

EDA using Data Visualization

Data visualization is the most effective tool in exploratory data analysis as it gives a birds eye view of the entire data in the most efficient way and amplifies the cognition. Python library package matplotlib, provides a very user friendly approach in data visualization. There are two approaches namely stateless and stateful.

Stateful and Stateless Approaches

In stateful approach, the program refers the pyplot directly and any changes made to the state reflect across all plots. Whereas in stateless approach is an object oriented approach, where an object with properties of figure and axes is created from pyplot and that object is referenced in creation of subsequent plots instead of directly calling pyplot function.

Pros and cons of each approach

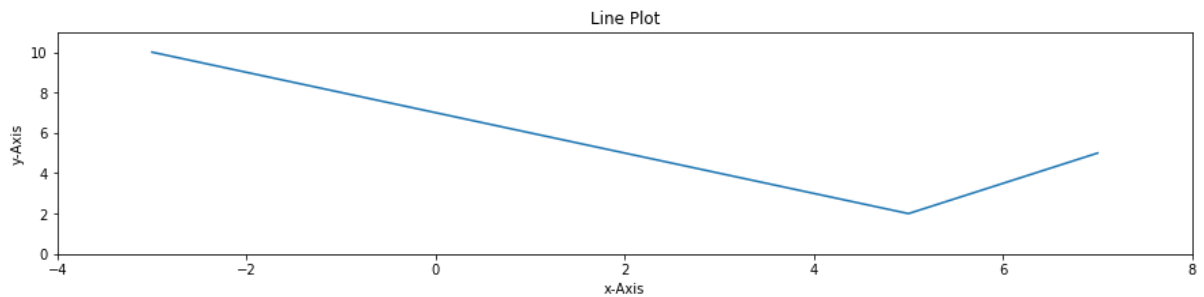
Pros and cons of each approach is dependent on the usage. For simple and fast plotting stateful is preferred but will become difficult to manage as the number of plots increases. On the other hand stateless approach has the overhead of creating an object but is easier to keep track of each plot and properties.

Before diving into the EDA, providing an example of each of these approaches for better understanding

```
In [1]: # import necessary packages
import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
%matplotlib inline

#Create two lists
x = [-3,5,7]
y = [10,2,5]

# Plot the above lists elements pairs on x-y axis using stateful approach
plt.figure(figsize=(15,3))
plt.plot(x,y)
plt.xlim(-4,8)
plt.ylim(0,11)
plt.xlabel("x-Axis")
plt.ylabel("y-Axis")
plt.title('Line Plot')
plt.show()
```



```
In [2]: %matplotlib inline
# Plot the above lists elements pairs on x-y axis using stateless approach

#create a pyplot object for figure and axis
fig,ax = plt.subplots(figsize=(15,3))

ax.plot(x,y)
ax.set_xlim(-4,8)
ax.set_ylim(-1,11)
ax.set_xlabel("x-Axis")
ax.set_ylabel("y-Axis")
ax.set_title('Line Plot')
fig.show()
```

```
In [4]: #Load data into a pandas dataframe
dirpath = "/Users/jacobdenny/Desktop/Babloo/Masters/Semester3/ADTA5340/D
ataFiles/"
dfIris = pd.read_csv(dirpath+"Iris.csv")

#get familiarized with the data
dfIris.head(5)
```

Out[4]:

	Id	SepalLengthCm	SepalWidthCm	PetalLengthCm	PetalWidthCm	Species
0	1	5.1	3.5	1.4	0.2	Iris-setosa
1	2	4.9	3.0	1.4	0.2	Iris-setosa
2	3	4.7	3.2	1.3	0.2	Iris-setosa
3	4	4.6	3.1	1.5	0.2	Iris-setosa
4	5	5.0	3.6	1.4	0.2	Iris-setosa

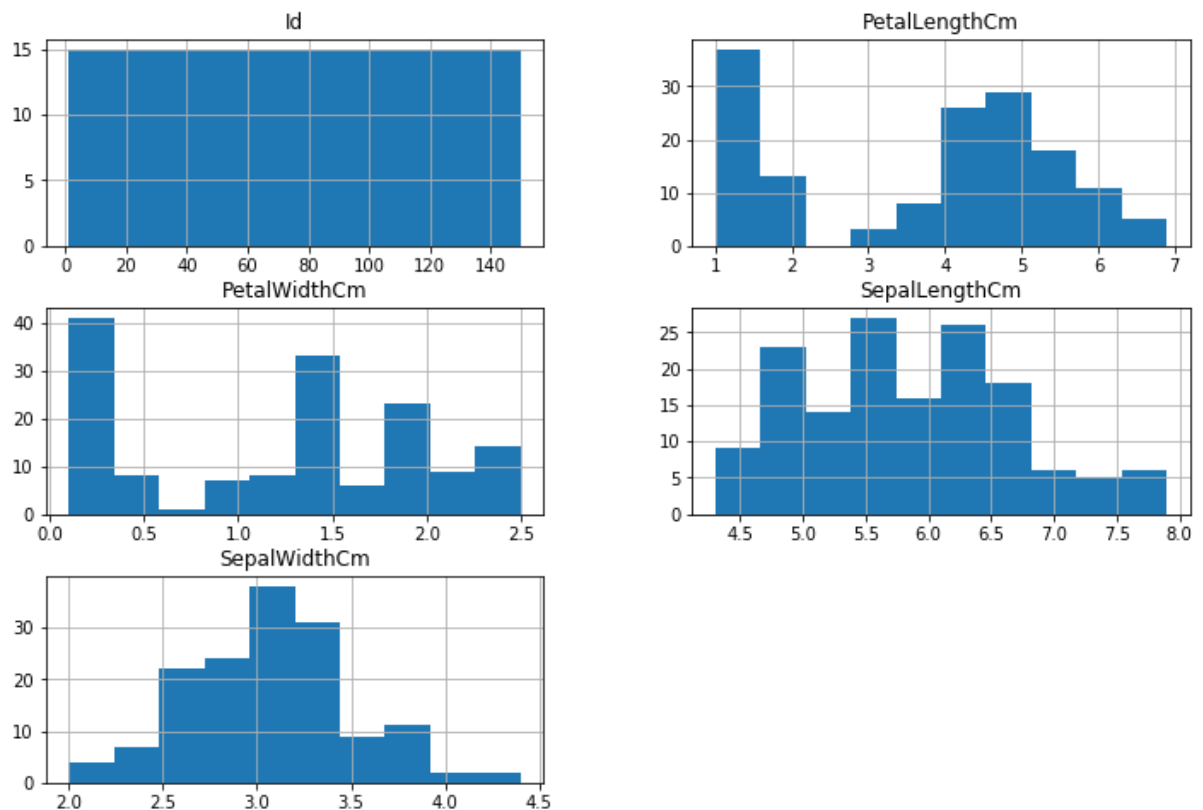
Univariate Plots

Plot each column to get an idea of the distribution. We will use the following

- Histogram
- Densityplot
- Boxplot

```
In [5]: #Histogram
dfIris.hist(figsize=(12,8))
plt.show
```

