

Using Notebooks for Data Exploration

In this lab we will be exploring the Iris dataset from British statistician and biologist Ronald Fisher (read more [here \(http://archive.ics.uci.edu/ml/machine-learning-databases/iris/iris.names\)](http://archive.ics.uci.edu/ml/machine-learning-databases/iris/iris.names)). This dataset is widely used to demonstrate Machine Learning classification techniques. The dataset contains measurements of three different species of Iris flowers collected by Ronald Fisher. The variables include the sepal length, the sepal width, the petal length, and the petal width. Please refer to the image below for a visual reference. The image file is provided courtesy of Wikimedia Commons, by Danielle Langlois under the Creative Commons Attribution-Share Alike 3.0 Unported license (https://commons.wikimedia.org/wiki/File:Iris_versicolor_3.jpg) (https://commons.wikimedia.org/wiki/File:Iris_versicolor_3.jpg)



To accomplish our task of understanding the data, we will need to use some Python libraries. Some of the packages we will use today include [Pandas \(https://towardsdatascience.com/a-quick-introduction-to-the-pandas-python-library-f1b678f34673\)](https://towardsdatascience.com/a-quick-introduction-to-the-pandas-python-library-f1b678f34673), [Scipy \(https://www.scipy.org/getting-started.html\)](https://www.scipy.org/getting-started.html), [Numpy \(https://docs.scipy.org/doc/numpy-1.15.0/user/whatisnumpy.html\)](https://docs.scipy.org/doc/numpy-1.15.0/user/whatisnumpy.html), and [Matplotlib \(https://matplotlib.org/\)](https://matplotlib.org/). If you want to learn more about how these different packages work with examples, click [here \(http://cs231n.github.io/python-numpy-tutorial/\)](http://cs231n.github.io/python-numpy-tutorial/).

```
In [2]: 1 #
2 # first, we need to import a bunch of Python libraries
3 #
4
5 # pandas is a very powerful data management package
6 import pandas as pd
7 from pandas.plotting import scatter_matrix
8
9 # these are some very standard math/science Python libraries
10 import scipy as sc
11 import numpy as np
12 import matplotlib
13 import matplotlib.pyplot as plt
14
15 # Scikit-learn is a widely used machine learning package
16 import sklearn as sk
17 from sklearn.datasets import load_iris
18
```

```
In [3]: 1 #
2 # We will now load a copy of the built-in iris dataset
3 #
4
5 #
6 # this will create a Scikit learn Bunches dataset structure for us
7 #
8 dataset = load_iris()
```

```
In [6]: 1 #
2 # For starters, let's look at the way our dataset is constructed
3 # Let's look at the column headers, or "keys"
4 #
5 print('Keys of the iris dataset')
6 for k in dataset.keys():
7     print('key = %s' % (k))
```

```
Keys of the iris dataset
key = data
key = target
key = target_names
key = DESCR
key = feature_names
key = filename
```

```
In [7]: 1 #
2 # using the type() method, we can double check that
3 # our variable dataset is a data type called a
4 # scikit-learn Bunch, which is kind of like a Panda's data frame
5 #
6 type(dataset)
```

```
Out[7]: sklearn.utils.Bunch
```

```
In [16]: 1 #
          2 # Let's take a peek at the actual numeric data
          3 #
          4 # note the way we access the dataset. We specify which attribute
          5 # we want ("data"), then specify which rows we want (here 1:5)
          6 #
          7 dataset['data'][1:5]
```

```
Out[16]: array([[4.9, 3. , 1.4, 0.2],
                [4.7, 3.2, 1.3, 0.2],
                [4.6, 3.1, 1.5, 0.2],
                [5. , 3.6, 1.4, 0.2]])
```

```
In [15]: 1 #
          2 # Let's look at some of the other parts of the data
          3 #
          4 # the "feature_names" corresponds to the column headers for the numeric
          5
          6 dataset['feature_names']
          7
```

```
Out[15]: ['sepal length (cm)',
          'sepal width (cm)',
          'petal length (cm)',
          'petal width (cm)']
```

```
In [72]: 1 #
          2 # Let's look at some of the other parts of the data
          3 #
          4 # the "target_names" corresponds to the known species for each row
          5
          6 dataset['target_names']
```

```
Out[72]: array(['setosa', 'versicolor', 'virginica'], dtype='<U10')
```

```

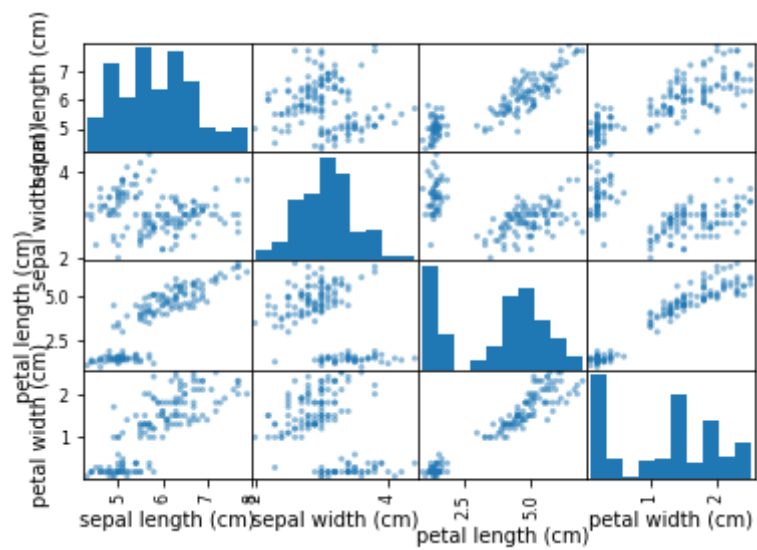
In [71]: 1 #
          2 # now let's visualize the data set
          3 # we will first create a Pandas data frame using the DataFrame() method
          4 # then we'll use the scatter_matrix() method to make a graphic
          5 #
          6 pdata = pd.DataFrame( dataset['data'], columns=dataset.feature_names )
          7 pd.plotting.scatter_matrix( pdata )
          8

