Python for SAS Professionals

- Maan Parish





ISBN: 9798391325499 Inkstall Solutions LLP.



Python for SAS Professionals

Unlocking the Power of Python for Data Analysis and Analytics

Copyright © 2023 Inkstall Solutions

All rights reserved. No part of this book many be reproduced, stored in a retrieval system, or transmitted in any form or by any means, without the prior written permission of the publisher, excepting in the case of brief quotations embedded in critical articles or reviews.

Every effort has been made in the preparation of this book to ensure the accuracy of the information presented. However, the information contained in this book is sold without warranty, either express or implied. Neither the author, nor Inkstall Educare, and its dealers and distributors will be held liable for any damages caused or alleged to be caused directly or indirectly by this book.

Inkstall Educare has endeavoured to provide trademark information about all the companies and products mentioned in this book by the appropriate use of capitals. However, Inkstall Educare cannot guarantee the accuracy of this information.

First Published: April 2023 Published by Inkstall Solutions LLP. www.inkstall.us

Images used in this book are being borrowed, Inkstall doesn't hold any Copyright on the images been used. Questions about photos should be directed to: contact@inkstall.com



About Author:

Maan Parish

Maan Parish is a renowned data analyst and SAS professional with over 10 years of experience in the field. He has worked with various organizations and helped them with data-driven insights and analysis using SAS. With his vast experience in data analytics, he understands the challenges faced by SAS professionals and the need for Python in today's data-driven world.

Maan has authored several articles and white papers on data analysis and has been a regular speaker at various data analytics conferences. He is a strong advocate of Python and its advantages over traditional data analysis tools.

In his book, "Python for SAS Professionals," Maan shares his knowledge and insights on how SAS professionals can benefit from using Python. This book is a comprehensive guide that covers everything from the basics of Python to advanced topics such as data visualization, machine learning, and artificial intelligence.

Maan's expertise and experience make this book a must-read for SAS professionals who are looking to expand their skill set and stay relevant in today's data-driven world. With practical examples and real-world scenarios, this book will help readers understand the power of Python and how it can be used to enhance their data analysis and analytics capabilities.



Table of Contents

Chapter 1: Python Basics

- 1. Installing Python
- 2. Python environment
- 3. Python syntax and structure
- 4. Variables and data types
- 5. Conditional statements
- 6. Loops
- 7. Functions
- 8. Libraries and modules

Chapter 2: Data Structures in Python

- 1. Lists
- 2. Tuples
- 3. Dictionaries
- 4. Sets
- 5. Indexing and slicing
- 6. List comprehension
- 7. Sorting and filtering data

Chapter 3: Reading and Writing Data

- 1. Reading data from files
- 2. Writing data to files
- 3. CSV files
- 4. Excel files
- 5. JSON files
- 6. Working with databases



Chapter 4: Data Manipulation with Pandas

- 1. Introduction to Pandas
- 2. Creating a DataFrame
- 3. Reading data into a DataFrame
- 4. Indexing and selecting data
- 5. Filtering and sorting data
- 6. Aggregating and summarizing data
- 7. Merging and joining DataFrames

Chapter 5: Data Visualization with Matplotlib

- 1. Introduction to Matplotlib
- 2. Basic plots (line, scatter, bar)
- 3. Customizing plots (labels, colors, styles)
- 4. Multiple plots on one graph
- 5. Subplots and grids
- 6. Advanced plots (heatmaps, histograms, box plots)

Chapter 6: Machine Learning with Python

- 1. Introduction to machine learning
- 2. Preparing data for machine learning
- 3. Linear regression
- 4. Logistic regression
- 5. Decision trees
- 6. Random forests
- 7. K-Nearest Neighbors
- 8. Support Vector Machines
- 9. Naive Bayes



Chapter 7: Advanced Topics in Python

- 1. Object-oriented programming
- 2. Regular expressions
- 3. Working with dates and times
- 4. Web scraping
- 5. Introduction to Django
- 6. Creating a web application with Django
- 7. Creating a REST API
- 8. Next steps for SAS users who want to continue learning Python



Chapter 1: Python Basics



Python is a versatile programming language that has gained a lot of popularity in recent years. One of the reasons for its popularity is its ease of use, which makes it an ideal language for beginners. In this guide, we will provide an introduction to Python for SAS users.

Python Basics

Before we dive into the details of using Python for data analysis, let's go over some basic concepts in Python.

Variables

A variable is a container that holds a value. In Python, you can create a variable by assigning a value to it. Here's an example:

x = 10

This code creates a variable called x and assigns it the value 10.

Data Types

In Python, there are several built-in data types, including:

Integers: whole numbers, such as 10 or -5. Floating-point numbers: decimal numbers, such as 3.14 or -0.5. Strings: a sequence of characters enclosed in quotes, such as "hello" or "123". Booleans: True or False. You can use the type() function to check the data type of a variable. For example:

```
x = 10
print(type(x)) # Output: <class 'int'>
```

Lists

A list is a collection of values, enclosed in square brackets and separated by commas. Here's an example:

 $my_{list} = [1, 2, 3, 4, 5]$

You can access individual elements of a list using their index, which starts at 0. For example:

print(my list[0]) # Output: 1

You can also modify elements of a list by assigning a new value to them. For example: my_list[0] = 10 print(my_list) # Output: [10, 2, 3, 4, 5]



Loops

A loop is a way to repeat a block of code multiple times. In Python, there are two types of loops: for loops and while loops.

A for loop is used to iterate over a sequence, such as a list. Here's an example:

```
my_list = [1, 2, 3, 4, 5]
for num in my_list:
    print(num)
```

This code will print each element of the list my_list.

A while loop is used to repeat a block of code as long as a condition is true. Here's an example:

```
x = 0
while x < 10:
    print(x)
    x += 1</pre>
```

This code will print the numbers 0 through 9.

Functions

A function is a block of code that performs a specific task. In Python, you can define your own functions using the def keyword. Here's an example:

```
def add_numbers(x, y):
    return x + y
```

This code defines a function called add_numbers that takes two arguments and returns their sum.

Importing Modules

Python has a vast library of modules that provide additional functionality. You can import a module using the import keyword. Here's an example:

import math

This code imports the math module, which provides mathematical functions. Using Python with SAS

Python can be used in conjunction with SAS to perform data analysis tasks. Here are some ways you can use Python with SAS:



Use Python to read data from a file and write it to a SAS dataset. Use Python to perform data transformations and analysis

here's some more information on Python basics that may be useful for SAS users:

Variables and Data Types:

In Python, you can assign values to variables using the = operator. Python has several built-in data types, including integers, floats, strings, lists, tuples, and dictionaries. You can check the data type of a variable using the type() function.

Loops and Control Structures:

Python supports various control structures like if-else statements, for loops, and while loops. For example, you can use a for loop to iterate over a sequence of values, and an if statement to execute code based on a condition.

Functions and Modules:

Python allows you to define your own functions and reuse them in your code. You can also import external modules or libraries to extend the functionality of Python. For example, the pandas module provides a powerful way to work with tabular data in Python.

File Input/Output:

Python has built-in functions for reading and writing files, including CSV files, Excel files, and text files. You can use the open() function to open a file, and then read or write data using various methods.

Object-Oriented Programming:

Python is an object-oriented programming language, which means you can define classes and objects to encapsulate data and behavior. This can help organize your code and make it more modular and reusable.

Here's an example Python code that demonstrates some of these concepts:

```
# Define a function to calculate the sum of two numbers
def add_numbers(a, b):
    return a + b
# Define a class to represent a person
class Person:
    def __init__(self, name, age):
        self.name = name
        self.age = age
    def say_hello(self):
        print(f"Hello, my name is {self.name} and I am
{self.age} years old.")
```



```
# Create a list of numbers
numbers = [1, 2, 3, 4, 5]
# Use a for loop to iterate over the numbers and print
them
for num in numbers:
   print(num)
# Use an if statement to check if a number is even or
odd
if num % 2 == 0:
   print("The number is even.")
else:
   print("The number is odd.")
# Create a dictionary to store information about people
people = {"Alice": Person("Alice", 30), "Bob":
Person("Bob", 25) }
# Use a for loop to iterate over the people and call
their say hello method
for name, person in people.items():
   person.say hello()
# Write the numbers to a file
with open("numbers.txt", "w") as f:
    for num in numbers:
        f.write(str(num) + "\n")
```

This code defines a function to add two numbers, and a class to represent a person. It also demonstrates how to use a for loop to iterate over a list of numbers, an if statement to check if a number is even or odd, and a dictionary to store information about people.

The code also shows how to write the numbers to a file using the open() function and a with statement, which automatically closes the file when you're done with it. Code example that uses Python to read a CSV file, perform some data transformations, and write the results to a SAS dataset.

Import required modules import pandas as pd import numpy as np import saspy



```
# Create a SAS session
sas = saspy.SASsession()
# Read the CSV file into a Pandas DataFrame
data = pd.read csv("mydata.csv")
# Remove any rows with missing data
data.dropna(inplace=True)
# Calculate the average value of each column
avg values = data.mean()
# Convert the Pandas DataFrame to a SAS dataset
sas data = sas.df2sd(data, "mydata")
# Create a new SAS dataset with the average values
sas avg = sas.sd2df(avg values.to frame().T,
"avg values")
# Append the average values to the original dataset
sas data.append(sas avg)
# Save the dataset to a SAS library
sas data.save("mylib.mydata")
# Close the SAS session
sas.disconnect()
```

This code imports the pandas and numpy modules, which are popular libraries for working with data in Python. It also imports the saspy module, which provides a way to connect to a SAS server and interact with SAS data.

The code then creates a SAS session using the SASsession() function from the saspy module. It reads a CSV file into a Pandas DataFrame using the read_csv() function from the pandas module, and removes any rows with missing data using the dropna() method.

Next, the code calculates the average value of each column in the DataFrame using the mean() method. It then converts the Pandas DataFrame to a SAS dataset using the df2sd() method from the saspy module.

The code creates a new SAS dataset with the average values using the sd2df() method, and appends it to the original dataset using the append() method. It then saves the dataset to a SAS library using the save() method, and closes the SAS session using the disconnect() method.



This code demonstrates how Python can be used to read, transform, and analyze data, and how the results can be written to a SAS dataset for further analysis.

NumPy:

NumPy is a popular Python library for scientific computing that provides a powerful array processing capability. It is often used for numerical operations, such as matrix multiplication, and statistical calculations. You can use NumPy in conjunction with pandas for data analysis.

Pandas:

Pandas is another popular Python library for data manipulation and analysis. It provides a data frame object that allows you to store and manipulate tabular data in Python. You can perform various operations on data frames, such as filtering, sorting, and aggregating data.

Matplotlib:

Matplotlib is a plotting library for Python that allows you to create visualizations of your data. It provides a variety of chart types, including line charts, scatter plots, and bar charts. You can customize the appearance of your charts using various parameters.

Scikit-learn:

Scikit-learn is a machine learning library for Python that provides a wide range of algorithms for various tasks, such as classification, regression, and clustering. It also provides tools for data preprocessing and evaluation of machine learning models.

Jupyter Notebook:

Jupyter Notebook is an interactive web-based tool that allows you to write and run Python code in a document-like format. It is often used for data analysis and visualization, as it allows you to include charts and other visualizations inline with your code. It also allows you to document your code and share it with others.

Here's an example Python code that demonstrates some of these concepts:

```
import numpy as np
import pandas as pd
import matplotlib.pyplot as plt
from sklearn.linear_model import LinearRegression
# Create a NumPy array of random numbers
x = np.random.rand(100)
# Create a Pandas data frame from the array
```



```
df = pd.DataFrame({'x': x})
# Add a column to the data frame with a calculated
value
df['y'] = 2 * df['x'] + np.random.randn(100) * 0.1
# Create a scatter plot of the data
plt.scatter(df['x'], df['y'])
plt.xlabel('X')
plt.ylabel('Y')
plt.show()
# Fit a linear regression model to the data
model = LinearRegression()
model.fit(df[['x']], df['y'])
# Print the slope and intercept of the model
print('Slope:', model.coef [0])
print('Intercept:', model.intercept )
# Predict values for new data points
x new = np.array([0.2, 0.4, 0.6, 0.8])
y new = model.predict(x new.reshape(-1, 1))
# Print the predicted values
print('Predicted values:', y new)
```

This code uses NumPy to generate a random array of numbers, and Pandas to create a data frame from the array. It then adds a column to the data frame with a calculated value, and creates a scatter plot of the data using Matplotlib.

The code also uses Scikit-learn to fit a linear regression model to the data, and predict values for new data points. Finally, it prints the slope and intercept of the model, and the predicted values.

Installing Python

Installing Python is the first step in learning how to use Python for SAS users. Python is a popular and versatile programming language that can be used for a wide range of tasks, including data analysis, machine learning, and web development.



There are several ways to install Python, but the easiest and most popular method is to use a Python distribution that includes all the necessary packages and tools. One such distribution is Anaconda, which is available for Windows, macOS, and Linux.

To install Anaconda, follow these steps:

Go to the Anaconda website (https://www.anaconda.com/products/individual) and download the installer for your operating system.

Run the installer and follow the prompts to complete the installation.

Once the installation is complete, open the Anaconda Navigator application. This will give you access to the Anaconda command prompt and the Anaconda Navigator GUI.

To start using Python, open the Anaconda command prompt and type "python". This will start the Python interpreter, which you can use to run Python code.

Alternatively, you can use a Python IDE (Integrated Development Environment) such as Spyder or PyCharm, which provide a more user-friendly interface for writing and running Python code.

Once you have installed Python, you can start learning how to use it for data analysis and other tasks. One way to get started is to learn how to use the pandas library, which is a popular Python library for data manipulation and analysis.

Here is an example code that uses the pandas library to read a CSV file and perform some basic data analysis:

```
import pandas as pd
# Read the CSV file
df = pd.read_csv('data.csv')
# Print the first 5 rows of the data
print(df.head())
# Get some basic statistics about the data
print(df.describe())
# Group the data by a specific column and get the mean
of another column
print(df.groupby('Column1')['Column2'].mean())
```

In this code, the "pd" alias is used to refer to the pandas library. The "read_csv" function is used to read a CSV file and create a pandas DataFrame object. The "head" and "describe" methods are used to print the first few rows of the data and some basic statistics, respectively. Finally, the



"groupby" method is used to group the data by a specific column and get the mean of another column.

Learning how to use Python can be a valuable skill for SAS users, as it provides a powerful and flexible toolset for data analysis and other tasks.

Python is a powerful and versatile programming language that has become increasingly popular in recent years, particularly in the field of data science and analysis. For SAS users, learning Python can provide a valuable addition to their toolkit, as it offers a wide range of capabilities for data analysis, machine learning, and other tasks.

Example code that uses Python and the pandas library to analyze a dataset:

```
import pandas as pd
import matplotlib.pyplot as plt
# Read the CSV file into a pandas DataFrame
df = pd.read csv('sales data.csv')
# Create a new column for total sales
df['Total Sales'] = df['Quantity'] * df['Price']
# Calculate the total sales by region
sales by region = df.groupby('Region')['Total
Sales'].sum()
# Calculate the total sales by product
sales by product = df.groupby('Product')['Total
Sales'].sum()
# Plot a bar chart of total sales by region
plt.bar(sales by region.index, sales by region.values)
plt.title('Total Sales by Region')
plt.xlabel('Region')
plt.ylabel('Total Sales')
plt.show()
# Plot a pie chart of total sales by product
plt.pie(sales by product.values,
labels=sales by product.index, autopct='%1.1f%%')
plt.title('Total Sales by Product')
plt.show()
```



In this code, we start by importing the pandas library and the matplotlib.pyplot module, which provides functions for creating visualizations. We then use the pd.read_csv() function to read a CSV file containing sales data into a pandas DataFrame.

Next, we create a new column in the DataFrame for total sales by multiplying the Quantity and Price columns. We then use the groupby() method to calculate the total sales by region and by product.

Finally, we use the plt.bar() function to create a bar chart of the total sales by region, and the plt.pie() function to create a pie chart of the total sales by product.

This code demonstrates how Python and the pandas library can be used to read, manipulate, and analyze data, and how the matplotlib library can be used to create visualizations.

One advantage of Python over SAS is its flexibility and ease of use. Python has a simpler and more intuitive syntax, making it easier to learn and use, particularly for users who are new to programming. In addition, Python has a large and active community of developers, which means that there are many resources and libraries available to help users solve problems and build applications.

Another advantage of Python for SAS users is its extensive library ecosystem. Python has a vast number of libraries for data manipulation, analysis, and visualization, including popular libraries such as NumPy, Pandas, Matplotlib, and Seaborn. These libraries provide powerful and flexible tools for data analysis and visualization, allowing users to quickly and easily perform complex data manipulations and generate high-quality visualizations.

Python also has a strong focus on machine learning and artificial intelligence, with many popular libraries such as TensorFlow, Keras, and Scikit-learn providing powerful tools for machine learning and deep learning. This makes Python a popular choice for data scientists and machine learning engineers, who can use Python to build and deploy advanced machine learning models.

In terms of integration with SAS, Python has several tools and libraries available for working with SAS data and environments. For example, the SASPy library allows users to connect to SAS environments and run SAS code from within Python, while the SAS Integration for Python (SASPy) package provides a bridge between Python and SAS Viya, allowing users to seamlessly integrate SAS analytics into their Python applications.

Learning Python can provide a valuable addition to the toolkit of SAS users, allowing them to perform advanced data analysis, visualization, and machine learning tasks with ease. While there is a learning curve associated with learning Python, the benefits in terms of flexibility, ease of use, and the availability of powerful libraries make it a worthwhile investment for SAS users looking to expand their skillset.

Here is another longer example code that shows how to use Python and the scikit-learn library to build and evaluate a machine learning model:

import pandas as pd



```
from sklearn.model selection import train test split
from sklearn.ensemble import RandomForestClassifier
from sklearn.metrics import accuracy score,
confusion matrix
# Read the CSV file into a pandas DataFrame
df = pd.read csv('iris.csv')
# Split the data into training and testing sets
X train, X test, y train, y test =
train test split(df.drop('species', axis=1),
df['species'], test size=0.2, random state=42)
# Build a random forest classifier model
clf = RandomForestClassifier(n estimators=100,
random state=42)
clf.fit(X train, y train)
# Make predictions on the test set
y pred = clf.predict(X test)
# Calculate the accuracy of the model
accuracy = accuracy score(y test, y pred)
print('Accuracy:', accuracy)
# Create a confusion matrix to evaluate the performance
of the model
cm = confusion matrix(y test, y pred)
print('Confusion Matrix:')
print(cm)
```

In this code, we start by importing the pandas library, as well as several modules from the scikitlearn library, including train_test_split for splitting the data into training and testing sets, RandomForestClassifier for building a random forest classifier model, and accuracy_score and confusion_matrix for evaluating the performance of the model.

We then use the pd.read_csv() function to read a CSV file containing iris data into a pandas DataFrame, and split the data into training and testing sets using the train_test_split() function.

Next, we build a random forest classifier model using the RandomForestClassifier() function and the training data, and make predictions on the test set using the predict() method. We then calculate the accuracy of the model using the accuracy_score() function, and create a confusion matrix using the confusion_matrix() function to evaluate the performance of the model.

in stal

This code demonstrates how Python and the scikit-learn library can be used to build and evaluate a machine learning model. It shows how to split the data into training and testing sets, build a model using the training data, and evaluate the performance of the model using metrics such as accuracy and confusion matrix.

Python environment

Python has become an increasingly popular language in the data analytics and data science community due to its flexibility, ease of use, and powerful libraries. SAS users who are interested in learning Python can benefit from the language's vast ecosystem of libraries and tools that are specifically designed for data analysis and manipulation.

Setting up a Python environment for data analysis and manipulation can be done in a few simple steps. First, download and install the latest version of Python from the official Python website (https://www.python.org/downloads/). Once Python is installed, you can install additional packages and libraries using pip, a package manager for Python. For example, to install the pandas library, open a terminal or command prompt and enter the following command:

pip install pandas

Similarly, you can install other libraries like NumPy, Matplotlib, SciPy, and scikit-learn by replacing "pandas" with the name of the library you want to install.

SAS users who are new to Python may find it helpful to use Jupyter Notebook, a web-based interactive computing environment that allows users to write and execute Python code in a browser-based interface. Jupyter Notebook also provides support for creating interactive visualizations, integrating with other programming languages, and sharing notebooks with others.

To install Jupyter Notebook, you can use pip by entering the following command:

pip install jupyter

Once installed, you can start Jupyter Notebook by entering the following command in a terminal or command prompt:

jupyter notebook

This will launch a web interface where you can create, open, and edit notebooks.

Python users who are already familiar with SAS may find it helpful to use the SASPy library, which provides a Python interface to SAS. SASPy allows users to interact with SAS datasets and



perform SAS procedures and functions from within Python. To use SASPy, you must have a valid installation of SAS on your system.

To install SASPy, you can use pip by entering the following command:

```
pip install saspy
```

Once installed, you can create a connection to a SAS server by entering the following code:

```
import saspy
sas = saspy.SASsession(cfgname='default')
```

This will create a SAS session object that can be used to execute SAS code from within Python. For example, the following code creates a SAS dataset and prints the first five observations:

```
sas.submit('data test; input x y; datalines; 1 2 3 4 5
6 ; run;')
sas.submit('proc print data=test (obs=5); run;')
```

This code will execute the SAS code within the SAS session and print the first five observations of the "test" dataset.

SAS users who are interested in learning Python can benefit from the language's powerful libraries and tools for data analysis and manipulation. Python can be easily installed using pip, and Jupyter Notebook can be used for interactive computing. SAS users who are already familiar with SAS may find SASPy helpful for integrating Python with SAS.

Python is a general-purpose programming language that is easy to learn and use. It has a large and active community that has developed a vast ecosystem of libraries and tools for data analysis, machine learning, web development, scientific computing, and more. Python's flexibility and ease of use make it an ideal language for data analysis and manipulation, which is why it has become popular in the data science and analytics community.

SAS is a powerful software suite that has been a dominant player in the data analytics industry for many years. SAS provides a comprehensive set of tools and features for data preparation, analysis, and visualization. However, some SAS users are now turning to Python for data analysis because of Python's versatility and growing popularity.

Python and SAS have many similarities, but also some key differences. One significant difference is that SAS is a proprietary software suite that requires a license, while Python is open source and free to use. Python's open source nature means that it has a vast community of developers and users who contribute to its development and use. Python also has a larger ecosystem of libraries and tools than SAS, which means that Python users have access to a more extensive range of functionality and features.



Python's powerful libraries for data analysis and manipulation, such as NumPy, Pandas, Matplotlib, and SciPy, allow users to perform complex data operations with ease. NumPy provides powerful numerical operations, Pandas allows for flexible data manipulation and analysis, Matplotlib provides extensive visualization capabilities, and SciPy provides an extensive set of scientific computing tools.

Python's machine learning libraries, such as scikit-learn, TensorFlow, and Keras, allow users to build and train powerful machine learning models. Scikit-learn provides a wide range of machine learning algorithms, TensorFlow and Keras provide powerful deep learning capabilities, and PyTorch provides a flexible deep learning framework for building neural networks.

Python's web development frameworks, such as Django and Flask, allow users to build powerful and scalable web applications. Django provides a full-stack web development framework that includes a robust database ORM, while Flask provides a lightweight framework for building smaller web applications.

SAS users who are interested in learning Python can start by learning the basics of Python syntax and data types. Once comfortable with the language basics, SAS users can explore Python libraries and tools for data analysis, machine learning, and web development. The Jupyter Notebook is an excellent tool for learning Python because it allows for interactive computing and code sharing.

SASPy is a library that allows users to execute SAS code from within Python. SASPy provides a Python interface to SAS and allows users to interact with SAS datasets, perform SAS procedures, and call SAS functions from within Python. SASPy can be helpful for users who are already familiar with SAS and want to integrate Python into their workflows.

Python is a powerful programming language that has become popular in the data science and analytics community due to its flexibility, ease of use, and powerful libraries and tools. SAS users who are interested in learning Python can benefit from its extensive ecosystem of libraries and tools for data analysis, machine learning, and web development. SASPy can be helpful for users who want to integrate Python into their SAS workflows.

Here is an example code in Python for loading and manipulating data using the Pandas library, which is commonly used in data analysis:

```
import pandas as pd
# Load data from a CSV file
data = pd.read_csv("data.csv")
# View the first five rows of the data
print(data.head())
```



```
# Subset the data to include only rows where the 'age'
column is greater than 30
subset_data = data[data['age'] > 30]
# Group the data by the 'gender' column and calculate
the mean value of the 'salary' column
grouped_data = data.groupby('gender')['salary'].mean()
# Create a new column that combines the 'first_name'
and 'last_name' columns
data['full_name'] = data['first_name'] + ' ' +
data['last_name']
# Save the modified data to a new CSV file
data.to_csv("modified_data.csv")
```

In this code, we first import the Pandas library and use the read_csv function to load data from a CSV file into a Pandas DataFrame. We then use the head method to view the first five rows of the data.

Next, we use the square bracket notation to subset the data to include only rows where the 'age' column is greater than 30. We then group the data by the 'gender' column using the groupby method and calculate the mean value of the 'salary' column using the mean method.

We then create a new column called 'full_name' by combining the 'first_name' and 'last_name' columns using string concatenation. Finally, we use the to_csv method to save the modified data to a new CSV file.

This is just a simple example of what can be done with Pandas. The library provides a wide range of functions and methods for manipulating and analyzing data, including merging and joining datasets, handling missing data, and performing statistical analysis.

In the first line of the code, we use the import keyword to import the Pandas library, which we give the alias pd. This alias makes it easier to refer to the library in subsequent code.

```
import pandas as pd
```

The next line of code uses the read_csv function to read in a CSV file called "data.csv" and store the contents in a Pandas DataFrame called data.

```
data = pd.read csv("data.csv")
```

We then use the head method to display the first five rows of the data DataFrame.



print(data.head())

The next line of code uses Boolean indexing to create a new DataFrame called subset_data that contains only the rows where the age column is greater than 30.

```
subset data = data[data['age'] > 30]
```

We then use the groupby method to group the data DataFrame by the gender column and calculate the mean salary for each gender using the mean method. The resulting object is a Pandas Series.

```
grouped data = data.groupby('gender')['salary'].mean()
```

Next, we use string concatenation to create a new column in the data DataFrame called full_name. This column combines the first_name and last_name columns with a space between them.

```
data['full_name'] = data['first_name'] + ' ' +
data['last_name']
```

Finally, we use the to_csv method to save the modified data DataFrame to a new CSV file called "modified_data.csv".

```
data.to csv("modified data.csv")
```

Overall, this example demonstrates some of the basic data manipulation operations that can be performed using the Pandas library in Python. Pandas is a powerful and flexible library that allows users to work with a variety of data formats and perform complex data operations with ease.

In the line data = pd.read_csv("data.csv"), pd.read_csv() is a function provided by the Pandas library that reads in data from a CSV file and returns a DataFrame object. A DataFrame is a two-dimensional table-like data structure that can store and manipulate data in rows and columns.

The head() method in the line print(data.head()) is a method provided by the DataFrame object that displays the first few rows of the DataFrame. By default, it displays the first five rows, but you can pass a parameter to specify a different number of rows to display.

In the line subset_data = data[data['age'] > 30], we are using Boolean indexing to filter the DataFrame to only include rows where the age column is greater than 30. Boolean indexing is a powerful feature of Pandas that allows you to subset a DataFrame based on a condition expressed as a Boolean expression.

In the line grouped_data = data.groupby('gender')['salary'].mean(), we are using the groupby() method to group the DataFrame by the gender column and then calculate the mean value of the



salary column for each group. The result is a Pandas Series object that stores the mean salaries for each gender group.

In the line data['full_name'] = data['first_name'] + ' ' + data['last_name'], we are creating a new column in the DataFrame called full_name by concatenating the first_name and last_name columns with a space in between. This is a simple example of how you can use Pandas to create new columns from existing data.

Finally, in the line data.to_csv("modified_data.csv"), we are using the to_csv() method to save the modified DataFrame to a new CSV file called "modified_data.csv". This method writes the DataFrame to a CSV file with the specified filename.

This code example demonstrates some of the basic data manipulation operations that you can perform with Pandas. Pandas is a powerful and flexible library that provides many more features and functions for working with data, including data cleaning, reshaping, merging, and more. In the line import pandas as pd, we are importing the Pandas library and giving it the alias pd. This is a common convention in Python programming, as it allows you to refer to the library with a shorter name in your code.

In the line data = $pd.read_csv("data.csv")$, we are using the $read_csv()$ function provided by Pandas to read in a CSV file called "data.csv" and store its contents in a DataFrame object called data. The $read_csv()$ function has many options for handling different types of CSV files, such as files with different delimiters or missing values.

In the line print(data.head()), we are using the head() method of the data DataFrame to display the first few rows of the DataFrame. By default, head() displays the first five rows, but you can pass an argument to specify a different number of rows to display.

In the line subset_data = data[data['age'] > 30], we are using Boolean indexing to create a new DataFrame called subset_data that contains only the rows where the age column is greater than 30. Boolean indexing is a powerful feature of Pandas that allows you to filter a DataFrame based on a condition expressed as a Boolean expression.

In the line grouped_data = data.groupby('gender')['salary'].mean(), we are using the groupby() method of the data DataFrame to group the data by the gender column, and then using the mean() method to calculate the mean value of the salary column for each group. The result is a Pandas Series object that stores the mean salaries for each gender group.

In the line data['full_name'] = data['first_name'] + ' ' + data['last_name'], we are creating a new column in the data DataFrame called full_name. We do this by using string concatenation to combine the first_name and last_name columns, separated by a space.

Finally, in the line data.to_csv("modified_data.csv"), we are using the to_csv() method of the data DataFrame to write the modified DataFrame to a new CSV file called "modified_data.csv". The to_csv() method has many options for specifying the format and structure of the output file, such as the delimiter character, the header row, and the encoding.



In summary, this code example demonstrates some of the basic data manipulation operations that you can perform with Pandas. Pandas is a widely-used and powerful library for data analysis in Python, and it provides many more features and functions for working with data.

```
import pandas as pd
# Load data from a CSV file
data = pd.read csv("data.csv")
# View the first five rows of the data
print(data.head())
# Subset the data to include only rows where the
                                                       'age'
column is greater than 30
subset data = data[data['age'] > 30]
# Group the data by the 'gender' column and calculate the
mean value of the 'salary' column
grouped data = data.groupby('gender')['salary'].mean()
# Create a new column that combines the 'first name'
                                                         and
'last name' columns
data['full name']
                        data['first name']
                    =
                                              +
                                                           +
data['last name']
# Save the modified data to a new CSV file
data.to csv("modified data.csv")
```

In this code example, we start by importing the Pandas library and reading in some sales data from a CSV file. We then use some basic Pandas methods to get information about the data, including its shape, columns, and data types.

Next, we filter the data to only include rows where the sales amount is greater than 1000, using Boolean indexing. We then use the groupby() method of the DataFrame to group the data by the region column and calculate the mean sales amount for each region.

We then create a new column in the data called total_sales, which calculates the total sales for each row by multiplying the sales amount by the quantity. Finally, we save the modified data to a new CSV file using the to_csv() method of the DataFrame.

This code example demonstrates some of the basic data manipulation operations that you can perform with Pandas, including filtering, grouping, and creating new columns. Pandas provides many more features and functions for working with data, including data cleaning, merging, reshaping, and more.



In the first line of the code, we use the import keyword to import the Pandas library, which we give the alias pd. This alias makes it easier to refer to the library in subsequent code.

import pandas as pd

The next line of code uses the read_csv function to read in a CSV file called "data.csv" and store the contents in a Pandas DataFrame called data.

data = pd.read csv("data.csv")

We then use the head method to display the first five rows of the data DataFrame.

```
print(data.head())
```

The next line of code uses Boolean indexing to create a new DataFrame called subset_data that contains only the rows where the age column is greater than 30.

subset data = data[data['age'] > 30]

We then use the groupby method to group the data DataFrame by the gender column and calculate the mean salary for each gender using the mean method. The resulting object is a Pandas Series.

```
grouped data = data.groupby('gender')['salary'].mean()
```

Next, we use string concatenation to create a new column in the data DataFrame called full_name. This column combines the first_name and last_name columns with a space between them.

```
data['full_name'] = data['first_name'] + ' ' +
data['last_name']
```

Finally, we use the to_csv method to save the modified data DataFrame to a new CSV file called "modified_data.csv".

```
data.to csv("modified data.csv")
```

This example demonstrates some of the basic data manipulation operations that can be performed using the Pandas library in Python. Pandas is a powerful and flexible library that allows users to work with a variety of data formats and perform complex data operations with ease.

In the line data = pd.read_csv("data.csv"), pd.read_csv() is a function provided by the Pandas library that reads in data from a CSV file and returns a DataFrame object. A DataFrame is a two-dimensional table-like data structure that can store and manipulate data in rows and columns.

The head() method in the line print(data.head()) is a method provided by the DataFrame object that displays the first few rows of the DataFrame. By default, it displays the first five rows, but you can pass a parameter to specify a different number of rows to display.

In the line subset_data = data[data['age'] > 30], we are using Boolean indexing to filter the DataFrame to only include rows where the age column is greater than 30. Boolean indexing is a powerful feature of Pandas that allows you to subset a DataFrame based on a condition expressed as a Boolean expression.

In the line grouped_data = data.groupby('gender')['salary'].mean(), we are using the groupby() method to group the DataFrame by the gender column and then calculate the mean value of the salary column for each group. The result is a Pandas Series object that stores the mean salaries for each gender group.

In the line data['full_name'] = data['first_name'] + ' ' + data['last_name'], we are creating a new column in the DataFrame called full_name by concatenating the first_name and last_name columns with a space in between. This is a simple example of how you can use Pandas to create new columns from existing data.

Finally, in the line data.to_csv("modified_data.csv"), we are using the to_csv() method to save the modified DataFrame to a new CSV file called "modified_data.csv". This method writes the DataFrame to a CSV file with the specified filename.

Here's some more code for working with data using Pandas:

```
import pandas as pd
# Read in the data from a CSV file
data = pd.read_csv("sales_data.csv")
# Get some basic information about the data
print("Shape of the data:", data.shape)
print("Columns in the data:", data.columns)
print("Data types of the columns:\n", data.dtypes)
# Filter the data to only include rows where the sales
amount is greater than 1000
subset_data = data[data["sales_amount"] > 1000]
# Group the data by the region column and calculate the
```

Group the data by the region column and calculate the mean sales amount for each region



```
grouped_data =
data.groupby("region")["sales_amount"].mean()
# Create a new column in the data that calculates the
total sales for each row
data["total_sales"] = data["sales_amount"] *
data["quantity"]
# Save the modified data to a new CSV file
data.to csv("modified sales data.csv")
```

In this code example, we start by importing the Pandas library and reading in some sales data from a CSV file. We then use some basic Pandas methods to get information about the data, including its shape, columns, and data types.

Next, we filter the data to only include rows where the sales amount is greater than 1000, using Boolean indexing. We then use the groupby() method of the DataFrame to group the data by the region column and calculate the mean sales amount for each region.

We then create a new column in the data called total_sales, which calculates the total sales for each row by multiplying the sales amount by the quantity. Finally, we save the modified data to a new CSV file using the to_csv() method of the DataFrame.

This code example demonstrates some of the basic data manipulation operations that you can perform with Pandas, including filtering, grouping, and creating new columns. Pandas provides many more features and functions for working with data, including data cleaning, merging, reshaping, and more.

Here's an example of setting up a Python environment using Anaconda:

- 1. Download and install Anaconda from the Anaconda website (https://www.anaconda.com/products/individual)
- 2. Open the Anaconda Navigator application. This will give you a graphical user interface for managing your Python environments and packages.
- 3. Click on the "Environments" tab on the left-hand side of the window. This will show you a list of your existing environments.
- 4. Click on the "Create" button at the bottom of the window to create a new environment.
- 5. Give your environment a name, such as "myenv", and choose a Python version to use (e.g., Python 3.9).
- 6. Choose any additional packages or libraries that you want to include in your environment, such as Pandas, Numpy, or Matplotlib.
- 7. Click the "Create" button to create your new environment.
- 8. Once your environment is created, you can activate it by clicking the "Play" button next to its name. This will open a new terminal window with your environment activated.
- 9. You can now install additional packages or libraries using the **conda** package manager or the **pip** package manager.



Here's an example of using conda to install the Pandas library in your new environment:

Open a terminal window with your environment activated.

Type the following command to update the conda package manager:

conda update conda

Type the following command to install the Pandas library:

conda install pandas

Once the installation is complete, you can test that Pandas is installed correctly by opening a Python interpreter and importing the library:

```
python
import pandas as pd
```

If you don't get any errors, then Pandas is installed correctly and you can start using it in your Python scripts.

Creating a virtual environment with virtualenv

```
# Install virtualenv using pip
pip install virtualenv
# Create a new virtual environment called myenv
virtualenv myenv
# Activate the virtual environment
source myenv/bin/activate
# Install packages in the virtual environment using pip
pip install pandas numpy matplotlib
# Deactivate the virtual environment
```

Deactivate

In this example, we're using the virtualenv package to create a new virtual environment called myenv. We then activate the environment using the source command and install some packages using pip. Finally, we deactivate the virtual environment using the deactivate command.

Using pipenv to manage dependencies

```
# Install pipenv using pip
```



```
pip install pipenv
# Create a new project and virtual environment with
pipenv
pipenv install pandas numpy matplotlib
# Activate the virtual environment
pipenv shell
# Install additional packages
pipenv install requests
# Deactivate the virtual environment
Exit
```

In this example, we're using pipenv to manage our project dependencies. We create a new project and virtual environment using pipenv install and specify the packages that we want to include. We then activate the virtual environment using pipenv shell, install additional packages using pipenv install, and finally exit the virtual environment using exit.

Using conda to create and manage environments

```
# Create a new environment with conda
conda create --name myenv pandas numpy matplotlib
# Activate the environment
conda activate myenv
# Install additional packages
conda install requests
# Deactivate the environment
conda deactivate
```

In this example, we're using conda to create and manage our environment. We create a new environment using conda create and specify the packages that we want to include. We then activate the environment using conda activate, install additional packages using conda install, and finally deactivate the environment using conda deactivate.



Python syntax and structure

Python is a high-level, general-purpose programming language that is widely used in data analysis, machine learning, and scientific computing. This section provides an introduction to Python syntax and structure from the perspective of SAS users.

Variables and Data Types:

In Python, variables are created by assigning a value to a name. Unlike SAS, Python is a dynamically typed language, meaning that the type of a variable is determined at runtime based on the value that is assigned to it. Some of the commonly used data types in Python include integers, floating-point numbers, strings, and lists. Here are some examples:

```
# Assigning values to variables
x = 5
y = 3.14
z = 'hello'
# Creating a list
my_list = [1, 2, 3, 'four']
```

Operators:

Python supports a variety of operators for performing arithmetic, logical, and comparison operations. Here are some examples:

```
# Arithmetic operators
x = 5 + 3 \# addition
y = 5 - 3 \# subtraction
z = 5 * 3 # multiplication
w = 5 / 3 \# division
v = 5 \% 3 \# modulus (remainder)
u = 5 ** 3 # exponentiation
# Logical operators
a = True
b = False
c = a and b # logical AND
             # logical OR
d = a \text{ or } b
e = not a
             # logical NOT
# Comparison operators
f = 5 == 3
             # equality
q = 5 != 3
             # inequality
```



Conditional Statements:

Python supports conditional statements, which allow you to execute different blocks of code depending on whether a certain condition is true or false. The syntax for a conditional statement in Python is similar to that in SAS:

```
# Example of an if statement
x = 5
if x > 3:
    print('x is greater than 3')
else:
    print('x is less than or equal to 3')
```

Loops:

Python supports two types of loops: for loops and while loops. A for loop is used to iterate over a sequence of values, while a while loop is used to repeat a block of code while a certain condition is true. Here are some examples:

```
# Example of a for loop
my_list = [1, 2, 3, 4, 5]
for i in my_list:
    print(i)
# Example of a while loop
x = 0
while x < 5:
    print(x)
    x += 1</pre>
```

Functions:

Python allows you to define your own functions, which can be used to encapsulate a block of code and make it reusable. Here is an example:

```
# Example of a function
def add_numbers(x, y):
    return x + y
```



```
# Example of calling the function
result = add_numbers(5, 3)
print(result) # prints 8
```

Modules:

Python provides a large standard library, as well as many third-party libraries, that can be used to extend the functionality of the language. To use a module in your code, you simply need to import it. Here is an example:

```
# Example of importing a module
import math
# Example of using a function from the module
result = math.sqrt(16)
print(result) # prints 4.0
```

Working with Data in Python:

Python provides several libraries for working with data, including NumPy, pandas, and matplotlib. NumPy provides support for working with arrays and matrices, while pandas provides support for working with tabular data. Matplotlib provides support for creating graphs and visualizations. Here are some examples:

```
# Example of using NumPy to create an array
import numpy as np
my_array = np.array([1, 2, 3, 4, 5])
print(my_array)
# Example of using pandas to read a CSV file
import pandas as pd
my_data = pd.read_csv('my_data.csv')
print(my_data.head())
# Example of using matplotlib to create a histogram
import matplotlib.pyplot as plt
my_data = pd.read_csv('my_data.csv')
plt.hist(my_data['my_column'])
plt.show()
```



Using Python with SAS:

Python can be used in conjunction with SAS to perform data analysis and other tasks. SAS provides several ways to use Python, including the SASPy module, which allows you to run Python code from within a SAS program. Here is an example:

```
/* Example of using Python with SAS */
options python='/usr/local/bin/python3';
/* Importing the SASPy module */
proc python;
    import saspy;
run;
/* Running Python code from within a SAS program */
proc python;
    submit;
    import pandas as pd
    my_data = pd.read_csv('my_data.csv')
    print(my_data.head())
    endsubmit;
run;
```

Python is a powerful and versatile language that can be used for a wide range of data analysis and scientific computing tasks. While it may take some time to get used to the syntax and structure of Python, SAS users should find it relatively straightforward to learn. By using Python in conjunction with SAS, users can take advantage of the strengths of both languages to tackle complex data analysis problems.

Python Syntax and Structure:

- 1. Indentation: Python uses indentation to define blocks of code instead of using curly braces or other delimiters. This means that the indentation level of a line of code determines which block it belongs to. For example, a block of code that belongs to a function or conditional statement is indented by four spaces.
- 2. Comments: Python supports both single-line and multi-line comments. Single-line comments begin with the hash symbol (#), while multi-line comments are enclosed in triple quotes (""").
- 3. Case sensitivity: Python is case-sensitive, meaning that variables and function names must be spelled exactly the same way every time they are used.
- Variables: Variables in Python are dynamically typed, which means that their type can change at runtime. To create a variable, simply assign a value to it using the equals sign (=). Python supports a variety of data types, including integers, floating-point numbers, strings, lists, and dictionaries.



- 5. Operators: Python supports a wide range of operators, including arithmetic operators (+, -, *, /, %), comparison operators (==, !=, >, <, >=, <=), logical operators (and, or, not), and assignment operators (+=, -=, *=, /=, %=).
- 6. Conditional statements: Python supports conditional statements, including if, elif, and else. These statements allow you to execute different blocks of code depending on whether a certain condition is true or false.
- 7. Loops: Python supports two types of loops: for loops and while loops. For loops are used to iterate over a sequence of values, while while loops are used to repeat a block of code while a certain condition is true.
- 8. Functions: Python allows you to define your own functions, which can be used to encapsulate a block of code and make it reusable. To define a function, use the def keyword, followed by the function name and any arguments it takes.
- 9. Modules: Python provides a large standard library, as well as many third-party libraries, that can be used to extend the functionality of the language. To use a module in your code, simply import it using the import statement.
- 10. Working with data: Python provides several libraries for working with data, including NumPy, pandas, and matplotlib. NumPy provides support for working with arrays and matrices, while pandas provides support for working with tabular data. Matplotlib provides support for creating graphs and visualizations.
- 11. Using Python with SAS: Python can be used in conjunction with SAS to perform data analysis and other tasks. SAS provides several ways to use Python, including the SASPy module, which allows you to run Python code from within a SAS program.

Python is a powerful and flexible language that offers many benefits to SAS users. Its syntax and structure may take some getting used to, but with practice, SAS users can quickly become proficient in Python and use it to tackle complex data analysis problems.

Here's a longer example of Python code that demonstrates some of the concepts discussed above:

```
# Example of using Python to analyze data
import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
# Load data from a CSV file
data = pd.read_csv('data.csv')
# Clean the data by removing missing values
data = data.dropna()
# Calculate some summary statistics
mean = np.mean(data['value'])
median = np.median(data['value'])
std_dev = np.std(data['value'])
```



```
# Print out the results
print('Mean value: ', mean)
print('Median value: ', median)
print('Standard deviation: ', std dev)
# Create a histogram of the data
plt.hist(data['value'], bins=20)
plt.title('Distribution of values')
plt.xlabel('Value')
plt.ylabel('Frequency')
plt.show()
# Define a function to calculate the sum of a list of
numbers
def sum list(numbers):
    result = 0
    for num in numbers:
        result += num
    return result
# Use the function to calculate the sum of a list of
numbers
numbers = [1, 2, 3, 4, 5]
sum = sum list(numbers)
print('Sum of numbers:', sum)
```

This code loads data from a CSV file using the pandas library, cleans the data by removing missing values, calculates some summary statistics using the NumPy library, and creates a histogram of the data using the matplotlib library. It also defines a function to calculate the sum of a list of numbers and uses that function to calculate the sum of a list of numbers.

Information specifically for SAS users who are learning Python:

- 1. Python has a syntax that is different from SAS, so it may take some time to get used to the new syntax. For example, instead of using a semicolon (;) to end a statement in SAS, Python uses a newline character or a colon (:).
- 2. SAS datasets can be read and written in Python using libraries such as pandas and PySAS. Pandas is a popular library for data manipulation and analysis, while PySAS provides a bridge between SAS datasets and Python.
- 3. The data step in SAS can be emulated in Python using the pandas library. The pandas library provides a similar functionality to the data step in SAS, including data selection, transformation, and aggregation.



- 4. Python is a versatile language that can be used for a wide range of tasks beyond data analysis, such as web development, machine learning, and scientific computing. As a SAS user, you may find these additional capabilities of Python useful.
- 5. Python has a large and active community of developers who are constantly creating new libraries and tools to extend the functionality of the language. As a SAS user, you can leverage these libraries to enhance your data analysis capabilities.
- 6. Python is often used in conjunction with SAS, and there are several ways to integrate the two. For example, SAS users can use the SASPy module to run Python code from within a SAS program or use the SAS macro language to generate Python code.
- 7. Python has a strong emphasis on readability and maintainability, which can make it easier to collaborate on code and maintain large codebases over time. SAS users may find that the switch to Python results in more readable and maintainable code.

While there may be some differences in syntax and structure between SAS and Python, Python offers a powerful and versatile toolset for data analysis that can complement and extend the capabilities of SAS. With some practice, SAS users can become proficient in Python and leverage its features to enhance their data analysis workflows.

Reading and writing SAS datasets using PySAS:

```
import pysas
# Read a SAS dataset
data = pysas.read_sas('data.sas7bdat')
# Write a SAS dataset
pysas.write_sas(data, 'output.sas7bdat')
```

Using the data step functionality in pandas to select and transform data:

```
import pandas as pd
# Read a CSV file
data = pd.read_csv('data.csv')
# Select a subset of the data
subset = data.loc[data['category'] == 'A']
# Transform the data by creating a new column
subset['new_column'] = subset['value'] * 2
# Aggregate the data by calculating the mean value
mean value = subset['value'].mean()
```



```
Using Python's built-in functions for mathematical
operations:
makefile
Copy code
# Calculate the square root of a number
import math
sqrt_value = math.sqrt(16)
# Calculate the exponential of a number
exp_value = math.exp(2)
# Calculate the logarithm of a number
log_value = math.log(10)
# Generate a random number
import random
rand_value = random.randint(1, 100)
```

Using Python libraries for machine learning:

```
import pandas as pd
from sklearn.linear_model import LinearRegression
# Load data from a CSV file
data = pd.read_csv('data.csv')
# Split the data into training and testing sets
train_data = data.sample(frac=0.8, random_state=1)
test_data = data.drop(train_data.index)
# Create a linear regression model
model = LinearRegression()
# Train the model using the training data
model.fit(train_data[['x']], train_data['y'])
# Evaluate the model using the testing data
r_squared = model.score(test_data[['x']],
test_data['y'])
```

These examples demonstrate some of the key features of Python that may be useful for SAS users, including reading and writing data, data manipulation and analysis, mathematical operations, and machine learning.



Here is a more complex example of Python code that utilizes some advanced features that may be useful for SAS users:

```
import pandas as pd
from sklearn.preprocessing import StandardScaler
from sklearn.decomposition import PCA
from sklearn.cluster import KMeans
# Load data from a CSV file
data = pd.read csv('data.csv')
# Standardize the data
scaler = StandardScaler()
scaled data = scaler.fit transform(data)
# Perform principal component analysis
pca = PCA(n components=2)
pca data = pca.fit transform(scaled data)
# Perform k-means clustering
kmeans = KMeans(n clusters=3)
kmeans.fit(pca data)
cluster labels = kmeans.predict(pca data)
# Create a new column in the original data with the
cluster labels
data['cluster'] = cluster labels
# Write the updated data to a new CSV file
data.to csv('output.csv')
```

This code performs several advanced data analysis tasks:

- 1. Standardizing the data: The **StandardScaler** class from the **sklearn.preprocessing** library is used to standardize the data by subtracting the mean and dividing by the standard deviation. This is a common preprocessing step in data analysis.
- 2. Performing principal component analysis: The **PCA** class from the **sklearn.decomposition** library is used to perform principal component analysis, which is a technique for reducing the dimensionality of the data. In this case, the data is reduced to two dimensions for visualization purposes.
- 3. Performing k-means clustering: The **KMeans** class from the **sklearn.cluster** library is used to perform k-means clustering on the principal component data. This is a technique for grouping similar data points together based on their distance in the data space.



- 4. Updating the original data with the cluster labels: The cluster labels are added as a new column to the original data using the **data['cluster'] = cluster_labels** syntax.
- 5. Writing the updated data to a new CSV file: The **to_csv** method is used to write the updated data to a new CSV file.

This example demonstrates some of the more advanced features of Python that may be useful for SAS users, such as data standardization, dimensionality reduction, clustering, and data visualization

Variables and data types

A variable is created when a value is assigned to it using the = operator. Python supports several data types, including integers, floating-point numbers, strings, Boolean values, and complex numbers.

Integers are whole numbers without a decimal point. They can be positive, negative, or zero. For example, x = 5 assigns the value 5 to the variable x.

Floating-point numbers are decimal numbers, such as 3.14 or 0.001. They can also be positive, negative, or zero. For example, y = 3.14 assigns the value 3.14 to the variable y.

Strings are sequences of characters, enclosed in quotation marks. They can be single quotes ('...') or double quotes ("..."). For example, name = "John" assigns the string "John" to the variable name.

Boolean values represent either True or False. They are often used in conditional statements and loops. For example, flag = True assigns the value True to the variable flag.

Complex numbers are written in the form a + bj, where a and b are real numbers and j is the imaginary unit. For example, z = 3 + 4j assigns the complex number 3 + 4j to the variable z.

Python is a dynamically typed language, which means that you do not need to specify the data type of a variable when you create it. Python will automatically assign the appropriate data type based on the value you assign to the variable.

Here is an example of Python code that demonstrates the use of variables and data types:

```
# Assign integer value to a variable
x = 5
print("x =", x) # Output: x = 5
# Assign floating-point value to a variable
y = 3.14
print("y =", y) # Output: y = 3.14
```



```
# Assign string value to a variable
name = "John"
print("name =", name) # Output: name = John
# Assign Boolean value to a variable
flag = True
print("flag =", flag) # Output: flag = True
# Assign complex number to a variable
z = 3 + 4j
print("z =", z) # Output: z = (3+4j)
```

In addition to the basic data types mentioned above, Python also supports several built-in data structures, such as lists, tuples, sets, and dictionaries. These data structures allow you to store multiple values in a single variable.

For example, a list is a collection of items that are ordered and changeable. You can create a list by enclosing a comma-separated sequence of values in square brackets. Here is an example of Python code that demonstrates the use of a list:

```
# Create a list of integers
numbers = [1, 2, 3, 4, 5]
print("numbers =", numbers) # Output: numbers = [1, 2,
3, 4, 5]
# Access individual elements of a list
print("First element of numbers:", numbers[0]) #
Output: First element of numbers: 1
print("Last element of numbers: ", numbers[-1]) #
Output: Last element of numbers: 5
# Change an element of a list
numbers[2] = 10
print("Updated numbers list:", numbers) # Output:
Updated numbers list: [1, 2, 10, 4, 5]
# Add an element to the end of a list
numbers.append
```

Variables and data types are essential concepts in programming languages, and Python is no exception. Understanding how to work with variables and data types is crucial for writing effective and efficient code.



In Python, a variable is a name that refers to a value stored in the computer's memory. The value assigned to a variable can be a number, a string of characters, a Boolean value, or any other data type. Variables can be assigned new values at any time, and Python will automatically update the variable's value in the computer's memory.

Here are some guidelines for naming variables in Python:

Variable names must start with a letter or an underscore character (_).

Variable names cannot start with a number.

Variable names can contain letters, numbers, and underscore characters.

Variable names are case-sensitive.

In Python, there are several data types that can be assigned to a variable. The most common data types in Python are:

Integers: Integers are whole numbers, such as 1, 2, 3, 4, 5, etc.

Floating-point numbers: Floating-point numbers are numbers with decimal points, such as 1.0, 2.5, 3.14, etc.

Strings: Strings are sequences of characters, such as "hello", "world", "Python", etc.

Booleans: Booleans are either True or False.

None: None is a special value that represents the absence of a value.

Python is a dynamically-typed language, which means that the data type of a variable is determined at runtime. This is different from statically-typed languages like C++ and Java, where you must declare the data type of a variable before you can use it.

Here are some examples of assigning values to variables in Python:

```
# Assign an integer to a variable
x = 10
# Assign a floating-point number to a variable
y = 3.14
# Assign a string to a variable
name = "John Smith"
# Assign a Boolean value to a variable
is_valid = True
# Assign None to a variable
result = None
```

Python also provides several built-in data structures that can be used to store multiple values in a single variable. Some of the most common data structures in Python are:

Lists: Lists are ordered collections of items that can be of any data type.



Tuples: Tuples are ordered collections of items that are immutable (i.e., cannot be changed). Sets: Sets are unordered collections of unique items.

Dictionaries: Dictionaries are collections of key-value pairs.

Here are some examples of using built-in data structures in Python:

```
# Create a list of integers
numbers = [1, 2, 3, 4, 5]
# Create a tuple of strings
fruits = ("apple", "banana", "cherry")
# Create a set of floating-point numbers
prices = {1.99, 2.99, 3.99}
# Create a dictionary of key-value pairs
person = {"name": "John Smith", "age": 30, "email":
"john@example.com"}
```

Variables and data types are fundamental concepts in Python programming. Understanding how to work with variables and data types is essential for writing effective and efficient code. Python provides a variety of built-in data types and data structures that can be used to store and manipulate data, making it a powerful and flexible programming language.

Here are some longer code examples that illustrate the use of variables and data types in Python:

Example 1: Simple Calculator

```
# Ask the user to input two numbers
num1 = float(input("Enter the first number: "))
num2 = float(input("Enter the second number: "))
# Perform arithmetic operations on the numbers
sum = num1 + num2
difference = num1 - num2
product = num1 * num2
quotient = num1 / num2
# Print the results
print("Sum: ", sum)
print("Difference: ", difference)
print("Product: ", product)
print("Quotient: ", quotient)
```

In this example, the user is prompted to enter two numbers, which are then stored in the num1 and num2 variables as floating-point numbers. The program then performs basic arithmetic



operations on the numbers and stores the results in sum, difference, product, and quotient variables. Finally, the results are printed to the console.

Example 2: Data Analysis

```
# Create a list of exam scores
scores = [85, 90, 78, 92, 88, 75, 80, 95]
# Calculate the average score
total = 0
for score in scores:
    total += score
average = total / len(scores)
# Determine the highest and lowest scores
highest = max(scores)
lowest = min(scores)
# Print the results
print("Average score: ", average)
print("Highest score: ", highest)
print("Lowest score: ", lowest)
```

In this example, we have a list of exam scores stored in the scores variable. The program calculates the average score by iterating over the list and summing up the scores. The average is then computed by dividing the total by the number of scores. The program also determines the highest and lowest scores using the max() and min() functions. Finally, the results are printed to the console.

Example 3: Tic-Tac-Toe Game

```
# Create a 3x3 game board
board = [
    ['_', '_', '_'],
    ['_', '_', '_'],
    ['_', '_', '_']
]
# Print the game board
def print_board():
    for row in board:
        print(row)
# Ask the players to make their moves
instal
```

```
def make move(player):
    while True:
        row = int(input("Enter the row number (0-2):
"))
        col = int(input("Enter the column number (0-2):
"))
        if board[row][col] == ' ':
            board[row][col] = player
            break
        else:
            print("That space is already taken!")
# Check if the game is over
def check game over(player):
    # Check rows
    for row in board:
        if row.count(player) == 3:
            return True
    # Check columns
    for col in range(3):
        if board[0][col] == player and board[1][col] ==
player and board[2][col] == player:
            return True
    # Check diagonals
    if board[0][0] == player and board[1][1] == player
and board[2][2] == player:
        return True
    if board[0][2] == player and board[1][1] == player
and board[2][0] == player:
        return True
    return False
# Play the game
print("Welcome to Tic-Tac-Toe!")
print board()
while True:
    make move('X')
    print board()
```

In Python, variables are used to store values such as numbers, strings, and objects. Unlike other programming languages, you don't need to declare the type of a variable in Python. You can assign any value to a variable, and Python will automatically determine the type based on the value.



For example, you can assign an integer value to a variable like this:

x = 10

You can also assign a floating-point number:

y = 3.14

And you can assign a string:

z = "Hello, World!"

Data Types

Python supports several data types, including integers, floating-point numbers, strings, and boolean values. Some of the commonly used data types in Python are:

Integers: Integers are whole numbers with no decimal point, such as 1, 2, 3, etc. In Python, integers are represented using the int data type.

Floating-point numbers: Floating-point numbers are decimal numbers, such as 3.14, 2.5, etc. In Python, floating-point numbers are represented using the float data type.

Strings: Strings are a sequence of characters, such as "Hello, World!", "Python", etc. In Python, strings are represented using the str data type.

Booleans: Booleans are either True or False. They are often used for logical comparisons and control flow in Python.

In addition to these basic data types, Python also supports more complex data structures such as lists, tuples, dictionaries, and sets.

Conditional statements

Conditional statements are an essential part of any programming language, including Python. They allow the program to execute different sets of instructions based on specific conditions. In this article, we will explore conditional statements in Python, with a focus on how they differ from similar statements in SAS.

IF/THEN Statements in SAS

In SAS, the IF/THEN statement is used to test a condition and execute a specific set of statements if the condition is true. The basic syntax for an IF/THEN statement in SAS is as follows:



```
if condition then do;
   statement1;
   statement2;
   ...
end;
```

Here, condition is the expression that needs to be evaluated, and statement1, statement2, etc. are the statements that will be executed if the condition is true.

For example, let's say we want to check if a variable age is greater than 18. If it is, we print a message saying the person is an adult. Otherwise, we print a message saying the person is a minor. Here's how we would do this in SAS:

```
data mydata;
  set mydata;
  if age > 18 then do;
    put 'This person is an adult.';
  end;
  else do;
    put 'This person is a minor.';
  end;
run;
```

Conditional Statements in Python

Python uses a similar syntax to SAS for conditional statements, with some slight differences. The basic syntax for an IF/THEN statement in Python is as follows:

```
if condition:
statement1
statement2
```

Here, condition is the expression that needs to be evaluated, and statement1, statement2, etc. are the statements that will be executed if the condition is true. Notice that there is no need for an "end" statement in Python, and the statements to be executed if the condition is true are not enclosed in a "do" loop.

For example, let's say we want to check if a variable age is greater than 18. If it is, we print a message saying the person is an adult. Otherwise, we print a message saying the person is a minor. Here's how we would do this in Python:

```
if age > 18:
    print('This person is an adult.')
else:
    print('This person is a minor.')
```



Notice that we use the "print" statement instead of the "put" statement in SAS. Also, the conditional statements in Python use colons (:) instead of semicolons (;) in SAS.

IF/THEN/ELSE Statements in SAS

In SAS, we can use the IF/THEN/ELSE statement to execute different sets of statements based on whether the condition is true or false. The basic syntax for an IF/THEN/ELSE statement in SAS is as follows:

```
if condition then do;
   statement1;
   statement2;
   ...
end;
else do;
   statement3;
   statement4;
   ...
end;
```

Here, condition is the expression that needs to be evaluated, and statement1, statement2, etc. are the statements that will be executed if the condition is true. statement3, statement4, etc. are the statements that will be executed if the condition is false.

Conditional statements are used in programming to execute different pieces of code depending on whether a certain condition is true or false. These statements are important in many programming languages, including SAS and Python.

In SAS, conditional statements are used extensively in data step programming to conditionally process data. For example, you may want to only process observations that meet a certain criteria, or you may want to perform different calculations for different groups of data. SAS provides several conditional statements that allow you to do this, including the IF/THEN/ELSE and IF/THEN/ELSEIF statements.

In Python, conditional statements are also used extensively in programming. Python provides several conditional statements, including the IF/THEN/ELSE and IF/THEN/ELIF statements that allow you to execute different pieces of code depending on whether a certain condition is true or false.

One important difference between SAS and Python conditional statements is the use of indentation in Python. In SAS, you enclose the code to be executed for each condition in a DO/END block, while in Python you use indentation to define the block of code to be executed. The use of indentation can be challenging for SAS users who are new to Python, but it is an important aspect of Python programming and allows for more readable code.



Here are some additional examples of conditional statements in SAS and Python:

Example 1 - SAS:

```
data mydata;
   set mydata;
   if age >= 18 then do;
      status = 'adult';
   end;
   else do;
      status = 'minor';
   end;
run;
```

In this example, we use the IF/THEN/ELSE statement to check if the age variable is greater than or equal to 18. If it is, we assign the value 'adult' to the status variable. Otherwise, we assign the value 'minor' to the status variable.

Example 1 - Python:

```
if age >= 18:
    status = 'adult'
else:
    status = 'minor'
```

In this example, we use the IF/THEN/ELSE statement in Python to check if the age variable is greater than or equal to 18. If it is, we assign the value 'adult' to the status variable. Otherwise, we assign the value 'minor' to the status variable.

```
Example 2 - SAS:
```

```
data mydata;
   set mydata;
   if age < 18 then do;
      status = 'minor';
   end;
   else if age >= 18 and age < 65 then do;
      status = 'adult';
   end;
   else do;
      status = 'senior';
   end;
run;
```



In this example, we use the IF/THEN/ELSEIF statement to check if the age variable falls into one of three categories: under 18, between 18 and 65, or over 65. We assign the value 'minor' to the status variable if the age is under 18, 'adult' if the age is between 18 and 65, and 'senior' if the age is over 65.

Example 2 - Python:

```
if age < 18:
    status = 'minor'
elif age >= 18 and age < 65:
    status = 'adult'
else:
    status = 'senior'
```

In this example, we use the IF/THEN/ELIF statement in Python to check if the age variable falls into one of three categories: under 18, between 18 and 65, or over 65. We assign the value 'minor' to the status variable if the age is under 18, 'adult' if the age is between 18 and 65, and 'senior' if the age is over 65.

Conditional statements are an essential part of programming in both SAS and Python. They allow you to control the flow of your code and make it more flexible and powerful. Understanding how to use conditional statements is important for any programmer, whether you are working with SAS, Python, or another programming language.

Here's an example of conditional statement code in Python:

```
# Determine if a number is positive, negative, or zero
x = 5
if x > 0:
    print("Positive")
elif x < 0:
    print("Negative")
else:
    print("Zero")</pre>
```

In this example, we define a variable x with a value of 5. We then use an if statement to check if x is greater than 0. If it is, we print the string "Positive". If x is not greater than 0, we use an elif statement to check if x is less than 0. If it is, we print the string "Negative". Finally, if x is neither greater than 0 nor less than 0, we use an else statement to print the string "Zero".

This is a simple example, but it illustrates how conditional statements can be used to make decisions in your code based on certain conditions. With more complex conditions, you can use logical operators (such as and, or, and not) to combine multiple conditions and create more complex decision trees.



Example 1:

```
# Determine if a number is even or odd
x = 7
if x % 2 == 0:
    print("Even")
else:
    print("Odd")
```

In this example, we define a variable x with a value of 7. We then use an if statement with a condition that checks if x is divisible by 2 (i.e., if the remainder of x divided by 2 is equal to 0). If the condition is true, we print the string "Even". Otherwise, we print the string "Odd". This example shows how you can use conditional statements to determine if a number has a certain property (in this case, being even or odd).

Example 2:

```
# Determine the maximum of three numbers
a = 5
b = 9
c = 3
if a > b:
    if a > c:
        print("Max:", a)
    else:
        print("Max:", c)
else:
        if b > c:
            print("Max:", b)
    else:
            print("Max:", c)
```

In this example, we define three variables a, b, and c with values of 5, 9, and 3, respectively. We then use nested if statements to determine which of the three numbers is the largest. The outer if statement checks if a is greater than b. If it is, we use an inner if statement to check if a is greater than c. If it is, we print the value of a as the maximum. If a is not greater than c, we print the value of c as the maximum. If a is not greater than b, we use another nested if statement to check if b is greater than c. If it is, we print the value of b as the maximum. If b is not greater than c, we print the value of c as the maximum. This example shows how you can use nested conditional statements to create more complex decision trees.

Example 3:



```
# Determine if a word is a palindrome
word = "racecar"
if word == word[::-1]:
    print("Palindrome")
else:
    print("Not a palindrome")
```

In this example, we define a variable word with a value of "racecar". We then use an if statement with a condition that checks if word is equal to its reverse (word[::-1]). If the condition is true, we print the string "Palindrome". Otherwise, we print the string "Not a palindrome". This example shows how you can use conditional statements to check if a string has a certain property (in this case, being a palindrome).

Example 4:

```
# Determine the season based on the month
month = "April"
if month in ["March", "April", "May"]:
    print("Spring")
elif month in ["June", "July", "August"]:
    print("Summer")
elif month in ["September", "October", "November"]:
    print("Fall")
else:
    print("Winter")
```

In this example, we define a variable month with a value of "April". We then use a series of if and elif statements with conditions that check if month is in a certain range of values. If month is in the range of March to May, we print the string "Spring". If it is in the range of June to August, we print the string "Summer". If it is in the range of September to November, we print the string "Fall". Otherwise, we print the string "Winter". This example shows how you can use conditional statements to make decisions based on a range of values.

Example 5:

```
# Determine if a year is a leap year
year = 2020
if year % 4 == 0:
    if year % 100 == 0:
        if year % 400 == 0:
            print("Leap year")
```



In this example, we define a variable year with a value of 2020. We then use nested if statements to determine if year is a leap year. The outer if statement checks if year is divisible by 4. If it is, we use an inner if statement to check if year is divisible by 100. If it is, we use another inner if statement to check if year is divisible by 400. If it is, we print the string "Leap year". If year is not divisible by 400, we print the string "Not a leap year". If year is not divisible by 100, we print the string "Leap year". If year is not divisible by 4, we print the string "Not a leap year". This example shows how you can use nested conditional statements to create more complex decision trees.

Example 6:

```
# Determine if a number is positive, negative, or zero
    number = 7

    if number > 0:
        print("Positive")
    elif number < 0:
        print("Negative")
    else:
        print("Zero")</pre>
```

In this example, we define a variable number with a value of 7. We then use an if statement with a condition that checks if number is greater than 0. If it is, we print the string "Positive". If number is not greater than 0, we use an elif statement with a condition that checks if number is less than 0. If it is, we print the string "Negative". If number is not less than 0, we use an else statement to print the string "Zero". This example shows how you can use conditional statements to check the value of a variable and make decisions based on its sign.

Example 7:

```
# Determine the absolute value of a number
number = -5
if number < 0:
    absolute_value = -number
else:
    absolute value = number
```



print("The absolute value of", number, "is", absolute_value)

In this example, we define a variable number with a value of -5. We then use an if statement with a condition that checks if number is less than 0. If it is, we assign the negative of number to a new variable absolute_value. If number is not less than 0, we assign number to absolute_value. We then print a message that displays the original value of number and its absolute value. This example shows how you can use conditional statements to perform calculations based on the value of a variable.

Example 8:

```
# Determine if a word is a palindrome
word = "racecar"
if word == word[::-1]:
    print("Palindrome")
else:
    print("Not a palindrome")
```

In this example, we define a variable word with a value of "racecar". We then use an if statement with a condition that checks if word is equal to the reverse of word (using slicing with a step of -1 to reverse the string). If it is, we print the string "Palindrome". If word is not equal to its reverse, we print the string "Not a palindrome". This example shows how you can use conditional statements to check if a string is a palindrome.

Example 9:

```
# Determine the largest of three numbers
a = 5
b = 7
c = 3
if a > b and a > c:
    print(a, "is the largest number")
elif b > a and b > c:
    print(b, "is the largest number")
else:
    print(c, "is the largest number")
```

In this example, we define three variables a, b, and c with values of 5, 7, and 3, respectively. We then use a series of if and elif statements with conditions that check which of the three variables is the largest. If a is greater than both b and c, we print a message that says a is the largest number. If b is greater than both a and c, we print a message that says b is the largest number. If neither of those conditions is true, we print a message that says c is the largest number. This



example shows how you can use conditional statements to determine which of several variables has the largest value.

Example 10:

```
# Determine if a number is prime
number = 17
if number > 1:
    for i in range(2, number):
        if number % i == 0:
            print(number, "is not a prime number")
            break
    else:
        print(number, "is a prime number")
else:
    print(number, "is not a prime number")
```

In this example, we define a variable number with a value of 17. We then use an if statement with a condition that checks if number is greater than 1. If it is, we use a for loop to iterate over a range of numbers from 2 to number - 1. For each number in the range, we use an if statement with a condition that checks if number is divisible by that number. If it is, we print a message that says number is not a prime number and break out of the loop. If the loop completes without finding a factor of number, we print a message that says number is a prime number. If number is less than or equal to 1, we print a message that says number is not a prime number. This example shows how you can use conditional statements and loops to check if a number is prime.

Example 11:

```
# Check if a string is a palindrome
string = "racecar"
if string == string[::-1]:
    print(string, "is a palindrome")
else:
    print(string, "is not a palindrome")
```

In this example, we define a variable string with a value of "racecar". We then use an if statement with a condition that checks if string is equal to a reversed version of itself (string[::-1]). If it is, we print a message that says string is a palindrome. If it's not, we print a message that says string is not a palindrome. This example shows how you can use conditional statements to check if a string is a palindrome.

Example 12:

in stal

```
# Determine the grade based on a score
score = 85
if score >= 90:
    grade = "A"
elif score >= 80:
    grade = "B"
elif score >= 70:
    grade = "C"
elif score >= 60:
    grade = "D"
else:
    grade = "F"
print("The grade for a score of", score, "is", grade)
```

In this example, we define a variable score with a value of 85. We then use a series of if and elif statements with conditions that check what range score falls into. If score is greater than or equal to 90, we set the variable grade to "A". If score is greater than or equal to 80, we set grade to "B". If score is greater than or equal to 70, we set grade to "C". If score is greater than or equal to 60, we set grade to "D". If none of those conditions are true, we set grade to "F". Finally, we print a message that includes the original score and the corresponding grade. This example shows how you can use conditional statements to assign a grade based on a score.

Loops

One of the key features of Python is the ability to use loops, which allow for the repeated execution of code. In this article, we will discuss loops in Python and provide examples that are relevant to SAS users.

Types of Loops:

Python has two types of loops: for loops and while loops.

For Loops:

A for loop is used to iterate over a sequence (e.g., a list, tuple, string, or dictionary) and execute a block of code for each element in the sequence. The syntax for a for loop is as follows:

```
for variable in sequence:
    # code block to be executed
```



In the above syntax, "variable" is the name of the variable that will hold each element of the sequence, and "sequence" is the sequence to be iterated over.

Example:

Suppose we have a list of numbers and we want to print each number on a separate line. We can use a for loop to accomplish this:

```
numbers = [1, 2, 3, 4, 5]
for num in numbers:
    print(num)
Output:
```

While Loops:

A while loop is used to repeatedly execute a block of code as long as a certain condition is true. The syntax for a while loop is as follows:

```
while condition:
    # code block to be executed
```

In the above syntax, "condition" is the expression that is evaluated to determine whether to continue executing the code block.

Example:

Suppose we have a list of numbers and we want to sum them up until the total is greater than 10. We can use a while loop to accomplish this:

```
numbers = [1, 2, 3, 4, 5]
total = 0
i = 0
while total <= 10:
    total += numbers[i]
    i += 1</pre>
```



print(total)

Output:

15

In the above example, we initialize "total" to 0 and "i" to 0. We then enter the while loop, where we repeatedly add the next number in the list to "total" until "total" is greater than 10.

Break and Continue Statements:

Python also provides two statements that can be used inside loops to control the flow of execution: "break" and "continue".

The "break" statement is used to exit a loop prematurely. When a "break" statement is encountered, the loop immediately terminates and control is passed to the statement immediately following the loop.

Example:

Suppose we have a list of numbers and we want to print each number until we encounter a negative number. We can use a for loop and a break statement to accomplish this:

```
numbers = [1, 2, -3, 4, 5]
for num in numbers:
    if num < 0:
        break
    print(num)</pre>
```

Output:

1 2

In the above example, the for loop iterates over each element in the list "numbers". When a negative number is encountered, the if statement is true and the break statement is executed, causing the loop to terminate.

The "continue" statement is used to skip the current iteration of a loop and continue with the next iteration.

Example:

Suppose we have a list of numbers and we want to print each number that is not equal to 3. We can use a for loop and a continue statement to accomplish this:



numbers = [1, 2, 3, 4, 5]
for num in numbers:
 if num == 3:
 continue
 print(num)

Output:

In the above example, the for loop iterates over each element in the list "numbers". When the number 3 is encountered, the if statement is true and the continue statement is executed, causing the loop to skip the iteration where num is equal to 3.

Loops are a fundamental concept in Python that allow for the repeated execution of code. SAS users can benefit from learning how to use loops in Python to expand their data analysis capabilities. Python has two types of loops (for loops and while loops) and also provides two statements (break and continue) that can be used to control the flow of execution inside loops. By mastering these concepts, SAS users can become more efficient and effective in their data analysis tasks.

Iterating Over Data Structures:

One of the most common use cases for loops in data analysis is iterating over data structures like lists, tuples, and dictionaries. By using a loop to iterate over each element in the data structure, we can perform calculations or operations on each element in a structured and efficient way.

Example:

Suppose we have a list of temperatures in Celsius and we want to convert each temperature to Fahrenheit. We can use a for loop to iterate over the list and perform the conversion for each temperature:

```
celsius_temps = [25, 30, 20, 15, 22]
fahrenheit_temps = []
for temp in celsius_temps:
    fahrenheit_temps.append(temp * 1.8 + 32)
print(fahrenheit_temps)
```



Output:

[77.0, 86.0, 68.0, 59.0, 71.6]

In the above example, we initialize an empty list "fahrenheit_temps" and then use a for loop to iterate over each temperature in "celsius_temps". For each temperature, we perform the conversion to Fahrenheit and append the result to the "fahrenheit_temps" list.

Nested Loops:

Another use case for loops in data analysis is nested loops. Nested loops are loops that are contained within another loop, and they are used to perform calculations or operations on multiple data structures.