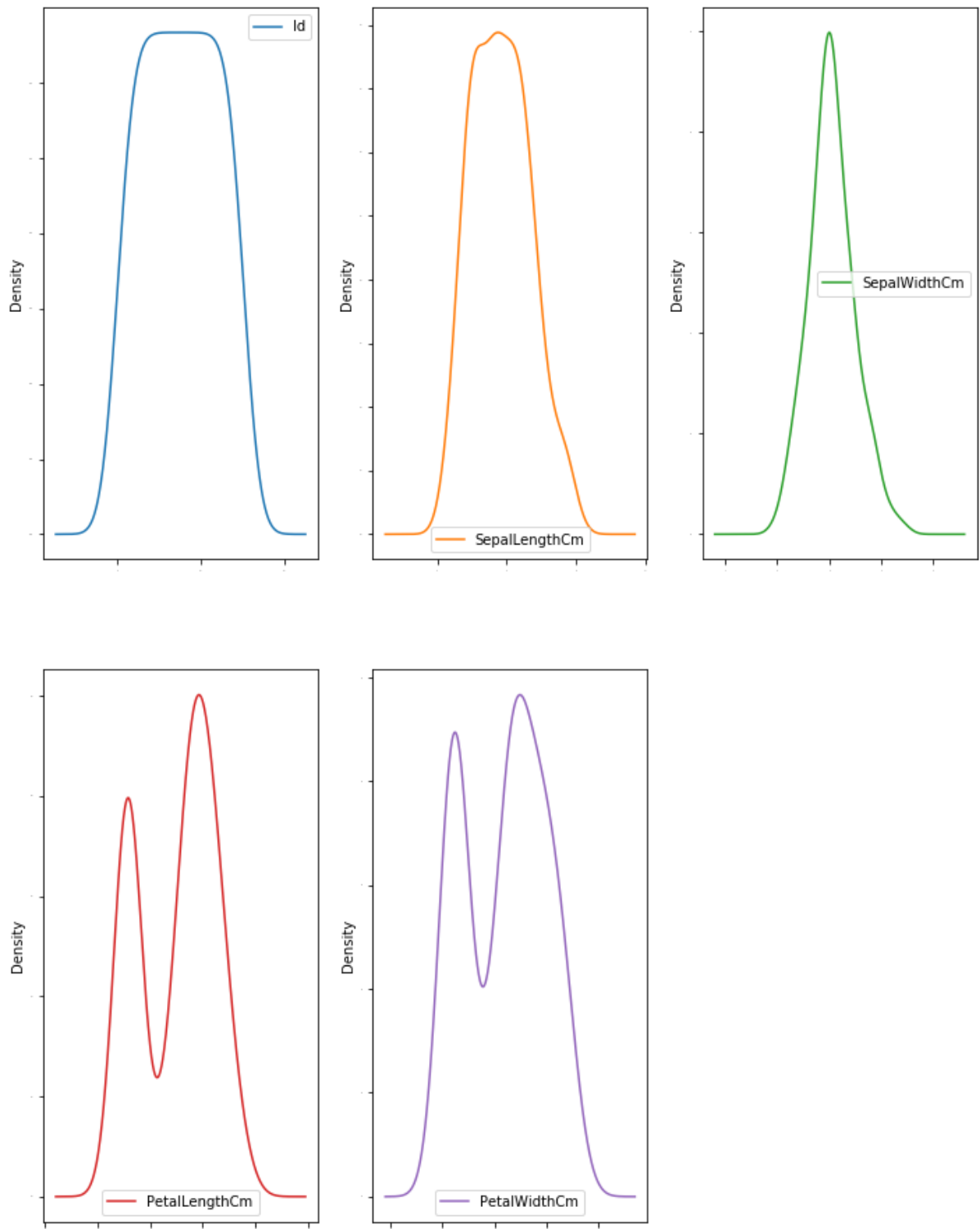
```
Out[5]: <function matplotlib.pyplot.show(*args, **kw)>
```



What Histogram Tells About Data

The histograms are displayed for each column of the data, which means the id column also gets plotted, which has no value in the data analysis and is ignored. Sepallength and sepalwidth shows almost normal distribution, petal lenght and width seems to be bimodal

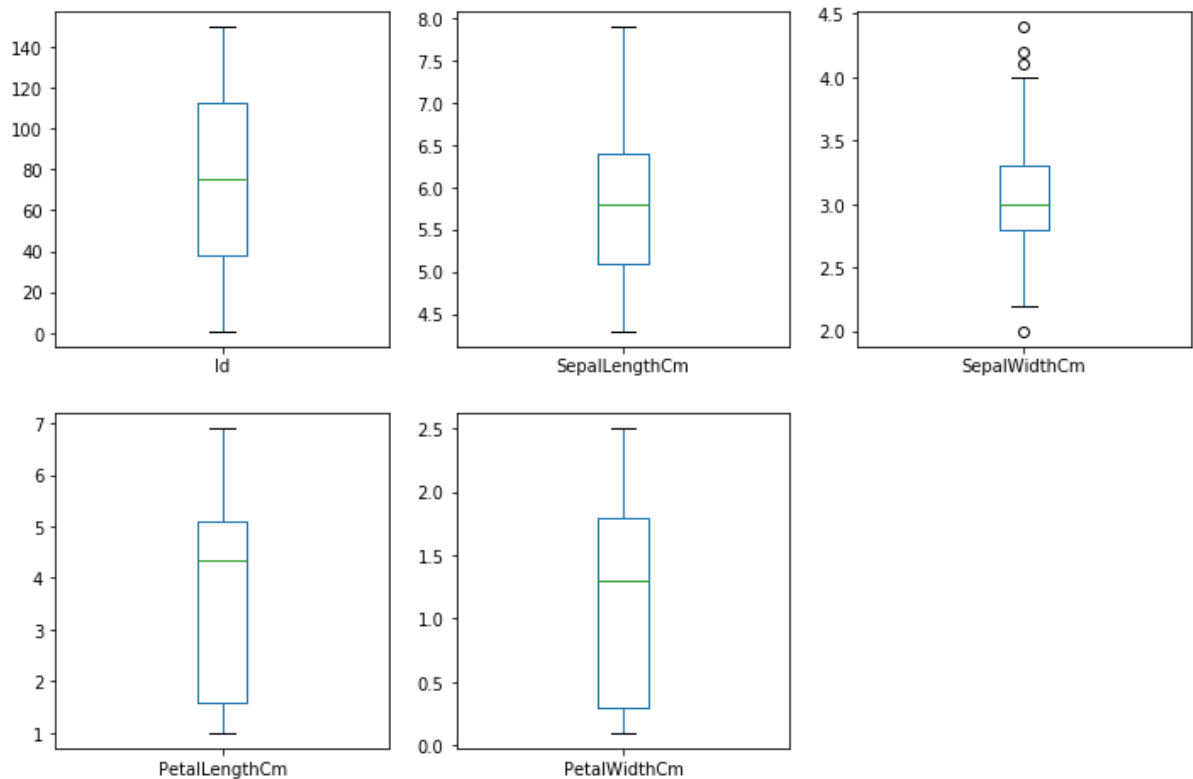
```
In [6]: #Density plot
dfIris.plot(kind="density", subplots = True, layout = (2,3), sharex = False,
legend = True, fontsize=1, figsize=(12,16) )
plt.show()
```



