

Module 3 : Data Exploration

Reference Book

- Reference the "Data Exploration" chapter in Book [Introduction to Data Mining](#) to understand some of the concepts introduced in this tutorial notebook.
- [Notebook Down](#)

Data Exploration

- Data exploration refers to the preliminary investigation of data in order to better understand its specific characteristics.
- Two motivations :
 1. To help users select the appropriate preprocessing and data analysis technique used.
 2. To make use of humans' abilities to recognize patterns in the data.

3.1 Summary Statistics

Introduction

Summary statistics are quantities, such as the mean and standard deviation, that capture various characteristics of a potentially large set of values with a single number or a small set of numbers. In this tutorial, we will use the Iris sample data, which contains information on 150 Iris flowers, 50 each from one of three Iris species: Setosa, Versicolour, and Virginica. Each flower is characterized by five attributes:

- sepal length in centimeters
- sepal width in centimeters
- petal length in centimeters
- petal width in centimeters
- class(Setosa,Versicolour,Virginica)

Target

- Load a CSV data file into a Pandas DataFrame object;
- Compute various summary statistics from the DataFrame.

(1) Download CSV File

UCI Machine Learning Repository : [Iris Dataset](#)

```
1 import pandas as pd
2
3 data = pd.read_csv("http://archive.ics.uci.edu/ml/machine-learning-
databases/iris/iris.data",header=None)
4 data.columns = ['sepal length','sepal width','petal length','petal width','class']
5
6 data
```

```
1 .dataframe tbody tr th {
2     vertical-align: top;
3 }
4
5 .dataframe thead th {
6     text-align: right;
7 }
```

	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
...
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	Iris-virginica
149	5.9	3.0	5.1	1.8	Iris-virginica

150 rows × 5 columns

```
1 data.head() # display first five row
```

```

1 .dataframe tbody tr th {
2     vertical-align: top;
3 }
4
5 .dataframe thead th {
6     text-align: right;
7 }

```

	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

(2) Calculation

For each quantitative attribute, calculate its average, standard deviation, minimum, and maximum values.

Functions

- `mean()`: get the average
- `std()` : get the standard deviation
- `min()` : get the minimum
- `max()` : get the maximum

```

1 from pandas.api.types import is_numeric_dtype
2 for col in data.columns:
3     if is_numeric_dtype(data[col]):
4         print('%s:' % (col))
5         print('\t Mean = %.2f' % data[col].mean())
6         print('\t Standard deviation = %.2f' % data[col].std())
7         print('\t Minimum = %.2f' % data[col].min())
8         print('\t Maximum = %.2f' % data[col].max())
9

```

```

1 sepal length:
2     Mean = 5.84
3     Standard deviation = 0.83
4     Minimum = 4.30
5     Maximum = 7.90
6 sepal width:
7     Mean = 3.05
8     Standard deviation = 0.43
9     Minimum = 2.00
10    Maximum = 4.40

```

```

11 petal length:
12     Mean = 3.76
13     Standard deviation = 1.76
14     Minimum = 1.00
15     Maximum = 6.90
16 petal width:
17     Mean = 1.20
18     Standard deviation = 0.76
19     Minimum = 0.10
20     Maximum = 2.50

```

(3) Count the Frequency

For the qualitative attribute (class), count the frequency for each of its distinct values.

```

1 data['class'].value_counts()

```

```

1 Iris-versicolor    50
2 Iris-setosa        50
3 Iris-virginica     50
4 Name: class, dtype: int64

```

(4) describe() function

- It is also possible to display the summary for all the attributes simultaneously in a table using the describe() function.
- If an attribute is quantitative, it will display its mean, standard deviation and various quantiles (including minimum, median, and maximum) values.
- If an attribute is qualitative, it will display its number of unique values and the top (most frequent) values.

```

1 data.describe(include='all')

```

```

1 .dataframe tbody tr th {
2     vertical-align: top;
3 }
4
5 .dataframe thead th {
6     text-align: right;
7 }

```

	sepal length	sepal width	petal length	petal width	class
count	150.000000	150.000000	150.000000	150.000000	150
unique	NaN	NaN	NaN	NaN	3
top	NaN	NaN	NaN	NaN	Iris-versicolor
freq	NaN	NaN	NaN	NaN	50
mean	5.843333	3.054000	3.758667	1.198667	NaN
std	0.828066	0.433594	1.764420	0.763161	NaN
min	4.300000	2.000000	1.000000	0.100000	NaN
25%	5.100000	2.800000	1.600000	0.300000	NaN
50%	5.800000	3.000000	4.350000	1.300000	NaN
75%	6.400000	3.300000	5.100000	1.800000	NaN
max	7.900000	4.400000	6.900000	2.500000	NaN

(5) Multivariate Statistics

- Compute the covariance and correlation between pairs of attributes.

