Homework 0

# CS 5785 Applied Machine Learning

Xialin Shen <xs293@cornell.edu>

An Le <aql6@cornell.edu>

24th August, 2017

# Question 1

Find and download the Iris Flowers dataset from the UC Irvine Machine Learning datasets archive at http://archive.ics.uci.edu/ml/datasets.html

Hint: The iris.names file describes the structure of the dataset. How many features/attributes are there per sample? How many different species are there, and how many samples of each species did Anderson record?

**ANSWERS:**

1. There are **4** attributes per sample, including "sepal length in cm", "sepal width in cm", "petal length in cm" and "petal width in cm"
2. There are **3** different species, including "Iris Setosa", "Iris Versicolour" and "Iris Virginica"
3. Anderson recorded **50** samples for each species.

# Question 2

Figure out how to parse the dataset you downloaded. Load the samples into an $N \times p$ array, where $N$ is the number of samples and $p$ is the number of attributes per sample. Additionally, create a $N$ -dimensional vector containing each sample's label (species).

**Xialin:**

**CODES:**

```python
 2 # -*- coding: utf-8 -*-
 3 """
 4 Created on Thu Aug 24 23:05:05 2017
 5
 6 @author: shenxialin (xs293@cornell.edu)
 7 """
 8 from matplotlib import pyplot as plt
 9 import numpy
10
11 # setup the notation for divide a line of data
12 split_token = ","
13
14 # there are 4 features in every sample
15 features_count = 4
16 features_name = ["sepal length", "sepal width", "petal length", "petal width"]
17
18 # there are total 3 different species
19 species = ["Iris-setosa","Iris-versicolor", "Iris-virginica"]
20 colors_options = ["r", "b", "g"]
21
22 # there are total 150 samples, and 50 samples per species
23 samples_count = 150
24 sample_set = [[0 for x in range(features_count)] for y in range(samples_count)]
25
26 labels = []
27 colors = []
28
29 index = 0
30 for line in open("./data/iris.data.txt"):
31     # print ("Line contains: " + line)
32
33     # get a line of data and extract features and label
34     attributesList = line.split(split_token)
35     features = attributesList[: features_count]
36     sample_set[index] = features
37     label = attributesList[len(attributesList) - 1].rstrip()
38     labels.append(label)
39
40     # assign different color for sample in different species
41     species_index = species.index(label)
42     if (species_index != ValueError):
43         colors.append(colors_options[species_index])
44
45     index += 1
```

**RESULTS:**

<sample_set>: a 150 x 4 dimensions matrix

<labels>: a list with 150 labels

**An:**

## Import data

```
In [139]:  1  import pandas as pd
           2  data_set = pd.read_csv('iris.data', header=None,
           3                    names=['Sepal Length','Sepal Width',
           4                           'Petal Length', 'Petal Width', 'Class'])
           5  print(data_set[:5])
```

```
   Sepal Length  Sepal Width  Petal Length  Petal Width        Class
0           5.1          3.5           1.4          0.2  Iris-setosa
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
```

```
In [140]:  1  print("Shape of dataset: {}".format(data_set.shape))
```

```
Shape of dataset: (150, 5)
```

```
In [141]:  1  data_set['Class'].value_counts()
```

```
Out[141]: Iris-setosa        50
          Iris-versicolor    50
          Iris-virginica     50
          Name: Class, dtype: int64
```

# Question 3

To visualize this dataset, we would have to build a *p*-dimensional scatterplot. Unfortunately, we only have 2D displays so we must reduce the dataset's dimensionality. The easiest way to view the set is to plot two attributes of the data against one another and repeat for each pair of attributes.