What Densityplots Tells About Data

The density plot is consistent with histogram and are displayed for each column of the data, which means the id column also gets plotted, which has no value in the data analysis and is ignored. Sepallength and sepalwidth shows almost normal distribution, petal length and width seems to be bimodal

```
In [7]: #Boxplot
dfIris.plot(kind="box", subplots = True, layout = (2,3), sharex = False,
sharey = False, figsize=(12,8) )
plt.show()
```



What Boxplots Tells About Data

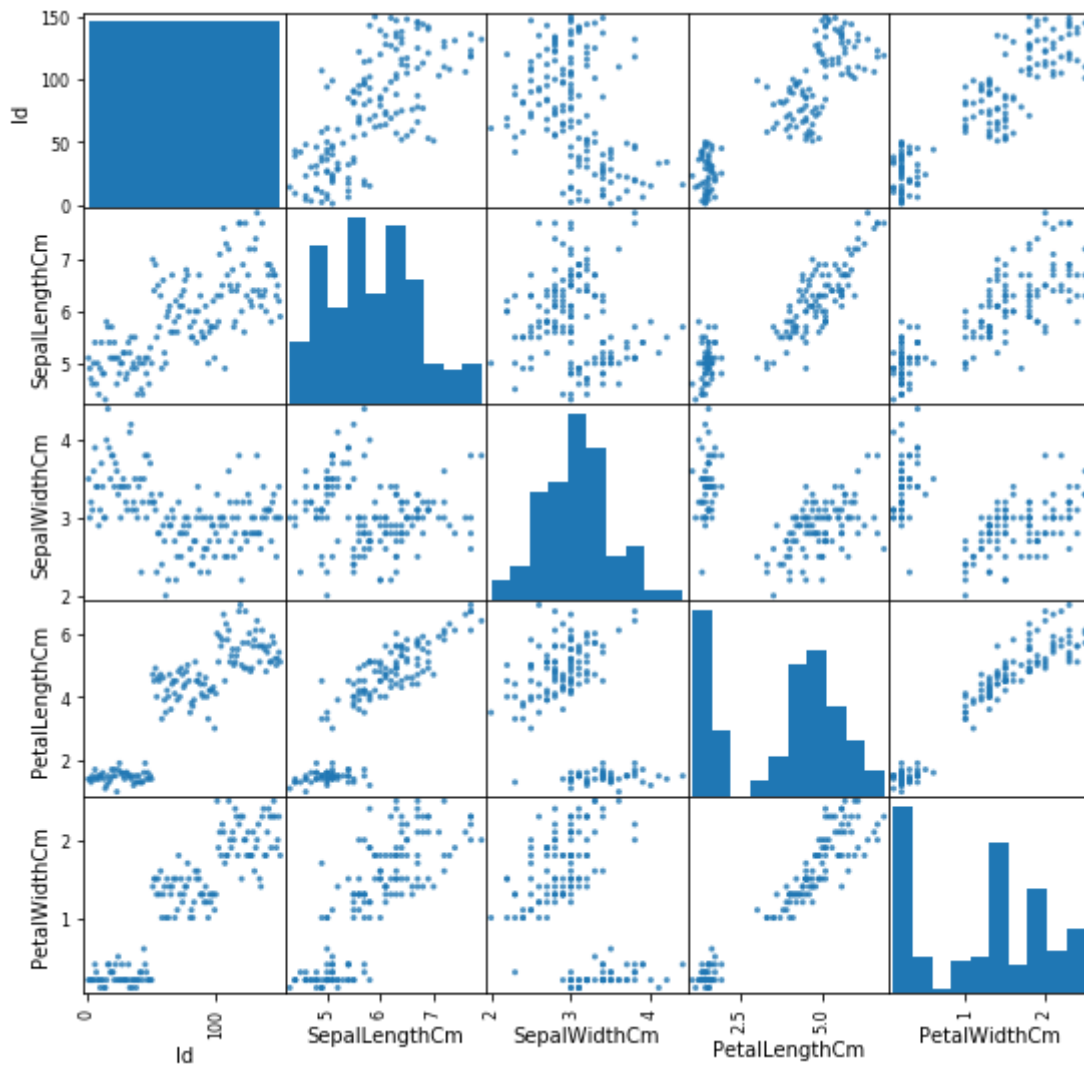
- Ignore the plot for id
- For sepal length the values range from 4.5cm to 8cm
- For sepal width the values range from 2cm to 4.5cm
- For petal length the values range from 1cm to 7cm
- For petal width the values range from 0 to 2.5cm
- Sepalwidth has evident outliers

Multivariate Analysis

The above plots shows characteristics of each column like distribution of data, range, dispersion or independent features of each column. Now we would like to see how each feature is affected by other features and we use multivariate analysis for this. Scatter matrix is the best technique for this analysis.


```
In [8]: #Scatter matrix
```

```
scatter_matrix(dfIris, alpha=0.8, figsize=(9,9))  
plt.show()
```



What scatter matrix tells us

- As sepal length increase petal length and width increase
- Sepal width exhibits almost null correlation with other features and hence has data points scattered all over
- Petal length and petal width also shows a positive linear correlation

```
In [ ]:
```

```
In [ ]:
```

```
In [ ]:
```

```
In [ ]:
```