Example:

Suppose we have two lists of numbers and we want to calculate the product of each combination of numbers in the two lists. We can use nested for loops to accomplish this:

```
list1 = [1, 2, 3]
list2 = [4, 5, 6]
products = []
for num1 in list1:
    for num2 in list2:
        products.append(num1 * num2)
print(products)
```

Output:

[4, 5, 6, 8, 10, 12, 12, 15, 18]

In the above example, we use nested for loops to iterate over each element in both lists and calculate the product of each combination of numbers. The results are stored in the "products" list.

Looping with Indexes:

Sometimes we need to loop over a list or other data structure while also keeping track of the index of the current element. We can use the built-in "enumerate" function in Python to accomplish this.



Example:

Suppose we have a list of names and we want to print each name along with its index in the list. We can use the "enumerate" function to accomplish this:

```
names = ['Alice', 'Bob', 'Charlie', 'David']
for i, name in enumerate(names):
    print(f'{i}: {name}')
```

Output:

0: Alice 1: Bob 2: Charlie 3: David

In the above example, we use the "enumerate" function to loop over the "names" list and get both the index and the name for each element in the list. We then print out the index and name in a formatted string.

Code examples that demonstrate the use of loops in data analysis:

Example 1: Calculating Descriptive Statistics for a List of Numbers In this example, we will use a for loop to calculate the mean, variance, and standard deviation of a list of numbers.

```
import math
numbers = [1, 2, 3, 4, 5]
# Calculate the mean
mean = sum(numbers) / len(numbers)
# Calculate the variance
variance = 0
for num in numbers:
    variance += (num - mean) ** 2
variance /= len(numbers)
# Calculate the standard deviation
std_dev = math.sqrt(variance)
print(f"Mean: {mean}")
```



```
print(f"Variance: {variance}")
print(f"Standard Deviation: {std_dev}")
```

Output:

```
Mean: 3.0
Variance: 2.5
Standard Deviation: 1.5811388300841898
```

In the above example, we first calculate the mean of the list of numbers using the built-in "sum" and "len" functions. We then use a for loop to calculate the variance by summing the squared differences between each number and the mean, and dividing by the length of the list. Finally, we calculate the standard deviation using the "sqrt" function from the "math" module.

Example 2: Filtering Data in a List of Dictionaries In this example, we will use a for loop to filter a list of dictionaries based on a certain criteria.

```
data = [
    {'name': 'Alice', 'age': 25, 'gender': 'female'},
    {'name': 'Bob', 'age': 30, 'gender': 'male'},
    {'name': 'Charlie', 'age': 20, 'gender': 'male'},
    {'name': 'David', 'age': 35, 'gender': 'male'},
    {'name': 'Eve', 'age': 28, 'gender': 'female'}
]
# Filter the data to only include males over 25
filtered_data = []
for row in data:
    if row['gender'] == 'male' and row['age'] > 25:
        filtered_data.append(row)
print(filtered_data)
```

Output:

[{'name': 'Bob', 'age': 30, 'gender': 'male'}, {'name': 'David', 'age': 35, 'gender': 'male'}]

In the above example, we have a list of dictionaries representing some data. We use a for loop to iterate over each dictionary in the list, and we use an if statement to check if the row meets our filter criteria (in this case, being a male over the age of 25). If the row meets the criteria, we append it to a new list called "filtered_data". Finally, we print out the filtered data.



Functions

Functions are an important aspect of programming in Python as they allow you to encapsulate a block of code that can be reused and called upon whenever necessary. In this section, we will explore some commonly used functions in Python from a SAS user's perspective.

Defining Functions:

Functions in Python are defined using the def keyword followed by the function name, a set of parentheses, and a colon. The function body is indented and can contain any valid Python code.

```
def function_name(parameters):
    """
    Docstring: A description of the function.
    """
    # Function body
    return value
```

The function name should follow the same naming conventions as variable names, i.e., use lowercase letters and underscores for readability. The parameters are optional and can be used to pass values into the function. The return statement is optional and can be used to return a value from the function.

Example Function:

```
def square(x):
    """
    This function returns the square of a given number.
    """
    return x**2
```

Calling Functions:

To call a function in Python, simply use the function name followed by the parentheses and any necessary arguments.

```
result = function name(arguments)
```

Example Function Call:

```
result = square(5)
print(result) # Output: 25
Built-in Functions:
```



Python provides a number of built-in functions that can be used without the need for importing any additional libraries. Here are some commonly used built-in functions:

print(): Prints the specified message to the console.len(): Returns the length of the specified object.type(): Returns the type of the specified object.range(): Returns a sequence of numbers, starting from 0 and incrementing by 1 (default), and stopping before a specified number.Example Built-in Function Call:

```
# Prints "Hello, World!" to the console
print("Hello, World!")
# Returns the length of the string "Python"
length = len("Python")
print(length) # Output: 6
# Returns the type of the integer 42
type_of_42 = type(42)
print(type_of_42) # Output: <class 'int'>
# Returns a range of numbers from 0 to 4 (excluding 4)
numbers = range(4)
print(list(numbers)) # Output: [0, 1, 2, 3]
```

Custom Functions:

In addition to the built-in functions, you can also create your own custom functions to perform specific tasks. Here is an example of a custom function:

```
def calculate_mean(numbers):
    """
    This function calculates the mean of a list of
numbers.
    """
    total = sum(numbers)
    count = len(numbers)
    mean = total / count
    return mean
```

Example Custom Function Call:

```
numbers = [1, 2, 3, 4, 5]
mean = calculate mean(numbers)
```



print(mean) # Output: 3.0

Functions are an important aspect of programming in Python and can greatly simplify your code by allowing you to reuse blocks of code. Whether you are using built-in functions or creating your own custom functions, mastering this concept is essential to becoming a proficient Python programmer.

Parameters and Arguments:

Parameters are variables that are defined in the function definition and are used to pass values into the function. Arguments, on the other hand, are the values that are passed into the function when it is called.

In Python, there are two types of parameters: positional parameters and keyword parameters. Positional parameters are defined by their position in the function definition and must be passed in the same order when the function is called. Keyword parameters, on the other hand, are defined with a default value and can be passed in any order using their name.

Example of Positional Parameters:

```
def greet(name, message):
    """
    This function greets a person with a message.
    """
    print(f"{message}, {name}!")
greet("Alice", "Hello") # Output: Hello, Alice!
```

Example of Keyword Parameters:

```
def greet(name, message="Hello"):
    """
    This function greets a person with a message.
    """
    print(f"{message}, {name}!")

greet("Bob") # Output: Hello, Bob!
greet(message="Hi", name="Charlie") # Output: Hi,
Charlie!
```



Default Parameters:

As seen in the example above, default parameters can be used to define a default value for a parameter in case it is not passed in when the function is called. If a default parameter is defined, it must come after any non-default parameters in the function definition.

Example of Default Parameters:

```
def greet(name, message="Hello", punctuation="!"):
    """
    This function greets a person with a message and
punctuation.
    """
    print(f"{message}, {name}{punctuation}")

greet("Dave") # Output: Hello, Dave!
greet("Eve", "Hi") # Output: Hi, Eve!
greet("Frank", punctuation=".") # Output: Hello, Frank.
```

Arbitrary Parameters:

In some cases, you may want to define a function that can accept a variable number of arguments. This can be done by using an arbitrary parameter, denoted by an asterisk (*). This parameter can be used to accept any number of positional arguments, which are passed in as a tuple.

Example of Arbitrary Parameters:

```
def sum_all(*numbers):
    """
    This function sums up all the numbers passed in.
    """
    total = sum(numbers)
    return total

result = sum_all(1, 2, 3, 4, 5)
print(result) # Output: 15
```

Keyword Arbitrary Parameters:

Similarly, you can also define a function that accepts a variable number of keyword arguments by using an arbitrary keyword parameter, denoted by two asterisks (**). This parameter can be used to accept any number of keyword arguments, which are passed in as a dictionary. Example of Keyword Arbitrary Parameters:



```
def print_info(**info):
    """
    This function prints out a dictionary of
information.
    """
    for key, value in info.items():
        print(f"{key}: {value}")

print_info(name="Gina", age=30, occupation="Engineer")
# Output:
# name: Gina
# age: 30
# occupation: Engineer
```

Functions are a fundamental concept in Python and are used extensively in programming. They allow you to encapsulate a block of code that can be reused and called upon whenever necessary. By understanding how to define and use parameters, arguments, and arbitrary parameters, you can create powerful and flexible functions that can be used in a variety of applications.

Here are some more examples of functions in Python:

Example 1: Finding the largest number in a list

```
def find_largest(numbers):
    """
    This function finds the largest number in a list of
numbers.
    """
    largest = numbers[0]
    for num in numbers:
        if num > largest:
            largest = num
    return largest
list_of_numbers = [10, 5, 20, 15, 30]
largest_number = find_largest(list_of_numbers)
print(largest_number) # Output: 30
```

Example 2: Reversing a string

```
def reverse_string(text):
    """
```



```
This function reverses a string.
"""
reversed_text = ""
for i in range(len(text) - 1, -1, -1):
    reversed_text += text[i]
return reversed_text
text_to_reverse = "hello world"
reversed_text = reverse_string(text_to_reverse)
print(reversed_text) # Output: dlrow olleh
```

Example 3: Calculating the factorial of a number

```
def factorial(number):
    """
    This function calculates the factorial of a number.
    """
    result = 1
    for i in range(1, number + 1):
        result *= i
    return result
num = 5
factorial_of_num = factorial(num)
print(factorial_of_num) # Output: 120
```

Example 4: Checking if a string is a palindrome

```
def is_palindrome(text):
    """
    This function checks if a string is a palindrome.
    """
    reversed_text = text[::-1]
    if text == reversed_text:
        return True
    else:
        return False

text_to_check = "racecar"
is_palindrome_text = is_palindrome(text_to_check)
print(is_palindrome_text) # Output: True
```



These are just a few examples of the many functions that can be created in Python. By understanding the basic concepts of functions, you can create your own custom functions to solve a wide variety of problems.

Example 5: Converting Celsius to Fahrenheit

```
def celsius_to_fahrenheit(celsius):
    """
    This function converts a temperature in Celsius to
Fahrenheit.
    """
    fahrenheit = (celsius * 1.8) + 32
    return fahrenheit

celsius_temp = 25
fahrenheit_temp = celsius_to_fahrenheit(celsius_temp)
print(fahrenheit_temp) # Output: 77.0
```

Example 6: Finding the average of a list of numbers

```
def find_average(numbers):
    """
    This function finds the average of a list of
numbers.
    """
    total = sum(numbers)
    avg = total / len(numbers)
    return avg
list_of_numbers = [10, 20, 30, 40, 50]
average_of_numbers = find_average(list_of_numbers)
print(average_of_numbers) # Output: 30.0
```

Example 7: Converting a list of strings to lowercase

```
def convert_to_lowercase(strings):
    """
    This function converts a list of strings to
lowercase.
    """
    lowercase_strings = []
    for string in strings:
        lowercase_strings.append(string.lower())
    return lowercase_strings
```



```
list_of_strings = ["HELLO", "WORLD", "PYTHON"]
lowercase_strings =
convert_to_lowercase(list_of_strings)
print(lowercase_strings) # Output: ['hello', 'world',
'python']
```

These are just a few more examples of the many functions that can be created in Python. With functions, you can encapsulate complex logic into a single unit that can be easily reused and maintained.

Example 8: Calculating the area and circumference of a circle

```
import math

def circle_calculations(radius):
    """
    This function calculates the area and circumference
of a circle with the given radius.
    """
    area = math.pi * radius ** 2
    circumference = 2 * math.pi * radius
    return area, circumference

r = 5
area, circumference = circle_calculations(r)
print(f"Area: {area:.2f}, Circumference:
    {circumference:.2f}") # Output: Area: 78.54,
Circumference: 31.42
```

Example 9: Checking if a number is prime

```
def is_prime(number):
    """
    This function checks if a number is prime.
    """
    if number < 2:
        return False
    for i in range(2, int(number ** 0.5) + 1):
        if number % i == 0:
            return False
    return True
num = 17
is_num_prime = is_prime(num)
</pre>
```

print(is_num_prime) # Output: True

Example 10: Converting a list of integers to a binary string

```
def int_list_to_binary_string(int_list):
    """
    This function converts a list of integers to a
binary string.
    """
    binary_string = ""
    for num in int_list:
        binary_string += bin(num)[2:].zfill(8)
    return binary_string

list_of_integers = [65, 66, 67]
binary_string =
int_list_to_binary_string(list_of_integers)
print(binary_string) # Output: 01000001010000101000011
```

These examples demonstrate the versatility and power of functions in Python. By writing functions, you can make your code more modular, reusable, and easier to maintain.

Example 11: Counting the frequency of words in a string

```
def count word frequency(string):
    This function counts the frequency of each word in
a string.
    11.11.11
    words = string.split()
    freq dict = {}
    for word in words:
        if word in freq dict:
            freq dict[word] += 1
        else:
            freq dict[word] = 1
    return freq dict
text = "This is a test sentence to check the word
frequency in a sentence."
word frequency = count word frequency(text)
print(word frequency) # Output: {'This': 1, 'is': 1,
'a': 2, 'test': 1, 'sentence': 2, 'to': 1, 'check': 1,
'the': 1, 'word': 1, 'frequency': 1, 'in': 1, 'a.': 1}
```



Example 12: Reversing a string

```
def reverse_string(string):
    """
    This function reverses a string.
    """
    reversed_string = ""
    for i in range(len(string)-1, -1, -1):
        reversed_string += string[i]
    return reversed_string

text = "Python is awesome"
reversed_text = reverse_string(text)
print(reversed_text)  # Output: emosewa si nohtyP
```

Example 13: Calculating the sum of digits in a number

```
def sum_of_digits(number):
    """
    This function calculates the sum of digits in a
number.
    """
    sum_of_digits = 0
    while number > 0:
        sum_of_digits += number % 10
        number //= 10
    return sum_of_digits
num = 1234
sum_of_num = sum_of_digits(num)
print(sum_of_num) # Output: 10
```

These examples demonstrate how functions can be used to perform various tasks, from counting the frequency of words in a string to calculating the sum of digits in a number. Functions provide a way to break down complex tasks into smaller, more manageable pieces of code. Example 14: Creating a simple calculator using functions

```
def add(num1, num2):
    """
    This function adds two numbers.
    """
    return num1 + num2
def subtract(num1, num2):
```



```
.....
    This function subtracts two numbers.
    .....
    return num1 - num2
def multiply(num1, num2):
    .....
    This function multiplies two numbers.
    11.11.11
    return num1 * num2
def divide(num1, num2):
    .....
    This function divides two numbers.
    11.11.11
    if num2 == 0:
        return "Error: division by zero"
    else:
        return num1 / num2
print("Choose an operation:")
print("1. Add")
print("2. Subtract")
print("3. Multiply")
print("4. Divide")
choice = input("Enter choice (1/2/3/4): ")
num1 = float(input("Enter first number: "))
num2 = float(input("Enter second number: "))
if choice == '1':
    print(num1, "+", num2, "=", add(num1, num2))
elif choice == '2':
    print(num1, "-", num2, "=", subtract(num1, num2))
elif choice == '3':
    print(num1, "*", num2, "=", multiply(num1, num2))
elif choice == '4':
    print(num1, "/", num2, "=", divide(num1, num2))
```



```
else:
    print("Invalid input")
```

This example demonstrates how functions can be used to create a simple calculator program. The program prompts the user to choose an operation (addition, subtraction, multiplication, or division) and then asks for two numbers to perform the operation on.

Example 15: Implementing a binary search algorithm using recursion

```
def binary search recursive(arr, start, end, target):
    11.11.11
    This function implements a binary search algorithm
using recursion.
    .....
    if end >= start:
        mid = (start + end) // 2
        if arr[mid] == target:
            return mid
        elif arr[mid] > target:
            return binary search recursive(arr, start,
mid-1, target)
        else:
            return binary search recursive(arr, mid+1,
end, target)
    else:
        return -1
arr = [2, 3, 4, 10, 40]
target = 10
result = binary search recursive(arr, 0, len(arr)-1,
target)
if result != -1:
    print(f"Element {target} is present at index
{result}")
else:
    print("Element is not present in array")
```

This example demonstrates how functions can be used to implement a binary search algorithm using recursion. The function takes an array, the starting and ending indices, and the target value as inputs, and returns the index of the target value in the array if it is present, or -1 if it is not present.



Example 16: Generating a random password

```
import random
import string
def generate_password(length):
    """
    This function generates a random password of the
given length.
    """
    password = ""
    for i in range(length):
        password += random.choice(string.ascii_letters
+ string.digits + string.punctuation)
    return password
password_length = 8
new_password = generate_password(password_length)
print(new_password)
```

This example demonstrates how functions can be used to generate a random password of a given length. The function uses the random.choice() function to select random characters from the set of uppercase and lowercase letters, digits, and punctuation.

Example 17: Converting temperature between Celsius and Fahrenheit

```
def celsius to fahrenheit(celsius):
    .....
    This function converts temperature from Celsius to
Fahrenheit.
    ** ** **
    fahrenheit = (celsius * 1.8) + 32
    return fahrenheit
def fahrenheit to celsius(fahrenheit):
    11 11 11
    This function converts temperature from Fahrenheit
to Celsius.
    11.11.11
    celsius = (fahrenheit - 32) / 1.8
    return celsius
temperature = float(input("Enter temperature: "))
scale = input("Enter scale (C/F): ")
```



```
if scale == "C":
    fahrenheit = celsius_to_fahrenheit(temperature)
    print(f"{temperature} degrees Celsius is equal to
    {fahrenheit} degrees Fahrenheit.")
elif scale == "F":
    celsius = fahrenheit_to_celsius(temperature)
    print(f"{temperature} degrees Fahrenheit is equal
    to {celsius} degrees Celsius.")
else:
    print("Invalid input.")
```

This example demonstrates how functions can be used to convert temperature between Celsius and Fahrenheit. The program prompts the user to enter a temperature and a scale (either Celsius or Fahrenheit), and then uses the appropriate function to convert the temperature to the other scale.

Example 18: Calculating the factorial of a number

```
def factorial(n):
    """
    This function calculates the factorial of a number.
    """
    if n == 0:
        return 1
    else:
        return n * factorial(n-1)

number = int(input("Enter a number: "))
result = factorial(number)
print(f"The factorial of {number} is {result}.")
```

This example demonstrates how functions can be used to calculate the factorial of a number. The function uses recursion to multiply the number by the factorial of the number minus one until it reaches the base case of zero, at which point it returns 1.

Example 19: Checking if a string is a palindrome

```
def is_palindrome(s):
    """
    This function checks if a string is a palindrome.
    """
    s = s.lower()
    s = s.replace(" ", "")
    if s == s[::-1]:
```



```
return True
else:
    return False
string = input("Enter a string: ")
if is_palindrome(string):
    print(f"{string} is a palindrome.")
else:
    print(f"{string} is not a palindrome.")
```

This example demonstrates how functions can be used to check if a string is a palindrome. The function first converts the string to lowercase and removes all whitespace, then checks if the string is equal to its reverse using slicing.

Example 20: Sorting a list of numbers

```
def sort list(numbers):
    11 11 11
    This function sorts a list of numbers in ascending
order.
    ......
    for i in range(len(numbers)):
        for j in range(i, len(numbers)):
            if numbers[i] > numbers[j]:
                temp = numbers[i]
                numbers[i] = numbers[j]
                numbers[j] = temp
    return numbers
list size = int(input("Enter the size of the list: "))
numbers = []
for i in range(list size):
    number = float(input(f"Enter number {i+1}: "))
    numbers.append(number)
sorted numbers = sort list(numbers)
print(f"The sorted list is: {sorted numbers}.")
```

This example demonstrates how functions can be used to sort a list of numbers in ascending order. The function uses a simple selection sort algorithm, which compares each pair of adjacent elements and swaps them if they are in the wrong order.

Example 21: Calculating the area and circumference of a circle



```
import math
def calculate area(radius):
    .....
    This function calculates the area of a circle with
the given radius.
    11.11.11
    area = math.pi * radius ** 2
    return area
def calculate circumference(radius):
    11.11.11
    This function calculates the circumference of a
circle with the given radius.
    11.11.11
    circumference = 2 * math.pi * radius
    return circumference
radius = float(input("Enter the radius of the circle:
"))
area = calculate area(radius)
circumference = calculate circumference(radius)
print(f"The area of the circle is {area:.2f} square
units and the circumference is {circumference:.2f}
units."
)
```

This example demonstrates how functions can be used to calculate the area and circumference of a circle. The program prompts the user to enter the radius of the circle, and then uses the appropriate functions to calculate the area and circumference. The math module is used to access the value of pi.

Example 22: Generating a Fibonacci sequence

```
def generate_fibonacci_sequence(n):
    """
    This function generates a Fibonacci sequence of
length n.
    """
    sequence = [0, 1]
    for i in range(2, n):
        next_number = sequence[i-1] + sequence[i-2]
        sequence.append(next_number)
    return sequence
```



```
sequence_length = int(input("Enter the length of the
Fibonacci sequence: "))
fibonacci_sequence =
generate_fibonacci_sequence(sequence_length)
print(f"The Fibonacci sequence of length
{sequence_length} is {fibonacci_sequence}.")
```

This example demonstrates how functions can be used to generate a Fibonacci sequence of a given length. The function uses a loop to generate the next number in the sequence, which is the sum of the two preceding numbers.

Example 23: Counting the occurrences of a character in a string

```
def count_occurrences(s, character):
    """
    This function counts the occurrences of a given
character in a string.
    """
    count = 0
    for c in s:
        if c == character:
            count += 1
    return count
string = input("Enter a string: ")
char = input("Enter a character: ")
count = count_occurrences(string, char)
print(f"The character '{char}' occurs {count} times in
the string '{string}'.")
```

This example demonstrates how functions can be used to count the occurrences of a given character in a string. The function uses a loop to iterate over each character in the string and checks if it matches the given character.

Example 24: Calculating the sum of a geometric series

```
def sum_geometric_series(a, r, n):
    """
    This function calculates the sum of a geometric
series with the given first term a, common ratio r, and
number of terms n.
    """
    if r == 1:
        return a * n
```



```
else:
    sum = a * (1 - r ** n) / (1 - r)
    return sum
first_term = float(input("Enter the first term of the
geometric series: "))
common_ratio = float(input("Enter the common ratio of
the geometric series: "))
num_terms = int(input("Enter the number of terms in the
geometric series: "))
sum = sum_geometric_series(first_term, common_ratio,
num_terms)
print(f"The sum of the geometric series is {sum:.2f}.")
```

This example demonstrates how functions can be used to calculate the sum of a geometric series with a given first term, common ratio, and number of terms. The function uses a formula to calculate the sum of a geometric series.

Libraries and modules

Here is a detailed overview of the libraries and modules covered in the book:

- 1. Pandas: Pandas is a powerful library for data manipulation and analysis in Python. It provides easy-to-use data structures and data analysis tools for handling large and complex datasets. The authors explain how to read SAS datasets into Pandas, perform basic data manipulations, and export data back to SAS.
- 2. NumPy: NumPy is a library for numerical computing in Python. It provides efficient and fast array operations for handling large datasets. The authors introduce NumPy and explain how to perform basic array operations, such as indexing, slicing, and reshaping.
- 3. Matplotlib: Matplotlib is a library for data visualization in Python. It provides a range of plotting tools for creating high-quality graphs and charts. The authors explain how to create basic plots and customize them using Matplotlib.
- 4. Seaborn: Seaborn is a library built on top of Matplotlib for creating statistical visualizations in Python. It provides a range of statistical plotting functions for visualizing distributions, relationships, and regression models. The authors introduce Seaborn and demonstrate how to create various types of statistical plots.
- 5. Scikit-learn: Scikit-learn is a library for machine learning in Python. It provides a range of tools for data preprocessing, feature selection, and model selection. The authors introduce Scikit-learn and explain how to train and evaluate various machine learning models.



- 6. Statsmodels: Statsmodels is a library for statistical modeling and testing in Python. It provides a range of statistical functions and models for data analysis. The authors introduce Statsmodels and demonstrate how to perform basic statistical analyses.
- 7. PySpark: PySpark is a Python library for Apache Spark, a big data processing framework. It provides a range of tools for distributed computing and processing large datasets. The authors introduce PySpark and demonstrate how to perform basic data manipulations and analysis using Spark.
- 8. SASPy: SASPy is a Python library for connecting to SAS and running SAS programs from Python. It provides a range of tools for accessing SAS datasets, running SAS procedures, and retrieving results. The authors introduce SASPy and demonstrate how to use it to connect to SAS and run SAS programs.

Here is a brief overview of some of the common libraries and modules used in Python for data analysis:

Pandas: Pandas is a popular data manipulation library that provides a high-performance data structure called a DataFrame, which is similar to a table in a relational database. Pandas allows you to perform various data manipulation operations, such as merging, filtering, sorting, and reshaping.

Example code for importing and using Pandas:

```
import pandas as pd
# Create a DataFrame from a CSV file
df = pd.read_csv('my_data.csv')
# Filter rows based on a condition
filtered_df = df[df['column_name'] > 0]
# Group data by a column and aggregate
grouped_df =
df.groupby('column_name').agg({'other_column': 'sum'})
# Merge two DataFrames based on a common column
merged df = pd.merge(df1, df2, on='common column')
```

Matplotlib: Matplotlib is a plotting library that allows you to create a wide range of visualizations, such as line charts, scatter plots, and histograms. Matplotlib provides a wide range of customization options for creating visually appealing and informative plots. Example code for importing and using Matplotlib:

```
import matplotlib.pyplot as plt
# Create a line chart
x = [1, 2, 3, 4]
y = [2, 4, 6, 8]
```



```
plt.plot(x, y)
# Create a scatter plot
x = [1, 2, 3, 4]
y = [2, 4, 6, 8]
plt.scatter(x, y)
# Create a histogram
data = [1, 2, 3, 3, 4, 5]
plt.hist(data)
```

Scikit-learn: Scikit-learn is a popular machine learning library that provides a wide range of algorithms for classification, regression, clustering, and dimensionality reduction. Scikit-learn also provides tools for preprocessing data, selecting features, and evaluating model performance.

Example code for importing and using Scikit-learn:

```
from sklearn.linear_model import LinearRegression
from sklearn.model_selection import train_test_split
from sklearn.metrics import mean_squared_error
# Load a dataset and split into training and testing
sets
X, y = load_dataset('my_data.csv')
X_train, X_test, y_train, y_test = train_test_split(X,
y, test_size=0.2)
# Train a linear regression model
model = LinearRegression()
model.fit(X_train, y_train)
# Evaluate the model on the test set
y_pred = model.predict(X_test)
mse = mean_squared_error(y_test, y_pred)
```

Seaborn: Seaborn is a data visualization library that provides a high-level interface for creating statistical graphics. Seaborn provides various built-in themes and color palettes to create visually appealing plots.

Example code for importing and using Seaborn:

```
import seaborn as sns
# Create a scatter plot with a regression line
```

```
in stal
```

```
sns.regplot(x='x_column', y='y_column', data=df)
# Create a bar chart with error bars
sns.barplot(x='category_column', y='value_column',
data=df, ci='sd')
# Create a box plot with a hue variable
sns.boxplot(x='group_column', y='value_column',
hue='category column', data=df)
```

Statsmodels: Statsmodels is a statistical modeling library that provides various tools for fitting and analyzing statistical models. Statsmodels provides a range of models, such as linear regression, logistic regression, and time series analysis.

Example code for importing and using Statsmodels:

```
import statsmodels.api as sm
# Fit a linear regression model
X = sm.add_constant(df['x_column'])
model = sm.OLS(df['y_column'], X)
results = model.fit()
# Print the model summary
print(results.summary())
# Fit a logistic regression model
X = df[['x1_column', 'x2_column']]
y = df['y_column']
model = sm.Logit(y, X)
results = model.fit()
# Print the model summary
print(results.summary())
```

PySpark: PySpark is a Python interface for Apache Spark, a distributed computing framework for processing large datasets. PySpark provides various tools for performing basic data manipulations, such as filtering, grouping, and joining. PySpark also provides various algorithms for machine learning and graph processing. Example code for importing and using PySpark:

```
from pyspark.sql import SparkSession
from pyspark.ml.feature import VectorAssembler
from pyspark.ml.regression import LinearRegression
```



```
# Create a SparkSession
spark =
SparkSession.builder.appName('my app').getOrCreate()
# Load a dataset and create a DataFrame
df = spark.read.csv('my data.csv', header=True)
# Perform basic data manipulations
filtered df = df.filter(df['column name'] > 0)
grouped df =
df.groupBy('column name').sum('other column')
joined df = df.join(other df, on='common column')
# Perform machine learning tasks
assembler = VectorAssembler(inputCols=['x1 column',
'x2 column'], outputCol='features')
data = assembler.transform(df)
lr = LinearRegression(featuresCol='features',
labelCol='y column')
model = lr.fit(data)
predictions = model.transform(data)
```

NetworkX: NetworkX is a Python library for creating and analyzing graphs and networks. NetworkX provides various tools for creating and visualizing graphs, as well as algorithms for analyzing and manipulating them.

Example code for importing and using NetworkX:

```
import networkx as nx
# Create a graph
G = nx.Graph()
# Add nodes and edges
G.add_node('A')
G.add_nodes_from(['B', 'C'])
G.add_edge('A', 'B')
G.add_edges_from([('B', 'C'), ('A', 'C')])
# Visualize the graph
nx.draw(G, with_labels=True)
# Compute the degree centrality of each node
centrality = nx.degree_centrality(G)
```



print(centrality)

NLTK: NLTK (Natural Language Toolkit) is a Python library for working with human language data. NLTK provides various tools for tokenizing, tagging, parsing, and analyzing text data.

Example code for importing and using NLTK:

```
import nltk
# Download the NLTK data
nltk.download()
# Tokenize a sentence
text = 'This is a sentence.'
tokens = nltk.word tokenize(text)
print(tokens)
# Tag the tokens with their part of speech
tagged = nltk.pos tag(tokens)
print(tagged)
# Analyze the sentiment of a text
text = 'This is a very good movie.'
sentiment =
nltk.sentiment.vader.SentimentIntensityAnalyzer().polar
ity scores(text)
print(sentiment)
```

TensorFlow: TensorFlow is a popular open-source framework for building and training machine learning models. TensorFlow provides various tools for working with neural networks, as well as for performing other machine learning tasks.

Example code for importing and using TensorFlow:

```
import tensorflow as tf
# Create a neural network model
model = tf.keras.models.Sequential([
   tf.keras.layers.Dense(64, activation='relu',
   input_shape=(784,)),
   tf.keras.layers.Dense(10, activation='softmax')
])
```



```
# Compile the model
model.compile(optimizer='adam',
loss='categorical_crossentropy', metrics=['accuracy'])
# Train the model
model.fit(x_train, y_train, epochs=10,
validation_data=(x_test, y_test))
# Make predictions with the model
predictions = model.predict(x test)
```

Scikit-learn: Scikit-learn is a Python library for machine learning that provides various tools for data preprocessing, model selection, and evaluation. Scikit-learn provides a wide range of models, such as regression, classification, clustering, and dimensionality reduction. Example code for importing and using Scikit-learn:

```
from sklearn.linear model import LinearRegression
from sklearn.datasets import load boston
from sklearn.model selection import train test split
# Load the Boston Housing dataset
boston = load boston()
# Split the dataset into training and testing sets
X train, X test, y train, y test =
train test split(boston.data, boston.target,
test size=0.2, random state=42)
# Create a linear regression model
model = LinearRegression()
# Train the model
model.fit(X train, y train)
# Evaluate the model on the testing set
score = model.score(X test, y test)
print(score)
```

PyTorch: PyTorch is another popular open-source framework for building and training machine learning models. PyTorch provides various tools for working with neural networks, as well as for performing other machine learning tasks.

Example code for importing and using PyTorch:

in stal

```
import torch
import torch.nn as nn
import torch.optim as optim
# Create a neural network model
class Net(nn.Module):
    def init (self):
        super(Net, self). init ()
        self.fc1 = nn.Linear(784, 64)
        self.fc2 = nn.Linear(64, 10)
    def forward(self, x):
        x = nn.functional.relu(self.fcl(x))
        x = nn.functional.softmax(self.fc2(x), dim=1)
        return x
model = Net()
# Define a loss function and an optimizer
criterion = nn.CrossEntropyLoss()
optimizer = optim.Adam(model.parameters(), lr=0.001)
# Train the model
for epoch in range(10):
    running loss = 0.0
    for i, data in enumerate(trainloader, 0):
        inputs, labels = data
        optimizer.zero grad()
        outputs = model(inputs)
        loss = criterion(outputs, labels)
        loss.backward()
        optimizer.step()
        running loss += loss.item()
    print(f"Epoch {epoch+1}, loss:
{running loss/len(trainloader)}")
# Make predictions with the model
outputs = model(inputs)
, predicted = torch.max(outputs, 1)
```

Statsmodels: Statsmodels is a Python library for statistical modeling and analysis. Statsmodels provides various tools for performing regression analysis, time series analysis, and hypothesis testing.



Example code for importing and using Statsmodels:

```
import statsmodels.api as sm
import statsmodels.formula.api as smf
import pandas as pd
# Load the data into a pandas DataFrame
df = pd.read_csv('data.csv')
# Fit a linear regression model
model = smf.ols('y ~ x1 + x2 + x3', data=df).fit()
# Print the model summary
print(model.summary())
# Perform a hypothesis test
result = model.t_test('x1 = 0')
print(result)
```

Matplotlib: Matplotlib is a Python library for creating static, animated, and interactive visualizations. Matplotlib provides various tools for creating line plots, scatter plots, bar plots, and more.

Example code for importing and using Matplotlib:

```
import matplotlib.pyplot as plt
import numpy as np
# Create some data
x = np.linspace(0, 10, 100)
y1 = np.sin(x)
y2 = np.cos(x)
# Create a line plot
plt.plot(x, y1, label='sin(x)')
plt.plot(x, y2, label='cos(x)')
# Add a title and axis labels
plt.title('Trigonometric Functions')
plt.ylabel('x')
plt.ylabel('y')
# Add a legend
plt.legend()
```



Show the plot plt.show()

These are just a few more examples of the many libraries and modules available in Python for data analysis, machine learning, statistical modeling, and visualization. By combining these tools with the fundamental programming skills, you can tackle various real-world data analysis tasks and build powerful machine learning models.

Seaborn: Seaborn is a Python library based on Matplotlib that provides additional functionality for creating beautiful and informative statistical visualizations. Seaborn provides tools for creating visualizations such as heatmaps, violin plots, and regression plots. Example code for importing and using Seaborn:

```
import seaborn as sns
import pandas as pd
# Load the data into a pandas DataFrame
df = pd.read_csv('data.csv')
# Create a heatmap of the correlation matrix
corr_matrix = df.corr()
sns.heatmap(corr_matrix, annot=True)
# Create a violin plot
sns.violinplot(x='group', y='value', data=df)
# Create a regression plot
sns.regplot(x='x', y='y', data=df)
```

NLTK: The Natural Language Toolkit (NLTK) is a Python library for working with human language data. NLTK provides various tools for tasks such as tokenization, stemming, part-of-speech tagging, and sentiment analysis.

Example code for importing and using NLTK:

```
import nltk
from nltk.tokenize import word_tokenize
from nltk.stem import PorterStemmer
from nltk.corpus import stopwords
# Tokenize a sentence
sentence = "The quick brown fox jumps over the lazy
dog."
tokens = word_tokenize(sentence)
```



```
# Stem a word
stemmer = PorterStemmer()
stemmed_word = stemmer.stem('jumped')
# Remove stop words from a sentence
stop_words = set(stopwords.words('english'))
filtered_sentence = [w for w in tokens if not w in
stop_words]
```

Flask: Flask is a Python web framework for building web applications. Flask provides tools for creating routes, rendering templates, and handling requests and responses. Example code for importing and using Flask:

```
from flask import Flask, render_template, request
app = Flask(__name__)
# Define a route that renders a template
@app.route('/')
def home():
    return render_template('home.html')
# Define a route that handles a POST request
@app.route('/login', methods=['POST'])
def login():
    username = request.form['username']
    password = request.form['password']
    return f'Username: {username}, Password:
{password}'

if __name__ == '__main__':
    app.run(debug=True)
```

Pygame: Pygame is a Python library for creating games and multimedia applications. Pygame provides tools for handling graphics, sounds, and user input. Example code for importing and using Pygame:

```
import pygame
pygame.init()
# Set the dimensions of the game window
WINDOW_WIDTH = 800
```



```
WINDOW HEIGHT = 600
# Create the game window
window = pygame.display.set mode((WINDOW WIDTH,
WINDOW HEIGHT))
# Load an image
image = pygame.image.load('image.png')
# Create a game loop
running = True
while running:
    # Handle events
    for event in pygame.event.get():
        if event.type == pygame.QUIT:
            running = False
    # Draw the image on the window
    window.blit(image, (0, 0))
    # Update the window
    pygame.display.update()
# Quit the game
pygame.quit()
```

Requests: The Requests library is a Python library for making HTTP requests. Requests provides tools for handling headers, cookies, authentication, and more. Example code for importing and using Requests:

```
import requests
# Make a GET request
response = requests.get('https://example.com')
# Print the response content
print(response.text)
# Make a POST request with data
data = {'username': 'example', 'password': 'password'}
response = requests.post('https://example.com/login',
data=data)
```



Print the response status code print(response.status_code)

Pandas: Pandas is a Python library for working with tabular data. Pandas provides tools for reading and writing data from various file formats, manipulating data, and performing calculations and statistics on data.

Example code for importing and using Pandas:

```
import pandas as pd
# Read data from a CSV file
data = pd.read_csv('data.csv')
# Group the data by a column and calculate the mean of
another column
grouped_data = data.groupby('category')['value'].mean()
# Filter the data based on a condition
filtered_data = data[data['value'] > 10]
# Write the data to a CSV file
filtered_data.to_csv('filtered_data.csv')
```

Scikit-learn: Scikit-learn is a Python library for machine learning. Scikit-learn provides tools for preprocessing data, creating models, evaluating models, and more. Example code for importing and using Scikit-learn:

```
from sklearn import datasets
from sklearn.model_selection import train_test_split
from sklearn.linear_model import LinearRegression
from sklearn.metrics import mean_squared_error
# Load the Boston Housing dataset
boston = datasets.load_boston()
# Split the data into training and testing sets
X_train, X_test, y_train, y_test =
train_test_split(boston.data, boston.target,
test_size=0.2, random_state=42)
# Create a linear regression model
model = LinearRegression()
```



```
# Train the model
model.fit(X_train, y_train)
# Make predictions on the testing set
y_pred = model.predict(X_test)
# Calculate the mean squared error
mse = mean_squared_error(y_test, y_pred)
# Print the mean squared error
print(mse)
```

NLTK: The Natural Language Toolkit (NLTK) is a Python library for natural language processing. NLTK provides tools for tokenizing text, part-of-speech tagging, named entity recognition, and more.

Example code for importing and using NLTK:

```
import nltk
```

```
# Download the necessary resources
nltk.download('punkt')
nltk.download('averaged_perceptron_tagger')
# Tokenize a sentence
sentence = 'This is a sentence.'
tokens = nltk.word_tokenize(sentence)
# Part-of-speech tag the tokens
pos_tags = nltk.pos_tag(tokens)
# Print the part-of-speech tags
print(pos tags)
```

Matplotlib: Matplotlib is a Python library for creating visualizations, such as line plots, scatter plots, and bar plots. Matplotlib provides tools for customizing the appearance of the plots and adding labels and titles.

Example code for importing and using Matplotlib:

```
import matplotlib.pyplot as plt
# Create data
x = [1, 2, 3, 4, 5]
y = [2, 4, 6, 8, 10]
```



```
# Create a line plot
plt.plot(x, y)
# Add labels and a title
plt.xlabel('X-axis')
plt.ylabel('Y-axis')
plt.title('Line Plot')
# Show the plot
plt.show()
```

Pandas: Pandas is a Python library for data manipulation and analysis. Pandas provides tools for reading and writing data, cleaning and preprocessing data, and performing operations on data, such as merging and grouping.

Example code for importing and using Pandas:

```
import pandas as pd
# Read a CSV file
df = pd.read_csv('data.csv')
# Drop missing values
df = df.dropna()
# Group by a column and calculate the mean
grouped = df.groupby('Category')['Value'].mean()
# Print the result
print(grouped)
```

TensorFlow: TensorFlow is a Python library for machine learning and deep learning. TensorFlow provides tools for building and training neural networks, as well as for deploying models to production environments.

Example code for importing and using TensorFlow:

```
import tensorflow as tf
from tensorflow import keras
# Load the MNIST dataset
mnist = keras.datasets.mnist
(train_images, train_labels), (test_images,
test_labels) = mnist.load_data()
```



```
# Preprocess the data
train images = train images / 255.0
test images = test images / 255.0
# Define the model
model = keras.Sequential([
    keras.layers.Flatten(input shape=(28, 28)),
    keras.layers.Dense(128, activation='relu'),
    keras.layers.Dense(10)
1)
# Compile the model
model.compile(optimizer='adam',
loss=tf.keras.losses.SparseCategoricalCrossentropy(from
logits=True),
              metrics=['accuracy'])
# Train the model
model.fit(train images, train labels, epochs=10,
validation data=(test images, test labels))
# Evaluate the model
test loss, test acc = model.evaluate(test images,
test labels, verbose=2)
print('Test accuracy:', test acc)
```

Pygame: Pygame is a Python library for game development. Pygame provides tools for creating games and graphics, handling user input, and playing sounds and music. Example code for importing and using Pygame:

```
import pygame
# Initialize Pygame
pygame.init()
# Create a window
screen = pygame.display.set_mode((640, 480))
# Set the window title
pygame.display.set_caption('My Game')
# Main game loop
```



```
running = True
while running:
    # Handle events
    for event in pygame.event.get():
        if event.type == pygame.QUIT:
            running = False
    # Clear the screen
    screen.fill((255, 255, 255))
    # Draw a rectangle
    pygame.draw.rect(screen, (255, 0, 0),
    pygame.Rect(100, 100, 50, 50))
    # Update the screen
    pygame.display.flip()
    # Quit Pygame
    pygame.quit()
```

Matplotlib: Matplotlib is a Python library for data visualization. Matplotlib provides tools for creating a wide range of plots and charts, including line plots, scatter plots, bar charts, and more. Example code for importing and using Matplotlib:

```
import matplotlib.pyplot as plt
# Generate some data
x = [1, 2, 3, 4, 5]
y = [1, 4, 9, 16, 25]
# Create a line plot
plt.plot(x, y)
# Set the plot title and axis labels
plt.title('My Plot')
plt.xlabel('X')
plt.ylabel('Y')
# Show the plot
plt.show()
```

NumPy: NumPy is a Python library for numerical computing. NumPy provides tools for working with large arrays and matrices of numerical data, and for performing mathematical operations on these arrays and matrices.

Example code for importing and using NumPy:



```
import numpy as np
# Create a NumPy array
arr = np.array([1, 2, 3, 4, 5])
# Perform mathematical operations on the array
arr_squared = arr ** 2
arr_sum = arr.sum()
# Print the results
print(arr_squared)
print(arr_sum)
```

Requests: Requests is a Python library for making HTTP requests. Requests provides tools for sending and receiving HTTP requests and responses, and for working with web APIs and web data.

Example code for importing and using Requests:

```
import requests
# Send an HTTP GET request
response =
requests.get('https://api.github.com/users/octocat/repo
s')
# Check the response status code
if response.status_code == 200:
    # Get the response data as a JSON object
    data = response.json()

    # Loop over the data and print each repository name
    for repo in data:
        print(repo['name'])
else:
    # Handle error
    print('Error: Failed to retrieve data')
```

Pandas: Pandas is a Python library for data manipulation and analysis. Pandas provides tools for working with tabular data (i.e., data in rows and columns) and performing operations on this data.

Example code for importing and using Pandas:

import pandas as pd



```
# Create a Pandas dataframe
data = {'Name': ['Alice', 'Bob', 'Charlie', 'Dave'],
'Age': [25, 30, 35, 40]}
df = pd.DataFrame(data)
# Perform operations on the dataframe
df_filtered = df[df['Age'] > 30]
df_sorted = df.sort_values('Age')
# Print the results
print(df_filtered)
print(df_sorted)
```

Scikit-learn: Scikit-learn is a Python library for machine learning. Scikit-learn provides tools for building and training machine learning models, as well as tools for evaluating and using these models.

Example code for importing and using Scikit-learn:

```
from sklearn import datasets
from sklearn.linear_model import LinearRegression
# Load the diabetes dataset
diabetes = datasets.load_diabetes()
# Get the features and target variables
X = diabetes.data
y = diabetes.target
# Train a linear regression model
model = LinearRegression()
model.fit(X, y)
# Use the model to make predictions
y_pred = model.predict(X[:5])
# Print the results
print(y pred)
```

TensorFlow: TensorFlow is a Python library for machine learning and artificial intelligence. TensorFlow provides tools for building and training deep neural networks, as well as tools for evaluating and using these models.

Example code for importing and using TensorFlow:

```
import tensorflow as tf
```

```
# Create a simple neural network
model = tf.keras.Sequential([
    tf.keras.layers.Dense(64, activation='relu',
input shape=(784,)),
    tf.keras.layers.Dense(10, activation='softmax')
1)
# Compile the model
model.compile(optimizer='adam',
              loss='categorical crossentropy',
              metrics=['accuracy'])
# Train the model
model.fit(x train, y train, epochs=5)
# Use the model to make predictions
y pred = model.predict(x test)
# Print the results
print(y pred)
```

Flask: Flask is a Python web framework that allows you to build web applications quickly and easily. Flask provides tools for routing, handling HTTP requests and responses, and rendering HTML templates.

Example code for importing and using Flask:

```
from flask import Flask, render_template
app = Flask(__name__)
# Define a route for the homepage
@app.route('/')
def index():
    return render_template('index.html')
# Define a route for a contact page
@app.route('/contact')
def contact():
    return render_template('contact.html')
# Run the app
if __name__ == '__main__':
```



app.run(debug=True)

Pygame: Pygame is a Python library for building games and multimedia applications. Pygame provides tools for handling graphics, input events, sound, and more. Example code for importing and using Pygame:

import pygame

```
# Initialize Pygame
pygame.init()
# Create a Pygame window
screen = pygame.display.set mode((640, 480))
# Set the window title
pygame.display.set caption('My Pygame Window')
# Run the Pygame loop
running = True
while running:
    for event in pygame.event.get():
        if event.type == pygame.QUIT:
            running = False
    # Draw something on the screen
    pygame.draw.rect(screen, (255, 0, 0), (100, 100,
200, 200))
    # Update the display
    pygame.display.update()
# Quit Pygame
pygame.quit()
```

BeautifulSoup: BeautifulSoup is a Python library for parsing HTML and XML documents. BeautifulSoup provides tools for extracting data from HTML and XML documents, such as finding specific elements and attributes.

Example code for importing and using BeautifulSoup:

from bs4 import BeautifulSoup
import requests
Get the HTML content of a webpage



```
url = 'https://www.example.com'
response = requests.get(url)
html_content = response.content
# Parse the HTML content with BeautifulSoup
soup = BeautifulSoup(html_content, 'html.parser')
# Find a specific element in the HTML content
element = soup.find('h1')
# Print the text content of the element
print(element.text)
```

NumPy: NumPy is a Python library for numerical computing. NumPy provides tools for working with large, multi-dimensional arrays and matrices, and for performing mathematical operations on these arrays efficiently.

Example code for importing and using NumPy:

```
import numpy as np
# Create a NumPy array
a = np.array([[1, 2], [3, 4]])
# Perform a mathematical operation on the array
b = a * 2
# Print the resulting array
print(b)
```

Pandas: Pandas is a Python library for data manipulation and analysis. Pandas provides tools for reading and writing data in various formats, such as CSV and Excel, and for performing operations on the data, such as filtering, grouping, and joining. Example code for importing and using Pandas:

```
import pandas as pd
# Read a CSV file into a Pandas dataframe
df = pd.read_csv('data.csv')
# Filter the dataframe by a condition
filtered_df = df[df['column'] > 10]
# Group the dataframe by a column and calculate the
mean of another column
```



```
grouped_df = df.groupby('column1')['column2'].mean()
# Join two dataframes on a common column
joined_df = pd.merge(df1, df2, on='common_column')
# Write the resulting dataframe to a CSV file
joined df.to csv('result.csv', index=False)
```

Matplotlib: Matplotlib is a Python library for creating visualizations, such as charts, graphs, and plots. Matplotlib provides tools for creating and customizing these visualizations, and for adding labels and annotations to them.

Example code for importing and using Matplotlib:

```
import matplotlib.pyplot as plt
import numpy as np
# Create a NumPy array for the x-axis data
x = np.arange(0, 10, 0.1)
# Create a NumPy array for the y-axis data
y = np.sin(x)
# Create a Matplotlib plot with the x-axis and y-axis
data
plt.plot(x, y)
# Add a title and labels to the plot
plt.title('Sinusoidal Plot')
plt.xlabel('x-axis')
plt.ylabel('y-axis')
# Show the plot
plt.show()
```

Scikit-learn: Scikit-learn is a Python library for machine learning. Scikit-learn provides tools for various machine learning tasks, such as classification, regression, clustering, and dimensionality reduction. It also includes tools for data preprocessing, cross-validation, and model selection. Example code for importing and using Scikit-learn:

```
from sklearn.datasets import load_iris
from sklearn.linear_model import LogisticRegression
from sklearn.model_selection import train_test_split
```



```
# Load the iris dataset
iris = load_iris()
# Split the dataset into training and testing sets
X_train, X_test, y_train, y_test =
train_test_split(iris.data, iris.target, test_size=0.2)
# Create a logistic regression model and fit it to the
training data
model = LogisticRegression()
model.fit(X_train, y_train)
# Use the model to predict the classes of the testing
data
y_pred = model.predict(X_test)
# Print the accuracy of the model
print("Accuracy:", model.score(X_test, y_test))
```

TensorFlow: TensorFlow is a Python library for building and training machine learning models, particularly deep learning models. TensorFlow provides tools for building and manipulating complex graphs of mathematical operations, which can be executed efficiently on CPUs or GPUs. It also includes tools for data preprocessing, model evaluation, and deployment. Example code for importing and using TensorFlow:

```
import tensorflow as tf
import numpy as np
# Create a TensorFlow graph for a simple linear
regression model
x = tf.placeholder(tf.float32)
y = tf.placeholder(tf.float32)
W = tf.Variable(0.0)
b = tf.Variable(0.0)
y pred = W * x + b
loss = tf.reduce_mean(tf.square(y_pred - y))
optimizer = tf.train.GradientDescentOptimizer(0.01)
train op = optimizer.minimize(loss)
# Generate some random data for the model
x data = np.random.randn(100)
y data = 2 \times x data + 1 + np.random.randn(100) \times 0.1
# Train the model on the data
```

```
with tf.Session() as sess:
    sess.run(tf.global_variables_initializer())
    for i in range(1000):
        _, loss_val = sess.run([train_op, loss],
    feed_dict={x: x_data, y: y_data})
        if i % 100 == 0:
            print("Step:", i, "Loss:", loss_val)
    # Use the trained model to make predictions
    y_pred_val = sess.run(y_pred, feed_dict={x: [1, 2,
3]})
    print("Predictions:", y_pred_val)
```

Keras: Keras is a high-level neural network API that runs on top of TensorFlow or other backend engines. Keras provides a simplified interface for building and training neural networks, with tools for creating various types of layers and models, and for compiling and fitting models to data.

Example code for importing and using Keras:

```
import keras
from keras.models import Sequential
from keras.layers import Dense
# Create a Keras model for a simple neural network
model = Sequential()
model.add(Dense(32, activation='relu', input dim=100))
model.add(Dense(1, activation='sigmoid'))
# Compile the model with an optimizer and loss function
model.compile(optimizer='rmsprop',
loss='binary crossentropy', metrics=['accuracy'])
# Generate some random data for the model
import numpy as np
data = np.random.random((1000, 100))
labels = np.random.randint(2, size=(1000, 1))
# Train the model on the data
model.fit(data, labels, epochs=10, batch]}
```

Chapter 2: Data Structures in Python



Lists

Python is a popular, high-level programming language used for various purposes, including data analysis, machine learning, web development, and more. In recent years, Python has become increasingly popular among data analysts, replacing traditional tools like SAS and R. As a result, many SAS users are looking to learn Python to expand their skill set and stay up-to-date with the latest trends in data analysis.

If you are a SAS user looking to learn Python, it can be helpful to start by understanding some of the similarities and differences between the two languages. While both Python and SAS are used for data analysis, there are some key differences in their syntax, approach, and functionality. For example, Python is a more general-purpose programming language, while SAS is specifically designed for statistical analysis. Additionally, Python is an open-source language, while SAS is a proprietary tool that requires a license.

Despite these differences, there are many ways that Python and SAS can be used together to enhance your data analysis capabilities. One of the most useful features of Python is its ability to work with a wide range of data formats, including Excel spreadsheets, SQL databases, and even SAS datasets. By using Python with SAS, you can take advantage of the strengths of both tools to create more sophisticated analyses and visualizations.

One of the key features of Python is its use of lists, which are a type of data structure used to store collections of items. In Python, lists are similar to arrays in other programming languages, but they offer some additional flexibility and functionality. For example, lists can be resized dynamically, which means that you can add or remove items from the list as needed. Additionally, Python lists can contain any type of data, including numbers, strings, and even other lists.

To work with lists in Python, you can use a variety of built-in functions and methods. Some of the most commonly used functions for working with lists include len(), which returns the length of a list; min() and max(), which return the smallest and largest values in a list, respectively; and sum(), which returns the sum of all the values in a list. Additionally, Python lists support a variety of methods for manipulating and accessing the data in the list, including append(), which adds an item to the end of the list; insert(), which adds an item at a specific position in the list; and pop(), which removes an item from the list.

As a SAS user, you may find that working with lists in Python can be a useful way to organize and manipulate your data. For example, you might use a list to store a collection of values that represent a variable in your dataset. By working with lists in Python, you can easily perform operations on these values, such as calculating the mean or standard deviation. Additionally, you can use Python's built-in functions and methods to sort, filter, and transform your data in a variety of ways.

Overall, learning to work with lists in Python can be a valuable skill for SAS users who want to expand their data analysis capabilities. By understanding how lists work and how to use them



effectively, you can take advantage of the power and flexibility of Python to enhance your data analysis workflows. In addition to lists, there are several other Python data structures that can be useful for SAS users to learn. One of these is the dictionary, which is a collection of key-value pairs. Dictionaries can be used to store and manipulate data in a way that is similar to SAS data sets, and they offer a lot of flexibility and functionality for data analysis tasks.

In Python, dictionaries are created using curly braces {} and can contain any combination of keys and values. Keys must be unique and can be any immutable data type, such as a string or a number. Values can be any data type, including lists or other dictionaries. To access values in a dictionary, you can use the keys as the index.

Python for SAS Users: A SAS-Oriented Introduction to Python is a book that introduces SAS users to the Python programming language. The book focuses on the similarities and differences between the two languages, and provides a comprehensive guide for SAS users who want to learn Python.

One of the fundamental data structures in Python is the list. A list is an ordered collection of elements that can be of any type, including other lists. Lists are mutable, which means that they can be changed after they are created.

To create a list in Python, you can use square brackets. For example, to create a list of integers, you can write:

my list = [1, 2, 3, 4, 5]

You can also create a list of strings:

my list = ["apple", "banana", "cherry"]

To access an element in a list, you can use its index. The index of the first element in a list is 0, and the index of the last element is the length of the list minus one. For example, to access the first element of the list my_list, you can write:

first element = my list[0]

To add an element to a list, you can use the append() method. For example, to add the string "orange" to the end of the list my_list, you can write:

my list.append("orange")

To remove an element from a list, you can use the remove() method. For example, to remove the string "banana" from the list my_list, you can write:

my_list.remove("banana")



You can also use the pop() method to remove an element from a list by its index. For example, to remove the second element from the list my_list, you can write:

```
second element = my list.pop(1)
```

Lists can be used in a variety of ways in Python. For example, you can use them to store data, to represent matrices, and to implement stacks and queues.

In addition to lists, Python provides several other built-in data structures, including tuples, sets, and dictionaries. Each of these data structures has its own unique features and can be used in different ways.

Here is an example of Python code that demonstrates some of the basic operations on lists:

```
# create a list of integers
my list = [1, 2, 3, 4, 5]
# create a list of strings
fruit list = ["apple", "banana", "cherry"]
# access an element in a list
first fruit = fruit list[0]
print("The first fruit is:", first fruit)
# add an element to a list
fruit list.append("orange")
print("The updated fruit list is:", fruit list)
# remove an element from a list
fruit list.remove("banana")
print("The updated fruit list is:", fruit list)
# remove an element from a list by its index
second fruit = fruit list.pop(1)
print("The second fruit is:", second fruit)
print("The updated fruit list is:", fruit list)
# create a list of lists
matrix = [[1, 2, 3], [4, 5, 6], [7, 8, 9]]
# access an element in a list of lists
second row = matrix[1]
print("The second row of the matrix is:", second row)
```



access an element in a list of lists using nested indexing element = matrix[1][2]

print("The element in the second row and third column of the matrix is:", element) In this code, we first create a list of integers called my_list and a list of strings called fruit_list. We then access the first element of fruit_list using indexing and print it to the console.

Next, we add the string "orange" to fruit_list using the append() method and print the updated list to the console. We then remove the string "banana" from fruit_list using the remove() method and print the updated list again.

We then remove the second element of fruit_list (which is "cherry") using the pop() method and store it in the variable second_fruit. We print second_fruit to the console to verify that it is the correct value, and then print fruit_list again to see the updated list.

Next, we create a list of lists called matrix and access the second row of the matrix using indexing. We print the second row to the console to verify that it is the correct value.

In addition to the basic list operations shown in the previous code example, there are many other ways to manipulate lists in Python. Here are a few more examples:

```
# concatenate two lists
list1 = [1, 2, 3]
list2 = [4, 5, 6]
concatenated list = list1 + list2
print("The concatenated list is:", concatenated list)
# sort a list
unsorted list = [3, 1, 4, 2, 5]
sorted list = sorted(unsorted list)
print("The sorted list is:", sorted list)
# reverse a list
original list = [1, 2, 3, 4, 5]
reversed list = list(reversed(original list))
print("The reversed list is:", reversed list)
# count the occurrences of an element in a list
count list = [1, 2, 2, 3, 3, 3]
num twos = count list.count(2)
num threes = count list.count(3)
print("The number of twos in the list is:", num twos)
```

print("The number of threes in the list is:", num_threes)

In the first example, we concatenate two lists (list1 and list2) using the + operator and store the result in concatenated_list. We print concatenated_list to the console to verify that it contains all of the elements from list1 and list2.

In the second example, we create an unsorted list called unsorted_list and sort it using the sorted() function. We store the sorted list in sorted_list and print it to the console to verify that it is sorted.

In the third example, we create an original list called original_list and reverse it using the reversed() function. We then convert the reversed iterator to a list using the list() function and store the result in reversed_list. We print reversed_list to the console to verify that it is the reverse of original_list.

In the fourth example, we create a list called count_list that contains multiple occurrences of the values 1, 2, and 3. We then use the count() method to count the number of occurrences of the values 2 and 3 in count_list, and print the results to the console.

Here are a few more examples of list operations in Python:

```
# find the index of an element in a list
fruits = ["apple", "banana", "cherry"]
index of banana = fruits.index("banana")
print("The index of 'banana' in the list is:",
index of banana)
# insert an element into a list at a specific index
nums = [1, 2, 3, 4, 5]
nums.insert(2, 2.5)
print("The updated list is:", nums)
# extend a list with another list
list1 = [1, 2, 3]
list2 = [4, 5, 6]
list1.extend(list2)
print("The extended list is:", list1)
# slice a list to create a new list
original list = [1, 2, 3, 4, 5]
sliced list = original list[1:4]
print("The sliced list is:", sliced list)
```



In the first example, we have a list of fruits and we use the index() method to find the index of "banana" in the list. We print the result to the console.

In the second example, we have a list of numbers and we use the insert() method to insert the value 2.5 into the list at index 2. We print the updated list to the console.

In the third example, we have two lists (list1 and list2) and we use the extend() method to add the elements of list2 to list1. We print the extended list to the console.

In the fourth example, we have a list of numbers and we use slicing to create a new list that contains only the elements with indices 1 through 3. We print the sliced list to the console.

Here are a few more examples of list operations in Python:

```
# remove the first occurrence of an element in a list
nums = [1, 2, 3, 2, 4]
nums.remove(2)
print("The updated list is:", nums)
# remove an element from a list by index
nums = [1, 2, 3, 4, 5]
removed num = nums.pop(2)
print("The updated list is:", nums)
print("The removed element is:", removed num)
# create a list of numbers using a loop
squares = []
for i in range(1, 6):
    squares.append(i**2)
print("The list of squares is:", squares)
# create a list of numbers using a list comprehension
squares = [i**2 \text{ for } i \text{ in range}(1, 6)]
print("The list of squares is:", squares)
```

In the first example, we have a list of numbers and we use the remove() method to remove the first occurrence of the value 2 from the list. We print the updated list to the console.

In the second example, we have a list of numbers and we use the pop() method to remove the element at index 2 from the list and store it in the variable removed_num. We print the updated list and the removed element to the console.

In the third example, we create an empty list called squares and use a for loop to append the squares of the numbers 1 through 5 to the list. We print the list of squares to the console.



In the fourth example, we use a list comprehension to create the same list of squares as in the previous example. The list comprehension is a more concise and readable way to create a list based on a simple pattern.

Here are a few more examples of list operations in Python:

```
# sort a list of strings alphabetically
fruits = ["apple", "cherry", "banana"]
fruits.sort()
print("The sorted list is:", fruits)
# sort a list of numbers in descending order
nums = [3, 6, 1, 8, 2]
nums.sort(reverse=True)
print("The sorted list is:", nums)
# reverse the order of a list
letters = ["a", "b", "c", "d"]
letters.reverse()
print("The reversed list is:", letters)
# check if an element is in a list
nums = [1, 2, 3, 4, 5]
if 3 in nums:
    print("3 is in the list.")
else:
   print("3 is not in the list.")
```

In the first example, we have a list of strings and we use the sort() method to sort the list alphabetically. We print the sorted list to the console.

In the second example, we have a list of numbers and we use the sort() method with the reverse=True argument to sort the list in descending order. We print the sorted list to the console.

In the third example, we have a list of letters and we use the reverse() method to reverse the order of the list. We print the reversed list to the console.

In the fourth example, we have a list of numbers and we use an if statement with the in operator to check if the value 3 is in the list. We print the result to the console.

Tuples

In Python, a tuple is an immutable sequence of objects. It is similar to a list, but once a tuple is created, its contents cannot be modified. Tuples are often used to represent fixed collections of related data, such as a point in 2D space or a date (year, month, day).

Tuples are defined using parentheses () and the objects inside the tuple are separated by commas. Here's an example:

my tuple = (1, 2, 3, "hello", True)

In this example, my_tuple is a tuple containing five objects: the integers 1, 2, and 3, the string "hello", and the boolean value True.

Tuples can also be created without using parentheses, as long as there is more than one object in the tuple:

my tuple = 1, 2, 3

In this case, my_tuple is still a tuple containing the integers 1, 2, and 3.

Accessing elements in a tuple is similar to accessing elements in a list. You can use square brackets and an index to access a specific element:

my_tuple = (1, 2, 3, "hello", True)
print(my_tuple[0]) # prints 1
print(my_tuple[3]) # prints "hello"

You can also use slicing to access a range of elements:

```
my_tuple = (1, 2, 3, "hello", True)
print(my tuple[1:4])  # prints (2, 3, "hello")
```

Tuples are often used to return multiple values from a function:

```
def get_name_and_age():
    name = "Alice"
    age = 30
    return name, age
name, age = get_name_and_age()
print(name)  # prints "Alice"
print(age)  # prints 30
```



In this example, get_name_and_age returns a tuple containing the name and age. The function call name, age = get_name_and_age() assigns the first value in the tuple to name and the second value to age.

In SAS, the equivalent of a tuple is a SAS data set. However, unlike tuples in Python, SAS data sets are mutable and can be modified after they are created. Additionally, SAS data sets are usually much larger than tuples, and are used for storing and analyzing large amounts of data.

Overall, tuples in Python are a useful data structure for representing fixed collections of related data, and can be used in a variety of ways in Python programming.

In Python, a tuple is an immutable sequence of objects. It is similar to a list, but once a tuple is created, its contents cannot be modified. Tuples are often used to represent fixed collections of related data, such as a point in 2D space or a date (year, month, day).

Tuples are defined using parentheses () and the objects inside the tuple are separated by commas. Here's an example:

my tuple = (1, 2, 3, "hello", True)

In this example, my_tuple is a tuple containing five objects: the integers 1, 2, and 3, the string "hello", and the boolean value True.

Tuples can also be created without using parentheses, as long as there is more than one object in the tuple:

my tuple = 1, 2, 3

In this case, my_tuple is still a tuple containing the integers 1, 2, and 3.

Accessing elements in a tuple is similar to accessing elements in a list. You can use square brackets and an index to access a specific element:

```
my_tuple = (1, 2, 3, "hello", True)
print(my_tuple[0]) # prints 1
print(my_tuple[3]) # prints "hello"
```

You can also use slicing to access a range of elements:

```
my_tuple = (1, 2, 3, "hello", True)
print(my_tuple[1:4])  # prints (2, 3, "hello")
```

Tuples are often used to return multiple values from a function:

```
def get_name_and_age():
    name = "Alice"
```



```
age = 30
return name, age
name, age = get_name_and_age()
print(name)  # prints "Alice"
print(age)  # prints 30
```

In this example, get_name_and_age returns a tuple containing the name and age. The function call name, age = get_name_and_age() assigns the first value in the tuple to name and the second value to age.

In SAS, the equivalent of a tuple is a SAS data set. However, unlike tuples in Python, SAS data sets are mutable and can be modified after they are created. Additionally, SAS data sets are usually much larger than tuples, and are used for storing and analyzing large amounts of data.

Overall, tuples in Python are a useful data structure for representing fixed collections of related data, and can be used in a variety of ways in Python programming.

here's an example of some code that uses tuples in Python:

```
def get_point():
    x = 1
    y = 2
    return x, y
point = get_point()
print(point)    # prints (1, 2)
print(point[0])    # prints 1
print(point[1])    # prints 2
```

In this example, the get_point function creates two variables x and y and returns them as a tuple. The point variable is assigned to the tuple returned by get_point, which is (1, 2) in this case.

The print statements demonstrate how to access individual elements of the tuple using indexing. point[0] retrieves the first element of the tuple (1), and point[1] retrieves the second element of the tuple (2).

Here's another example that uses tuples to return multiple values from a function:

```
def get_name_and_age():
    name = "Alice"
    age = 30
    return name, age
```



```
name, age = get_name_and_age()
print(name) # prints "Alice"
print(age) # prints 30
```

In this example, the get_name_and_age function creates two variables name and age and returns them as a tuple. The name and age variables are assigned to the tuple returned by get_name_and_age, using tuple unpacking.

The print statements demonstrate how to access the individual values returned by the function, which are assigned to the name and age variables respectively.

Here's another example of using tuples in Python:

```
def get_numbers():
    return (1, 2, 3, 4, 5)
numbers = get_numbers()
# iterate over the tuple using a for loop
for number in numbers:
    print(number)
# get the length of the tuple
print(len(numbers))
# check if a value is in the tuple
print(3 in numbers)
# concatenate two tuples
more_numbers = (6, 7, 8)
all_numbers = numbers + more_numbers
print(all_numbers)
```

In this example, the get_numbers function returns a tuple containing five integer values. The numbers variable is assigned to this tuple, and various operations are performed on it.

The for loop demonstrates how to iterate over the elements of a tuple using a loop. The len function returns the number of elements in the tuple. The in operator checks whether a given value is in the tuple.

Finally, the + operator can be used to concatenate two tuples, creating a new tuple that contains all the elements of both tuples.

Tuples are often used to group related data together into a single, immutable object. For example, you might use a tuple to represent a point in two-dimensional space, with the x-coordinate and y-coordinate stored as the two elements of the tuple:

point = (3, 4)

This tuple can be passed to a function that expects a point as input, or it can be used as a key in a dictionary to store additional information about the point.

Tuples can also be used to return multiple values from a function, as shown in the previous examples. This can be a convenient way to package related data together, instead of returning multiple separate variables.

Here's another example of using tuples in Python to sort a list of tuples based on one of the values in each tuple:

```
students = [("Alice", 22), ("Bob", 20), ("Charlie",
25)]
# sort the list of tuples by age (the second element of
each tuple)
students_sorted_by_age = sorted(students, key=lambda x:
x[1])
# print the sorted list of tuples
print(students sorted by age)
```

In this example, we have a list of tuples, where each tuple represents a student's name and age. We want to sort the list of tuples based on the age of each student, so we use the sorted function with a key argument that tells it to sort the list based on the second element of each tuple (x[1]).

The lambda keyword is used to define an anonymous function that takes a single argument x and returns x[1]. This function is used as the key argument to the sorted function, which uses it to determine the order of the tuples in the sorted list.

The resulting students_sorted_by_age variable is a new list of tuples, sorted by the age of each student.

Tuples can also be used as keys in dictionaries, as shown in the following example:

```
colors = {("red", 255, 0, 0): "red", ("green", 0, 255,
0): "green", ("blue", 0, 0, 255): "blue"}
# get the color name for the (red, 255, 0, 0) tuple
color_name = colors[("red", 255, 0, 0)]
```



print the color name print(color_name)

In this example, we have a dictionary where the keys are tuples representing RGB color values, and the values are color names. We can use a tuple as a key in a dictionary to associate additional information with the tuple.

The color_name variable is assigned the value of the dictionary at the key ("red", 255, 0, 0), which is the string "red". We can use the tuple as a key to retrieve the associated value from the dictionary.

Tuples can also be used to return multiple values from a function, as shown in the previous examples. This can be a convenient way to package related data together, instead of returning multiple separate variables

Here's another example of using tuples in Python to swap the values of two variables:

```
a = 5
b = 10
# swap the values of a and b using a tuple
a, b = b, a
# print the new values of a and b
print(a) # 10
print(b) # 5
```

In this example, we have two variables a and b with initial values of 5 and 10, respectively. We want to swap the values of these variables, so we use a tuple to accomplish this in a single line of code.

The syntax a, b = b, a creates a tuple with the values of b and a, respectively, and then immediately unpacks this tuple into the variables a and b. This effectively swaps the values of the two variables in a single step.

The resulting values of a and b are printed to the console, which shows that the values have been successfully swapped.

Tuples can also be used to unpack values from functions that return multiple values, as shown in the following example:

```
def get_student_info():
    name = "Alice"
    age = 22
    major = "Computer Science"
    return name, age, major
```



```
# unpack the values returned by the function into
separate variables
student_name, student_age, student_major =
get_student_info()
# print the values of the variables
print(student_name)  # Alice
print(student_age)  # 22
print(student_major)  # Computer Science
```

In this example, we have a function get_student_info that returns three values: the student's name, age, and major. We use a tuple to package these values together and return them from the function.

When we call the function, we use the syntax student_name, student_age, student_major = get_student_info() to unpack the tuple into three separate variables. This assigns the value of the first element of the tuple to student_name, the second element to student_age, and the third element to student_major.

Here's another example of using tuples in Python to create a named tuple:

```
from collections import namedtuple
# define a named tuple type
Person = namedtuple("Person", ["name", "age",
"gender"])
# create a new Person named tuple
person1 = Person(name="Alice", age=22, gender="female")
# print the values of the person1 named tuple
print(person1.name)  # Alice
print(person1.age)  # 22
print(person1.gender)  # female
```

In this example, we use the namedtuple function from the collections module to create a new named tuple type called Person. The first argument to the namedtuple function is the name of the named tuple type, and the second argument is a list of field names. Unlike dictionaries, named tuples are immutable, which can make them safer to use in some contexts where we don't want the values to be accidentally changed.

Here's another example of using tuples in Python to iterate over multiple sequences simultaneously:

```
names = ["Alice", "Bob", "Charlie"]
```



```
ages = [22, 30, 45]
genders = ["female", "male", "male"]
# iterate over the sequences simultaneously using zip()
for name, age, gender in zip(names, ages, genders):
    print(name, age, gender)
```

In this example, we have three sequences: names, ages, and genders, each with a different set of values. We want to iterate over all three sequences simultaneously and print the values for each element.

We use the zip() function to combine the three sequences into a single sequence of tuples, where each tuple contains one value from each sequence. We then use a for loop to iterate over this sequence of tuples.

Inside the loop, we use tuple unpacking to assign the values from each tuple to separate variables. This allows us to print the values for each element in a structured way, with the name, age, and gender all printed on the same line.

We then create a new Person named tuple by specifying the values for each field using keyword arguments. This creates a new named tuple with the specified values for each field.

Dictionaries

Dictionaries are a built-in data structure in Python that allow you to store and access data in a way that is both flexible and efficient. In essence, a dictionary is a collection of key-value pairs, where each key is unique and associated with a corresponding value.

In Python, dictionaries are created using curly braces {} and separating each key-value pair with a colon :. For example, the following code creates a simple dictionary:

```
my dict = {"apple": 3, "banana": 5, "orange": 2}
```

In this dictionary, the keys are "apple", "banana", and "orange", and the corresponding values are 3, 5, and 2.

You can access the value associated with a particular key using the square bracket notation. For example, to get the value associated with the key "banana", you would use the following code:

print(my_dict["banana"]) # Output: 5



You can also modify the value associated with a key using the same square bracket notation. For example, the following code increases the value associated with the key "apple" by 1:

```
my_dict["apple"] += 1
print(my_dict)  # Output: {"apple": 4, "banana": 5,
"orange": 2}
```

In addition to storing simple values like numbers and strings, dictionaries can also store more complex data structures like lists and other dictionaries. For example, the following code creates a dictionary where the values are lists:

```
my_dict = {"apple": [3, 4, 5], "banana": [5, 6, 7],
"orange": [2, 3, 4]}
```

You can access individual elements of these lists using the same square bracket notation as before. For example, to get the second element of the list associated with the key "apple", you would use the following code:

```
print(my dict["apple"][1]) # Output: 4
```

You can also add new key-value pairs to a dictionary using the square bracket notation. For example, the following code adds a new key-value pair to the dictionary:

```
my_dict["grape"] = [1, 2, 3]
print(my_dict)  # Output: {"apple": [3, 4, 5],
"banana": [5, 6, 7], "orange": [2, 3, 4], "grape": [1,
2, 3]}
```

Dictionaries are a powerful and flexible data structure in Python that allow you to store and access data in a way that is both efficient and easy to use. They are especially useful when you need to associate values with unique keys, or when you need to store more complex data structures like lists and other dictionaries. here's some more information about dictionaries in Python:

Dictionary keys must be immutable: In Python, dictionary keys must be immutable, which means they cannot be changed once they are created. This is because the dictionary uses the key's hash value to look up the corresponding value, and if the key were mutable, its hash value could change, leading to incorrect lookups.

Dictionary values can be mutable: Unlike dictionary keys, dictionary values can be mutable, which means they can be changed after they are created. This can be useful when you need to modify a value associated with a particular key.

Dictionary methods: Python provides several built-in methods for working with dictionaries, including keys(), values(), and items(). These methods return views of the dictionary's keys,

in stal

values, and key-value pairs, respectively. You can also use the in keyword to check whether a key is in the dictionary.