Create every possible scatterplot from all pairs of two attributes. (For example, one scatterplot would graph petal length vs sepal width, another would graph petal length vs. sepal length, and so on). Within each scatterplot, the color of each dot should correspond with the sample species.
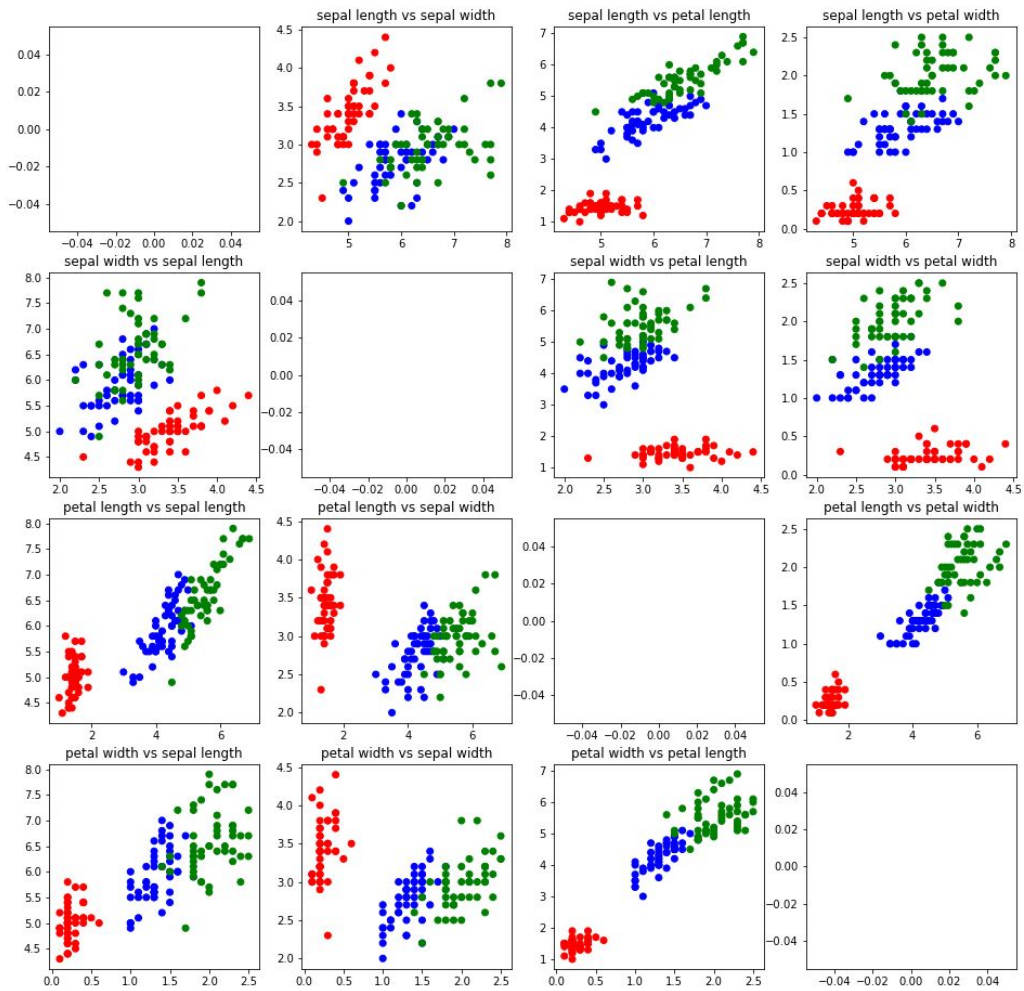
**Xialin:**