In []:

Supervised Linear Regression

Creating a model defining the outcome variable as a linear combination of a set of independent features. Steps involve:

- Preprocess data
- EDA
- Split data into test and train set
- Build and train model
- Test/evaluate model

```
In [9]: # import scikit module for linear regression
from sklearn.linear_model import LinearRegression
# import scikit module for splitting data into training and testing datasets
from sklearn.model_selection import train_test_split
# import scikit module for kfold cross validation for evaluation and cross validation
from sklearn.model_selection import KFold
from sklearn.model_selection import cross_val_score
```

```
In [10]: #Load data into a pandas dataframe
dirpath = "/Users/jacobdenny/Desktop/Babloo/Masters/Semester3/ADTA5340/DataFiles/"
df = pd.read_csv(dirpath+"housing_boston_w_hdrs.csv")

# Not assigning column headers, as this dataset has headers
#get familiarized with the data
df.head(5)
```

Out[10]:

	CRIM	ZN	INDUS	CHAS	NOX	RM	AGE	DIS	RAD	TAX	PTRATIO	B	LSTAT
0	0.00632	18.0	2.31	0	0.538	6.575	65.2	4.0900	1	296	15.3	396.90	4.98
1	0.02731	0.0	7.07	0	0.469	6.421	78.9	4.9671	2	242	17.8	396.90	9.14
2	0.02729	0.0	7.07	0	0.469	7.185	61.1	4.9671	2	242	17.8	392.83	4.03
3	0.03237	0.0	2.18	0	0.458	6.998	45.8	6.0622	3	222	18.7	394.63	2.94
4	0.06905	0.0	2.18	0	0.458	7.147	54.2	6.0622	3	222	18.7	396.90	5.33

Preprocessing Data

We are trying to predict the median value of owner occupied houses in 1000 dollars(MEDV) using linear regression with 'RM', 'AGE', 'DIS', 'RAD', 'PTRATIO' as predictors. Hence the subset is extracted from dataframe and is stored into another data frame. Also, from domain knowledge we know that the average value of a home, the pupil teacher ratio and average number of rooms per dwelling cannot be zero. If any row cell contain zero for these three columns that needs to be marked as missing. So preprocessing includes:

- Extract the columns under consideration
- Data cleaning
- Check for invalid Values

```
In [11]: # Extract a sub-dataset from the original one --> dataframe: df2
df2 = df[['RM', 'AGE', 'DIS', 'RAD', 'PTRATIO', 'MEDV']]
#Data cleaning
#Replace 0 value with NaN for columns 'RM', 'PTRATIO', 'MEDV'
df[['RM', 'PTRATIO', 'MEDV']] = df[['RM', 'PTRATIO', 'MEDV']].replace(0,
np.NaN)
#Verify there are no missing values in the dataset
print(df.isnull().sum())
```

```
CRIM      0
ZN        0
INDUS     0
CHAS      0
NOX       0
RM        0
AGE       0
DIS       0
RAD       0
TAX       0
PTRATIO   0
B         0
LSTAT     0
MEDV      0
dtype: int64
```

There are no missing values in the dataset. Now we are ready to move to Exploratory data analysis

Exploratory Data Analysis

In this step we are getting to know our data. The number of columns and rows, what each column represent etc. The following are explored in this section:

- Shape of data set
- Type of each value in the dataset
- View a snippet of the dataset
- Descriptive statistics of each column
- Univariate plots
- Multivariate plots

```
In [12]: #SHAPE
# Get the dimensions or Shape of the dataset
# i.e. number of records/rows x number of variables/columns
print(df2.shape)

(452, 6)
```

We have 6 variables and 452 rows.

```
In [13]: #TYPE
# Get the data types of all variables/attributes of the data set
# The results show
print(df2.dtypes)

RM          float64
AGE          float64
DIS          float64
RAD           int64
PTRATIO      float64
MEDV         float64
dtype: object
```

All the features and the dependent variable(MEDV) are of type floating point except RAD which is integer type.

```
In [14]: #View of a snippet of Dataset
# Get the first five records
print(df2.head(5))
```

	RM	AGE	DIS	RAD	PTRATIO	MEDV
0	6.575	65.2	4.0900	1	15.3	24.0
1	6.421	78.9	4.9671	2	17.8	21.6
2	7.185	61.1	4.9671	2	17.8	34.7
3	6.998	45.8	6.0622	3	18.7	33.4
4	7.147	54.2	6.0622	3	18.7	36.2

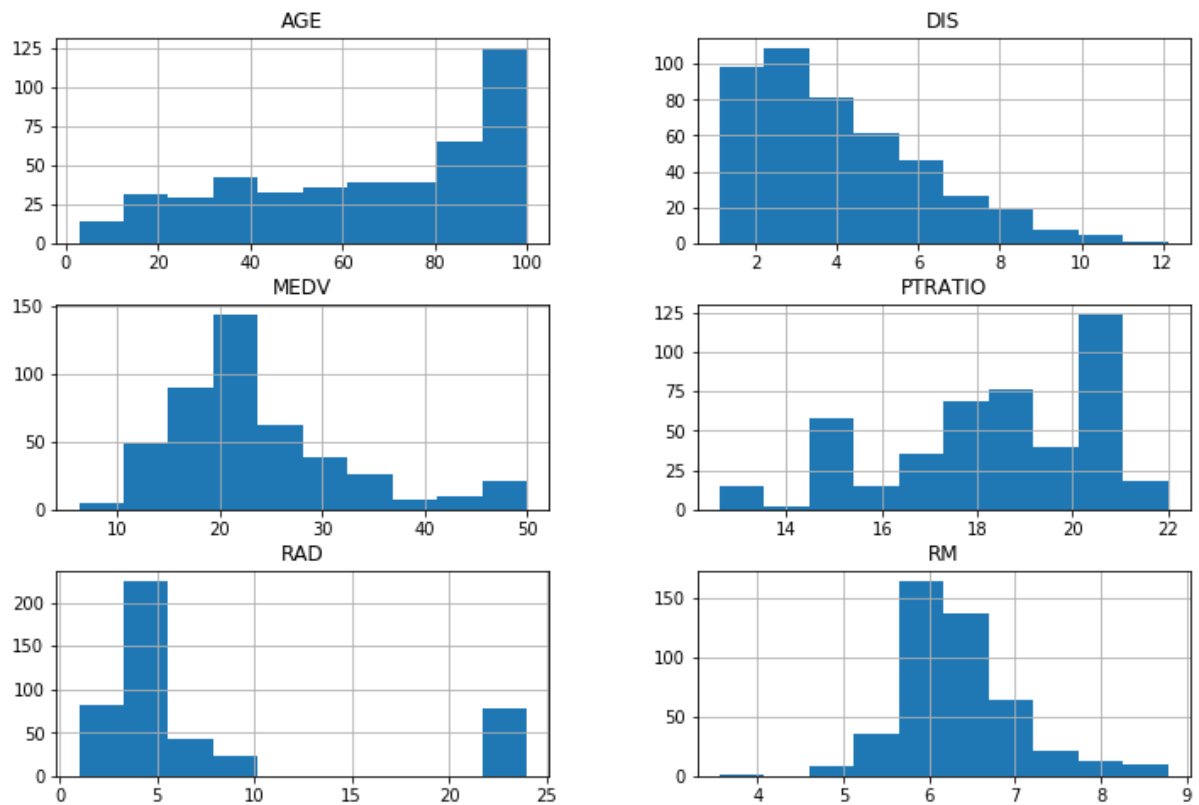
```
In [15]: #Summary Statistics
# Get the summary statistics of the numeric variables/attributes of the
dataset
print(df2.describe())
```