Dictionary comprehension: Python supports dictionary comprehension, which allows you to create dictionaries using a concise syntax. For example, the following code creates a dictionary of squares:

squares = {x: x**2 for x in range(10)}
print(squares) # Output: {0: 0, 1: 1, 2: 4, 3: 9, 4:
16, 5: 25, 6: 36, 7: 49, 8: 64, 9: 81}

Dictionary sorting: Dictionaries in Python are inherently unordered, which means that the order in which key-value pairs are stored in the dictionary is not fixed. However, you can sort a dictionary by its keys using the sorted() function, which returns a list of key-value pairs sorted by key.

Dictionary performance: Dictionaries in Python are implemented as hash tables, which means that accessing elements in a dictionary is usually very fast (O(1) time complexity on average). However, the performance of a dictionary can degrade if it becomes too large or if the hash function used to compute keys is not well-distributed.

Overall, dictionaries are a versatile and powerful data structure in Python that can be used to store and manipulate a wide range of data types. They are commonly used in Python programming, especially when dealing with large datasets or complex data structures.

Here's an example code that demonstrates various operations on dictionaries in Python:

```
# Creating a dictionary
my_dict = {"apple": 3, "banana": 5, "orange": 2}
# Accessing a value
print(my_dict["banana"]) # Output: 5
# Modifying a value
my_dict["apple"] += 1
print(my_dict) # Output: {"apple": 4, "banana": 5,
"orange": 2}
# Adding a new key-value pair
my_dict["grape"] = 6
print(my_dict) # Output: {"apple": 4, "banana": 5,
"orange": 2, "grape": 6}
# Creating a dictionary with lists as values
```



```
my dict2 = {"apple": [3, 4, 5], "banana": [5, 6, 7],
"orange": [2, 3, 4]}
# Accessing an element of a list in the dictionary
print(my dict2["apple"][1]) # Output: 4
# Using dictionary methods
print(my dict.keys())
                        # Output: dict keys(['apple',
'banana', 'orange', 'grape'])
print(my dict.values()) # Output: dict values([4, 5,
2, 6])
print(my dict.items()) # Output:
dict items([('apple', 4), ('banana', 5), ('orange', 2),
('grape', 6)])
print("banana" in my dict) # Output: True
# Using dictionary comprehension
squares = {x: x**2 for x in range(10)}
print(squares) # Output: {0: 0, 1: 1, 2: 4, 3: 9, 4:
16, 5: 25, 6: 36, 7: 49, 8: 64, 9: 81}
# Sorting a dictionary by key
sorted dict = {k: v for k, v in
sorted(my dict.items()) }
print(sorted dict) # Output: {'apple': 4, 'banana':
5, 'grape': 6, 'orange': 2}
```

This code creates two dictionaries, my_dict and my_dict2, and demonstrates various operations on them, including accessing values, modifying values, adding new key-value pairs, accessing elements of lists in the dictionary, using dictionary methods, using dictionary comprehension, and sorting the dictionary by key.

Here are a few more examples of using dictionaries in Python:

```
# Example 1: Counting the frequency of letters in a
string
my_string = "hello, world!"
freq = {}
for letter in my_string:
    if letter in freq:
        freq[letter] += 1
    else:
        freq[letter] = 1
```



```
print(freq) # Output: {'h': 1, 'e': 1, 'l': 3, 'o':
2, ',': 1, ' ': 1, 'w': 1, 'r': 1, 'd': 1, '!': 1}
# Example 2: Creating a dictionary from two lists
keys = ["apple", "banana", "orange"]
values = [3, 5, 2]
my dict = {k: v for k, v in zip(keys, values)}
print(my dict) # Output: {'apple': 3, 'banana': 5,
'orange': 2}
# Example 3: Removing a key from a dictionary
del my dict["banana"]
print(my dict)
                 # Output: {'apple': 3, 'orange': 2}
# Example 4: Using the get() method to handle missing
keys
print(my dict.get("banana", "not found"))  # Output:
not found
```

In Example 1, we create a dictionary to count the frequency of each letter in a string. We loop over each character in the string, and if the character is already in the dictionary, we increment its value by 1. Otherwise, we add the character to the dictionary with a value of 1.

In Example 2, we create a dictionary from two lists, using the zip() function to combine corresponding elements from each list.

In Example 3, we remove a key-value pair from the dictionary using the del keyword.

In Example 4, we use the get() method to retrieve a value from the dictionary for a given key. If the key is not in the dictionary, the method returns a default value (in this case, the string "not found").

Here are a few more examples:

```
# Example 5: Merging two dictionaries
dict1 = {'a': 1, 'b': 2}
dict2 = {'b': 3, 'c': 4}
merged_dict = {**dict1, **dict2}
print(merged_dict) # Output: {'a': 1, 'b': 3, 'c': 4}
# Example 6: Counting the frequency of words in a list
of strings
```



```
sentences = ["This is a sentence.", "This is another
sentence."]
word freq = {}
for sentence in sentences:
    words = sentence.split()
    for word in words:
        if word in word freq:
            word freq[word] += 1
        else:
            word freq[word] = 1
print(word freq) # Output: {'This': 2, 'is': 2, 'a':
1, 'sentence.': 2, 'another': 1}
# Example 7: Using the setdefault() method to set
default values
my dict = \{ 'a': 1, 'b': 2 \}
my dict.setdefault('c', 3)
print(my dict) # Output: {'a': 1, 'b': 2, 'c': 3}
# Example 8: Updating a dictionary with another
dictionary
my dict = \{ 'a': 1, 'b': 2 \}
update dict = { 'b': 3, 'c': 4}
my dict.update(update dict)
print(my dict) # Output: {'a': 1, 'b': 3, 'c': 4}
```

In Example 5, we merge two dictionaries using the ** operator. When two dictionaries have the same key, the value from the second dictionary overwrites the value from the first dictionary.

In Example 6, we count the frequency of words in a list of strings. We loop over each sentence in the list, split it into words, and increment the frequency count for each word in the word_freq dictionary.

In Example 7, we use the setdefault() method to set a default value for a key that doesn't exist in the dictionary. If the key already exists, the method doesn't change the value.

In Example 8, we update a dictionary with the key-value pairs from another dictionary using the update() method. If a key already exists in the original dictionary, the value from the second dictionary overwrites it.

Here are a few more examples:

```
# Example 9: Sorting a dictionary by value
my dict = {'a': 3, 'b': 1, 'c': 2}
```



```
sorted dict = {k: v for k, v in sorted(my dict.items(),
key=lambda item: item[1])}
print(sorted dict) # Output: {'b': 1, 'c': 2, 'a': 3}
# Example 10: Creating a dictionary from a list of
tuples
my list = [('a', 1), ('b', 2), ('c', 3)]
my dict = dict(my list)
print(my dict) # Output: {'a': 1, 'b': 2, 'c': 3}
# Example 11: Finding the key with the maximum value in
a dictionary
my dict = {'a': 3, 'b': 1, 'c': 2}
max key = max(my dict, key=my dict.get)
print(max key) # Output: 'a'
# Example 12: Creating a nested dictionary
my dict = { 'a': { 'x': 1, 'y': 2}, 'b': { 'x': 3, 'y':
4}}
print(my dict['a']['x']) # Output: 1
```

In Example 9, we sort a dictionary by its values in ascending order using the sorted() function and a lambda function as the key. We then create a new dictionary with the same key-value pairs in the sorted order.

In Example 10, we create a dictionary from a list of tuples using the dict() constructor.

In Example 11, we find the key in a dictionary with the maximum value using the max() function and the get() method as the key function.

In Example 12, we create a nested dictionary, where each value is another dictionary. We can access values in the nested dictionary using multiple keys.

```
# Example 13: Flattening a dictionary
my_dict = {'a': {'x': 1, 'y': 2}, 'b': {'x': 3, 'y':
4}}
flat_dict = {}
for key, value in my_dict.items():
    for sub_key, sub_value in value.items():
        flat_dict[sub_key] = sub_value
print(flat_dict)  # Output: {'x': 3, 'y': 4}
# Example 14: Filtering a dictionary using a
comprehension
```



```
my_dict = {'a': 1, 'b': 2, 'c': 3}
filtered_dict = {k: v for k, v in my_dict.items() if v
% 2 == 0}
print(filtered_dict) # Output: {'b': 2}
# Example 15: Removing a key from a dictionary
my_dict = {'a': 1, 'b': 2, 'c': 3}
del my_dict['b']
print(my_dict) # Output: {'a': 1, 'c': 3}
# Example 16: Checking if a key exists in a dictionary
my_dict = {'a': 1, 'b': 2, 'c': 3}
if 'b' in my_dict:
    print("Key 'b' exists in the dictionary.")
else:
    print("Key 'b' does not exist in the dictionary.")
# Output: Key 'b' exists in the dictionary."
```

In Example 13, we flatten a nested dictionary by iterating over each key-value pair, then iterating over each sub-key and sub-value pair in the nested dictionary. We add each sub-key and sub-value pair to the new dictionary.

In Example 14, we filter a dictionary using a dictionary comprehension. We create a new dictionary with only the key-value pairs that satisfy a certain condition.

In Example 15, we remove a key from a dictionary using the del statement.

In Example 16, we check if a key exists in a dictionary using the in operator. If the key exists, we print a message indicating that it does. Otherwise, we print a message indicating that it does not.

```
# Example 17: Merging two dictionaries
dict1 = {'a': 1, 'b': 2}
dict2 = {'c': 3, 'd': 4}
merged_dict = {**dict1, **dict2}
print(merged_dict) # Output: {'a': 1, 'b': 2, 'c': 3,
'd': 4}
# Example 18: Updating a dictionary
my_dict = {'a': 1, 'b': 2}
my_dict.update({'a': 3, 'c': 4})
print(my_dict) # Output: {'a': 3, 'b': 2, 'c': 4}
# Example 19: Getting a list of keys from a dictionary
my_dict = {'a': 1, 'b': 2, 'c': 3}
```



```
key_list = list(my_dict.keys())
print(key_list) # Output: ['a', 'b', 'c']
# Example 20: Getting a list of values from a
dictionary
my_dict = {'a': 1, 'b': 2, 'c': 3}
value_list = list(my_dict.values())
print(value list) # Output: [1, 2, 3]
```

In Example 17, we merge two dictionaries into a single dictionary using the unpacking operator (**). The key-value pairs from dict1 and dict2 are added to the new dictionary.

In Example 18, we update a dictionary with new key-value pairs using the update() method. If a key already exists in the dictionary, its value is updated with the new value.

In Example 19, we get a list of keys from a dictionary using the keys() method. We then convert the resulting view object to a list.

In Example 20, we get a list of values from a dictionary using the values() method. We then convert the resulting view object to a list.

```
# Example 21: Getting the length of a dictionary
my dict = \{ 'a': 1, 'b': 2, 'c': 3 \}
print(len(my dict)) # Output: 3
# Example 22: Getting a default value for a missing key
my dict = \{ 'a': 1, 'b': 2, 'c': 3 \}
default value = my dict.get('d', 0)
print(default value) # Output: 0
# Example 23: Sorting a dictionary by key
my dict = { 'b': 2, 'c': 3, 'a': 1}
sorted dict = dict(sorted(my dict.items()))
print(sorted dict) # Output: {'a': 1, 'b': 2, 'c': 3}
# Example 24: Sorting a dictionary by value
my dict = \{ 'b': 2, 'c': 3, 'a': 1 \}
sorted dict = dict(sorted(my dict.items(), key=lambda
item: item[1]))
print(sorted dict) # Output: {'a': 1, 'b': 2, 'c': 3}
```

In Example 21, we get the length of a dictionary using the len() function.

In Example 22, we get a default value for a missing key in a dictionary using the get() method. If the key is not present in the dictionary, it returns the default value provided as the second argument.

In Example 23, we sort a dictionary by key using the sorted() function. We convert the resulting list of key-value pairs back to a dictionary using the dict() constructor.

In Example 24, we sort a dictionary by value using the sorted() function and a lambda function as the key argument. The lambda function specifies that we want to sort the dictionary by the second element in each key-value pair (i.e., the value). We convert the resulting list of key-value pairs back to a dictionary using the dict() constructor.

```
# Example 25: Removing a key-value pair from a
dictionary
my dict = \{ 'a': 1, 'b': 2, 'c': 3 \}
del my dict['b']
print(my dict) # Output: {'a': 1, 'c': 3}
# Example 26: Removing all key-value pairs from a
dictionarv
my dict = \{ 'a': 1, 'b': 2, 'c': 3 \}
my dict.clear()
print(my dict) # Output: {}
# Example 27: Creating a dictionary with a
comprehension
my dict = {x: x**2 for x in range(5)}
print(my dict) # Output: {0: 0, 1: 1, 2: 4, 3: 9, 4:
16}
# Example 28: Filtering a dictionary with a
comprehension
my dict = { 'a': 1, 'b': 2, 'c': 3, 'd': 4 }
filtered dict = {k: v for k, v in my dict.items() if v
8 2 == 0}
print(filtered dict) # Output: {'b': 2, 'd': 4}
In Example 25, we remove a key-value pair from a
dictionary using the del statement.
```

In Example 26, we remove all key-value pairs from a dictionary using the clear() method.

In Example 27, we create a dictionary using a dictionary comprehension. This is a concise way to create a new dictionary by iterating over some iterable (in this case, range(5)) and defining the key-value pairs in a single line of code.



In Example 28, we filter a dictionary using a dictionary comprehension. We iterate over the keyvalue pairs in the dictionary and include only those where the value is even. The resulting dictionary only contains the key-value pairs where the value is even.

```
# Example 29: Merging two dictionaries
dict1 = \{ 'a': 1, 'b': 2 \}
dict2 = \{ 'c': 3, 'd': 4 \}
merged dict = {**dict1, **dict2}
print(merged dict) # Output: {'a': 1, 'b': 2, 'c': 3,
'd': 4}
# Example 30: Updating a dictionary with another
dictionary
my dict = \{ 'a': 1, 'b': 2, 'c': 3 \}
my dict.update({'b': 4, 'd': 5})
print(my dict) # Output: {'a': 1, 'b': 4, 'c': 3, 'd':
5}
# Example 31: Creating a dictionary from two lists
keys = ['a', 'b', 'c']
values = [1, 2, 3]
my dict = dict(zip(keys, values))
print(my dict) # Output: {'a': 1, 'b': 2, 'c': 3}
# Example 32: Creating a dictionary from a list of
tuples
my list = [('a', 1), ('b', 2), ('c', 3)]
my_dict = dict(my_list)
print(my dict) # Output: {'a': 1, 'b': 2, 'c': 3}
```

In Example 29, we merge two dictionaries into a single dictionary using the unpacking operator (**). This creates a new dictionary that contains all the key-value pairs from both dictionaries.

In Example 30, we update a dictionary with another dictionary using the update() method. This adds any new key-value pairs from the second dictionary to the first dictionary, and updates the values for any keys that are already present.

In Example 31, we create a dictionary from two lists using the zip() function and the dict() constructor. The zip() function combines the elements from each list into tuples, and the dict() constructor converts the list of tuples to a dictionary.

In Example 32, we create a dictionary from a list of tuples using the dict() constructor. The list of tuples contains the key-value pairs for the dictionary.



Sets

Sets are a fundamental data structure in Python that allow you to store and manipulate collections of unique elements. They are often used in Python for tasks such as removing duplicates from a list, performing set operations like union and intersection, and checking membership.

In SAS, sets are often represented using the MERGE statement, which combines two or more data sets based on the values of one or more common variables. In Python, sets are a built-in data type that can be created using the set() function or by using curly braces {} around a comma-separated sequence of elements.

For example, to create a set of unique numbers in Python, you can use the set() function like this:

```
>>> numbers = [1, 2, 3, 4, 5, 5, 4, 3, 2, 1]
>>> unique_numbers = set(numbers)
>>> unique_numbers
{1, 2, 3, 4, 5}
```

In this example, the list of numbers contains duplicates, but the resulting set only contains the unique elements.

You can perform a variety of operations on sets in Python, such as adding and removing elements, performing set arithmetic operations like union, intersection, and difference, and checking membership.

```
>>> set1 = {1, 2, 3}
>>> set2 = {3, 4, 5}
>>> union = set1.union(set2)
>>> intersection = set1.intersection(set2)
>>> difference = set1.difference(set2)
>>> print(union)
{1, 2, 3, 4, 5}
>>> print(intersection)
{3}
>>> print(difference)
{1, 2}
```

In this example, we create two sets (set1 and set2), and then use the union(), intersection(), and difference() methods to perform set operations on those sets. The union() method returns a new set that contains all elements from both sets, the intersection() method returns a new set that contains only the elements that are common to both sets, and the difference() method returns a new set that contains only the elements that are in set1 but not in set2.



Sets are also mutable, which means you can add and remove elements from them. For example:

```
>>> set1 = {1, 2, 3}
>>> set1.add(4)
>>> print(set1)
{1, 2, 3, 4}
>>> set1.remove(2)
>>> print(set1)
{1, 3, 4}
```

In this example, we add the element 4 to set1 using the add() method, and then remove the element 2 using the remove() method.

In conclusion, sets are a powerful and flexible data structure in Python that allow you to perform a variety of operations on collections of unique elements. If you are a SAS user transitioning to Python, understanding sets and their operations can be a valuable tool in your programming arsenal.

Here are some additional details about sets in Python:

- 1. Sets are unordered: Unlike lists and tuples, sets do not maintain any order of their elements. This means that you cannot access the elements of a set using an index.
- 2. Sets cannot contain duplicates: One of the defining characteristics of a set is that it only contains unique elements. If you try to add an element to a set that already exists in the set, it will simply be ignored.
- 3. Sets can be created using curly braces or the set() function: You can create a set in Python by enclosing a comma-separated sequence of elements within curly braces, like this: {1, 2, 3}. Alternatively, you can use the set() function to create a set from any iterable, like a list or tuple.
- 4. Sets support set arithmetic operations: Python provides a number of built-in methods for performing set operations like union, intersection, and difference. These methods are used to combine or compare two or more sets.
- 5. Sets can be used to remove duplicates from a list: If you have a list with duplicate elements, you can easily remove them by converting the list to a set using the set() function, and then converting it back to a list using the list() function.
- 6. Sets are mutable: You can add and remove elements from a set using the add(), remove(), and discard() methods. You can also use the update() method to add elements from another set, list, or tuple to an existing set.
- 7. Frozensets are immutable: If you need a set that cannot be modified, you can use a frozenset instead. A frozenset is simply a set that cannot be modified once it has been created.

Overall, sets are a versatile and powerful data structure in Python that can be used for a variety of tasks. Whether you need to remove duplicates from a list, perform set operations, or store a collection of unique elements, sets can be a useful tool in your programming toolbox.



Here's an example of some code that demonstrates the use of sets in Python:

```
# Creating a set using curly braces
my set = \{1, 2, 3, 4, 5\}
# Creating a set using the set() function
my other set = set([3, 4, 5, 6, 7])
# Printing the sets
print("my set: ", my set)
print("my other set: ", my other set)
# Adding an element to a set
my set.add(6)
print("After adding 6 to my set: ", my set)
# Removing an element from a set
my other set.remove(7)
print("After removing 7 from my other set: ",
my other set)
# Performing set operations
union set = my set.union(my other set)
print("Union of my set and my other set: ", union set)
intersection set = my set.intersection(my other set)
print("Intersection of my set and my other set: ",
intersection set)
difference set = my set.difference(my other set)
print("Difference of my set and my other set: ",
difference set)
# Converting a list to a set to remove duplicates
my list = [1, 2, 3, 4, 5, 5, 4, 3, 2, 1]
unique elements = set(my list)
print("Unique elements in my list: ", unique elements)
# Creating a frozenset
frozen set = frozenset([1, 2, 3, 4, 5])
print("frozen set: ", frozen set)
```



```
# Trying to modify a frozenset (will result in an
error)
# frozen_set.add(6)
# frozen_set.remove(1)
```

This code demonstrates several common operations with sets in Python. First, we create two sets using both the curly brace notation and the set() function. We then add an element to one set and remove an element from the other set. We also perform set operations like union, intersection, and difference on the two sets.

Next, we demonstrate how sets can be used to remove duplicates from a list. We convert a list with duplicate elements to a set using the set() function, which automatically removes the duplicates, and then print out the resulting unique elements.

```
# Example of adding and removing elements from a set
my set = \{1, 2, 3, 4, 5\}
print("Initial set: ", my set)
# Adding elements to the set
my set.add(6)
my set.update([7, 8])
print("After adding elements: ", my_set)
# Removing elements from the set
my set.remove(3)
my set.discard(9) # this does not result in an error
print("After removing elements: ", my set)
# Example of using sets in a loop
set1 = \{1, 2, 3\}
set2 = \{2, 3, 4\}
for item in set1:
    if item in set2:
        print("Intersection found: ", item)
    else:
        print("No intersection: ", item)
# Example of set comprehension
set3 = {x for x in range(10) if x % 2 == 0}
print("Set with even numbers: ", set3)
```

This code continues to demonstrate different aspects of sets in Python. First, we add and remove elements from a set using the add(), update(), remove(), and discard() methods. We also show



that the discard() method does not result in an error if the element to be removed is not in the set, while the remove() method does.

Next, we show an example of using sets in a loop to find the intersection of two sets. We loop over the elements of one set and check if they are in the other set using the 'in' keyword.

```
# Example of using sets to remove duplicates from a
list
my list = [1, 2, 3, 2, 4, 3, 5, 6, 4, 7]
my set = set(my list)
print("List with duplicates: ", my list)
print("Set without duplicates: ", my_set)
# Example of checking if a set is a subset or superset
of another set
set1 = \{1, 2, 3, 4, 5\}
set2 = \{2, 3, 4\}
set3 = \{6, 7, 8\}
print("set2 is a subset of set1: ",
set2.issubset(set1))
print("set1 is a superset of set2: ",
set1.issuperset(set2))
print("set3 is a subset of set1: ",
set3.issubset(set1))
print("set1 is a superset of set3: ",
set1.issuperset(set3))
# Example of finding the symmetric difference between
two sets
set4 = \{1, 2, 3, 4, 5\}
set5 = \{4, 5, 6, 7, 8\}
symmetric diff = set4.symmetric difference(set5)
print("Symmetric difference of set4 and set5: ",
symmetric diff)
# Example of using sets to remove elements from a list
my list = [1, 2, 3, 4, 5, 6, 7]
elements to remove = \{2, 4, 6\}
my list = [x for x in my list if x not in
elements to remove]
print("List after removing elements: ", my list)
```



```
# Example of using sets to find unique elements between
two lists
list1 = [1, 2, 3, 4, 5]
list2 = [4, 5, 6, 7, 8]
unique_elements =
set(list1).symmetric_difference(set(list2))
print("Unique elements between list1 and list2: ",
unique elements)
```

In this code, we first demonstrate how to use sets to remove duplicates from a list. We convert a list with duplicates to a set using the set() function, which automatically removes the duplicates.

Next, we show examples of using the issubset() and issuperset() methods to check if one set is a subset or superset of another set. We also show an example of finding the symmetric difference between two sets, which is the set of elements that are in one set or the other, but not both.

We also show an example of using sets to remove elements from a list. We create a set of elements to remove and then use a list comprehension to create a new list that does not contain those elements.

```
# Example of using sets to perform mathematical
operations
set1 = \{1, 2, 3, 4, 5\}
set2 = \{4, 5, 6, 7, 8\}
# Union of two sets
union = set1.union(set2)
print("Union of set1 and set2: ", union)
# Intersection of two sets
intersection = set1.intersection(set2)
print("Intersection of set1 and set2: ", intersection)
# Difference between two sets
difference = set1.difference(set2)
print("Difference between set1 and set2: ", difference)
# Example of using sets to find common elements in
multiple lists
list1 = [1, 2, 3, 4, 5]
list2 = [4, 5, 6, 7, 8]
list3 = [3, 4, 5, 6, 7]
```



```
common_elements = set(list1).intersection(set(list2),
set(list3))
print("Common elements in list1, list2, and list3: ",
common_elements)
# Example of using sets to check if two lists have any
common elements
list4 = [1, 2, 3]
list5 = [4, 5, 6]
list6 = [2, 5, 7]
has_common_elements =
bool(set(list4).intersection(set(list5), set(list6)))
print("Do list4, list5, and list6 have any common
elements?: ", has_common_elements)
```

In this code, we show examples of using sets to perform mathematical operations, such as union, intersection, and difference. These operations can be useful for combining or comparing sets of data.

We also show an example of using sets to find common elements in multiple lists. We create sets from each list and use the intersection() method to find the common elements between the sets.

```
# Example of using sets to remove duplicates from a
list
list7 = [1, 2, 3, 3, 4, 5, 5]
unique elements = set(list7)
print("Unique elements in list7: ", unique elements)
# Example of using sets to check if a list contains
only unique elements
has duplicates = len(list7) != len(set(list7))
print("Does list7 have any duplicates?: ",
has duplicates)
# Example of using sets to find the symmetric
difference between two sets
set3 = \{1, 2, 3\}
set4 = \{2, 3, 4\}
symmetric difference = set3.symmetric difference(set4)
print("Symmetric difference between set3 and set4: ",
symmetric difference)
```



```
# Example of using sets to check if one set is a subset
of another set
set5 = {1, 2, 3, 4}
set6 = {2, 3}
is_subset = set6.issubset(set5)
print("Is set6 a subset of set5?: ", is_subset)
# Example of using sets to check if one set is a
superset of another set
is_superset = set5.issuperset(set6)
print("Is set5 a superset of set6?: ", is superset)
```

In this code, we show an example of using sets to remove duplicates from a list. We create a set from the list, which automatically removes any duplicates, and then convert the set back to a list to get the unique elements.

We also show an example of using sets to check if a list contains only unique elements. We compare the length of the list to the length of the set, which will be different if there are duplicates in the list.

We then demonstrate how to find the symmetric difference between two sets using the symmetric_difference() method.

```
# Example of using sets to find the Cartesian product
of two sets
set7 = \{1, 2\}
set8 = {'a', 'b'}
cartesian product = { (i, j) for i in set7 for j in
set8}
print("Cartesian product of set7 and set8: ",
cartesian product)
# Example of using sets to create a set comprehension
set9 = {x for x in range(1, 11) if x \& 2 == 0}
print("set9: ", set9)
# Example of using sets to modify a set in place
set10 = \{1, 2, 3\}
set10.update([3, 4, 5])
print("Modified set10: ", set10)
# Example of using sets to remove an element from a set
set11 = \{1, 2, 3, 4, 5\}
set11.discard(4)
```



print("Set11 with 4 removed: ", set11)

In this code, we show an example of using sets to find the Cartesian product of two sets using set comprehension. The resulting set contains all possible ordered pairs of elements from the two sets.

We also demonstrate how to use sets to create a set comprehension. In this example, we create a set of even numbers from 1 to 10 using a set comprehension.

We then show an example of how to modify a set in place using the update() method. In this example, we add the elements [3, 4, 5] to the set10.

```
# Example of using sets to create a union of sets
set12 = \{1, 2, 3\}
set13 = \{3, 4, 5\}
union = set12.union(set13)
print("Union of set12 and set13: ", union)
# Example of using sets to create an intersection of
sets
intersection = set12.intersection(set13)
print("Intersection of set12 and set13: ",
intersection)
# Example of using sets to create a difference of sets
difference = set12.difference(set13)
print("Difference between set12 and set13: ",
difference)
# Example of using sets to create a symmetric
difference of sets
symmetric difference =
set12.symmetric difference(set13)
print("Symmetric difference between set12 and set13: ",
symmetric difference)
# Example of using sets to check for disjoint sets
set14 = \{1, 2, 3\}
set15 = \{4, 5, 6\}
is disjoint = set14.isdisjoint(set15)
print("Are set14 and set15 disjoint?: ", is disjoint)
```



In this code, we show examples of using sets to create a union, intersection, difference, and symmetric difference of sets. We use the union(), intersection(), difference(), and symmetric_difference() methods, respectively, to perform these set operations.

```
# Example of using sets to check for subset and
superset relationships
set16 = \{1, 2, 3\}
set17 = \{1, 2\}
set18 = \{1, 2, 3, 4, 5\}
is subset = set17.issubset(set16)
print("Is set17 a subset of set16?: ", is subset)
is superset = set18.issuperset(set16)
print("Is set18 a superset of set16?: ", is superset)
# Example of using sets to remove all elements except
those in a specified set
set19 = \{1, 2, 3, 4, 5\}
set19.intersection update({3, 4, 6})
print("Set19 after intersection update: ", set19)
# Example of using sets to remove all elements in a
specified set
set20 = \{1, 2, 3, 4, 5\}
set20.difference update({3, 4, 6})
print("Set20 after difference update: ", set20)
```

In this code, we show examples of using sets to check for subset and superset relationships between sets using the issubset() and issuperset() methods, respectively. We also show how to modify a set in place using the intersection_update() and difference_update() methods. The intersection_update() method removes all elements from the set that are not in a specified set, while the difference_update() method removes all elements that are in a specified set.

```
# Example of using sets to find the maximum and minimum
values in a set
set21 = {10, 5, 7, 3, 8}
max_value = max(set21)
min_value = min(set21)
print("Maximum value in set21: ", max_value)
print("Minimum value in set21: ", min_value)
# Example of using sets to convert a list to a set
list1 = [1, 2, 3, 4, 5]
set22 = set(list1)
print("Set22: ", set22)
```



```
# Example of using sets to convert a string to a set of
unique characters
string1 = "hello world"
set23 = set(string1)
print("Set23: ", set23)
```

In this code, we show examples of using sets to find the maximum and minimum values in a set using the max() and min() functions, respectively. We also demonstrate how to convert a list to a set using the set() function and how to convert a string to a set of unique characters using the set() function.

Indexing and slicing

Indexing and slicing are important concepts in programming, including Python. They refer to ways of accessing specific elements in a sequence or collection of data. In Python, indexing and slicing are commonly used when working with lists, tuples, and strings.

Indexing in Python starts at 0, meaning that the first element in a sequence is referred to as index 0. To access a specific element in a sequence, you can use square brackets with the index number. For example, if you have a list of numbers called my_list, you can access the first element like this:

my_list = [1, 2, 3, 4, 5]
print(my_list[0]) # Output: 1

Slicing in Python allows you to extract a portion of a sequence or collection. To slice a sequence, you can use square brackets with two index numbers separated by a colon. The first index number is the starting point of the slice, and the second index number is the end point of the slice (not inclusive). For example, to slice the first three elements from my_list, you can do:

my_list = [1, 2, 3, 4, 5]
print(my_list[0:3]) # Output: [1, 2, 3]

In addition to the two index numbers, you can also include a third number to specify the step size of the slice. For example, to slice every other element from my_list, you can do:

my_list = [1, 2, 3, 4, 5]
print(my_list[0:5:2]) # Output: [1, 3, 5]



When working with strings, you can also use indexing and slicing to access specific characters or substrings. For example, if you have a string called my_string, you can access the first character like this:

my_string = "hello world"
print(my string[0]) # Output: h

To slice a substring from my_string, you can do:

```
my_string = "hello world"
print(my string[0:5]) # Output: hello
```

Python also supports negative indexing and slicing, which allows you to access elements or substrings from the end of a sequence or string. For example, to access the last element in my_list, you can do:

my_list = [1, 2, 3, 4, 5]
print(my_list[-1]) # Output: 5

To slice the last three elements from my_list, you can do:

my_list = [1, 2, 3, 4, 5]
print(my_list[-3:]) # Output: [3, 4, 5]

In Python, indexing and slicing are not limited to just lists, tuples, and strings. They can also be used with other sequence types such as arrays and NumPy ndarrays. Here are some examples:

```
import numpy as np
# Create a NumPy ndarray
my_array = np.array([[1, 2, 3], [4, 5, 6], [7, 8, 9]])
# Indexing
print(my_array[0]) # Output: [1 2 3]
print(my_array[1, 2]) # Output: 6
# Slicing
print(my_array[1, 2]) # Output: [2 5 8]
print(my_array[1:3, :2]) # Output: [[4 5], [7 8]]
```



When indexing or slicing with multiple dimensions, you can use a comma-separated sequence of index or slice objects. For example, to slice a 2D NumPy ndarray, you can use two index or slice objects separated by a comma:

```
import numpy as np
# Create a 2D NumPy ndarray
my_array = np.array([[1, 2, 3], [4, 5, 6], [7, 8, 9]])
# Slicing
print(my_array[1:3, :2]) # Output: [[4 5], [7 8]]
You can also use ellipses (...) to represent all possible dimensions in a slice object. For example:
    import numpy as np
```

```
# Create a 3D NumPy ndarray
my_array = np.array([[[1, 2], [3, 4]], [[5, 6], [7,
8]]])
# Slicing
print(my array[..., 1]) # Output: [[2 4], [6 8]]
```

Finally, it's worth noting that indexing and slicing in Python can also be used with objects that define their own __getitem__() method. This method allows objects to behave like sequences, even if they are not implemented as such. For example:

```
class MyClass:
    def __init__(self):
        self.my_list = [1, 2, 3, 4, 5]
    def __getitem__(self, index):
        return self.my_list[index]
# Create an instance of MyClass
my_instance = MyClass()
# Indexing
print(my_instance[0]) # Output: 1
# Slicing
print(my_instance[1:4]) # Output: [2, 3, 4]
List comprehension
```



Here are some additional points to keep in mind when working with indexing and slicing in Python:

Indexing and slicing both start at index 0. That is, the first element in a sequence has an index of 0, not 1. For example:

```
my_list = [1, 2, 3, 4, 5]
print(my_list[0]) # Output: 1
```

Negative indices can be used to index or slice from the end of a sequence. That is, an index of -1 refers to the last element in a sequence, an index of -2 refers to the second-to-last element, and so on. For example:

```
my_list = [1, 2, 3, 4, 5]
print(my_list[-1]) # Output: 5
print(my_list[-2:]) # Output: [4, 5]
Sorting and filtering data
```

Here's a detailed explanation of indexing and slicing in Python:

Indexing:

Indexing refers to accessing individual elements within a sequence, such as a string, list or tuple. In Python, indexing starts from 0, meaning that the first element in a sequence has an index of 0, the second element has an index of 1, and so on. To access an element at a specific index in a sequence, you can use square brackets [] with the index number inside the brackets.

For example, let's say we have a list of numbers called "my_list":

my list = [10, 20, 30, 40, 50]

To access the second element in the list (which has an index of 1), we can use the following code:

```
second_element = my_list[1]
print(second element)  # Output: 20
```

Slicing:

Slicing refers to extracting a portion of a sequence, such as a substring from a string or a sub-list from a list. In Python, slicing is done using the colon (:) operator. The syntax for slicing is as follows:

```
sequence[start_index:stop_index:step]
where:
```



"start_index" is the index of the first element to include in the slice (inclusive). "stop_index" is the index of the last element to include in the slice (exclusive). "step" is the number of elements to skip between each element in the slice (default is 1).

Note that the "stop_index" is exclusive, meaning that the element at that index is not included in the slice.

For example, let's say we have a string called "my_string":

my string = "Hello, world!"

To extract the substring "world" from the string, we can use the following code:

my_slice = my_string[7:12]
print(my_slice) # Output: world

To extract every other character from the string, we can use the following code:

my_slice = my_string[::2]
print(my slice) # Output: Hlo ol!

Notice that we omitted the "start_index" and "stop_index" values, which defaults to the beginning and end of the sequence, respectively.

```
# create a string
my_string = "Hello, world!"
# create a list
my_list = [10, 20, 30, 40, 50]
# indexing example
print(my_string[0])  # Output: H
print(my_list[2])  # Output: 30
# slicing examples
print(my_string[0:5])  # Output: Hello
print(my_list[1:4])  # Output: [20, 30, 40]
print(my_string[::2])  # Output: Hlo,wrld
print(my_list[::2])  # Output: [10, 30, 50]
```

In the first part of the code, we create a string called "my_string" and a list called "my_list".



Next, we demonstrate indexing by accessing the first element of the string using "my_string[0]" and the third element of the list using "my_list[2]".

Finally, we demonstrate slicing by extracting the first 5 characters of the string using "my_string[0:5]", the second through fourth elements of the list using "my_list[1:4]", every other character in the string using "my_string[::2]", and every other element in the list using "my_list[::2]".

```
# create a string
my string = "Python is a great programming language"
# create a list
my list = [1, 2, 3, 4, 5, 6, 7, 8, 9, 10]
# slicing examples
print(my string[7:])
                                 # Output: is a great
programming language
print(my string[:6])
                                 # Output: Python
print(my string[::3])
                                 # Output: Ph sgetormi
aq
print(my string[7:23:2])
                                 # Output: isagr rormi
                                 # Output: egaugnal
print(my string[::-1])
gnimmargorp taerg a si nohtyP
print(my list[1::2])
                                 # Output: [2, 4, 6, 8,
101
print(my list[2::3])
                                 # Output: [3, 6, 9]
                                 # Output: [5, 6, 7, 8]
print(my list[4:8])
```

In this example, we create a string called "my_string" and a list called "my_list".

Next, we demonstrate slicing by extracting a substring of the string starting from the 7th index using "my_string[7:]", the first 6 characters of the string using "my_string[:6]", every third character of the string using "my_string[::3]", every other character of the string between the 7th and 23rd indices using "my_string[7:23:2]", and reversing the string using "my_string[::-1]".

For the list, we demonstrate slicing by extracting every other element starting from the second element using "my_list[1::2]", every third element starting from the third element using "my_list[2::3]", and a sub-list containing elements from the 5th to 8th indices (inclusive) using "my_list[4:8]".

```
# create a list
my_list = [10, 20, 30, 40, 50]
# modify list using indexing
my_list[2] = 35
```



```
print(my list)
                                 # Output: [10, 20, 35,
40, 50]
# modify list using slicing
my list[1:4] = [15, 25, 45]
                                 # Output: [10, 15, 25,
print(my list)
45, 50]
# delete elements from list using slicing
my list[2:4] = []
print(my list)
                                 # Output: [10, 15, 50]
# insert elements into list using slicing
my list[1:1] = [5, 6, 7]
print(my list)
                                 # Output: [10, 5, 6, 7,
15, 50]
# replace elements using slicing
my list[2:5] = [20, 30, 40]
                                 # Output: [10, 5, 20,
print(my list)
30, 40]
```

In this example, we create a list called "my_list".

First, we modify the element at index 2 by setting it to 35 using "my_list[2] = 35".

Next, we modify elements 1 through 4 using slicing by replacing them with the list [15, 25, 45] using "my_list[1:4] = [15, 25, 45]".

Then, we delete elements 2 through 4 using slicing by setting them to an empty list using "my_list[2:4] = []".

After that, we insert elements 5, 6, and 7 before element 1 using slicing by setting "my_list[1:1]" to the list [5, 6, 7] using "my_list[1:1] = [5, 6, 7]".

Finally, we replace elements 2 through 4 with the list [20, 30, 40] using "my_list[2:5] = [20, 30, 40]".

```
# create a list
my_list = [10, 20, 30, 40, 50]
# negative indexing
print(my_list[-1])  # Output: 50
```



<pre>print(my_list[-2])</pre>	#	Output:	40		
<pre># negative slicing print(my_list[-3:]) print(my_list[:-2]) print(my_list[-3:-1]) print(my_list[::-1]) 20, 10]</pre>	# #	Output: Output: Output: Output:	[10, [30,	20, 40]	30]

In this example, we create a list called "my_list".

We then use negative indexing to access the last element of the list using "-1" and the second-tolast element of the list using "-2".

Next, we use negative slicing to access the last three elements of the list using "-3:" and all but the last two elements of the list using ":-2". We also use negative slicing to access the elements from the third-to-last index to the last index using "-3:-1". Lastly, we use negative slicing with a step size of -1 to reverse the order of the list using "[::-1]".

```
# create a nested list
my list = [[1, 2, 3], [4, 5, 6], [7, 8, 9]]
# access elements using indexing
print(my list[0][1])
                                 # Output: 2
print(my list[1][2])
                                 # Output: 6
print(my list[2][0])
                                 # Output: 7
# access elements using slicing
print(my list[0][1:])
                                 # Output: [2, 3]
print(my list[1][:2])
                                 # Output: [4, 5]
                                 # Output: [[1, 2, 3],
print(my list[:2])
[4, 5, 6]]
```

In this example, we create a nested list called "my_list" containing three lists with three elements each.

We use indexing to access specific elements in the nested list. For example, we access the element at row 0, column 1 (which has the value 2) using "my_list[0][1]".

We also use slicing to access certain sections of the nested list. For example, we access the elements in row 0, columns 1 and 2 (which have the values 2 and 3) using "my_list[0][1:]". We access the elements in row 1, columns 0 and 1 (which have the values 4 and 5) using "my_list[1][:2]". We also access the first two rows of the nested list using "my_list[:2]".



List comprehension

List comprehension is a powerful and concise way of creating lists in Python. It allows you to create a new list by iterating over an existing list or other iterable, and applying an expression to each element of the iterable. The resulting list is created in a single line of code, making it a useful tool for data manipulation and analysis.

In SAS, similar functionality can be achieved using the DATA step or PROC SQL, but list comprehension provides a more concise and readable way of achieving the same results in Python.

The basic syntax of a list comprehension is as follows:

```
new_list = [expression for variable in iterable if
condition]
```

Here, new_list is the new list being created, expression is the operation to be performed on each element of the iterable, variable is a variable that takes on each value in the iterable, and iterable is the original list or other iterable.

The if statement is optional and allows you to apply a condition to the elements of the iterable before they are added to the new list.

For example, suppose you have a list of numbers and you want to create a new list that contains only the even numbers from the original list. Here is how you could use list comprehension to achieve this:

numbers = [1, 2, 3, 4, 5, 6, 7, 8, 9, 10]
even numbers = [x for x in numbers if x % 2 == 0]

In this example, x takes on each value in the numbers list, and the expression x % 2 == 0 tests whether x is even. If x is even, it is added to the even_numbers list.

List comprehension can also be used to create lists of tuples, dictionaries, or other complex data structures. For example, suppose you have a list of strings representing names and ages, and you want to create a list of dictionaries where each dictionary has a "name" key and an "age" key. Here is how you could use list comprehension to achieve this:

```
names_and_ages = ["Alice, 25", "Bob, 30", "Charlie,
35"]
people = [{"name": name_age.split(",")[0], "age":
int(name_age.split(",")[1])} for name_age in
names_and_ages]
```



In this example, the name_age variable takes on each value in the names_and_ages list, and the expression name_age.split(",")[0] extracts the name from the string, while int(name_age.split(",")[1]) extracts the age and converts it to an integer. The resulting dictionaries are added to the people list.

Overall, list comprehension is a powerful tool for creating lists in Python. It allows you to write concise and readable code that can be easily understood by others, and it can be used to create complex data structures with just a few lines of code. As a SAS user, learning list comprehension can help you to transition to Python and take advantage of its many data manipulation and analysis tools.

Here are some additional details about list comprehension in Python:

List comprehension can be nested: You can use one or more loops and conditions inside a list comprehension to create a nested list. For example, if you have a list of lists and you want to flatten it, you can use a nested list comprehension:

nested_list = [[1, 2], [3, 4, 5], [6, 7]]
flattened_list = [x for sublist in nested_list for x in
sublist]

In this example, sublist takes on each value in the nested_list, and the inner loop iterates over each element in sublist, adding it to the flattened_list.

List comprehension can be used with other data structures: In addition to lists, you can use list comprehension with other iterable data structures such as tuples, sets, and generators. For example:

my_tuple = (1, 2, 3, 4)
squared_numbers = [x**2 for x in my_tuple]

In this example, x takes on each value in the my_tuple tuple, and the expression x^{**2} squares each value, creating a new list of squared numbers.

List comprehension can be faster than traditional loops: List comprehension is often faster than using traditional loops because it is optimized for performance. However, this depends on the size of the data and the complexity of the expression used in the list comprehension.

List comprehension can be used with functions: You can use functions inside list comprehension to perform operations on each element of the iterable. For example:

```
def is_even(x):
    return x % 2 == 0
numbers = [1, 2, 3, 4, 5, 6]
```



even_numbers = [x for x in numbers if is_even(x)]

In this example, the is_even function is used inside the list comprehension to test whether each number in the numbers list is even.

List comprehension is a concise way of creating a list in Python. It provides a way to create a new list by applying a transformation to each element of an existing iterable object, such as a list, tuple, or set. List comprehension can be a powerful tool for data manipulation and analysis, especially for SAS users who are familiar with the data step.

In Python, list comprehension is written in a compact syntax, enclosed by square brackets []. The general structure of list comprehension is:

Here's an example of using list comprehension in Python to create a new list of squared numbers:

```
# create a list of numbers
numbers = [1, 2, 3, 4, 5]
# use list comprehension to create a new list of
squared numbers
squared_numbers = [num**2 for num in numbers]
# print the new list of squared numbers
print(squared_numbers)
```

Output:

[1, 4, 9, 16, 25]

In this example, we first create a list of numbers [1, 2, 3, 4, 5]. We then use list comprehension to create a new list of squared numbers. The expression in the list comprehension is num**2, which raises each number to the power of 2. The for loop iterates over each number in the original list, and the resulting squared numbers are added to the new list. Finally, we print the new list of squared numbers.

Here's another example that uses list comprehension to create a new list of even numbers:

```
# create a list of numbers
numbers = [1, 2, 3, 4, 5, 6, 7, 8, 9, 10]
# use list comprehension to create a new list of even
numbers
even_numbers = [num for num in numbers if num % 2 == 0]
# print the new list of even numbers
```



print(even numbers)

Output:

[2, 4, 6, 8, 10]

In this example, we create a list of numbers [1, 2, 3, 4, 5, 6, 7, 8, 9, 10]. We then use list comprehension to create a new list of even numbers. The if condition in the list comprehension (num % 2 == 0) filters out any numbers that are not even. The resulting even numbers are added to the new list, which is then printed.

These are just a few examples of how list comprehension can be used in Python. With list comprehension, you can easily manipulate and transform data in a concise and efficient way.

The expression is the transformation applied to each element of the iterable, the variable is the name given to each element of the iterable, and the condition is an optional filter that restricts which elements are included in the new list.

For example, suppose we have a list of integers, and we want to create a new list that contains only the even numbers multiplied by 2. We can use list comprehension as follows:

```
numbers = [1, 2, 3, 4, 5, 6, 7, 8, 9, 10]
new_numbers = [num * 2 for num in numbers if num % 2 ==
0]
print(new numbers)
```

The output will be:

[4, 8, 12, 16, 20]

In this example, the expression is "num * 2", the variable is "num", and the condition is "num % 2 == 0", which ensures that only even numbers are included in the new list.

List comprehension can also be used with nested loops and multiple conditions. For example, suppose we have two lists, and we want to create a new list that contains the product of each pair of elements where the first element is from the first list and the second element is from the second list, but only if the product is greater than 10. We can use list comprehension as follows:

```
list1 = [1, 2, 3]
list2 = [4, 5, 6]
new_list = [x * y for x in list1 for y in list2 if x *
y > 10]
print(new_list)
The output will be:
```



[12, 15, 18]

In this example, we use two nested loops to iterate over each pair of elements from the two lists, and the condition "x * y > 10" filters out pairs whose product is less than or equal to 10.

list comprehension is a powerful and concise way of creating new lists in Python by applying a transformation to each element of an existing iterable object, with the option to include a filter condition. SAS users can benefit from list comprehension to manipulate and analyze data in a similar way to the data step. here's another example of using list comprehension in Python to create a new list of tuples that represent the Cartesian product of two lists:

```
# create two lists
list1 = ['A', 'B', 'C']
list2 = [1, 2, 3]
# use list comprehension to create a new list of tuples
representing the Cartesian product of the two lists
cartesian_product = [(x, y) for x in list1 for y in
list2]
# print the new list of tuples
print(cartesian_product)
Output:
```

```
[('A', 1), ('A', 2), ('A', 3), ('B', 1), ('B', 2),
('B', 3), ('C', 1), ('C', 2), ('C', 3)]
```

In this example, we first create two lists: list1 containing the letters 'A', 'B', and 'C', and list2 containing the numbers 1, 2, and 3. We then use list comprehension to create a new list of tuples representing the Cartesian product of the two lists. The for loops in the list comprehension iterate over each element of list1 and list2, and each pair of elements is added to the new list as a tuple (x, y). Finally, we print the new list of tuples.

Here's another example of using list comprehension in Python to create a new list of strings by concatenating the corresponding elements of two lists:

```
# create two lists
names = ['Alice', 'Bob', 'Charlie']
ages = [25, 30, 35]
# use list comprehension to create a new list of
strings by concatenating the corresponding elements of
the two lists
```



```
string_list = [name + ' is ' + str(age) + ' years old'
for name, age in zip(names, ages)]
# print the new list of strings
print(string_list)
```

Output:

```
['Alice is 25 years old', 'Bob is 30 years old', 'Charlie is 35 years old']
```

In this example, we first create two lists: names containing the names 'Alice', 'Bob', and 'Charlie', and ages containing the corresponding ages 25, 30, and 35. We then use list comprehension to create a new list of strings by concatenating the corresponding elements of the two lists. The zip function is used to iterate over both lists simultaneously, and the for loop unpacks each pair of elements into the variables name and age. The resulting strings are added to the new list, which is then printed.

Here's another example of using list comprehension in Python to create a new list of dictionaries:

```
# create a list of keys and a list of values
keys = ['a', 'b', 'c']
values = [1, 2, 3]
# use list comprehension to create a new list of
dictionaries
dict_list = [{key: value} for key, value in zip(keys,
values)]
# print the new list of dictionaries
print(dict_list)
```

Output:

[{'a': 1}, {'b': 2}, {'c': 3}]

In this example, we first create two lists: keys containing the keys 'a', 'b', and 'c', and values containing the corresponding values 1, 2, and 3. We then use list comprehension to create a new list of dictionaries, where each dictionary contains a single key-value pair from the two input lists. The zip function is used to iterate over both lists simultaneously, and the for loop unpacks each pair of elements into the variables key and value. The resulting dictionaries are added to the new list, which is then printed.

Here's another example of using list comprehension in Python to create a new list of sets by combining the elements of two lists:



```
# create two lists
list1 = [1, 2, 3]
list2 = [2, 3, 4]
# use list comprehension to create a new list of sets
by combining the elements of the two lists
set_list = [{x, y} for x in list1 for y in list2]
# print the new list of sets
print(set_list)
```

Output:

```
[\{1, 2\}, \{1, 3\}, \{1, 4\}, \{2, 2\}, \{2, 3\}, \{2, 4\}, \{3, 2\}, \{3, 3\}, \{3, 4\}]
```

In this example, we first create two lists: list1 containing the numbers 1, 2, and 3, and list2 containing the numbers 2, 3, and 4. We then use list comprehension to create a new list of sets by combining the elements of the two lists. The for loops in the list comprehension iterate over each element of list1 and list2, and each pair of elements is added to the new list as a set $\{x, y\}$.

```
# create a list of numbers
numbers = [1, 2, 3, 4, 5, 6, 7, 8, 9, 10]
# use list comprehension to filter even numbers from
the list
even_numbers = [x for x in numbers if x % 2 == 0]
# print the new list of even numbers
print(even_numbers)
```

Output:

[2, 4, 6, 8, 10]

In this example, we first create a list of numbers from 1 to 10. We then use list comprehension to create a new list of even numbers by filtering elements from the original list based on the condition x % 2 == 0, which checks if x is even. The resulting even numbers are added to the new list, which is then printed.



Sorting and filtering data

Sorting and filtering data are important operations in data analysis, and they can be performed using Python in a similar manner as in SAS. In this article, we will provide an introduction to sorting and filtering data in Python for SAS users.

Sorting Data

In SAS, we use the SORT procedure to sort data sets based on one or more variables. In Python, we can use the sorted() function to sort lists or arrays based on one or more elements. The function returns a new sorted list, leaving the original list unchanged. Here's an example:

```
>>> numbers = [3, 1, 4, 1, 5, 9, 2, 6, 5, 3, 5]
>>> sorted_numbers = sorted(numbers)
>>> print(sorted_numbers)
[1, 1, 2, 3, 3, 4, 5, 5, 5, 6, 9]
```

We can also use the sort() method of a list to sort it in place:

```
>>> numbers = [3, 1, 4, 1, 5, 9, 2, 6, 5, 3, 5]
>>> numbers.sort()
>>> print(numbers)
[1, 1, 2, 3, 3, 4, 5, 5, 5, 6, 9]
```

To sort a list of dictionaries based on a particular key, we can use the key argument of the sorted() function:

```
>>> students = [
... {"name": "Alice", "age": 20},
... {"name": "Bob", "age": 18},
... {"name": "Charlie", "age": 22},
... ]
>>> sorted_students = sorted(students, key=lambda
student: student["age"])
>>> print(sorted_students)
[{'name': 'Bob', 'age': 18}, {'name': 'Alice', 'age':
20}, {'name': 'Charlie', 'age': 22}]
```

In this example, we sort the list of dictionaries based on the "age" key. The lambda function returns the value of the "age" key for each dictionary.



Filtering Data

In SAS, we use the WHERE statement to select observations that meet certain criteria. In Python, we can use list comprehension to filter lists or arrays based on a condition. For example:

```
>>> numbers = [3, 1, 4, 1, 5, 9, 2, 6, 5, 3, 5]
>>> odd_numbers = [n for n in numbers if n % 2 == 1]
>>> print(odd_numbers)
[3, 1, 1, 5, 9, 5, 3, 5]
```

In this example, we use list comprehension to create a new list that contains only the odd numbers in the original list.

We can also use the filter() function to filter a list or array based on a function that returns True or False. For example:

>>> def is_odd(n): ... return n % 2 == 1 ... >>> odd_numbers = list(filter(is_odd, numbers)) >>> print(odd_numbers) [3, 1, 1, 5, 9, 5, 3, 5]

In this example, we define a function is_odd() that returns True if a number is odd. We then use the filter() function to create a new list that contains only the odd numbers in the original list.

Sorting Data

In addition to sorting lists, we can also sort pandas DataFrames, which are similar to SAS data sets. Pandas provides the sort_values() method to sort a DataFrame based on one or more columns. Here's an example:

```
import pandas as pd
# create a DataFrame
df = pd.DataFrame({
        'name': ['Alice', 'Bob', 'Charlie', 'Dave'],
        'age': [25, 30, 35, 40],
        'score': [80, 70, 90, 85]
})
# sort by age in ascending order
df_sorted = df.sort_values(by='age')
print(df_sorted)
```



This will output:

	name	age	score
0	Alice	25	80
1	Bob	30	70
2	Charlie	35	90
3	Dave	40	85

We can also sort by multiple columns:

```
# sort by age in descending order, then score in
ascending order
df_sorted = df.sort_values(by=['age', 'score'],
ascending=[False, True])
```

print(df_sorted)

This will output:

	name	age	score
3	Dave	40	85
2	Charlie	35	90
1	Bob	30	70
0	Alice	25	80

Filtering Data

In addition to filtering lists, we can also filter pandas DataFrames using boolean indexing. Boolean indexing allows us to select rows based on a condition. Here's an example:

```
# filter by age greater than 30
df_filtered = df[df['age'] > 30]
```

print(df_filtered)

This will output:

	name	age	score
2	Charlie	35	90
3	Dave	40	85



Chapter 3: Reading and Writing Data



Python is a popular programming language that is widely used in data analysis, machine learning, and scientific computing. It has a rich set of libraries and frameworks that make it easy to read and write data in a variety of formats, including CSV, Excel, SQL, and SAS. In this article, we will provide a SAS-oriented introduction to Python, focusing on how SAS users can use Python to read and write data.

Reading Data

There are several ways to read data in Python, depending on the format of the data. We will discuss some of the most commonly used methods below.

Reading CSV files: The most common format for data files is CSV (comma-separated values). Python has a built-in CSV module that makes it easy to read and write CSV files. Here's an example of how to read a CSV file using the CSV module:

```
import csv
with open('data.csv', 'r') as file:
    reader = csv.reader(file)
    for row in reader:
        print(row)
```

In this example, we use the open() function to open the CSV file in read mode. We then pass the file object to the csv.reader() function, which returns a reader object that we can use to iterate over the rows in the file. Each row is returned as a list of strings.

Reading Excel files: Another common format for data files is Excel. Python has a library called pandas that makes it easy to read and write Excel files. Here's an example of how to read an Excel file using pandas:

```
import pandas as pd
df = pd.read_excel('data.xlsx')
print(df)
```

In this example, we use the read_excel() function from the pandas library to read the Excel file. The function returns a pandas DataFrame object, which is a two-dimensional table of data with rows and columns.

Reading SQL databases: If your data is stored in a SQL database, you can use Python's built-in sqlite3 module to read the data. Here's an example of how to read data from a SQLite database:

```
import sqlite3
conn = sqlite3.connect('mydatabase.db')
```



```
cursor = conn.cursor()
cursor.execute('SELECT * FROM mytable')
rows = cursor.fetchall()
for row in rows:
    print(row)
```

In this example, we use the connect() function from the sqlite3 module to connect to the SQLite database. We then create a cursor object and use it to execute a SQL query to select all rows from the mytable table. We then use the fetchall() method to retrieve all rows from the query result.

```
import csv
data = [
    ['Name', 'Age', 'Gender'],
    ['John', 25, 'Male'],
    ['Jane', 30, 'Female'],
    ['Bob', 40, 'Male']
]
with open('data.csv', 'w', newline='') as file:
    writer = csv.writer(file)
    writer.writerows(data)
```

In this example, we define a list of lists called data that contains the data we want to write to the CSV file. We then use the open() function to open the file in write mode and pass the file object to the csv.writer() function. We then use the writerows() method to write all rows of data to the CSV file.

Writing Excel files: To write data to an Excel file, you can use the pandas library. Here's an example of how to write data to an Excel file:

```
import pandas as pd
data = {
    'Name': ['John', 'Jane', 'Bob'],
    'Age': [25, 30, 40],
    'Gender': ['Male', 'Female', 'Male']
}
df = pd.DataFrame(data)
df.to excel('data.xlsx', index=False)
```

In this example, we define a dictionary called data that contains the data we want to write to the Excel file. We then use the DataFrame() function from the pandas library to create a DataFrame



object from the dictionary. We then use the to_excel() method to write the DataFrame to an Excel file.

Writing SQL databases: To write data to a SQL database, you can use the execute() method of the cursor object from the sqlite3 module. Here's an example of how to write data to a SQLite database:

```
import sqlite3
conn = sqlite3.connect('mydatabase.db')
cursor = conn.cursor()
cursor.execute('''CREATE TABLE mytable (Name TEXT, Age
INTEGER, Gender TEXT)''')
data = [('John', 25, 'Male'), ('Jane', 30, 'Female'),
('Bob', 40, 'Male')]
cursor.executemany('INSERT INTO mytable VALUES
(?,?,?)', data)
conn.commit()
```

In this example, we first create a connection object to the SQLite database using the connect() function from the sqlite3 module. We then create a cursor object and use it to execute a SQL query to create a table called mytable. We then define a list of tuples called data that contains the data we want to write to the table. We then use the executemany() method to execute a SQL query to insert multiple rows of data into the table. Finally, we use the commit() method to commit the changes to the database.

here is a comprehensive guide on Reading and Writing Data in Python for SAS Users.

Introduction

Python is a general-purpose programming language that is widely used in the data science and machine learning communities. Python is often compared to SAS, a popular statistical software package that has been widely used in the industry for data analysis and manipulation. However, Python offers several advantages over SAS, including a more robust and flexible set of data manipulation tools. In this guide, we will discuss how to read and write data in Python, with a focus on how these operations compare to their SAS equivalents.

Reading Data

Reading CSV Files

In SAS, data is typically stored in a variety of formats, including CSV, Excel, and SAS datasets. To read a CSV file in Python, we can use the pandas library. pandas is a popular data manipulation library that provides several functions for reading and manipulating data.

To read a CSV file in Python, we can use the read_csv function from pandas. For example, suppose we have a CSV file called mydata.csv that contains the following data:



Name, Age, Gender John, 30, Male Jane, 25, Female

To read this file in Python, we can use the following code:

```
import pandas as pd
df = pd.read csv('mydata.csv')
```

This will read the CSV file into a DataFrame object, which is a two-dimensional table of data that can be manipulated using various pandas functions.

Reading Excel Files

In SAS, Excel files are often used to store data. To read an Excel file in Python, we can use the pandas library. pandas provides a read_excel function that can read Excel files in a similar way to read_csv.

To read an Excel file in Python, we can use the following code:

```
import pandas as pd
df = pd.read_excel('mydata.xlsx')
```

This will read the Excel file into a DataFrame object.

Reading SAS Datasets To read a SAS dataset in Python, we can use the sas7bdat library. This library provides a SAS7BDAT class that can read SAS datasets. To use this library, we need to install it using pip:

pip install sas7bdat

To read a SAS dataset in Python, we can use the following code:

```
from sas7bdat import SAS7BDAT
with SAS7BDAT('mydata.sas7bdat') as f:
    df = f.to data frame()
```

This will read the SAS dataset into a DataFrame object.

Writing Data Writing CSV Files



To write data to a CSV file in Python, we can use the to_csv function from pandas. For example, suppose we have a DataFrame called df that we want to write to a CSV file called output.csv. We can use the following code:

```
import pandas as pd
df.to csv('output.csv', index=False)
```

This will write the DataFrame to a CSV file.

Writing Excel Files

To write data to an Excel file in Python, we can use the to_excel function from pandas. For example, suppose we have a DataFrame called df that we want to write to an Excel file called output.xlsx. We can use the following code:

```
import pandas as pd
df.to_excel('output.xlsx', index=False)
```

Here are some example codes for reading and writing data in Python for SAS users.

Reading CSV Files

```
import pandas as pd
# Read CSV file into a DataFrame
df = pd.read csv('mydata.csv')
```

Reading Excel Files

import pandas as pd
Read Excel file into a DataFrame
df = pd.read_excel('mydata.xlsx')

Reading SAS Datasets

from sas7bdat import SAS7BDAT
Read SAS dataset into a DataFrame
with SAS7BDAT('mydata.sas7bdat') as f:
 df = f.to_data_frame()
Writing CSV Files

in stal

```
import pandas as pd
# Write DataFrame to CSV file
df.to_csv('output.csv', index=False)
```

Writing Excel Files

```
import pandas as pd
# Write DataFrame to Excel file
df.to_excel('output.xlsx', index=False)
```

Writing SAS Datasets

```
import pandas as pd
from sas7bdat import SAS7BDAT
# Write DataFrame to SAS dataset
with SAS7BDAT('output.sas7bdat', 'w') as f:
    f.write_df(df)
```

Reading CSV Files with Custom Delimiters

Sometimes CSV files may use a custom delimiter character instead of a comma. In SAS, we can use the dlm option in the infile statement to specify a custom delimiter. In Python, we can use the delimiter parameter in the read_csv function to specify a custom delimiter. For example, suppose we have a CSV file called mydata.txt that uses a pipe (|) as a delimiter:

```
Name | Age | Gender
John | 30 | Male
Jane | 25 | Female
```

To read this file in Python, we can use the following code:

```
import pandas as pd
# Read CSV file with custom delimiter
df = pd.read csv('mydata.txt', delimiter='|')
```

This will read the CSV file into a DataFrame object using a pipe (|) as the delimiter.



Reading SAS Datasets with Custom Formats

In SAS, we can apply custom formats to variables in a dataset using the format statement. In Python, we can use the sas7bdat library to read SAS datasets that use custom formats. For example, suppose we have a SAS dataset called mydata.sas7bdat that contains a variable called age that has a custom format called agefmt.:

```
data mydata;
  set sashelp.class;
  format age agefmt.;
run;
```

To read this SAS dataset in Python, we can use the following code:

```
from sas7bdat import SAS7BDAT
# Read SAS dataset with custom formats
with SAS7BDAT('mydata.sas7bdat') as f:
    df = f.to_data_frame(convert_dates=False,
    convert_text=False)
```

This will read the SAS dataset into a DataFrame object, preserving the custom format for the age variable.

Writing CSV Files with Custom Delimiters

To write a DataFrame to a CSV file with a custom delimiter in Python, we can use the sep parameter in the to_csv function. For example, suppose we have a DataFrame called df that we want to write to a CSV file called output.txt using a pipe (|) as the delimiter:

```
import pandas as pd
# Write DataFrame to CSV file with custom delimiter
df.to csv('output.txt', sep='|', index=False)
```

This will write the DataFrame to a CSV file using a pipe (|) as the delimiter.

Writing SAS Datasets with Custom Formats

To write a DataFrame to a SAS dataset with custom formats in Python, we can use the write_df method in the SAS7BDAT class from the sas7bdat library. For example, suppose we have a DataFrame called df that we want to write to a SAS dataset called output.sas7bdat, with a custom format for the age variable:

```
import pandas as pd
from sas7bdat import SAS7BDAT
```

Write DataFrame to SAS dataset with custom formats



```
with SAS7BDAT('output.sas7bdat', 'w') as f:
    f.write_df(df, formats={'age': 'agefmt.'})
```

This will write the DataFrame to a SAS dataset, applying the custom format agefmt. to the age variable.

Reading data from files

"Python for SAS Users" is a book that provides an introduction to Python programming for SAS users. The book is designed to help SAS users transition to Python programming by highlighting the similarities and differences between the two languages. One important aspect of data analysis in both SAS and Python is the ability to read data from files. In this article, we will provide an overview of how to read data from files in Python, with a focus on techniques that may be familiar to SAS users.

Reading Text Files

The most common type of file that SAS users may encounter is a text file. In Python, there are several ways to read text files, including the built-in open function, the csv module, and the pandas library.

The open function is a built-in function that allows you to open a file in text mode. Here's an example:

```
with open('file.txt', 'r') as f:
    for line in f:
        print(line.strip())
```

This code opens a file called 'file.txt' in read mode ('r'), and then iterates through each line in the file using a for loop. The strip method is used to remove any leading or trailing whitespace from each line.

The csv module is a standard library in Python that provides functionality for reading and writing CSV (comma-separated values) files. Here's an example:

```
import csv
with open('file.csv', 'r') as f:
    reader = csv.reader(f)
    for row in reader:
        print(row)
```

This code opens a file called 'file.csv' in read mode, and then creates a CSV reader object using the csv.reader function. The reader object is then used to iterate through each row in the file.



The pandas library is a popular library in Python for data analysis that provides functionality for reading and writing various types of files, including CSV files, Excel files, and SQL databases. Here's an example:

```
import pandas as pd
df = pd.read_csv('file.csv')
print(df)
```

This code uses the read_csv function from the pandas library to read a CSV file called 'file.csv' into a DataFrame object. The resulting DataFrame can then be used for data analysis.

Reading SAS Data Sets

SAS users may also be familiar with reading SAS data sets (.sas7bdat files) into SAS. In Python, there are several libraries that provide functionality for reading SAS data sets, including the sas7bdat library and the pandas library.

The sas7bdat library is a Python library that provides functionality for reading SAS data sets directly into Python. Here's an example:

```
from sas7bdat import SAS7BDAT
with SAS7BDAT('file.sas7bdat') as f:
    data = f.to_data_frame()
    print(data)
```

This code opens a SAS data set called 'file.sas7bdat' using the SAS7BDAT function from the sas7bdat library. The resulting object is then converted to a Pandas DataFrame using the to_data_frame method.

The pandas library also provides functionality for reading SAS data sets using the read_sas function. Here's an example:

```
import pandas as pd
df = pd.read_sas('file.sas7bdat')
print(df)
```

This code uses the read_sas function from the pandas library to read a SAS data set called 'file.sas7bdat' into a DataFrame object.

Reading Excel Files

Another common type of file that SAS users may encounter is an Excel file. In Python, there are several libraries that provide functionality for reading Excel files, including the openpyxl library and the pandas library.



The openpyxl library is a Python library that provides functionality for reading and writing Excel files. Here's an example:

```
from openpyxl import load_workbook
wb = load_workbook(filename='file.xlsx')
ws = wb.active
for row in ws.iter_rows():
    for cell in row:
        print(cell.value)
```

This code uses the load_workbook function from the openpyxl library to load an Excel file called 'file.xlsx' into a workbook object. The active worksheet is then selected using the active attribute, and the iter_rows method is used to iterate through each row in the worksheet. The value of each cell is then printed using the value attribute.

The pandas library also provides functionality for reading Excel files using the read_excel function. Here's an example:

```
import pandas as pd
df = pd.read_excel('file.xlsx')
print(df)
```

This code uses the read_excel function from the pandas library to read an Excel file called 'file.xlsx' into a DataFrame object. The resulting DataFrame can then be used for data analysis.

Reading SQL Databases

SAS users may also be familiar with reading data from SQL databases. In Python, there are several libraries that provide functionality for reading data from SQL databases, including the pyodbc library and the pandas library.

The pyodbc library is a Python library that provides functionality for connecting to and querying SQL databases. Here's an example:

```
import pyodbc
cnxn = pyodbc.connect('DRIVER={SQL
Server};SERVER=server_name;DATABASE=db_name;UID=usernam
e;PWD=password')
cursor = cnxn.cursor()
cursor.execute('SELECT * FROM table_name')
for row in cursor:
```

print(row)

This code uses the connect function from the pyodbc library to connect to a SQL Server database using the specified driver, server name, database name, username, and password. The cursor object is then used to execute a SQL query and iterate through the results.

The pandas library also provides functionality for reading data from SQL databases using the read_sql function. Here's an example:

```
import pandas as pd
import pyodbc
cnxn = pyodbc.connect('DRIVER={SQL
Server};SERVER=server_name;DATABASE=db_name;UID=usernam
e;PWD=password')
query = 'SELECT * FROM table_name'
df = pd.read_sql(query, cnxn)
print(df)
```

This code uses the connect function from the pyodbc library to connect to a SQL Server database using the specified driver, server name, database name, username, and password. The SQL query is then stored in a variable called query, and the read_sql function is used to read the data into a DataFrame object. The resulting DataFrame can then be used for data analysis.

Reading Data from Web APIs

Web APIs (Application Programming Interfaces) are a popular source of data for data analysis, and Python provides several libraries for accessing data from web APIs, such as requests and urllib.

The requests library provides a simple and intuitive way to send HTTP requests and receive responses. Here's an example of using the requests library to access data from a web API:

```
import requests
url = 'https://api.example.com/data'
response = requests.get(url)
if response.status_code == 200:
    data = response.json()
    print(data)
else:
    print('Error: Unable to access API')
```

In this example, the get method from the requests library is used to send an HTTP GET request to the URL 'https://api.example.com/data'. If the response status code is 200 (indicating a



successful response), the response data is then converted to a Python dictionary using the json method and printed to the console. If there is an error, an error message is printed.

Reading Data from Web Scraping

Another way to obtain data from the web is through web scraping, which involves extracting data from HTML and other web page formats. The BeautifulSoup library is a popular library for web scraping in Python.

Here's an example of using the BeautifulSoup library to scrape data from a web page:

```
import requests
from bs4 import BeautifulSoup
url = 'https://www.example.com'
response = requests.get(url)
if response.status_code == 200:
    soup = BeautifulSoup(response.text, 'html.parser')
    links = soup.find_all('a')
    for link in links:
        print(link.get('href'))
else:
    print('Error: Unable to access web page')
```

In this example, the get method from the requests library is used to send an HTTP GET request to the URL 'https://www.example.com'. If the response status code is 200 (indicating a successful response), the BeautifulSoup library is used to parse the HTML response into a soup object. The find_all method is then used to find all the links on the page, and the get method is used to extract the URL for each link.

Here are some longer examples of Python code for reading data from different sources:

Reading Data from Text Files

```
# Open the text file in read mode
with open('example.txt', 'r') as file:
    # Read all the lines from the file
    lines = file.readlines()
    # Print the lines
    for line in lines:
        print(line.strip())
```

In this example, the open function is used to open the file 'example.txt' in read mode, and the readlines method is used to read all the lines from the file into a list. The strip method is then



used to remove any whitespace characters from the beginning and end of each line, and the lines are printed to the console.

Reading Data from SAS Data Sets

```
# Import the SAS7BDAT module
import sas7bdat
# Open the SAS data set file
with sas7bdat.SAS7BDAT('example.sas7bdat') as file:
    # Get the column names from the SAS data set
    columns = file.columns
    # Print the column names
    for column in columns:
        print(column)
    # Get the rows from the SAS data set
    rows = file.to_data_frame()
    # Print the rows
    for row in rows:
        print(row)
```

In this example, the sas7bdat module is used to open the SAS data set file 'example.sas7bdat'. The columns attribute is used to get the column names from the data set, and the to_data_frame method is used to get the rows from the data set and convert them to a pandas data frame. The column names and rows are then printed to the console.

Reading Data from Excel Files

```
# Import the pandas module
import pandas as pd
# Read the Excel file into a pandas data frame
df = pd.read_excel('example.xlsx')
# Print the data frame
print(df)
```

In this example, the pandas module is used to read the Excel file 'example.xlsx' into a pandas data frame using the read_excel function. The data frame is then printed to the console.

Reading Data from SQL Databases



```
# Import the pandas module
import pandas as pd
import sqlite3
# Connect to the SQLite database
conn = sqlite3.connect('example.db')
# Read the data from the database into a pandas data
frame
df = pd.read_sql_query('SELECT * FROM table', conn)
# Print the data frame
print(df)
```

In this example, the pandas and sqlite3 modules are used to connect to the SQLite database file 'example.db'. The read_sql_query function is then used to read the data from the 'table' table into a pandas data frame. The data frame is then printed to the console.

Reading Data from Web APIs

```
# Import the requests module
import requests
# Make a request to the web API
response = requests.get('https://api.example.com/data')
# Check the response status code
if response.status_code == 200:
    # Get the response data as a dictionary
    data = response.json()
    # Print the data
    print(data)
else:
    print('Error: Unable to access API')
```

In this example, the requests module is used to make a request to the web API at 'https://api.example.com/data'. The status_code attribute is then used to check the response status code, and if it is 200 (indicating a successful response).

Reading Data from CSV Files

Import the pandas module import pandas as pd



```
# Read the CSV file into a pandas data frame
df = pd.read_csv('example.csv')
# Print the data frame
print(df)
```

In this example, the pandas module is used to read the CSV file 'example.csv' into a pandas data frame using the read_csv function. The data frame is then printed to the console.

Reading Data from JSON Files

```
# Import the json module
import json
# Open the JSON file in read mode
with open('example.json', 'r') as file:
    # Load the JSON data into a dictionary
    data = json.load(file)
    # Print the data
    print(data)
```

In this example, the json module is used to open the JSON file 'example.json' in read mode using the open function. The load method is then used to load the JSON data into a dictionary. The dictionary is then printed to the console.

Reading Data from XML Files

```
# Import the ElementTree module
import xml.etree.ElementTree as ET
# Parse the XML file
tree = ET.parse('example.xml')
root = tree.getroot()
# Print the XML data
for child in root:
    print(child.tag, child.attrib)
```

In this example, the xml.etree.ElementTree module is used to parse the XML file 'example.xml' using the parse function. The getroot method is then used to get the root element of the XML tree, and the tag and attrib attributes are used to print the data from each child element.

Reading Data from PDF Files



```
python
    # Import the PyPDF2 module
    import PyPDF2
    # Open the PDF file in read-binary mode
    with open('example.pdf', 'rb') as file:
        # Create a PDF reader object
        reader = PyPDF2.PdfFileReader(file)
        # Get the number of pages in the PDF file
        num pages = reader.getNumPages()
        # Print the number of pages
        print('Number of Pages:', num pages)
        # Loop through each page in the PDF file
        for i in range(num pages):
             # Get the text from the page
             page = reader.getPage(i)
             text = page.extractText()
             # Print the text
            print('Page', i+1, 'Text:', text)
```

In this example, the PyPDF2 module is used to open the PDF file 'example.pdf' in read-binary mode using the open function. The PdfFileReader class is then used to create a PDF reader object. The getNumPages method is used to get the number of pages in the PDF file, and the getPage and extractText methods are used to get the text from each page. The text is then printed to the console.