```

```

Out[71]: array([[<matplotlib.axes._subplots.AxesSubplot object at 0x0000029211CD35
C0>,
                  <matplotlib.axes._subplots.AxesSubplot object at 0x0000029212E1C5
C0>,
                  <matplotlib.axes._subplots.AxesSubplot object at 0x0000029212E428
28>,
                  <matplotlib.axes._subplots.AxesSubplot object at 0x0000029212E6DA
90>],
                [<matplotlib.axes._subplots.AxesSubplot object at 0x0000029212E97C
F8>,
                  <matplotlib.axes._subplots.AxesSubplot object at 0x0000029212EBDF
60>,
                  <matplotlib.axes._subplots.AxesSubplot object at 0x0000029212EF12
08>,
                  <matplotlib.axes._subplots.AxesSubplot object at 0x0000029212F1A4
A8>],
                [<matplotlib.axes._subplots.AxesSubplot object at 0x0000029212F1A4
E0>,
                  <matplotlib.axes._subplots.AxesSubplot object at 0x0000029212F6A9
40>,
                  <matplotlib.axes._subplots.AxesSubplot object at 0x0000029212F94B
A8>,
                  <matplotlib.axes._subplots.AxesSubplot object at 0x0000029212FBEE
10>],
                [<matplotlib.axes._subplots.AxesSubplot object at 0x0000029212FF00
B8>,
                  <matplotlib.axes._subplots.AxesSubplot object at 0x000002921301B3
20>,
                  <matplotlib.axes._subplots.AxesSubplot object at 0x00000292130415
88>,
                  <matplotlib.axes._subplots.AxesSubplot object at 0x00000292112BFE
B8>]],
              dtype=object)