	RM	AGE	DIS	RAD	PTRATIO	
MEDV						
count	452.000000	452.000000	452.000000	452.000000	452.000000	452.000000
mean	6.343538	65.557965	4.043570	7.823009	18.247124	23.750442
std	0.666808	28.127025	2.090492	7.543494	2.200064	8.808602
min	3.561000	2.900000	1.129600	1.000000	12.600000	6.300000
25%	5.926750	40.950000	2.354750	4.000000	16.800000	18.500000
50%	6.229000	71.800000	3.550400	5.000000	18.600000	21.950000
75%	6.635000	91.625000	5.401100	7.000000	20.200000	26.600000
max	8.780000	100.000000	12.126500	24.000000	22.000000	50.000000

The above table gives the summary statistics for each of the columns. This includes number of records, mean, standard deviation, minimum 25 percentile 50 percentile 75 percentile and the maximum values of each column namely RM, AGE, DIS, RAD, PTRATIO and MEDV

Histograms

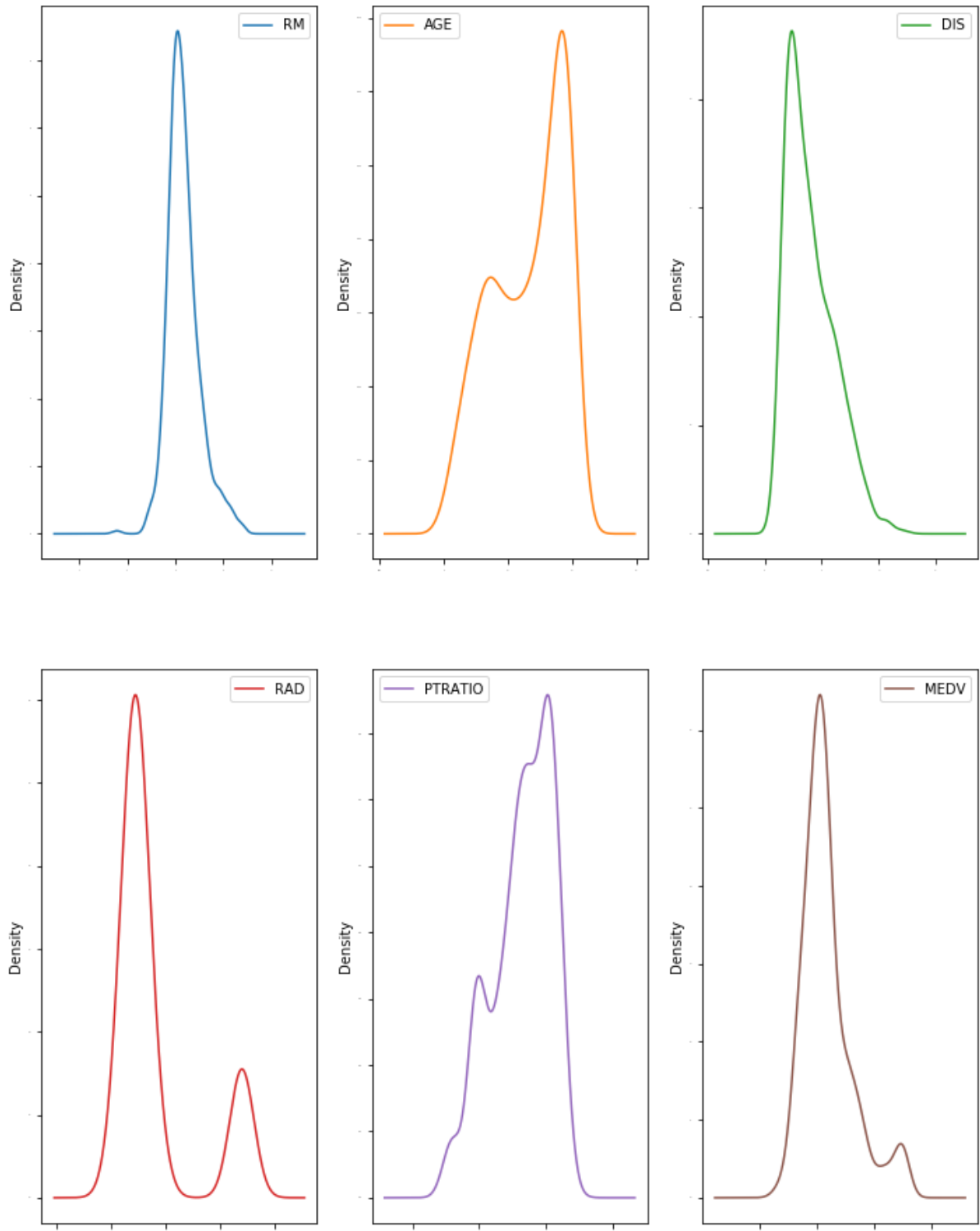
```
In [16]: # Plot histogram for each numeric  
df2.hist(figsize=(12, 8))  
plt.show()
```



MEDV, PTRATIO and RM has a distribution that is almost normal but the rest are not.

