# Part-11

# Python Fundamentals: Mastering Concepts and Techniques

## The __name__ Variable

The __name__ variable is a special variable in Python that holds the name of the current module. If the file is executed directly (i.e., not imported from another module), __name__ will have the value '__main__'. If the file is imported from another module, __name__ will hold the name of the imported module.

- The __name__ variable is used to determine the entry point of a program
- When a Python script is run directly, __name__ is set to '__main__'
- When a Python script is imported as a module, __name__ is set to the name of the module

## NumPy Arrays

NumPy arrays are a fundamental data structure in Python's scientific computing ecosystem. They offer several advantages over standard Python lists.

- NumPy arrays are grids of values, all of the same data type
- NumPy arrays are indexed by tuples of non-negative integers
- The number of dimensions in a NumPy array is called its rank
- The shape of a NumPy array is a tuple of integers giving the size of the array along each dimension

## Matrices vs. Arrays

While related, matrices and arrays are distinct data structures in Python:

- Matrices are a two-dimensional representation of data from linear algebra
- Matrices come with a powerful set of mathematical operations
- Arrays are a more general sequence of objects of similar data type
- An array within another array forms a matrix-like structure

# Finding the Indices of the N Maximum Values in a NumPy Array

To find the indices of the N maximum values in a NumPy array:

1. Import the NumPy library as `np`
2. Create the NumPy array
3. Sort the array in descending order
4. Use negative indexing to get the indices of the N maximum values

```
import numpy as np
arr = np.array([1, 2, 3, 4, 5])
indices_of_max = np.argsort(arr)[-n:]
```

# Combining Training and Test Sets

When working with machine learning datasets, it's common to have separate training and test sets. To combine these sets:

1. Create the training and test sets as NumPy arrays or arrays of arrays
2. Use `np.concatenate()` to combine the sets horizontally (column-wise)
3. Use `np.vstack()` to combine the sets vertically (row-wise)

```
train_set = np.array([1, 2, 3])
test_set = np.array([[0, 1, 2], [1, 2, 3]])
combined_set = np.vstack((train_set, test_set))
```

# Importing a Decision Tree Classifier

To import a Decision Tree Classifier from the scikit-learn library:

```
from sklearn.tree import DecisionTreeClassifier
```

The `DecisionTreeClassifier` class is part of the `sklearn.tree` module.

# Accessing Data from a Google Spreadsheet

To access data stored in a Google Spreadsheet:

1. Upload the data to a Google Spreadsheet and share it publicly

2. Get the public link to the spreadsheet

3. Use the `pandas` library to read the data from the link

```
import pandas as pd
import io
import requests
url = 'https://docs.google.com/spreadsheets/d/[spreadsheet_id]/export?format=csv&id=[spreadsheet_id]'
data = pd.read_csv(io.StringIO(requests.get(url).content))
```

# Data Frame Views vs. Copies

When working with Pandas data frames, it's important to understand the difference between views and copies:

- `DF_name` and `DF_location` are both copies of the original data frame
- Pandas data frames are often returned as views, not copies, unless explicitly specified

# Handling Encoding Errors when Reading CSV Files

When encountering a `UnicodeEncodeError`
while reading a CSV file with Pandas, the solution is to specify the correct encoding:

```
pd.read_csv('temp.csv', encoding='utf-8')
```

The `encoding='utf-8'` parameter tells Pandas to use the UTF-8 encoding when reading the file.

# Setting Line Width in Plots

To set the line width in a Matplotlib plot, use the `linewidth` or `lw` parameter:

```
plt.plot(x, y, linewidth=3) or plt.plot(x, y, lw=3)
```

# Resetting the Index of a Pandas Data Frame

To reset the index of a Pandas data frame to a given list:

```
df = df.reindex(new_index_list)
```

The `reindex()` method allows you to change the index of a data frame to a new set of labels.

# Copying Objects in Python

Python provides two ways to copy objects:

1. `copy.copy()` for shallow copying
2. `copy.deepcopy()`for deep copying

Shallow copying creates a new object that references the same elements as the original, while deep copying creates a new object with completely independent elements.

# Range vs. XRange

The main difference between `range()` and `xrange()` (in Python 2) is:

- `range()`returns a Python list object
- `xrange()` returns an `xrange` object, which generates values as needed, rather than creating the entire list upfront
- This makes `xrange()`more memory-efficient for large ranges.

# Checking if a Pandas Data Frame is Empty

To check if a Pandas data frame is empty, use the `empty` attribute:

```
if df.empty:
  print("The data frame is empty.")
else:
  print("The data frame is not empty.")
```

# Sorting a NumPy Array by the N-1 Column

To sort a NumPy array by the N-1 column:

```
import numpy as np
X = np.array([[1, 2, 3], [0, 5, 2], [2, 3, 4]])
sorted_indices = X[:, -2].argsort()
sorted_X = X[sorted_indices]
```

The `argsort()` method returns the indices that would sort the array, which can then be used to reorder the array.

# Creating Series from Lists, NumPy Arrays, and Dictionaries

To create a Pandas Series from various data structures:

```
import pandas as pd
import numpy as np
```

```
my_list = ['a', 'b', 'c', 'd', 'e', 'f', 'g']
my_array = np.arange(26)
my_dict = dict(zip(my_list, my_array))
series_from_list = pd.Series(my_list)
series_from_array = pd.Series(my_array)
series_from_dict = pd.Series(my_dict)
```

# Finding Items Not Common to Both Series

To find the items that are not common to both Pandas Series A and B:

```
import pandas as pd
series_1 = pd.Series([1, 2, 3, 4, 5, 4, 5, 6, 7, 8])
series_2 = pd.Series([4, 5, 6, 7, 8, 9, 10])
unique_items = series_1[~series_1.isin(series_2)]
```

The `isin()` method is used to check which elements in `series_1` are also present in `series_2`, and the ~operator is used to negate the result, giving us the unique items.