```



```
In [54]: 1 #
2 # let's do some basic stats
3 # the data set has 50 observations per flower type
4 # so we can look at blocks of 50 rows at a time
5 #
6 for i in range(0,3):
7     lo = 50*i
8     hi = lo + 50
9     name = dataset.target_names[i]
10    print('data for flower type %s' % (name))
11    print(pdata[lo:hi].describe())
12
```

data for flower type setosa

	sepal length (cm)	sepal width (cm)	petal length (cm) \
count	50.00000	50.000000	50.000000
mean	5.00600	3.428000	1.462000
std	0.35249	0.379064	0.173664
min	4.30000	2.300000	1.000000
25%	4.80000	3.200000	1.400000
50%	5.00000	3.400000	1.500000
75%	5.20000	3.675000	1.575000
max	5.80000	4.400000	1.900000

	petal width (cm)
count	50.000000
mean	0.246000
std	0.105386
min	0.100000
25%	0.200000
50%	0.200000
75%	0.300000
max	0.600000

data for flower type versicolor

	sepal length (cm)	sepal width (cm)	petal length (cm) \
count	50.000000	50.000000	50.000000
mean	5.936000	2.770000	4.260000
std	0.516171	0.313798	0.469911
min	4.900000	2.000000	3.000000
25%	5.600000	2.525000	4.000000
50%	5.900000	2.800000	4.350000
75%	6.300000	3.000000	4.600000
max	7.000000	3.400000	5.100000

	petal width (cm)
count	50.000000
mean	1.326000
std	0.197753
min	1.000000
25%	1.200000
50%	1.300000
75%	1.500000
max	1.800000

data for flower type virginica

	sepal length (cm)	sepal width (cm)	petal length (cm) \
count	50.00000	50.000000	50.000000
mean	6.58800	2.974000	5.552000

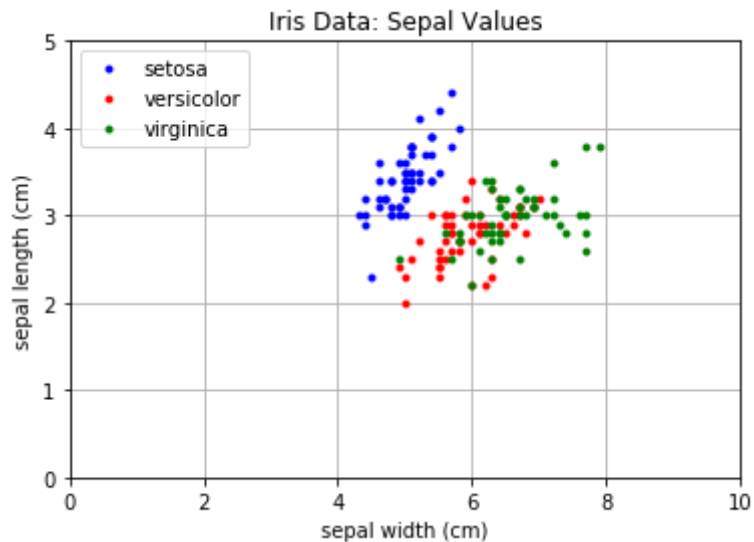
std	0.63588	0.322497	0.551895
min	4.90000	2.200000	4.500000
25%	6.22500	2.800000	5.100000
50%	6.50000	3.000000	5.550000
75%	6.90000	3.175000	5.875000
max	7.90000	3.800000	6.900000

