SIT719 Task 2.1 Sanket Thakur

December 5, 2020

```
import pandas as pd
     import numpy as np
[4]: iris_data = pd.read_csv("https://archive.ics.uci.edu/ml/
      →machine-learning-databases/iris/iris.data", header=None)
[5]: iris_data
[5]:
            0
                 1
                       2
                            3
                                             4
                    1.4 0.2
     0
          5.1
               3.5
                                  Iris-setosa
     1
          4.9
               3.0
                    1.4 0.2
                                  Iris-setosa
     2
          4.7
               3.2
                    1.3 0.2
                                  Iris-setosa
          4.6
     3
               3.1
                    1.5
                          0.2
                                  Iris-setosa
                    1.4 0.2
          5.0 3.6
                                  Iris-setosa
     145
               3.0
                   5.2
                          2.3
                               Iris-virginica
          6.7
                               Iris-virginica
     146
         6.3
               2.5
                    5.0
                          1.9
     147
          6.5
                          2.0
                               Iris-virginica
               3.0
                    5.2
     148
         6.2
               3.4
                    5.4
                          2.3
                               Iris-virginica
               3.0
                    5.1
                          1.8
                               Iris-virginica
     [150 rows x 5 columns]
[8]: iris_data.
      -columns=['sepal_length', 'sepal_width', 'petal_length', 'petal_width', 'class']
[9]: iris data
[9]:
          sepal_length
                        sepal_width petal_length petal_width
                                                                            class
                   5.1
                                 3.5
                                                1.4
                                                             0.2
                                                                      Iris-setosa
                   4.9
                                 3.0
                                                             0.2
     1
                                                1.4
                                                                      Iris-setosa
     2
                   4.7
                                 3.2
                                                1.3
                                                             0.2
                                                                      Iris-setosa
                                 3.1
                                                             0.2
                                                                      Iris-setosa
     3
                   4.6
                                                1.5
     4
                   5.0
                                 3.6
                                                1.4
                                                             0.2
                                                                      Iris-setosa
     145
                   6.7
                                 3.0
                                                5.2
                                                             2.3
                                                                   Iris-virginica
     146
                   6.3
                                 2.5
                                                5.0
                                                             1.9
                                                                   Iris-virginica
```