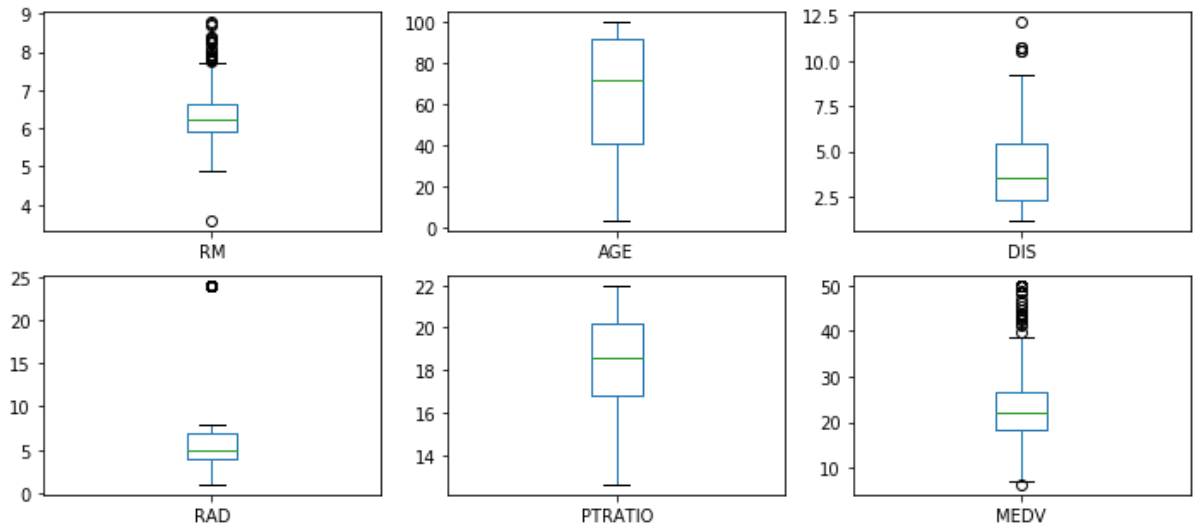
Density Plots

```
In [17]: # Density plots
# IMPORTANT NOTES: 5 numeric variables -> at Least 5 plots -> Layout (2,
3): 2 rows, each row with 3 plots
df2.plot(kind='density', subplots=True, layout=(2, 3), sharex=False, leg
end=True, fontsize=1,
figsize=(12, 16))
plt.show()
```



Consistent with the histogram, density plots also gives an almost normal distribution for MEDV, PTRATIO and RM. But DIS and AGE are bimodal.

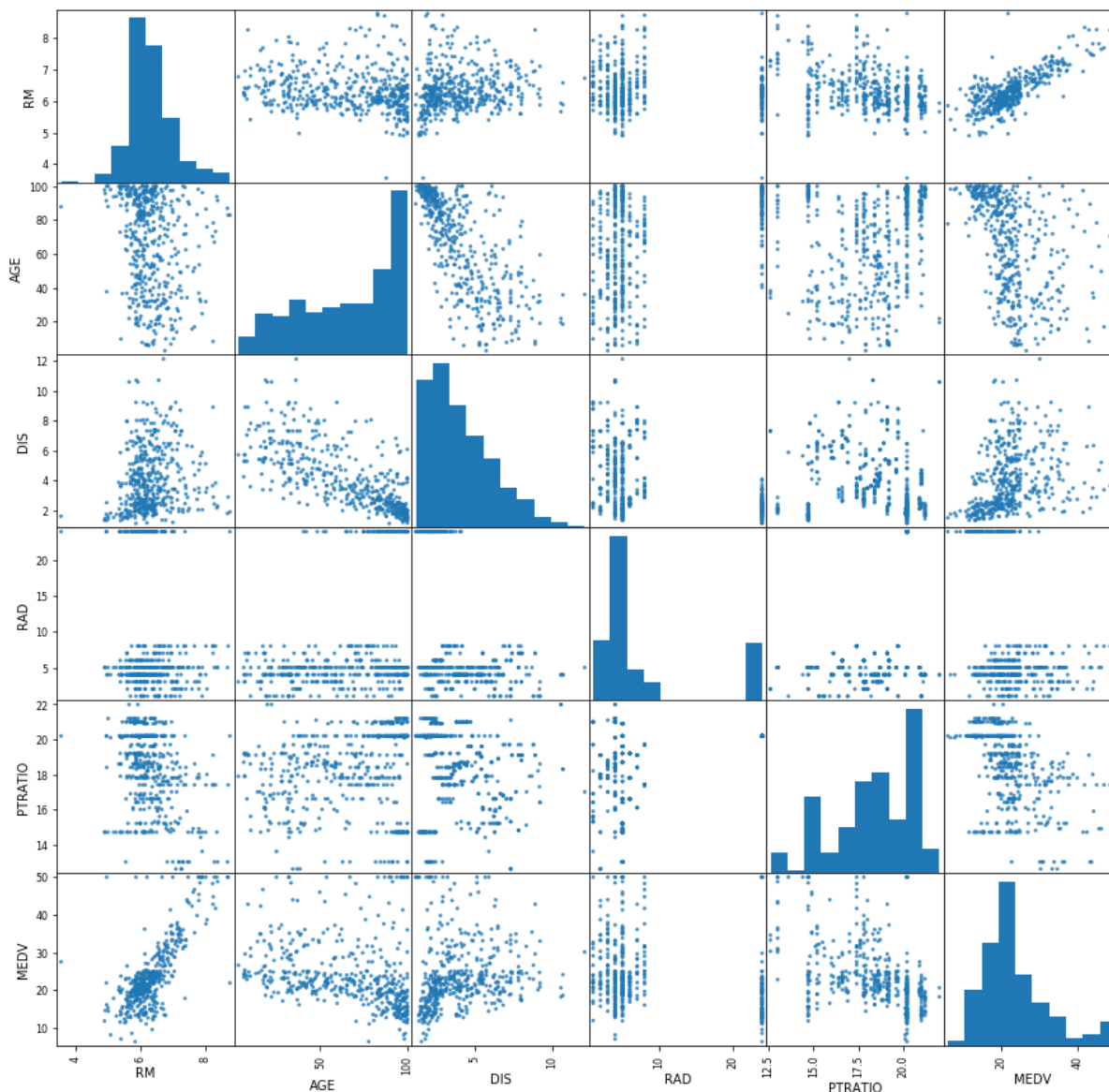
```
In [18]: df2.plot(kind='box', subplots=True, layout=(3,3), sharex=False, figsize=(12,8))
plt.show()
```



- RAD has a single outlier
- MEDV, RM and DIS has outliers at both ends
- Min value of RM is 1 and max at 9
- Min value of AGE is 4 and max just below 100
- Min value of DIS is 1 and max just below 12.5
- Min value of RAD is 1 and max just below 10 with an outlier at 25
- Min value of PTRATIO is 5 and max just below 22
- Min value of MEDV is 41 and max just above 50