Function

- `cov()` : calculate the covariance
- `corr()`: calculate the correlation

```
1 print("Covariance:")
2 data.cov()
```

```
1 Covariance:
```

```
1 .dataframe tbody tr th {
2     vertical-align: top;
3 }
4
5 .dataframe thead th {
6     text-align: right;
7 }
```

	sepal length	sepal width	petal length	petal width
sepal length	0.685694	-0.039268	1.273682	0.516904
sepal width	-0.039268	0.188004	-0.321713	-0.117981
petal length	1.273682	-0.321713	3.113179	1.296387
petal width	0.516904	-0.117981	1.296387	0.582414

```
1 print("Correlation:")
2 data.corr()
```

```
1 Correlation:
```

```
1 .dataframe tbody tr th {
2     vertical-align: top;
3 }
4
5 .dataframe thead th {
6     text-align: right;
7 }
```

	sepal length	sepal width	petal length	petal width
sepal length	1.000000	-0.109369	0.871754	0.817954
sepal width	-0.109369	1.000000	-0.420516	-0.356544
petal length	0.871754	-0.420516	1.000000	0.962757
petal width	0.817954	-0.356544	0.962757	1.000000

3.2 Data Visualization

- Data visualization is the display of information in a graphic or tabular format.
- Target : the characteristics of the data and the relationships among data items or attributes can be analyzed or reported.

(1) Histogram

Function : hist()

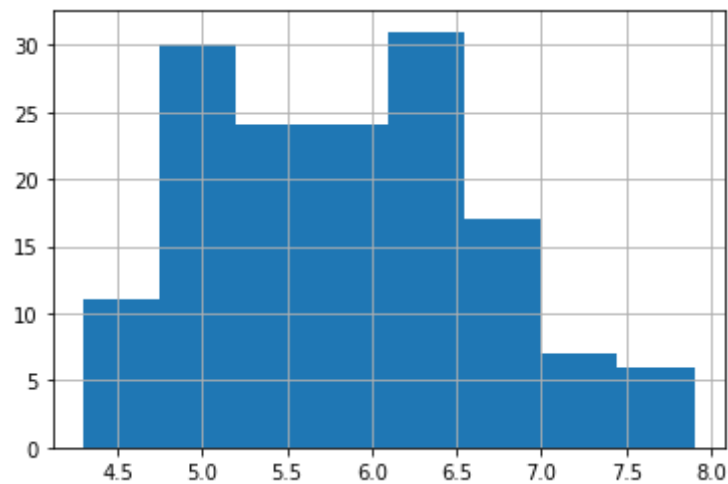
- Create the histogram for target attribute;
- Count the frequency for each bin.

parameter :

- bin : the number of separate bins

```
1 # First, we will display the histogram for the sepal length attribute
2 # by discretizing it into 8 separate bins and counting the frequency
3 # for each bin.
4
5 %matplotlib inline
6 data['sepal length'].hist(bins = 8)
```

```
1 <matplotlib.axes._subplots.AxesSubplot at 0x1a586045548>
```



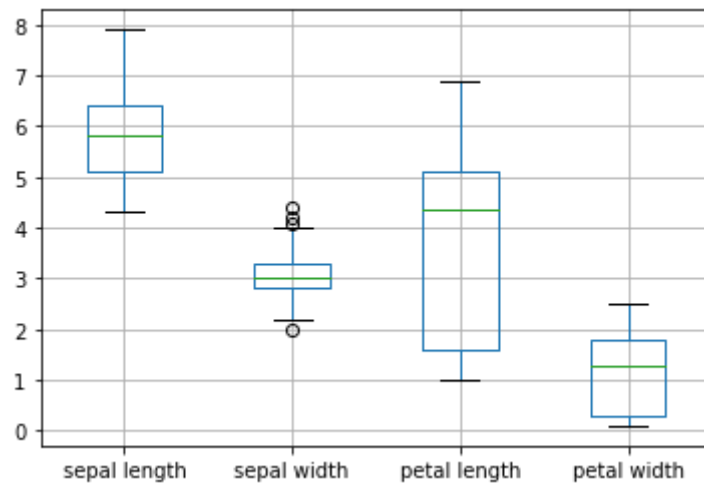
(2) Box plot

Function : boxplot()

- Show the distribution of values for each attribute.

```
1 data.boxplot()
```

```
1 <matplotlib.axes._subplots.AxesSubplot at 0x1a58862cf08>
```



(3) Scatter Plot

For each pair of attributes, we can use a scatter plot to visualize their joint distribution.