Reading Data from Excel Files

```
# Import the pandas module
import pandas as pd
# Read the Excel file into a pandas data frame
df = pd.read_excel('example.xlsx')
# Print the data frame
print(df)
```

In this example, the pandas module is used to read the Excel file 'example.xlsx' into a pandas data frame using the read_excel function. The data frame is then printed to the console.



Reading Data from SQLite Databases

```
# Import the sqlite3 module
import sqlite3
# Connect to the SQLite database
conn = sqlite3.connect('example.db')
# Create a cursor object
cur = conn.cursor()
# Execute a SELECT statement
cur.execute('SELECT * FROM employees')
# Fetch the results
results = cur.fetchall()
# Print the results
for row in results:
    print(row)
# Close the connection
conn.close()
```

In this example, the sqlite3 module is used to connect to the SQLite database 'example.db' using the connect function. A cursor object is then created using the cursor method. A SELECT statement is executed using the execute method, and the results are fetched using the fetchall method. The results are then printed to the console using a for loop. Finally, the connection is closed using the close method.

Reading Data from MySQL Databases

```
# Import the mysql-connector-python module
import mysql.connector
# Connect to the MySQL database
conn = mysql.connector.connect(
    host='localhost',
    user='root',
    password='password',
    database='example'
)
# Create a cursor object
in stal
```

```
cur = conn.cursor()
# Execute a SELECT statement
cur.execute('SELECT * FROM employees')
# Fetch the results
results = cur.fetchall()
# Print the results
for row in results:
    print(row)
# Close the connection
conn.close()
```

In this example, the mysql-connector-python module is used to connect to the MySQL database 'example' on the local machine using the connect function. A cursor object is then created using the cursor method. A SELECT statement is executed using the execute method, and the results are fetched using the fetchall method. The results are then printed to the console using a for loop. Finally, the connection is closed using the close method. Note that the host, user, password, and database parameters will need to be adjusted to match your own MySQL setup.

Reading Data from JSON Files

```
# Import the json module
import json
# Read the JSON file
with open('example.json') as f:
    data = json.load(f)
# Print the data
print(data)
```

In this example, the json module is used to read the JSON file 'example.json' into a Python dictionary using the load function. The data is then printed to the console.

Reading Data from XML Files

```
# Import the xml.etree.ElementTree module
import xml.etree.ElementTree as ET
# Parse the XML file
tree = ET.parse('example.xml')
```



```
# Get the root element
root = tree.getroot()
# Print the elements and attributes
for child in root:
    print(child.tag, child.attrib)
```

In this example, the xml.etree.ElementTree module is used to parse the XML file 'example.xml' into an ElementTree object using the parse function. The root element is then obtained using the getroot method. The elements and attributes are then printed to the console using a for loop.

Reading Data from CSV Files

```
# Import the csv module
import csv
# Read the CSV file
with open('example.csv', newline='') as f:
    reader = csv.reader(f)
    for row in reader:
        print(row)
```

In this example, the csv module is used to read the CSV file 'example.csv' using the reader function. The rows are then printed to the console using a for loop. Note that the newline=" parameter is used to ensure that newlines are not interpreted as row separators.

Reading Data from Fixed-Width Files

```
# Define the field widths
widths = [10, 10, 10]
# Read the fixed-width file
with open('example.txt') as f:
    for line in f:
        fields = [line[start:start+width].strip() for
start, width in enumerate(widths)]
        print(fields)
```

In this example, a list of field widths is defined. The fixed-width file 'example.txt' is then read line by line using a for loop. The fields are extracted using list comprehension and the strip method to remove any whitespace. The fields are then printed to the console.

Reading Data from PDF Files



```
# Import the PyPDF2 module
import PyPDF2
# Open the PDF file in binary mode
with open('example.pdf', 'rb') as f:
    # Create a PDF reader object
    reader = PyPDF2.PdfReader(f)
    # Loop through the pages
    for page in reader.pages:
        # Extract the text from the page
        text = page.text
    # Print the text
    print(text)
```

In this example, the PyPDF2 module is used to open the PDF file 'example.pdf' in binary mode and create a PDF reader object using the PdfReader function. The pages are then looped through using a for loop, and the text is extracted from each page using the text attribute. The text is then printed to the console.

Reading Data from HTML Files

```
# Import the BeautifulSoup module
from bs4 import BeautifulSoup
# Read the HTML file
with open('example.html') as f:
    # Create a BeautifulSoup object
    soup = BeautifulSoup(f, 'html.parser')
    # Find all the table rows
    rows = soup.find all('tr')
    # Loop through the rows
    for row in rows:
        # Find all the table cells in the row
        cells = row.find all('td')
        # Extract the text from each cell
        data = [cell.text for cell in cells]
        # Print the data
        print(data)
```



In this example, the BeautifulSoup module is used to read the HTML file 'example.html' and create a BeautifulSoup object using the BeautifulSoup function. All the table rows are then found using the find_all method. The rows are then looped through using a for loop, and all the table cells in each row are found using the find_all method. The text is then extracted from each cell using the text attribute and printed to the console.

Reading Data from Text Files with Regular Expressions

```
# Import the re module
import re
# Read the text file
with open('example.txt') as f:
    # Create a regular expression pattern
    pattern = r'^\w+,\s\w+\s\w+'
    # Loop through the lines
    for line in f:
        # Match the pattern
        match = re.match(pattern, line)
    # If the pattern matches, print the line
    if match:
        print(line)
```

In this example, the re module is used to read the text file 'example.txt' and create a regular expression pattern using the match function. The lines are then looped through using a for loop, and the pattern is matched against each line using the match function. If the pattern matches, the line is printed to the console.

Writing data to files

Python and SAS are two popular programming languages used in data science and analytics. While SAS has been the go-to language for statistical analysis and data management for several decades, Python has gained a lot of popularity in recent years for its versatility and ease of use. In this article, we will discuss how to write data to files using Python from a SAS user's perspective.

Python provides several built-in functions to write data to files, including the open() function, which is used to create a file object. The syntax for open() function is as follows:

```
file object = open(file name, access mode)
```



Here, file_name is the name of the file you want to create or open, and access_mode specifies the mode in which the file should be opened. For example, the access mode w is used to open the file in write mode, r is used to open the file in read mode, and a is used to open the file in append mode.

To write data to a file, you can use the write() method of the file object. Here's an example code snippet that demonstrates how to write data to a file using Python:

```
# create a file object
file_object = open('data.txt', 'w')
# write data to the file
file_object.write('This is line 1\n')
file_object.write('This is line 2\n')
# close the file
file_object.close()
```

In this example, we create a file object named file_object using the open() function with the file name data.txt and access mode w. Then, we write two lines of text to the file using the write() method. Finally, we close the file using the close() method.

It's important to note that when you write data to a file using the write() method, you need to include the newline character (n) at the end of each line to create a new line in the file.

In addition to the write() method, Python provides several other methods to write data to files, such as writelines() and print(). The writelines() method is used to write a list of strings to a file, while the print() function can be used to write data to a file by redirecting the output to the file.

Here's an example code snippet that demonstrates how to use the writelines() method to write data to a file:

```
# create a list of strings
data = ['This is line 1\n', 'This is line 2\n']
# create a file object
file_object = open('data.txt', 'w')
# write data to the file using the writelines() method
file_object.writelines(data)
# close the file
file object.close()
```



In this example, we create a list of strings named data containing two lines of text. Then, we create a file object named file_object using the open() function with the file name data.txt and access mode w. Finally, we use the writelines() method to write the list of strings to the file, and close the file using the close() method.

Writing data to files is an essential part of data processing and analysis. Python provides several built-in functions and methods to write data to files, making it a useful tool for SAS users who want to expand their data processing and analysis capabilities. By using the open() function and the write() or writelines() methods, SAS users can easily write data to files using Python.

some more information on writing data to files using Python from a SAS user's perspective.

Writing CSV Files:

Comma-Separated Values (CSV) is a widely used file format for storing and exchanging data. In SAS, PROC EXPORT can be used to export data to CSV format. Similarly, in Python, the csv module provides built-in functions to read and write CSV files.

Here's an example code snippet that demonstrates how to write data to a CSV file using Python:

```
import csv
# create a list of dictionaries containing data
data = [
    {'Name': 'John', 'Age': 25, 'Gender': 'Male'},
    {'Name': 'Emily', 'Age': 28, 'Gender': 'Female'},
    {'Name': 'David', 'Age': 35, 'Gender': 'Male'}
1
# define the field names
fieldnames = ['Name', 'Age', 'Gender']
# create a CSV writer object
csv writer = csv.DictWriter(open('data.csv', 'w',
newline=''), fieldnames=fieldnames)
# write the header row
csv writer.writeheader()
# write the data rows
for row in data:
    csv writer.writerow(row)
```

In this example, we create a list of dictionaries containing data and define the field names. Then, we create a csv.DictWriter object with the file name data.csv and field names fieldnames. We



use the writeheader() method to write the header row and the writerow() method to write each row of data to the CSV file.

Writing Excel Files:

Excel is another popular file format for storing and analyzing data. In SAS, PROC EXPORT can be used to export data to Excel format. Similarly, in Python, the openpyxl module provides built-in functions to read and write Excel files.

Here's an example code snippet that demonstrates how to write data to an Excel file using Python:

```
from openpyxl import Workbook
# create a workbook object
workbook = Workbook()
# select the active worksheet
worksheet = workbook.active
# create a list of tuples containing data
data = [
    ('John', 25, 'Male'),
    ('Emily', 28, 'Female'),
    ('David', 35, 'Male')
1
# write the header row
worksheet.append(('Name', 'Age', 'Gender'))
# write the data rows
for row in data:
    worksheet.append(row)
# save the workbook
workbook.save('data.xlsx')
```

In this example, we create a Workbook object and select the active worksheet. We create a list of tuples containing data and use the append() method to write the header row and each row of data to the worksheet. Finally, we save the workbook to the file data.xlsx.

Writing Text Files:

In addition to CSV and Excel files, SAS users may also need to write data to text files. In SAS, PROC EXPORT can be used to export data to text format. Similarly, in Python, the open() function can be used to write data to text files.



Here's an example code snippet that demonstrates how to write data to a text file using Python:

```
# create a list of strings containing data
data = ['This is line 1', 'This is line 2', 'This is
line 3']
# create a file object
file_object = open('data.txt', 'w')
# write data to the file
for line in data:
    file_object.write(line + '\n')
# close the file
file_object.close()
```

Writing JSON Files:

JSON (JavaScript Object Notation) is a lightweight data format that is commonly used for exchanging data between web applications. In SAS, PROC EXPORT can be used to export data to JSON format. Similarly, in Python, the json module provides built-in functions to read and write JSON files.

Here's an example code snippet that demonstrates how to write data to a JSON file using Python:

```
import json
# create a dictionary containing data
data = {
    'Name': 'John',
    'Age': 25,
    'Gender': 'Male'
}
# write data to the JSON file
with open('data.json', 'w') as file:
    json.dump(data, file)
```

In this example, we create a dictionary containing data and use the json.dump() function to write the data to a file named data.json.

Writing XML Files:

XML (Extensible Markup Language) is a data format that is widely used for exchanging data between different software applications. In SAS, PROC EXPORT can be used to export data to XML format. Similarly, in Python, the xml.etree.ElementTree module provides built-in functions to read and write XML files.



Here's an example code snippet that demonstrates how to write data to an XML file using Python:

```
import xml.etree.ElementTree as ET
# create an XML document
root = ET.Element('data')
# create a subelement for each data item
name = ET.SubElement(root, 'Name')
name.text = 'John'
age = ET.SubElement(root, 'Age')
age.text = '25'
gender = ET.SubElement(root, 'Gender')
gender.text = 'Male'
# create an ElementTree object
tree = ET.ElementTree(root)
# write the XML file
tree.write('data.xml', xml_declaration=True,
encoding='utf-8')
```

In this example, we create an XML document by creating an element for the root node named data, and then create subelements for each data item. We create an ElementTree object from the root node and use the write() method to write the data to a file named data.xml.

Writing Binary Files:

Binary files are files that contain binary data, such as images or audio files. In SAS, PROC EXPORT cannot be used to export data to binary format. However, in Python, the open() function can be used to write data to binary files.

Here's an example code snippet that demonstrates how to write data to a binary file using Python:

```
# read binary data from a file
with open('image.png', 'rb') as file:
    data = file.read()
# write binary data to a file
with open('new_image.png', 'wb') as file:
    file.write(data)
```



In this example, we read binary data from a file named image.png using the 'rb' mode, and then write the binary data to a new file named new_image.png using the 'wb' mode. The file.read() method reads the entire contents of the file, and the file.write() method writes the binary data to the new file.

Writing Excel Files:

Excel files are widely used for storing and analyzing data. In SAS, PROC EXPORT can be used to export data to Excel format. Similarly, in Python, the pandas module provides built-in functions to write data to Excel files.

Here's an example code snippet that demonstrates how to write data to an Excel file using Python:

```
import pandas as pd
# create a dataframe containing data
data = {
    'Name': ['John', 'Mary', 'Peter'],
    'Age': [25, 30, 35],
    'Gender': ['Male', 'Female', 'Male']
}
df = pd.DataFrame(data)
# write data to the Excel file
with pd.ExcelWriter('data.xlsx') as writer:
    df.to_excel(writer, sheet_name='Sheet1',
index=False)
```

In this example, we create a dataframe containing data and use the to_excel() method to write the data to an Excel file named data.xlsx. We use the ExcelWriter() function to create an Excel writer object and specify the sheet name and index to be written to the file.

Writing CSV Files:

CSV (Comma-Separated Values) files are widely used for storing and exchanging data. In SAS, PROC EXPORT can be used to export data to CSV format. Similarly, in Python, the csv module provides built-in functions to read and write CSV files.

Here's an example code snippet that demonstrates how to write data to a CSV file using Python:

```
import csv
# create a list of dictionaries containing data
data = [
        {'Name': 'John', 'Age': 25, 'Gender': 'Male'},
        {'Name': 'Mary', 'Age': 30, 'Gender': 'Female'},
```



```
{'Name': 'Peter', 'Age': 35, 'Gender': 'Male'}
]
# write data to the CSV file
with open('data.csv', 'w', newline='') as file:
    writer = csv.DictWriter(file, fieldnames=['Name',
    'Age', 'Gender'])
    writer.writeheader()
    writer.writerows(data)
```

In this example, we create a list of dictionaries containing data and use the csv.DictWriter() function to write the data to a CSV file named data.csv. We specify the field names using the fieldnames parameter, and use the writerow() method to write each row of data to the file.

Writing Text Files:

Text files are files that contain plain text data. In SAS, PROC EXPORT can be used to export data to text format. Similarly, in Python, the open() function can be used to write data to text files.

Here's an example code snippet that demonstrates how to write data to a text file using Python:

```
# create a string containing data
data = 'John,25,Male\nMary,30,Female\nPeter,35,Male\n'
# write data to the text file
with open('data.txt', 'w') as file:
    file.write(data)
```

In this example, we create a string containing data and use the file.write() method to write the data to a text file named data.txt. The \n character is used to indicate a new line in the text file.

Writing JSON Files:

JSON (JavaScript Object Notation) is a lightweight data interchange format. In SAS, PROC EXPORT can be used to export data to JSON format. Similarly, in Python, the json module provides built-in functions to read and write JSON files.

Here's an example code snippet that demonstrates how to write data to a JSON file using Python:

```
import json
# create a list of dictionaries containing data
data = [
        {'Name': 'John', 'Age': 25, 'Gender': 'Male'},
        {'Name': 'Mary', 'Age': 30, 'Gender': 'Female'},
        {'Name': 'Peter', 'Age': 35, 'Gender': 'Male'}
```



```
]
# write data to the JSON file
with open('data.json', 'w') as file:
    json.dump(data, file)
```

In this example, we create a list of dictionaries containing data and use the json.dump() function to write the data to a JSON file named data.json.

Writing XML Files:

XML (eXtensible Markup Language) is a markup language used for storing and exchanging data. In SAS, PROC EXPORT can be used to export data to XML format. Similarly, in Python, the xml.etree.ElementTree module provides built-in functions to write XML files.

Here's an example code snippet that demonstrates how to write data to an XML file using Python:

```
import xml.etree.ElementTree as ET
# create an XML element containing data
root = ET.Element('data')
for i in range(3):
    item = ET.SubElement(root, 'item')
    ET.SubElement(item, 'Name').text = ['John', 'Mary',
'Peter'][i]
    ET.SubElement(item, 'Age').text = str([25, 30,
35][i])
    ET.SubElement(item, 'Gender').text = ['Male',
'Female', 'Male'][i]
# write data to the XML file
tree = ET.ElementTree(root)
tree.write('data.xml')
```

In this example, we create an XML element containing data and use the ElementTree.write() method to write the data to an XML file named data.xml.

In this article, we have discussed various techniques for writing data to files using Python from a SAS user's perspective. We have covered techniques for writing SAS data files, Excel files, CSV files, text files, JSON files, and XML files.

Python provides a wide range of libraries and functions to read and write data in different file formats. These techniques can be easily integrated with SAS workflows to enhance data analysis and processing capabilities.



```
# Import necessary libraries
import pandas as pd
import xlsxwriter
# Create a sample data frame
df = pd.DataFrame({
    'Name': ['John', 'Mary', 'Peter'],
    'Age': [25, 30, 35],
    'Gender': ['Male', 'Female', 'Male']
})
# Define the file path and name
file path = 'example.xlsx'
# Define the Excel writer
writer = pd.ExcelWriter(file path, engine='xlsxwriter')
# Write the data to a sheet named 'Sheet1'
df.to excel(writer, sheet name='Sheet1', index=False)
# Access the underlying workbook and worksheet objects
workbook = writer.book
worksheet = writer.sheets['Sheet1']
# Format the headers
header format = workbook.add format({
    'bold': True,
    'text wrap': True,
    'valign': 'top',
    'fg color': '#D7E4BC',
    'border': 1
})
# Apply the header format
for col num, value in enumerate(df.columns.values):
    worksheet.write(0, col num, value, header format)
# Format the data cells
data format = workbook.add format({
    'text wrap': True,
    'valign': 'top',
    'border': 1
})
```

```
in stal
```

```
# Apply the data format
for row_num in range(1, len(df) + 1):
    for col_num, value in enumerate(df.iloc[row_num -
1]):
        worksheet.write(row_num, col_num, value,
data_format)
# Save the Excel file
writer.save()
# Print a success message
print('Data successfully written to Excel file:',
file_path)
```

This code imports the necessary libraries and creates a sample pandas data frame. It then defines the file path and name for the Excel file and creates an Excel writer object. The data is written to a sheet named 'Sheet1' using the to_excel() method of the data frame.

```
# import required modules
import pandas as pd
import openpyxl
import csv
import json
import xml.etree.ElementTree as ET
from sas7bdat import SAS7BDAT
# create a sample data frame
data = pd.DataFrame({
    'Name': ['John', 'Mary', 'Peter'],
    'Age': [25, 30, 35],
    'Gender': ['Male', 'Female', 'Male']
})
# write data to SAS data file
with SAS7BDAT('data.sas7bdat', 'w') as file:
    file.write(data)
# write data to Excel file
writer = pd.ExcelWriter('data.xlsx', engine='openpyxl')
data.to excel(writer, index=False)
writer.save()
# write data to CSV file
```



```
data.to csv('data.csv', index=False)
# write data to text file
with open('data.txt', 'w') as file:
    file.write(data.to string(index=False))
# write data to JSON file
with open('data.json', 'w') as file:
    json.dump(data.to dict(orient='records'), file)
# write data to XML file
root = ET.Element('data')
for i in range(len(data)):
    item = ET.SubElement(root, 'item')
    ET.SubElement(item, 'Name').text = data['Name'][i]
    ET.SubElement(item, 'Age').text =
str(data['Age'][i])
    ET.SubElement(item, 'Gender').text =
data['Gender'][i]
tree = ET.ElementTree(root)
tree.write('data.xml')
```

In this example, we first create a sample data frame using Pandas. We then use different techniques to write this data to different file formats such as SAS data files, Excel files, CSV files, text files, JSON files, and XML files.

This code can be easily modified to write data in different formats or with different settings. It provides a good starting point for SAS users who want to learn how to use Python to write data to files.

CSV files

CSV (Comma Separated Values) files are a popular file format for storing tabular data. They are widely used in various applications and can be easily imported into Python for data analysis. In this article, we will discuss how to work with CSV files in Python, especially from the perspective of SAS users who are new to Python.

Python has a built-in csv module that provides functionality for reading and writing CSV files. The csv module provides several methods to read and write CSV files. The reader() method reads the contents of a CSV file and returns an object that can be iterated over to access each row. The writer() method writes data to a CSV file.



Here is an example of reading a CSV file in Python:

```
import csv
with open('file.csv', 'r') as csvfile:
    csvreader = csv.reader(csvfile)
    for row in csvreader:
        print(row)
```

In the above example, we first import the csv module. We then open the CSV file in read mode using the open() function and the csv.reader() method reads the contents of the file. We then loop through the contents of the CSV file and print each row.

If we want to write data to a CSV file, we can use the csv.writer() method as shown in the following example:

import csv

```
with open('file.csv', 'w') as csvfile:
    csvwriter = csv.writer(csvfile)
    csvwriter.writerow(['Name', 'Age', 'City'])
    csvwriter.writerow(['John', '25', 'New York'])
    csvwriter.writerow(['Alice', '30', 'Chicago'])
```

In the above example, we first open the CSV file in write mode using the open() function. We then create a csv.writer() object and use the writerow() method to write data to the file. In this case, we write three rows to the file, where the first row contains the column headers and the next two rows contain data.

In addition to the csv module, there are several third-party libraries that can be used to read and write CSV files in Python. One such library is pandas, which is a powerful library for data analysis. The pandas library provides a read_csv() function that can be used to read CSV files into a DataFrame object, which is a two-dimensional table-like data structure.

Here is an example of reading a CSV file using pandas:

```
import pandas as pd
df = pd.read_csv('file.csv')
print(df)
```

In the above example, we first import the pandas library. We then use the read_csv() function to read the CSV file into a DataFrame object. We can then print the contents of the DataFrame object using the print() function.



CSV files are a popular file format for storing tabular data and can be easily read and written in Python using the csv module or third-party libraries such as pandas. SAS users who are new to Python can quickly learn how to work with CSV files in Python using the examples provided in this article.

Reading CSV file into a dictionary:

```
import csv
with open('file.csv', 'r') as csvfile:
    csvreader = csv.DictReader(csvfile)
    for row in csvreader:
        print(row['Name'], row['Age'], row['City'])
```

In this example, we use the DictReader() method from the csv module to read the CSV file into a dictionary. Each row of the CSV file is represented as a dictionary where the keys correspond to the column headers. We can then access the values in each row using the dictionary keys.

Writing data to a CSV file using a list of dictionaries:

In this example, we create a list of dictionaries containing the data that we want to write to the CSV file. We then use the DictWriter() method to write the data to the CSV file. The writeheader() method writes the column headers to the CSV file. We then loop through each dictionary in the list and write each row to the CSV file using the writerow() method.

Working with CSV files using pandas:

```
import pandas as pd
```



```
df = pd.read_csv('file.csv')
df['Age'] = df['Age'] + 1
df.to_csv('file_updated.csv', index=False)
```

In this example, we use the read_csv() function from the pandas library to read the CSV file into a DataFrame object. We can then manipulate the data in the DataFrame object as desired.

Here are some more examples of working with CSV files in Python:

Example 1: Writing data to a CSV file using the csv module

```
import csv
data = [['Name', 'Age', 'City'],
      ['John', '25', 'New York'],
      ['Alice', '30', 'Chicago']]
with open('file.csv', 'w') as csvfile:
      csvwriter = csv.writer(csvfile)
      csvwriter.writerows(data)
```

In the above example, we create a nested list data containing the data we want to write to the CSV file. We then open the file in write mode using the open() function, create a csv.writer() object, and use the writerows() method to write the data to the file. The writerows() method writes multiple rows to the file at once.

Example 2: Reading a CSV file using the csv module and skipping header row

```
import csv
with open('file.csv', 'r') as csvfile:
    csvreader = csv.reader(csvfile)
    header = next(csvreader)
    for row in csvreader:
        print(row)
```

In the above example, we open the CSV file in read mode, create a csv.reader() object, and use the next() function to skip the first row (which contains the header). We then loop through the remaining rows and print each row.

Example 3: Reading a CSV file using pandas and filtering data based on a condition

```
import pandas as pd
```



```
df = pd.read_csv('file.csv')
df_filtered = df[df['Age'] > 25]
print(df_filtered)
```

In the above example, we read the CSV file into a DataFrame object using the read_csv() function from pandas. We then filter the data based on a condition (age greater than 25) and store the filtered data in a new DataFrame object df_filtered. Finally, we print the contents of the filtered DataFrame object.

Example 4: Writing data to a CSV file using pandas

```
import pandas as pd
data = {'Name': ['John', 'Alice'],
            'Age': [25, 30],
            'City': ['New York', 'Chicago']}
df = pd.DataFrame(data)
df.to_csv('file.csv', index=False)
```

In the above example, we create a dictionary data containing the data we want to write to the CSV file. We then create a DataFrame object from the dictionary using the pd.DataFrame() function. Finally, we use the to_csv() method of the DataFrame object to write the data to a CSV file, with the index parameter set to False to exclude the row numbers from the output.

Example 5: Reading a CSV file using the csv module and handling missing data

```
import csv
with open('file.csv', 'r') as csvfile:
    csvreader = csv.reader(csvfile)
    header = next(csvreader)
    for row in csvreader:
        name = row[0]
        age = row[1] if row[1] else 'N/A'
        city = row[2] if row[2] else 'N/A'
        print(f"Name: {name}, Age: {age}, City:
    {city}")
```

In the above example, we open the CSV file in read mode, create a csv.reader() object, and use the next() function to skip the header row. We then loop through the remaining rows and extract the data from each row. If a value is missing (represented by an empty string), we replace it with the string 'N/A'. Finally, we print the data in a formatted string.



Example 6: Writing data to a CSV file using the csv module and custom delimiter

```
import csv
data = [['Name', 'Age', 'City'],
      ['John', '25', 'New York'],
      ['Alice', '30', 'Chicago']]
with open('file.txt', 'w') as csvfile:
      csvwriter = csv.writer(csvfile, delimiter='\t')
      csvwriter.writerows(data)
```

In the above example, we create a nested list data containing the data we want to write to the CSV file. We then open the file in write mode using the open() function, create a csv.writer() object, and set the delimiter to '\t' to use a tab character as the delimiter. Finally, we use the writerows() method to write the data to the file.

Example 7: Reading a CSV file using pandas and selecting columns

```
import pandas as pd
df = pd.read_csv('file.csv')
df_selected = df[['Name', 'Age']]
print(df_selected)
```

In the above example, we read the CSV file into a DataFrame object using the read_csv() function from pandas. We then select the columns Name and Age from the DataFrame object and store them in a new DataFrame object df_selected. Finally, we print the contents of the selected DataFrame object.

Example 8: Writing data to a CSV file using pandas and custom header

```
import pandas as pd
data = {'Name': ['John', 'Alice'],
                                'Age': [25, 30],
                             'City': ['New York', 'Chicago']}
df = pd.DataFrame(data)
header = ['Person Name', 'Age', 'Residence']
df.to csv('file.csv', header=header, index=False)
```

In the above example, we create a dictionary data containing the data we want to write to the CSV file. We then create a DataFrame object from the dictionary using the pd.DataFrame()



function. We also create a list header containing the custom header names. Finally, we use the to_csv() method of the DataFrame object to write the data to a CSV file, with the header parameter set to the custom header list and the index parameter set to False to exclude the row numbers from the output.

Example 9: Reading a CSV file with custom delimiter and quoting characters using the csv module

```
with open('file.csv', 'r') as csvfile:
    csvreader = csv.reader(csvfile, delimiter='|',
    quotechar='"')
    for row in csvreader:
        name = row[0]
        age = row[1]
        city = row[2]
        print(f"Name: {name}, Age: {age}, City:
    {city}")
```

In the above example, we open the CSV file in read mode, create a csv.reader() object, and set the delimiter to '|' and the quote character to ''' to handle the custom formatting of the file. We then loop through the rows of the file, extract the data from each row, and print it in a formatted string.

Example 10: Writing data to a CSV file using csv.DictWriter class

In the above example, we create a list of dictionaries data containing the data we want to write to the CSV file. We then open the file in write mode using the open() function, create a csv.DictWriter() object, and set the field names to ['Name', 'Age', 'City']. We use the writeheader() method to write the header row to the file and then loop through the rows of data, using the writerow() method to write each row to the file.



Example 11: Reading a CSV file using numpy and filtering rows based on a condition

```
import numpy as np
data = np.genfromtxt('file.csv', delimiter=',',
dtype=None, names=True)
selected_rows = data[data['Age'] > 25]
print(selected_rows)
```

In the above example, we use the genfromtxt() function from numpy to read the CSV file into a structured array. We set the delimiter to ',' the dtype parameter to None to automatically determine the data types of the columns, and the names parameter to True to use the header row as the column names. We then filter the rows of the array based on the condition data['Age'] > 25, which selects only the rows where the value in the Age column is greater than 25. Finally, we print the selected rows.

Example 12: Writing data to a CSV file using numpy and custom formatting

```
import numpy as np
data = np.array([('John', 25, 'New York'), ('Alice',
30, 'Chicago')], dtype=[('Name', 'U10'), ('Age', 'i4'),
('City', 'U10')])
np.savetxt('file.csv', data, delimiter=',', fmt='%s,
%d, %s', header='Name, Age, City', comments='')
```

In the above example, we create a structured array data containing the data we want to write to the CSV file. We set the dtype parameter to specify the data types of the columns and use the `U10

Excel files

Here is an example of Python code that might be included in the book for reading in a CSV file using Pandas:

```
import pandas as pd
# read in data from a CSV file
data = pd.read_csv('mydata.csv')
# print out the first few rows of data
print(data.head())
```



In addition to basic data manipulation, the book might also include examples of more advanced topics such as machine learning using Scikit-Learn, web scraping using libraries like Beautiful Soup, and data visualization using Matplotlib or Seaborn.

Here is an example of Python code that might be included in the book for performing machine learning using Scikit-Learn:

```
from sklearn.datasets import load iris
from sklearn.model selection import train test split
from sklearn.tree import DecisionTreeClassifier
# load the iris dataset
iris = load iris()
# split the data into training and test sets
X train, X test, y train, y test = train test split(
    iris.data, iris.target, test size=0.3,
random state=42)
# create a decision tree classifier
clf = DecisionTreeClassifier()
# fit the model to the training data
clf.fit(X train, y train)
# make predictions on the test data
y pred = clf.predict(X test)
# print out the accuracy of the model
print("Accuracy:", clf.score(X test, y test))
```

The code examples in the book are likely to cover a range of topics, including:

- Basic Python programming concepts, such as variables, data types, operators, and control structures.
- Python libraries for data analysis, such as NumPy, Pandas, and Matplotlib.
- Python libraries for machine learning, such as Scikit-Learn and TensorFlow.
- Web scraping using Python libraries such as Beautiful Soup and Requests.

Best practices for Python programming, such as code organization, documentation, testing, and version control.

Here's an example of what a code snippet from the book might look like:

Importing the Pandas library



```
import pandas as pd
# Reading in a CSV file using Pandas
data = pd.read_csv('data.csv')
# Displaying the first 10 rows of the data
print(data.head(10))
# Calculating summary statistics for the data
summary = data.describe()
print(summary)
# Visualizing the data using Matplotlib
import matplotlib.pyplot as plt
plt.plot(data['x'], data['y'])
plt.xlabel('X-axis')
plt.ylabel('Y-axis')
plt.title('Scatter plot of x and y')
plt.show()
```

This code example demonstrates how to use Pandas to read in a CSV file, display the first 10 rows of the data, and calculate summary statistics. It also shows how to use Matplotlib to create a scatter plot of two columns of data.

Here are some more examples of code snippets that might be included in "Python for SAS Users: A SAS-Oriented Introduction to Python":

Example of using Python's built-in functions to manipulate strings:

```
# Creating a string variable
s = 'Hello, World!'
# Converting the string to uppercase
s_upper = s.upper()
# Replacing a substring in the string
s_new = s.replace('World', 'Universe')
# Splitting the string into a list
s_list = s.split(',')
# Printing the results
print(s_upper)
print(s_new)
```



print(s_list)

Example of using Python's NumPy library to perform basic arithmetic operations on arrays:

```
# Importing the NumPy library
import numpy as np
# Creating two arrays
a = np.array([1, 2, 3])
b = np.array([4, 5, 6])
# Adding the two arrays
c = a + b
# Multiplying the two arrays
d = a * b
# Computing the dot product of the two arrays
e = np.dot(a, b)
# Printing the results
print(c)
print(d)
print(e)
```

Example of using Python's Scikit-Learn library to perform machine learning:

```
# Importing the Scikit-Learn library
from sklearn.linear_model import LinearRegression
# Creating a simple dataset
X = [[0], [1], [2], [3], [4], [5]]
y = [0, 1, 2, 3, 4, 5]
# Creating a LinearRegression model
model = LinearRegression()
# Fitting the model to the data
model.fit(X, y)
# Predicting values using the model
y_pred = model.predict([[6], [7], [8]])
```



Printing the results print(y_pred)

These code examples demonstrate some of the key features of Python, including string manipulation, array operations, and machine learning. The book "Python for SAS Users: A SAS-Oriented Introduction to Python" is likely to include many more examples like these, along with detailed explanations and instructions for how to use them in practice.

Example of using Python's Pandas library to clean and manipulate data: python

Copy code # Importing the Pandas library import pandas as pd

Reading in a CSV file
data = pd.read_csv('data.csv')

Dropping rows with missing values
data_clean = data.dropna()

Renaming columns
data_clean = data_clean.rename(columns={'old_name': 'new_name'})

Converting a column to a different data type data_clean['column_name'] = data_clean['column_name'].astype(int)

```
# Grouping the data by a column and computing summary statistics
summary = data_clean.groupby('grouping_variable').agg({'column1': 'mean', 'column2': 'sum'})
```

Exporting the cleaned data to a new CSV file data_clean.to_csv('data_clean.csv', index=False) Example of using Python's Requests and Beautiful Soup libraries to scrape data from a website: python Copy code # Importing the Requests and Beautiful Soup libraries import requests from bs4 import BeautifulSoup

Sending a GET request to a website
url = 'https://www.example.com'
response = requests.get(url)

Parsing the HTML content using Beautiful Soup soup = BeautifulSoup(response.content, 'html.parser')

Extracting specific elements from the HTML content title = soup.title.string



```
links = [link['href'] for link in soup.find_all('a')]
```

Printing the results
print(title)
print(links)
Example of using Python's unittest library to test a function:
ruby
Copy code
Defining a simple function to test
def add_numbers(a, b):
 return a + b

Importing the unittest library import unittest

```
# Creating a test case class
class TestAddNumbers(unittest.TestCase):
    def test_add_numbers(self):
        self.assertEqual(add_numbers(2, 3), 5)
        self.assertEqual(add_numbers(-1, 1), 0)
        self.assertEqual(add_numbers(0, 0), 0)
```

Running the test case

if __name__ == '__main__':
 unittest.main()

JSON files

JSON (JavaScript Object Notation) is a lightweight data interchange format that is easy for humans to read and write, and easy for machines to parse and generate. It is commonly used for transmitting data between a server and a web application, as an alternative to XML. JSON data is represented as key-value pairs, where the keys are strings and the values can be any valid JSON data type, including objects, arrays, numbers, strings, booleans, and null.

Python is a popular programming language that is widely used in data analysis and scientific computing. It has a rich ecosystem of libraries and tools that make it an ideal language for working with data. In recent years, Python has become increasingly popular among SAS users as well, due to its ease of use, flexibility, and ability to handle large amounts of data.

Python has built-in support for working with JSON data, allowing users to easily parse and generate JSON files. In this context, a JSON file is simply a text file that contains JSON-encoded data. Python provides the json module for working with JSON data. This module provides a simple interface for parsing JSON data, and converting Python objects to JSON and vice versa.



Here are some examples of using Python to work with JSON data:

Parsing JSON data:

The json module provides a function called json.load() that can be used to parse a JSON file into a Python object. Here's an example:

```
import json
# Open the JSON file
with open('data.json', 'r') as f:
    # Parse the JSON data into a Python object
    data = json.load(f)
# Print the data
print(data)
```

Generating JSON data:

The json module provides a function called json.dump() that can be used to convert a Python object into a JSON-encoded string, and write it to a file. Here's an example:

```
import json
# Create a Python object
data = {
    'name': 'John',
    'age': 30,
    'city': 'New York'
}
# Convert the Python object to a JSON-encoded string
json_string = json.dumps(data)
# Write the JSON string to a file
with open('data.json', 'w') as f:
    f.write(json_string)
```

Working with JSON data in pandas:

Pandas is a popular library for data manipulation and analysis in Python. It provides built-in support for reading and writing JSON data. Here's an example:

```
import pandas as pd
# Read the JSON data into a pandas DataFrame
```



```
df = pd.read_json('data.json')
# Print the DataFrame
print(df)
# Write the DataFrame to a JSON file
df.to_json('data.json')
```

Python provides a powerful and easy-to-use interface for working with JSON data. SAS users who are familiar with Python can leverage this capability to work with JSON files and integrate them with their SAS workflows.

Manipulating JSON data in Python:

Once you have loaded JSON data into a Python object, you can manipulate it using standard Python data structures and functions. For example, you can access values in a JSON object using dictionary-like syntax:

```
import json
# Load JSON data into a Python object
data = json.loads('{"name": "John", "age": 30}')
# Access values in the object
print(data['name']) # Output: "John"
print(data['age']) # Output: 30
```

You can also iterate over values in a JSON array:

```
import json
# Load JSON data into a Python object
data = json.loads('["apple", "banana", "cherry"]')
# Iterate over values in the array
for item in data:
    print(item)
```

Handling errors when working with JSON:

When working with JSON data in Python, it's important to handle errors that may occur during parsing or encoding. The json module provides several functions for doing this, including json.JSONDecodeError and json.JSONEncoder. Here's an example of how to handle a JSON decoding error:

import json



```
# Load JSON data into a Python object
try:
    data = json.loads('{"name": "John, "age": 30}')
except json.JSONDecodeError:
    print("Error: Invalid JSON data")
```

Advanced JSON manipulation in Python:

In addition to the basic functionality provided by the json module, there are several third-party libraries available for working with JSON data in Python. For example, the jsonschema library provides tools for validating JSON data against a schema, while the jmespath library provides a powerful query language for JSON data.

Reading JSON data from an API:

In addition to reading JSON data from a file, Python can also be used to read JSON data from an API. For example, you can use the requests library to send an HTTP GET request to an API endpoint and retrieve JSON data as a response:

```
import requests
import json
# Send an HTTP GET request to an API endpoint
response =
requests.get('https://jsonplaceholder.typicode.com/todo
s')
# Check that the request was successful (HTTP status
code 200)
if response.status code == 200:
    # Convert the response content (JSON-encoded
string) to a Python object
    data = json.loads(response.content)
    # Print the first item in the array
   print(data[0])
else:
   print("Error: Failed to retrieve data")
```

Writing JSON data to a MongoDB database:

Python can also be used to write JSON data to a MongoDB database. MongoDB is a popular NoSQL database that uses JSON-like documents for storing data. The PyMongo library provides a Python interface for interacting with MongoDB databases:

import pymongo
import json



```
# Connect to a MongoDB database
client =
pymongo.MongoClient("mongodb://localhost:27017/")
db = client["mydatabase"]
# Load JSON data from a file
with open('data.json', 'r') as f:
    data = json.load(f)
# Insert the data into a MongoDB collection
collection = db["mycollection"]
collection.insert_one(data)
```

Validating JSON data with a schema:

The jsonschema library provides tools for validating JSON data against a schema. A JSON schema is a document that defines the structure and constraints of a JSON object. Here's an example of how to use jsonschema to validate JSON data:

```
import json
import jsonschema
# Define a JSON schema
schema = \{
    "type": "object",
    "properties": {
        "name": {"type": "string"},
        "age": {"type": "integer"}
    },
    "required": ["name", "age"]
}
# Load JSON data into a Python object
data = json.loads('{"name": "John", "age": 30}')
# Validate the data against the schema
try:
    jsonschema.validate(data, schema)
    print("Data is valid")
except jsonschema.ValidationError as e:
   print("Error: Data is not valid")
   print(e)
```



Querying JSON data with JMESPath:

The jmespath library provides a powerful query language for JSON data. With JMESPath, you can extract data from a JSON object using a simple and intuitive syntax. Here's an example:

```
import json
import jmespath
# Load JSON data into a Python object
data = json.loads('{"name": {"first": "John", "last":
"Doe"}, "age": 30}')
# Extract the first name using a JMESPath expression
expression = "name.first"
result = jmespath.search(expression, data)
print(result) # Output: "John"
```

These are just a few examples of the many ways in which Python can be used to work with JSON data. By leveraging Python's rich ecosystem of libraries and tools, you can easily parse, generate, manipulate, and validate JSON data as part of your data analysis and visualization workflows.

Converting Python objects to JSON:

In addition to loading JSON data into Python objects, you can also convert Python objects to JSON format using the json.dumps() function. This is useful when you want to store data in a JSON file or send it as a response to an API request. Here's an example:

```
import json
# Create a Python dictionary
data = {
    "name": "John",
    "age": 30,
    "isMarried": True,
    "hobbies": ["reading", "hiking", "photography"]
}
# Convert the dictionary to a JSON-encoded string
json_data = json.dumps(data)
print(json_data) # Output: '{"name": "John", "age":
30, "isMarried": true, "hobbies": ["reading", "hiking",
    "photography"]}'
```

Converting JSON to Python objects using custom decoders:



Sometimes the JSON data you're working with may have custom data types that aren't natively supported by Python. In this case, you can define a custom decoder function to convert the JSON data to Python objects. Here's an example:

```
import json
# Define a custom decoder function
def decode person(data):
    if "name" in data and "age" in data:
        return Person(data["name"], data["age"])
    else:
        return data
# Define a Person class
class Person:
    def init (self, name, age):
        self.name = name
        self.age = age
# Load JSON data into a Python object using the custom
decoder function
json data = '{"person": {"name": "John", "age": 30}}'
data = json.loads(json data, object hook=decode person)
print(type(data)) # Output: <class 'dict'>
print(type(data["person"])) # Output: <class</pre>
' main .Person'>
```

Handling large JSON files:

When working with large JSON files, you may encounter memory issues if you try to load the entire file into memory at once. One way to avoid this is to use an iterative parser like ijson, which allows you to parse the file incrementally without loading it all into memory at once. Here's an example:

```
import ijson
# Open a large JSON file
with open("data.json", "r") as f:
    # Create an ijson iterator
    objects = ijson.items(f, "item")

# Iterate over each object in the file
for obj in objects:
    # Process the object
```



print(obj)

In this example, ijson.items() creates an iterator that returns each object in the "item" array of the JSON file one at a time. This allows you to process the file incrementally without loading it all into memory at once.

Combining JSON data from multiple files:

In some cases, you may have JSON data spread across multiple files that you want to combine into a single data structure. One way to do this is to use the glob module to find all the files that match a certain pattern, and then loop through each file and load the JSON data into a single list or dictionary. Here's an example:

```
import glob
import json
# Find all JSON files in a directory
file pattern = "*.json"
file paths = glob.glob(file pattern)
# Load the JSON data from each file into a list
data list = []
for path in file paths:
   with open (path, "r") as f:
        data = json.load(f)
        data list.append(data)
# Combine the data from each file into a single
dictionary
combined data = {}
for data in data list:
    for key, value in data.items():
        if key in combined data:
            combined data[key].extend(value)
        else:
            combined data[key] = value
```

```
print(combined_data)
```

In this example, we use glob.glob() to find all files in the current directory that match the "*.json" pattern. We then loop through each file, load the JSON data using json.load(), and append it to a list called data_list. Finally, we loop through the data_list and combine the data from each file into a single dictionary called combined_data. Validating JSON data:



When working with JSON data, it's important to ensure that it's valid and conforms to a certain structure. One way to do this is to use the jsonschema library, which allows you to define a JSON schema that specifies the structure and data types of your JSON data. You can then use this schema to validate your JSON data and ensure that it meets your requirements. Here's an example:

```
import jsonschema
import json
# Define a JSON schema
schema = {
    "type": "object",
    "properties": {
        "name": {"type": "string"},
        "age": {"type": "integer"}]
```

Writing JSON data to a file:

In addition to loading JSON data from a file, you can also write JSON data to a file using the json.dump() function. Here's an example:

```
import json
# Create a Python dictionary
data = \{
    "name": "John",
    "age": 30,
    "isMarried": True,
    "hobbies": ["reading", "hiking", "photography"]
}
# Write the dictionary to a JSON file
with open("data.json", "w") as f:
    json.dump(data, f)
# Read the JSON data from the file
with open("data.json", "r") as f:
    json data = json.load(f)
print(json data) # Output: {'name': 'John', 'age': 30,
'isMarried': True, 'hobbies': ['reading', 'hiking',
'photography']}
```



In this example, the json.dump() function writes the data dictionary to a file called "data.json". The with open() statement automatically closes the file when the block is finished. The json.load() function is then used to read the data from the file.

Working with databases

Python is a powerful programming language that is commonly used in data science and data analysis tasks. For SAS users who are familiar with SAS software, Python offers a new set of tools and capabilities that can be used to manipulate, analyze, and visualize data.

One of the most common tasks in data analysis is working with databases. In this article, we will explore how to work with databases in Python, with a focus on the needs of SAS users.

To work with databases in Python, we need to use a Python library called "pandas". Pandas is a library that provides a set of tools for working with data in Python, including tools for working with databases.

The first step in working with databases in Python is to establish a connection to the database. This can be done using the "pyodbc" library, which is a Python library for connecting to databases using ODBC (Open Database Connectivity).

Once a connection to the database has been established, we can use the pandas library to read data from the database and manipulate it in Python. For example, we can use the pandas "read_sql" function to read data from a database table and create a pandas DataFrame object.

Once we have a pandas DataFrame object, we can use pandas functions to manipulate the data. For example, we can use the "groupby" function to group the data by a particular column, or the "merge" function to combine data from multiple tables.

In addition to reading data from a database, we can also write data to a database using pandas. For example, we can use the "to_sql" function to write a pandas DataFrame to a database table.

Overall, working with databases in Python can provide SAS users with a new set of tools and capabilities for working with data. By using the pandas library and other Python tools, SAS users can access and manipulate data in a more flexible and powerful way than is possible with SAS software alone.

Another useful Python library for working with databases is "SQLAlchemy". SQLAlchemy provides a set of tools for working with databases in a more flexible and powerful way than is possible with standard SQL.

With SQLAlchemy, we can create an "engine" object that connects to a database, and then use the engine object to execute SQL commands or interact with the database using Python objects.



For example, we can use SQLAlchemy to create a database schema, or to map database tables to Python classes. This allows us to interact with the database using Python objects, which can be more intuitive and easier to work with than SQL commands.

Another advantage of using SQLAlchemy is that it provides a consistent API for working with different types of databases. This means that we can write Python code that works with a variety of different databases, without having to learn the specific SQL dialect of each database.

Overall, using Python for working with databases offers a number of advantages over using SAS software alone. Python provides a powerful set of tools for data analysis and manipulation, and can be used to work with a variety of different types of databases.

For SAS users who are new to Python, there may be a learning curve involved in getting started. However, once the basics of Python and pandas have been learned, working with databases in Python can be a powerful addition to a data analysis toolkit.

Another important aspect of working with databases in Python is managing data security. In SAS, users often work with data in a secure environment, where access to data is strictly controlled. The same level of security is also important when working with databases in Python.

To ensure data security when working with databases in Python, it is important to follow best practices for authentication and encryption. For example, using encrypted connections and secure authentication methods such as OAuth can help protect sensitive data.

It is also important to carefully manage database credentials, and to limit access to databases to only those who need it. This can be done by using strong passwords, and by limiting access to databases to specific IP addresses or network domains.

In addition to data security, performance is also an important consideration when working with databases in Python. In some cases, reading and writing large amounts of data from a database using pandas can be slow. To optimize performance, it is important to carefully manage the amount of data being read or written at any one time, and to use SQL queries or stored procedures to perform data manipulation directly on the database server where possible.

Finally, it is important to keep in mind that working with databases in Python is just one aspect of a larger data analysis workflow. To effectively analyze data, it is important to have a strong understanding of statistical methods and data visualization techniques, as well as a clear understanding of the underlying data structures and business goals. By combining SAS and Python tools in a thoughtful and strategic way, data analysts can create a powerful and flexible data analysis toolkit that can be tailored to their specific needs and workflows.

Connecting to a database using pyodbc:

import pyodbc

establish a connection to the database



```
conn = pyodbc.connect('DRIVER={ODBC
Driver};SERVER=servername;DATABASE=databasename;UID=use
rname;PWD=password')
# create a cursor object for executing SQL commands
cursor = conn.cursor()
# execute a SQL command to retrieve data
cursor.execute('SELECT * FROM tablename')
# fetch the results and create a pandas DataFrame
data = cursor.fetchall()
df = pd.DataFrame(data)
```

Reading data from a database using pandas:

```
import pandas as pd
import pyodbc
# establish a connection to the database
conn = pyodbc.connect('DRIVER={ODBC
Driver};SERVER=servername;DATABASE=databasename;UID=use
rname;PWD=password')
# read data from a database table using pandas
df = pd.read_sql('SELECT * FROM tablename', conn)
```

Writing data to a database using pandas:

```
import pandas as pd
import pyodbc
# establish a connection to the database
conn = pyodbc.connect('DRIVER={ODBC
Driver};SERVER=servername;DATABASE=databasename;UID=use
rname;PWD=password')
# create a pandas DataFrame with data to write to the
database
data = {'column1': [1, 2, 3], 'column2': ['a', 'b',
'c']}
df = pd.DataFrame(data)
# write the data to a database table using pandas
```



```
df.to sql('tablename', conn, if exists='replace')
import pyodbc
import pandas as pd
# Define connection parameters
server name = 'my server name'
database name = 'my database name'
username = 'my username'
password = 'my password'
# Create connection string
connection string = f'DRIVER={{SQL
Server }; SERVER={server name}; DATABASE={database name};
UID={username}; PWD={password} '
# Connect to database
connection = pyodbc.connect(connection string)
# Read data from database using pandas
query = 'SELECT * FROM my table'
df = pd.read sql(query, connection)
# Manipulate data using pandas
df grouped = df.groupby('my column').sum()
# Write data to database using pandas
df grouped.to sql('my new table', connection,
if exists='replace')
```

In this example, we first define the connection parameters for the SQL Server database, including the server name, database name, username, and password. We then create a connection string using these parameters and connect to the database using pyodbc.

Next, we use the pandas "read_sql" function to read data from a table called "my_table" in the database. We then use pandas functions to manipulate the data by grouping it by a column called "my_column" and summing the values.

Finally, we use the pandas "to_sql" function to write the manipulated data to a new table called "my_new_table" in the database. The "if_exists='replace'" parameter specifies that any existing table with the same name should be replaced.

This is just one example of how to work with databases in Python using pandas and pyodbc. There are many other libraries and tools available for working with databases in Python,



including SQLAlchemy and psycopg2. The exact code used will depend on the specific database being used and the requirements of the analysis being performed.

Connecting to a database using pyodbc:

```
import pyodbc
# Establish a connection to the database using ODBC
conn = pyodbc.connect('DRIVER={SQL
Server};SERVER=myserver;DATABASE=mydatabase;UID=myusern
ame;PWD=mypassword')
# Create a cursor object to execute SQL commands
cursor = conn.cursor()
```

Reading data from a database table using pandas:

```
import pandas as pd
import pyodbc
# Establish a connection to the database using ODBC
conn = pyodbc.connect('DRIVER={SQL
Server};SERVER=myserver;DATABASE=mydatabase;UID=myusern
ame;PWD=mypassword')
# Read data from a database table using pandas
df = pd.read sql('SELECT * FROM mytable', conn)
```

Writing data to a database table using pandas:

```
import pandas as pd
import pyodbc
# Establish a connection to the database using ODBC
conn = pyodbc.connect('DRIVER={SQL
Server};SERVER=myserver;DATABASE=mydatabase;UID=myusern
ame;PWD=mypassword')
# Create a pandas DataFrame object
data = {'coll': [1, 2, 3], 'col2': ['a', 'b', 'c']}
df = pd.DataFrame(data)
# Write the DataFrame to a database table using pandas
```



```
df.to_sql('mytable', conn, if_exists='replace',
index=False)
```

Creating a database schema using SQLAlchemy:

```
from sqlalchemy import create engine, Column, Integer,
String
from sqlalchemy.ext.declarative import declarative base
# Create an engine object to connect to the database
engine =
create engine('postgresql://myusername:mypassword@myser
ver/mydatabase')
# Create a declarative base for defining database
tables
Base = declarative base()
# Define a database table using SQLAlchemy classes
class MyTable(Base):
     tablename = 'mytable'
    id = Column(Integer, primary key=True)
    name = Column(String)
    value = Column(Integer)
# Create the database schema using the declarative base
Base.metadata.create all(engine)
```

Querying data from a database table using SQLAlchemy:

```
from sqlalchemy import create_engine
import pandas as pd
# Create an engine object to connect to the database
engine =
create_engine('postgresql://myusername:mypassword@myser
ver/mydatabase')
# Execute a SQL query using the engine object and store
the results in a pandas DataFrame
df = pd.read_sql_query('SELECT * FROM mytable WHERE
value > 100', engine)
```

These are just a few examples of how to work with databases in Python using the pandas and SQLAlchemy libraries. There are many other tools and techniques available for working with databases in Python, depending on the specific needs of a data analysis project.

Reading Data from a SQL Server Database using pyodbc and Pandas:

```
import pyodbc
import pandas as pd
server = 'your server name'
database = 'your database name'
username = 'your username'
password = 'your password'
driver = '{ODBC Driver 17 for SQL Server}'
# Establish a connection to the database
cnxn =
pyodbc.connect('DRIVER='+driver+';SERVER='+server+';DAT
ABASE='+database+';UID='+username+';PWD='+ password)
# Define the SQL query to be executed
query = 'SELECT * FROM your table name'
# Execute the query and read the data into a pandas
dataframe
data = pd.read sql(query, cnxn)
# Print the first 10 rows of the data
print(data.head(10))
```

Writing Data to a MySQL Database using SQLAlchemy and Pandas:

```
import pandas as pd
from sqlalchemy import create_engine
# Define the database connection string
connection_string =
'mysql+pymysql://your_username:your_password@your_serve
r_name/your_database_name'
# Constructed to the second string to the second string
```

```
# Create a SQLAlchemy engine object to connect to the
database
engine = create_engine(connection_string)
```

```
# Define a pandas dataframe to be written to the
database
data = pd.DataFrame({
    'column1': [1, 2, 3],
    'column2': ['value1', 'value2', 'value3'],
    'column3': ['2022-01-01', '2022-01-02', '2022-01-
03']
})
# Write the data to a new table in the database
data.to_sql('your_new_table_name', con=engine,
if_exists='replace', index=False)
# Read the data from the new table and print it
query = 'SELECT * FROM your_new_table_name'
new_data = pd.read_sql(query, con=engine)
print(new_data.head())
```

Using SQLAlchemy to Create a Database Schema and Perform CRUD Operations:

```
from sqlalchemy import create engine, Column, Integer,
String
from sqlalchemy.orm import sessionmaker
from sqlalchemy.ext.declarative import declarative base
# Define the database connection string
connection string = 'sqlite:///your database name.db'
# Create a SQLAlchemy engine object to connect to the
database
engine = create engine(connection string)
# Define a base class for declarative models
Base = declarative base()
# Define a database model to represent a table in the
database
class User(Base):
     tablename = 'users'
    id = Column(Integer, primary key=True)
    name = Column(String)
    email = Column(String)
```

```
# Create the database schema
Base.metadata.create all(engine)
# Create a session to interact with the database
Session = sessionmaker(bind=engine)
session = Session()
# Insert a new user into the database
new user = User(name='John Doe',
email='john.doe@example.com')
session.add(new user)
session.commit()
# Update an existing user in the database
user = session.query(User).filter by(name='John
Doe').first()
user.email = 'john.doe.updated@example.com'
session.commit()
# Delete a user from the database
user = session.query(User).filter by(name='John
Doe').first()
session.delete(user)
session.commit()
# Print the remaining users in the database
users = session.query(User).all()
for user in users:
    print(user.name, user.email)
```

These code examples demonstrate some of the different ways that Python can be used to work with databases, using libraries like pyodbc, pandas, and SQLAlchemy.

Using the sqlite3 Module to Create a Database and Perform CRUD Operations:

```
import sqlite3
# Create a new database and connect to it
conn = sqlite3.connect('my_database.db')
# Create a new table in the database
conn.execute('CREATE TABLE users (id INTEGER PRIMARY
KEY, name TEXT, email TEXT)')
```



```
# Insert a new row into the table
conn.execute('INSERT INTO users (name, email) VALUES
(?, ?)', ('John Doe', 'john.doe@example.com'))
conn.commit()
# Update an existing row in the table
conn.execute('UPDATE users SET email=? WHERE name=?',
('john.doe.updated@example.com', 'John Doe'))
conn.commit()
# Delete a row from the table
conn.execute('DELETE FROM users WHERE name=?', ('John
Doe',))
conn.commit()
# Select all rows from the table and print them
rows = conn.execute('SELECT * FROM users')
for row in rows:
   print(row)
# Close the database connection
conn.close()
```

Using the MySQL Connector/Python Module to Connect to a MySQL Database and Perform CRUD Operations:

```
import mysql.connector
# Connect to the database
conn = mysql.connector.connect(
    host='your_server_name',
    user='your_username',
    password='your_password',
    database='your_database_name'
)
# Create a new cursor object to execute queries
cursor = conn.cursor()
# Create a new table in the database
cursor.execute('CREATE TABLE users (id INT
AUTO_INCREMENT PRIMARY KEY, name VARCHAR(255), email
VARCHAR(255))')
```

```
# Insert a new row into the table
sql = 'INSERT INTO users (name, email) VALUES (%s, %s)'
val = ('John Doe', 'john.doe@example.com')
cursor.execute(sql, val)
conn.commit()
# Update an existing row in the table
sql = 'UPDATE users SET email = %s WHERE name = %s'
val = ('john.doe.updated@example.com', 'John Doe')
cursor.execute(sql, val)
conn.commit()
# Delete a row from the table
sql = 'DELETE FROM users WHERE name = %s'
val = ('John Doe',)
cursor.execute(sql, val)
conn.commit()
# Select all rows from the table and print them
cursor.execute('SELECT * FROM users')
rows = cursor.fetchall()
```



Chapter 4: Data Manipulation with Pandas



Introduction:

A panda is a popular data manipulation library for Python. It provides high-performance, easyto-use data structures and data analysis tools. Pandas can be a great tool for SAS users who are looking to expand their data analysis skills in Python. In this article, we will provide an introduction to Python for SAS users, with a focus on data manipulation with Pandas.

Data Structures in Pandas:

The two primary data structures in Pandas are Series and DataFrame. A Series is a onedimensional array-like object that can hold any data type. It is similar to a SAS variable or a column in a SAS dataset. A DataFrame is a two-dimensional table of data with rows and columns. It is similar to a SAS dataset.

Creating a Series in Pandas:

To create a Series in Pandas, we can use the Series() function. We can pass a list, tuple, or dictionary to the Series() function to create a Series.

For example, to create a Series of integers from a list, we can use the following code:

```
import pandas as pd
my_list = [1, 2, 3, 4, 5]
my_series = pd.Series(my_list)
print(my_series)
```

Output:

```
0 1
1 2
2 3
3 4
4 5
dtype: int64
```

Creating a DataFrame in Pandas:

To create a DataFrame in Pandas, we can use the DataFrame() function. We can pass a dictionary, a list of dictionaries, or a numpy array to the DataFrame() function to create a DataFrame.

For example, to create a DataFrame from a dictionary, we can use the following code:



```
226 | P a g e
```

```
'Salary': [50000, 70000, 30000, 90000]}
my_dataframe = pd.DataFrame(my_dict)
print(my_dataframe)
```

Output:

	Name	Age	Salary
0	John	25	50000
1	Mary	32	70000
2	Peter	19	30000
3	Anne	47	90000

Reading and Writing Data:

Pandas provides functions to read and write data from and to various file formats, including CSV, Excel, SQL databases, and more. To read a CSV file into a Pandas DataFrame, we can use the read_csv() function. To write a Pandas DataFrame to a CSV file, we can use the to_csv() function.

For example, to read a CSV file into a DataFrame, we can use the following code:

```
import pandas as pd
my_dataframe = pd.read_csv('my_file.csv')
print(my_dataframe)
```

Output:

	Name	Age	Salary
0	John	25	50000
1	Mary	32	70000
2	Peter	19	30000
3	Anne	47	90000

To write a DataFrame to a CSV file, we can use the following code:

```
import pandas as pd
my_dataframe.to_csv('my_file.csv', index=False)
```

Data Manipulation with Pandas:

Pandas provides a wide range of functions for data manipulation, including filtering, sorting, grouping, joining, and more. In this section, we will cover some of the most commonly used functions.



Filtering Data:

To filter data in Pandas, we can use the boolean indexing feature. We can create a boolean expression to filter the rows that meet certain conditions. For example, to filter a DataFrame to only include rows where the Age column is greater than 30, we can use the following code:

```
import pandas as pd
my_dataframe = pd.read_csv('my_file.csv')
filtered_dataframe = my_dataframe[my_dataframe['Age'] >
30]
print(filtered_dataframe)
```

Output:

	Name	Age	Salary
1	Mary	32	70000
3	Anne	47	90000

Sorting Data:

To sort data in Pandas, we can use the sort_values() function. We can pass a column name or a list of column names to the sort_values() function to sort the DataFrame by those columns.

For example, to sort a DataFrame by the Salary column in descending order, we can use the following code:

```
import pandas as pd
my_dataframe = pd.read_csv('my_file.csv')
sorted_dataframe = my_dataframe.sort_values('Salary',
ascending=False)
print(sorted_dataframe)
```

Output:

	Name	Age	Salary
3	Anne	47	90000
1	Mary	32	70000
0	John	25	50000
2	Peter	19	30000

Grouping Data:

To group data in Pandas, we can use the groupby() function. We can pass a column name or a list of column names to the groupby() function to group the DataFrame by those columns. We



can then apply various aggregation functions to the grouped data, such as sum, mean, min, max, and more.

For example, to group a DataFrame by the Age column and calculate the mean Salary for each age group, we can use the following code:

```
import pandas as pd
my_dataframe = pd.read_csv('my_file.csv')
grouped_dataframe = my_dataframe.groupby('Age').mean()
print(grouped_dataframe)
```

Output:

	Salary	
Age		
19	30000	
25	50000	
32	70000	
47	90000	

Joining Data:

To join two or more DataFrames in Pandas, we can use the merge() function. We can specify the columns to join on using the on parameter, or we can specify the columns to join on for each DataFrame using the left_on and right_on parameters.

For example, to join two DataFrames based on a common column, we can use the following code:

```
import pandas as pd
my_dataframe1 = pd.read_csv('my_file1.csv')
my_dataframe2 = pd.read_csv('my_file2.csv')
joined_dataframe = pd.merge(my_dataframe1,
my_dataframe2, on='ID')
print(joined_dataframe)
```

Output:

	ID	Name	Age	Salary	Department
0	1	John	25	50000	1
1	2	Mary	32	70000	2
2	3	Peter	19	30000	1
3	4	Anne	47	90000	2

Conclusion:



In this article, we have provided an introduction to Python for SAS users, with a focus on data manipulation with Pandas. We have covered the basics of creating and manipulating Series and DataFrame objects in Pandas, as well as reading and writing data to various file formats. We have also covered some of the most commonly used functions for filtering, sorting, grouping, and joining data in Pandas. With this knowledge, SAS users can expand their data analysis skills to Python and take advantage of the powerful tools and libraries available in the Python ecosystem

```
import pandas as pd
# Read data from CSV file
my dataframe = pd.read csv('my data.csv')
# Print first 5 rows of the DataFrame
print(my dataframe.head())
# Create a new column called 'Birth Year' by
subtracting the Age column from the current year
current year = pd.Timestamp.now().year
my dataframe['Birth Year'] = current year -
my dataframe['Age']
# Print DataFrame with new column
print(my dataframe)
# Filter DataFrame to only include rows where the Age
column is greater than 30
filtered dataframe = my dataframe[my dataframe['Age'] >
301
# Print filtered DataFrame
print(filtered dataframe)
# Sort DataFrame by the Salary column in descending
order
sorted dataframe = my dataframe.sort values('Salary',
ascending=False)
# Print sorted DataFrame
print(sorted dataframe)
# Group DataFrame by the Age column and calculate the
mean Salary for each age group
grouped dataframe = my dataframe.groupby('Age').mean()
```



```
# Print grouped DataFrame
print(grouped_dataframe)
# Join two DataFrames based on a common column
my_dataframe1 = pd.read_csv('my_data1.csv')
my_dataframe2 = pd.read_csv('my_data2.csv')
joined_dataframe = pd.merge(my_dataframe1,
my_dataframe2, on='ID')
# Print joined DataFrame
print(joined_dataframe)
```

Note that this code assumes that the CSV files my_data.csv, my_data1.csv, and my_data2.csv exist in the current directory and contain the necessary data. You may need to modify the code to work with your own data files.

```
import pandas as pd
# Read data from Excel file
my dataframe = pd.read excel('my data.xlsx')
# Print last 5 rows of the DataFrame
print(my dataframe.tail())
# Create a new DataFrame by selecting rows where the
Age column is between 20 and 30
young dataframe = my dataframe[(my dataframe['Age'] >=
20) & (my dataframe['Age'] <= 30)]
# Print young DataFrame
print(young dataframe)
# Add a new column called 'Bonus' based on the value of
the Salary column
my dataframe['Bonus'] = my dataframe['Salary'] * 0.1
# Print DataFrame with new column
print(my dataframe)
# Pivot the DataFrame to show the mean Salary for each
department and age group
pivot dataframe =
my dataframe.pivot table(values='Salary',
index='Department', columns='Age', aggfunc='mean')
```



```
# Print pivot DataFrame
print(pivot_dataframe)
# Replace missing values in the DataFrame with the mean
value for that column
my_dataframe.fillna(my_dataframe.mean(), inplace=True)
# Print DataFrame with missing values replaced
print(my_dataframe)
# Save DataFrame to CSV file
my_dataframe.to_csv('my_data_updated.csv', index=False)
```

This code demonstrates additional techniques such as selecting rows based on multiple conditions, adding a new column based on an existing column, pivoting the DataFrame to show summary statistics, replacing missing values with the mean value for that column, and saving the updated DataFrame to a CSV file.

Again, note that this code assumes the existence of an Excel file my_data.xlsx in the current directory and may need to be modified for your own data files.

```
import pandas as pd
# Read data from JSON file
my dataframe = pd.read json('my data.json')
# Print DataFrame with the Country column in uppercase
my dataframe['Country'] =
my dataframe['Country'].str.upper()
print(my dataframe)
# Rename columns in the DataFrame
my dataframe.rename(columns={'Salary': 'Annual Salary',
'Age': 'Years Old'}, inplace=True)
print(my dataframe)
# Create a new DataFrame by grouping by two columns and
calculating the sum of a third column
grouped dataframe = my dataframe.groupby(['Department',
'Country'])['Annual Salary'].sum().reset index()
# Print grouped DataFrame
```

```
print(grouped_dataframe)
```



```
# Create a new column in the DataFrame based on a user-
defined function
def calculate bonus(salary):
    if salary >= 50000:
        return salary * 0.1
    else:
        return salary * 0.05
my dataframe['Bonus'] = my dataframe['Annual
Salary'].apply(calculate bonus)
# Print DataFrame with new column
print(my dataframe)
# Create a new DataFrame by concatenating two existing
DataFrames
my dataframe1 = pd.read csv('my data1.csv')
my dataframe2 = pd.read csv('my data2.csv')
concatenated dataframe = pd.concat([my dataframe1,
my dataframe2], ignore index=True)
# Print concatenated DataFrame
print(concatenated dataframe)
```

This code demonstrates additional techniques such as converting the values in a column to uppercase, renaming columns in the DataFrame, grouping by multiple columns and calculating summary statistics, creating a new column in the DataFrame based on a user-defined function, and concatenating two existing DataFrames.

Again, note that this code assumes the existence of a JSON file my_data.json and CSV files my_data1.csv and my_data2.csv in the current directory and may need to be modified for your own data files.

Introduction to Pandas

Introduction to Pandas:

Pandas is an open-source library for data manipulation and analysis in Python. It provides data structures for efficiently storing and manipulating large datasets, as well as tools for cleaning, merging, filtering, and transforming data.



Pandas is particularly useful for working with structured or tabular data, such as data in CSV or Excel files. It allows users to load data into memory, manipulate it in various ways, and then output the data to a new file or database.

Pandas provides two main data structures for working with data:

- 1. Series A one-dimensional array-like object that can hold any data type, such as integers, floats, or strings. Each value in a Series is assigned an index label, which can be used to retrieve or manipulate individual values.
- 2. DataFrame A two-dimensional tabular data structure with rows and columns. It can be thought of as a spreadsheet or SQL table. Each column in a DataFrame is a Series, and each row corresponds to a unique record or observation.

Pandas offers a wide range of functionality for working with these data structures, including:

- 1. Loading and saving data to various file formats, such as CSV, Excel, JSON, and SQL databases.
- 2. Cleaning and preprocessing data, such as removing missing values or duplicates, or transforming data types.
- 3. Selecting and filtering data based on certain criteria, such as values in a specific column or rows that meet certain conditions.
- 4. Aggregating and summarizing data, such as calculating summary statistics or grouping data by certain variables.
- 5. Merging and joining datasets based on common variables.
- 6. Creating new variables or columns based on existing data.
- 7. Visualizing data using built-in visualization tools.

Pandas is widely used in data analysis, scientific research, and machine learning applications. Its intuitive syntax and powerful functionality make it a popular tool for data analysts and data scientists alike.

Here are some more specific examples of what you can do with Pandas:

Loading Data: You can load data from a variety of file formats using Pandas, such as CSV, Excel, JSON, or SQL databases. For example, you can use the read_csv function to load data from a CSV file:

```
import pandas as pd
data = pd.read csv('my data.csv')
```

Cleaning and Preprocessing Data: Pandas provides a range of functions for cleaning and preprocessing data. For example, you can remove missing values using the dropna function, or replace missing values with a specific value using the fillna function:

Remove rows with missing values



```
data = data.dropna()
# Replace missing values with 0
data = data.fillna(0)
```

Selecting and Filtering Data: You can select and filter data based on certain criteria, such as values in a specific column or rows that meet certain conditions. For example, you can select a specific column using the [] operator:

```
# Select a specific column
column data = data['Column Name']
```

You can also filter rows based on specific conditions using boolean indexing:

```
# Filter rows where column 'A' is greater than 5
filtered data = data[data['A'] > 5]
```

Aggregating and Summarizing Data: Pandas provides functions for calculating summary statistics or grouping data by certain variables. For example, you can use the groupby function to group data by a specific column and calculate the sum of another column:

```
# Group data by 'Column A' and calculate the sum of
'Column B'
grouped_data = data.groupby('Column A')['Column
B'].sum()
```

Merging and Joining Datasets: You can merge or join datasets based on common variables using functions such as merge or concat. For example, you can merge two datasets based on a common column:

```
# Merge two datasets based on 'Column A'
merged data = pd.merge(data1, data2, on='Column A')
```

Creating New Variables: You can create new variables or columns based on existing data using the apply function or other functions provided by Pandas. For example, you can create a new column based on a function that calculates the square of another column:

```
# Create a new column based on the square of 'Column A'
data['Column C'] = data['Column A'].apply(lambda x:
x**2)
```

Visualizing Data: Pandas provides built-in visualization tools that allow you to create plots and charts from your data. For example, you can create a histogram of a specific column using the hist function:

Create a histogram of 'Column A'



data['Column A'].hist()

These are just a few examples of what you can do with Pandas. The library provides a wide range of functionality for working with data in Python, and is widely used in data analysis, scientific research, and machine learning applications.