Multivariate Analysis : Scatter Matrix


```
In [19]: # scatter plot matrix
scatter_matrix(df2, alpha=0.8, figsize=(15, 15))
plt.show()
```



- There is an evident linear relation between MEDV and RM
- MEDV and AGE is showing an inverse linear relation
- DIS and MEDV on the other hand has a direct linear relation

Separate Dataset into Input & Output NumPy Arrays

```
In [20]: # Store dataframe values into a numpy array
array = df2.values
# separate array into input and output components by slicing
# For X (input)[: , 5] --> all the rows, columns from 0 - 4 (5 - 1)
X = array[:,0:5]
# For Y (output)[: , 5] --> all the rows, column index 5 (Last column)
Y = array[:,5]
```

- Independent variables are 'RM', 'AGE', 'DIS', 'RAD' and 'PTRATIO' columns 0 to 4
- Dependent variable is 'MEDV' column 5

Split Input/Output Arrays into Training/Testing Datasets

```
In [21]: # Split the dataset --> training sub-dataset: 67%; test sub-dataset:
test_size = 0.33
# Selection of records to include in which sub-dataset must be done rand
omly
# use this seed for randomization
seed = 7
# Split the dataset (both input & output) into training/testing datasets
X_train, X_test, Y_train, Y_test = train_test_split(X, Y, test_size=test
_size,
random_state=seed)
```

- 33% of the data is test data
- 67% of the data is train data
- Split is based on random sampling

Build and Train the Model

```
In [22]: # Build the model
model = LinearRegression()
# Train the model using the training sub-dataset
model.fit(X_train, Y_train)
# Print out the coefficients and the intercept
# print intercept and coefficients

print ("Intercept",model.intercept_)
#print (model.coef_)

# If we want to print out the list of the coefficients with their corresponding variable name
# pair the feature names with the coefficients
names_2 = ['RM', 'AGE', 'DIS', 'RAD', 'PTRATIO']
coeffs_zip = zip(names_2, model.coef_)
# Convert iterator into set
coeffs = set(coeffs_zip)
#print (coeffs)
for coef in coeffs:
    print(coef, "\n")

Intercept -4.53602172360959
('DIS', -0.800386625412262)

('RAD', -0.14269109757375334)

('AGE', -0.08160589380401152)

('PTRATIO', -0.8641387261012642)

('RM', 8.457012650599342)
```

The model build by the algorithm is

- $MEDV = -4.53602172360959 + (8.45701265 \text{ RM}) - (0.08160589 \text{ AGE}) - (0.80038663 \text{ DIS}) - (0.1426911 \text{ RAD}) - (0.86413873 * \text{PTRATIO})$
 - As per the model for unit increase in RM, with all other features held constant MEDV increases 8.45701265
 - As per the model for unit increase in AGE, with all other features held constant MEDV decreases 0.08160589
 - As per the model for unit increase in DIS, with all other features held constant MEDV decreases 0.80038663
 - As per the model for unit increase in RAD, with all other features held constant MEDV decreases 0.1426911
 - As per the model for unit increase in PTRATIO, with all other features held constant MEDV decreases 0.86413873

Calculate R-Squared

```
In [23]: R_squared = model.score(X_test, Y_test)
         print(R_squared)
```

```
0.4921278466441376
```

The model has successfully explained 49 % of variance in MEDV using variables 'RM', 'AGE', 'DIS', 'RAD' and 'PTRATIO'

Prediction

We have trained the model. Let's use the trained model to predict the "Median value of owner-occupied homes in 1000 dollars" (MEDV)_

The suburb area has the following predictors:

- RM: average number of rooms per dwelling = 6.0
- AGE: proportion of owner-occupied units built prior to 1940 = 55
- DIS: weighted distances to five Boston employment centers = 5
- RAD: index of accessibility to radial highways = 2
- PTRATIO: pupil-teacher ratio by town = 16

```
In [24]: Predicted_Value = model.predict([[6.0, 55, 5, 2, 16]])
         print("Predicted Median value of owner-occupied homes in 1000 dollars for the above house is :", Predicted_Value[0], "Thousand" )
```

```
Predicted Median value of owner-occupied homes in 1000 dollars for the above house is : 23.60419508093679 Thousand
```

Evaluate/Validate Algorithm/Model Using K-Fold Cross-Validation

```
In [25]: # Evaluate the algorithm
# Specify the K-size
num_folds = 10
# Fix the random seed
# must use the same seed value so that the same subsets can be obtained
# for each time the process is repeated
seed = 7
# Split the whole data set into folds
kfold = KFold(n_splits=num_folds, random_state=seed)
# For Linear regression, we can use MSE (mean squared error) value
# to evaluate the model/algorithm
scoring = 'neg_mean_squared_error'
# Train the model and run K-fold cross-validation to validate/evaluate the model
results = cross_val_score(model, X, Y, cv=kfold, scoring=scoring)
# Print out the evaluation results
# Result: the average of all the results obtained from the k-fold cross-validation
print(results.mean())
```

-31.1777676938244

After we train we evaluate Use K-Fold to determine if the model is acceptable We pass the whole set because the system will divide for us -31 avg of all error (mean of square errors) this value would traditionally be positive value, but scikit reports as neg Square root would be between 5 and 6

With the test train split we calculated an R square of 49 but that is applicable to only that sample. Where as in K-fold validation with number of folds set to 10, we have 10 different accuracies which averages to 31. 17777. K-fold validation describes the population better than a predicting on single sample from population and that decreases the accuracy to 31.2.

Logistic Regression

Logistic regression is a classification technique.

Dataset

Predicted attribute: class of iris plant. Number of Instances: 150 (50 in each of three classes) Number of predictors: 4 numeric, predictive attributes and the class Attribute Information:

1. sepal length in cm
2. sepal width in cm
3. petal length in cm
4. petal width in cm
5. class: Iris Setosa, Iris Versicolour, Iris Virginica

Import necessary packages

```
In [26]: # Import scikit-Learn module for the algorithm/model: Logistic Regression
from sklearn.linear_model import LogisticRegression

# Import scikit-Learn module classification report to later use for information about how the system
# try to classify / label each record
from sklearn.metrics import classification_report
```

```
In [27]: #Load data into a pandas dataframe
# Already loaded into dfIris during visualization exercise
print(dfIris.head(5))
dfIris=dfIris.drop(['Id'], axis=1)
print("After dropping column Id \n",dfIris.head(5))
```