# Keeping the Top 2 Most Frequent Values in a Series

To keep the top 2 most frequent values in a Pandas Series and replace everything else as "Other":

```
import pandas as pd
import numpy as np
series = pd.Series(np.random.randint(1, 6, 12))
top_2_values = series.value_counts().head(2).indexseries[~series.isin(top_2_values)] = 'Other'
```

The `value_counts()`method is used to get the frequency of each value, and the `head(2)` method is used to get the top 2 most frequent values.

The `isin()`method is then used to identify the values that are not in the top 2, and these are replaced with the "Other" label.

# Finding Positions of Numbers Divisible by 3 in a Series

To find the positions of numbers that are multiples of 3 in a Pandas Series:

```
import pandas as pd
import numpy as np
series = pd.Series(np.random.randint(1, 11, 10))
positions_of_multiples_of_3 = series[series % 3 == 0].index
```

The key steps are:

1. Create a Pandas Series

2. Use the modulo operator % to check which values are divisible by 3

3. Get the index of the resulting boolean Series to find the positions of the multiples of 3.

# Computing Euclidean Distance Between Two Series

To compute the Euclidean distance between two Pandas Series:

```
import pandas as pd
import numpy as np
p = pd.Series([1, 2, 3, 4, 5])
q = pd.Series([2, 3, 4, 5, 6])
# Solution1: Using Pandaseuclidean_distance_1 = ((p - q) ** 2).sum() ** 0.5
# Solution2: Using NumPyeuclidean_distance_2 = np.linalg.norm(p - q)
```

Both solutions compute the Euclidean distance between the two series, with the second solution using the more conci se `np.linalg.norm()` function.

# Reversing the Rows of a Pandas Data Frame

To reverse the rows of a Pandas data frame:

```
import numpy as np
import pandas as pd
df = pd.DataFrame(np.arange(25).reshape(5, 5))
reversed_df = df.iloc[::-1, :]
```

The key steps are:

1. Create a sample Pandas data frame

2. Use the `iloc` indexer with a step of `-1` to reverse the rows

3. Keep all the columns by using the `:` notation

# Overfitting when Splitting Data into Train and Test Sets

It is possible to overfit a model when splitting data into train and test sets, especially if you:

- Retune the model parameters after seeing the test set performance
- Train new models with different parameters until you get the desired result on the test set

This is a common beginner mistake, as it can lead to overly optimistic performance on the test set that does not generalize well to new, unseen data.

# Seaborn: A Data Visualization Library Built on Matplotlib and Pandas

Seaborn is a popular data visualization library in Python that is built on top of Matplotlib and Pandas. It provides a high-level interface for drawing informative and attractive statistical graphics.

Seaborn is often used in conjunction with Pandas and Matplotlib to create complex, publication-quality visualizations with minimal code.

# Table of Contents

| Section | Description |
|---|---|
| The __name__ Variable | Explains the purpose and behavior of the __name__ variable in Python |
| NumPy Arrays | Discusses the key characteristics and features of NumPy arrays |
| Matrices vs. Arrays | Differentiates between matrices and arrays in Python |
| Finding the Indices of the N Maximum Values in a NumPy Array | Demonstrates how to find the indices of the N maximum values in a NumPy array |
| Combining Training and Test Sets | Explains how to combine training and test sets using NumPy |
| Importing a Decision Tree Classifier | Shows how to import a Decision Tree Classifier from the scikit-learn library |
| Accessing Data from a Google Spreadsheet | Describes the process of accessing data stored in a Google Spreadsheet using Pandas |
| Data Frame Views vs. Copies | Differentiates between views and copies of Pandas data frames |

| | |
|---|---|
| Handling Encoding Errors when Reading CSV Files | Provides a solution for handling `UnicodeEncodeError` when reading CSV files with Pandas |
| Setting Line Width in Plots | Demonstrates how to set the line width in Matplotlib plots |
| Resetting the Index of a Pandas Data Frame | Explains how to reset the index of a Pandas data frame to a given list |
| Copying Objects in Python | Discusses the differences between shallow and deep copying in Python |
| Range vs. XRange | Compares the `range()` and `xrange()` functions in Python |
| Checking if a Pandas Data Frame is Empty | Shows how to check if a Pandas data frame is empty |
| Sorting a NumPy Array by the N-1 Column | Demonstrates how to sort a NumPy array by the N-1 column |
| Creating Series from Lists, NumPy Arrays, and Dictionaries | Explains how to create Pandas Series from various data structures |
| Finding Items Not Common to Both Series | Describes how to find the items that are not common to two Pandas Series |
| Keeping the Top 2 Most Frequent Values in a Series | Illustrates how to keep the top 2 most frequent values in a Pandas Series |
| Finding Positions of Numbers Divisible by 3 in a Series | Shows how to find the positions of numbers that are multiples of 3 in a Pandas Series |
| Computing Euclidean Distance Between Two Series | Provides two solutions for computing the Euclidean distance between two Pandas Series |
| Reversing the Rows of a Pandas Data Frame | Demonstrates how to reverse the rows of a Pandas data frame |
| Overfitting when Splitting Data into Train and Test Sets | Discusses the risk of overfitting when splitting data into train and test sets |
| Seaborn: A Data Visualization Library Built on Matplotlib and Pandas | Introduces Seaborn, a data visualization library built on top of Matplotlib and Pandas |