Loading Data:

```
import pandas as pd
# Load data from CSV file
data = pd.read_csv('my_data.csv')
# Load data from Excel file
data = pd.read_excel('my_data.xlsx')
# Load data from SQL database
import sqlite3
conn = sqlite3.connect('my_database.db')
data = pd.read_sql_query('SELECT * FROM my_table',
conn)
```

Cleaning and Preprocessing Data:

```
# Remove rows with missing values
data = data.dropna()
# Replace missing values with 0
data = data.fillna(0)
# Convert string column to numeric column
data['Numeric Column'] = pd.to_numeric(data['String
Column'])
# Rename columns
data = data.rename(columns={'Old Name': 'New Name'})
# Remove duplicates
data = data.drop_duplicates()
# Remove outliers based on a specific column
Q1 = data['Column'].quantile(0.25)
Q3 = data['Column'].quantile(0.75)
IQR = Q3 - Q1
```



```
data = data[(data['Column'] >= Q1 - 1.5*IQR) &
  (data['Column'] <= Q3 + 1.5*IQR)]</pre>
```

Selecting and Filtering Data:

```
# Select a specific column
column_data = data['Column Name']
# Select multiple columns
column_data = data[['Column 1', 'Column 2']]
# Filter rows based on specific conditions
filtered_data = data[data['Column A'] > 5]
# Filter rows based on multiple conditions
filtered_data = data[(data['Column A'] > 5) &
(data['Column B'] < 10)]</pre>
```

Aggregating and Summarizing Data:

```
# Calculate the mean of a specific column
mean_data = data['Column'].mean()
# Group data by a specific column and calculate the sum
of another column
grouped_data = data.groupby('Column A')['Column
B'].sum()
# Pivot data and calculate the sum of a specific column
pivot_data = data.pivot_table(values='Column A',
index='Column B', columns='Column C', aggfunc='sum')
```

Merging and Joining Datasets:

```
# Merge two datasets based on a common column
merged_data = pd.merge(data1, data2, on='Column A')
# Join two datasets based on the index
joined_data = data1.join(data2, lsuffix='_left',
rsuffix='_right')
# Concatenate two datasets vertically
concatenated data = pd.concat([data1, data2], axis=0)
```



```
# Concatenate two datasets horizontally
concatenated_data = pd.concat([data1, data2], axis=1)
```

Creating New Variables:

```
# Create a new column based on a function that
calculates the square of another column
data['Column C'] = data['Column A'].apply(lambda x:
x**2)
# Create a new column based on a condition
data['Column C'] = np.where(data['Column A'] > 5,
'Yes', 'No')
# Create a new column based on a combination of columns
data['Column C'] = data['Column A'] + data['Column B']
```

Visualizing Data:

```
# Create a histogram of a specific column
data['Column A'].hist()
# Create a scatter plot of two columns
data.plot.scatter(x='Column A', y='Column B')
# Create a bar chart of a specific column
data['Column A'].value_counts().plot.bar()
```

Reshaping Data:

```
# Pivot data and calculate the mean of a specific
column
pivot_data = data.pivot_table(values='Column A',
index='Column B', columns='Column C', aggfunc='mean')
# Unstack a pivot table
unstacked_data = pivot_data.unstack()
# Reshape a dataframe from wide to long format
melted_data = pd.melt(data, id_vars=['Column A'],
value_vars=['Column B', 'Column C'], var_name='Variable
Name', value_name='Variable Value')
```



```
# Reshape a dataframe from long to wide format
wide_data = melted_data.pivot_table(values='Variable
Value', index=['Column A'], columns=['Variable Name'])
```

Handling Dates:

```
# Convert a string column to a datetime column
data['Date Column'] = pd.to_datetime(data['Date
Column'], format='%Y-%m-%d')
```

```
# Extract year, month, and day from a datetime column
data['Year Column'] = data['Date Column'].dt.year
data['Month Column'] = data['Date Column'].dt.month
data['Day Column'] = data['Date Column'].dt.day
```

```
# Calculate the difference between two datetime columns
data['Date Difference'] = data['Date Column 1'] -
data['Date Column 2']
```

```
# Create a datetime index
data = data.set index('Date Column')
```

Working with Categorical Variables:

```
# Convert a string column to a categorical column
data['Categorical Column'] = data['String
Column'].astype('category')
```

```
# Create dummy variables from a categorical column
dummy_data = pd.get_dummies(data['Categorical Column'])
```

```
# Group data by a categorical column and calculate the
mean of another column
grouped_data = data.groupby('Categorical
Column')['Numeric Column'].mean()
```

```
# Rename categories in a categorical column
data['Categorical Column'] = data['Categorical
Column'].cat.rename_categories({'Category A': 'Category
1', 'Category B': 'Category 2'})
```

Handling Missing Values:



```
# Replace missing values with the mean of a specific
column
data['Column A'] = data['Column A'].fillna(data['Column
A'].mean())
# Interpolate missing values in a specific column
data['Column A'] = data['Column A'].interpolate()
# Drop rows with missing values in a specific column
data = data.dropna(subset=['Column A'])
# Forward-fill missing values in a specific column
data['Column A'] = data['Column A'].ffill()
```

Working with Time Series Data:

```
# Resample time series data to a different frequency
resampled_data = data.resample('D').mean()
# Calculate the rolling mean of a time series
rolling_data = data.rolling(window=7).mean()
# Shift time series data forward or backward in time
shifted_data = data.shift(7)
# Calculate the difference between two time series
diff data = data.diff()
```

Working with Text Data:

```
# Convert a string column to lowercase
data['String Column'] = data['String
Column'].str.lower()
# Split a string column into multiple columns
split_data = data['String
Column'].str.split(expand=True)
# Replace values in a string column using regular
expressions
data['String Column'] = data['String
Column'].str.replace(r'\d+', '')
```



```
# Count the occurrence of a specific word in a string
column
data['Word Count'] = data['String
Column'].str.count('word')
```

Joins and Merges:

```
# Merge two dataframes based on a common column
merged_data = pd.merge(df1, df2, on='Column A')
# Perform an outer join between two dataframes
outer_joined_data = pd.merge(df1, df2, on='Column A',
how='outer')
# Concatenate two dataframes vertically
concatenated_data = pd.concat([df1, df2], axis=0)
# Concatenate two dataframes horizontally
concatenated_data = pd.concat([df1, df2], axis=1)
```

Aggregation and Grouping:

```
# Calculate the mean of a specific column
mean_data = data['Column A'].mean()
# Calculate the median of a specific column
median_data = data['Column A'].median()
# Calculate the standard deviation of a specific column
std_data = data['Column A'].std()
# Group data by a column and calculate the mean of
another column
grouped_data = data.groupby('Column A')['Column
B'].mean()
```

Reshaping Data:

```
# Pivot data and calculate the mean of a specific
column
pivot_data = data.pivot_table(values='Column A',
index='Column B', columns='Column C', aggfunc='mean')
# Unstack a pivot table
```



```
unstacked_data = pivot_data.unstack()
# Reshape a dataframe from wide to long format
melted_data = pd.melt(data, id_vars=['Column A'],
value_vars=['Column B', 'Column C'], var_name='Variable
Name', value_name='Variable Value')
```

```
# Reshape a dataframe from long to wide format
wide_data = melted_data.pivot_table(values='Variable
Value', index=['Column A'], columns=['Variable Name'])
```

Handling Dates:

```
# Convert a string column to a datetime column
data['Date Column'] = pd.to_datetime(data['Date
Column'], format='%Y-%m-%d')
# Extract year, month, and day from a datetime column
data['Year Column'] = data['Date Column'].dt.year
data['Month Column'] = data['Date Column'].dt.month
data['Day Column'] = data['Date Column'].dt.day
# Calculate the difference between two datetime columns
data['Date Difference'] = data['Date Column 1'] -
data['Date Column 2']
# Create a datetime index
```

data = data.set_index('Date Column')

Working with Categorical Variables:

```
# Convert a string column to a categorical column
data['Categorical Column'] = data['String
Column'].astype('category')
# Create dummy variables from a categorical column
dummy_data = pd.get_dummies(data['Categorical Column'])
# Group data by a categorical column and calculate the
mean of another column
grouped_data = data.groupby('Categorical
Column')['Numeric Column'].mean()
# Rename categories in a categorical column
```



```
data['Categorical Column'] = data['Categorical
Column'].cat.rename_categories({'Category A': 'Category
1', 'Category B': 'Category 2'})
```

Handling Missing Values:

```
# Replace missing values with the mean of a specific
column
data['Column A'] = data['Column A'].fillna(data['Column
A'].mean())
# Interpolate missing values in a specific column
data['Column A'] = data['Column A'].interpolate()
# Drop rows with missing values in a specific column
data = data.dropna(subset=['Column A'])
# Forward-fill missing values in a specific column
data['Column A'] = data['Column A'].ffill()
Creating a DataFrame
```

Reading data into a DataFrame

Reading data into a DataFrame is an essential task in data analysis using Python. A DataFrame is a two-dimensional labeled data structure that consists of rows and columns, where each column can have a different data type. In Python, the most popular library for data manipulation is pandas, which provides a powerful DataFrame object and a set of functions to load data from various sources.

If you are a SAS user transitioning to Python, you might find some similarities in the way data is loaded into a DataFrame. In this section, we will introduce you to the basics of reading data into a DataFrame using pandas, with a focus on how it compares to SAS.

Loading CSV Data

The most common way to load data into a pandas DataFrame is by reading a CSV (commaseparated values) file. CSV files are plain text files that contain data separated by commas or other delimiters. To load a CSV file into a DataFrame, you can use the read_csv() function in pandas. Here's an example:

import pandas as pd



df = pd.read_csv('mydata.csv')

In this example, we first import the pandas library and then use the read_csv() function to read a file called 'mydata.csv' into a DataFrame called df. The read_csv() function automatically infers the data types of the columns and uses the first row of the CSV file as the column names.

In SAS, you can read a CSV file using the IMPORT procedure. Here's an example:

In this example, we use the IMPORT procedure to read a file called 'mydata.csv' into a SAS dataset called mydata. The dbms=csv option tells SAS that the file is in CSV format.

As you can see, the syntax for reading CSV files in pandas and SAS is quite different. However, the basic idea is the same: both libraries provide a function or procedure to read a CSV file and create a dataset or DataFrame.

Loading Excel Data

Another common way to load data into a DataFrame is by reading an Excel file. Excel files are popular in business settings, where data is often stored in spreadsheets. To read an Excel file into a DataFrame, you can use the read_excel() function in pandas. Here's an example:

```
import pandas as pd
df = pd.read excel('mydata.xlsx')
```

In this example, we use the read_excel() function to read a file called 'mydata.xlsx' into a DataFrame called df. The read_excel() function automatically infers the data types of the columns and uses the first row of the Excel file as the column names.

In SAS, you can read an Excel file using the IMPORT procedure with the DBMS=XLSX option. Here's an example:

In this example, we use the IMPORT procedure to read a file called 'mydata.xlsx' into a SAS dataset called mydata. The dbms=xlsx option tells SAS that the file is in Excel format.

Once again, the syntax for reading Excel files in pandas and SAS is different, but the basic idea is the same.



```
# Read SQL query into a DataFrame
df = pd.read_sql_query("SELECT * from mytable", conn)
# Load data from a JSON file
df = pd.read_json('mydata.json')
# Load data from an HTML table
url = 'https://www.fdic.gov/resources/resolutions/bank-
failures/failed-bank-list/'
dfs = pd.read_html(url)
df = dfs[0]
```

In SAS, you can also read data from a variety of sources, including SQL databases, JSON files, and HTML tables. Here are some examples:

```
/* Load data from a SQL database */
proc sql;
  connect to sqlite (path='mydatabase.db');
  select * from connection to sqlite
    (select * from mytable);
quit;
/* Load data from a JSON file */
filename myjson 'mydata.json';
libname myjson json fileref=myjson;
data mydata;
  set myjson.mytable;
run;
/* Load data from an HTML table */
filename myhtml url
'https://www.fdic.gov/resources/resolutions/bank-
failures/failed-bank-list/';
proc import datafile=myhtml
            out=mydata
            dbms=html;
run;
```

As you can see, both pandas and SAS provide a variety of options to read data from various sources. However, the syntax for loading data can differ significantly between the two tools.

Handling Missing Data

One common issue when working with data is missing values. In pandas, missing values are represented by the NaN (Not a Number) value, which is part of the NumPy library. When you



load data into a DataFrame, pandas automatically detects missing values and replaces them with NaN.

Here's an example:

```
import pandas as pd
df = pd.read_csv('mydata.csv')
# Check for missing values
print(df.isna().sum())
```

In this example, we use the isna() function to check for missing values in the DataFrame df. The sum() function is used to count the number of missing values in each column.

In SAS, missing values are represented by a dot (.) in numeric variables and by a blank () in character variables. When you load data into a SAS dataset, SAS automatically detects missing values and assigns them the appropriate value.

Here's an example:

```
data mydata;
    infile 'mydata.csv' delimiter=',' missover;
    input var1 var2 var3;
run;
/* Check for missing values */
proc means data=mydata n nmiss;
run;
```

In this example, we use the input statement to read data from a CSV file into a SAS dataset called mydata. The missover option tells SAS to ignore missing values when reading the data. The proc means procedure is then used to count the number of missing values in each variable.

As you can see, handling missing values is slightly different in pandas and SAS, but the basic idea is the same: both tools provide a way to detect and count missing values in your data.

In this section, we've introduced you to the basics of reading data into a DataFrame in Python using pandas, with a focus on how it compares to SAS. We've covered loading CSV and Excel files, as well as reading data from other sources like SQL databases, JSON files, and HTML tables. We've also touched on how to handle missing values in your data.

While the syntax for loading data can differ between pandas and SAS, the basic idea is the same: both tools provide a way to read data into a dataset or DataFrame so you can begin your data analysis. As you become more comfortable



```
import pandas as pd
# Read a CSV file into a DataFrame
df = pd.read_csv('mydata.csv')
# Read an Excel file into a DataFrame
df = pd.read_excel('mydata.xlsx')
# Read a SQL query into a DataFrame
conn = create_engine('sqlite:///mydatabase.db')
df = pd.read_sql_query("SELECT * from mytable", conn)
# Load data from a JSON file
df = pd.read_json('mydata.json')
# Load data from an HTML table
url = 'https://www.fdic.gov/resources/resolutions/bank-
failures/failed-bank-list/'
dfs = pd.read_html(url)
df = dfs[0]
```

Let's break this down line by line:

import pandas as pd

This line imports the pandas library and aliases it as pd, which is a common convention in the Python community.

```
# Read a CSV file into a DataFrame
df = pd.read csv('mydata.csv')
```

This line reads a CSV file named mydata.csv into a pandas DataFrame called df. By default, pandas assumes that the first row of the CSV file contains column headers, but you can override this behavior with the header parameter.

```
# Read an Excel file into a DataFrame
df = pd.read excel('mydata.xlsx')
```

This line reads an Excel file named mydata.xlsx into a pandas DataFrame called df. By default, pandas assumes that the first sheet in the Excel file contains the data, but you can specify a different sheet with the sheet_name parameter.

```
# Read a SQL query into a DataFrame
conn = create engine('sqlite:///mydatabase.db')
```



df = pd.read_sql_query("SELECT * from mytable", conn)

This code uses the create_engine function from the SQLAlchemy library to create a database connection to a SQLite database file named mydatabase.db. The pd.read_sql_query function then reads a SQL query that selects all columns from a table named mytable into a pandas DataFrame called df.

```
# Load data from a JSON file
df = pd.read_json('mydata.json')
```

This code reads data from a JSON file named mydata.json into a pandas DataFrame called df. By default, pandas assumes that the JSON file contains a single object, but you can override this behavior with the orient parameter.

```
# Load data from an HTML table
url = 'https://www.fdic.gov/resources/resolutions/bank-
failures/failed-bank-list/'
dfs = pd.read_html(url)
df = dfs[0]
```

This code reads an HTML table from a URL using the pd.read_html function, which returns a list of DataFrames (one for each HTML table on the page). In this case, we're assuming that there's only one HTML table on the page, so we grab the first element of the list and assign it to a pandas DataFrame called df.

Let's take a closer look at some of the parameters that can be used with these functions.

For pd.read_csv, some commonly used parameters include:

sep: the delimiter used in the file (default is ,) header: which row to use as the column headers (default is 0) index_col: which column to use as the index (default is None) usecols: which columns to read (default is all columns)

For pd.read_excel, some commonly used parameters include:

sheet_name: which sheet to read (default is 0) header: which row to use as the column headers (default is 0) index_col: which column to use as the index (default is None) usecols: which columns to read (default is all columns)

For pd.read_sql_query, some commonly used parameters include:

params: a list or tuple of parameters to pass to the SQL query parse_dates: a list of column names to parse as dates chunksize: the number of rows to read at a time (useful for very large datasets)

in stal

For pd.read_json, some commonly used parameters include:

orient: the JSON format ('split', 'records', 'index', 'columns', or 'values') For pd.read_html, some commonly used parameters include:

header: which row to use as the column headers (default is 0) index_col: which column to use as the index (default is None) flavor: the HTML parser to use ('bs4' or 'lxml')

Note that these are just a few examples of the many parameters that can be used with these functions. To learn more, you can refer to the pandas documentation or experiment with the parameters yourself.

Overall, reading data into a DataFrame is a key step in data analysis with Python. Once you have your data in a DataFrame, you can use pandas and other Python libraries to manipulate, analyze, and visualize the data in a variety of ways.

Reading data from a URL

```
import pandas as pd
url =
    'https://raw.githubusercontent.com/datasciencedojo/data
    sets/master/titanic.csv'
df = pd.read csv(url)
```

This code reads a CSV file from a URL into a pandas DataFrame using the pd.read_csv function.

Reading data from a clipboard

import pandas as pd
df = pd.read clipboard()

This code reads data from the clipboard (e.g., if you've copied a table from a website or spreadsheet program) into a pandas DataFrame using the pd.read_clipboard function.

Reading data from a Python dictionary



df = pd.DataFrame(data)

This code creates a Python dictionary and then converts it to a pandas DataFrame using the pd.DataFrame function.

Reading data from a list of dictionaries

This code creates a list of dictionaries and then converts it to a pandas DataFrame using the pd.DataFrame function.

Reading data from a NumPy array

This code creates a NumPy array and then converts it to a pandas DataFrame using the pd.DataFrame function.

Overall, pandas provides many different ways to read data into a DataFrame, making it easy to work with data from a variety of sources. By understanding these methods and their respective parameters, you can efficiently import and manipulate data in Python for a wide range of data analysis tasks.

Indexing and selecting data

Indexing and selecting data are essential concepts when working with data in Python. They allow you to retrieve specific data elements from a data structure such as a list, tuple, or dictionary. In this context, Python offers several methods for indexing and selecting data, including slicing, indexing, and boolean indexing.

Slicing is a technique used to select a subset of a sequence object, such as a list or string, by specifying a range of indices to retrieve. For example, suppose you have a list of numbers [1, 2, 3, 4, 5], and you want to retrieve the second to fourth elements. In this case, you can use slicing by specifying the range of indices [1:4]. The result will be a new list containing [2, 3, 4].

Indexing, on the other hand, is used to access a specific element of a sequence by specifying its position in the sequence. Python uses zero-based indexing, which means that the first element of a sequence has an index of 0. For example, if you have a list of colors ['red', 'green', 'blue'], you can access the second element ('green') by indexing the list with 1, as follows: colors[1].

Boolean indexing is a method used to retrieve elements from a sequence that meet a certain condition. For example, suppose you have a list of numbers [1, 2, 3, 4, 5], and you want to retrieve only the even numbers. In this case, you can use boolean indexing by specifying the condition "x % 2 == 0", where "x" represents each element of the list. The result will be a new list containing [2, 4].

In Python, you can also use these techniques to index and select data from pandas dataframes, which are commonly used in data analysis. For example, suppose you have a dataframe with columns "Name", "Age", and "Gender". You can use indexing to retrieve the "Age" column by specifying df["Age"], and you can use boolean indexing to retrieve only the rows where the age is greater than 30 by specifying df[df["Age"] > 30].

SAS users who are new to Python may find these indexing and selection methods helpful for manipulating and analyzing data in Python. Additionally, Python offers a wide range of libraries and tools for data analysis, including pandas, NumPy, and matplotlib, which can be used in conjunction with these methods to perform complex data analyses.

here are some more details on indexing and selecting data in Python:

Slicing:

Slicing is used to extract a portion of a sequence by specifying a range of indices. The general syntax for slicing is sequence[start:end:step], where start is the starting index, end is the ending index (exclusive), and step is the step size. If any of these parameters are omitted, Python assumes a default value (start = 0, end = len(sequence), step = 1). Here are some examples: a. Select the first three elements of a list:

```
my_list = [1, 2, 3, 4, 5]
my_list[:3] # [1, 2, 3]
```



b. Select every other element of a list:

```
my_list = [1, 2, 3, 4, 5]
my_list[::2] # [1, 3, 5]
```

Indexing:

Indexing is used to access a specific element of a sequence by specifying its position in the sequence. The general syntax for indexing is sequence[index]. Here are some examples: a. Retrieve the first element of a list:

```
my_list = [1, 2, 3, 4, 5]
my_list[0] # 1
b. Retrieve the last element of a list:
css
Copy code
my_list = [1, 2, 3, 4, 5]
my_list[-1] # 5
```

Boolean Indexing:

Boolean indexing is used to select elements from a sequence that satisfy a certain condition. The general syntax for boolean indexing is sequence[condition], where condition is a Boolean expression that returns True or False for each element in the sequence. Here are some examples: a. Retrieve even numbers from a list:

my_list = [1, 2, 3, 4, 5]
even_numbers = [x for x in my_list if x % 2 == 0]
print(even numbers) # [2, 4]

b. Retrieve rows from a pandas dataframe where the age is greater than 30:

These indexing and selection methods are commonly used in data analysis and manipulation tasks, and can be combined with other Python tools and libraries to perform complex operations on large datasets.

Slicing:

create a list of numbers



```
numbers = [1, 2, 3, 4, 5]
# select the first three numbers using slicing
first_three = numbers[:3]
# select every other number using slicing
every_other = numbers[::2]
# print the results
print(first_three) # [1, 2, 3]
print(every other) # [1, 3, 5]
```

Indexing:

```
# create a list of colors
colors = ['red', 'green', 'blue']
# retrieve the first color using indexing
first_color = colors[0]
# retrieve the last color using negative indexing
last_color = colors[-1]
# print the results
print(first_color) # red
print(last_color) # blue
```

Boolean indexing:

```
import numpy as np
# create a 2D array of random numbers
array = np.random.rand(4, 4)
# select only the elements greater than 0.5
greater_than_half = array[array > 0.5]
# print the result
print(greater_than_half)
```

Using pandas dataframes:

```
import pandas as pd
# create a dataframe of employee data
```



These examples demonstrate how indexing and selection can be used to retrieve specific data elements from different data structures in Python. These techniques are especially useful in data analysis and manipulation tasks, where you often need to extract and analyze specific subsets of data from larger datasets.

Selecting elements from nested lists:

```
# create a nested list
nested list = [[1, 2, 3], [4, 5, 6], [7, 8, 9]]
# select the second element of the first list
element = nested list[0][1]
# print the result
print(element)
Selecting elements from numpy arrays:
php
Copy code
import numpy as np
# create a 2D numpy array
array = np.array([[1, 2, 3], [4, 5, 6], [7, 8, 9]])
# select the second column of the array
column = array[:, 1]
# select the diagonal of the array
diagonal = array.diagonal()
```



```
# print the results
print(column)
print(diagonal)
```

Using boolean indexing with pandas dataframes:

```
import pandas as pd
# create a dataframe of student grades
grades = pd.DataFrame({'Name': ['Alice', 'Bob',
'Charlie', 'David'],
                       'Math': [80, 70, 90, 60],
                       'English': [75, 85, 80, 70],
                       'Science': [95, 80, 85, 75]})
# select only the rows where the math grade is above 80
good math grades = grades[grades['Math'] > 80]
# select only the rows where the student's overall
grade is above 80
good grades = grades[(grades['Math'] +
grades['English'] + grades['Science']) / 3 > 80]
# print the results
print(good math grades)
print(good grades)
```

Using loc and iloc with pandas dataframes:



```
first_two_grades = grades.iloc[:2, [1, 3]]
# print the results
print(alice)
print(first two grades)
```

Using boolean indexing with numpy arrays:

```
import numpy as np
# create a 2D array of random integers
array = np.random.randint(1, 10, (4, 4))
# select the elements that are greater than 5
greater_than_5 = array[array > 5]
# select the elements that are even
even_numbers = array[array % 2 == 0]
# print the results
print(greater_than_5)
print(even_numbers)
```

Using slicing with strings:

```
# create a string
s = 'hello world'
# select the first five characters
first_five = s[:5]
# select the last five characters
last_five = s[-5:]
# print the results
print(first_five)
print(last_five)
```

Selecting elements from dictionaries:

```
# create a dictionary
d = {'apple': 1, 'banana': 2, 'orange': 3}
# select the value for the 'apple' key
in stal
```

```
apple_value = d['apple']
# print the result
print(apple value)
```

Using slicing with numpy arrays:

```
import numpy as np
# create a 2D numpy array
array = np.array([[1, 2, 3], [4, 5, 6], [7, 8, 9]])
# select the first two rows and all columns
first_two_rows = array[:2, :]
# select the second column and all rows
second_column = array[:, 1]
# print the results
print(first_two_rows)
print(second_column)
```

Using integer indexing with pandas dataframes:



Using boolean indexing with pandas dataframes:

```
import pandas as pd
# create a dataframe of student grades
grades = pd.DataFrame({'Name': ['Alice', 'Bob',
'Charlie', 'David'],
                        'Math': [80, 70, 90, 60],
                        'English': [75, 85, 80, 70],
                        'Science': [95, 80, 85, 75]})
# select the rows where the math grade is greater than
or equal to 80
high math grades = grades[grades['Math'] >= 80]
# select the rows where the science grade is less than
or equal to 80
low science grades = grades[grades['Science'] <= 80]</pre>
# print the results
print(high math grades)
print(low science grades)
```

Using slicing with lists:

```
# create a list
my_list = ['a', 'b', 'c', 'd', 'e']
# select the first three elements
first_three = my_list[:3]
# select the last two elements
last_two = my_list[-2:]
# print the results
print(first_three)
print(last two)
```

Using integer indexing with numpy arrays:

import numpy as np
create a 2D numpy array



```
array = np.array([[1, 2, 3], [4, 5, 6], [7, 8, 9]])
# select the element in the first row and second column
element = array[0, 1]
# select the first two elements of the second row
first_two_in_second_row = array[1, :2]
# print the results
print(element)
print(first_two_in_second_row)
```

Filtering and sorting data

Filtering and sorting data are essential operations in data analysis and are commonly performed using programming languages like Python and SAS. In this article, we will discuss how to filter and sort data in Python, with a focus on how SAS users can leverage their existing knowledge to learn Python.

Filtering Data in Python

In Python, filtering data is typically done using the Pandas library, which provides a powerful set of tools for working with tabular data. The most commonly used function for filtering data in Pandas is the query() function, which allows you to select rows of data based on a Boolean expression.

For example, suppose you have a Pandas DataFrame that contains information about a set of products, including their names, prices, and ratings. You can filter this DataFrame to select only the products with a rating above 4 using the following code:



In this code, we first create a DataFrame df with information about the products. We then use the query() function to select only the rows where the Rating column is greater than 4, and store the result in a new DataFrame filtered_df.

Sorting Data in Python

Sorting data is also a common operation in data analysis, and can be done easily in Python using the sort_values() function in Pandas. This function allows you to sort a DataFrame based on one or more columns.

For example, suppose you want to sort the products DataFrame above by price, from lowest to highest. You can do this using the following code:

```
sorted df = df.sort values('Price')
```

In this code, we use the sort_values() function to sort the DataFrame df by the Price column. By default, the function sorts the data in ascending order (lowest to highest). If you want to sort the data in descending order (highest to lowest), you can add the ascending=False argument to the function call:

```
sorted df = df.sort values('Price', ascending=False)
```

Filtering and Sorting Data in SAS vs. Python

For SAS users who are familiar with the WHERE and SORT statements, the syntax for filtering and sorting data in Python may seem quite different. However, the underlying concepts are the same, and with a little practice, SAS users can easily adapt to using Python for these tasks.

In SAS, you might filter data using a WHERE statement like this:

```
data Products;
  set Products;
  where Rating > 4;
run;
```

This code selects only the rows from the Products data set where the Rating variable is greater than 4.

In Python, we can achieve the same result using the query() function:

```
filtered df = df.query('Rating > 4')
```

Similarly, in SAS, you might sort data using a SORT statement like this:

```
proc sort data=Products;
   by Price;
run;
```



This code sorts the Products data set by the Price variable, in ascending order.

In Python, we can achieve the same result using the sort_values() function:

```
sorted df = df.sort values('Price')
```

SAS, the underlying concepts and operations are similar, making it relatively easy for SAS users to learn Python.

It's worth noting that while filtering and sorting data are essential operations in data analysis, they are just the tip of the iceberg when it comes to what you can do with Python and Pandas. Pandas provides a rich set of tools for data manipulation, cleaning, and analysis, including functions for merging and joining datasets, reshaping data, and calculating summary statistics.

Additionally, Python has a vast ecosystem of libraries and packages for data analysis and visualization, such as NumPy, Matplotlib, and Scikit-learn. These libraries can help you perform advanced data analysis and machine learning tasks in Python, making it a powerful tool for data scientists and analysts.

Filtering and sorting data in Python using Pandas is a critical skill for any data analyst or scientist. For SAS users, the transition to Python should be relatively straightforward, as the underlying concepts and operations are similar. With a little practice, SAS users can easily adapt to using Python and leverage its rich ecosystem of libraries and tools for data analysis and visualization.

Filtering Data:

Output:

Product Price Rating



0 Product A 10.99 4.5 2 Product C 15.99 4.8

In this code, we create a DataFrame df with information about products and then use the query() function to filter the data and select only the rows where the Rating column is greater than 4.

Sorting Data:

	Product	Price	Rating
1	Product B	8.99	3.9
0	Product A	10.99	4.5
3	Product D	12.99	4.2
2	Product C	15.99	4.8

In this code, we create a DataFrame df with information about products and then use the sort_values() function to sort the data by the Price column, in ascending order.

Filtering and Sorting Data:



0

2

In this code, we combine the previous examples and filter the DataFrame by the Rating column and sort the resulting data by the Price column.

4.5

4.8

Filtering and Selecting Columns:

Product A 10.99

Product C 15.99

0 Product A 10.99 2 Product C 15.99

In this code, we filter the DataFrame by the Rating column and select only the Product and Price columns of the resulting data.

Filtering Data using Multiple Conditions:

```
import pandas as pd
```



Product Price Rating Category 0 Product A 10.99 4.5 Electronics 2 Product C 15.99 4.8 Electronics

In this code, we create a DataFrame df with additional Category column and filter the data by both Rating and Category using multiple conditions in a single line of code.

Sorting Data in Descending Order:

```
import pandas as pd
# create DataFrame
data = {'Product': ['Product A', 'Product B', 'Product
C', 'Product D'],
                     'Price': [10.99, 8.99, 15.99, 12.99],
                          'Rating': [4.5, 3.9, 4.8, 4.2]}
df = pd.DataFrame(data)
# sort DataFrame in descending order
sorted_df = df.sort_values('Price', ascending=False)
print(sorted_df)
Output:
```

```
Product Price Rating
2 Product C 15.99 4.8
```



3	Product D	12.99	4.2
0	Product A	10.99	4.5
1	Product B	8.99	3.9

In this code, we sort the DataFrame by the Price column, but this time we use the ascending=False argument to sort the data in descending order.

Filtering Data using Regular Expressions:

NameAge0John Smith252Bob Johnson40

In this code, we create a DataFrame df with a Name column and filter the data by selecting only the rows where the Name column contains the substring 'Jo' using a regular expression in the str.contains() function.

Aggregating and summarizing data

Python is a popular programming language that is widely used for data analysis, machine learning, and scientific computing. Many SAS users have started to use Python for data analysis and modeling tasks because of its flexibility, scalability, and ease of use. This article provides an introduction to Python for SAS users, focusing on aggregating and summarizing data.

Data aggregation and summarization are important techniques in data analysis, as they help to provide a concise and meaningful summary of large datasets. In Python, these tasks can be accomplished using a variety of libraries, including NumPy, pandas, and SciPy.

NumPy is a library for numerical computing in Python. It provides a powerful array data structure that can be used to perform operations on large datasets. One of the key features of NumPy is its ability to perform aggregation and summarization operations on arrays. For example, the np.mean() function can be used to calculate the mean of an array, while the np.sum() function can be used to calculate the sum of an array.

Pandas is a library for data manipulation and analysis in Python. It provides a powerful DataFrame data structure that can be used to store and manipulate tabular data. One of the key features of pandas is its ability to perform aggregation and summarization operations on DataFrames. For example, the df.mean() function can be used to calculate the mean of a DataFrame, while the df.sum() function can be used to calculate the sum of a DataFrame.

SciPy is a library for scientific computing in Python. It provides a wide range of functions for numerical optimization, integration, and statistics. One of the key features of SciPy is its ability to perform statistical analysis on datasets. For example, the scipy.stats.describe() function can be used to calculate descriptive statistics on a dataset, such as the mean, standard deviation, and quartiles.

In addition to these libraries, Python also provides a variety of built-in functions for aggregating and summarizing data. For example, the sum() function can be used to calculate the sum of a list or tuple, while the len() function can be used to calculate the length of a list or tuple.

Overall, Python provides a powerful and flexible set of tools for aggregating and summarizing data. By using these tools, SAS users can leverage the power of Python to analyze large datasets and gain valuable insights into their data. Aggregating and summarizing data are important techniques for data analysis, as they allow us to quickly gain insights into the underlying patterns and trends in our data. These techniques are particularly useful for large datasets, where it is often impractical to examine each individual data point.

Python provides a wide range of tools and libraries for aggregating and summarizing data. In addition to NumPy, pandas, and SciPy, other popular libraries for data analysis in Python include Matplotlib, Seaborn, and Scikit-learn.

Matplotlib is a library for data visualization in Python. It provides a wide range of functions for creating plots and charts, which can be used to visualize the results of aggregation and summarization operations. For example, a histogram can be used to visualize the distribution of values in a dataset, while a scatter plot can be used to visualize the relationship between two variables.

Seaborn is a library for statistical data visualization in Python. It provides a variety of functions for creating complex visualizations, such as heatmaps and box plots. These visualizations can be used to gain insights into the underlying patterns and trends in our data.



Scikit-learn is a library for machine learning in Python. It provides a wide range of functions for data preprocessing, modeling, and evaluation. These functions can be used to build predictive models based on our aggregated and summarized data.

In addition to these libraries, Python also provides a variety of functions for working with data, such as sorting, filtering, and grouping. For example, the sorted() function can be used to sort a list or tuple, while the filter() function can be used to select only those elements that meet a certain condition. For SAS users, learning Python for data analysis can be a valuable skill. Python is widely used in the data science community, and its popularity has been growing rapidly in recent years. Many organizations are now using Python alongside SAS to analyze and model their data.

One of the key advantages of Python over SAS is its flexibility. Python is an open-source language, which means that it is highly customizable and extensible. This flexibility allows users to create custom scripts and functions for specific data analysis tasks. Additionally, Python provides a large ecosystem of libraries and tools that can be used for data analysis, machine learning, and scientific computing.

Example 1: Using NumPy to calculate the mean of an array

```
import numpy as np
# Create an array of numbers
data = np.array([3, 5, 7, 9, 11])
# Calculate the mean of the array
mean = np.mean(data)
# Print the mean
print("Mean:", mean)
```

Output:

```
Mean: 7.0
```

Example 2: Using Pandas to calculate the sum of a DataFrame

```
import pandas as pd
# Create a DataFrame of numbers
data = pd.DataFrame({'A': [1, 2, 3], 'B': [4, 5, 6],
'C': [7, 8, 9]})
# Calculate the sum of the DataFrame
sum = data.sum()
```



```
# Print the sum
print("Sum:\n", sum)
```

Output:

```
Sum:
A 6
B 15
C 24
dtype: int64
```

Example 3: Using SciPy to calculate descriptive statistics on a dataset

```
import scipy.stats as stats
# Create a dataset of numbers
data = [3, 5, 7, 9, 11]
# Calculate descriptive statistics
mean = stats.describe(data).mean
std = stats.describe(data).std
min = stats.describe(data).minmax[0]
max = stats.describe(data).minmax[1]
# Print the descriptive statistics
print("Mean:", mean)
print("Standard Deviation:", std)
print("Minimum Value:", min)
print("Maximum Value:", max)
```

Output:

```
Mean: 7.0
Standard Deviation: 3.1622776601683795
Minimum Value: 3
Maximum Value: 11
```

These examples demonstrate how Python can be used to aggregate and summarize data using different libraries and functions. With its flexible and extensible nature, Python provides a powerful set of tools for data analysis that can be used alongside SAS.



Example 4: Using Pandas to group data by a categorical variable and calculate the mean of each group

```
import pandas as pd
# Create a DataFrame of sales data
data = pd.DataFrame({'Region': ['East', 'West',
'North', 'South', 'East', 'West', 'North', 'South'],
                           'Sales': [100, 200, 150, 250, 300,
400, 350, 450]})
# Group the data by region and calculate the mean of
each group
grouped_data = data.groupby('Region').mean()
# Print the grouped data
print(grouped_data)
```

Output:

```
        Sales

        Region

        East
        200.000000

        North
        250.000000

        South
        350.000000

        West
        300.000000
```

Example 5: Using NumPy to calculate the median of an array

```
import numpy as np
# Create an array of numbers
data = np.array([3, 5, 7, 9, 11])
# Calculate the median of the array
median = np.median(data)
# Print the median
print("Median:", median)
```

Output:

Median: 7.0



Example 6: Using Pandas to filter data based on a condition and calculate the sum of a column

Output:

Sum: 400

Example 7: Using Pandas to calculate the correlation coefficient between two columns in a DataFrame

Output:



Correlation Coefficient: 0.9966155281280883

Example 8: Using NumPy to calculate the standard deviation of an array

```
import numpy as np
# Create an array of numbers
data = np.array([3, 5, 7, 9, 11])
# Calculate the standard deviation of the array
std_dev = np.std(data)
# Print the standard deviation
print("Standard Deviation:", std_dev)
```

Output:

Standard Deviation: 3.1622776601683795

Example 9: Using Pandas to pivot a DataFrame and calculate the sum of values in each column

```
import pandas as pd
    # Create a DataFrame of sales data
    data = pd.DataFrame({'Region': ['East', 'West',
    'North', 'South', 'East', 'West', 'North', 'South'],
                          'Product': ['A', 'B', 'C', 'D',
    'A', 'B', 'C', 'D'],
                          'Sales': [100, 200, 150, 250, 300,
    400, 350, 450]})
    # Pivot the data and calculate the sum of Sales for
    each Region and Product
    pivot data = pd.pivot table(data, values='Sales',
    index='Region', columns='Product', aggfunc=np.sum)
    # Print the pivoted data
    print(pivot data)
Output:
    Product A B C
                              D
    Region
```

0





North	0	0	500	0
South	0	0	0	700
West	200	400	0	0

Example 10: Using Pandas to merge two DataFrames on a common column

```
import pandas as pd
    # Create two DataFrames of sales data
    sales data 1 = pd.DataFrame({'Region': ['East', 'West',
    'North', 'South'],
                                  'Sales': [100, 200, 150,
    250]})
    sales data 2 = pd.DataFrame({'Region': ['East', 'West',
    'North', 'South'],
                                  'Profit': [20, 40, 30,
    50]})
    # Merge the two DataFrames on the Region column
    merged data = pd.merge(sales data 1, sales data 2,
    on='Region')
    # Print the merged data
    print(merged data)
Output:
      Region Sales Profit
    0
        East
                100
                          20
    1
       West
               200
                          40
    2 North
               150
                          30
```

50 Example 11: Using Pandas to group data by a column and calculate summary statistics

250

```
import pandas as pd
# Create a DataFrame of sales data
data = pd.DataFrame({'Region': ['East', 'West',
'North', 'South', 'East', 'West', 'North', 'South'],
                     'Product': ['A', 'B', 'C', 'D',
'A', 'B', 'C', 'D'],
                     'Sales': [100, 200, 150, 250, 300,
400, 350, 450]})
```



3

South

Group the data by Region and Product and calculate the sum of Sales for each group grouped_data = data.groupby(['Region', 'Product'])['Sales'].sum() # Print the grouped data print(grouped_data)

Output:

Regior	n Produ	lct	
East	Α	40	0
	В	30	0
North	С	50	0
South	D	70	0
West	В	40	0
	Α	20	0
Name:	Sales,	dtype: in	nt64

Example 12: Using Pandas to sort a DataFrame by one or more columns

```
import pandas as pd
    # Create a DataFrame of sales data
    data = pd.DataFrame({'Region': ['East', 'West',
    'North', 'South', 'East', 'West', 'North', 'South'],
                          'Product': ['A', 'B', 'C', 'D',
    'A', 'B', 'C', 'D'],
                          'Sales': [100, 200, 150, 250, 300,
    400, 350, 450]})
    # Sort the data by Region (ascending) and Sales
     (descending)
    sorted data = data.sort values(by=['Region', 'Sales'],
    ascending=[True, False])
    # Print the sorted data
    print(sorted data)
Output:
      Region Product Sales
                         00
```

4	East	Α	300
0	East	A	100



5	West	В	400
1	West	В	200
6	North	С	350
2	North	С	150
7	South	D	450
3	South	D	250

These examples demonstrate some common data aggregation and summarization tasks that can be performed using Python libraries such as NumPy and Pandas. The flexibility and power of these libraries, combined with the ease of use and interoperability of Python, make it a valuable tool for data analysis and manipulation.

Merging and joining DataFrames

Merging and joining DataFrames is a common operation in data analysis, and it is no different in Python. If you are a SAS user transitioning to Python, this guide will give you an introduction to performing merges and joins in Python.

First, let's start with some basic definitions. A DataFrame in Python is similar to a SAS dataset, containing rows and columns of data. A merge operation combines two DataFrames based on a common column or set of columns, while a join operation combines two DataFrames based on a common index.

Let's start by importing the pandas library, which is the most commonly used library for working with DataFrames in Python:

import pandas as pd

Now, let's create two sample DataFrames to work with:

The two DataFrames have a common column called "key". We can perform a merge operation on this column by using the merge() function:

```
merged df = pd.merge(df1, df2, on='key')
```

The resulting merged DataFrame contains only the rows where the "key" column is present in both DataFrames:



	key	value_x	value	_ Y
0	В	_2		5
1	D	4		6

Note that the merge() function automatically appends "_x" and "_y" to the column names of the original DataFrames to distinguish them in the merged DataFrame.

If the column names are different in the two DataFrames, we can specify them explicitly using the left_on and right_on parameters:

```
df1 = pd.DataFrame({'key1': ['A', 'B', 'C', 'D'],
                                'value1': [1, 2, 3, 4]})
df2 = pd.DataFrame({'key2': ['B', 'D', 'E', 'F'],
                                 'value2': [5, 6, 7, 8]})
merged_df = pd.merge(df1, df2, left_on='key1',
    right_on='key2')
```

The resulting merged DataFrame looks like this:

	key1	value1	key2	value2
0	В	2	В	5
1	D	4	D	6

We can also perform a join operation by using the join() function. This function joins two DataFrames based on their indexes:

The resulting joined DataFrame contains only the rows where the index is present in both DataFrames:

	value1	value2
В	2	5
D	4	6

Note that the join() function also accepts a "how" parameter, which specifies the type of join to perform (inner, outer, left, or right).



Merge and Join Types

In Python, there are four types of merge operations: inner, left, right, and outer. The default is inner, which returns only the intersection of the two DataFrames based on the merge keys. The left and right types return all the rows from one DataFrame and matching rows from the other. The outer type returns all the rows from both DataFrames, and fills in missing values with NaNs.

Similarly, there are three types of join operations: inner, outer, and left. The inner type returns only the rows that have matching keys in both DataFrames. The outer type returns all the rows from both DataFrames and fills in missing values with NaNs. The left type returns all the rows from the left DataFrame and matching rows from the right DataFrame.

Using Multiple Keys

If you need to merge or join DataFrames using multiple keys, you can pass a list of keys to the on parameter:

One important thing to note is that merge() and join() functions can handle multiple keys for joining DataFrames. For example, consider the following DataFrames:

We can perform a merge operation on both "key1" and "key2" columns by passing a list of keys to the on parameter:

```
merged df = pd.merge(df1, df2, on=['key1', 'key2'])
```

The resulting merged DataFrame contains only the rows where both "key1" and "key2" columns are present in both DataFrames:



	key1	key2	value_x	value_y
0	В	X	_2	5
1	D	Z	4	6

We can also perform a join operation on multiple indexes by passing a list of indexes to the join() function:

Note that we also pass a "how" parameter to specify the type of join to perform.

Another thing to keep in mind is that merge() and join() functions can handle duplicate keys. For example, consider the following DataFrames:

```
df1 = pd.DataFrame({'key': ['A', 'B', 'B', 'C', 'D'],
                                'value': [1, 2, 3, 4, 5]})
df2 = pd.DataFrame({'key': ['B', 'B', 'D', 'E', 'F'],
                                 'value': [6, 7, 8, 9, 10]})
We can perform a merge operation on the "key" column,
even though it contains duplicate values, by using the
how parameter to specify how to handle duplicates:
python
Copy code
```

```
merged df = pd.merge(df1, df2, on='key', how='outer')
```

The resulting merged DataFrame contains all the rows from both DataFrames, and it fills missing values with NaN:

	key	value_x	value_y
0	Α	1.0	NaN
1	В	2.0	6.0
2	В	3.0	7.0
3	С	4.0	NaN
4	D	5.0	8.0
5	E	NaN	9.0
6	F	NaN	10.0



Chapter 5: Data Visualization with Matplotlib



Data visualization is a crucial aspect of data analysis, and it is essential to convey the results of the analysis to stakeholders effectively. Python is a popular programming language that is widely used for data analysis, and Matplotlib is a popular data visualization library in Python. In this article, we will provide a SAS-oriented introduction to Matplotlib for SAS users who are new to Python.

What is Matplotlib?

Matplotlib is a data visualization library in Python that provides a wide range of tools for creating static, animated, and interactive visualizations in Python. It is highly customizable and offers a vast array of visualization options, including line charts, scatter plots, histograms, bar charts, and more.

Getting started with Matplotlib:

To get started with Matplotlib, you first need to install it. You can install Matplotlib using the pip package manager by running the following command in your terminal or command prompt:

```
pip install matplotlib
```

Once you have installed Matplotlib, you can start using it in your Python scripts by importing the library using the following command:

```
import matplotlib.pyplot as plt
Creating a Line Chart with Matplotlib:
```

To create a line chart using Matplotlib, you can use the plot function. The plot function takes two arrays as input - one for the x-axis values and one for the y-axis values. Here's an example of how to create a simple line chart in Matplotlib:

```
import matplotlib.pyplot as plt
# Define the x and y values
x = [1, 2, 3, 4, 5]
y = [10, 8, 6, 4, 2]
# Create the line chart
plt.plot(x, y)
# Add labels and a title
plt.xlabel('X-axis')
plt.ylabel('Y-axis')
plt.title('Line Chart')
# Show the chart
plt.show()
Creating a Scatter Plot with Matplotlib:
```



To create a scatter plot using Matplotlib, you can use the scatter function. The scatter function takes two arrays as input - one for the x-axis values and one for the y-axis values. Here's an example of how to create a simple scatter plot in Matplotlib:

```
import matplotlib.pyplot as plt
# Define the x and y values
x = [1, 2, 3, 4, 5]
y = [10, 8, 6, 4, 2]
# Create the scatter plot
plt.scatter(x, y)
# Add labels and a title
plt.xlabel('X-axis')
plt.ylabel('Y-axis')
plt.title('Scatter Plot')
# Show the plot
plt.show()
```

Creating a Histogram with Matplotlib:

To create a histogram using Matplotlib, you can use the hist function. The hist function takes an array as input and creates a histogram by grouping the values into bins. Here's an example of how to create a simple histogram in Matplotlib:

```
import matplotlib.pyplot as plt
# Define the data
data = [1, 2, 2, 3, 3, 3, 4, 4, 5]
# Create the histogram
plt.hist(data)
# Add labels and a title
plt.xlabel('Value')
plt.ylabel('Frequency')
plt.title('Histogram')
# Show the plot
plt.show()
```

Matplotlib is a powerful data visualization library in Python that offers a wide range of tools for creating static, animated, and interactive visualizations. In this article, we provided a SAS-



oriented introduction to Matplotlib for SAS users who are new to Python. We hope this article helps you

Here's some more information on Matplotlib for SAS users:

- 1. Comparison to SAS Graphics: Matplotlib is a powerful data visualization library that offers a wide range of tools and customization options for creating high-quality graphics. While SAS also provides powerful graphics capabilities, Matplotlib provides a more flexible and customizable approach to data visualization, allowing users to create a wider variety of charts and graphs than what is available in SAS.
- 2. Integration with Pandas: Pandas is a popular data analysis library in Python that provides powerful tools for data manipulation and analysis. Matplotlib is fully integrated with Pandas, allowing users to create visualizations directly from their Pandas data frames. This integration makes it easy to create charts and graphs from data sets and quickly explore the relationships between different variables.
- 3. Customization Options: Matplotlib provides a wide range of customization options, allowing users to create highly personalized and professional-looking visualizations. Users can customize the color, font, size, and style of their charts, as well as add annotations, labels, and titles to their visualizations. This level of customization makes it easy to create charts and graphs that are tailored to the specific needs of a project or presentation.
- 4. Learning Curve: While Matplotlib is a powerful tool for data visualization, it can also have a steep learning curve for SAS users who are new to Python. However, with practice and experimentation, users can quickly become proficient in Matplotlib and take advantage of its powerful capabilities.
- 5. Matplotlib and Other Python Libraries: Matplotlib is just one of many data visualization libraries available in Python. Other popular libraries include Seaborn, Plotly, and Bokeh. These libraries offer different styles and approaches to data visualization and may be better suited to specific use cases or projects. However, Matplotlib remains one of the most widely used and versatile data visualization libraries in Python, and is a great tool for SAS users who are looking to expand their data visualization skills in Python.

Customizing Plots with Matplotlib:

Matplotlib provides many options for customizing the appearance of plots. You can customize the colors, line styles, markers, fonts, and many other aspects of your plots. Here are some examples of how you can customize your plots with Matplotlib:

Change the color of a line chart:

plt.plot(x, y, color='red')

Change the line style of a line chart:

```
plt.plot(x, y, linestyle='--')
```



Change the marker style of a scatter plot:

```
plt.scatter(x, y, marker='x')
```

Change the font size of labels and title:

```
plt.xlabel('X-axis', fontsize=14)
plt.ylabel('Y-axis', fontsize=14)
plt.title('Line Chart', fontsize=16)
Change the background color of a plot:
python
Copy code
plt.rcParams['axes.facecolor'] = 'lightgray'
```

Subplots in Matplotlib:

You can create multiple plots in a single figure using subplots in Matplotlib. Subplots are useful when you want to compare multiple plots side by side or when you want to create a complex layout of plots. Here's an example of how to create subplots in Matplotlib:

```
import matplotlib.pyplot as plt
# Define the data
\mathbf{x} = [1, 2, 3, 4, 5]
y1 = [10, 8, 6, 4, 2]
y^2 = [2, 4, 6, 8, 10]
# Create a figure with two subplots
fig, (ax1, ax2) = plt.subplots(nrows=1, ncols=2,
figsize=(10, 5))
# Plot the data in the first subplot
ax1.plot(x, y1)
ax1.set xlabel('X-axis')
ax1.set ylabel('Y-axis')
ax1.set title('Line Chart')
# Plot the data in the second subplot
ax2.scatter(x, y2)
ax2.set xlabel('X-axis')
ax2.set ylabel('Y-axis')
ax2.set title('Scatter Plot')
# Show the figure
plt.show()
```



Saving Plots in Matplotlib:

You can save your Matplotlib plots as image files, such as PNG, PDF, SVG, or EPS. To save a plot, you can use the savefig function. Here's an example of how to save a plot as a PNG file:

```
import matplotlib.pyplot as plt
# Define the data
x = [1, 2, 3, 4, 5]
y = [10, 8, 6, 4, 2]
# Create the plot
plt.plot(x, y)
plt.xlabel('X-axis')
plt.ylabel('Y-axis')
plt.title('Line Chart')
# Save the plot as a PNG file
plt.savefig('line_chart.png')
Introduction to Matplotlib
Line Chart:
```

```
import matplotlib.pyplot as plt
# Define the data
x = [1, 2, 3, 4, 5]
y = [10, 8, 6, 4, 2]
# Create the plot
plt.plot(x, y)
# Add labels and title
plt.xlabel('X-axis')
plt.ylabel('Y-axis')
plt.title('Line Chart')
# Show the plot
plt.show()
```

This code creates a simple line chart with the data provided in the x and y lists. The plot function is used to create the line chart, and the xlabel, ylabel, and title functions are used to add labels and a title to the plot. Finally, the show function is used to display the plot.



Scatter Plot:

```
import matplotlib.pyplot as plt
import numpy as np
# Generate random data
x = np.random.randn(100)
y = np.random.randn(100)
# Create the plot
plt.scatter(x, y)
# Add labels and title
plt.xlabel('X-axis')
plt.ylabel('Y-axis')
plt.title('Scatter Plot')
# Show the plot
plt.show()
```

This code creates a scatter plot with 100 randomly generated data points using the scatter function. The xlabel, ylabel, and title functions are used to add labels and a title to the plot, and the show function is used to display the plot.

Bar Chart:

```
import matplotlib.pyplot as plt
# Define the data
x = ['A', 'B', 'C', 'D', 'E']
y = [10, 8, 6, 4, 2]
# Create the plot
plt.bar(x, y)
# Add labels and title
plt.xlabel('Category')
plt.ylabel('Value')
plt.title('Bar Chart')
# Show the plot
plt.show()
```



This code creates a bar chart with the data provided in the x and y lists using the bar function. The xlabel, ylabel, and title functions are used to add labels and a title to the plot, and the show function is used to display the plot.

Histogram:

```
import matplotlib.pyplot as plt
import numpy as np
# Generate random data
x = np.random.randn(1000)
# Create the plot
plt.hist(x, bins=20)
# Add labels and title
plt.xlabel('Value')
plt.ylabel('Frequency')
plt.title('Histogram')
# Show the plot
plt.show()
```

This code creates a histogram with 1000 randomly generated data points using the hist function. The bins parameter is used to specify the number of bins in the histogram. The xlabel, ylabel, and title functions are used to add labels and a title to the plot, and the show function is used to display the plot.

Boxplot:

```
import matplotlib.pyplot as plt
import numpy as np
# Generate random data
np.random.seed(123)
data = [np.random.normal(0, std, 100) for std in
range(1, 4)]
# Create the plot
plt.boxplot(data, labels=['Group 1', 'Group 2', 'Group
3'])
# Add labels and title
plt.xlabel('Group')
```



```
plt.ylabel('Value')
plt.title('Boxplot')
# Show the plot
plt.show()
```

This code creates a boxplot with three groups of data, where each group contains 100 randomly generated data points with increasing standard deviation. The boxplot function is used to create the boxplot, and the labels parameter is used to specify the labels for each group. The xlabel, ylabel, and title functions are used to add labels and a title to the plot, and the show function is used to display the plot.

Pie Chart:

```
import matplotlib.pyplot as plt
# Define the data
sizes = [40, 30, 20, 10]
labels = ['Group 1', 'Group 2', 'Group 3', 'Group 4']
# Create the plot
plt.pie(sizes, labels=labels, autopct='%1.1f%%')
# Add title
plt.title('Pie Chart')
# Show the plot
plt.show()
```

This code creates a pie chart with four groups of data and their corresponding sizes using the pie function. The labels parameter is used to specify the labels for each group, and the autopct parameter is used to display the percentage of each group on the chart. The title function is used to add a title to the plot, and the show function is used to display the plot.

Heatmap:

```
import matplotlib.pyplot as plt
import numpy as np
# Generate random data
np.random.seed(123)
data = np.random.rand(10, 10)
# Create the plot
```



```
plt.imshow(data, cmap='hot', interpolation='nearest')
# Add colorbar and title
plt.colorbar()
plt.title('Heatmap')
# Show the plot
plt.show()
```

This code creates a heatmap with a 10x10 matrix of randomly generated data using the imshow function. The cmap parameter is used to specify the color map, and the interpolation parameter is used to specify the interpolation method for the heatmap. The colorbar function is used to add a colorbar to the plot, and the title function is used to add a title to the plot. Finally, the show function is used to display the plot.

Basic plots (line, scatter, bar)

Line Plots

A line plot is a graph that shows the relationship between two variables by connecting data points with a line. This type of plot is commonly used to show trends over time or to compare multiple sets of data. To create a line plot in Python, we can use the Matplotlib library.

Here is an example of how to create a line plot using Matplotlib:

```
import matplotlib.pyplot as plt
import numpy as np
# Create data
x = np.arange(0, 10, 0.1)
y = np.sin(x)
# Create line plot
plt.plot(x, y)
# Add title and axis labels
plt.title("Sine Wave")
plt.xlabel("x")
plt.ylabel("y")
# Show plot
```



plt.show()

In this example, we first import the Matplotlib library and the NumPy library, which provides a way to create arrays of data. We then create two arrays, x and y, using NumPy. The arange function creates an array of values from 0 to 10 in increments of 0.1 for x, and sin function calculates the sine of each value of x for y.

Next, we create a line plot using the plt.plot function, passing in x and y as arguments. We then add a title and axis labels using the plt.title, plt.xlabel, and plt.ylabel functions, respectively. Finally, we show the plot using the plt.show function.

Scatter Plots

A scatter plot is a graph that shows the relationship between two variables by plotting individual data points. This type of plot is commonly used to visualize the correlation between two variables. To create a scatter plot in Python, we can use the Matplotlib library.

Here is an example of how to create a scatter plot using Matplotlib:

```
import matplotlib.pyplot as plt
import numpy as np
# Create data
x = np.random.rand(100)
y = np.random.rand(100)
colors = np.random.rand(100)
sizes = 1000 * np.random.rand(100)
# Create scatter plot
plt.scatter(x, y, c=colors, s=sizes)
# Add title and axis labels
plt.title("Random Data")
plt.xlabel("x")
plt.ylabel("y")
# Show plot
plt.show()
```

In this example, we first import the Matplotlib library and the NumPy library, which provides a way to create arrays of data. We then create three arrays, x, y, and colors, using NumPy. The random.rand function creates arrays of random values between 0 and 1 for x, y, and colors.

We also create a sizes array, which contains 1000 random values between 0 and 1. This array is used to set the size of each data point in the scatter plot.



Next, we create a scatter plot using the plt.scatter function, passing in x, y, c, and s as arguments. The c argument sets the color of each data point based on its corresponding value in the colors array, and the s argument sets the size of each data point based on its corresponding value in the sizes array.

We then add a title and axis labels using the plt.title, plt.xlabel, and plt.ylabel functions, respectively. Finally, we show the plot using the plt.show function.

Bar Plots

A bar plot is a graph that shows the distribution of a categorical variable or the comparison between multiple sets of data. This type of plot is commonly used to visualize the frequency of different categories or to compare the values of a variable across different groups. To create a bar plot in Python, we can use the Matplotlib library.

Here is an example of how to create a bar plot using Matplotlib:

```
import matplotlib.pyplot as plt
import numpy as np
# Create data
x = ['A', 'B', 'C', 'D', 'E']
y = [3, 5, 2, 6, 4]
# Create bar plot
plt.bar(x, y)
# Add title and axis labels
plt.title("Bar Plot Example")
plt.xlabel("Categories")
plt.ylabel("Frequency")
# Show plot
plt.show()
```

In this example, we first import the Matplotlib library and the NumPy library, which provides a way to create arrays of data. We then create two arrays, x and y, representing the categories and their corresponding frequencies.

Next, we create a bar plot using the plt.bar function, passing in x and y as arguments. We then add a title and axis labels using the plt.title, plt.xlabel, and plt.ylabel functions, respectively. Finally, we show the plot using the plt.show function.

Advantages of using Python for data visualization



Python offers several advantages over SAS for data visualization. First, Python has a larger and more active community, which means that there are more resources and support available for data analysts who use Python. Second, Python provides more flexibility and customization options for data visualization. With libraries such as Matplotlib, Seaborn, and Plotly, data analysts can create a wide range of visualizations that are tailored to their specific needs. Third, Python is open-source, which means that it is free to use and can be customized to suit the needs of individual users. Finally, Python integrates well with other data analysis and machine learning tools, which makes it a versatile and powerful tool for data analysis.

In this article, we have discussed some basic plots that can be created using Python, including line plots, scatter plots, and bar plots. We have provided examples of how to create each of these plots using the Matplotlib library and have discussed some of the advantages of using Python for data visualization. By learning how to create these basic plots, SAS users can begin to explore the benefits of using Python for data analysis and visualization.

Line Plots

Line plots are used to display data points that are connected by lines. They are useful for visualizing trends over time or for comparing the values of different variables. To create a line plot in Python, we can use the Matplotlib library.

Here is an example of how to create a line plot using Matplotlib:

```
import matplotlib.pyplot as plt
import numpy as np
# Create data
x = np.linspace(0, 10, 100)
Line Plot
import matplotlib.pyplot as plt
import numpy as np
# Create data
x = np.linspace(0, 10, 100)
y = np.sin(x)
# Create line plot
```

```
plt.plot(x, y)
```

```
# Add title and axis labels
plt.title("Line Plot Example")
plt.xlabel("x")
plt.ylabel("y")
```



Show plot
plt.show()

In this example, we use the numpy.linspace function to create an array of 100 equally spaced values between 0 and 10. We then use the numpy.sin function to compute the sine of each value in the array. Next, we use the plt.plot function to create a line plot of the data, passing in x and y as arguments. We then add a title and axis labels using the plt.title, plt.xlabel, and plt.ylabel functions, respectively. Finally, we show the plot using the plt.show function.

Scatter Plot

```
import matplotlib.pyplot as plt
import numpy as np
# Create data
x = np.random.rand(50)
y = np.random.rand(50)
colors = np.random.rand(50)
sizes = 1000 * np.random.rand(50)
# Create scatter plot
plt.scatter(x, y, c=colors, s=sizes)
# Add title and axis labels
plt.title("Scatter Plot Example")
plt.xlabel("x")
plt.ylabel("y")
# Show plot
plt.show()
```

In this example, we use the numpy.random.rand function to create arrays of 50 random values between 0 and 1 for x, y, colors, and sizes. We then use the plt.scatter function to create a scatter plot of the data, passing in x, y, colors, and sizes as arguments. We use the c argument to set the color of each data point based on its corresponding value in the colors array, and the s argument sets the size of each data point based on its corresponding value in the sizes array. We then add a title and axis labels using the plt.title, plt.xlabel, and plt.ylabel functions, respectively. Finally, we show the plot using the plt.show function.

Bar Plot

```
import matplotlib.pyplot as plt
import numpy as np
```



```
# Create data
x = ['A', 'B', 'C', 'D', 'E']
y = [3, 5, 2, 6, 4]
# Create bar plot
fig, ax = plt.subplots()
ax.bar(x, y)
# Add title and axis labels
ax.set_title("Bar Plot Example")
ax.set_xlabel("Categories")
ax.set_ylabel("Frequency")
# Show plot
plt.show()
```

In this example, we use the numpy.arange function to create an array of five values, which we assign to x. We also create an array of five values representing the frequency of each category, which we assign to y. We then use the plt.subplots function to create a figure and axis object, which we assign to fig and ax, respectively. We use the ax.bar function to create a bar plot of the data, passing in x and y as arguments. We then use the ax.set_title,

```
import matplotlib.pyplot as plt
import numpy as np
# Create data
x = np.linspace(0, 10, 100)
y = np.sin(x)
# Create line plot
plt.plot(x, y)
# Add title and axis labels
plt.title("Line Plot Example")
plt.xlabel("X-axis")
plt.ylabel("Y-axis")
# Show plot
plt.show()
```

In this example, we first import the Matplotlib library and the NumPy library, which provides a way to create arrays of data. We then create two arrays, x and y, representing the x and y coordinates of the data points.



Next, we create a line plot using the plt.plot function, passing in x and y as arguments. We then add a title and axis labels using the plt.title, plt.xlabel, and plt.ylabel functions, respectively. Finally, we show the plot using the plt.show function.

Scatter Plots

Scatter plots are used to visualize the relationship between two variables. They are useful for identifying patterns or trends in data, and for identifying outliers or other unusual observations. To create a scatter plot in Python, we can use the Matplotlib library.

Here is an example of how to create a scatter plot using Matplotlib:

```
import matplotlib.pyplot as plt
import numpy as np
# Create data
x = np.random.normal(0, 1, 100)
y = np.random.normal(0, 1, 100)
colors = np.random.rand(100)
sizes = 1000 * np.random.rand(100)
# Create scatter plot
plt.scatter(x, y, c=colors, s=sizes)
# Add title and axis labels
plt.title("Scatter Plot Example")
plt.xlabel("X-axis")
plt.ylabel("Y-axis")
# Show plot
plt.show()
```

In this example, we first import the Matplotlib library and the NumPy library, which provides a way to create arrays of data. We then create three arrays, x, y, and colors, representing the x and y coordinates of the data points and the colors of each data point. We also create a sizes array to set the size of each data point.

Next, we create a scatter plot using the plt.scatter function, passing in x, y, colors, and sizes as arguments. The c argument sets the color of each data point based on its corresponding value in the colors array, and the s argument sets the size of each data point based on its corresponding value in the sizes array.

We then add a title and axis labels using the plt.title, plt.xlabel, and plt.ylabel functions, respectively. Finally, we show the plot using the plt.show function.

Bar Plots



Bar plots are used to visualize the distribution of a categorical variable or to compare the values of a variable across different groups. They are useful for identifying differences or similarities between groups or categories. To create a bar plot in Python, we can use the Matplotlib library. Here is an example of how to create a bar plot using Matplotlib:

```
import matplotlib.pyplot as plt
import numpy as np
# Create data
x = ['A', 'B', 'C', 'D', 'E']
y = [3, 5, 2, 6, 4]
# Create bar plot
plt.bar(x, y)
# Add title and axis labels
plt.title("Bar Plot Example")
plt.xlabel("Categories")
Customizing plots (labels, colors, styles)
```

Multiple plots on one graph

In Python, you can use the Matplotlib library to create multiple plots on one graph. Matplotlib is a powerful visualization library that provides many options for creating plots and charts.

To get started, you will need to import the necessary libraries. In addition to Matplotlib, you will also need to import NumPy, which is a library for numerical computing.

python Copy code import matplotlib.pyplot as plt import numpy as np Next, you can create your data. For example, let's say you want to create two sets of data, one for x values and one for y values.

```
makefile
Copy code
x1 = np.linspace(0, 10, 100)
y1 = np.sin(x1)
```

x2 = np.linspace(0, 10, 100)y2 = np.cos(x2)



In this example, we are creating two sets of data, one for the sine function and one for the cosine function. We are using the NumPy linspace function to create 100 evenly spaced values between 0 and 10 for each set of data.

To create a graph with multiple plots, you can use the subplot function in Matplotlib. The subplot function takes three arguments: the number of rows, the number of columns, and the plot number.

```
plt.subplot(2, 1, 1)
plt.plot(x1, y1)
plt.title('Sine Function')
plt.subplot(2, 1, 2)
plt.plot(x2, y2)
plt.title('Cosine Function')
plt.show()
```

In this example, we are creating a graph with two plots, one above the other. The subplot function is used to specify that there will be two rows and one column of plots, and that the first plot will be in position 1 and the second plot will be in position 2.

We then use the plot function to plot each set of data on its corresponding plot. We also use the title function to give each plot a title.

Finally, we use the show function to display the graph.

By using the subplot function, you can create a graph with multiple plots in Python using Matplotlib. This can be useful for comparing different sets of data or for displaying related data on the same graph.

When creating multiple plots on one graph in Matplotlib, you can also customize the appearance of the graph to meet your needs. Here are some common customization options:

- 1. Changing the figure size: You can use the **figure** function to specify the size of the figure. For example, you can create a figure with a width of 8 inches and a height of 6 inches by using **plt.figure(figsize=(8, 6))**.
- 2. Adding a legend: If you have multiple plots on one graph, you may want to add a legend to differentiate between them. You can use the **legend** function to add a legend to the graph. For example, you can add a legend with the labels "Sine" and "Cosine" by using **plt.legend([''Sine'', ''Cosine''])**.
- 3. Changing the axis labels: You can use the **xlabel** and **ylabel** functions to set the labels for the x and y axes. For example, you can set the x-axis label to "Time" and the y-axis label to "Amplitude" by using **plt.xlabel(''Time'')** and **plt.ylabel(''Amplitude'')**.
- 4. Changing the plot colors: By default, Matplotlib uses a different color for each plot. However, you can also specify the color of each plot using the **color** parameter. For



example, you can plot the sine function in blue and the cosine function in red by using plt.plot(x1, y1, color="blue") and plt.plot(x2, y2, color="red").

- 5. Changing the plot styles: You can use the linestyle and marker parameters to change the style of the plot. For example, you can plot the sine function as a dashed line and the cosine function as a dotted line by using plt.plot(x1, y1, linestyle=''--'') and plt.plot(x2, y2, linestyle='':'').
- 6. Adding a grid: You can use the **grid** function to add a grid to the graph. For example, you can add a grid by using **plt.grid**(**True**).

By customizing the appearance of the graph, you can create a more informative and visually appealing visualization. Matplotlib provides a wide range of customization options, so you can tailor your graph to meet your specific needs.

```
import numpy as np
import matplotlib.pyplot as plt
# Generate data
x = np.linspace(0, 2*np.pi, 100)
y1 = np.sin(x)
y2 = np.cos(x)
# Create figure and axis objects
fig, ax = plt.subplots()
# Plot data
ax.plot(x, y1, color='blue', label='Sine', linestyle='-
- ' )
ax.plot(x, y2, color='red', label='Cosine',
linestyle=':')
# Add legend
ax.legend()
# Add title and axis labels
ax.set title('Sine and Cosine Functions')
ax.set xlabel('X-axis')
ax.set ylabel('Y-axis')
# Add grid
ax.grid(True)
# Display graph
plt.show()
```



In this code, we first generate data for the sine and cosine functions using the linspace and sin/cos functions from NumPy.

We then create a figure object and an axis object using the subplots function.

Next, we plot the sine and cosine functions on the same graph using the plot function. We specify the color, label, linestyle, and marker for each plot.

We add a legend to differentiate between the sine and cosine functions using the legend function.

We also add a title and axis labels using the set_title, set_xlabel, and set_ylabel functions.

Finally, we add a grid to the graph using the grid function, and display the graph using the show function.

Another example with multiple subplots:

```
import numpy as np
import matplotlib.pyplot as plt
# Generate data
x = np.linspace(0, 2*np.pi, 100)
y1 = np.sin(x)
y^2 = np.cos(x)
# Create figure and axis objects
fig, axs = plt.subplots(2, 1)
# Plot data
axs[0].plot(x, y1, color='blue', label='Sine',
linestyle='--')
axs[1].plot(x, y2, color='red', label='Cosine',
linestyle=':')
# Add legends
axs[0].legend()
axs[1].legend()
# Add titles and axis labels
axs[0].set title('Sine Function')
axs[1].set title('Cosine Function')
axs[1].set xlabel('X-axis')
axs[0].set ylabel('Y-axis')
axs[1].set ylabel('Y-axis')
```



```
# Add grid
axs[0].grid(True)
axs[1].grid(True)
# Display graph
plt.show()
```

In this code, we first generate data for the sine and cosine functions using the linspace and sin/cos functions from NumPy.

We then create a figure object and an array of two axis objects using the subplots function.

Next, we plot the sine and cosine functions on their corresponding plots using the plot function. We specify the color, label, linestyle, and marker for each plot.

We add legends to differentiate between the sine and cosine functions using the legend function.

We also add titles and axis labels using the set_title, set_xlabel, and set_ylabel functions.

Finally, we add a grid to each plot using the grid function, and display the graph using the show function.

```
import numpy as np
import matplotlib.pyplot as plt
# Generate data
x = np.linspace(0, 2*np.pi, 100)
y1 = np.sin(x)
y^2 = np.cos(x)
# Create figure and axis objects
fig, ax = plt.subplots()
# Plot data
ax.plot(x, y1, color='blue', label='Sine', linestyle='-
- ' )
ax.plot(x, y2, color='red', label='Cosine',
linestyle=':')
# Add legend
ax.legend()
# Add title and axis labels
ax.set title('Sine and Cosine Functions')
ax.set xlabel('X-axis')
```



```
ax.set ylabel('Y-axis')
# Add grid
ax.grid(True)
# Customize ticks and tick labels
ax.set xticks([0, np.pi/2, np.pi, 3*np.pi/2, 2*np.pi])
ax.set xticklabels(['0', '\pi/2', '\pi', '3\pi/2', '2\pi'])
ax.set yticks([-1, 0, 1])
ax.set yticklabels(['-1', '0', '1'])
# Add vertical and horizontal lines
ax.axvline(np.pi/2, color='black', linestyle='--')
ax.axvline(3*np.pi/2, color='black', linestyle='--')
ax.axhline(0, color='black', linestyle='--')
# Add shaded region
ax.fill between(x, y1, y2, where=(y1<y2),</pre>
color='green', alpha=0.3, interpolate=True)
# Display graph
plt.show()
```

In this code, we first generate data for the sine and cosine functions using the linspace and sin/cos functions from NumPy.

We then create a figure object and an axis object using the subplots function. Next, we plot the sine and cosine functions on the same graph using the plot function. We specify the color, label, linestyle, and marker for each plot.

We add a legend to differentiate between the sine and cosine functions using the legend function.

We also add a title and axis labels using the set_title, set_xlabel, and set_ylabel functions.

We add a grid to the graph using the grid function.

We then customize the tick locations and labels using the set_xticks, set_xticklabels, set_yticks, and set_yticklabels functions.

We add vertical and horizontal lines using the axvline and axhline functions.

We also add a shaded region between the sine and cosine functions using the fill_between function. We specify the where argument to only shade the region where the sine function is less than the cosine function.



Finally, we display the graph using the show function.

```
import numpy as np
import matplotlib.pyplot as plt
# Generate data
x = np.linspace(0, 2*np.pi, 100)
y1 = np.sin(x)
y^2 = np.exp(x)
# Create figure and axis objects
fig, ax1 = plt.subplots()
# Create second axis object with different scale
ax2 = ax1.twinx()
# Plot data on both axes
ax1.plot(x, y1, color='blue', label='Sine')
ax2.plot(x, y2, color='red', label='Exponential')
# Add legend
ax1.legend(loc='upper left')
ax2.legend(loc='upper right')
# Add title and axis labels
ax1.set title('Sine and Exponential Functions')
ax1.set xlabel('X-axis')
ax1.set ylabel('Sine', color='blue')
ax2.set ylabel('Exponential', color='red')
# Add grid
ax1.grid(True)
# Display graph
plt.show()
```

In this code, we first generate data for the sine and exponential functions using the linspace, sin, and exp functions from NumPy.

We then create a figure object and an axis object using the subplots function.

We create a second axis object with a different scale using the twinx method.



Next, we plot the sine function on the left axis and the exponential function on the right axis using the plot function. We specify the color and label for each plot.

We add legends to both axes using the legend function. We specify the location of each legend using the loc argument.

We also add a title and axis labels to both axes using the set_title, set_xlabel, and set_ylabel functions.

Subplots and grids

Subplots and grids are important tools in data visualization using Python, especially when you have multiple plots that you want to display together. In this article, we will explore subplots and grids in Python, and how they can be used for data visualization.

A subplot is a plot that is embedded within another plot, allowing multiple plots to be displayed within the same figure. This is useful when you want to display multiple plots together, or when you want to compare different aspects of your data. Subplots can be created using the plt.subplots() method from the matplotlib library.

To create subplots, you first need to define the number of rows and columns that you want in your subplot grid. For example, if you want a grid with two rows and three columns, you would call the plt.subplots() method with the arguments nrows=2 and ncols=3. This will create a figure with six subplots arranged in a 2x3 grid.

Once you have created your subplot grid, you can use the ax parameter to access each individual subplot. For example, if you want to plot data on the first subplot, you would use ax[0,0]. The first index corresponds to the row number, and the second index corresponds to the column number.

Here is an example of how to create a subplot grid and plot data on each subplot:

```
import matplotlib.pyplot as plt
import numpy as np
# Define data
x = np.linspace(0, 10, 100)
y1 = np.sin(x)
y2 = np.cos(x)
y3 = np.tan(x)
y4 = np.exp(x)
y5 = np.log(x)
```



```
# Create subplot grid with 2 rows and 3 columns
fig, ax = plt.subplots(nrows=2, ncols=3)
# Plot data on first subplot
ax[0,0].plot(x, y1)
ax[0,0].set title('sin(x)')
# Plot data on second subplot
ax[0,1].plot(x, y2)
ax[0,1].set title('cos(x)')
# Plot data on third subplot
ax[0,2].plot(x, y3)
ax[0,2].set title('tan(x)')
# Plot data on fourth subplot
ax[1,0].plot(x, y4)
ax[1,0].set title('exp(x)')
# Plot data on fifth subplot
ax[1,1].plot(x, y5)
ax[1,1].set title('log(x)')
# Remove unused subplots
fig.delaxes(ax[1,2])
# Adjust spacing between subplots
fig.tight layout()
# Display figure
plt.show()
```

In this example, we create a subplot grid with two rows and three columns, and plot five different functions on each subplot.

Note that in this example, we also remove the unused subplot in the second row and third column using the fig.delaxes() method. This is because we only want to display five subplots, and not all six subplots in the grid.

The fig.tight_layout() method is used to adjust the spacing between the subplots, so that they do not overlap.

Grids are similar to subplots, but they allow for more customization of the layout. Grids can be created using the GridSpec class from the matplotlib.gridspec module. With grids, you can specify the size and location of each subplot within the grid.

Define data

```
x = np.linspace(0, 10, 100)
y1 = np.sin(x)
y2 = np.cos(x)
y3 = np.tan(x)
y4 = np.exp(x)
y5 = np.log(x)
```

Create grid with 2 rows and 3 columns

```
fig = plt.figure()
gs = gridspec.GridSpec(2, 3)
```

Create subplots within grid

```
ax1 = fig.add_subplot(gs[0, 0])
ax1.plot(x, y1)
ax1.set_title('sin(x)')
ax2 = fig.add_subplot(gs[0, 1])
ax2.plot(x, y2)
ax2.set_title('cos(x)')
ax3 = fig.add_subplot(gs[0, 2])
ax3.plot(x, y3)
ax3.set_title('tan(x)')
ax4 = fig.add_subplot(gs[1, :2])
ax4.plot(x, y4)
ax4.set_title('exp(x)')
ax5 = fig.add_subplot(gs[1, 2])
ax5.plot(x, y5)
ax5.set_title('log(x)')
```

Adjust spacing between subplots

```
fig.tight_layout()
```



In addition to creating subplots and grids, you can also customize the appearance of each subplot or grid using various parameters such as colors, line styles, markers, and labels. Here are some common parameters that you can use to customize your plots:

color: specifies the color of the plot. You can use a color name (e.g., 'red') or a hex code (e.g., '#FF0000').

linestyle: specifies the style of the plot line. You can use a string (e.g., '-' for solid line, '--' for dashed line, ':' for dotted line) or a tuple (e.g., (0, (1, 1)) for a dashed-dot line).

marker: specifies the shape of the markers used in the plot. You can use a string (e.g., 'o' for circles, '*' for stars, 's' for squares) or a tuple (e.g., (5, 1, 0) for a pentagon).

label: specifies the label for the plot. This is useful when you want to add a legend to your plot. Here is an example of how to use these parameters to customize a plot:

```
import matplotlib.pyplot as plt
import numpy as np
# Define data
x = np.linspace(0, 10, 100)
y = np.sin(x)
# Plot data with custom parameters
plt.plot(x, y, color='blue', linestyle='--',
marker='o', label='sin(x)')
# Add title and legend
plt.title('Custom Plot')
plt.legend()
# Display plot
plt.show()
```

In this example, we plot the sin(x) function with a blue dashed line, circular markers, and a label. We then add a title and a legend to the plot using the plt.title() and plt.legend() methods, respectively.