**CODES (CON'T):**

```python
47 # setup figure size and title
48 plt.figure(figsize=(16, 16))
49 plt.suptitle('Iris Data', fontsize=20, fontweight='bold')
50
51 plot_index = 1
52 for x in range(features_count):
53     for y in range(features_count):
54         plt.subplot(features_count,features_count,plot_index)
55         plot_index += 1
56
57         if (x == y):
58             plt.plot([], [])
59             continue
60
61         xs = numpy.array(sample_set)[:,x]
62         ys = numpy.array(sample_set)[:,y]
63
64         plt.title(features_name[x] + " vs " + features_name[y])
65         plt.scatter(xs, ys, c=colors)
66
67 plt.savefig("plot.png")
```
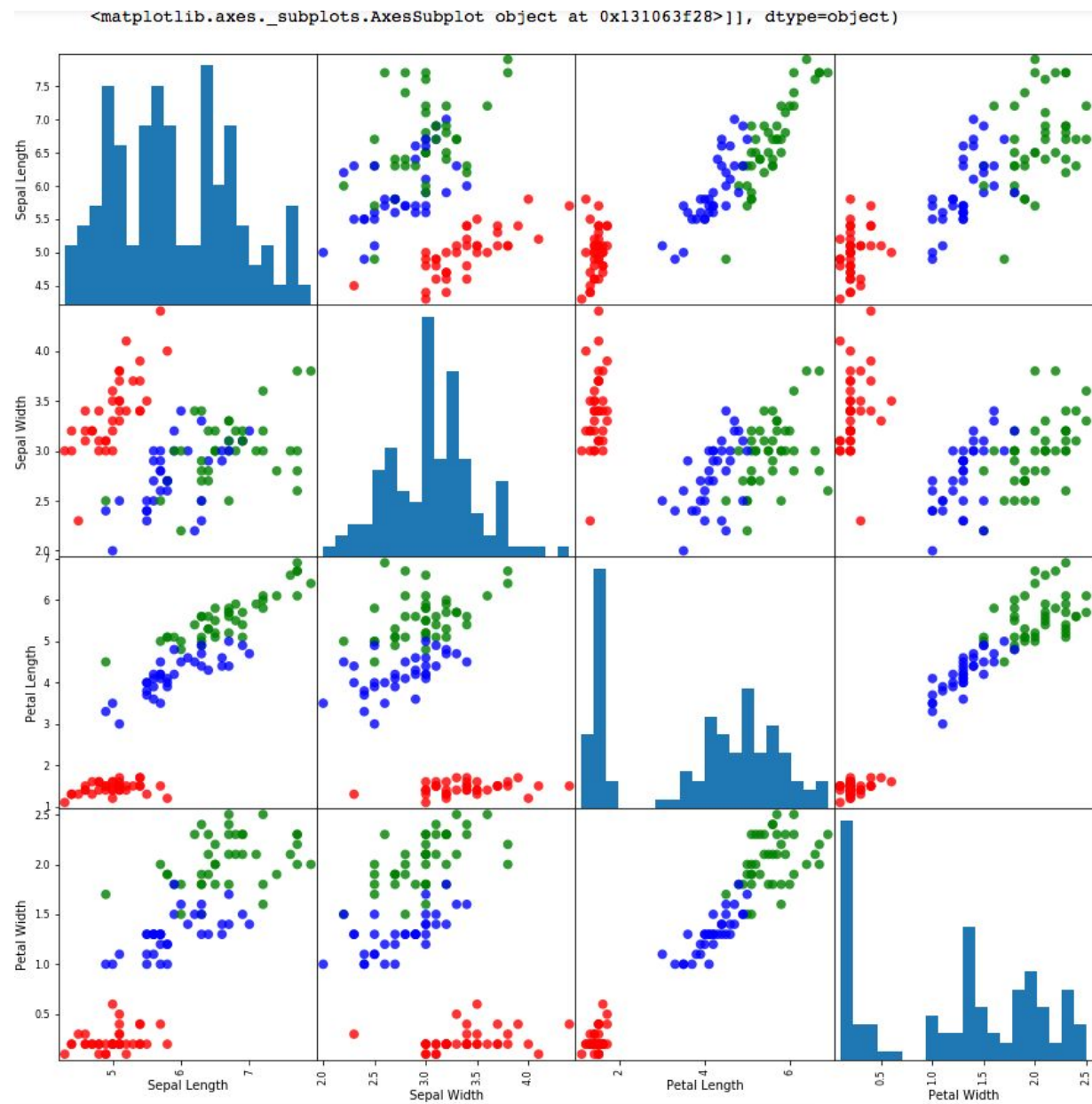
**RESULTS:**

**Iris Data**

## Convert to narray and split

```
In [142]:  1  column_names = ["Sepal Length","Sepal Width","Petal Length","Petal Width"];
           2  data = data_set.as_matrix(column_names)
           3  target = data_set.as_matrix(["Class"])
           4
           5  from sklearn.model_selection import train_test_split
           6  X_train, X_test, y_train, y_test = train_test_split(data, target, random_state = 0)
```

## Plot

```
In [143]:   1  %matplotlib inline
            2  import matplotlib.pyplot as plt
            3  import numpy as np
            4
            5  iris_dataset = pd.DataFrame(X_train, columns=column_names)
            6
            7  colors = []
            8  for e in y_train:
            9      if e=='Iris-setosa':
           10          colors.append('r')
           11      elif e=='Iris-versicolor':
           12          colors.append('b')
           13      elif e=='Iris-virginica':
           14          colors.append('g')
           15
           16  pd.plotting.scatter_matrix(iris_dataset, figsize=(15,15),
           17                             marker='o', hist_kwds={'bins': 20}, s=60,
           18                             alpha=.8, c=colors)
```

**RESULTS:**



`<matplotlib.axes._subplots.AxesSubplot object at 0x131063f28>]], dtype=object)`

## CODE

[Xialin](#)

[An Le](#)

- END -