	petal width (cm)
count	50.00000
mean	2.02600
std	0.27465
min	1.40000
25%	1.80000
50%	2.00000
75%	2.30000
max	2.50000

```

In [76]: 1 #
2 # let's do some basic visualizations
3 # the data set has 50 observations per flower type
4 # so we can look at blocks of 50 rows at a time
5 #
6
7 #
8 # make a list of blocks of data
9 #
10 blocks = []
11 for i in range(0,3):
12     lo = 50*i
13     hi = lo + 50
14     block = pdata[lo:hi]
15     blocks.append( block )
16
17 #
18 # first render the data as dots
19 #
20 plt.plot(blocks[0]['sepal length (cm)'], blocks[0]['sepal width (cm)'],
21 plt.plot(blocks[1]['sepal length (cm)'], blocks[1]['sepal width (cm)'],
22 plt.plot(blocks[2]['sepal length (cm)'], blocks[2]['sepal width (cm)'],
23
24 #
25 # set up plot embellishments
26 #
27 plt.ylim(0,5)
28 plt.xlim(0,10)
29 plt.title('Iris Data: Sepal Values')
30 plt.ylabel('sepal length (cm)')
31 plt.xlabel('sepal width (cm)')
32 plt.grid()
33 plt.legend(('setosa', 'versicolor', 'virginica'))
34 plt.show()
35

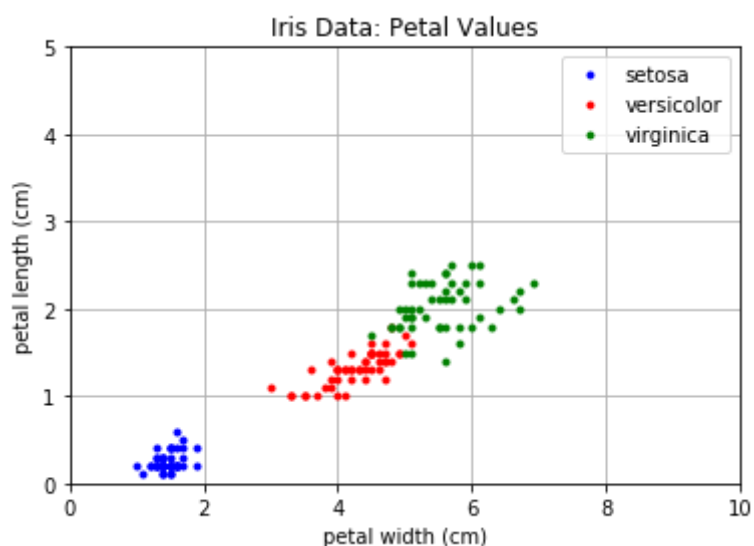
```




```

In [77]: 1 #
2 # let's do some basic visualizations
3 # the data set has 50 observations per flower type
4 # so we can look at blocks of 50 rows at a time
5 #
6
7 #
8 # make a list of blocks of data
9 #
10 blocks = []
11 for i in range(0,3):
12     lo = 50*i
13     hi = lo + 50
14     block = pdata[lo:hi]
15     blocks.append( block )
16
17 #
18 # first render the data as dots
19 #
20 plt.plot(blocks[0]['petal length (cm)'], blocks[0]['petal width (cm)'],
21 plt.plot(blocks[1]['petal length (cm)'], blocks[1]['petal width (cm)'],
22 plt.plot(blocks[2]['petal length (cm)'], blocks[2]['petal width (cm)'],
23
24 #
25 # set up plot embellishments
26 #
27 plt.ylim(0,5)
28 plt.xlim(0,10)
29 plt.title('Iris Data: Petal Values')
30 plt.ylabel('petal length (cm)')
31 plt.xlabel('petal width (cm)')
32 plt.grid()
33 plt.legend(('setosa', 'versicolor', 'virginica'))
34 plt.show()

```



```
In [ ]: 1 #  
        2 # now use the data to answer some questions  
        3 #  
        4 # 1. what is the mean petal size of each species of Iris  
        5 # 2. what is the mean sepal size of each species of Iris?  
        6 # 3. could you distinguish species using only sepal length and width? W  
        7 # 4. could you distinguish species using only petal length and width? W
```