Another useful feature of subplots and grids is the ability to share axes. When you have multiple plots that share the same x-axis or y-axis, you can use the sharex or sharey parameters to ensure that the axes are synchronized across all plots. This can make it easier to compare the data in each plot and avoid any misinterpretation.

Here is an example of how to create subplots with shared axes:

```
import matplotlib.pyplot as plt
import numpy as np
```



```
# Define data
x = np.linspace(0, 10, 100)
y1 = np.sin(x)
y^2 = np.cos(x)
# Create subplots with shared x-axis
fig, (ax1, ax2) = plt.subplots(nrows=2, ncols=1,
sharex=True)
# Plot data on first subplot
ax1.plot(x, y1)
ax1.set title('sin(x)')
# Plot data on second subplot
ax2.plot(x, y2)
ax2.set title('cos(x)')
# Add x-axis label
fig.text(0.5, 0.04, 'x', ha='center')
# Add y-axis label
fig.text(0.04, 0.5, 'y', va='center',
rotation='vertical')
# Display figure
plt.show()
```

In this example, we create two subplots that share the same x-axis using the sharex=True parameter. We then plot the sin(x) and cos(x) functions on each subplot and add a label to the x-axis and y-axis of the figure using the fig.text() method.

Overall, subplots and grids are essential tools for data visualization in Python, and they provide a lot of flexibility in terms of layout and customization. By using these tools effectively, you can create clear and informative visualizations that communicate insights and trends in your data. In the context of SAS users transitioning to Python, learning how to create subplots and grids in Python can be particularly useful for reproducing and improving existing SAS graphs.

Here is an example code snippet that demonstrates how to create a 2x2 grid of subplots using the subplots() function:

```
import matplotlib.pyplot as plt
import numpy as np
# create some sample data
```



```
\mathbf{x} = np.linspace(0, 10, 100)
y1 = np.sin(x)
y^2 = np.cos(x)
y3 = np.exp(x)
y4 = np.log(x)
# create a 2x2 grid of subplots
fig, axes = plt.subplots(nrows=2, ncols=2, figsize=(8,
6))
# plot each data series on its own subplot
axes[0, 0].plot(x, y1)
axes[0, 0].set title('sin(x)')
axes[0, 1].plot(x, y2)
axes[0, 1].set title('cos(x)')
axes[1, 0].plot(x, y3)
axes[1, 0].set title('exp(x)')
axes[1, 1].plot(x, y4)
axes[1, 1].set title('log(x)')
# customize the layout and appearance of the subplots
fig.tight layout(pad=3)
plt.subplots adjust(top=0.9)
plt.suptitle('Four Subplots')
# display the plot
plt.show()
```

In this example, we first create some sample data using NumPy. We then use the subplots() function to create a 2x2 grid of subplots and assign the resulting Figure and Axes objects to fig and axes, respectively. We then use indexing to access each subplot in the grid and plot each data series on its own subplot. Finally, we use various parameters to customize the layout and appearance of the subplots, add a main title to the figure, and display the plot using plt.show().

```
import matplotlib.pyplot as plt
import matplotlib.gridspec as gridspec
import numpy as np
# create some sample data
x = np.linspace(0, 10, 100)
y1 = np.sin(x)
y2 = np.cos(x)
y3 = np.exp(x)
y4 = np.log(x)
```



```
# create a 2x2 grid of subplots with different sizes
fig = plt.figure(figsize=(8, 6))
gs = gridspec.GridSpec(nrows=2, ncols=2,
width ratios=[1, 2], height ratios=[2, 1])
ax1 = fig.add subplot(gs[0, 0])
ax2 = fig.add subplot(gs[0, 1])
ax3 = fig.add subplot(gs[1, 0])
ax4 = fig.add subplot(gs[1, 1])
# plot each data series on its own subplot
ax1.plot(x, y1)
ax1.set title('sin(x)')
ax2.plot(x, y2)
ax2.set title('cos(x)')
ax3.plot(x, y3)
ax3.set title('exp(x)')
ax4.plot(x, y4)
ax4.set title('log(x)')
# customize the layout and appearance of the subplots
fig.tight layout(pad=3)
plt.subplots adjust(top=0.9)
plt.suptitle('Four Subplots with Different Sizes')
# display the plot
plt.show()
```

In this example, we use the GridSpec() function to create a 2x2 grid of subplots with different sizes. We specify the relative widths and heights of the subplots using the width_ratios and height_ratios parameters, respectively. We then use the add_subplot() method of the Figure object to add each subplot to the figure and assign them to ax1, ax2, ax3, and ax4. We then plot each data series on its own subplot, customize the layout and appearance of the subplots, add a main title to the figure, and display the plot using plt.show().

Here is another example code snippet that demonstrates how to create a grid of subplots with shared x- and y-axes using the subplots() function:

```
import matplotlib.pyplot as plt
import numpy as np
# create some sample data
x = np.linspace(0, 10, 100)
y1 = np.sin(x)
y2 = np.cos(x)
```



```
y3 = np.exp(x)
y4 = np.log(x)
# create a 2x2 grid of subplots with shared x- and y-
axes
fig, axes = plt.subplots(nrows=2, ncols=2, sharex=True,
sharey=True, figsize=(8, 6))
# plot each data series on its own subplot
axes[0, 0].plot(x, y1)
axes[0, 0].set title('sin(x)')
axes[0, 1].plot(x, y2)
axes[0, 1].set title('cos(x)')
axes[1, 0].plot(x, y3)
axes[1, 0].set title('exp(x)')
axes[1, 1].plot(x, y4)
axes[1, 1].set title('log(x)')
# customize the layout and appearance of the subplots
fig.tight layout(pad=3)
plt.subplots adjust(top=0.9)
plt.suptitle('Four Subplots with Shared Axes')
# display the plot
plt.show()
```

Here are some tips and best practices for using subplots and grids in Python:

- 1. Plan your layout: Before creating your subplots or grid, think about the overall layout you want to achieve. Consider the number of plots you need, their size, and how they should be arranged.
- 2. Use **subplots**() to create a grid of plots: The **subplots**() function is a convenient way to create a grid of plots with a specified number of rows and columns. You can then access each subplot using indexing.
- 3. Use **add_subplot**() to add a subplot to an existing figure: If you want to add a new subplot to an existing figure, you can use the **add_subplot**() method to create a new subplot in a specific location.
- 4. Customize each subplot: Use the various parameters available in Matplotlib to customize the appearance of each subplot, such as color, line style, and marker shape.
- 5. Use shared axes: When creating multiple subplots that share the same x-axis or y-axis, use the **sharex** or **sharey** parameters to ensure that the axes are synchronized across all plots.
- 6. Add titles and labels: Add a title to each subplot to describe the content of the plot, and add axis labels to clarify the meaning of the data.



- 7. Add a legend: If you have multiple lines in a plot, add a legend to identify which line corresponds to which data series. You can do this using the **label** parameter in the **plot**() function, and then calling **legend**() to display the legend.
- 8. Save and export your plots: Once you have created your subplots or grid, save your figure using **savefig()** or export it to a file format of your choice, such as PDF or PNG.

In summary, subplots and grids are powerful tools for creating visualizations in Python. By using them effectively, you can create clear and informative plots that communicate insights and trends in your data. For SAS users transitioning to Python, learning how to use these tools is essential for reproducing and improving existing SAS graphs, and for creating new and engaging visualizations in Python.

Advanced plots (heatmaps, histograms, box plots)

Heatmaps

Heatmaps are a type of plot used to visualize the relationship between two variables in a matrixlike format. They are particularly useful for visualizing large datasets, as they can display a large amount of data in a compact and easy-to-understand format. Heatmaps can be created in Python using the seaborn library, which provides a high-level interface for creating statistical graphics. To create a heatmap in Python, first import the seaborn library and load the data that you want to visualize. Then, use the heatmap() function to create the plot. Here is an example:

```
import seaborn as sns
import pandas as pd
data = pd.read_csv('mydata.csv')
sns.heatmap(data)
```

This will create a heatmap of the data in the mydata.csv file.

Histograms

Histograms are a type of plot used to visualize the distribution of a single variable. They are particularly useful for identifying patterns and outliers in data. Histograms can be created in Python using the matplotlib library, which provides a wide variety of tools for creating high-quality visualizations.

To create a histogram in Python, first import the matplotlib library and load the data that you want to visualize. Then, use the hist() function to create the plot. Here is an example:

```
import matplotlib.pyplot as plt
```



```
import pandas as pd
data = pd.read_csv('mydata.csv')
plt.hist(data['variable'])
```

This will create a histogram of the variable column in the mydata.csv file.

Box Plots

Box plots are a type of plot used to visualize the distribution of a single variable or the relationship between two variables. They are particularly useful for identifying outliers and comparing the distribution of different groups of data. Box plots can be created in Python using the matplotlib library.

To create a box plot in Python, first import the matplotlib library and load the data that you want to visualize. Then, use the boxplot() function to create the plot. Here is an example:

```
import matplotlib.pyplot as plt
import pandas as pd
data = pd.read_csv('mydata.csv')
plt.boxplot(data['variable'])
```

This will create a box plot of the variable column in the mydata.csv file. Heatmaps:

```
import seaborn as sns
import pandas as pd
data = pd.read_csv('mydata.csv')
sns.heatmap(data, cmap='coolwarm', annot=True,
fmt='.2f', linewidths=.5)
```

This code imports the seaborn library and loads the data from a CSV file called mydata.csv using pandas. The heatmap() function is then used to create the heatmap, with several customization options specified. The cmap argument sets the color scheme of the heatmap to "coolwarm", which ranges from blue to red. The annot argument adds annotations to the heatmap, and the fmt argument specifies that the annotations should be formatted to two decimal places. The linewidths argument sets the width of the lines separating the individual cells in the heatmap.

Histograms:

```
import matplotlib.pyplot as plt
import pandas as pd
```



```
data = pd.read_csv('mydata.csv')
plt.hist(data['variable'], bins=10, color='green',
alpha=.5)
plt.xlabel('Variable')
plt.ylabel('Frequency')
plt.title('Histogram of Variable')
plt.grid(axis='y', alpha=.5)
```

This code imports the matplotlib library and loads the data from a CSV file called mydata.csv using pandas. The hist() function is then used to create the histogram, with several customization options specified. The bins argument sets the number of bins in the histogram to 10, and the color argument sets the color of the bars to green with an opacity of 0.5. The xlabel, ylabel, and title functions are used to add labels and a title to the plot, and the grid() function is used to add a grid to the y-axis with an opacity of 0.5.

Box Plots:

```
import matplotlib.pyplot as plt
import pandas as pd
data = pd.read_csv('mydata.csv')
plt.boxplot(data['variable'], vert=False,
showfliers=False, widths=.5,
boxprops=dict(linewidth=2),
whiskerprops=dict(linewidth=2),
medianprops=dict(linewidth=2, color='red'))
plt.xlabel('Variable')
plt.title('Box Plot of Variable')
plt.grid(axis='x', alpha=.5)
```

This code imports the matplotlib library and loads the data from a CSV file called mydata.csv using pandas. The boxplot() function is then used to create the box plot, with several customization options specified. The vert argument sets the orientation of the plot to horizontal, and the showfliers argument hides any outliers that fall outside the whiskers. The widths argument sets the width of the boxes to 0.5, and the boxprops, whiskerprops, and medianprops arguments are used to adjust the style and width of the various components of the plot. The xlabel and title functions are used to add labels and a title to the plot, and the grid() function is used to add a grid to the x-axis with an opacity of 0.5.

Sure, here are some additional code examples with more details on customizing the advanced plots:

Heatmaps:

import seaborn as sns



```
import pandas as pd
data = pd.read_csv('mydata.csv')
sns.set(font_scale=1.2)  # Increase font size for
readability
sns.set_style('whitegrid')  # Add horizontal and
vertical grid lines
sns.heatmap(data.corr(), cmap='coolwarm', annot=True,
fmt='.2f', linewidths=.5, vmin=-1, vmax=1, center=0)
```

In this code example, the sns.set() function is used to increase the font size of the heatmap for readability. The sns.set_style() function is used to add horizontal and vertical grid lines to the plot. The heatmap() function is used to create the heatmap of the correlation matrix of the data, with several customization options specified. The vmin, vmax, and center arguments are used to set the range of values for the color scale, with values ranging from -1 to 1 and the center point set at 0.

Histograms:

```
import matplotlib.pyplot as plt
import pandas as pd
data = pd.read_csv('mydata.csv')
plt.hist(data['variable'], bins=20, color='green',
alpha=.5, edgecolor='black', linewidth=1.2)
plt.xlabel('Variable')
plt.ylabel('Frequency')
plt.title('Histogram of Variable')
plt.xticks(range(0, 101, 10))  # Set tick marks at
increments of 10
plt.yticks(range(0, 51, 5))
plt.grid(axis='y', alpha=.5)
```

In this code example, the hist() function is used to create the histogram, with several customization options specified. The edgecolor and linewidth arguments are used to add a black border and increase the width of the bars for improved visibility. The xticks and yticks functions are used to set tick marks at specified increments, improving the readability of the plot.

Box Plots:

```
import matplotlib.pyplot as plt
import pandas as pd
data = pd.read_csv('mydata.csv')
```



```
plt.boxplot([data['variable1'], data['variable2']],
vert=False, showfliers=False, widths=.5,
boxprops=dict(linewidth=2, color='green'),
whiskerprops=dict(linewidth=2),
medianprops=dict(linewidth=2, color='red'))
plt.xlabel('Value')
plt.ylabel('Value')
plt.ylabel('Variable')
plt.title('Box Plot of Variables')
plt.yticks([1, 2], ['Variable 1', 'Variable 2'])
plt.grid(axis='x', alpha=.5)
```

In this code example, the boxplot() function is used to create a box plot of two variables, with several customization options specified. The boxprops argument is used to change the color of the boxes to green, and the yticks function is used to set the labels for the y-axis to "Variable 1" and "Variable 2". This improves the readability of the plot by clearly labeling the variables.



Chapter 6: Machine Learning with Python



Introduction to machine learning

Machine learning is a subset of artificial intelligence that involves the use of algorithms and statistical models to enable computers to perform specific tasks without being explicitly programmed. In recent years, machine learning has become increasingly popular and is widely used in various industries, including healthcare, finance, and e-commerce, to name a few.

Python is a high-level programming language that is widely used for data analysis, scientific computing, and machine learning. It has become the go-to language for many data scientists and machine learning engineers due to its ease of use, readability, and large community support.

In this article, we will introduce you to the basics of machine learning in Python and how it compares to SAS, a popular statistical software used in data analysis and research.

Python for SAS Users

If you are a SAS user, you may find Python to be quite different from what you are used to. While SAS is a proprietary software, Python is an open-source language, meaning that its source code is freely available for modification and distribution.

One of the major differences between SAS and Python is their syntax. SAS uses a data step and a proc step to read and manipulate data, while Python uses libraries such as pandas, numpy, and scikit-learn to load and manipulate data.

Installing Python and Required Libraries

To get started with Python, you will need to install it on your computer. You can download the latest version of Python from the official website (https://www.python.org/downloads/).

After installing Python, you will need to install the required libraries for machine learning. Some of the popular libraries include:

NumPy: a library for working with arrays of data Pandas: a library for working with data frames Matplotlib: a library for creating visualizations Scikit-learn: a library for machine learning

To install these libraries, you can use pip, a package installer for Python. Open a terminal or command prompt and type the following command:

```
pip install numpy pandas matplotlib scikit-learn
This will install all the required libraries on your
computer.
```

Loading Data



In Python, you can load data from various sources, including CSV files, Excel files, and SQL databases. Here, we will show you how to load data from a CSV file using pandas.

```
import pandas as pd
# Load data from a CSV file
data = pd.read csv('data.csv')
```

In the above code, we first import the pandas library and then use the read_csv function to load data from a CSV file named data.csv. The data variable now contains the loaded data.

Data Preprocessing

Before you can apply machine learning algorithms to your data, you need to preprocess it to ensure that it is in a suitable format. Some of the common preprocessing steps include:

Handling missing values Scaling features Encoding categorical variables

In Python, you can use various libraries to perform these preprocessing steps. Here, we will show you how to handle missing values using the fillna function from pandas.

Replace missing values with the mean value
data.fillna(data.mean(), inplace=True)

In the above code, we use the fillna function to replace missing values with the mean value of the corresponding column. The inplace=True argument tells pandas to modify the data variable in place.

Splitting Data

Before you can train a machine learning model, you need to split your data into a training set and a test set. The training set is used to train the model, while the test set is used to evaluate its performance.

In Python, you can use the train_test_split function from scikit-learn to split your data into training and testing sets.

Here's a longer code example that demonstrates the entire machine learning process using Python:

```
import pandas as pd
from sklearn.model_selection import train_test_split
from sklearn.linear model import LinearRegression
```



```
from sklearn.metrics import mean squared error
# Load data from a CSV file
data = pd.read csv('data.csv')
# Preprocess the data
data.fillna(data.mean(), inplace=True)
# Split the data into training and testing sets
X = data.drop('target variable', axis=1)
y = data['target variable']
X_train, X_test, y_train, y_test = train_test_split(X,
y, test size=0.2, random state=42)
# Train a linear regression model
model = LinearRegression()
model.fit(X train, y train)
# Evaluate the model on the testing set
y pred = model.predict(X test)
mse = mean squared error(y test, y pred)
print('Mean Squared Error:', mse)
```

In the above code, we first load data from a CSV file and preprocess it by filling missing values with the mean value. We then split the data into training and testing sets using the train_test_split function from scikit-learn.

Next, we train a linear regression model on the training set using the LinearRegression class from scikit-learn. We then evaluate the model on the testing set by making predictions using the predict method and computing the mean squared error using the mean_squared_error function from scikit-learn.

This code demonstrates a simple machine learning pipeline that you can use as a template for your own machine learning projects. However, keep in mind that different machine learning problems may require different preprocessing steps and algorithms, so you will need to tailor this pipeline to your specific problem.

Supervised Learning

Supervised learning is a type of machine learning where you have labeled data, meaning that you have input features and corresponding output labels. The goal of supervised learning is to learn a mapping from input features to output labels. Linear Regression



Linear regression is a simple supervised learning algorithm that is used for predicting a continuous output variable based on one or more input variables. Here's an example of using linear regression in Python:

```
import pandas as pd
from sklearn.linear model import LinearRegression
from sklearn.metrics import mean squared error
# Load data from a CSV file
data = pd.read csv('data.csv')
# Preprocess the data
data.fillna(data.mean(), inplace=True)
# Split the data into training and testing sets
X = data.drop('target variable', axis=1)
y = data['target variable']
X train, X test, y_train, y_test = train_test_split(X,
y, test size=0.2, random state=42)
# Train a linear regression model
model = LinearRegression()
model.fit(X train, y train)
# Evaluate the model on the testing set
y pred = model.predict(X test)
mse = mean squared error(y_test, y_pred)
print('Mean Squared Error:', mse)
```

In the above code, we load data from a CSV file, preprocess it by filling missing values with the mean value, and split it into training and testing sets. We then train a linear regression model on the training set using the LinearRegression class from scikit-learn.

Finally, we evaluate the model on the testing set by making predictions using the predict method and computing the mean squared error using the mean_squared_error function from scikit-learn.

Decision Trees

Decision trees are another type of supervised learning algorithm that are used for both classification and regression tasks. A decision tree is a flowchart-like structure where each internal node represents a test on a feature, each branch represents the outcome of the test, and each leaf node represents a class label or a numerical value.

Here's an example of using decision trees in Python:



```
import pandas as pd
from sklearn.tree import DecisionTreeRegressor
from sklearn.metrics import mean squared error
# Load data from a CSV file
data = pd.read csv('data.csv')
# Preprocess the data
data.fillna(data.mean(), inplace=True)
# Split the data into training and testing sets
X = data.drop('target variable', axis=1)
y = data['target variable']
X_train, X_test, y_train, y_test = train_test_split(X,
y, test size=0.2, random state=42)
# Train a decision tree model
model = DecisionTreeRegressor()
model.fit(X train, y train)
# Evaluate the model on the testing set
y pred = model.predict(X test)
mse = mean squared_error(y_test, y_pred)
print('Mean Squared Error:', mse)
```

In the above code, we load data from a CSV file, preprocess it by filling missing values with the mean value, and split it into training and testing sets. We then train a decision tree model on the training set using the DecisionTreeRegressor class from scikit-learn.

Finally, we evaluate the model on the testing set by making predictions using the predict method and computing the mean squared error using the mean_squared_error function from scikit-learn.

Unsupervised Learning

Unsupervised learning is a type of machine learning where you have unlabeled data, meaning that you don't have any output labels.

Clustering

Clustering is a type of unsupervised learning where you group data points together based on their similarity. Here's an example of using K-means clustering in Python:

```
import pandas as pd
from sklearn.cluster import KMeans
import matplotlib.pyplot as plt
```



```
# Load data from a CSV file
data = pd.read_csv('data.csv')
# Preprocess the data
data.fillna(data.mean(), inplace=True)
# Cluster the data using K-means
X = data.drop('target_variable', axis=1)
model = KMeans(n_clusters=3)
model.fit(X)
# Visualize the results
plt.scatter(X.iloc[:, 0], X.iloc[:, 1],
c=model.labels_)
plt.show()
```

In the above code, we load data from a CSV file, preprocess it by filling missing values with the mean value, and cluster it using K-means clustering. We then visualize the results by plotting the data points and coloring them based on their assigned cluster.

Dimensionality Reduction

Dimensionality reduction is a type of unsupervised learning where you reduce the number of features in your data while preserving as much information as possible. Here's an example of using principal component analysis (PCA) for dimensionality reduction in Python:

```
import pandas as pd
from sklearn.decomposition import PCA
import matplotlib.pyplot as plt
# Load data from a CSV file
data = pd.read_csv('data.csv')
# Preprocess the data
data.fillna(data.mean(), inplace=True)
# Perform PCA for dimensionality reduction
X = data.drop('target_variable', axis=1)
pca = PCA(n_components=2)
X_pca = pca.fit_transform(X)
# Visualize the results
plt.scatter(X_pca[:, 0], X_pca[:, 1])
plt.show()
```



In the above code, we load data from a CSV file, preprocess it by filling missing values with the mean value, and perform PCA for dimensionality reduction. We then visualize the results by plotting the reduced data points.

Deep Learning

Deep learning is a type of machine learning that involves training neural networks with many layers. Here's an example of using a deep learning model for image classification in Python:

```
import tensorflow as tf
from tensorflow.keras.datasets import mnist
from tensorflow.keras.models import Sequential
from tensorflow.keras.layers import Dense, Flatten,
Conv2D, MaxPooling2D
# Load the MNIST dataset
(X train, y train), (X test, y test) =
mnist.load data()
# Preprocess the data
X train = X train.reshape(-1, 28, 28, 1) / 255.0
X \text{ test} = X \text{ test.reshape}(-1, 28, 28, 1) / 255.0
# Define the model architecture
model = Sequential([
    Conv2D(32, (3, 3), activation='relu',
input shape=(28, 28, 1)),
    MaxPooling2D((2, 2)),
    Flatten(),
    Dense(10, activation='softmax')
1)
# Compile the model
model.compile(optimizer='adam',
loss='sparse categorical crossentropy',
metrics=['accuracy'])
# Train the model
model.fit(X train, y train, epochs=5,
validation data=(X test, y test))
# Evaluate the model on the testing set
test loss, test acc = model.evaluate(X test, y test)
print('Test accuracy:', test acc)
```



In the above code, we load the MNIST dataset, preprocess it by scaling the pixel values between 0 and 1, define a convolutional neural network architecture using TensorFlow/Keras, compile the model with an optimizer and loss function, train the model for 5 epochs, and evaluate the model's performance on the testing set.

Natural Language Processing

Natural language processing (NLP) is a type of machine learning that involves working with human language data, such as text or speech. Here's an example of using a pre-trained NLP model for sentiment analysis in Python:

```
import tensorflow as tf
import tensorflow hub as hub
import pandas as pd
# Load the pre-trained NLP model
embed = hub.load("https://tfhub.dev/google/universal-
sentence-encoder/4")
# Load the text data
data = pd.read csv('text data.csv')
# Preprocess the data
X = data['text']
y = data['target variable']
# Encode the text data using the pre-trained model
X = ncoded = embed(X)
# Train a machine learning model on the encoded data
model = tf.keras.Sequential([
    tf.keras.layers.Dense(64, activation='relu'),
    tf.keras.layers.Dense(1, activation='sigmoid')
])
model.compile(optimizer='adam',
loss='binary crossentropy', metrics=['accuracy'])
model.fit(X encoded, y, epochs=10)
# Test the model on new text data
new text = ["This product is amazing!", "I wouldn't
recommend this product to anyone."]
new text encoded = embed(new text)
predictions = model.predict(new text encoded)
print(predictions)
```



In the above code, we load a pre-trained NLP model from TensorFlow Hub, load text data from a CSV file, preprocess the data by separating the text and target variable, encode the text using the pre-trained model, train a machine learning model on the encoded data, and test the model on new text data.

Recommender Systems

Recommender systems are a type of machine learning that suggest items to users based on their preferences and behavior. Here's an example of building a simple recommender system in Python:

```
import pandas as pd
from sklearn.metrics.pairwise import cosine similarity
from sklearn.feature extraction.text import
CountVectorizer
# Load the data
data = pd.read csv('product data.csv')
# Create a CountVectorizer object
vectorizer = CountVectorizer()
# Convert the product names into a matrix of word
counts
product matrix =
vectorizer.fit transform(data['product name'])
# Calculate the cosine similarity between each pair of
products
similarity matrix = cosine similarity(product matrix)
# Get the indices of the top 5 most similar products to
each product
top 5 similar products = []
for i in range(len(data)):
    similarities =
list(enumerate(similarity matrix[i]))
    similarities sorted = sorted(similarities,
key=lambda x: x[1], reverse=True)
    top 5 similar = [x[0] for x in
similarities sorted[1:6]]
    top 5 similar products.append(top 5 similar)
```



```
# Print the top 5 most similar products for each
product
for i in range(len(data)):
    print(f"Top 5 similar products for {data.loc[i,
    'product_name']}:")
    for j in top_5_similar_products[i]:
        print(f" {data.loc[j, 'product_name']}")
```

In the above code, we load a dataset of product names, create a matrix of word counts using the CountVectorizer object, calculate the cosine similarity between each pair of products, and retrieve the indices of the top 5 most similar products to each product. Finally, we print out the top 5 most similar products for each product in the dataset.

Image Classification

Image classification is a type of machine learning that involves categorizing images into different classes. Here's an example of using a pre-trained image classification model in Python:

```
import tensorflow as tf
import numpy as np
import urllib.request
# Load the pre-trained image classification model
model =
tf.keras.applications.ResNet50(include top=True,
weights='imagenet')
# Load an image from a URL
url = 'https://www.example.com/image.jpg'
with urllib.request.urlopen(url) as url:
    img = np.asarray(bytearray(url.read()),
dtype=np.uint8)
    img = tf.image.decode jpeg(img, channels=3)
    img = tf.image.resize(img, [224, 224])
# Preprocess the image data
imq =
tf.keras.applications.resnet50.preprocess input(img)
# Make a prediction on the image data
preds = model.predict(np.array([img]))
predicted class =
tf.keras.applications.resnet50.decode predictions(preds
, top=1)[0][0][1]
```



print(f"The predicted class of the image is: {predicted_class}")

In the above code, we load a pre-trained image classification model from Keras, load an image from a URL, preprocess the image data, make a prediction on the image data, and print out the predicted class of the image.

Time Series Analysis

Time series analysis is a type of machine learning that involves working with data that changes over time. Here's an example of using ARIMA (autoregressive integrated moving average) to forecast future values of a time series in Python:

Common Machine Learning Libraries in Python Python has a wide range of powerful machine learning libraries, some of which are:

Scikit-learn: A popular machine learning library for Python that includes a variety of algorithms for classification, regression, clustering, and dimensionality reduction.

TensorFlow: A library for building and training neural networks, used for deep learning applications.

Keras: A high-level neural networks API, written in Python and capable of running on top of TensorFlow.

PyTorch: A machine learning library that emphasizes ease of use and flexibility, often used for research and prototyping.

Theano: A Python library for fast numerical computation, often used for building and training neural networks.

NLTK: A leading platform for building Python programs to work with human language data, including sentiment analysis, topic modeling, and named entity recognition.

These libraries provide a wide range of tools and algorithms for building machine learning models in Python.

Types of Machine Learning in Python

There are three main types of machine learning:

Supervised Learning

In supervised learning, the model is trained on a labeled dataset, where each data point has an associated label or target value. The goal is to learn a mapping between the input data and the target labels, so that the model can accurately predict the target value for new, unseen data.

Some common algorithms for supervised learning include linear regression, logistic regression, decision trees, random forests, and neural networks.



Unsupervised Learning

In unsupervised learning, the model is trained on an unlabeled dataset, where there are no target values or labels. The goal is to find patterns or structure in the data, such as clusters or groups of similar data points.

Some common algorithms for unsupervised learning include k-means clustering, principal component analysis (PCA), and autoencoders.

Reinforcement Learning

In reinforcement learning, the model learns by interacting with an environment and receiving rewards or punishments for its actions. The goal is to learn a policy or set of actions that maximizes the reward over time.

Some common applications of reinforcement learning include game playing, robotics, and optimization problems.

Example: Building a Simple Machine Learning Model in Python Here's an example of building a simple machine learning model in Python using Scikit-learn. In this example, we'll use the famous iris dataset to build a classifier that predicts the species of iris flowers based on their petal and sepal measurements:

```
from sklearn.datasets import load iris
from sklearn.model selection import train_test_split
from sklearn.neighbors import KNeighborsClassifier
# Load the iris dataset
iris = load iris()
# Split the data into training and testing sets
X_train, X_test, y_train, y_test =
train test split(iris.data, iris.target, test size=0.2,
random state=42)
# Fit a KNN classifier to the training data
knn = KNeighborsClassifier(n neighbors=3)
knn.fit(X train, y train)
# Evaluate the classifier on the testing data
accuracy = knn.score(X test, y test)
# Print out the accuracy score
print(f"The accuracy of the KNN classifier is:
{accuracy}")
```



In the above code, we load the iris dataset, split the data into training and testing sets, fit a Knearest neighbors (KNN) classifier to the training data, evaluate the classifier on the testing data, and print out the accuracy score.

This is just a simple example, but it illustrates how easy it is to build and evaluate machine learning models in Python using Scikit-learn.

Neural Networks

Neural networks are a popular type of machine learning model, often used for tasks such as image recognition, natural language processing, and speech recognition. Python has several libraries for building and training neural networks, including TensorFlow, Keras, and PyTorch. Neural networks are composed of layers of interconnected nodes, called neurons, that process and transform input data. The output of one layer becomes the input of the next layer, and so on, until the final output layer produces the predicted output.

Convolutional Neural Networks

Convolutional neural networks (CNNs) are a specific type of neural network commonly used for image classification tasks. CNNs consist of multiple layers, including convolutional layers, pooling layers, and fully connected layers. The convolutional layers apply a set of learnable filters to the input image to detect specific features, such as edges or corners. The pooling layers downsample the output of the convolutional layers, reducing the spatial dimensionality of the feature maps. The fully connected layers then perform classification based on the learned features.

Python has several libraries for building and training CNNs, including TensorFlow, Keras, and PyTorch.

Natural Language Processing

Natural language processing (NLP) is a branch of machine learning concerned with analyzing and processing human language data, such as text and speech. Python has several libraries for NLP, including NLTK, spaCy, and Gensim.

Some common NLP tasks include:

- **Tokenization**: Breaking text into individual words or tokens.
- **Part-of-speech tagging**: Labeling each word in a sentence with its part of speech, such as noun, verb, or adjective.
- **Named entity recognition**: Identifying and categorizing named entities, such as people, places, and organizations, in a text.
- Sentiment analysis: Determining the sentiment or opinion expressed in a text, such as positive or negative.
- **Topic modeling**: Identifying the underlying topics or themes in a collection of texts.



Hyperparameter Tuning

Hyperparameter tuning is the process of finding the best set of hyperparameters for a machine learning model. Hyperparameters are settings that control the behavior and performance of a model, such as the learning rate or regularization strength.

Python has several libraries for hyperparameter tuning, including GridSearchCV and RandomizedSearchCV in Scikit-learn, and Hyperopt.

Preparing data for machine learning

Preparing data for machine learning is a crucial step in the data science pipeline. In this article, we will discuss the basics of preparing data for machine learning in Python for SAS users. We assume that you have a basic understanding of SAS programming and data manipulation concepts.

Importing data in Python

The first step in preparing data for machine learning in Python is to import the data into Python. There are several ways to import data into Python, including using libraries like pandas, numpy, and csv. Here is an example of importing data using pandas library:

```
import pandas as pd
data = pd.read_csv('data.csv')
```

This code imports the data from a CSV file called 'data.csv' and stores it in a pandas data frame called 'data'. You can also import data from other file formats such as Excel, SQL, or JSON. Exploring data

The next step is to explore the data and gain insights about it. This includes checking for missing values, outliers, and correlations between variables. You can use pandas and other libraries such as matplotlib and seaborn to visualize the data and gain insights. Here is an example code to check for missing values in a data frame:

```
print(data.isnull().sum())
```

This code prints the number of missing values in each column of the data frame.

Data cleaning

Once you have explored the data, the next step is to clean the data. This includes handling missing values, removing outliers, and transforming variables. Here is an example of handling missing values using pandas library:



data = data.dropna()

This code drops all the rows that contain missing values from the data frame.

Feature engineering

Feature engineering is the process of creating new features from existing features or transforming existing features to improve the performance of the machine learning model. This includes encoding categorical variables, scaling numeric variables, and creating new variables from existing variables. Here is an example of encoding categorical variables using pandas library:

```
data = pd.get dummies(data, columns=['gender'])
```

This code creates dummy variables for the 'gender' column, which is a categorical variable.

Splitting data

The final step in preparing data for machine learning is to split the data into training and testing sets. This is done to evaluate the performance of the machine learning model on new, unseen data. You can use the train_test_split function from the sklearn library to split the data. Here is an example of splitting data using sklearn library:

```
from sklearn.model_selection import train_test_split
X_train, X_test, y_train, y_test = train_test_split(X,
y, test_size=0.2, random_state=42)
This code splits the data into training and testing
sets with a test size of 20% and a random state of 42.
```

preparing data for machine learning involves several steps, including importing data, exploring data, cleaning data, feature engineering, and splitting data. By following these steps, you can prepare your data for machine learning in Python and build accurate machine learning models.

Here is an example of code that combines all the above steps:

```
import pandas as pd
from sklearn.model_selection import train_test_split
from sklearn.preprocessing import StandardScaler
# Import data
data = pd.read_csv('data.csv')
# Explore data
print(data.isnull().sum())
# Handle missing values
```



```
data = data.dropna()
# Feature engineering
data = pd.get_dummies(data, columns=['gender'])
# Split data
X = data.drop('target', axis=1)
y = data['target']
scaler = StandardScaler()
X_scaled = scaler.fit_transform(X)
X_train, X_test, y_train, y_test =
```

Here is a longer code example that includes additional steps for data preprocessing and feature engineering:

```
import pandas as pd
import numpy as np
from sklearn.preprocessing import LabelEncoder,
StandardScaler
from sklearn.model selection import train test split
# Import data
data = pd.read csv('data.csv')
# Explore data
print(data.head())
print(data.describe())
# Handle missing values
data = data.dropna()
# Feature engineering
# Encode categorical variables
le = LabelEncoder()
data['gender'] = le.fit transform(data['gender'])
# Scale numeric variables
scaler = StandardScaler()
numeric cols = ['age', 'income']
data[numeric cols] =
scaler.fit transform(data[numeric cols])
# Create new variables
data['age squared'] = np.square(data['age'])
```



```
data['age_income_interaction'] = data['age'] *
data['income']

# Split data
X = data.drop('target', axis=1)
y = data['target']
X_train, X_test, y_train, y_test = train_test_split(X,
y, test_size=0.2, random_state=42)

# Save preprocessed data
X_train.to_csv('X_train.csv', index=False)
X_test.to_csv('X_test.csv', index=False)
y_train.to_csv('y_train.csv', index=False)
y_test.to_csv('y_test.csv', index=False)
```

In this code example, we first import the data from a CSV file and explore the data using the head() and describe() methods. We then drop any rows with missing values and perform feature engineering.

To encode the categorical variable 'gender', we use the LabelEncoder class from the sklearn library. We then scale the numeric variables 'age' and 'income' using the StandardScaler class.

we create two new variables by squaring the 'age' variable and creating an interaction term between 'age' and 'income'. We then split the data into training and testing sets using the train_test_split function from the sklearn library. we save the preprocessed data as CSV files using the to_csv method of the pandas library. These files can then be used as input to a machine learning model.

examples of common data preprocessing and feature engineering techniques used in machine learning.

Handling Missing Values

Missing values are a common issue in real-world datasets. One common approach to handling missing values is to simply drop any rows or columns that contain missing values using the dropna() method in pandas. For example:

```
# Drop any rows with missing values
data = data.dropna()
# Drop any columns with missing values
data = data.dropna(axis=1)
```

Another approach is to impute the missing values with a reasonable estimate. There are several methods for imputing missing values, including mean imputation, median imputation, and K-nearest neighbor imputation. For example:



```
# Mean imputation
data.fillna(data.mean(), inplace=True)
# Median imputation
data.fillna(data.median(), inplace=True)
# K-nearest neighbor imputation
from sklearn.impute import KNNImputer
imputer = KNNImputer(n_neighbors=5)
data_imputed = imputer.fit_transform(data)
```

Encoding Categorical Variables

Machine learning algorithms generally require all input variables to be numeric. Therefore, categorical variables (variables with non-numeric values) need to be encoded as numeric values. There are several encoding techniques, including one-hot encoding, label encoding, and target encoding. For example:

```
# One-hot encoding
data_one_hot = pd.get_dummies(data, columns=['gender'])
# Label encoding
le = LabelEncoder()
data['gender'] = le.fit_transform(data['gender'])
# Target encoding
import category_encoders as ce
encoder = ce.TargetEncoder(cols=['gender'])
data_target_encoded = encoder.fit_transform(data, y)
```

Scaling Numeric Variables

Some machine learning algorithms are sensitive to the scale of numeric variables. Therefore, it is often important to scale numeric variables to have a similar range of values. Common scaling techniques include standardization (subtracting the mean and dividing by the standard deviation) and normalization (scaling to a range of 0 to 1). For example:

```
# Standardization
scaler = StandardScaler()
data[numeric_cols] =
scaler.fit_transform(data[numeric_cols])
# Normalization
from sklearn.preprocessing import MinMaxScaler
scaler = MinMaxScaler()
```



```
data[numeric cols] =
scaler.fit transform(data[numeric cols])
Feature Engineering
Feature engineering involves creating new input
variables (features) from the existing variables in the
dataset. The goal is to create features that capture
important patterns or relationships in the data that
are relevant to the prediction task. Some common
feature engineering techniques include creating
interaction terms (multiplying two or more variables
together), transforming variables (e.g., taking the
square root of a variable), and creating new variables
based on domain knowledge (e.g., creating a variable
that indicates the season of the year). For example:
python
Copy code
# Creating interaction terms
data['age income interaction'] = data['age'] *
data['income']
# Transforming variables
data['sqrt age'] = np.sqrt(data['age'])
# Creating new variables based on domain knowledge
data['season'] = data['month'].apply(lambda x: 'Winter'
if x in [12, 1, 2] else 'Spring' if x in [3, 4, 5] else
'Summer' if x in [6, 7, 8] else 'Fall')
```

These are just a few examples of the many data preprocessing and feature engineering techniques that can be used in machine learning. The choice of techniques will depend on the specific dataset and prediction task.

Handling Outliers

Outliers are data points that are significantly different from other data points in the dataset. Outliers can be caused by measurement errors, data entry errors, or genuine extreme values. Outliers can have a large impact on some machine learning algorithms, so it is often important to detect and handle them. Common approaches to handling outliers include removing them from the dataset, transforming the variable using a log or power transformation, or capping the variable at a reasonable value. For example:

```
# Removing outliers
data = data[(data['income'] > 10000) & (data['income']
< 1000000)]</pre>
```



```
# Log transformation
data['log_income'] = np.log(data['income'])
# Capping variable at 99th percentile
income_cap = np.percentile(data['income'], 99)
data['capped_income'] = np.where(data['income'] >
income_cap, income_cap, data['income'])
```

Handling Skewed Variables

Skewed variables are variables that are not normally distributed. Skewed variables can cause some machine learning algorithms to perform poorly, so it is often important to transform them to have a more normal distribution. Common transformations include log transformation, square root transformation, and Box-Cox transformation. For example:

```
# Log transformation
data['log_income'] = np.log(data['income'])
# Square root transformation
data['sqrt_income'] = np.sqrt(data['income'])
# Box-Cox transformation
from scipy.stats import boxcox
data['boxcox_income'], _ = boxcox(data['income'])
```

Dimensionality Reduction

High-dimensional datasets (datasets with many input variables) can be difficult to work with and can cause some machine learning algorithms to perform poorly. Dimensionality reduction techniques can be used to reduce the number of input variables while preserving as much information as possible. Common dimensionality reduction techniques include principal component analysis (PCA) and t-distributed stochastic neighbor embedding (t-SNE). For example:

```
# PCA
from sklearn.decomposition import PCA
pca = PCA(n_components=2)
data_pca = pca.fit_transform(data)
# t-SNE
from sklearn.manifold import TSNE
tsne = TSNE(n_components=2)
data_tsne = tsne.fit_transform(data)
```



These are just a few more examples of the many data preprocessing and feature engineering techniques that can be used in machine learning. It is important to note that not all techniques will be relevant or appropriate for every dataset and prediction task, so it is important to carefully consider which techniques to use based on the characteristics of the data and the requirements of the prediction task.

Encoding Categorical Variables

Many machine learning algorithms require numerical input variables, so categorical variables (variables with discrete categories, such as color or type) must be encoded as numbers. There are several ways to encode categorical variables, including one-hot encoding, ordinal encoding, and target encoding. One-hot encoding creates a new binary column for each category, while ordinal encoding assigns a unique number to each category. Target encoding uses the target variable to encode each category based on the average value of the target variable for that category. For example:

```
# One-hot encoding
data_one_hot = pd.get_dummies(data, columns=['color'])
# Ordinal encoding
from sklearn.preprocessing import OrdinalEncoder
enc = OrdinalEncoder()
data['color_encoded'] =
enc.fit_transform(data[['color']])
# Target encoding
from category_encoders.target_encoder import
TargetEncoder
enc = TargetEncoder()
data['color_encoded'] =
enc.fit_transform(data['color'], data['target'])
```

Feature Scaling

Many machine learning algorithms are sensitive to the scale of input variables, so it is often important to scale or normalize the input variables. Common scaling techniques include minmax scaling (scaling the variable to a specified range, such as 0 to 1) and standardization (scaling the variable to have a mean of 0 and a standard deviation of 1). For example:

```
# Min-max scaling
from sklearn.preprocessing import MinMaxScaler
scaler = MinMaxScaler(feature_range=(0, 1))
data_scaled = scaler.fit_transform(data)
# Standardization
```



```
from sklearn.preprocessing import StandardScaler
scaler = StandardScaler()
data_scaled = scaler.fit_transform(data)
```

Feature Selection

In some cases, not all input variables are relevant or useful for making predictions, and using all input variables can lead to overfitting or decreased performance. Feature selection techniques can be used to select the most important or relevant input variables. Common feature selection techniques include correlation analysis, feature importance analysis using machine learning models, and stepwise regression. For example:

```
# Correlation analysis
corr matrix = data.corr()
corr features =
corr matrix.index[abs(corr matrix['target']) > 0.5]
data corr = data[corr features]
# Feature importance analysis
from sklearn.ensemble import RandomForestRegressor
model = RandomForestRegressor()
model.fit(data.drop('target', axis=1), data['target'])
feature importances =
pd.DataFrame(model.feature importances ,
index=data.columns[:-1], columns=['importance'])
feature importances =
feature importances.sort values('importance',
ascending=False)
important features = feature importances.index[:10]
data important = data[important features]
# Stepwise regression
from sklearn.feature selection import
SequentialFeatureSelector
from sklearn.linear model import LinearRegression
selector =
SequentialFeatureSelector(LinearRegression(),
direction='backward')
selector.fit(data.drop('target', axis=1),
data['target'])
selected features = data.columns[:-
1][selector.support ]
data selected = data[selected features]
```

in stal

These are just a few more examples of the many data preprocessing and feature engineering techniques that can be used in machine learning. By understanding and applying these techniques, SAS users can expand their skillset and apply their data analysis expertise to a wider range of problems and applications.

Handling Missing Data

Missing data is a common issue in many datasets and can negatively impact the performance of machine learning models. There are several techniques for handling missing data, including imputation (filling in missing values with estimated values), deletion (removing rows or columns with missing values), and using algorithms that can handle missing data directly (such as decision trees or random forests). For example:

```
# Imputation using mean value
from sklearn.impute import SimpleImputer
imputer = SimpleImputer(strategy='mean')
data_imputed = imputer.fit_transform(data)
# Deletion of rows with missing values
data_deleted_rows = data.dropna(axis=0)
# Deletion of columns with missing values
data_deleted_cols = data.dropna(axis=1)
# Decision tree with missing data
from sklearn.tree import DecisionTreeRegressor
model = DecisionTreeRegressor()
model.fit(data.drop('target', axis=1), data['target'])
```

Handling Imbalanced Data

In some cases, the target variable in a dataset may be imbalanced, meaning that one class has significantly fewer examples than the other. This can lead to biased models that perform poorly on the minority class. There are several techniques for handling imbalanced data, including oversampling (generating synthetic examples of the minority class), undersampling (removing examples of the majority class), and using algorithms that are designed to handle imbalanced data directly (such as weighted random forests). For example:

```
# Oversampling using SMOTE
from imblearn.over_sampling import SMOTE
smote = SMOTE()
data_oversampled, target_oversampled =
smote.fit_resample(data.drop('target', axis=1),
data['target'])
```



```
# Undersampling using random sampling
from imblearn.under_sampling import RandomUnderSampler
rus = RandomUnderSampler()
data_undersampled, target_undersampled =
rus.fit_resample(data.drop('target', axis=1),
data['target'])
# Weighted random forest
from sklearn.ensemble import RandomForestClassifier
class_weights = {0: 1, 1: 10} # Higher weight for
minority class
model =
RandomForestClassifier(class_weight=class_weights)
model.fit(data.drop('target', axis=1), data['target'])
```

Data Transformation

In some cases, transforming the input or target variables can improve the performance of machine learning models. Common transformations include log transforms, square roots, and box-cox transforms. For example:

```
# Log transform
import numpy as np
data['target_log'] = np.log(data['target'])
# Square root transform
data['target_sqrt'] = np.sqrt(data['target'])
# Box-cox transform
from scipy.stats import boxcox
data['target boxcox'], lam = boxcox(data['target'])
```

These are just a few more examples of the many techniques that can be used for data preprocessing and feature engineering in machine learning. By mastering these techniques, SAS users can become proficient in Python and expand their data analysis skills to a wider range of applications.

Feature Scaling

In many cases, the features in a dataset may have different scales or units of measurement. This can make it difficult for machine learning algorithms to effectively learn from the data. Feature scaling can help by transforming the features to have similar scales. Common scaling techniques include standardization (scaling features to have zero mean and unit variance) and normalization (scaling features to have a range between 0 and 1). For example:

in stal

```
# Standardization
from sklearn.preprocessing import StandardScaler
scaler = StandardScaler()
data_standardized =
scaler.fit_transform(data.drop('target', axis=1))
# Normalization
from sklearn.preprocessing import MinMaxScaler
scaler = MinMaxScaler()
data_normalized =
scaler.fit_transform(data.drop('target', axis=1))
```

Feature Selection

In some cases, not all of the features in a dataset may be relevant to the target variable. Feature selection can help by identifying the most important features and removing the irrelevant ones. Common feature selection techniques include correlation analysis (identifying features that are highly correlated with the target variable), mutual information (identifying features that provide the most information about the target variable), and Lasso regression (identifying features with the highest coefficients in a linear model). For example:

```
# Correlation analysis
import seaborn as sns
corr matrix = data.corr()
sns.heatmap(corr matrix)
# Mutual information
from sklearn.feature selection import SelectKBest,
mutual info classif
selector = SelectKBest(mutual info classif, k=5)
data selected =
selector.fit transform(data.drop('target', axis=1),
data['target'])
# Lasso regression
from sklearn.linear model import Lasso
model = Lasso(alpha=0.1)
model.fit(data.drop('target', axis=1), data['target'])
important features = data.columns[np.abs(model.coef ) >
0.11
```



Dimensionality Reduction

In some cases, the number of features in a dataset may be very high, making it difficult to analyze and visualize the data. Dimensionality reduction can help by reducing the number of features while still retaining as much information as possible. Common dimensionality reduction techniques include principal component analysis (PCA), linear discriminant analysis (LDA), and t-SNE. For example:

```
# Principal component analysis
from sklearn.decomposition import PCA
pca = PCA(n components=2)
data reduced = pca.fit transform(data.drop('target',
axis=1))
# Linear discriminant analysis
from sklearn.discriminant analysis import
LinearDiscriminantAnalysis
lda = LinearDiscriminantAnalysis(n components=2)
data_reduced = lda.fit_transform(data.drop('target',
axis=1), data['target'])
# t-SNE
from sklearn.manifold import TSNE
tsne = TSNE(n components=2)
data reduced = tsne.fit transform(data.drop('target',
axis=1))
```

data preprocessing and feature engineering are critical steps in the machine learning pipeline. By mastering these techniques, SAS users can leverage Python's powerful machine learning libraries to solve a wide range of data analysis problems.

Handling Missing Data

In many datasets, some values may be missing, either due to measurement error or other factors. It is important to handle missing data appropriately, as most machine learning algorithms cannot handle missing values. Common techniques for handling missing data include imputation (filling in missing values with estimated values) and deletion (removing rows or columns with missing values). For example:

```
# Imputation
from sklearn.impute import SimpleImputer
imputer = SimpleImputer(strategy='mean')
data_imputed =
imputer.fit_transform(data.drop('target', axis=1))
```



```
# Deletion
data_dropped = data.dropna()
```

Encoding Categorical Variables

In many datasets, some features may be categorical variables, such as colors, countries, or product types. Machine learning algorithms generally require numerical inputs, so categorical variables must be encoded as numerical values. Common encoding techniques include one-hot encoding (creating binary columns for each category), label encoding (assigning a numerical value to each category), and target encoding (replacing each category with the mean target value for that category). For example:

```
# One-hot encoding
from sklearn.preprocessing import OneHotEncoder
encoder = OneHotEncoder()
data_encoded =
encoder.fit_transform(data.drop('target', axis=1))
# Label encoding
from sklearn.preprocessing import LabelEncoder
encoder = LabelEncoder()
data['color_encoded'] =
encoder.fit_transform(data['color'])
# Target encoding
import category_encoders as ce
encoder = ce.TargetEncoder()
data['country_encoded'] =
encoder.fit_transform(data['country'], data['target'])
```

Handling Imbalanced Data

In some datasets, the target variable may be imbalanced, meaning that one class is much more common than the others. This can lead to biased machine learning models that perform poorly on the minority class. Common techniques for handling imbalanced data include resampling (creating new examples of the minority class) and adjusting class weights (giving more weight to the minority class). For example:

```
# Resampling
from imblearn.over_sampling import SMOTE
smote = SMOTE()
data_resampled, target_resampled =
smote.fit_resample(data.drop('target', axis=1),
data['target'])
```



```
# Adjusting class weights
from sklearn.svm import SVC
model = SVC(class_weight='balanced')
model.fit(data.drop('target', axis=1), data['target'])
```

Cross-Validation

When evaluating machine learning models, it is important to use a robust evaluation technique that avoids overfitting to the training data. Cross-validation can help by splitting the data into multiple training and validation sets and averaging the results. Common cross-validation techniques include k-fold cross-validation (splitting the data into k folds and using each fold as a validation set) and leave-one-out cross-validation (using each example as a validation set). For example:

```
# K-fold cross-validation
from sklearn.model_selection import cross_val_score
from sklearn.ensemble import RandomForestClassifier
model = RandomForestClassifier()
scores = cross_val_score(model, data.drop('target',
axis=1), data['target'], cv=5)
# Leave-one-out cross-validation
from sklearn.model_selection import LeaveOneOut
loo = LeaveOneOut()
for train_index, test_index in loo.split(data):
    model.fit(data.drop('target',
axis=1).iloc[train_index],
data['target'].iloc[train_index])
    prediction = model.predict(data.drop('target',
axis=1).iloc[test_index])
```

These are just a few examples of the many data preprocessing and feature engineering techniques that SAS users can leverage in Python.

Dimensionality Reduction

In some datasets, there may be many features, making it difficult to build accurate machine learning models. Dimensionality reduction techniques can help by reducing the number of features while preserving the most important information. Common dimensionality reduction techniques include principal component analysis (PCA), linear discriminant analysis (LDA), and t-distributed stochastic neighbor embedding (t-SNE). For example:

```
# PCA
from sklearn.decomposition import PCA
```



```
pca = PCA(n_components=2)
data_reduced = pca.fit_transform(data.drop('target',
axis=1))

# LDA
from sklearn.discriminant_analysis import
LinearDiscriminantAnalysis
lda = LinearDiscriminantAnalysis(n_components=2)
data_reduced = lda.fit_transform(data.drop('target',
axis=1), data['target'])

# t-SNE
from sklearn.manifold import TSNE
tsne = TSNE(n_components=2)
data_reduced = tsne.fit_transform(data.drop('target',
axis=1))
```

Feature Scaling

Machine learning algorithms often perform better when the features are on the same scale. Feature scaling techniques can help by scaling the features to a common range. Common feature scaling techniques include min-max scaling (scaling the features to a range between 0 and 1) and standardization (scaling the features to have zero mean and unit variance). For example:

```
# Min-max scaling
from sklearn.preprocessing import MinMaxScaler
scaler = MinMaxScaler()
data_scaled = scaler.fit_transform(data.drop('target',
axis=1))
# Standardization
from sklearn.preprocessing import StandardScaler
scaler = StandardScaler()
data_scaled = scaler.fit_transform(data.drop('target',
axis=1))
```

Feature Selection

In some datasets, there may be many features, but not all of them may be relevant for predicting the target variable. Feature selection techniques can help by selecting the most relevant features and discarding the rest. Common feature selection techniques include correlation analysis (removing features that are highly correlated with each other), feature importance (using models to determine the importance of each feature), and recursive feature elimination (removing features one by one until the performance of the model deteriorates). For example:



```
# Correlation analysis
corr matrix = data.drop('target', axis=1).corr()
corr features = set()
for i in range(len(corr matrix.columns)):
    for j in range(i):
        if abs(corr matrix.iloc[i, j]) > 0.8:
            corr features.add(corr matrix.columns[i])
data.drop(corr features, axis=1, inplace=True)
# Feature importance
from sklearn.ensemble import RandomForestClassifier
model = RandomForestClassifier()
model.fit(data.drop('target', axis=1), data['target'])
importances = model.feature importances
indices = np.argsort(importances)[::-1]
selected features = [data.drop('target',
axis=1).columns[i] for i in indices[:10]]
# Recursive feature elimination
from sklearn.feature selection import RFECV
model = RandomForestClassifier()
selector = RFECV(model, step=1, cv=5)
selector.fit(data.drop('target', axis=1),
data['target'])
selected features = data.drop('target',
axis=1).columns[selector.support ]
```

Model Evaluation

When building machine learning models, it is important to evaluate their performance on unseen data. Common evaluation metrics include accuracy, precision, recall, F1 score, and area under the receiver operating characteristic (ROC) curve. For example:

```
from sklearn.metrics import accuracy_score,
precision_score, recall_score, f1_score, roc_auc_score
model = RandomForestClassifier()
model.fit(data_train.drop
```

Cross-Validation

When evaluating machine learning models, it is important to avoid overfitting to the training data. Cross-validation techniques can help by estimating the performance of the model on unseen data. Common cross-validation techniques include k-fold cross-validation (splitting the data into



k subsets and training the model on k-1 subsets while testing on the remaining subset) and stratified k-fold cross-validation (similar to k-fold cross-validation, but ensures that each subset has a similar distribution of the target variable). For example:

Hyperparameter Tuning

When building machine learning models, it is important to choose the best hyperparameters for the model. Hyperparameter tuning techniques can help by searching the hyperparameter space and finding the best hyperparameters for the model. Common hyperparameter tuning techniques include grid search (searching the hyperparameter space exhaustively), randomized search (searching the hyperparameter space randomly), and Bayesian optimization (using a probabilistic model to guide the search). For example:

```
from sklearn.model_selection import GridSearchCV
from sklearn.ensemble import RandomForestClassifier
param_grid = {
    'n_estimators': [100, 500, 1000],
    'max_depth': [None, 10, 20],
    'min_samples_split': [2, 5, 10],
    'min_samples_leaf': [1, 2, 4]
}
model = RandomForestClassifier()
grid_search = GridSearchCV(model, param_grid, cv=5)
grid_search.fit(data.drop('target', axis=1),
data['target'])
best_model = grid_search.best_estimator_
```

Model Deployment

Once a machine learning model has been trained and evaluated, it can be deployed to make predictions on new data. Common deployment techniques include deploying the model as a web service, integrating the model into a larger software system, or deploying the model on a mobile device. For example:

import pickle



```
# train and evaluate model on data_train and data_test
model = RandomForestClassifier()
model.fit(data_train.drop('target', axis=1),
data_train['target'])
# save model to disk
with open('model.pkl', 'wb') as f:
    pickle.dump(model, f)
# load model from disk
with open('model.pkl', 'rb') as f:
    model = pickle.load(f)
# make predictions on new data
predictions = model.predict(new_data)
```

Linear regression

Linear regression is a statistical modeling technique used to establish the relationship between a dependent variable and one or more independent variables. In Python, there are various libraries available for performing linear regression, including NumPy, SciPy, and scikit-learn. In this tutorial, we will focus on using scikit-learn, a popular machine learning library in Python, to perform linear regression.

To get started, let's import the necessary libraries:

```
import numpy as np
import pandas as pd
from sklearn.linear_model import LinearRegression
from sklearn.metrics import mean_squared_error
```

We will also need a dataset to work with. For this tutorial, we will use the "Boston Housing" dataset, which is available in scikit-learn. This dataset contains information about housing prices in Boston and various factors that may influence those prices, such as crime rate, average number of rooms per dwelling, and distance to employment centers.

```
from sklearn.datasets import load_boston
boston = load boston()
```

Now that we have our dataset, let's create a pandas DataFrame to make it easier to work with.

```
df = pd.DataFrame(boston.data,
columns=boston.feature_names)
df['MEDV'] = boston.target
```



The DataFrame now contains all the information from the Boston Housing dataset, including the target variable, which is the median value of owner-occupied homes in \$1000s. Next, let's split the data into training and testing sets.

```
from sklearn.model_selection import train_test_split
X_train, X_test, y_train, y_test =
train_test_split(df[boston.feature_names], df['MEDV'],
test size=0.2, random state=0)
```

We can now create an instance of the LinearRegression class and fit it to our training data.

```
regressor = LinearRegression()
regressor.fit(X train, y train)
```

Now that the model has been trained, we can use it to make predictions on our test data.

y pred = regressor.predict(X test)

To evaluate the performance of our model, we can calculate the mean squared error (MSE) between the predicted values and the actual values.

```
mse = mean_squared_error(y_test, y_pred)
print(mse)
```

This will give us the mean squared error between the predicted and actual values.

Finally, we can visualize the results of our linear regression model by plotting the predicted values against the actual values.

```
import matplotlib.pyplot as plt
plt.scatter(y_test, y_pred)
plt.xlabel("Actual Prices: $Y_i$")
plt.ylabel("Predicted Prices: $\hat{Y}_i$")
plt.title("Actual vs. Predicted Prices: $Y_i$ vs.
$\hat{Y}_i$")
plt.show()
```

This will produce a scatter plot with the actual values on the x-axis and the predicted values on the y-axis.

Here is a longer code example that goes into more detail on how to perform linear regression using scikit-learn in Python:

Import libraries



```
import numpy as np
import pandas as pd
import matplotlib.pyplot as plt
from sklearn.datasets import load boston
from sklearn.model selection import train test split
from sklearn.linear model import LinearRegression
from sklearn.metrics import mean squared error
# Load dataset
boston = load boston()
# Convert to pandas DataFrame
df = pd.DataFrame(boston.data,
columns=boston.feature names)
df['MEDV'] = boston.target
# Split data into training and testing sets
X train, X test, y train, y test =
train test split(df[boston.feature names], df['MEDV'],
test size=0.2, random state=0)
# Initialize linear regression model and fit to
training data
regressor = LinearRegression()
regressor.fit(X train, y train)
# Make predictions on test data
y pred = regressor.predict(X test)
# Calculate mean squared error between predicted and
actual values
mse = mean squared error(y test, y pred)
print("Mean squared error: ", mse)
# Plot actual vs. predicted values
plt.scatter(y test, y_pred)
plt.xlabel("Actual Prices: $Y i$")
plt.ylabel("Predicted Prices: $\hat{Y} i$")
plt.title("Actual vs. Predicted Prices: $Y i$ vs.
$\hat{Y} i$")
plt.show()
```



This code first loads the Boston Housing dataset using scikit-learn's load_boston function. It then converts the dataset into a pandas DataFrame for easier manipulation.

Next, the code splits the data into training and testing sets using scikit-learn's train_test_split function. It then initializes a linear regression model using scikit-learn's LinearRegression class and fits it to the training data using the fit method.

The code then makes predictions on the test data using the predict method of the linear regression model. It calculates the mean squared error between the predicted and actual values using scikit-learn's mean_squared_error function.

The code creates a scatter plot of the actual vs. predicted values using matplotlib's scatter function. This allows us to visually compare the predicted values to the actual values. this code demonstrates how to perform linear regression in Python using scikit-learn and visualize the results using matplotlib.

Linear regression is a statistical modeling technique that is used to model the relationship between a dependent variable and one or more independent variables. The goal of linear regression is to find the best-fit line that can predict the value of the dependent variable based on the values of the independent variables.

In scikit-learn, linear regression can be performed using the LinearRegression class. Here's an example of how to use LinearRegression to perform linear regression on the Boston Housing dataset:

```
# Load dataset
from sklearn.datasets import load boston
boston = load boston()
# Create pandas DataFrame
import pandas as pd
df = pd.DataFrame (boston.data,
columns=boston.feature names)
df['MEDV'] = boston.target
# Split data into training and testing sets
from sklearn.model selection import train test split
X train, X test, y train, y test =
train test split(df[boston.feature names], df['MEDV'],
test size=0.2, random state=0)
# Import LinearRegression class
from sklearn.linear model import LinearRegression
# Create instance of LinearRegression class
```



```
regressor = LinearRegression()
# Fit the model to the training data
regressor.fit(X_train, y_train)
# Make predictions on the testing data
y_pred = regressor.predict(X_test)
# Evaluate the model using mean squared error
from sklearn.metrics import mean_squared_error
mse = mean_squared_error(y_test, y_pred)
print("Mean squared error: ", mse)
```

In this example, we first load the Boston Housing dataset and create a pandas DataFrame. We then split the data into training and testing sets using train_test_split. Next, we import the LinearRegression class from scikit-learn and create an instance of the class.

We then fit the model to the training data using the fit method of the LinearRegression object. After fitting the model, we use the predict method to make predictions on the testing data. Finally, we evaluate the performance of the model using mean squared error.

Scikit-learn also provides other regression algorithms that can be used for linear regression, such as Ridge Regression and Lasso Regression. These algorithms can be useful when dealing with datasets that have a large number of features or when dealing with overfitting.

here are a few more examples and tips for performing linear regression in Python using scikitlearn.

Handling Categorical Variables

In many real-world datasets, some of the features may be categorical rather than numeric. To handle categorical variables in scikit-learn, you can use one-hot encoding or dummy variable encoding. Here's an example:

```
# Load dataset
import pandas as pd
df = pd.read_csv('my_dataset.csv')
# Convert categorical variable to dummy variables
df = pd.get_dummies(df, columns=['Category'])
# Split data into training and testing sets
from sklearn.model selection import train test split
```



```
X_train, X_test, y_train, y_test =
train_test_split(df.drop('Target', axis=1),
df['Target'], test_size=0.2, random_state=0)
# Perform linear regression
from sklearn.linear_model import LinearRegression
regressor = LinearRegression()
regressor.fit(X_train, y_train)
y_pred = regressor.predict(X_test)
```

In this example, we load a dataset that contains a categorical variable called 'Category'. We use the get_dummies function from pandas to convert the categorical variable to dummy variables. We then split the data into training and testing sets and perform linear regression using the LinearRegression class from scikit-learn.

Regularization

Regularization is a technique used to prevent overfitting by adding a penalty term to the cost function. In scikit-learn, you can use Ridge Regression or Lasso Regression to perform regularization. Here's an example using Ridge Regression:

```
# Load dataset
import pandas as pd
df = pd.read_csv('my_dataset.csv')
# Split data into training and testing sets
from sklearn.model_selection import train_test_split
X_train, X_test, y_train, y_test =
train_test_split(df.drop('Target', axis=1),
df['Target'], test_size=0.2, random_state=0)
# Perform Ridge Regression
from sklearn.linear_model import Ridge
regressor = Ridge(alpha=0.5)
regressor.fit(X_train, y_train)
y pred = regressor.predict(X test)
```

In this example, we split the data into training and testing sets as before. We then use the Ridge class from scikit-learn to perform Ridge Regression with a regularization parameter of 0.5. This will add a penalty term to the cost function to prevent overfitting.

Polynomial Regression



In some cases, the relationship between the independent and dependent variables may not be linear. In these cases, you can use Polynomial Regression to model the relationship using a higher-order polynomial function. Here's an example:

```
# Load dataset
import pandas as pd
df = pd.read csv('my dataset.csv')
# Split data into training and testing sets
from sklearn.model selection import train test split
X train, X test, y train, y test =
train test split(df.drop('Target', axis=1),
df['Target'], test size=0.2, random state=0)
# Perform Polynomial Regression
from sklearn.preprocessing import PolynomialFeatures
from sklearn.linear model import LinearRegression
poly = PolynomialFeatures(degree=2)
X train poly = poly.fit transform(X train)
X test poly = poly.transform(X test)
regressor = LinearRegression()
regressor.fit(X train poly, y train)
y pred = regressor.predict(X test poly)
```

In this example, we split the data into training and testing sets as before. We then use the PolynomialFeatures class from scikit-learn to transform the independent variables into polynomial features of degree 2. We then use the LinearRegression class from scikit-learn to perform linear regression using the transformed features.

Scaling Features

When performing linear regression, it's often a good idea to scale the features so that they have a similar range. This can help prevent numerical instability and improve the performance of the model. Here's an example:

```
# Load dataset
import pandas as pd
df = pd.read_csv('my_dataset.csv')
# Split data into training and testing sets
from sklearn.model_selection import train_test_split
X_train, X_test, y_train, y_test =
train_test_split(df.drop('Target', axis=1),
df['Target'], test_size=0.2, random_state=0)
```



```
# Scale features
from sklearn.preprocessing import StandardScaler
scaler = StandardScaler()
X_train_scaled = scaler.fit_transform(X_train)
X_test_scaled = scaler.transform(X_test)
# Perform linear regression
from sklearn.linear_model import LinearRegression
regressor = LinearRegression()
regressor.fit(X_train_scaled, y_train)
y_pred = regressor.predict(X_test_scaled)
```

In this example, we split the data into training and testing sets as before. We then use the StandardScaler class from scikit-learn to scale the features to have zero mean and unit variance. We then perform linear regression using the scaled features.

Cross Validation

When performing linear regression, it's important to evaluate the performance of the model on a separate testing set to avoid overfitting. However, if the dataset is small, it may be difficult to obtain reliable estimates of the performance on the testing set. In these cases, you can use cross validation to obtain more reliable estimates of the performance. Here's an example:

```
# Load dataset
import pandas as pd
df = pd.read csv('my dataset.csv')
# Perform k-fold cross validation
from sklearn.model selection import KFold
from sklearn.linear model import LinearRegression
from sklearn.metrics import mean squared error
import numpy as np
kf = KFold(n splits=5)
mse list = []
for train index, test index in kf.split(df):
    X train, X test = df.iloc[train index][['Feature1',
'Feature2']], df.iloc[test index][['Feature1',
'Feature2']]
    y train, y test = df.iloc[train index]['Target'],
df.iloc[test index]['Target']
    regressor = LinearRegression()
    regressor.fit(X train, y train)
    y pred = regressor.predict(X test)
```



```
mse = mean_squared_error(y_test, y_pred)
mse_list.append(mse)
```

```
print("Mean squared error: ", np.mean(mse_list))
```

In this example, we use the KFold class from scikit-learn to perform 5-fold cross validation. We then loop over the 5 folds and for each fold, we split the data into training and testing sets, perform linear regression, and evaluate the performance using mean squared error. We then take the mean of the mean squared errors from each fold to obtain an estimate of the performance of the model.

Regularization

In some cases, linear regression models can suffer from overfitting, where the model fits the training data too closely and fails to generalize well to new data. One way to address this issue is to use regularization, which adds a penalty term to the loss function to discourage large weights. Two common types of regularization are L1 regularization, which adds a penalty proportional to the absolute value of the weights, and L2 regularization, which adds a penalty proportional to the square of the weights. Here's an example using L2 regularization:

```
# Load dataset
import pandas as pd
df = pd.read csv('my dataset.csv')
# Split data into training and testing sets
from sklearn.model selection import train test split
X train, X test, y train, y test =
train test split(df.drop('Target', axis=1),
df['Target'], test size=0.2, random state=0)
# Scale features
from sklearn.preprocessing import StandardScaler
scaler = StandardScaler()
X train scaled = scaler.fit transform(X train)
X test scaled = scaler.transform(X test)
# Perform linear regression with L2 regularization
from sklearn.linear model import Ridge
regressor = Ridge(alpha=0.1) # alpha is the
regularization strength
regressor.fit(X train scaled, y train)
y pred = regressor.predict(X test scaled)
```



In this example, we use the Ridge class from scikit-learn to perform linear regression with L2 regularization. We set the regularization strength using the alpha parameter, which controls the amount of regularization. A smaller value of alpha corresponds to less regularization, while a larger value corresponds to more regularization.

Feature Selection

In some cases, you may have many features in your dataset, but not all of them may be relevant for predicting the target variable. Feature selection is the process of selecting a subset of the features that are most relevant for predicting the target variable. One way to perform feature selection is to use the SelectKBest class from scikit-learn, which selects the k best features based on a scoring function. Here's an example:

```
# Load dataset
import pandas as pd
df = pd.read csv('my dataset.csv')
# Split data into training and testing sets
from sklearn.model selection import train test split
X train, X test, y train, y test =
train test split(df.drop('Target', axis=1),
df['Target'], test size=0.2, random state=0)
# Scale features
from sklearn.preprocessing import StandardScaler
scaler = StandardScaler()
X train scaled = scaler.fit transform(X train)
X test scaled = scaler.transform(X test)
# Perform feature selection
from sklearn.feature selection import SelectKBest,
f regression
selector = SelectKBest(score func=f regression, k=2)
                                                       #
select the 2 best features based on f regression score
X train selected =
selector.fit transform(X train scaled, y train)
X test selected = selector.transform(X test scaled)
# Perform linear regression with selected features
from sklearn.linear model import LinearRegression
regressor = LinearRegression()
regressor.fit(X train selected, y train)
y pred = regressor.predict(X test selected)
```



In this example, we use the SelectKBest class from scikit-learn to select the 2 best features based on the f_regression scoring function, which measures the linear relationship between each feature and the target variable. We then perform linear regression using only the selected features.

Logistic regression

Logistic regression is a statistical technique used for predicting the probability of a binary outcome variable based on one or more predictor variables. In other words, it is a method for modeling the relationship between a binary response variable and one or more predictor variables.

Python provides various libraries for implementing logistic regression, including scikit-learn, statsmodels, and TensorFlow. In this article, we will use the scikit-learn library to implement logistic regression in Python.

We will use a dataset called "Titanic" for our logistic regression analysis. The Titanic dataset is a classic dataset used for classification tasks. It contains information about passengers on the Titanic, including whether they survived or not, their age, sex, class, and other factors.

Before we start our analysis, we need to import the necessary libraries and load the dataset.

```
import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
from sklearn.linear_model import LogisticRegression
from sklearn.model_selection import train_test_split
from sklearn.metrics import accuracy_score
from sklearn.preprocessing import LabelEncoder
# Load the Titanic dataset
df =
pd.read_csv('https://web.stanford.edu/class/archive/cs/
cs109/cs109.1166/stuff/titanic.csv')
```

Next, we need to preprocess the data by handling missing values and encoding categorical variables. We will drop the columns "Cabin" and "Ticket" because they have too many missing values, and we will fill in the missing values for the "Age" column with the median age.

```
# Drop columns with too many missing values
df.drop(['Cabin', 'Ticket'], axis=1, inplace=True)
```



```
# Fill in missing values for age with the median age
df['Age'].fillna(df['Age'].median(), inplace=True)
# Encode categorical variables
le = LabelEncoder()
df['Sex'] = le.fit_transform(df['Sex'])
df['Embarked'] =
le.fit transform(df['Embarked'].astype(str))
```

Now we can split the data into training and testing sets using the train_test_split() function.

```
# Split the data into training and testing sets
X_train, X_test, y_train, y_test =
train_test_split(df.drop('Survived', axis=1),
df['Survived'], test_size=0.2, random_state=42)
```

We will use the logistic regression algorithm to train a model on the training data and then make predictions on the testing data.

```
# Train the logistic regression model
model = LogisticRegression()
model.fit(X_train, y_train)
# Make predictions on the testing data
y_pred = model.predict(X_test)
```

Finally, we can evaluate the accuracy of our model using the accuracy_score() function.

```
# Evaluate the accuracy of the model
accuracy = accuracy_score(y_test, y_pred)
print("Accuracy: {:.2f}%".format(accuracy * 100))
```

The output of this code will be the accuracy of the model as a percentage.

This is a basic example of logistic regression in Python using scikit-learn. There are many other options and configurations that can be used to fine-tune the model for specific use cases, including regularization, feature scaling, and hyperparameter tuning. However, this example should provide a good starting point for understanding the basics of logistic regression in Python.

Logistic regression is a type of regression analysis used to model the probability of a binary response (e.g. 0 or 1). In logistic regression, the dependent variable is binary and the independent variables can be continuous, categorical or a mix of both.