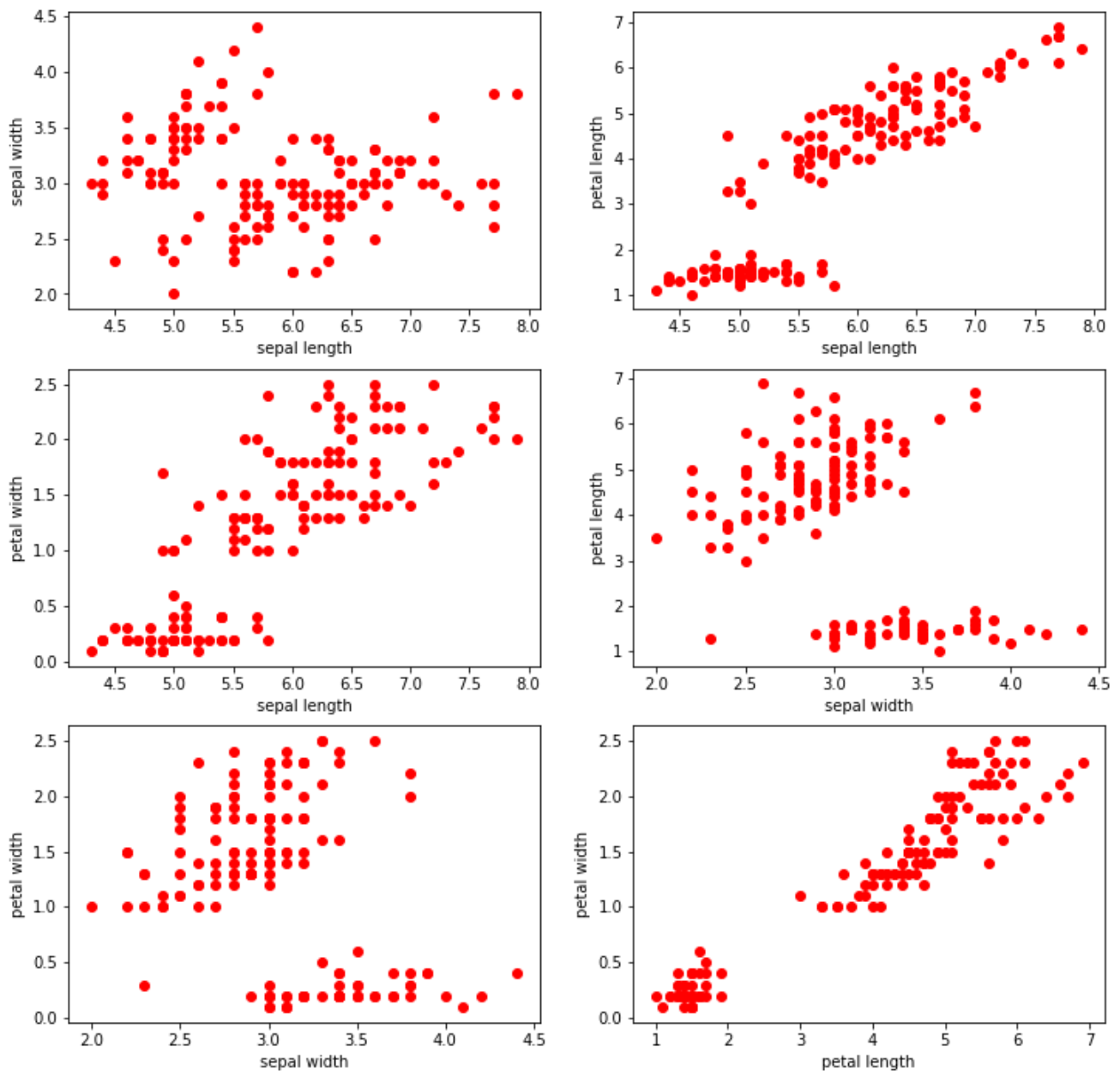
Function

- `subplots()`

```

1  import matplotlib.pyplot as plt
2
3  fig, axes = plt.subplots(3, 2, figsize=(12, 12))
4  index = 0
5  for i in range(3):
6      for j in range(i+1, 4):
7          ax1 = int(index/2)
8          ax2 = index % 2
9          axes[ax1][ax2].scatter(data[data.columns[i]], data[data.columns[j]], color='red')
10         axes[ax1][ax2].set_xlabel(data.columns[i])
11         axes[ax1][ax2].set_ylabel(data.columns[j])
12         index = index + 1

```

(4) Paralle Corrdinates

- Parallel coordinates can be used to display all the data points simultaneously.
- Parallel coordinates have one coordinate axis for each attribute, but the different axes are parallel to one other instead of perpendicular, as is traditional.
- Furthermore, an object is represented as a line instead of as a point. In the example below, the distribution of values for each class can be identified in a separate color.

```
1 from pandas.plotting import parallel_coordinates
2 %matplotlib inline
3
4 parallel_coordinates(data, 'class')
```

```
1 <matplotlib.axes._subplots.AxesSubplot at 0x1a588ac5208>
```