	147	6.5	3.0	5.2	2.0	Iris-virginica								
	148	6.2	3.4	5.4	2.3	_								
	140					•								
	149	5.9	3.0	5.1	1.8	Iris-virginica								
	[150 rows x 5 columns]													
[10]:	<pre>iris_data.isnull().tail()</pre>													
[10]:	se	epal_length	sepal_width	petal_length	petal_width	class								
	145	False	False	False	False									
	146	False	False	False	False	False								
	147	False	False	False	False	False								
	148	False	False	False	False	False								
	149	False	False	False	False	False								
[11]:	iris_da	ata.loc[9:29	, 'petal_leng	th']= pd.np.Na	ıN									
	<pre><ipython-input-11-9f54047bc343>:1: FutureWarning: The pandas.np module is</ipython-input-11-9f54047bc343></pre>													
	deprecated and will be removed from pandas in a future version. Import numpy													
	directly instead													
	iris_data.loc[9:29, 'petal_length'] = pd.np.NaN													
	1115_(uata.100[9.2	e, petal_ler	igui]- pa.np.i	van									
[13]:	_	ata.isnull()	es, petal_ler	igui 1- puinpii	van									
	iris_da	nta.isnull()	•			class								
[13]: [13]:	iris_da	ata.isnull()	sepal_width	petal_length	petal_width									
	iris_da	ata.isnull() epal_length False	sepal_width False	petal_length False	petal_width False	False								
	iris_da	ata.isnull() epal_length False False	sepal_width False False	petal_length False False	petal_width False False	False False								
	iris_da	epal_length False False False	sepal_width False False False	petal_length False False False	petal_width False False False	False False False								
	iris_da se 0 1 2 3	epal_length False False False False	sepal_width False False False False	petal_length False False False False	petal_width False False False False	False False False False								
	iris_da se 0 1 2 3 4	epal_length False False False False False False False	sepal_width False False False False False	petal_length False False False	petal_width False False False False False	False False False False								
	iris_da se 0 1 2 3 4	epal_length False False False False False False False	sepal_width False False False False False	petal_length False False False False False False	petal_width False False False False False False	False False False False False								
	iris_da se 0 1 2 3 4 145	epal_length False False False False False False False False False	sepal_width False False False False False False	petal_length False False False False False False	petal_width False False False False False False False	False False False False False False								
	iris_da se 0 1 2 3 4 145 146	epal_length False	sepal_width False False False False False False False	petal_length False False False False False False False False	petal_width False False False False False False False False	False False False False False False False								
	iris_da se 0 1 2 3 4 145 146 147	epal_length False	sepal_width False	petal_length False	petal_width False	False False False False False False False False False								
	iris_da se 0 1 2 3 4 145 146 147 148	eta.isnull() epal_length False	sepal_width False	petal_length False	petal_width False False	False								
	iris_da se 0 1 2 3 4 145 146 147	epal_length False	sepal_width False	petal_length False	petal_width False	False False False False False False False False False								
	iris_da se 0 1 2 3 4 145 146 147 148 149	eta.isnull() epal_length False	sepal_width False	petal_length False	petal_width False False	False								
	iris_da se 0 1 2 3 4 145 146 147 148 149	epal_length False	sepal_width False	petal_length False False	petal_width False False	False								
[13]:	iris_da se 0 1 2 3 4 145 146 147 148 149 [150 ro	epal_length False	sepal_width False	petal_length False False	petal_width False False	False								

2

1.4

0.2

class

Iris-setosa

sepal_length sepal_width petal_length petal_width

3.5

[15]: iris_data.fillna(10.0)

5.1

[15]:

0

1	4.9	3.0	1.4		0.2	Iris-setosa
2	4.7	3.2	1.3		0.2	Iris-setosa
3	4.6	3.1	1.5		0.2	Iris-setosa
4	5.0	3.6	1.4		0.2	Iris-setosa
• •	•••	•••	•••	•••		•••
145	6.7	3.0	5.2		2.3	Iris-virginica
146	6.3	2.5	5.0		1.9	<pre>Iris-virginica</pre>
147	6.5	3.0	5.2		2.0	Iris-virginica
148	6.2	3.4	5.4		2.3	Iris-virginica
149	5.9	3.0	5.1		1.8	Iris-virginica

[150 rows x 5 columns]

[]:

SIT719 Security and Privacy Issues in Analytics

Pass Task 2.1: Basic scripting with python

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.

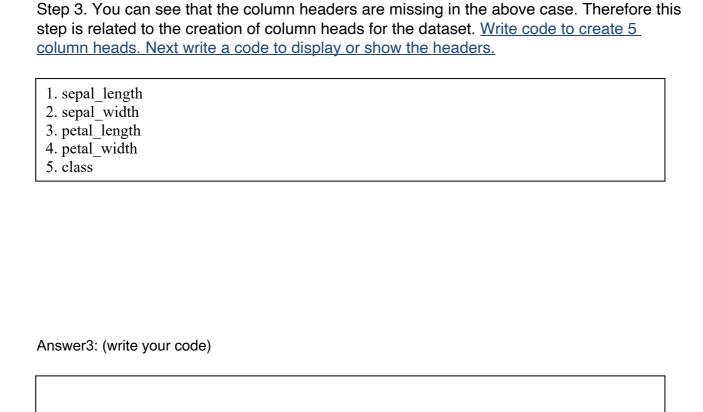
https://archive.ics.uci.edu/ml/machine-learning-databases/iris/iris.data

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)