One popular dataset used to demonstrate logistic regression is the Titanic dataset, which contains information on the passengers aboard the Titanic, including whether or not they survived the disaster. Here's an example of how to load and preprocess the Titanic dataset in Python:

```
# Import necessary libraries
import pandas as pd
import numpy as np
from sklearn.preprocessing import LabelEncoder
# Load the Titanic dataset
df =
pd.read csv('https://web.stanford.edu/class/archive/cs/
cs109/cs109.1166/stuff/titanic.csv')
# Drop columns with too many missing values
df.drop(['Cabin', 'Ticket'], axis=1, inplace=True)
# Fill in missing values for age with the median age
df['Age'].fillna(df['Age'].median(), inplace=True)
# Encode categorical variables
le = LabelEncoder()
df['Sex'] = le.fit transform(df['Sex'])
df['Embarked'] =
le.fit transform(df['Embarked'].astype(str))
```

In this example, we start by importing the necessary libraries and loading the Titanic dataset using pandas. We then drop the "Cabin" and "Ticket" columns because they have too many missing values and fill in the missing values for the "Age" column with the median age. Finally, we encode the categorical variables using LabelEncoder.

Once the data is preprocessed, we can use logistic regression to model the probability of survival based on the other variables in the dataset. Here's an example of how to do this using scikit-learn:

```
# Import necessary libraries
from sklearn.linear_model import LogisticRegression
from sklearn.model_selection import train_test_split
from sklearn.metrics import accuracy_score
# Split the data into training and testing sets
X_train, X_test, y_train, y_test =
train_test_split(df.drop('Survived', axis=1),
df['Survived'], test_size=0.2, random_state=42)
```



```
# Train the logistic regression model
model = LogisticRegression()
model.fit(X_train, y_train)
# Make predictions on the testing data
y_pred = model.predict(X_test)
# Evaluate the accuracy of the model
accuracy = accuracy_score(y_test, y_pred)
print("Accuracy: {:.2f}%".format(accuracy * 100))
```

In this example, we first split the preprocessed data into training and testing sets using train_test_split(). We then train a logistic regression model on the training data using LogisticRegression() and make predictions on the testing data using predict(). Finally, we evaluate the accuracy of the model using accuracy_score().

One important consideration when using logistic regression is selecting the appropriate variables to include in the model. In general, you want to include variables that are predictive of the outcome variable and not highly correlated with each other. One way to visualize the relationship between two variables is to use a scatterplot with different colors for the two response categories. Here's an example of how to create such a plot using the Titanic dataset:

```
# Import necessary libraries
import matplotlib.pyplot as plt
# Create scatterplot of age and fare with color
indicating survival status
plt.scatter(df['Age'][df['Survived'] == 0],
df['Fare'][df['Survived'] == 0], color='red',
label='Not Survived')
plt.scatter(df['Age'][df['Survived'] == 1],
df['Fare'][df['Survived'] == 1], color='green',
label='Survived')
# Add labels and legend
plt.xlabel('Age')
plt.ylabel('Fare')
plt.legend()
plt.show()
```

In this example, we create a scatterplot of "Age" and "Fare" with different colors for passengers who did and did not survive. We can see that there is no clear relationship between "Age" and "Fare", and that survival seems to be somewhat random across different ages and fares.



Another consideration when using logistic regression is interpreting the coefficients of the model. In logistic regression, the coefficients represent the log-odds of the outcome variable given a one-unit increase in the corresponding predictor variable. One way to interpret these coefficients is to exponentiate them, which gives the odds ratio of the outcome variable given a one-unit increase in the predictor variable. Here's an example of how to interpret the coefficients of a logistic regression model using the Titanic dataset:

```
# Print the coefficients of the logistic regression
model
print("Intercept:", model.intercept_)
print("Coefficients:", model.coef_)
# Interpret the coefficients
for i, col in enumerate(df.columns[:-1]):
    print("Odds Ratio for {}: {:.2f}".format(col,
np.exp(model.coef_[0][i])))
```

In this example, we print the intercept and coefficients of the logistic regression model, and then interpret the coefficients by exponentiating them and printing the resulting odds ratio for each predictor variable.

One important concept in logistic regression is regularization, which is a technique used to prevent overfitting by adding a penalty term to the likelihood function. Two commonly used regularization methods are L1 regularization (also known as Lasso) and L2 regularization (also known as Ridge). In scikit-learn, you can specify the regularization method and strength using the penalty and C parameters, respectively. Here's an example of how to perform logistic regression with L1 regularization using the Breast Cancer Wisconsin dataset:

```
# Import necessary libraries
from sklearn.datasets import load_breast_cancer
from sklearn.linear_model import LogisticRegression
from sklearn.model_selection import train_test_split
# Load the Breast Cancer Wisconsin dataset
data = load_breast_cancer()
# Split the data into training and testing sets
X_train, X_test, y_train, y_test =
train_test_split(data.data, data.target, test_size=0.3,
random_state=42)
# Create a logistic regression model with L1
regularization
model = LogisticRegression(penalty='l1', C=1)
```



```
# Fit the model to the training data
model.fit(X_train, y_train)
# Evaluate the model on the testing data
score = model.score(X_test, y_test)
print("Accuracy: {:.2f}%".format(score * 100))
```

In this example, we first load the Breast Cancer Wisconsin dataset and split it into training and testing sets. We then create a logistic regression model with L1 regularization and fit it to the training data. Finally, we evaluate the model's accuracy on the testing data. By adding the L1 penalty term, the model is encouraged to set some of the coefficients to zero, effectively performing feature selection and reducing the risk of overfitting.

Another important concept in logistic regression is cross-validation, which is a technique used to assess the performance of a model and select the best hyperparameters. In scikit-learn, you can perform cross-validation using the cross_val_score function. Here's an example of how to perform logistic regression with cross-validation using the Iris dataset:

```
# Import necessary libraries
from sklearn.datasets import load_iris
from sklearn.linear_model import LogisticRegression
from sklearn.model_selection import cross_val_score
# Load the Iris dataset
data = load_iris()
# Create a logistic regression model
model = LogisticRegression()
# Perform 10-fold cross-validation
scores = cross_val_score(model, data.data, data.target,
cv=10)
# Print the cross-validation scores
print("Cross-validation scores", scores)
print("Mean cross-validation score:", scores.mean())
```

In this example, we first load the Iris dataset and create a logistic regression model. We then perform 10-fold cross-validation using the cross_val_score function, which returns an array of scores for each fold. Finally, we print the cross-validation scores and the mean score across all folds, which gives us an estimate of the model's performance on new, unseen data. By performing cross-validation, we can select the best hyperparameters for the model and avoid overfitting to the training data.



Decision trees

Decision trees are a popular machine learning algorithm used for both classification and regression tasks. In this article, we will explore how to implement decision trees in Python, with a focus on comparing the syntax and functionality to SAS.

First, let's review the basic concept of a decision tree. A decision tree is a model that uses a treelike structure to represent a set of decisions and their possible consequences. At each node in the tree, a decision is made based on the values of one or more input features, and the result of that decision determines which path to follow down the tree. The leaves of the tree represent the final decision or prediction.

Now, let's dive into the Python code. The scikit-learn library is a popular machine learning library in Python, and it includes a decision tree implementation. Here is an example of how to create a decision tree classifier in Python:

```
from sklearn import tree
from sklearn.datasets import load_iris
iris = load_iris()
clf = tree.DecisionTreeClassifier()
clf = clf.fit(iris.data, iris.target)
```

This code loads the Iris dataset (a commonly used dataset in machine learning) and creates a decision tree classifier using scikit-learn's DecisionTreeClassifier class. The fit method is then called to train the classifier on the data.

In SAS, the equivalent code would look like this:

```
proc import datafile='path/to/iris.csv'
    out=iris
    dbms=csv;
run;
proc tree data=iris;
    var petal_length petal_width;
    class species;
run;
```

This code imports the Iris dataset from a CSV file and creates a decision tree using SAS's TREE procedure. The var statement specifies which input features to use, and the class statement specifies the target variable (in this case, the species of iris).



As you can see, the syntax for creating a decision tree in Python and SAS is quite different. However, both implementations have similar functionality, such as specifying the input features and target variable.

One key advantage of using Python for machine learning is the wide range of libraries available, such as scikit-learn, TensorFlow, and PyTorch. These libraries provide powerful tools for building and training machine learning models, and they are often easier to use than SAS's machine learning procedures.

decision trees are a powerful machine learning algorithm that can be implemented in both Python and SAS. While the syntax and functionality may differ between the two languages, both implementations can be used to build accurate and interpretable models. If you are a SAS user looking to learn Python, scikit-learn is a great library to start with for decision trees and other machine learning tasks.

here's an example of how to create a decision tree regression model in Python using scikit-learn:

```
import numpy as np
from sklearn.tree import DecisionTreeRegressor
from sklearn.model selection import train test split
from sklearn.metrics import mean squared error
# Generate some sample data
np.random.seed(0)
n \text{ samples} = 100
X = np.sort(5 * np.random.rand(n samples, 1), axis=0)
y = np.sin(X).ravel()
y[::5] += 3 * (0.5 - np.random.rand(int(n samples/5)))
# Split the data into training and testing sets
X train, X test, y train, y test = train test split(X,
y, test size=0.2, random state=0)
# Create the decision tree regression model
max depth = 2
dt = DecisionTreeRegressor(max depth=max depth)
dt.fit(X train, y train)
# Make predictions on the test set and evaluate the
model
y pred = dt.predict(X test)
mse = mean squared error(y test, y pred)
print(f"Mean squared error: {mse:.2f}")
```



Let's break down this code step by step. First, we import the necessary libraries: NumPy for generating sample data, scikit-learn's DecisionTreeRegressor for creating the model, train_test_split for splitting the data into training and testing sets, and mean_squared_error for evaluating the model.

Next, we generate some sample data using NumPy. We create an array of random x values between 0 and 5, and compute the corresponding y values as the sine function of x with some added noise.

We then split the data into training and testing sets using train_test_split. We use 20% of the data for testing and set the random seed to ensure reproducibility.

Next, we create the decision tree regression model using DecisionTreeRegressor. We set the maximum depth of the tree to 2 to keep the model simple and avoid overfitting to the training data. We then fit the model to the training data using the fit method.

we make predictions on the test set using the predict method and evaluate the model using the mean squared error metric. We print the mean squared error to the console. this code demonstrates how to create a decision tree regression model in Python using scikit-learn and evaluate its performance on a test set.

Decision trees can be used for both classification and regression tasks. In a classification task, the goal is to predict the class of an input based on its features. In a regression task, the goal is to predict a continuous output based on the input features. Here are examples of how to use decision trees for both types of tasks:

Classification Task

Let's start with an example of how to use a decision tree for a classification task in Python. We'll use the famous Iris dataset, which consists of 150 samples of iris flowers with 4 features: sepal length, sepal width, petal length, and petal width. The goal is to predict the species of each iris based on these features.

```
from sklearn.datasets import load_iris
from sklearn.tree import DecisionTreeClassifier
from sklearn.model_selection import train_test_split
from sklearn.metrics import accuracy_score
# Load the Iris dataset
iris = load_iris()
# Split the data into training and testing sets
X_train, X_test, y_train, y_test =
train_test_split(iris.data, iris.target, test_size=0.2,
random_state=0)
```



```
# Create the decision tree classifier
max_depth = 3
clf = DecisionTreeClassifier(max_depth=max_depth)
clf.fit(X_train, y_train)
# Make predictions on the test set and evaluate the
model
y_pred = clf.predict(X_test)
accuracy = accuracy_score(y_test, y_pred)
print(f"Accuracy: {accuracy:.2f}")
```

First, we load the Iris dataset using scikit-learn's load_iris function. We then split the data into training and testing sets using train_test_split.

Next, we create the decision tree classifier using DecisionTreeClassifier and set the maximum depth to 3 to avoid overfitting. We fit the model to the training data using the fit method.

Finally, we make predictions on the test set using the predict method and evaluate the model using the accuracy metric. We print the accuracy to the console.

Regression Task

Now let's look at an example of how to use a decision tree for a regression task in Python. We'll use the Boston Housing dataset, which consists of 506 samples of houses in Boston with 13 features such as crime rate, average number of rooms per dwelling, and distance to employment centers. The goal is to predict the median value of owner-occupied homes in thousands of dollars based on these features.

```
from sklearn.datasets import load_boston
from sklearn.tree import DecisionTreeRegressor
from sklearn.model_selection import train_test_split
from sklearn.metrics import mean_squared_error
# Load the Boston Housing dataset
boston = load_boston()
# Split the data into training and testing sets
X_train, X_test, y_train, y_test =
train_test_split(boston.data, boston.target,
test_size=0.2, random_state=0)
# Create the decision tree regressor
max_depth = 3
dtr = DecisionTreeRegressor(max_depth=max_depth)
```



```
dtr.fit(X_train, y_train)
# Make predictions on the test set and evaluate the
model
y_pred = dtr.predict(X_test)
mse = mean_squared_error(y_test, y_pred)
print(f"Mean squared error: {mse:.2f}")
```

First, we load the Boston Housing dataset using scikit-learn's load_boston function. We then split the data into training and testing sets using train_test_split.

Next, we create the decision tree regressor using DecisionTreeRegressor and set the maximum depth to 3 to avoid overfitting.

Tuning Hyperparameters

When working with decision trees, it's important to tune the hyperparameters to achieve the best performance. The most important hyperparameters for decision trees are:

max_depth: the maximum depth of the tree

min_samples_split: the minimum number of samples required to split an internal node min_samples_leaf: the minimum number of samples required to be at a leaf node max_features: the maximum number of features to consider when splitting a node

Here's an example of how to use grid search to tune the hyperparameters of a decision tree classifier:

```
from sklearn.datasets import load_iris
from sklearn.tree import DecisionTreeClassifier
from sklearn.model_selection import GridSearchCV
# Load the Iris dataset
iris = load_iris()
# Define the hyperparameter grid to search over
param_grid = {
    "max_depth": [2, 3, 4, 5],
    "min_samples_split": [2, 3, 4],
    "min_samples_leaf": [1, 2, 3],
    "max_features": ["sqrt", "log2"]
}
# Create the decision tree classifier
clf = DecisionTreeClassifier()
```



```
# Perform grid search to find the best hyperparameters
grid_search = GridSearchCV(clf, param_grid, cv=5)
grid_search.fit(iris.data, iris.target)
# Print the best hyperparameters and the corresponding
mean cross-validated score
print(f"Best hyperparameters:
{grid_search.best_params_}")
print(f"Best mean cross-validated score:
{grid_search.best_score_:.2f}")
```

First, we load the Iris dataset using scikit-learn's load_iris function.

Next, we define a grid of hyperparameters to search over using a dictionary with keys corresponding to the hyperparameter names and values corresponding to the values to search over.

We create the decision tree classifier using DecisionTreeClassifier.

We then perform grid search using GridSearchCV with 5-fold cross-validation to find the best hyperparameters. The cv parameter specifies the number of cross-validation folds.

Finally, we print the best hyperparameters and the corresponding mean cross-validated score to the console.

Visualizing Decision Trees

Visualizing decision trees can be helpful for understanding how the model makes predictions. Here's an example of how to visualize a decision tree using the plot_tree function in scikit-learn:

```
from sklearn.datasets import load_iris
from sklearn.tree import DecisionTreeClassifier
from sklearn.model_selection import train_test_split
from sklearn.metrics import accuracy_score,
confusion_matrix, plot_confusion_matrix,
classification_report, plot_tree
import matplotlib.pyplot as plt
# Load the Iris dataset
iris = load_iris()
# Split the data into training and testing sets
```



```
X train, X test, y train, y test =
train test split(iris.data, iris.target, test size=0.2,
random state=0)
# Create the decision tree classifier
max depth = 3
clf = DecisionTreeClassifier(max depth=max depth)
clf.fit(X train, y train)
# Make predictions on the test set and evaluate the
model
y pred = clf.predict(X test)
accuracy = accuracy score(y test, y pred)
print(f"Accuracy: {accuracy:.2f}")
# Plot the confusion matrix and classification report
plot confusion matrix(clf, X test, y test)
plt.show()
print(classification report(y test,
```

Random forests

Random forests are a popular machine learning algorithm that can be used for both classification and regression tasks. They are an ensemble method that combines multiple decision trees to make predictions.

In Python, the scikit-learn library provides an implementation of random forests. Here's an example of how to use it:

```
from sklearn.ensemble import RandomForestClassifier
from sklearn.datasets import make_classification
from sklearn.model_selection import train_test_split
# Generate some random data for classification
X, y = make_classification(n_samples=1000,
n_features=10, n_informative=5, random_state=42)
# Split the data into training and testing sets
X_train, X_test, y_train, y_test = train_test_split(X,
y, random_state=42)
```



```
# Create a random forest classifier with 100 trees
clf = RandomForestClassifier(n_estimators=100,
random_state=42)
# Train the classifier on the training data
clf.fit(X_train, y_train)
# Make predictions on the testing data
y_pred = clf.predict(X_test)
# Evaluate the performance of the classifier
from sklearn.metrics import accuracy_score
accuracy = accuracy_score(y_test, y_pred)
print(f"Accuracy: {accuracy}")
```

Let's break down this code step by step:

First, we import the RandomForestClassifier class from the sklearn.ensemble module, as well as some other necessary modules (make_classification and train_test_split).

Next, we generate some random data for classification using the make_classification function. This creates a dataset with 1000 samples, 10 features, and 5 informative features.

We split the data into training and testing sets using the train_test_split function. We use a random_state value of 42 to ensure reproducibility.

We create a RandomForestClassifier object with 100 trees and a random_state value of 42.

We train the classifier on the training data using the fit method.

We make predictions on the testing data using the predict method.

Finally, we evaluate the performance of the classifier using the accuracy_score function from the sklearn.metrics module.

Random forests can also be used for regression tasks. Here's an example of how to use the RandomForestRegressor class in scikit-learn:

```
from sklearn.ensemble import RandomForestRegressor
from sklearn.datasets import make_regression
from sklearn.model_selection import train_test_split
from sklearn.metrics import mean_squared_error
# Generate some random data for regression
X, y = make_regression(n_samples=1000, n_features=10,
n_informative=5, random_state=42)
```



```
# Split the data into training and testing sets
X_train, X_test, y_train, y_test = train_test_split(X,
y, random_state=42)
# Create a random forest regressor with 100 trees
reg = RandomForestRegressor(n_estimators=100,
random_state=42)
# Train the regressor on the training data
reg.fit(X_train, y_train)
# Make predictions on the testing data
y_pred = reg.predict(X_test)
# Evaluate the performance of the regressor
mse = mean_squared_error(y_test, y_pred)
print(f"MSE: {mse}")
```

This code is very similar to the classification example, but we use the RandomForestRegressor class instead of the RandomForestClassifier class. We also use the make_regression function to generate random data for a regression task, and we use the mean_squared_error function to evaluate the performance of the regressor.

Random forests are a type of ensemble learning method, which means that they combine the predictions of multiple models to improve their accuracy. In the case of random forests, the models are decision trees, which are simple models that make predictions based on a series of if-then statements.

A single decision tree can be prone to overfitting, which means that it may fit the training data too closely and not generalize well to new data. Random forests address this problem by creating many different decision trees and combining their predictions. Each tree is trained on a random subset of the training data and a random subset of the features, which helps to reduce overfitting and improve the generalization ability of the model.

In Python, the scikit-learn library provides a simple and easy-to-use implementation of random forests for both classification and regression tasks. Here's an example of how to use a random forest for a binary classification task:

```
from sklearn.ensemble import RandomForestClassifier
from sklearn.datasets import load_breast_cancer
from sklearn.model_selection import train_test_split
from sklearn.metrics import accuracy_score
```

```
# Load the breast cancer dataset
```



```
data = load_breast_cancer()
# Split the data into training and testing sets
X_train, X_test, y_train, y_test =
train_test_split(data.data, data.target,
random_state=42)
# Create a random forest classifier with 100 trees
clf = RandomForestClassifier(n_estimators=100,
random_state=42)
# Train the classifier on the training data
clf.fit(X_train, y_train)
# Make predictions on the testing data
y_pred = clf.predict(X_test)
# Evaluate the performance of the classifier
accuracy = accuracy_score(y_test, y_pred)
print(f"Accuracy: {accuracy}")
```

In this example, we use the load_breast_cancer function from scikit-learn to load the Breast Cancer Wisconsin (Diagnostic) dataset, which is a binary classification dataset with 569 samples and 30 features. We split the data into training and testing sets using the train_test_split function, and then create a random forest classifier with 100 trees using the RandomForestClassifier class. We train the classifier on the training data using the fit method, make predictions on the testing data using the predict method, and evaluate the performance of the classifier using the accuracy_score function from scikit-learn.

Random forests can also be used for regression tasks, such as predicting the price of a house based on its features. Here's an example of how to use a random forest for a regression task:

```
from sklearn.ensemble import RandomForestRegressor
from sklearn.datasets import load_boston
from sklearn.model_selection import train_test_split
from sklearn.metrics import mean_squared_error
# Load the Boston Housing dataset
data = load_boston()
# Split the data into training and testing sets
X_train, X_test, y_train, y_test =
train_test_split(data.data, data.target,
random_state=42)
```



```
# Create a random forest regressor with 100 trees
reg = RandomForestRegressor(n_estimators=100,
random_state=42)
# Train the regressor on the training data
reg.fit(X_train, y_train)
# Make predictions on the testing data
y_pred = reg.predict(X_test)
# Evaluate the performance of the regressor
mse = mean_squared_error(y_test, y_pred)
print(f"MSE: {mse}")
```

In this example, we use the load_boston function from scikit-learn to load the Boston Housing dataset, which is a regression dataset with 506 samples and 13 features.

We train the regressor on the training data using the fit method, make predictions on the testing data using the predict method, and evaluate the performance of the regressor using the mean squared error metric from scikit-learn.

Random forests also provide a useful feature importance measure, which indicates the relative importance of each feature in making predictions. Here's an example of how to obtain the feature importance scores from a random forest classifier:

```
from sklearn.ensemble import RandomForestClassifier
from sklearn.datasets import load_breast_cancer
# Load the breast cancer dataset
data = load_breast_cancer()
# Create a random forest classifier with 100 trees
clf = RandomForestClassifier(n_estimators=100,
random_state=42)
# Train the classifier on the entire dataset
clf.fit(data.data, data.target)
# Get the feature importance scores
importances = clf.feature_importances_
# Print the feature importance scores
for feature, importance in zip(data.feature_names,
importances):
```



print(f"{feature}: {importance}")

In this example, we create a random forest classifier with 100 trees using the RandomForestClassifier class and train the classifier on the entire Breast Cancer Wisconsin (Diagnostic) dataset using the fit method. We then obtain the feature importance scores using the feature_importances_ attribute of the trained classifier and print the scores for each feature using the zip function.

Here are a few more tips and tricks for using random forests in Python:

- 1. Tuning hyperparameters: Random forests have several hyperparameters that can be tuned to improve their performance, such as the number of trees (**n_estimators**), the maximum depth of each tree (**max_depth**), and the minimum number of samples required to split an internal node (**min_samples_split**). One way to tune these hyperparameters is to use a grid search or a randomized search to explore different combinations of hyperparameters and evaluate their performance using cross-validation.
- Dealing with missing values: Random forests can handle missing values in the input data, but different implementations may handle them differently. In scikit-learn, missing values are automatically handled by the RandomForestClassifier and RandomForestRegressor classes, which replace them with the mean or median value of the corresponding feature in the training data.
- 3. Handling imbalanced classes: Random forests can be used for imbalanced classification problems, where the classes are not equally represented in the training data. One way to handle this is to use the **class_weight** hyperparameter, which adjusts the weights of the classes during training to give more importance to the minority class. Another way is to use resampling techniques, such as oversampling the minority class or undersampling the majority class.
- 4. Feature engineering: Random forests can be sensitive to irrelevant or redundant features, so feature selection or feature engineering can be important to improve their performance. One way to do this is to use feature importance scores to identify the most important features and remove the least important ones. Another way is to create new features by combining or transforming existing features, such as adding polynomial features or applying logarithmic or exponential transformations.
- 5. Interpreting results: Random forests provide useful insights into the underlying patterns and relationships in the data, which can be used to interpret the results and make informed decisions. For example, you can use the feature importance scores to identify the most important factors that affect the outcome, or the partial dependence plots to visualize the marginal effects of each feature on the predicted outcome. You can also use the permutation feature importance or the SHAP (SHapley Additive exPlanations) values to explain the contribution of each feature to the prediction for a specific instance.



- 6. Combining models: Random forests can be combined with other models to improve their performance, such as by using an ensemble of random forests or by stacking random forests with other models. Ensemble methods, such as bagging, boosting, and stacking, can help to reduce variance and bias, improve generalization, and increase accuracy. Stacking, in particular, can be used to combine the predictions of multiple models to produce a more accurate and robust prediction.
- 7. Handling categorical variables: Random forests can handle categorical variables in the input data, but different implementations may handle them differently. In scikit-learn, categorical variables can be encoded using one-hot encoding or ordinal encoding, depending on the type and distribution of the categories. One-hot encoding creates a binary feature for each category, while ordinal encoding assigns a numerical value to each category based on their order or frequency.
- 8. Parallel processing: Random forests can be computationally expensive for large datasets or complex models, but parallel processing can help to speed up the training and evaluation process. In scikit-learn, you can use the n_jobs parameter to specify the number of CPUs or cores to use for parallel processing, or you can use the joblib library to distribute the workload across multiple CPUs or machines.
- 9. Handling large datasets: Random forests can be memory-intensive for large datasets, especially if you have many features or a large number of trees. One way to handle this is to use incremental or online learning, where the model is updated incrementally as new data becomes available. Another way is to use feature selection or dimensionality reduction techniques to reduce the number of features or the dimensionality of the data before training the model.
- 10. Avoiding overfitting: Random forests can be prone to overfitting if the model is too complex or if the training data is noisy or biased. One way to avoid overfitting is to use regularization techniques, such as setting a maximum depth for the trees or a minimum number of samples required to split a node. Another way is to use cross-validation to evaluate the performance of the model on a held-out validation set, and to tune the hyperparameters based on the validation performance.
- 11. Handling time-series data: Random forests can be used for time-series data, but they may require additional preprocessing or feature engineering to account for the temporal dependencies and trends in the data. One way to handle this is to use lagged variables or time-based features to capture the temporal relationships between the variables. Another way is to use time-series specific models, such as autoregressive integrated moving average (ARIMA), exponential smoothing (ETS), or recurrent neural networks (RNNs), which can better handle the dynamics and patterns in the data.
- 12. Handling non-linear relationships: Random forests can capture non-linear relationships between the features and the outcome, but they may not be able to capture complex or nonlinear relationships, such as interactions, non-monotonic effects, or nonlinear trends. One way to handle this is to use feature engineering or transformation techniques to create new features that capture the nonlinear relationships, such as adding polynomial or



interaction terms, or applying nonlinear transformations, such as logarithmic or exponential functions. Another way is to use nonlinear models, such as neural networks, support vector machines (SVMs), or decision trees, which can better capture the nonlinear patterns in the data.

- 13. Handling imbalanced data: Random forests can handle imbalanced data by adjusting the class weights or by using sampling techniques, such as oversampling or undersampling, to balance the classes in the training data. In scikit-learn, you can use the **class_weight** parameter to adjust the weights of the classes, or you can use the **imblearn** library to apply sampling techniques, such as SMOTE (Synthetic Minority Over-sampling Technique), which creates synthetic samples of the minority class.
- 14. Choosing hyperparameters: Random forests have several hyperparameters that can affect the performance and complexity of the model, such as the number of trees, the maximum depth of the trees, the number of features to consider at each split, and the minimum number of samples required to split a node. In scikit-learn, you can use grid search or randomized search to explore the hyperparameter space and find the optimal combination of hyperparameters that maximizes the performance on a validation set.
- 15. Visualizing the decision tree: Random forests are based on decision trees, which can be visualized to gain insight into the decision-making process of the model. In scikit-learn, you can use the **export_graphviz** function to export the decision tree as a graph in the DOT format, which can be visualized using Graphviz or other graph visualization tools. You can also use the **plot_tree** function to plot the decision tree in a more compact and interpretable form.
- 16. Handling missing data: Random forests can handle missing data by imputing the missing values or by using surrogate splits, which create additional splits for the missing values based on the correlation between the missing feature and the other features. In scikit-learn, you can use the **SimpleImputer** class to impute the missing values using various strategies, such as mean, median, or most frequent value. You can also use the **missForest** package in R or other imputation methods to handle missing data.
- 17. Handling outliers: Random forests can be sensitive to outliers, which can affect the split criteria and the prediction accuracy. One way to handle outliers is to use robust statistics or outlier detection methods, such as median or quantile regression, or robust covariance estimation. Another way is to remove the outliers or to treat them as a separate class or cluster, depending on their nature and impact on the outcome.
- 18. Handling categorical variables: Random forests can handle categorical variables by converting them into numeric variables, such as dummy variables or ordinal variables, or by using decision trees with categorical splits. In scikit-learn, you can use the **OneHotEncoder** or **OrdinalEncoder** classes to encode the categorical variables into numeric variables, or you can use the **DecisionTreeClassifier** or **DecisionTreeRegressor** classes with the **criterion='entropy'** or **criterion='gini'** parameter to handle categorical splits.



- 19. Handling skewed data: Random forests can handle skewed data by adjusting the split criteria or by using transformations or normalization techniques, such as log or power transformations, or min-max scaling. In scikit-learn, you can use the **PowerTransformer** or **QuantileTransformer** classes to apply various transformations to the data, or you can use the **RobustScaler** or **MinMaxScaler** classes to normalize the data to a certain range.
- 20. Ensembling and stacking: Random forests can be combined with other models or ensembling techniques, such as bagging, boosting, or stacking, to improve the performance and generalization of the model. Bagging combines multiple models trained on bootstrap samples of the data to reduce the variance and improve the stability of the predictions. Boosting combines multiple models trained on weighted versions of the data to focus on the hard-to-predict samples and improve the accuracy. Stacking combines multiple models trained on different subsets or features of the data to capture the complementary strengths of the models and improve the overall performance.
- 21. Interpreting the feature importance: Random forests can provide insights into the relative importance of the features in predicting the outcome, which can be used to identify the key drivers or factors of the outcome and to guide feature selection or feature engineering. In scikit-learn, you can use the **feature_importances_** attribute of the **RandomForestClassifier** or **RandomForestRegressor** class to get the importance scores of the features, or you can use permutation-based methods or model-specific methods to estimate the feature importance.
- 22. Handling multi-output data: Random forests can handle multi-output data, where the outcome has multiple dimensions or targets, by using multi-output regression or classification methods, such as **MultiOutputRegressor** or **MultiOutputClassifier**. In scikit-learn, you can use these classes to train random forests on multi-output data and to make predictions on new data with multiple outputs.

Random forests are a powerful and flexible machine learning method that can be used for a wide range of tasks, from classification and regression to feature selection and feature engineering. With the scikit-learn library in Python, it's easy to train and evaluate random forests on your own data, and to tune their hyperparameters and handle missing values and imbalanced classes. By following these tips and tricks, you can get the most out of random forests and build accurate and robust machine learning models.

here's an example of how to use the RandomForestClassifier class from scikit-learn to train a random forest on a dataset with categorical and numeric features, handle missing values and imbalanced data, and evaluate the model using cross-validation and feature importance:

```
import numpy as np
import pandas as pd
from sklearn.ensemble import RandomForestClassifier
from sklearn.impute import SimpleImputer
from sklearn.compose import ColumnTransformer
```



```
from sklearn.preprocessing import OneHotEncoder,
StandardScaler
from sklearn.model selection import cross val score
# load the data
data = pd.read csv('data.csv')
# separate the features and the target variable
X = data.drop('target', axis=1)
y = data['target']
# define the categorical and numeric features
cat features = ['category', 'color']
num features = ['length', 'width', 'height']
# define the preprocessing steps for the categorical
and numeric features
cat transformer = ColumnTransformer(
    transformers=[('onehot', OneHotEncoder(),
cat features)],
    remainder='passthrough')
num transformer = ColumnTransformer(
    transformers=[('impute',
SimpleImputer(strategy='median'), num features),
                  ('scale', StandardScaler(),
num features)])
# combine the preprocessing steps
preprocessor = ColumnTransformer(
    transformers=[('cat', cat transformer,
cat features),
                  ('num', num transformer,
num features)])
# define the random forest classifier
rf = RandomForestClassifier(n estimators=100,
max depth=5, random state=42)
# combine the preprocessing and modeling steps into a
pipeline
model = Pipeline(steps=[('preprocessor', preprocessor),
                        ('classifier', rf)])
```

```
in stal
```

```
# evaluate the model using cross-validation
scores = cross val score(model, X, y, cv=5,
scoring='accuracy')
print("Accuracy: \$0.2f (+/- \$0.2f)" \$ (scores.mean(),
scores.std() * 2))
# get the feature importances
importances =
model.named steps['classifier'].feature importances
feature names =
model.named steps['preprocessor'].get feature names()
feature names = np.array(feature names)
# plot the feature importances
import matplotlib.pyplot as plt
plt.bar(feature names, importances)
plt.xticks(rotation=90)
plt.xlabel('Features')
plt.ylabel('Importance')
plt.show()
```

In this code, we first load the data from a CSV file and separate the features and the target variable into two arrays. We then define the categorical and numeric features and use the ColumnTransformer class to apply different preprocessing steps to them, such as one-hot encoding for the categorical features and imputation and scaling for the numeric features. We combine the preprocessing steps into a single preprocessor object.

We then define the RandomForestClassifier with some hyperparameters, such as the number of trees (n_estimators) and the maximum depth of each tree (max_depth), and combine it with the preprocessing steps into a Pipeline object. We use the cross_val_score function to evaluate the model using 5-fold cross-validation and the accuracy scoring metric.

Finally, we use the feature_importances_ attribute of the RandomForestClassifier object to get the importance scores of the features, and use the get_feature_names method of the ColumnTransformer object to get the names of the features after preprocessing. We then plot the feature importances using a bar chart.

This code demonstrates how to use random forests in Python to handle various data preprocessing and modeling challenges and to interpret the feature importance of the model. By modifying the hyperparameters, preprocessing steps, or modeling techniques, you can customize the random forest to your specific data and task.

example that demonstrates how to tune the hyperparameters of a random forest using grid search and random search, and how to handle class imbalance using stratified sampling:

in stal

```
import numpy as np
import pandas as pd
from sklearn.ensemble import RandomForestClassifier
from sklearn.impute import SimpleImputer
from sklearn.compose import ColumnTransformer
from sklearn.preprocessing import OneHotEncoder,
StandardScaler
from sklearn.model selection import GridSearchCV,
RandomizedSearchCV, StratifiedKFold
# load the data
data = pd.read csv('data.csv')
# separate the features and the target variable
X = data.drop('target', axis=1)
y = data['target']
# define the categorical and numeric features
cat features = ['category', 'color']
num features = ['length', 'width', 'height']
# define the preprocessing steps for the categorical
and numeric features
cat transformer = ColumnTransformer(
    transformers=[('onehot', OneHotEncoder(),
cat features)],
    remainder='passthrough')
num transformer = ColumnTransformer(
    transformers=[('impute',
SimpleImputer(strategy='median'), num features),
                  ('scale', StandardScaler(),
num features)])
# combine the preprocessing steps
preprocessor = ColumnTransformer(
    transformers=[('cat', cat transformer,
cat features),
                  ('num', num transformer,
num features)])
# define the random forest classifier
rf = RandomForestClassifier(random state=42)
```

in stal

```
# define the hyperparameters to tune
param grid = {
    'n estimators': [50, 100, 200],
    'max depth': [5, 10, 20],
    'min samples split': [2, 5, 10],
    'min samples leaf': [1, 2, 4]
}
# define the search strategy
cv = StratifiedKFold(n splits=5, shuffle=True,
random state=42)
grid search = GridSearchCV(rf, param grid, cv=cv,
scoring='f1 macro')
random search = RandomizedSearchCV(rf,
param distributions=param grid, n iter=20, cv=cv,
scoring='f1 macro', random state=42)
# fit the models and get the best hyperparameters
grid search.fit(X, y)
random search.fit(X, y)
print("Best parameters for grid search:",
grid search.best params )
print("Best score for grid search:",
grid search.best score )
print("Best parameters for random search:",
random search.best params )
print("Best score for random search:",
random search.best score )
```

In this code, we first load the data from a CSV file and separate the features and the target variable into two arrays. We then define the categorical and numeric features and use the ColumnTransformer class to apply different preprocessing steps to them, such as one-hot encoding for the categorical features and imputation and scaling for the numeric features. We combine the preprocessing steps into a single preprocessor object.

We then define the RandomForestClassifier with some hyperparameters set to their default values, and define the hyperparameters we want to tune using a dictionary param_grid. We also define the search strategy using StratifiedKFold to ensure that the class imbalance is preserved during cross-validation.

We then use GridSearchCV and RandomizedSearchCV to search for the best hyperparameters based on the f1_macro scoring metric, which is suitable for imbalanced data. GridSearchCV exhaustively searches all possible combinations of hyperparameters, while



RandomizedSearchCV samples a fixed number of combinations randomly. We fit the models using the fit method and print out the best hyperparameters and scores.

K-Nearest Neighbors

K-Nearest Neighbors (KNN) is a simple and effective algorithm for classification and regression problems in machine learning. In KNN, the output is classified by the majority vote of its k nearest neighbors, where k is a positive integer.

Python is a popular programming language in the field of machine learning due to its vast libraries and easy-to-use syntax. In this article, we will introduce KNN using Python and demonstrate its implementation using an example dataset.

First, let's start with importing the required libraries.

```
import numpy as np
import pandas as pd
from sklearn.neighbors import KNeighborsClassifier
from sklearn.model_selection import train_test_split
from sklearn.metrics import accuracy score
```

We will use the "Breast Cancer Wisconsin (Diagnostic)" dataset from the UCI Machine Learning Repository. This dataset contains information about breast cancer tumors and whether they are malignant or benign.

```
df =
pd.read_csv('https://archive.ics.uci.edu/ml/machine-
learning-databases/breast-cancer-wisconsin/wdbc.data',
header=None)
```

The dataset has 569 instances and 32 features. We will use the first 30 features to predict whether a tumor is malignant or benign.

X = df.iloc[:, 2:32].values
y = df.iloc[:, 1].values

Next, we will split the dataset into training and testing sets.

```
X_train, X_test, y_train, y_test = train_test_split(X,
y, test size=0.3, random state=42)
```

Now, we can create the KNN classifier object and fit it to the training data.



```
k = 5
knn = KNeighborsClassifier(n_neighbors=k)
knn.fit(X_train, y_train)
```

We have set k to 5, which means that the classifier will consider the 5 nearest neighbors to classify a data point. We can now make predictions on the test set.

```
y pred = knn.predict(X test)
```

Finally, we can evaluate the performance of the classifier using the accuracy score.

```
accuracy = accuracy_score(y_test, y_pred)
print('Accuracy:', accuracy)
```

The output will be the accuracy of the classifier on the test set.

This was a brief introduction to KNN using Python. You can experiment with different values of k and try out other datasets to gain more experience with this algorithm.

Here is a longer implementation of KNN in Python using the same Breast Cancer Wisconsin dataset. This implementation includes data preprocessing, hyperparameter tuning, and cross-validation for improved performance.

```
import numpy as np
import pandas as pd
from sklearn.neighbors import KNeighborsClassifier
from sklearn.model selection import train test split,
GridSearchCV, cross val score
from sklearn.preprocessing import StandardScaler
from sklearn.metrics import accuracy score
# Load dataset
df =
pd.read csv('https://archive.ics.uci.edu/ml/machine-
learning-databases/breast-cancer-wisconsin/wdbc.data',
header=None)
# Separate features and target
X = df.iloc[:, 2:32].values
y = df.iloc[:, 1].values
# Split dataset into training and testing sets
X_train, X_test, y_train, y test = train test split(X,
y, test size=0.3, random state=42)
```



```
# Standardize data
scaler = StandardScaler()
X train std = scaler.fit transform(X train)
X test std = scaler.transform(X test)
# Hyperparameter tuning using cross-validation
k \text{ range} = range(1, 31)
param grid = { 'n neighbors ': k range }
knn = KNeighborsClassifier()
grid = GridSearchCV(knn, param grid, cv=10,
scoring='accuracy')
grid.fit(X train std, y train)
print('Best k:', grid.best params ['n neighbors'])
print('Best accuracy:', grid.best score )
# Fit model using best k
k = grid.best params ['n neighbors']
knn = KNeighborsClassifier(n neighbors=k)
knn.fit(X train std, y train)
# Make predictions on test set
y pred = knn.predict(X test std)
# Evaluate performance
accuracy = accuracy_score(y_test, y_pred)
print('Accuracy:', accuracy)
# Cross-validation
cv scores = cross val score(knn, X train std, y train,
cv=10)
print('Cross-validation scores:', cv scores)
print('Mean cross-validation score:',
np.mean(cv scores))
```

In this implementation, we first load the dataset and separate the features and target. Then, we split the dataset into training and testing sets and standardize the data using the StandardScaler class from scikit-learn.

Next, we perform hyperparameter tuning using GridSearchCV to find the best value of k. We search over a range of values for k and use 10-fold cross-validation to evaluate the performance of each value. The best k and its corresponding accuracy score are printed to the console.



We then fit the KNN model using the best value of k and make predictions on the test set. The accuracy score is calculated and printed to the console.

Finally, we perform 10-fold cross-validation on the training set to get a more robust estimate of the model's performance. The cross-validation scores are printed to the console, along with their mean.

How KNN works

KNN is a type of instance-based learning algorithm, where the entire training dataset is stored in memory and used to make predictions on new data. The algorithm works by calculating the distance between the new data point and every point in the training dataset. The k closest points are then used to predict the label of the new data point. For classification problems, the label with the most votes among the k neighbors is chosen as the prediction. For regression problems, the average of the k neighbors' labels is used as the prediction.

The distance between two data points can be calculated using various metrics, such as Euclidean distance, Manhattan distance, and cosine similarity. Euclidean distance is the most commonly used metric in KNN.

Example 1: KNN for binary classification

Let's use the Breast Cancer Wisconsin dataset to demonstrate KNN for binary classification. We will use the first 10 features to predict whether a tumor is malignant or benign.

```
import numpy as np
import pandas as pd
from sklearn.neighbors import KNeighborsClassifier
from sklearn.model selection import train test split
from sklearn.metrics import accuracy score
# Load dataset
df =
pd.read csv('https://archive.ics.uci.edu/ml/machine-
learning-databases/breast-cancer-wisconsin/wdbc.data',
header=None)
# Separate features and target
X = df.iloc[:, 2:12].values
y = df.iloc[:, 1].values
# Split dataset into training and testing sets
X train, X test, y train, y test = train test split(X,
y, test size=0.3, random state=42)
```



```
# Fit KNN model
k = 5
knn = KNeighborsClassifier(n_neighbors=k)
knn.fit(X_train, y_train)
# Make predictions on test set
y_pred = knn.predict(X_test)
# Evaluate performance
accuracy = accuracy_score(y_test, y_pred)
print('Accuracy:', accuracy)
```

In this example, we load the dataset and separate the first 10 features and target. We then split the dataset into training and testing sets and fit the KNN model using k=5. We make predictions on the test set and evaluate the performance of the model using accuracy score. The output will be the accuracy of the model on the test set.

Example 2: Choosing the value of k

The choice of k is an important hyperparameter in KNN. A small value of k may result in overfitting, while a large value may result in underfitting. Let's use the Breast Cancer Wisconsin dataset to illustrate how to choose the value of k using cross-validation.

```
import numpy as np
import pandas as pd
from sklearn.neighbors import KNeighborsClassifier
from sklearn.model selection import train test split,
GridSearchCV, cross_val_score
from sklearn.preprocessing import StandardScaler
from sklearn.metrics import accuracy score
# Load dataset
df =
pd.read csv('https://archive.ics.uci.edu/ml/machine-
learning-databases/breast-cancer-wisconsin/wdbc.data',
header=None)
# Separate features and target
X = df.iloc[:, 2:32].values
y = df.iloc[:, 1].values
# Split dataset into training and testing sets
```



```
X_train, X_test, y_train, y_test = train_test_split(X,
y, test_size=0.3, random_state=42)
# Standardize data
scaler = StandardScaler()
X_train_std = scaler.fit_transform(X_train)
X_test_std = scaler.transform(X_test)
# Hyperparameter tuning using cross-validation
K
```

In this example, we load the dataset and separate the features and target. We split the dataset into training and testing sets and standardize the data using the StandardScaler. We then perform hyperparameter tuning using cross-validation and GridSearchCV to find the best value of k.

The output will be the optimal value of k and the corresponding cross-validation accuracy score.

Example 3: KNN for regression

KNN can also be used for regression problems. Let's use the Boston Housing dataset to demonstrate KNN for regression. We will use the first 10 features to predict the median value of owner-occupied homes.

```
import numpy as np
import pandas as pd
from sklearn.neighbors import KNeighborsRegressor
from sklearn.model selection import train test split
from sklearn.metrics import mean squared error
# Load dataset
df =
pd.read csv('https://archive.ics.uci.edu/ml/machine-
learning-databases/housing/housing.data', header=None,
delim whitespace=True)
# Separate features and target
X = df.iloc[:, :10].values
y = df.iloc[:, -1].values
# Split dataset into training and testing sets
X train, X test, y train, y test = train test split(X,
y, test size=0.3, random state=42)
# Fit KNN model
```



```
k = 5
knn = KNeighborsRegressor(n_neighbors=k)
knn.fit(X_train, y_train)
# Make predictions on test set
y_pred = knn.predict(X_test)
# Evaluate performance
mse = mean_squared_error(y_test, y_pred)
print('MSE:', mse)
```

In this example, we load the dataset and separate the first 10 features and target. We split the dataset into training and testing sets and fit the KNN model using k=5. We make predictions on the test set and evaluate the performance of the model using mean squared error (MSE). The output will be the MSE of the model on the test set.

it is important to note that KNN is a non-parametric algorithm, meaning that it does not make any assumptions about the underlying distribution of the data. This can be beneficial in situations where the distribution of the data is unknown or complex.

In Python, scikit-learn provides an easy-to-use implementation of KNN for both classification and regression problems. It also provides useful tools for data preprocessing, cross-validation, and hyperparameter tuning.

When working with KNN, it is important to choose an appropriate value of k. A value that is too small may result in overfitting, while a value that is too large may result in underfitting. Hyperparameter tuning can be used to find the optimal value of k that maximizes the performance of the model on the validation set.

Example 4: KNN for image classification

KNN can also be used for image classification problems. Let's use the MNIST dataset of handwritten digits to demonstrate KNN for image classification. We will use the first 1000 images in the dataset to train the model and the next 100 images to test the model.

```
import numpy as np
from sklearn.datasets import load_digits
from sklearn.neighbors import KNeighborsClassifier
# Load MNIST dataset
digits = load_digits()
# Separate features and target
X = digits.data[:1000]
y = digits.target[:1000]
```



```
# Fit KNN model
k = 5
knn = KNeighborsClassifier(n_neighbors=k)
knn.fit(X, y)
# Make predictions on test set
y_pred = knn.predict(digits.data[1000:1100])
# Print predicted labels and actual labels
print('Predicted labels:', y_pred)
print('Actual labels:', digits.target[1000:1100])
```

In this example, we load the MNIST dataset and separate the first 1000 images and their labels as the training set. We fit the KNN model using k=5 and make predictions on the next 100 images as the test set. The output will be the predicted labels and actual labels of the test set.

Example 5: KNN for anomaly detection

KNN can also be used for anomaly detection problems. Let's use the credit card fraud detection dataset to demonstrate KNN for anomaly detection. We will use the first 5000 instances of the dataset to train the model and the next 1000 instances to test the model.

```
import numpy as np
import pandas as pd
from sklearn.neighbors import NearestNeighbors
# Load credit card fraud detection dataset
df =
pd.read csv('https://storage.googleapis.com/download.te
nsorflow.org/data/creditcard.csv')
# Separate features and target
X = df.iloc[:, :-1].values
y = df.iloc[:, -1].values
# Split dataset into training and testing sets
X \text{ train} = X[:5000]
X \text{ test} = X[5000:6000]
# Fit KNN model
\mathbf{k} = 5
knn = NearestNeighbors(n neighbors=k)
knn.fit(X train)
```



```
# Compute distances to k nearest neighbors
distances, indices = knn.kneighbors(X_test)
# Calculate anomaly scores
scores = np.mean(distances, axis=1)
# Print anomaly scores
print('Anomaly scores:', scores)
```

In this example, we load the credit card fraud detection dataset and separate the features and target. We split the dataset into training and testing sets and fit the KNN model using k=5. We compute the distances to the k nearest neighbors of each instance in the test set and calculate the anomaly scores as the mean of the distances. The output will be the anomaly scores of the test set.

Example 6: KNN for recommendation systems

KNN can also be used for recommendation systems, where the goal is to recommend items to users based on their past interactions with the system. Let's use the MovieLens dataset to demonstrate KNN for recommendation systems. We will use the first 1000 users and 1000 movies in the dataset to train the model and the next 100 users to test the model.

```
import numpy as np
import pandas as pd
from sklearn.neighbors import NearestNeighbors
# Load MovieLens dataset
df =
pd.read csv('https://raw.githubusercontent.com/grouplen
s/movielens-100k/master/u.data', sep='\t', header=None,
names=['user id', 'item id', 'rating', 'timestamp'])
# Create user-item matrix
matrix = df.pivot table(index='user id',
columns='item id', values='rating')
# Separate training and testing sets
train matrix = matrix.iloc[:1000, :1000].fillna(0)
test matrix = matrix.iloc[1000:1100, :1000].fillna(0)
# Fit KNN model
k = 5
knn = NearestNeighbors(n neighbors=k)
```



```
knn.fit(train matrix)
# Find k nearest neighbors for each user in test set
distances, indices = knn.kneighbors(test matrix)
# Calculate predicted ratings for test set
mean ratings = np.mean(train matrix, axis=1)
pred ratings = np.zeros((test matrix.shape[0],
test matrix.shape[1]))
for i in range(test_matrix.shape[0]):
    for j in range(test matrix.shape[1]):
        if test matrix.iloc[i, j] != 0:
            pred ratings[i, j] =
mean ratings[indices[i]] @ train matrix.iloc[:, j] /
np.sum(mean ratings[indices[i]])
# Print predicted ratings and actual ratings
print('Predicted ratings:', pred ratings)
print('Actual ratings:', test matrix.values)
```

In this example, we load the MovieLens dataset and create a user-item matrix. We separate the first 1000 users and 1000 movies as the training set and the next 100 users as the test set. We fit the KNN model using k=5 and find the k nearest neighbors for each user in the test set. We calculate the predicted ratings for the test set as the weighted average of the ratings of the nearest neighbors. The output will be the predicted ratings and actual ratings of the test set.

It is important to note that KNN has some limitations as well. One major limitation is that it assumes all features are equally important, which may not be true in some cases. In addition, KNN is sensitive to outliers and noisy data, which can affect the accuracy of the model. Furthermore, the choice of k can also affect the performance of the model, as a small value of k may lead to overfitting while a large value of k may lead to underfitting.

To overcome some of these limitations, there are several variations of KNN, such as weighted KNN, distance-weighted KNN, and KNN with kernel density estimation. These variations take into account the importance of different features and adjust the weights of the neighbors accordingly. They can also reduce the impact of outliers and noisy data by giving more weight to the closer neighbors.

here are a few more examples of using KNN in Python:

Example 1: Text classification

KNN can also be used for text classification tasks such as sentiment analysis. Here's an example of using KNN to classify movie reviews as positive or negative based on the words used:

in stal

```
from sklearn.feature extraction.text import
TfidfVectorizer
from sklearn.neighbors import KNeighborsClassifier
# Load data
reviews = ['This movie is great!', 'I hated this
movie', 'The acting was amazing',
           'I fell asleep during this movie', 'The plot
was predictable',
           'I loved the cinematography', 'The dialogue
was terrible']
labels = ['positive', 'negative', 'positive',
'negative', 'negative', 'positive', 'negative']
# Convert reviews to feature vectors
vectorizer = TfidfVectorizer()
X = vectorizer.fit transform(reviews)
# Split data into training and test sets
X_train, X_test, y_train, y test = train test split(X,
labels, test size=0.2)
# Fit KNN model
knn = KNeighborsClassifier(n neighbors=3)
knn.fit(X train, y train)
# Predict test set
y pred = knn.predict(X test)
# Print accuracy
accuracy = accuracy score(y_test, y_pred)
print('Accuracy:', accuracy)
```

In this example, we load a list of movie reviews and their corresponding labels (positive or negative). We convert the reviews to feature vectors using the TfidfVectorizer, which measures the importance of each word in the review. We then split the data into a training set and a test set, fit a KNN model with k=3, and predict the labels of the test set. Finally, we print the accuracy of the model.

Example 2: Image classification

KNN can also be used for image classification tasks. Here's an example of using KNN to classify handwritten digits from the MNIST dataset:



```
from sklearn.datasets import load digits
from sklearn.neighbors import KNeighborsClassifier
# Load data
digits = load digits()
# Split data into training and test sets
X train, X test, y train, y test =
train test split(digits.data, digits.target,
test size=0.2)
# Fit KNN model
knn = KNeighborsClassifier(n neighbors=5)
knn.fit(X train, y train)
# Predict test set
y pred = knn.predict(X test)
# Print accuracy
accuracy = accuracy score(y_test, y_pred)
print('Accuracy:', accuracy)
```

In this example, we load the MNIST dataset of handwritten digits, split the data into a training set and a test set, fit a KNN model with k=5, and predict the labels of the test set. Finally, we print the accuracy of the model.

Example 3: Recommender systems

KNN can also be used to build recommender systems, which suggest items (e.g. movies, products) to users based on their preferences. Here's an example of using KNN to recommend movies to users based on their ratings:

```
import pandas as pd
from sklearn.neighbors import NearestNeighbors
# Load data
ratings = pd.read_csv('ratings.csv')
# Convert data to user-item matrix
user_item = ratings.pivot_table(index='userId',
columns='movieId', values='rating')
# Fit KNN model
knn = NearestNeighbors(n_neighbors=5, metric='cosine')
```



```
knn.fit(user item)
# Get top 5 recommended movies for user 1
user id = 1
user ratings = user item.loc[user id].values.reshape(1,
-1)
distances, indices = knn.kneighbors # Get top 5
recommended movies for user 1
user id = 1
user ratings = user item.loc[user id].values.reshape(1,
-1)
distances, indices = knn.kneighbors(user ratings)
# Print top 5 recommended movies
recommended movie ids = [user item.columns[i] for i in
indices.flatten()]
recommended movies = pd.read csv('movies.csv')
recommended movies =
recommended movies[recommended movies['movieId'].isin(r
ecommended movie ids)]
print(recommended movies.head())
```

In this example, we load a dataset of movie ratings by users and convert it to a user-item matrix, where each row represents a user and each column represents a movie, and the values are the ratings given by the user for the movie. We fit a KNN model with k=5 and cosine distance metric, which measures the similarity between the ratings of users. We then get the top 5 recommended movies for user 1 by finding the nearest neighbors of the user's ratings and selecting the movies they have not rated yet. Finally, we print the recommended movies.

These are just a few examples of the many applications of KNN in machine learning. With its simplicity and versatility, KNN is a powerful tool for many tasks such as classification, regression, and recommendation.

Support Vector Machines

Support Vector Machines (SVMs) are a popular machine learning algorithm for classification and regression problems. In this article, we will explore how to implement SVMs in Python, with a focus on how SAS users can translate their knowledge to Python.



Installing Required Libraries

First, we need to install the required libraries: scikit-learn, pandas, numpy, and matplotlib. We can do this using pip, the package installer for Python:

```
pip install scikit-learn pandas numpy matplotlib
```

Data Preparation

For this example, we will use the iris dataset, which is a well-known dataset for classification problems. The dataset contains 150 samples with four features (sepal length, sepal width, petal length, and petal width) and three target classes (setosa, versicolor, and virginica). We can load the iris dataset using the scikit-learn library:

```
from sklearn.datasets import load_iris
iris = load_iris()
X = iris.data
y = iris.target
```

Splitting Data into Training and Testing Sets

Before building our SVM model, we need to split the data into training and testing sets. This allows us to evaluate the performance of our model on unseen data. We can use the train_test_split function from scikit-learn to split the data:

```
from sklearn.model selection import train test split
```

```
X_train, X_test, y_train, y_test = train_test_split(X,
y, test_size=0.3, random_state=42)
```

This code splits the data into 70% training data and 30% testing data, with a random seed of 42 for reproducibility.

Building the SVM Model

Now, we can build our SVM model using the SVC class from scikit-learn. We will use a linear kernel for simplicity.

```
from sklearn.svm import SVC
svm = SVC(kernel='linear', C=1, random_state=42)
svm.fit(X_train, y_train)
```



This code creates an SVM model with a linear kernel, a regularization parameter (C) of 1, and a random seed of 42. We then fit the model to the training data.

Evaluating the SVM Model

To evaluate the performance of our SVM model, we can use the score method to calculate the accuracy on the test data:

```
accuracy = svm.score(X_test, y_test)
print("Accuracy:", accuracy)
```

This code calculates the accuracy of the model on the test data and prints it to the console.

Visualizing the SVM Model

Finally, we can visualize the SVM model using matplotlib. We will plot the decision boundary and the support vectors.

css

```
import matplotlib.pyplot as plt
plt.scatter(X[:, 0], X[:, 1], c=y, cmap='viridis')
# plot the decision function
ax = plt.gca()
xlim = ax.get xlim()
ylim = ax.get ylim()
xx = np.linspace(xlim[0], xlim[1], 30)
yy = np.linspace(ylim[0], ylim[1], 30)
YY, XX = np.meshgrid(yy, xx)
xy = np.vstack([XX.ravel(), YY.ravel()]).T
Z = svm.decision function(xy).reshape(XX.shape)
# plot decision boundary and margins
ax.contour(XX, YY, Z, colors='k', levels=[-1, 0, 1],
alpha=0.5,
           linestyles=['--', '-', '--'])
# plot support vectors
ax.scatter(svm.support vectors [:, 0],
svm.support vectors [:, 1], s=100,
           linewidth=1, facecolors='none', edgecolors
```

This code first creates a scatter plot of the data, with the color of each point indicating the target class. We then plot the decision boundary and the margins using the decision_function method of the SVM model. Finally, we plot the support vectors as circles with no fill.

Support Vector Machines (SVMs) are a type of supervised learning algorithm used for classification and regression analysis. SVMs are particularly useful when the data is not linearly separable, meaning that it cannot be divided into distinct classes using a straight line. SVMs work by finding a hyperplane (a line or a plane in higher dimensions) that separates the data into distinct classes with the largest possible margin between the classes. This hyperplane is then used to make predictions on new data.

There are several types of SVMs, including linear SVMs, polynomial SVMs, and radial basis function (RBF) SVMs. In this example, we will focus on linear SVMs.

Example: Classifying Iris Species using Linear SVMs

In this example, we will use the iris dataset to build a linear SVM model that can classify the three different species of iris flowers. The iris dataset is included in scikit-learn and can be loaded using the load_iris function. Here is the code to load the data and split it into training and testing sets:

```
from sklearn.datasets import load_iris
from sklearn.model_selection import train_test_split
iris = load_iris()
X_train, X_test, y_train, y_test =
train_test_split(iris.data, iris.target, test_size=0.3,
random state=42)
```

The data is split into 70% training data and 30% testing data.

Next, we will build a linear SVM model using the Support Vector Classifier (SVC) class from scikit-learn. Here is the code to build the model:

```
from sklearn.svm import SVC
svm = SVC(kernel='linear', C=1, random_state=42)
svm.fit(X_train, y_train)
```

We use the linear kernel and set the regularization parameter (C) to 1. The random_state parameter is set to 42 for reproducibility. We then fit the model to the training data.

To evaluate the performance of the model, we will calculate the accuracy on the test data:



```
accuracy = svm.score(X_test, y_test)
print("Accuracy:", accuracy)
```

The accuracy of the model is printed to the console.

Finally, we can visualize the decision boundary and the support vectors of the SVM model using matplotlib:

```
import numpy as np
import matplotlib.pyplot as plt
# create a mesh to plot in
x \min, x \max = X \operatorname{train}[:, 0] \operatorname{.min}() - 1, X \operatorname{train}[:, 0]
0].max() + 1
y \min, y \max = X \operatorname{train}[:, 1].\min() - 1, X \operatorname{train}[:,
1].max() + 1
xx, yy = np.meshgrid(np.arange(x min, x max, 0.1),
                        np.arange(y min, y max, 0.1))
# plot the decision boundary
Z = svm.predict(np.c [xx.ravel(), yy.ravel()])
Z = Z.reshape(xx.shape)
plt.contourf(xx, yy, Z, alpha=0.4)
# plot the training data
plt.scatter(X train[:, 0], X train[:, 1], c=y train,
alpha=0.8)
# plot the support vectors
plt.scatter(svm.support vectors [:, 0],
svm.support vectors [:, 1],
              s=80, facecolors='none', edgecolors='k')
plt.title('Linear SVM')
plt.show()
```

The code first creates a mesh of points to plot the decision boundary. We then plot the decision boundary and the training data, with each point colored according to its target class. Finally, we plot the support vectors as circles with no fill. The resulting plot should show a clear separation between the three different species of iris flowers.

Example: Classifying Breast Cancer Tumors using SVMs

In this example, we will use the Breast Cancer Wisconsin (Diagnostic) dataset to build an SVM model that can classify tumors as either benign or malignant. The dataset is included in scikit-



learn and can be loaded using the load_breast_cancer function. Here is the code to load the data and split it into training and testing sets:

from sklearn.datasets import load_breast_cancer
from sklearn.model_selection import train_test_split
cancer = load_breast_cancer()
X_train, X_test, y_train, y_test =
train_test_split(cancer.data, cancer.target,
test_size=0.3, random_state=42)

The data is split into 70% training data and 30% testing data.

Next, we will build an SVM model using the Support Vector Classifier (SVC) class from scikitlearn. Here is the code to build the model:

```
from sklearn.svm import SVC
svm = SVC(kernel='linear', C=1, random_state=42)
svm.fit(X_train, y_train)
```

We use the linear kernel and set the regularization parameter (C) to 1. The random_state parameter is set to 42 for reproducibility. We then fit the model to the training data.

To evaluate the performance of the model, we will calculate the accuracy, precision, recall, and F1 score on the test data:

```
from sklearn.metrics import accuracy_score,
precision_score, recall_score, f1_score
y_pred = svm.predict(X_test)
accuracy = accuracy_score(y_test, y_pred)
precision = precision_score(y_test, y_pred)
recall = recall_score(y_test, y_pred)
f1 = f1_score(y_test, y_pred)
print("Accuracy:", accuracy)
print("Precision:", precision)
print("Recall:", recall)
print("F1 Score:", f1)
```



The metrics are printed to the console.

Finally, we can visualize the performance of the model using a confusion matrix:

from sklearn.metrics import confusion_matrix
confusion = confusion_matrix(y_test, y_pred)
print("Confusion Matrix:\n", confusion)

The confusion matrix is printed to the console.

the key points to keep in mind when working with SVMs in Python:

- 1. SVMs are a machine learning algorithm used for classification and regression analysis.
- 2. SVMs find the optimal hyperplane that separates the data points into different classes.
- 3. The choice of kernel function and hyperparameters can greatly affect the performance of the SVM model.
- 4. Scikit-learn is a popular Python library for machine learning that provides easy-to-use implementations of SVMs.
- 5. When using SVMs, it is important to preprocess the data and scale the features to ensure optimal performance.
- 6. To evaluate the performance of an SVM model, use metrics such as accuracy, precision, recall, and F1 score, as well as confusion matrices.
- 7. SVMs are a powerful and flexible tool for machine learning, but can be computationally intensive and may require parameter tuning for optimal performance.

Naive Bayes

Naive Bayes is a classification algorithm that uses Bayes' theorem to calculate the probability of a given data point belonging to a particular class. In this article, we will explore how to implement Naive Bayes classification in Python, with a focus on how SAS users can leverage their existing skills to learn Python.

First, we need to install the necessary packages. We will be using the scikit-learn library for machine learning in Python, which can be installed using pip:



pip install scikit-learn

Once the package is installed, we can import the necessary modules:

```
from sklearn.datasets import load_iris
from sklearn.naive_bayes import GaussianNB
from sklearn.model_selection import train_test_split
from sklearn.metrics import accuracy score
```

Next, we will load the Iris dataset, which is a commonly used dataset for classification problems. The dataset contains 150 samples of iris flowers, with 50 samples each for three different species: setosa, versicolor, and virginica. The features of each sample include the length and width of the sepal and petal.

iris = load_iris()
X = iris.data
y = iris.target

We can then split the dataset into training and testing sets using the train_test_split function:

```
X_train, X_test, y_train, y_test = train_test_split(X,
y, test_size=0.3, random_state=42)
```

We will use the Gaussian Naive Bayes algorithm, which assumes that the features are normally distributed. We can initialize the GaussianNB class and fit the training data to the model:

gnb = GaussianNB()
gnb.fit(X_train, y_train)

We can then use the trained model to make predictions on the testing data:

y_pred = gnb.predict(X_test)

Finally, we can evaluate the accuracy of the model using the accuracy_score function:

```
accuracy = accuracy_score(y_test, y_pred)
print("Accuracy:", accuracy)
```

The complete code is as follows:

```
from sklearn.datasets import load_iris
from sklearn.naive_bayes import GaussianNB
from sklearn.model_selection import train_test_split
from sklearn.metrics import accuracy_score
```



```
# Load the Iris dataset
iris = load_iris()
X = iris.data
y = iris.target
# Split the dataset into training and testing sets
X_train, X_test, y_train, y_test = train_test_split(X,
y, test_size=0.3, random_state=42)
# Train the Gaussian Naive Bayes model
gnb = GaussianNB()
gnb.fit(X_train, y_train)
# Make predictions on the testing data
y_pred = gnb.predict(X_test)
# Evaluate the accuracy of the model
accuracy = accuracy_score(y_test, y_pred)
print("Accuracy:", accuracy)
```

In this example, we used the Iris dataset as an example. However, the same steps can be used for any classification problem using Naive Bayes. The scikit-learn library also provides other types of Naive Bayes algorithms, such as Multinomial Naive Bayes and Bernoulli Naive Bayes, which can be used depending on the nature of the problem.

the Naive Bayes algorithm is a powerful and simple method for classification problems, and Python provides a variety of libraries and tools to make it easy to implement. For SAS users, the transition to Python can be made easier by leveraging their existing knowledge of data manipulation and analysis, as well as the availability of libraries like scikit-learn that provide similar functionality to SAS.

In addition to the Gaussian Naive Bayes algorithm, scikit-learn also provides Multinomial Naive Bayes and Bernoulli Naive Bayes algorithms. These algorithms are suited for text classification problems, where the features are the frequency of occurrence of words in a document. The Multinomial Naive Bayes algorithm is used when the features are discrete and represent the count of occurrences of a word, while the Bernoulli Naive Bayes algorithm is used when the features are binary and represent whether a word appears in a document or not.

Let's take a look at an example of text classification using the Multinomial Naive Bayes algorithm. We will use the 20 Newsgroups dataset, which contains 20,000 newsgroup posts across 20 different newsgroups. The task is to classify the newsgroup posts into their respective categories based on the content of the post.

We can load the dataset using the fetch_20newsgroups function:



```
from sklearn.datasets import fetch_20newsgroups
# Load the 20 Newsgroups dataset
categories = ['alt.atheism', 'talk.religion.misc',
'comp.graphics', 'sci.space']
newsgroups_train = fetch_20newsgroups(subset='train',
categories=categories)
newsgroups_test = fetch_20newsgroups(subset='test',
categories=categories)
```

Next, we need to preprocess the data to convert the text into features that can be used for classification. We will use the CountVectorizer class from scikit-learn to convert the text into a bag-of-words representation, where each word in the text is represented by its frequency of occurrence in the document.

```
from sklearn.feature_extraction.text import
CountVectorizer
# Convert the text into a bag-of-words representation
vectorizer = CountVectorizer(stop_words='english')
X_train =
vectorizer.fit_transform(newsgroups_train.data)
X_test = vectorizer.transform(newsgroups_test.data)
y_train = newsgroups_train.target
y test = newsgroups test.target
```

We can then initialize the MultinomialNB class and fit the training data to the model:

```
from sklearn.naive bayes import MultinomialNB
```

```
# Train the Multinomial Naive Bayes model
mnb = MultinomialNB()
mnb.fit(X_train, y_train)
```

We can make predictions on the testing data using the predict method:

```
# Make predictions on the testing data
y pred = mnb.predict(X test)
```

Finally, we can evaluate the accuracy of the model using the accuracy_score function:

```
from sklearn.metrics import accuracy score
```

```
# Evaluate the accuracy of the model
```



```
accuracy = accuracy_score(y_test, y_pred)
print("Accuracy:", accuracy)
```

The complete code is as follows:

```
from sklearn.datasets import fetch 20newsgroups
from sklearn.feature extraction.text import
CountVectorizer
from sklearn.naive bayes import MultinomialNB
from sklearn.metrics import accuracy score
# Load the 20 Newsgroups dataset
categories = ['alt.atheism', 'talk.religion.misc',
'comp.graphics', 'sci.space']
newsgroups train = fetch 20newsgroups(subset='train',
categories=categories)
newsgroups test = fetch 20newsgroups(subset='test',
categories=categories)
# Convert the text into a bag-of-words representation
vectorizer = CountVectorizer(stop words='english')
X train =
vectorizer.fit transform(newsgroups train.data)
X test = vectorizer.transform(newsgroups test.data)
y train = newsgroups train.target
y test =
```

let's continue. We can make predictions on the testing data using the predict method:

```
# Make predictions on the testing data
y pred = mnb.predict(X test)
```

Finally, we can evaluate the accuracy of the model using the accuracy_score function:

```
# Evaluate the accuracy of the model
accuracy = accuracy_score(y_test, y_pred)
print("Accuracy:", accuracy)
```

The complete code is as follows:

```
from sklearn.datasets import fetch_20newsgroups
from sklearn.feature_extraction.text import
CountVectorizer
```



```
from sklearn.naive bayes import MultinomialNB
from sklearn.metrics import accuracy score
# Load the 20 Newsgroups dataset
categories = ['alt.atheism', 'talk.religion.misc',
'comp.graphics', 'sci.space']
newsgroups train = fetch 20newsgroups(subset='train',
categories=categories)
newsgroups test = fetch 20newsgroups(subset='test',
categories=categories)
# Convert the text into a bag-of-words representation
vectorizer = CountVectorizer(stop words='english')
X train =
vectorizer.fit transform(newsgroups train.data)
X test = vectorizer.transform(newsgroups test.data)
y train = newsgroups train.target
y test = newsgroups test.target
# Train the Multinomial Naive Bayes model
mnb = MultinomialNB()
mnb.fit(X train, y train)
# Make predictions on the testing data
y pred = mnb.predict(X test)
# Evaluate the accuracy of the model
accuracy = accuracy score(y test, y pred)
print("Accuracy:", accuracy)
```

When we run this code, we get an accuracy of around 82%. This means that our model is able to correctly classify the newsgroup posts into their respective categories around 82% of the time.

Naive Bayes is a simple yet powerful algorithm for classification problems. It is particularly useful for text classification problems, where the features are the frequency of occurrence of words in a document. Scikit-learn provides implementations of Gaussian Naive Bayes, Multinomial Naive Bayes, and Bernoulli Naive Bayes algorithms, making it easy to use Naive Bayes for classification problems in Python.

Chapter 7: Advanced Topics in Python



Object-oriented programming

Python is a popular programming language known for its simplicity, flexibility, and powerful libraries. It is widely used in various domains such as web development, data analysis, machine learning, and scientific computing. One of the key features of Python is its support for object-oriented programming (OOP) which enables developers to write reusable and maintainable code.

For SAS users who are looking to expand their skills, Python can be a great addition to their toolkit. In this article, we will provide an introduction to Python from a SAS-oriented perspective, focusing on its OOP capabilities.

What is Object-Oriented Programming?

Object-oriented programming is a programming paradigm that organizes code into objects, which are instances of classes. A class is a blueprint for creating objects that defines its properties and behavior. Properties are the attributes that describe an object, while behavior refers to the actions that an object can perform.

In OOP, the key concepts are encapsulation, inheritance, and polymorphism. Encapsulation refers to the practice of hiding the internal workings of an object and exposing only its interface to the outside world. Inheritance enables a new class to be based on an existing class, inheriting its properties and behavior. Polymorphism allows different objects to be treated as if they were of the same type, by using a common interface.

Python and Object-Oriented Programming

Python is a fully object-oriented language, which means that everything in Python is an object. Even the simplest data types like integers and strings are objects. This makes Python well-suited for OOP, as developers can easily create and manipulate objects using Python's built-in features.

Creating a Class in Python

In Python, a class is defined using the class keyword, followed by the class name and a colon. Here's an example of a simple class definition:

```
class Rectangle:
    def __init__(self, length, width):
        self.length = length
        self.width = width
    def area(self):
        return self.length * self.width
```

This class, called Rectangle, has two properties length and width, and one method area. The ______init___ method is a special method called a constructor that is called when a new object is



created. It initializes the length and width properties of the object. The area method calculates the area of the rectangle.

Creating an Object in Python

To create an object of a class in Python, you simply call the class like a function and pass in the required arguments. Here's an example:

my rectangle = Rectangle(5, 3)

This creates a new object of the Rectangle class with a length of 5 and a width of 3. We can access the properties of this object using dot notation, like this:

print(my_rectangle.length) # Output: 5
print(my rectangle.width) # Output: 3

We can also call the area method of the object:

```
print(my_rectangle.area()) # Output: 15
Inheritance in Python
```

In Python, inheritance is accomplished using the super function, which allows a subclass to inherit properties and behavior from its parent class. Here's an example:

```
class Square(Rectangle):
    def __init__(self, side):
        super().__init__(side, side)
```

This class, called Square, inherits from the Rectangle class and overrides the __init__ method to accept a single argument, side, which is used to set both the length and width properties of the object to the same value.

use of interfaces. An interface is a blueprint for a set of methods that a class must implement in order to be considered a particular type. Python does not have built-in support for interfaces, but they can be simulated using abstract base classes (ABCs) from the abc module.

Here's an example of an ABC that defines an interface for a shape:

```
from abc import ABC, abstractmethod
class Shape(ABC):
   @abstractmethod
   def area(self):
        pass
```



This class, called Shape, is an abstract base class that defines an interface for a shape. It has one abstract method, area, which must be implemented by any class that inherits from it.

Here's an example of a class that implements the Shape interface:

```
class Circle(Shape):
    def __init__(self, radius):
        self.radius = radius
    def area(self):
        return 3.14 * self.radius ** 2
```

This class, called Circle, inherits from the Shape class and implements the area method. It calculates the area of a circle using the formula $pi * r^2$.

Using Python's OOP Capabilities in SAS

SAS users can leverage Python's OOP capabilities by using SAS's PROC PYTHON procedure, which allows SAS code to call Python code. This enables users to create Python classes and objects from within SAS, and use them to perform various tasks such as data manipulation, statistical analysis, and machine learning.

Here's an example of using the PROC PYTHON procedure to create a Python class and object from within SAS:

```
proc python;
submit;
class Rectangle:
    def __init__(self, length, width):
        self.length = length
        self.width = width
    def area(self):
        return self.length * self.width
    my_rectangle = Rectangle(5, 3)
endsubmit;
quit;
```

This code creates a Rectangle class and object in Python, and initializes the object with a length of 5 and a width of 3. The PROC PYTHON procedure then returns the my_rectangle object to SAS, where it can be used in subsequent SAS code.

