

SIT719 Security and Privacy Issues in Analytics

Pass Task 2.1: Basic scripting with python

Overview

Python is an amazingly versatile programming language and extremely popular among the data science people. This powerful tool will give you access to a wide variety of data science libraries which will help you to develop your script easily. By the end of week 2, you will be familiar with basic python scripting. Please see the weekly resources for some basic operations.

If you are new to python scripting, you might follow the below references:

- Python for Data Analysis: Data Wrangling with Pandas, NumPy, and IPython by Wes McKinney, O'Reilly Media, Inc.

Because of the evolving nature of the open-source tools like Python and its libraries, it is always wise to look for the updated learning material from the python library website tutorials, user guides and manuals. For example, the user guide of the pandas data frame can be obtained from the below link:

https://pandas.pydata.org/pandas-docs/stable/user_guide/index.html

Similarly, numpy can be learned based on the material presented in the following links:

<https://docs.scipy.org/doc/numpy/user/basics.html>

<https://docs.scipy.org/doc/numpy/user/quickstart.html>

This is a Pass task, so you **MUST** complete the task and submit the evidence of your work to Ontrack.

Submit the following files to Ontrack:

- A screenshot of the output you obtained by executing the python program (in Section 1)
- Some reflections on what you got out of this experience of learning fundamental concepts of python scripting (see Section 2)

Section 1

Instructions: In this task, you will be asked to perform some basic python operations using pandas and numpy libraries. Please write the code, execute and take a screenshot of the results of the completed outputs.

Step 1. Import the pandas and numpy libraries

Answer1: (This one has been done for you)

```
In [140]: import pandas as pd
...: import numpy as np
```

Step 2. Import the popular 'iris' dataset from the below address. And then check the header of the dataset.

<https://archive.ics.uci.edu/ml/machine-learning-databases/iris/iris.data>

Answer2: (This one has also been done for you)

```
In [141]: url = 'https://archive.ics.uci.edu/ml/machine-learning-databases/iris/iris.data'

In [142]: iris = pd.read_csv(url)

In [143]: iris.head()
Out[143]:
5.1 3.5 1.4 0.2 Iris-setosa
0 4.9 3.0 1.4 0.2 Iris-setosa
1 4.7 3.2 1.3 0.2 Iris-setosa
2 4.6 3.1 1.5 0.2 Iris-setosa
3 5.0 3.6 1.4 0.2 Iris-setosa
4 5.4 3.9 1.7 0.4 Iris-setosa
```

Step 3. You can see that the column headers are missing in the above case. Therefore this step is related to the creation of column heads for the dataset. [Write code to create 5 column heads.](#)
[Next write a code to display or show the headers.](#)

```
1. sepal_length
2. sepal_width
3. petal_length
4. petal_width
5. class
```

Answer3: (write your code)

Step 4. [Write a code to check if there are any missing values in the dataframe?](#)

Answer4: (write your code)

Hints: there is no missing values but check it thorough the code

Step 5. [Write a code to set the values of the rows 10 to 29 of the column 'petal_length' to NaN.](#)

Answer5: (write your code)

Step 6. [Now again, check if there is any missing values \(NaN\) in the dataframe? Count, how many missing values.](#)

Answer6: (write your code)

Hints: this time you will have missing values.

Step 7. [Substitute the NaN values to 10.0](#)

Answer7: (write your code)

Section 2

Create a short document that summarizes some key concepts and usages of important python libraries like 'pandas', 'numpy', 'matplotlib' you have learned in this week (doc size max 200 words). In your response you need to address the following queries.

- [Importance of python libraries for data analysis](#)
- [Some common functionalities and usages related to dataframe manipulation \(for example, NaN check, slicing dataset using iloc\). Just show 2 examples.](#)
- [Demonstration of a sample visualization example using matplotlib library.](#)

OUTPUT TO SUBMIT

Take screenshots of section 1 as a proof that you have completed. Then combine those screenshots and answers of section 2 in a **single PDF document** and upload using OnTrack system.