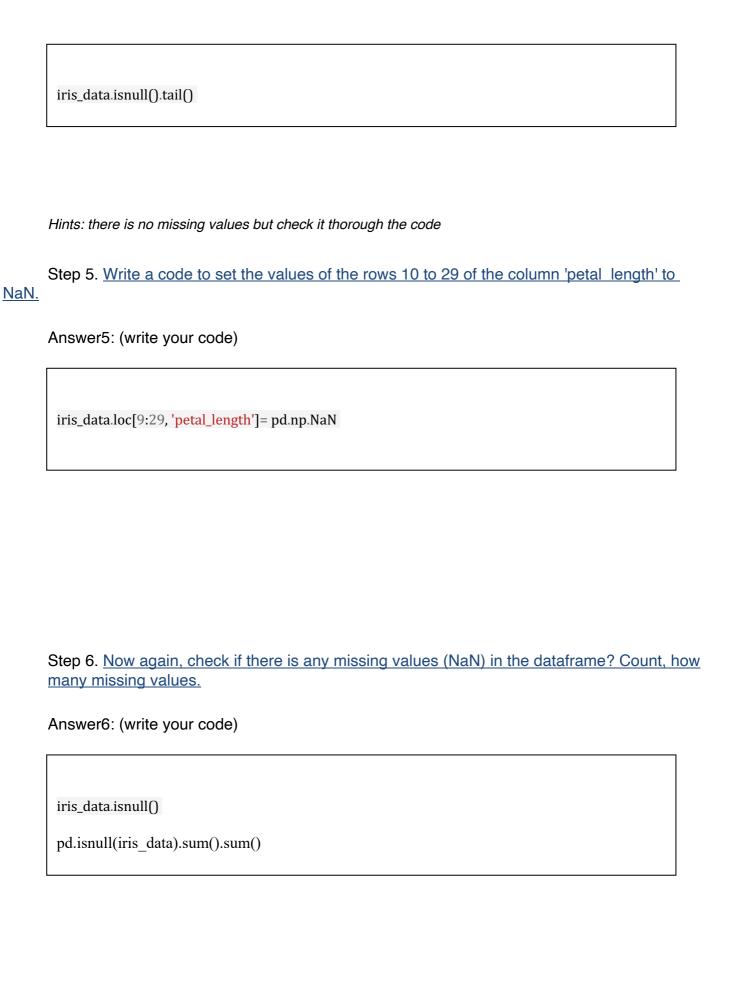
In [143]: iris.head()
Out[143]:
5.1 3.5 1.4 0.2 Iris-setosa
0 4.9 3.0 1.4 0.2 Iris-setosa
1 4.7 3.2 1.3 0.2 Iris-setosa
2 4.6 3.1 1.5 0.2 Iris-setosa
3 5.0 3.6 1.4 0.2 Iris-setosa
4 5.4 3.9 1.7 0.4 Iris-setosa
```



Step 4. Write a code to check if there are any missing values in the dataframe?

iris_data. @-columns=['sepal_length','sepal_width','petal_length','petal_width','class']

Answer4: (write your code)



Hints: this time you will have missing values.

Answer7: (write your code)

```
iris_data.fillna(10.0)
```

Section 2

Numpy is an open source library written in C++, with Python being the basic package for scientific computing.

It includes a range of methods for most machine learning tasks, algorithms
Pandas is an open source data analysis and manipulation tool that is quick, efficient,
versatile and easy to use, built on top of the programming language of Python
Matplotlib is a Python programming language plotting library and its NumPy numeric
al mathematics extension

We can use slicing and indexing in python for creating a subset of the given data frame. It can be also used to operate on a particular section in the given dataset.

Visualising data gives us a very intuitive insight of the data we want to work with. Matplotlib is one of the best visualization tools available in Python.

```
In [4]: from matplotlib import pyplot as plt

In [9]: x = [1,5,6]
y = [2,3,4]
plt.plot(x,y)
plt.plot(y,x)
plt.title('test plot')
plt.xlabel('x')
plt.ylabel('x')
plt.show()
```