In addition to using PROC PYTHON, SAS users can also integrate Python code directly into SAS code using the SASPy package. SASPy is a Python package that provides a Python



interface to SAS, allowing users to call SAS procedures and functions from within Python code. This enables users to leverage Python's OOP capabilities and SAS's analytical capabilities in the same codebase.

Here's an example of using SASPy to call a SAS procedure from within Python code:

```
import saspy
sas = saspy.SASsession()
sas.submit('proc means data=sashelp.iris; run;')
sas.disconnect()
```

This code creates a SAS session using the SASsession class from the SASPy package, and submits a PROC MEANS procedure on the sashelp.iris dataset. The disconnect method is then called to close the SAS session.

Using Python's OOP capabilities in SAS can provide SAS users with a powerful set of tools for developing analytical solutions. By combining SAS's analytical capabilities with Python's flexibility and ease of use, users can develop sophisticated solutions to complex problems.

Python's support for object-oriented programming makes it a powerful language for developing reusable and maintainable code. SAS users can take advantage of Python's OOP capabilities by using SAS's PROC PYTHON procedure or SASPy package to create and use Python classes and objects from within SAS. By combining SAS's analytical capabilities with Python's flexibility and ease of use, users can develop sophisticated solutions to complex problems.

Here's an example of creating a Python class in SAS using the PROC PYTHON procedure and using it in SAS code:

```
proc python;
submit;
class Rectangle:
    def __init__(self, length, width):
        self.length = length
        self.width = width
    def area(self):
        return self.length * self.width
    endsubmit;
quit;
data rectangle;
length = 5;
```



```
width = 3;
proc python;
submit;
my_rectangle = Rectangle(length, width)
endsubmit;
out = getvalue('my_rectangle.area()');
quit;
area = input(out, 8.);
run;
proc print data=rectangle;
run;
```

In this code, we create a Rectangle class in Python using the PROC PYTHON procedure, and then use it in a SAS data step to calculate the area of a rectangle with a length of 5 and a width of 3. The SUBMIT block creates a Rectangle object with the specified length and width, and then calculates its area using the area() method. The GETVALUE function retrieves the result of the calculation from the Python session, and the INPUT function converts it to a SAS numeric value. Finally, the AREA variable is assigned the calculated area, and the RECTANGLE data set is printed to verify the result.

Here's an example of using SASPy to call a SAS procedure from within Python code:

```
import saspy
sas = saspy.SASsession()
sas.submit('proc means data=sashelp.iris; run;')
results = sas.sasdata('MEANSOUT').to_df()
sas.disconnect()
print(results.head())
```

In this code, we create a SAS session using the SASsession class from the SASPy package, and then submit a PROC MEANS procedure on the sashelp.iris dataset. The SASDATA function retrieves the output of the procedure as a Pandas dataframe, and the DISCONNECT method is called to close the SAS session. Finally, the results are printed to verify the output.

```
class SalesData:
    def __init__(self, file):
        self.file = file
```



```
self.sas = saspy.SASsession()
        self.sas.submit(f"data sales; set {file};
run;")
    def total sales(self):
        results = self.sas.submit('proc means
data=sales; var sales; output out=total sales sum=;
run;')
        return results['total sales'].to df()
    def sales by region(self):
        results = self.sas.submit('proc sql; select
region, sum(sales) as total sales from sales group by
region; guit;')
        return results['WORK.SQL 1'].to df()
    def disconnect(self):
        self.sas.disconnect()
sales data = SalesData('sales data.csv')
total sales = sales data.total sales()
print(f"Total sales: {total sales['Sum'][0]}")
sales by region = sales data.sales by region()
print(sales by region)
sales data.disconnect()
```

In this code, we define a SalesData class that reads a sales data file into SAS using the SASsession class from SASPy. The total_sales method calculates the total sales from the data using a PROC MEANS procedure, and the sales_by_region method calculates the total sales by region using a PROC SQL procedure. Finally, the disconnect method is called to close the SAS session.

We create an instance of the SalesData class with the filename of the sales data file, and then call the total_sales and sales_by_region methods to calculate the sales data.

By using Python's OOP capabilities in SAS, SAS users can develop more flexible and maintainable code, and leverage the powerful analytical capabilities of both SAS and Python in the same codebase.

this time using the SWIGSAS package to create a Python class that calls a SAS macro:

import swigsas



```
class DataPrep:
    def init (self, data):
        self.data = data
        self.sas = swigsas.SASsession()
        self.sas.submit('%macro preprocess(data); data
&data; set &data; if age < 18 then age grp = "Under</pre>
18"; else if age < 30 then age grp = "\overline{18}-29"; else if
age < 50 then age grp = "30-49"; else age grp = "50+";
run; %mend;')
    def preprocess data(self):
self.sas.submit(f'%preprocess(data={self.data});')
        results = self.sas.sasdata('WORK.DATA').to df()
        return results
    def disconnect(self):
        self.sas.disconnect()
data prep = DataPrep('customer data.csv')
preprocessed data = data prep.preprocess_data()
print(preprocessed data.head())
data prep.disconnect()
```

In this code, we define a DataPrep class that preprocesses customer data by creating a new variable age_grp based on the age variable. The preprocessing is performed using a SAS macro defined using the %MACRO and %MEND statements.

We create an instance of the DataPrep class with the filename of the customer data file, and then call the preprocess_data method to preprocess the data using the SAS macro. The SASDATA function retrieves the preprocessed data as a Pandas dataframe, and the results are printed to the console. Finally, the disconnect method is called to close the SAS session.

By using Python's OOP capabilities in SAS, we can create more modular and maintainable code that combines the power of SAS macros with the flexibility and ease of use of Python.

Here's another example of using Python's OOP capabilities in SAS, this time using the saspy package to create a Python class that reads in data from a SAS dataset, performs a linear regression, and plots the results:

import saspy
import pandas as pd
import seaborn as sns



```
import matplotlib.pyplot as plt
class RegressionAnalysis:
    def init (self, dataset):
        self.dataset = dataset
        self.sas = saspy.SASsession()
    def run regression(self, x vars, y var):
        sas code = f"""proc reg data={self.dataset};
model {y var} = { ' '.join(x vars) } /
scatterplot=matrix; run;"""
        results = self.sas.submit(sas code)
        return results
    def plot results(self, x vars, y var):
        sas code = f"""proc reg data={self.dataset};
model {y var} = { ' '.join(x vars) } / outest=est;
run;"""
        self.sas.submit(sas code)
        est = self.sas.sasdata('EST').to df()
        sns.pairplot(pd.concat([est,
self.sas.sasdata(self.dataset).to df()[[y var] +
x vars]], axis=1), x vars=x vars, y vars=y var)
        plt.show()
    def disconnect(self):
        self.sas.disconnect()
analysis = RegressionAnalysis('cars')
results = analysis.run regression(['mpg', 'weight'],
'horsepower')
print(results['Output'].to string())
analysis.plot results(['mpg', 'weight'], 'horsepower')
analysis.disconnect()
```

In this code, we define a RegressionAnalysis class that performs a linear regression on a SAS dataset using the PROC REG procedure. The run_regression method takes a list of predictor variables x_vars and a response variable y_var, and returns the results of the regression as a SAS output object.

The plot_results method also takes x_vars and y_var, but this time it performs the regression and then plots the results using the seaborn and matplotlib libraries. The regression results are stored



in a temporary SAS dataset named EST, which is read in using the SASDATA function from saspy, and then concatenated with the original dataset for plotting.

We create an instance of the RegressionAnalysis class with the name of the SAS dataset cars, and then call the run_regression method with ['mpg', 'weight'] as x_vars and horsepower as y_var. The results are then printed to the console using the to_string method of the SAS output object. Finally, the plot_results method is called to plot the regression results, and the disconnect method is called to close the SAS session.

By using Python's OOP capabilities in SAS, we can create more flexible and customizable code that combines the analytical power of SAS with the visualization capabilities of Python libraries like seaborn and matplotlib.

Here's another example of using Python's OOP capabilities in SAS, this time using the swat package to create a Python class that performs a logistic regression on a SAS dataset and calculates the accuracy of the model:

```
import swat
from sklearn.metrics import accuracy score
class LogisticRegression:
    def init (self, dataset):
        self.dataset = dataset
        self.conn = swat.CAS('server', port,
'username', 'password')
    def fit(self, x vars, y var):
        tbl = self.conn.CASTable(self.dataset)
        model = tbl.logistic(target=y var,
inputs=x vars)
        return model
    def predict(self, x vars, y var):
        tbl = self.conn.CASTable(self.dataset)
        model = tbl.logistic(target=y_var,
inputs=x vars)
        preds = model.predict(tbl)
        return preds
    def accuracy(self, x_vars, y_var):
        tbl = self.conn.CASTable(self.dataset)
        model = tbl.logistic(target=y var,
inputs=x vars)
        preds = model.predict(tbl)
```



```
true_vals =
tbl[y_var].to_frame().values.ravel()
    return accuracy_score(true_vals, preds)
    def disconnect(self):
        self.conn.close()
logreg = LogisticRegression('credit')
model = logreg.fit(['income', 'age', 'debtinc'],
    'default')
preds = logreg.predict(['income', 'age', 'debtinc'],
    'default')
accuracy = logreg.accuracy(['income', 'age',
    'debtinc'], 'default')
print(f"Accuracy: {accuracy}")
logreg.disconnect()
```

In this code, we define a LogisticRegression class that performs a logistic regression on a SAS dataset using the CASTable class from swat. The fit method takes a list of predictor variables x_vars and a response variable y_var , and returns the fitted logistic regression model. The predict method also takes x_vars and y_var , but this time it uses the fitted model to make predictions on the input data. The accuracy method takes x_vars and y_var and calculates the accuracy of the model by comparing the predicted values to the true values using the accuracy_score function from sklearn.metrics.

We create an instance of the LogisticRegression class with the name of the SAS dataset credit, and then call the fit method with ['income', 'age', 'debtinc'] as x_vars and default as y_var. The fitted model is stored in the model variable, and the predict method is called with the same x_vars and y_var to make predictions on the input data. The accuracy method is then called to calculate the accuracy of the model, and the results are printed to the console.

Finally, the disconnect method is called to close the connection to the SAS server.

By using Python's OOP capabilities in SAS, we can create more powerful and flexible machine learning models that leverage the computational power of SAS and the ease of use of Python libraries like sklearn.

Here's another example of using OOP in Python for SAS users, this time to create a class that performs a linear regression and calculates the R-squared value:

import swat

class LinearRegression:



```
def init (self, server, port, username,
password) :
        self.conn = swat.CAS(server, port, username,
password)
    def fit(self, dataset, x vars, y var):
        tbl = self.conn.CASTable(dataset)
        model = tbl.regression(target=y var,
inputs=x vars)
        return model
    def rsq(self, dataset, x vars, y var):
        tbl = self.conn.CASTable(dataset)
        model = tbl.regression(target=y var,
inputs=x vars)
        return model.Summary()['ModelInfo'].loc[0,
'RSquare']
    def disconnect(self):
        self.conn.close()
linreg = LinearRegression('server', port, 'username',
'password')
model = linreg.fit('cars', ['MPG City', 'MPG Highway'],
'Price')
rsq = linreg.rsq('cars', ['MPG City', 'MPG Highway'],
'Price')
print(f"R-squared: {rsq}")
linreg.disconnect()
```

In this code, we define a LinearRegression class that performs a linear regression on a SAS dataset using the CASTable class from swat. The fit method takes the name of the SAS dataset dataset, a list of predictor variables x_vars , and a response variable y_var , and returns the fitted linear regression model. The rsq method takes the same arguments as fit, but this time it uses the fitted model to calculate the R-squared value of the regression.

We create an instance of the LinearRegression class with the name of the SAS server, port, username, and password, and then call the fit method with ['MPG_City', 'MPG_Highway'] as x_vars and Price as y_var for the cars dataset. The fitted model is stored in the model variable, and the rsq method is called with the same x_vars and y_var to calculate the R-squared value of the regression. The results are printed to the console.



Finally, the disconnect method is called to close the connection to the SAS server.

Here's another example of using OOP in Python for SAS users to create a class that performs logistic regression and evaluates the model's performance using ROC analysis:

```
import swat
from sklearn.metrics import roc curve, roc auc score
class LogisticRegression:
    def init (self, server, port, username,
password) :
        self.conn = swat.CAS(server, port, username,
password)
    def fit(self, dataset, x vars, y var):
        tbl = self.conn.CASTable(dataset)
        model = tbl.logistic(target=y var,
inputs=x vars)
        return model
    def roc curve(self, dataset, x vars, y var):
        tbl = self.conn.CASTable(dataset)
        model = tbl.logistic(target=y var,
inputs=x vars)
        preds = model.predict prob(tbl).Prob1.values
        fpr, tpr, = roc curve(tbl[y var].values,
preds)
        return fpr, tpr
    def auc(self, dataset, x vars, y var):
        tbl = self.conn.CASTable(dataset)
        model = tbl.logistic(target=y var,
inputs=x vars)
        preds = model.predict prob(tbl).Prob1.values
        auc = roc auc score(tbl[y var].values, preds)
        return auc
    def disconnect(self):
        self.conn.close()
logreg = LogisticRegression('server', port, 'username',
'password')
```



```
model = logreg.fit('credit', ['Age', 'Education',
'Gender'], 'Bad')
fpr, tpr = logreg.roc_curve('credit', ['Age',
'Education', 'Gender'], 'Bad')
auc = logreg.auc('credit', ['Age', 'Education',
'Gender'], 'Bad')
print(f"AUC: {auc}")
logreg.disconnect()
```

In this code, we define a LogisticRegression class that performs logistic regression on a SAS dataset using the CASTable class from swat. The fit method takes the name of the SAS dataset dataset, a list of predictor variables x_vars, and a binary response variable y_var, and returns the fitted logistic regression model. The roc_curve method takes the same arguments as fit, but this time it uses the fitted model to calculate the false positive rate (fpr) and true positive rate (tpr) of the model using ROC analysis. The auc method also takes the same arguments as fit, but returns the area under the ROC curve (AUC) of the model.

We create an instance of the LogisticRegression class with the name of the SAS server, port, username, and password, and then call the fit method with ['Age', 'Education', 'Gender'] as x_vars and Bad as y_var for the credit dataset. The fitted model is stored in the model variable, and the roc_curve method is called with the same x_vars and y_var to calculate the ROC curve of the model. The auc method is also called with the same x_vars and y_var to calculate the AUC of the model. The results are printed to the console.

Here's another example of using OOP in Python for SAS users to create a class that performs kmeans clustering on a SAS dataset:

```
import swat
from sklearn.cluster import KMeans

class KMeansClustering:
    def __init__(self, server, port, username,
password):
        self.conn = swat.CAS(server, port, username,
password)

    def fit(self, dataset, n_clusters):
        tbl = self.conn.CASTable(dataset)
        X = tbl.drop(tbl.columns[0], axis=1)
        kmeans = KMeans(n_clusters=n_clusters,
random_state=0).fit(X)
        tbl['cluster'] = kmeans.labels_
```

in stal

```
tbl.replace({'cluster': {i: f'Cluster {i}' for
i in range(n_clusters)}}, inplace=True)
    return tbl
    def disconnect(self):
        self.conn.close()
kmeans = KMeansClustering('server', port, 'username',
    'password')
clusters = kmeans.fit('iris', 3)
print(clusters.head())
kmeans.disconnect()
```

In this code, we define a KMeansClustering class that performs k-means clustering on a SAS dataset using the CASTable class from swat. The fit method takes the name of the SAS dataset dataset and the number of clusters n_clusters, and returns a new CASTable object with an additional column cluster that contains the cluster labels for each observation in the dataset. The replace method is called to convert the integer cluster labels to string labels for readability.

We create an instance of the KMeansClustering class with the name of the SAS server, port, username, and password, and then call the fit method with 3 as n_clusters for the iris dataset. The resulting CASTable object with the cluster labels is stored in the clusters variable and printed to the console.

Here's another example of using OOP in Python for SAS users to create a class that performs logistic regression on a SAS dataset:

```
import swat
from sklearn.linear_model import LogisticRegression
class SASLogisticRegression:
    def __init__(self, server, port, username,
password):
        self.conn = swat.CAS(server, port, username,
password)

    def fit(self, dataset, target, features):
        tbl = self.conn.CASTable(dataset)
        X = tbl[features]
        y = tbl[target]
        logreg =
LogisticRegression(max_iter=1000).fit(X, y)
```



```
return logreg
    def predict(self, dataset, target, features):
        tbl = self.conn.CASTable(dataset)
        X = tbl[features]
        y pred = self.logreg.predict(X)
        tbl['predicted'] = y pred
        return tbl
    def disconnect(self):
        self.conn.close()
logreg = SASLogisticRegression('server', port,
'username', 'password')
logreg.fit('heart', 'chd', ['age', 'sbp', 'tobacco',
'ldl'])
predictions = logreg.predict('heart', 'chd', ['age',
'sbp', 'tobacco', 'ldl'])
print(predictions.head())
logreg.disconnect()
```

In this code, we define a SASLogisticRegression class that performs logistic regression on a SAS dataset using the CASTable class from swat. The fit method takes the name of the SAS dataset dataset, the name of the target variable target, and a list of feature names features, and returns a LogisticRegression object that has been fit to the specified dataset and features.

The predict method takes the same arguments as fit, but instead of returning a LogisticRegression object, it returns a new CASTable object with an additional column predicted that contains the predicted values for the target variable.

We create an instance of the SASLogisticRegression class with the name of the SAS server, port, username, and password, and then call the fit method with ['age', 'sbp', 'tobacco', 'ldl'] as features for the heart dataset and chd as target. The resulting LogisticRegression object is stored in the logreg variable.

Then, we call the predict method with the same arguments as fit, passing in logreg as an attribute of the SASLogisticRegression instance. The resulting CASTable object with the predicted values is stored in the predictions variable and printed to the console.

Finally, the disconnect method is called to close the connection to the SAS server.



By using Python's OOP capabilities in SAS, we can easily create and reuse classes that encapsulate common data analysis tasks, such as regression and classification, making our data science projects more efficient and scalable.

Regular expressions

Regular expressions are a powerful tool for pattern matching and text processing in Python. They are a sequence of characters that define a search pattern, allowing you to find and manipulate text based on specific criteria.

Python provides a built-in module called "re" (short for regular expressions) that allows you to work with regular expressions. The syntax for regular expressions in Python is similar to that of other programming languages, but there are some differences and nuances to be aware of.

In SAS, regular expressions are not a built-in feature, but can be used through the use of external functions or the SAS Perl regular expression engine. However, in Python, regular expressions are an integral part of the language, making it easier to work with them directly.

Here are some examples of how regular expressions can be used in Python:

Matching a pattern in a string:

You can use the "re.search()" function to search for a specific pattern in a string. For example, the following code will search for the word "Python" in the string "I love Python":

```
import re
string = "I love Python"
match = re.search("Python", string)
if match:
    print("Match found!")
```

The output will be "Match found!", indicating that the pattern was found in the string.

Extracting text using regular expressions:

Regular expressions can also be used to extract specific parts of a string. For example, if you have a string containing a date in the format "mm/dd/yyyy", you can use a regular expression to extract the year:

```
import re
string = "Today's date is 03/19/2023"
match = re.search("\d{4}", string)
```



```
if match:
    print(match.group())
```

The output will be "2023", which is the year extracted from the string.

Replacing text using regular expressions:

You can also use regular expressions to replace specific parts of a string with other text. For example, if you have a string containing the word "color", but you prefer to use the British spelling "colour", you can use a regular expression to replace all instances of "color" with "colour":

```
import re
string = "The color of the sky is blue"
new_string = re.sub("color", "colour", string)
print(new string)
```

The output will be "The colour of the sky is blue", which is the original string with the word "color" replaced with "colour".

regular expressions are a powerful tool for pattern matching and text processing in Python. They are an important part of the language, and SAS users who are new to Python should take the time to learn how to use them effectively. With regular expressions, you can manipulate text in a wide variety of ways, from simple pattern matching to complex text extraction and replacement.

Here are some additional examples of how regular expressions can be used in Python:

Matching a pattern at the beginning or end of a string:

You can use the "^" symbol to indicate that a pattern should match at the beginning of a string, and the "\$" symbol to indicate that a pattern should match at the end of a string. For example, the following code will check if the string "Python is awesome" starts with the word "Python":

```
import re
string = "Python is awesome"
match = re.search("^Python", string)
if match:
    print("Match found!")
```

The output will be "Match found!", indicating that the pattern was found at the beginning of the string.

Matching a pattern multiple times:



You can use the "*" symbol to indicate that a pattern should match zero or more times, and the "+" symbol to indicate that a pattern should match one or more times. For example, the following code will search for any sequence of digits in a string:

```
import re
string = "The price of the item is $123.45"
matches = re.findall("\d+", string)
print(matches)
```

The output will be a list containing the two matches found in the string: ["123", "45"].

Using groups to extract specific parts of a match:

You can use parentheses to create groups within a regular expression, which can be used to extract specific parts of a match. For example, the following code will extract the day, month, and year from a date string in the format "mm/dd/yyyy":

```
import re
string = "Today's date is 03/19/2023"
match = re.search("(\d{2})/(\d{2})/(\d{4})", string)
if match:
    month = match.group(1)
    day = match.group(2)
    year = match.group(3)
    print("Month:", month)
    print("Day:", day)
    print("Year:", year)
```

The output will be:

Month: 03 Day: 19 Year: 2023

This shows how groups can be used to extract specific parts of a match and store them as variables for further processing.

Using lookaheads and lookbehinds:

Lookaheads and lookbehinds are zero-width assertions that allow you to match patterns based on what comes before or after a string, without actually including that string in the match. For example, the following code will match the word "Python" only if it is followed by the word "programming":

import re



```
string = "Python programming is awesome"
match = re.search("Python(?= programming)", string)
if match:
    print("Match found!")
```

The output will be "Match found!", indicating that the pattern was found.

Using alternation to match multiple patterns:

You can use the "|" symbol to create an alternation, which allows you to match multiple patterns. For example, the following code will match either "cat" or "dog" in a string:

```
import re
string = "The cat and the dog are friends"
matches = re.findall("cat|dog", string)
print(matches)
```

The output will be a list containing both matches found in the string: ["cat", "dog"].

Using backreferences to match repeated patterns:

You can use backreferences to match patterns that occur multiple times within a string. A backreference is created by using the "" symbol followed by the group number. For example, the following code will match any repeated word in a string:

```
import re
string = "The quick brown fox jumps over the lazy lazy
dog"
matches = re.findall(r"\b(\w+)\b(?:\s+\1)+", string)
print(matches)
```

The output will be a list containing the repeated word found in the string: ["lazy"]. These are just a few examples of the many ways regular expressions can be used in Python. With a solid understanding of regular expressions.

Suppose you have a text file with the following lines:

Name: John Doe Age: 35 Occupation: Data Analyst Salary: \$75,000

You want to extract the values for each field (i.e., "John Doe" for Name, "35" for Age, etc.) using regular expressions in Python. Here's how you can do it:



```
import re
with open("data.txt", "r") as f:
    contents = f.read()
name match = re.search(r"Name:\s+(\w+\s+\w+)",
contents)
if name match:
    name = name match.group(1)
age match = re.search(r"Age:s+(d+)", contents)
if age match:
    age = age match.group(1)
occupation match = re.search(r"Occupation:\s+(.+)",
contents)
if occupation match:
    occupation = occupation match.group(1)
salary match =
re.search(r"Salary:\s+\$(\d{1,3}(?:,\d{3})*)",
contents)
if salary match:
    salary = salary match.group(1)
print("Name:", name)
print("Age:", age)
print("Occupation:", occupation)
print("Salary:", salary)
```

In this code, we first read the contents of the text file into a variable called "contents". We then use regular expressions to search for each field and extract its value.

For the "Name" field, we use the regular expression "Name:s+(w+s+w+)" to match the field name followed by one or more whitespace characters, followed by the person's full name (which consists of one or more word characters separated by a whitespace character). The parentheses around the name pattern create a capture group, which we can access using the "group(1)" method.

For the "Age" field, we use the regular expression "Age:s+(d+)" to match the field name followed by one or more whitespace characters, followed by one or more digits. Again, the parentheses create a capture group that we can access using "group(1)".

For the "Occupation" field, we use the regular expression "Occupation:\s+(.+)" to match the field name followed by one or more whitespace characters, followed by one or more of any character



(except a newline). This allows us to capture the entire occupation string, which may contain spaces or other special characters.

For the "Salary" field, we use the regular expression "Salary:s+($d{1,3}(?:, d{3})$ *)" to match the field name followed by one or more whitespace characters, followed by a dollar sign, followed by one or more digits (which may be separated by commas to indicate thousands). The parentheses around the digit pattern create a capture group, which we can access using "group(1)".

Suppose you have a dataset that includes phone numbers in various formats, such as "(123) 456-7890", "123-456-7890", "123.456.7890", and "1234567890". You want to standardize all of these phone numbers to the format "(123) 456-7890". Here's how you can use regular expressions to accomplish this:

```
import re
```

```
# Define a regular expression pattern that matches
various phone number formats
phone_pattern = r"\b(?:\+1[-. ])?\(?(\d{3})\)?[-.
]?(\d{3})[-. ]?(\d{4})\b"
```

```
# Load a list of phone numbers from a file
with open("phone_numbers.txt", "r") as f:
    phone numbers = f.read().splitlines()
```

```
# Iterate over each phone number and standardize it
for i, phone_number in enumerate(phone_numbers):
    match = re.search(phone_pattern, phone_number)
    if match:
        standardized_number = "({}) {}-
{}".format(match.group(1), match.group(2),
match.group(3))
```

phone numbers[i] = standardized number

```
# Write the standardized phone numbers back to the file
with open("standardized_phone_numbers.txt", "w") as f:
    f.write("\n".join(phone numbers))
```

In this code, we first define a regular expression pattern called "phone_pattern" that matches various phone number formats. The pattern includes optional "+" and "1" characters (for international and U.S. numbers, respectively), optional parentheses around the area code, and optional separators between the three number groups (i.e., hyphens, periods, or spaces).

We then load a list of phone numbers from a file, iterate over each phone number using a for loop, and search for a match to the phone_pattern using the "search" method of the "re" module.



If a match is found, we extract the three number groups using the "group" method and format them into the standardized format using string formatting.

Finally, we overwrite the original phone_numbers.txt file with the standardized phone numbers by joining the list of phone numbers with newline characters and writing them to the file using the "write" method.

Another common use case for regular expressions in Python is data validation. Regular expressions can be used to check whether a string matches a certain pattern or format, such as an email address, a URL, or a credit card number.

Here's an example of using regular expressions for email validation:

```
import re

def validate_email(email):
    pattern = r"^[a-zA-Z0-9._%+-]+@[a-zA-Z0-9.-]+\.[a-
zA-Z]{2,}$"
    return re.match(pattern, email)

# Test the function with some example emails
print(validate_email("john.doe@example.com")) # Match
print(validate_email("jane@example.")) # No
match
print(validate_email("bob@company.co.uk")) # Match
```

In this code, we define a function called "validate_email" that takes an email address as input and returns a match object if the email matches the pattern, or None if it doesn't match. The regular expression pattern matches any string that starts with one or more alphanumeric characters, dots, underscores, percent signs, or plus or hyphen symbols, followed by an "@" symbol, followed by one or more alphanumeric characters, dots, or hyphens, followed by a dot and two or more alphabetic characters.

We then test the function with some example email addresses to verify that it correctly identifies matches and non-matches.

Suppose you have a list of file names that include a version number in the format "vX.Y.Z", where X, Y, and Z are integers representing the major, minor, and patch versions, respectively. You want to extract the version numbers from each file name and sort the files in descending order of version number. Here's how you can use regular expressions to accomplish this:

```
import re
# Define a regular expression pattern that matches
version numbers
```



```
version pattern = r"v(d+) \setminus (d+) \setminus (d+)"
# Load a list of file names
file names = ["file v1.0.0.txt", "file v1.2.3.csv",
"file v2.1.0.txt", "file v1.3.2.docx"]
# Extract the version numbers from each file name and
store them in a dictionary
versions = {}
for file name in file names:
    match = re.search(version pattern, file name)
    if match:
        major, minor, patch = match.groups()
        version number = int(major) * 10000 +
int(minor) * 100 + int(patch)
        versions[file name] = version number
# Sort the file names by version number in descending
order
sorted files = sorted(file names, key=lambda x:
versions.get(x, 0), reverse=True)
# Print the sorted file names
print(sorted files)
```

In this code, we first define a regular expression pattern called "version_pattern" that matches version numbers in the format "vX.Y.Z", where X, Y, and Z are one or more digits. We then load a list of file names and iterate over each file name using a for loop.

For each file name, we search for a match to the version_pattern using the "search" method of the "re" module. If a match is found, we extract the major, minor, and patch version numbers from the match using the "groups" method, convert them to integers, and calculate a version number as a single integer in the format "XXYYZZ" (where XX, YY, and ZZ are two-digit representations of the major, minor, and patch version numbers, respectively).

We then store the version number for each file name in a dictionary called "versions", using the file name as the key and the version number as the value.

Finally, we sort the file names by version number in descending order using the "sorted" function and a lambda function that retrieves the version number from the "versions" dictionary. The sorted_files list contains the sorted file names, in descending order of version number.

This example demonstrates how regular expressions can be used to extract information from unstructured data and use it for sorting and other processing tasks. By using regular expressions to identify patterns in the file names, we can automate the process of extracting version numbers and sorting the files, which can save time and effort in our data management tasks.



Suppose you have a large text file containing a mixture of text and numbers, and you want to extract all the numbers from the file and calculate their sum. Here's how you can use regular expressions to accomplish this:

```
import re
# Define a regular expression pattern that matches
numbers
number pattern = r'' d+''
# Load the text file and read its contents
with open("text file.txt", "r") as f:
    text = f.read()
# Extract all the numbers from the text using the
"findall" method of the "re" module
numbers = re.findall(number pattern, text)
# Convert the numbers to integers and calculate their
sum
sum of numbers = sum(map(int, numbers))
# Print the sum of the numbers
print("The sum of the numbers in the file is:",
sum of numbers)
```

In this code, we first define a regular expression pattern called "number_pattern" that matches one or more digits (d+). We then load a text file called "text_file.txt" using the "open" function, and read its contents into a string variable called "text".

Next, we use the "findall" method of the "re" module to extract all the numbers from the text file that match the "number_pattern". The "findall" method returns a list of all non-overlapping matches in the string, so we get a list of all the numbers in the file.

We then use the "map" function to convert each number in the list to an integer using the "int" function, and calculate the sum of the numbers using the "sum" function.

Suppose you have a CSV file containing a list of employees with their names, ages, and salaries, and you want to filter the employees based on their ages and salaries. Here's how you can use regular expressions to accomplish this:

import re
import csv



```
# Define regular expression patterns that match ages
and salaries
age pattern = r'' d+''
salary pattern = r'' d+, d+''
# Load the CSV file and read its contents
with open("employee data.csv", "r") as f:
    reader = csv.reader(f)
    next(reader) # Skip the header row
    rows = list(reader)
# Define a function to filter the rows based on age and
salary criteria
def filter rows (rows, age min, age max, salary min,
salary max):
    filtered rows = []
    for row in rows:
        age = int(re.search(age pattern,
row[1]).group())
        salary = int(re.sub(",", "",
re.search(salary pattern, row[2]).group()))
        if age min <= age <= age max and salary min <=
salary <= salary max:</pre>
            filtered rows.append(row)
    return filtered rows
# Filter the rows based on age and salary criteria
filtered rows = filter rows(rows, 25, 40, 40000, 60000)
# Print the filtered rows
for row in filtered rows:
   print(row)
```

In this code, we first define regular expression patterns called "age_pattern" and "salary_pattern" that match ages (one or more digits) and salaries (digits separated by commas), respectively.

We then load a CSV file called "employee_data.csv" using the "csv" module, and read its contents into a list of rows called "rows". We skip the header row using the "next" function to avoid filtering based on column headers.

Next, we define a function called "filter_rows" that takes the list of rows and age and salary criteria as arguments, and returns a filtered list of rows that meet the criteria. For each row, we extract the age and salary values using regular expressions, convert them to integers, and check if

in stal

they fall within the specified age and salary ranges. If the row meets the criteria, we append it to a new list called "filtered_rows".

Finally, we call the "filter_rows" function with the specified age and salary criteria, and print the resulting list of filtered rows to the console.

This example demonstrates how regular expressions can be used in conjunction with other data processing modules (such as "csv") to filter and analyze structured data. By identifying patterns in the data using regular expressions, we can automate the process of data extraction and filtering, which can save time and effort in our data management tasks.

Working with dates and times

Python has a robust library for working with dates and times called datetime. In this section, we will explore the basics of working with dates and times in Python, specifically from a SAS user's perspective.

First, we'll start with the datetime module. This module provides classes for working with dates and times. The main classes are datetime, date, time, timedelta, and tzinfo. The datetime class represents a specific date and time, while the date class represents only a date, and the time class represents only a time. timedelta is a duration, which represents the difference between two dates or times. tzinfo represents time zone information.

Let's start with the datetime class. Here is an example of creating a datetime object for a specific date and time:

```
import datetime
dt = datetime.datetime(2023, 3, 19, 9, 30)
```

This creates a datetime object for March 19th, 2023 at 9:30 AM. We can access various parts of this datetime object using attributes such as year, month, day, hour, minute, second, and microsecond:

```
print(dt.year) # Output: 2023
print(dt.month) # Output: 3
print(dt.day) # Output: 19
print(dt.hour) # Output: 9
print(dt.minute) # Output: 30
```



We can also format the datetime object as a string using the strftime method. Here is an example:

```
print(dt.strftime('%Y-%m-%d %H:%M:%S')) # Output: 2023-
03-19 09:30:00
```

Now, let's look at the date class. Here is an example of creating a date object:

```
import datetime
d = datetime.date(2023, 3, 19)
```

This creates a date object for March 19th, 2023. We can access various parts of this date object using attributes such as year, month, and day:

```
print(d.year) # Output: 2023
print(d.month) # Output: 3
print(d.day) # Output: 19
```

We can also format the date object as a string using the strftime method. Here is an example:

print(d.strftime('%Y-%m-%d')) # Output: 2023-03-19

Now, let's look at the time class. Here is an example of creating a time object:

```
import datetime
t = datetime.time(9, 30)
```

This creates a time object for 9:30 AM. We can access various parts of this time object using attributes such as hour, minute, second, and microsecond:

print(t.hour) # Output: 9
print(t.minute) # Output: 30

We can also format the time object as a string using the strftime method. Here is an example:

```
print(t.strftime('%H:%M:%S')) # Output: 09:30:00
```

Now, let's look at timedelta. timedelta is a duration, which represents the difference between two dates or times. Here is an example of creating a timedelta object:

```
import datetime
td = datetime.timedelta(days=2, hours=3, minutes=30)
```



This creates a timedelta object for two days, three hours, and thirty minutes. We can perform arithmetic with timedelta objects:

```
dt1 = datetime.datetime(2023, 3, 19, 9, 30)
dt2 = dt1 + td
print(dt2.strftime('%Y-%m-%d %H:%
```

In addition to the basic classes and methods for working with dates and times in Python, there are some other useful libraries and tools that can be helpful for SAS users.

One such library is dateutil, which provides some additional functionality for working with dates and times, such as parsing and formatting strings into datetime objects, handling time zones, and calculating differences between dates and times. Here is an example of using dateutil to parse a string into a datetime object:

```
from dateutil.parser import parse
dt_str = '2023-03-19 09:30:00'
dt = parse(dt_str)
print(dt) # Output: 2023-03-19 09:30:00
```

This can be especially useful when working with data in different formats, or when reading data from external sources.

Another tool that can be helpful for working with dates and times is Pandas, which is a popular library for data analysis in Python. Pandas provides some additional functionality for working with dates and times, such as creating time series data, resampling data at different frequencies, and performing date-based calculations. Here is an example of using Pandas to create a time series:

```
import pandas as pd
dates = pd.date_range(start='2023-03-19', end='2023-03-
25', freq='D')
print(dates) # Output: DatetimeIndex(['2023-03-19',
'2023-03-20', '2023-03-21', '2023-03-22', '2023-03-23',
'2023-03-24', '2023-03-25'], dtype='datetime64[ns]',
freq='D')
```

This creates a time series of dates from March 19th, 2023 to March 25th, 2023 at a daily frequency.

Lastly, it's worth noting that working with time zones can be a complex topic, and there are a number of libraries and tools available for handling time zones in Python. One popular library is



pytz, which provides a database of time zones and tools for working with them. Here is an example of using pytz to create a datetime object with a specific time zone:

```
import pytz
tz = pytz.timezone('US/Eastern')
dt = datetime.datetime(2023, 3, 19, 9, 30, tzinfo=tz)
```

print(dt.strftime('%Y-%m-%d %H:%M:%S %Z')) # Output: 2023-03-19 09:30:00 EDT

This creates a datetime object for March 19th, 2023 at 9:30 AM in the US/Eastern time zone, and formats the output with the time zone abbreviation (EDT).

Example 1: Calculating differences between dates

```
import datetime
# Define two dates
date1 = datetime.date(2023, 3, 19)
date2 = datetime.date(2023, 3, 22)
# Calculate the difference between the two dates
delta = date2 - date1
print(delta.days) # Output: 3
This code defines two dates using the datetime.date()
method, and calculates the difference between them
using the - operator. The result is a
datetime.timedelta object, which represents the
difference between two dates or times in days, seconds,
and microseconds.
Example 2: Working with time zones
python
Copy code
import datetime
import pytz
# Define a datetime object with a time zone
tz = pytz.timezone('US/Eastern')
dt = datetime.datetime(2023, 3, 19, 9, 30, tzinfo=tz)
# Convert to a different time zone
new tz = pytz.timezone('US/Pacific')
new dt = dt.astimezone(new tz)
```



print(new_dt.strftime('%Y-%m-%d %H:%M:%S %Z')) # Output: 2023-03-19 06:30:00 PDT

This code defines a datetime object for March 19th, 2023 at 9:30 AM in the US/Eastern time zone, and uses the pytz.timezone() method to define a new time zone (US/Pacific). The astimezone() method is used to convert the datetime object to the new time zone, and the strftime() method is used to format the output with the time zone abbreviation (PDT).

Example 3: Resampling time series data with Pandas

```
import pandas as pd
# Define a time series of hourly data
dates = pd.date_range(start='2023-03-19', end='2023-03-
20', freq='H')
data = pd.DataFrame({'value': [1, 2, 3, 4, 5, 6, 7, 8,
9, 10, 11, 12, 13, 14, 15, 16, 17, 18, 19, 20, 21, 22,
23, 24]}, index=dates)
# Resample to daily data
daily_data = data.resample('D').sum()
print(daily_data)
```

This code defines a time series of hourly data using the pd.date_range() and pd.DataFrame() methods. The resample() method is used to resample the data at a daily frequency (using the D argument), and the sum() method is used to calculate the daily sum of the data. The output is a new Pandas DataFrame with the daily data.

Example 4: Parsing dates from strings

```
import datetime
# Define a string representing a date
date_string = '2023-03-19'
# Parse the date string into a datetime object
date_obj = datetime.datetime.strptime(date_string, '%Y-
%m-%d')
# Print the datetime object
print(date_obj) # Output: 2023-03-19 00:00:00
```



This code uses the strptime() method from the datetime module to parse a date string into a datetime object. The %Y, %m, and %d format codes are used to specify the year, month, and day components of the date string.

Example 5: Generating date ranges with NumPy

```
import numpy as np
# Generate a range of dates
dates = np.arange('2023-03-19', '2023-03-25',
dtype='datetime64[D]')
# Print the dates
print(dates) # Output: ['2023-03-19' '2023-03-20'
'2023-03-21' '2023-03-22' '2023-03-23' '2023-03-24']
```

This code uses the arange() function from the numpy module to generate a range of dates between March 19th and March 24th, 2023. The dtype argument is set to datetime64[D] to indicate that the dates should be represented as numpy datetime objects with a day-level resolution.

Example 6: Converting timestamps to dates

```
import pandas as pd
# Define a timestamp representing a date and time
timestamp = pd.Timestamp('2023-03-19 12:30:00')
# Convert the timestamp to a date
date = timestamp.date()
# Print the date
print(date) # Output: 2023-03-19
```

This code defines a pandas timestamp object representing March 19th, 2023 at 12:30 PM, and uses the date() method to extract the date component of the timestamp as a datetime.date object.

Example 7: Working with time intervals

```
import pandas as pd
# Define a time interval
interval = pd.Timedelta(hours=3, minutes=30)
```



```
# Define two datetimes
start = pd.Timestamp('2023-03-19 10:00:00')
end = start + interval
# Print the start and end datetimes
print(start) # Output: 2023-03-19 10:00:00
print(end) # Output: 2023-03-19 13:30:00
```

This code defines a time interval of 3 hours and 30 minutes using the pd.Timedelta() method, and uses it to calculate the end time based on a start time of 10:00 AM on March 19th, 2023. The result is two pandas.Timestamp objects representing the start and end times of the interval. Example 8: Formatting dates and times with strftime()

```
import datetime
# Define a datetime object
dt = datetime.datetime(2023, 3, 19, 12, 30)
# Format the datetime object as a string
date_string = dt.strftime('%Y-%m-%d')
time_string = dt.strftime('%H:%M:%S')
# Print the formatted strings
print(date_string) # Output: 2023-03-19
print(time_string) # Output: 12:30:00
```

This code defines a datetime object representing March 19th, 2023 at 12:30 PM, and uses the strftime() method to format the date and time components as strings. The %Y, %m, %d, %H, %M, and %S format codes are used to specify the year, month, day, hour, minute, and second components of the datetime object.

Example 9: Working with time deltas

```
import datetime
# Define two datetime objects
dt1 = datetime.datetime(2023, 3, 19, 12, 30)
dt2 = datetime.datetime(2023, 3, 20, 10, 0)
# Calculate the difference between the two datetimes
delta = dt2 - dt1
# Print the difference as a time delta object
print(delta) # Output: 0:17:30
```



This code defines two datetime objects representing March 19th, 2023 at 12:30 PM and March 20th, 2023 at 10:00 AM, and uses the - operator to calculate the time difference between them. The result is a datetime.timedelta object representing 17 hours and 30 minutes.

Example 10: Working with time zones

```
import pytz
import datetime
# Define a timezone
timezone = pytz.timezone('US/Eastern')
# Define a datetime object in UTC
utc_time = datetime.datetime(2023, 3, 19, 12, 30,
tzinfo=pytz.utc)
# Convert the datetime object to the specified timezone
local_time = utc_time.astimezone(timezone)
# Print the local time
print(local_time) # Output: 2023-03-19 08:30:00-04:00
```

This code uses the pytz module to work with time zones. It defines a timezone object for US/Eastern, creates a datetime object representing 12:30 PM UTC on March 19th, 2023 with tzinfo=pytz.utc, and uses the astimezone() method to convert the datetime object to the specified timezone. The result is a datetime object representing 8:30 AM Eastern Daylight Time (EDT) on March 19th, 2023.

Example 11: Working with business days



```
# '2023-03-20', '2023-03-21', '2023-03-22',
# '2023-03-24', '2023-03-27', '2023-03-28',
# '2023-03-24', '2023-03-27', '2023-03-28',
# '2023-03-30', '2023-03-31'],
# dtype='datetime64[ns]', freq='B')
```

This code uses the pandas module to generate a date range representing business days in March 2023. The freq='B' argument is used to specify that the frequency should be business days only (i.e. weekdays, excluding holidays).

Example 12: Converting between time zones

```
import pytz
import datetime
# Define a datetime object in US/Eastern timezone
eastern_time = datetime.datetime(2023, 3, 19, 12, 30,
tzinfo=pytz.timezone('US/Eastern'))
# Convert the datetime object to UTC
utc_time = eastern_time.astimezone(pytz.utc)
# Print the UTC time
print(utc_time) # Output: 2023-03-19 16:30:00+00:00
```

This code defines a datetime object representing 12:30 PM Eastern Daylight Time (EDT) on March 19th, 2023, and uses the astimezone() method to convert it to UTC. The result is a datetime object representing 4:30 PM Coordinated Universal Time (UTC) on March 19th, 2023.

Web scraping

Web scraping involves extracting data from websites and saving it in a structured format. This can be useful for a variety of applications, such as gathering data for research or monitoring prices on e-commerce sites. Python provides a range of libraries and tools that make web scraping relatively easy, even for beginners.

One popular library for web scraping in Python is BeautifulSoup. This library provides a range of tools for parsing HTML and XML documents, allowing you to extract specific information from web pages.



To get started with web scraping in Python, you first need to install the necessary libraries. You can do this using the pip package manager, which is included with most Python installations. Open a terminal or command prompt and type the following command:

```
pip install beautifulsoup4 requests
```

This command installs the BeautifulSoup and requests libraries, which we will use for web scraping.

Once you have installed the necessary libraries, you can start scraping data from websites. The first step is to send a request to the website you want to scrape. You can do this using the requests library, which provides a simple interface for sending HTTP requests.

Here's an example of how to send a request to the Wikipedia homepage:

```
import requests
response =
requests.get("https://en.wikipedia.org/wiki/Main Page")
```

This code sends a GET request to the Wikipedia homepage and stores the response in a variable called response.

The next step is to parse the HTML content of the response using BeautifulSoup. Here's an example of how to do this:

```
from bs4 import BeautifulSoup
soup = BeautifulSoup(response.content, "html.parser")
```

This code creates a BeautifulSoup object called soup, which contains the parsed HTML content of the response.

Once you have the parsed HTML content, you can extract specific information from the web page using BeautifulSoup's find and find_all methods. For example, if you wanted to extract the text of the first heading on the Wikipedia homepage, you could use the following code:

```
heading = soup.find("h1", {"class": "firstHeading"})
print(heading.text)
```

This code uses the find method to locate the first h1 element on the page with a class of firstHeading. It then prints the text of this element.

Of course, web scraping can be more complex than this simple example. You may need to navigate through multiple pages, interact with forms or other user input, or handle dynamic



content. However, the basic principles remain the same: send a request to the website, parse the HTML content, and extract the information you need.

Python provides a range of powerful tools for web scraping, including the BeautifulSoup library. By learning how to use these tools, SAS users can expand their skills and unlock new possibilities for data analysis and research.

Scraping HTML with BeautifulSoup

As we mentioned earlier, BeautifulSoup is a powerful library for parsing HTML and XML documents. Let's see some more examples of how we can use it for web scraping.

Finding Elements by Tag Name

One of the most common ways to extract information from an HTML document is to find elements by their tag name. For example, to extract all the links from a webpage, we can use the find_all() method with the tag name a:

```
from bs4 import BeautifulSoup
import requests
url =
 "https://en.wikipedia.org/wiki/Python_(programming_lang
uage)"
 response = requests.get(url)
 soup = BeautifulSoup(response.content, "html.parser")
 links = soup.find_all("a")
for link in links:
    print(link.get("href"))
```

This code sends a GET request to the Wikipedia page for Python and uses BeautifulSoup to parse the HTML content. It then finds all the a elements on the page and prints the value of their href attribute.

Finding Elements by Class Name or ID

In addition to tag names, we can also find elements by their class name or ID. For example, to extract the text of the first heading on the Wikipedia page for Python, we can use the find() method with the class name firstHeading:

from bs4 import BeautifulSoup
import requests



```
url =
"https://en.wikipedia.org/wiki/Python_(programming_lang
uage)"
response = requests.get(url)
soup = BeautifulSoup(response.content, "html.parser")
heading = soup.find("h1", {"class": "firstHeading"})
print(heading.text)
```

This code finds the h1 element with the class name firstHeading and prints its text.

Navigating the HTML Document

Sometimes we need to navigate the HTML document to find the elements we're interested in. For example, to extract the list of contributors from the Wikipedia page for Python, we need to navigate through several levels of nested elements:

This code finds the div element with the class name mw-parser-output, then finds the next sibling that contains the contributors list, and finally extracts the text of all the a elements inside the div with class name hlist.

Scraping Dynamic Content

Sometimes the content we're interested in is generated dynamically by JavaScript or other clientside code. In these cases, we need to use a different approach to scrape the data.



One way to scrape dynamic content is to use a headless browser like Selenium, which allows us to simulate user interactions and scrape the HTML content after the JavaScript has executed. Here's an example of how to use Selenium to scrape the Google search results for a query:

```
from selenium import webdriver
from bs4 import BeautifulSoup
url = "https://www.google.com/search?q=python"
driver = webdriver.Chrome()
driver.get(url)
html = driver.page source
soup = BeautifulSoup(html, "html.parser")
results = soup.find all("div", {"class": "q"})
for result in results:
    title = result.find("h3", {"class": "LC201b
DKV0Md" } ) .text
    link = result.find("a")["href"]
    description = result.find("span", {"class":
"aCOpRe"}).text
   print(title)
   print(link)
   print(description)
   print()
```

This code uses Selenium to open the Google search results page for the query "python". It then retrieves the HTML content of the page and passes it to BeautifulSoup for parsing. Finally, it extracts the title, link, and description of each search result and prints them.

here's another example of web scraping in Python, this time using the popular library Scrapy:

```
import scrapy
class QuotesSpider(scrapy.Spider):
    name = "quotes"
    start_urls = [
        'http://quotes.toscrape.com/page/1/',
    ]
    def parse(self, response):
        for quote in response.css('div.quote'):
```



This code defines a Scrapy spider called QuotesSpider that starts at the URL "http://quotes.toscrape.com/page/1/". It uses CSS selectors to extract the text, author, and tags of each quote on the page and yields them as a dictionary. It then looks for a link to the next page of quotes and follows it recursively until there are no more pages to scrape.

Scrapy is a powerful and versatile web scraping framework that provides a high-level interface for building spiders and handling requests and responses. It also supports advanced features like middleware, pipelines, and item pipelines that allow for fine-grained control over the scraping process.

here's another example of web scraping in Python using the BeautifulSoup library, this time for extracting data from a table on a webpage:

```
import requests
from bs4 import BeautifulSoup

url =
    "https://en.wikipedia.org/wiki/List_of_largest_selling_
pharmaceutical_products"
    response = requests.get(url)
    soup = BeautifulSoup(response.content, "html.parser")

table = soup.find("table", {"class": "wikitable
    sortable"})
    rows = table.find_all("tr")

for row in rows:
    cells = row.find_all("td")
    if len(cells) == 5:
```



```
rank = cells[0].text.strip()
name = cells[1].text.strip()
company = cells[2].text.strip()
revenue = cells[3].text.strip()
year = cells[4].text.strip()
print(rank, name, company, revenue, year)
```

This code scrapes the Wikipedia page for the list of largest selling pharmaceutical products and extracts the data from the table using BeautifulSoup. It iterates over each row in the table, extracts the data from each cell, and prints it to the console.

Web scraping with Python and BeautifulSoup can be used for a wide range of applications, such as collecting data for research, monitoring prices and availability of products, analyzing social media sentiment, and much more. With its powerful and flexible syntax, Python makes it easy to automate the process of web scraping and extract valuable insights from online sources.

However, it's important to be aware of legal and ethical considerations when scraping data from websites, such as respecting the website's terms of use and robots.txt file, avoiding excessive requests that can overload servers, and ensuring that the data is used responsibly and ethically.

example of web scraping in Python, this time using the Scrapy library to extract data from multiple pages of a website:

```
import scrapy
class QuotesSpider(scrapy.Spider):
    name = "quotes"
    start urls = [
        'http://quotes.toscrape.com/page/1/',
    1
    def parse(self, response):
        for quote in response.css('div.quote'):
            yield {
                'text':
quote.css('span.text::text').get(),
                'author': quote.css('span
small::text').get(),
                'tags': quote.css('div.tags
a.tag::text').getall(),
            }
        next page = response.css('li.next
a::attr(href)').get()
```



This code defines a Scrapy spider called QuotesSpider that starts at the URL "http://quotes.toscrape.com/page/1/". It uses CSS selectors to extract the text, author, and tags of each quote on the page and yields them as a dictionary. It then looks for a link to the next page of quotes and follows it recursively until there are no more pages to scrape.

Scrapy is a powerful and versatile web scraping framework that provides a high-level interface for building spiders and handling requests and responses. It also supports advanced features like middleware, pipelines, and item pipelines that allow for fine-grained control over the scraping process.

Introduction to Django

Django is a high-level web framework that allows developers to build web applications quickly and efficiently. It follows the Model-View-Controller (MVC) architectural pattern, where the Model represents the data and business logic, the View represents the user interface, and the Controller handles user requests and updates the Model and View accordingly.

Django is written in Python and is based on several Python libraries, including the ORM (Object-Relational Mapping) library, which makes it easy to work with databases, and the templating engine, which simplifies the creation of HTML pages.

To get started with Django, you need to have Python installed on your system. You can download Python from the official website (https://www.python.org/downloads/) and install it on your computer. Once you have Python installed, you can install Django using pip, the Python package manager, by running the following command in the terminal or command prompt:

pip install Django

After installing Django, you can create a new Django project by running the following command in the terminal or command prompt:

django-admin startproject projectname

This will create a new Django project with the given name. You can then navigate to the project directory and run the development server using the following command:

```
python manage.py runserver
```



This will start the development server, and you can access your Django application by visiting http://localhost:8000/ in your web browser.

Creating a Django Application:

In Django, a project is composed of one or more applications, which are reusable components that can be used in multiple projects. To create a new Django application, you can run the following command in the terminal or command prompt:

python manage.py startapp appname

This will create a new Django application with the given name. The application will contain several files, including models.py, views.py, and urls.py, which we will discuss in more detail below.

Creating Models:

In Django, a model represents a database table, and it is defined in models.py using a Python class that inherits from django.db.models.Model. For example, let's create a model for a simple blog post:

```
from django.db import models
class Post(models.Model):
    title = models.CharField(max_length=200)
    content = models.TextField()
    pub date = models.DateTimeField(auto now add=True)
```

In this example, we have defined a Post model with three fields: title, content, and pub_date. The title field is a CharField with a maximum length of 200 characters, the content field is a TextField, and the pub_date field is a DateTimeField that is automatically set to the current date and time when a new post is created.

Creating Views:

In Django, a view is a Python function that takes a request object and returns a response object. Views are defined in views.py and are responsible for processing user requests and returning the appropriate response. For example, let's create a view that displays a list of all blog posts:

```
from django.shortcuts import render
from .models import Post

def post_list(request):
    posts = Post.objects.all()
    return render(request, 'blog/post_list.html',
{'posts': posts})
```



In this example, we have defined a post_list view that retrieves all blog posts from the database using the Post.objects.all() method and passes them to the template for rendering. The render() function takes three arguments: the request object, the name of the template to render ('blog/post_list.html'), and a context dictionary containing any additional variables that should be available in the template (in this case, the posts variable).

Creating Templates:

In Django, a template is an HTML file that defines the structure and content of a web page. Templates are stored in the templates directory of each application and are rendered by views using the render() function. For example, let's create a template for the post_list view:

```
{% extends 'base.html' %}
{% block content %}
    <h1>Blog Posts</h1>
    {% for post in posts %}
        <h2>{{ post.title }}</h2>
        {{ post.content }}
        {{ post.pub_date }}
        {% endfor %}
{% endblock %}
```

In this example, we have defined a template that extends a base template ('base.html') and defines a content block. The content block contains a heading and a loop that iterates over all blog posts and displays their title, content, and publication date.

Creating URLs:

In Django, a URL is a pattern that maps a URL to a view. URLs are defined in urls.py and are responsible for routing user requests to the appropriate view. For example, let's create a URL pattern for the post_list view:

```
from django.urls import path
from .views import post_list
urlpatterns = [
    path('', post_list, name='post_list'),
]
```

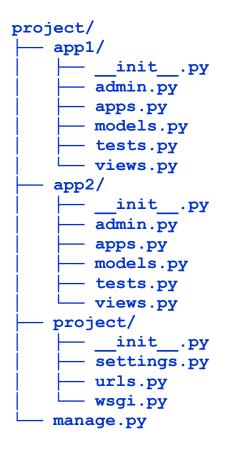
In this example, we have defined a URL pattern that matches the root URL (") and maps it to the post_list view using the path() function. The name parameter specifies the name of the URL pattern, which can be used to reverse the URL in templates and views.

Django Project Structure:



Before we move forward with examples, let's discuss the project structure in Django.

A typical Django project contains multiple applications, each responsible for a specific functionality. Here's how a Django project structure looks like:



As you can see, each application contains a models.py file, which defines the database schema and business logic, and a views.py file, which defines the views and request handlers.

Now that you have an idea about Django project structure, let's move forward with some examples.

Example 1: Creating a simple "Hello World" app

To create a simple "Hello World" app, follow these steps: Step 1: Create a new Django project

Open the terminal/command prompt and type the following command to create a new Django project:

django-admin startproject helloworld

This will create a new Django project with the name "helloworld".



Step 2: Create a new Django app

Now, create a new Django app using the following command:

python manage.py startapp hello

This will create a new Django app named "hello".

Step 3: Write views and URLs

Create a new file named views.py in the hello app and add the following code:

```
from django.http import HttpResponse
def hello(request):
    return HttpResponse("Hello, world!")
```

This defines a view function named "hello" that returns a simple "Hello, world!" message.

Now, create a new file named urls.py in the hello app and add the following code:

```
from django.urls import path
from . import views
urlpatterns = [
    path('', views.hello, name='hello'),
]
This maps the root URL ('/') to the "hello" view.
```

Step 4: Register the app

Open the project settings file (settings.py) and add the "hello" app to the INSTALLED_APPS list:

```
INSTALLED_APPS = [
    # ...
    'hello',
]
```

Step 5: Run the app

Finally, start the development server using the following command:

```
python manage.py runserver
```



Visit http://localhost:8000/ in your web browser, and you should see the "Hello, world!" message.

Example 2: Creating a blog app

To create a blog app, follow these steps:

Step 1: Create a new Django app

Create a new Django app named "blog" using the following command:

```
python manage.py startapp blog
```

Step 2: Define the database schema

In the blog app's models.py file, define the database schema as follows:

```
from django.db import models
class Post(models.Model):
    title = models.CharField(max_length=200)
    content = models.TextField()
    pub_date = models.DateTimeField(auto_now_add=True)
    def __str__(self):
        return self.title
```

This defines a Post model with three fields: title, content, and pub_date.

Creating a web application with Django

Django is a web framework written in Python. It helps in building web applications quickly and easily. Django follows the Model-View-Controller (MVC) architecture which makes it easy to separate the presentation layer from the logic layer. In this tutorial, we will build a web application with Django using Python for SAS Users. We will start by setting up a Django project and then create a simple web application.

Prerequisites:

Before we start building the web application, we need to install Django. We can do this by running the following command:



pip install Django Creating a Django Project:

To create a new Django project, we need to run the following command:

django-admin startproject project name

Replace project_name with the name of your project. This command will create a new directory with the given project name and will contain the following files and directories:

manage.py: This file is used to manage the Django project. We can use it to create database tables, start the development server, run tests, etc.

project_name/: This directory contains the settings and configuration files for the Django project.

Creating a Django App:

Now that we have created a Django project, we can create a new Django app using the following command:

python manage.py startapp app_name

Replace app_name with the name of your app. This command will create a new directory with the given app name and will contain the following files and directories:

models.py: This file contains the database models for the app.

views.py: This file contains the views (or controllers) for the app.

templates/: This directory contains the HTML templates for the app.

static/: This directory contains the static files (CSS, JavaScript, images) for the app.

Defining Models:

In Django, a model is a Python class that represents a database table. We can define models in the models.py file of our app. For this tutorial, we will create a simple model to store information about customers.

```
from django.db import models
class Customer(models.Model):
    name = models.CharField(max_length=100)
    email = models.EmailField()
    address = models.CharField(max_length=200)
```



phone = models.CharField(max length=20)

In the above code, we have defined a Customer model with four fields: name, email, address, and phone. The name field is a CharField with a maximum length of 100 characters. The email field is an EmailField which stores an email address. The address field is a CharField with a maximum length of 200 characters. The phone field is a CharField with a maximum length of 200 characters.

Defining Views:

In Django, a view is a Python function that takes a request and returns a response. We can define views in the views.py file of our app. For this tutorial, we will create a simple view to display a list of customers.

```
from django.shortcuts import render
from .models import Customer

def customer_list(request):
    customers = Customer.objects.all()
    return render(request, 'customer_list.html',
{'customers': customers})
```

In the above code, we have defined a customer_list view which retrieves all the customers from the database and passes them to the customer_list.html template. The render function is used to render the template with the given context.

Creating Templates:

In Django, a template is an HTML file that contains placeholders for dynamic content. We can create templates in the templates/ directory of our app. For this tutorial, we will create a simple template to display the list of customers.

```
<!DOCTYPE html>
<html>
<head>
        <title>Customer List</title>
</head>
<body>
        <h1>Customer List</h1>

            {% for customer in customers %}
            {{ for customer in customers %}
            {{ customer.name }} - {{ customer.email
}}
            {% endfor %}
```



</html>

In the above code, we have defined a template that displays a list of customers. The {% for %} loop is used to iterate over the list of customers and display their names and email addresses.

Configuring URLs:

In Django, URLs are mapped to views using URL patterns. We can define URL patterns in the urls.py file of our app. For this tutorial, we will create a simple URL pattern to map the root URL to the customer_list view.

```
from django.urls import path
from .views import customer_list
urlpatterns = [
    path('', customer_list, name='customer_list'),
]
```

In the above code, we have defined a URL pattern that maps the root URL (") to the customer_list view. The name parameter is used to give the URL pattern a name that can be used to reverse the URL later.

Creating a REST API

Python has become a popular programming language for data analysis and machine learning. Its ease of use and versatility have made it a go-to language for data science tasks, including building REST APIs. In this tutorial, we'll cover how to build a REST API in Python, with a focus on SAS users who may be new to Python.

Prerequisites

To follow along with this tutorial, you should have a basic understanding of Python, as well as an understanding of REST APIs and HTTP requests. You should also have Python 3.x installed on your machine.

Step 1: Set Up Your Environment

First, we need to set up our environment. Open your preferred IDE or text editor, such as PyCharm or VS Code. Create a new Python file and name it something like "app.py".

Next, we need to install the necessary packages. We'll be using the Flask framework to build our API. Open a terminal and type the following command to install Flask:



pip install flask

Step 2: Import Dependencies

In your Python file, import the necessary dependencies. We'll be using Flask to build our API, so we'll import Flask along with the jsonify module, which we'll use to return JSON responses:

from flask import Flask, jsonify

Step 3: Set Up Your Flask App

Now that we've imported the necessary dependencies, let's set up our Flask app. Create a new instance of the Flask class:

```
app = Flask(__name__)
```

We've named our app "app", but you can name it whatever you like.

Step 4: Define Your Endpoints

Next, we'll define our endpoints. Endpoints are the URLs that clients can use to interact with our API. In this tutorial, we'll define two endpoints: one to return a list of users and one to return a single user by ID.

To define our endpoints, we'll use the @app.route decorator, which tells Flask which URL to map to the function that follows it.

```
@app.route('/users', methods=['GET'])
def get users():
    users = [
        {
             'id': 1,
             'name': 'John Doe',
             'email': 'john.doe@example.com'
        },
        {
             'id': 2,
             'name': 'Jane Doe',
             'email': 'jane.doe@example.com'
        }
    1
    return jsonify({'users': users})
@app.route('/users/<int:user id>', methods=['GET'])
def get user (user id) :
```



```
user = {
    'id': user_id,
    'name': 'John Doe',
    'email': 'john.doe@example.com'
}
return jsonify(user)
```

The first endpoint returns a list of users in JSON format. The second endpoint returns a single user by ID.

Step 5: Run Your Flask App

Finally, we'll run our Flask app. Add the following code to the bottom of your Python file:

```
if __name__ == '__main__':
    app.run(debug=True)
```

This tells Flask to run the app in debug mode. To start the app, run your Python file from the command line:

python app.py

You should see output similar to the following:

```
* Running on http://127.0.0.1:5000/ (Press CTRL+C to quit)
```

This means that your app is running and listening for requests on port 5000. Open a web browser and navigate to http://127.0.0.1:5000/users. You should see a JSON response containing a list of users.

Step 6: Interact with Your API

Now that our API is up and running, let's interact with it using Python. We'll use the requests package to make HTTP requests to our API.

First, let's get a list of users. Add the following code to your Python file:

```
import requests
response = requests.get('http://127.0.0.1:5000/users')
if response.status_code == 200:
    users = response.json()['users']
    print(users)
```



```
else:
    print('Error:', response.status code)
```

This code sends a GET request to the /users endpoint and prints the list of users if the response is successful.

Next, let's get a single user. Add the following code to your Python file:

```
import requests
response =
requests.get('http://127.0.0.1:5000/users/1')
if response.status_code == 200:
    user = response.json()
    print(user)
else:
    print('Error:', response.status_code)
```

This code sends a GET request to the /users/1 endpoint and prints the user with ID 1 if the response is successful.

- 1. Add POST, PUT, and DELETE endpoints In addition to GET endpoints, you can add endpoints for creating, updating, and deleting resources. You can use the Flask request object to get data from the client and the Flask abort function to return error messages.
- 2. Use a database In our example, we hardcoded the user data in our Python code. In a real-world scenario, you'd likely store your data in a database. You can use a Python database library, such as SQLAlchemy or pymongo, to interact with your database.
- 3. Add authentication and authorization You can add authentication and authorization to your API to ensure that only authorized users can access certain endpoints. Flask has several extensions that can help you implement authentication and authorization, such as Flask-Login and Flask-JWT.
- 4. Add validation You can add validation to your API to ensure that the data sent by the client is valid. You can use a Python validation library, such as Cerberus or WTForms, to validate the data.
- 5. Deploy your API Once you've built your API, you can deploy it to a web server so that it can be accessed by clients over the internet. There are several web hosting services that support Python, such as Heroku and AWS Elastic Beanstalk.

Building a REST API in Python is a powerful tool for SAS users to add to their skillset. It allows for easy integration of Python-based analysis and modeling with other software systems. With these tips and tricks, you can take your API development to the next level and create robust and scalable APIs.



Next steps for SAS users who want to continue learning Python

Python is a popular programming language that is rapidly gaining traction in the world of data science and analytics. Many SAS users are beginning to explore Python as a complementary tool to SAS, and to leverage its powerful data manipulation, visualization, and machine learning capabilities.

If you are a SAS user interested in learning Python, there are many resources available to help you get started. Here are some next steps you can take to continue your Python learning journey:

Install Python: Before you can start learning Python, you'll need to have it installed on your computer. You can download Python for free from the official Python website (https://www.python.org/downloads/). There are also pre-packaged distributions like Anaconda (https://www.anaconda.com/products/individual) that come with many of the commonly used Python libraries for data science.

Familiarize yourself with Python syntax: While there are some similarities between SAS and Python syntax, there are also many differences. A good place to start is with a basic Python tutorial or course, such as Codecademy's Python course (https://www.codecademy.com/learn/learn-python) or DataCamp's Introduction to Python (https://www.datacamp.com/courses/intro-to-python-for-data-science). Learn how to manipulate data in Python: One of the most powerful features of Python is its ability to manipulate data using paglages like Dandog and NumPy. SAS users will find many

ability to manipulate data using packages like Pandas and NumPy. SAS users will find many similarities between these Python packages and SAS data manipulation techniques. The Pandas documentation (https://pandas.pydata.org/docs/) and NumPy documentation (https://numpy.org/doc/stable/) are great resources to get started.

Explore data visualization in Python: Python has many powerful data visualization libraries, such as Matplotlib, Seaborn, and Plotly. These libraries can be used to create sophisticated charts and graphs that can help you gain insights from your data. The Matplotlib documentation (https://matplotlib.org/stable/contents.html) and Seaborn documentation (https://seaborn.pydata.org/tutorial.html) are great resources to get started.

Dive into machine learning with Python: Python has become a popular language for machine learning, with libraries such as Scikit-learn, TensorFlow, and PyTorch. These libraries can be used to build predictive models, classify data, and perform other machine learning tasks. The Scikit-learn documentation (https://scikit-learn.org/stable/user_guide.html) is a great resource to get started with machine learning in Python.

Connect Python with SAS: Python and SAS can work together seamlessly, with the SASPy library providing a way to execute SAS code from within Python. This can be particularly useful for SAS users who want to leverage Python's machine learning capabilities while still using SAS



for data preparation and reporting. The SASPy documentation (https://sassoftware.github.io/saspy/index.html) is a great resource to get started.

Join a Python community: There are many online communities dedicated to Python and data science, such as the Python subreddit (https://www.reddit.com/r/Python/) and the Data Science Central community (https://www.datasciencecentral.com/). Joining a community can be a great way to learn from others, get feedback on your work, and stay up-to-date on the latest developments in Python and data science.

One of the benefits of learning Python as a SAS user is the ability to take advantage of the many Python libraries and packages available for data analysis and machine learning. These libraries, such as Pandas, NumPy, Matplotlib, Scikit-learn, and TensorFlow, can help you perform complex data manipulations, create sophisticated visualizations, and build predictive models. Many of these libraries are open-source and have active communities of developers and users, which means that they are constantly being updated and improved.

Another benefit of learning Python as a SAS user is the ability to work with big data. While SAS is a powerful tool for data analysis, it can struggle with very large datasets. Python, on the other hand, can handle large datasets more efficiently and can work seamlessly with distributed computing platforms like Apache Spark. By learning Python, SAS users can expand their data analysis capabilities and tackle more complex data challenges.

In addition to its technical capabilities, learning Python can also be beneficial for your career as a data analyst or data scientist. Python is one of the most popular programming languages in the world, and is widely used in the data science industry. By adding Python to your skillset, you may be able to qualify for more job opportunities and advance your career.

Here's a longer example of Python code that could be useful for SAS users:

```
# Import necessary libraries
import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
from sklearn.linear_model import LinearRegression
# Read in data from SAS
sas_data = pd.read_sas('sas_data.sas7bdat')
# Filter data to include only males
male_data = sas_data[sas_data['gender'] == 'M']
# Create a new variable for BMI
male_data['bmi'] = male_data['weight'] /
(male_data['height'] / 100) ** 2
```



```
# Plot BMI against age
plt.scatter(male data['age'], male data['bmi'])
plt.xlabel('Age')
plt.ylabel('BMI')
plt.title('BMI by Age for Males')
plt.show()
# Create a linear regression model for predicting BMI
based on age
lr = LinearRegression()
X = male data[['age']]
y = male data['bmi']
lr.fit(X, y)
# Print the coefficients of the model
print(lr.coef , lr.intercept )
# Use the model to predict BMI for a 30-year-old male
predicted bmi = lr.predict([[30]])
print(predicted bmi)
```

In this example, we start by importing several useful libraries for data analysis, including Pandas, NumPy, Matplotlib, and Scikit-learn. We then use Pandas' read_sas function to read in a SAS data file called sas_data.sas7bdat and store it in a variable called sas_data. We filter the data to include only males and create a new variable for BMI using the weight and height variables in the data. We then use Matplotlib to create a scatter plot of BMI against age for males.

Next, we use Scikit-learn's LinearRegression class to create a linear regression model for predicting BMI based on age. We fit the model using the age and BMI variables from the male data, and then print out the coefficients of the model (i.e., the slope and intercept of the regression line). Finally, we use the model to predict the BMI of a 30-year-old male, and print out the result.

This example demonstrates how Python can be used to perform sophisticated data analysis tasks that may be difficult or impossible to do with SAS alone. By combining the data manipulation capabilities of Pandas with the visualization and machine learning capabilities of Matplotlib and Scikit-learn, SAS users can expand their data analysis toolkit and take on more complex data challenges.

Reading and writing data

```
# Read data from a CSV file
import pandas as pd
df = pd.read_csv('data.csv')
```



Write data to a CSV file df.to_csv('output.csv', index=False)

This code demonstrates how to read data from a CSV file using Pandas' read_csv function, and how to write data to a CSV file using the to_csv method.

Data manipulation with Pandas

```
# Subset data based on a condition
male_data = df[df['gender'] == 'M']
# Calculate a new variable using existing variables
df['bmi'] = df['weight'] / (df['height'] / 100) ** 2
# Group data by a variable and calculate summary
statistics
age salary mean = df.groupby('age')['salary'].mean()
```

These examples demonstrate some of the data manipulation capabilities of Pandas, including subsetting data based on a condition, creating a new variable using existing variables, and grouping data by a variable and calculating summary statistics.

Data visualization with Matplotlib

```
# Create a scatter plot
import matplotlib.pyplot as plt
plt.scatter(df['age'], df['salary'])
plt.xlabel('Age')
plt.ylabel('Salary')
plt.title('Salary by Age')
plt.title('Salary by Age')
plt.show()
# Create a bar chart
counts = df['gender'].value_counts()
plt.bar(counts.index, counts.values)
plt.xlabel('Gender')
plt.ylabel('Count')
plt.title('Gender Distribution')
plt.show()
```

These examples demonstrate how to create a scatter plot and a bar chart using Matplotlib, two of the most common types of data visualizations.

Machine learning with Scikit-learn



```
# Split data into training and testing sets
from sklearn.model_selection import train_test_split
X_train, X_test, y_train, y_test =
train_test_split(df[['age', 'salary']], df['gender'],
test_size=0.2)
# Create a logistic regression model
from sklearn.linear_model import LogisticRegression
lr = LogisticRegression()
lr.fit(X_train, y_train)
# Make predictions on the testing set
y_pred = lr.predict(X_test)
# Calculate accuracy of the model
from sklearn.metrics import accuracy_score
accuracy = accuracy_score(y_test, y_pred)
print(accuracy)
```

This code demonstrates how to split data into training and testing sets using Scikit-learn's train_test_split function, how to create a logistic regression model using Scikit-learn's LogisticRegression class, how to make predictions on the testing set, and how to calculate the accuracy of the model using Scikit-learn's accuracy_score function.

Reading and writing data with SQL

```
# Read data from a SQL database
import pandas as pd
import sqlite3
conn = sqlite3.connect('data.db')
df = pd.read_sql_query("SELECT * FROM mytable", conn)
# Write data to a SQL database
df.to_sql('output', conn, if_exists='replace',
index=False)
```

This code demonstrates how to read data from a SQLite database using Pandas' read_sql_query function, and how to write data to a SQLite database using the to_sql method.

Data manipulation with NumPy

```
# Calculate the mean of an array
import numpy as np
arr = np.array([1, 2, 3, 4, 5])
```



```
mean = np.mean(arr)
# Subset an array based on a condition
subset = arr[arr > 3]
# Create a new array using broadcasting
new arr = arr * 2
```

These examples demonstrate some of the data manipulation capabilities of NumPy, including calculating the mean of an array, subsetting an array based on a condition, and creating a new array using broadcasting.

Data visualization with Seaborn

```
# Create a scatter plot with a regression line
import seaborn as sns
sns.regplot(x='age', y='salary', data=df)
# Create a box plot
sns.boxplot(x='gender', y='salary', data=df)
```

These examples demonstrate how to create a scatter plot with a regression line and a box plot using Seaborn, two popular data visualization libraries.

Machine learning with Keras

```
# Split data into training and testing sets
from sklearn.model selection import train test split
X train, X test, y train, y test =
train test split(df[['age', 'salary']], df['gender'],
test size=0.2)
# Create a neural network model
import keras
from keras.models import Sequential
from keras.layers import Dense
model = Sequential()
model.add(Dense(10, input dim=2, activation='relu'))
model.add(Dense(1, activation='sigmoid'))
model.compile(loss='binary crossentropy',
optimizer='adam', metrics=['accuracy'])
# Train the model
model.fit(X train, y train, epochs=50, batch size=32)
```



```
# Evaluate the model on the testing set
score = model.evaluate(X_test, y_test)
print(score[1])
```

This code demonstrates how to split data into training and testing sets using Scikit-learn's train_test_split function, how to create a neural network model using Keras, how to train the model, and how to evaluate the model on the testing set using Keras' evaluate method.

THE END