	Id	SepalLengthCm	SepalWidthCm	PetalLengthCm	PetalWidthCm	Species
0	1	5.1	3.5	1.4	0.2	Iris-setosa
1	2	4.9	3.0	1.4	0.2	Iris-setosa
2	3	4.7	3.2	1.3	0.2	Iris-setosa
3	4	4.6	3.1	1.5	0.2	Iris-setosa
4	5	5.0	3.6	1.4	0.2	Iris-setosa

After dropping column Id

	SepalLengthCm	SepalWidthCm	PetalLengthCm	PetalWidthCm	Species
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

Preprocess Dataset

Clean Data: Find & Mark Missing Values

- Zero values are invalid in these columns as petal and sepal width cannot have 0 as a valid value. If they exist, we need to mark them as missing value or numpy.NaN

```
In [28]: # mark zero values as missing or NaN
dfIris[['SepalLengthCm' , 'SepalWidthCm' , 'PetalLengthCm' , 'PetalWidthCm' ]] = dfIris[['SepalLengthCm' , 'SepalWidthCm' , 'PetalLengthCm' , 'PetalWidthCm' ]].replace(0,np.NaN)
# count the number of NaN values in each column
print (dfIris.isnull().sum())

SepalLengthCm    0
SepalWidthCm     0
PetalLengthCm    0
PetalWidthCm     0
Species          0
dtype: int64
```

There are no missing values in the dataset. Now we are ready to move to Exploratory data analysis

Exploratory Data Analysis

In this step we are getting to know our data. The number of columns and rows, what each column represent etc. The following are explored in this section:

- Shape of data set
- Type of each value in the dataset
- View a snippet of the dataset
- Descriptive statistics of each column
- Univariate plots
- Multivariate plots

```
In [29]: # get the dimensions or shape of the dataset
# i.e. number of records / rows X number of variables / columns
print(dfIris.shape)

(150, 5)
```

We have 5 variables and 150 rows as we have deleted the column labeled Id.

```
In [30]: #get the data types of all the variables / attributes in the data set
print(dfIris.dtypes)
#return the first five records / rows of the data set
print(dfIris.head(5))
```

```
SepalLengthCm    float64
SepalWidthCm     float64
PetalLengthCm    float64
PetalWidthCm     float64
Species          object
dtype: object
```

	SepalLengthCm	SepalWidthCm	PetalLengthCm	PetalWidthCm	Species
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

All the features are of type floating point except species which is object type.

```
In [31]: #return the summary statistics of the numeric variables / attributes in
         the data set
print(dfIris.describe())
```

	SepalLengthCm	SepalWidthCm	PetalLengthCm	PetalWidthCm
count	150.000000	150.000000	150.000000	150.000000
mean	5.843333	3.054000	3.758667	1.198667
std	0.828066	0.433594	1.764420	0.763161
min	4.300000	2.000000	1.000000	0.100000
25%	5.100000	2.800000	1.600000	0.300000
50%	5.800000	3.000000	4.350000	1.300000
75%	6.400000	3.300000	5.100000	1.800000
max	7.900000	4.400000	6.900000	2.500000

The above table gives the summary statistics for each of the numeric columns.

This includes number of records, mean, standard deviation, minimum 25 percentile 50 percentile 75 percentile and the maximum values of each column namely SepalLengthCm, SepalWidthCm, PetalLengthCm, PetalWidthCm

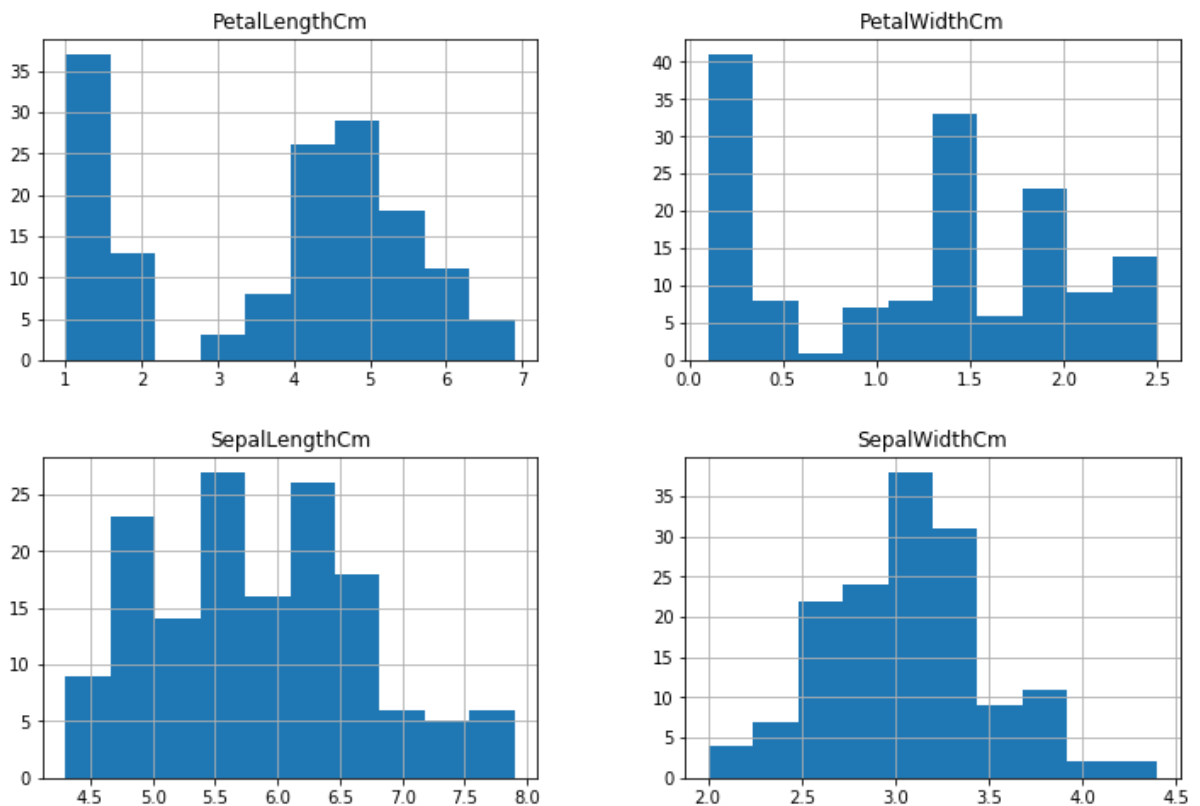
Univariate Plots

Plot each column to get an idea of the distribution. We will use the following

- Histogram
- Densityplot
- Boxplot

```
In [32]: #Histogram
dfIris.hist(figsize=(12,8))
plt.show
```

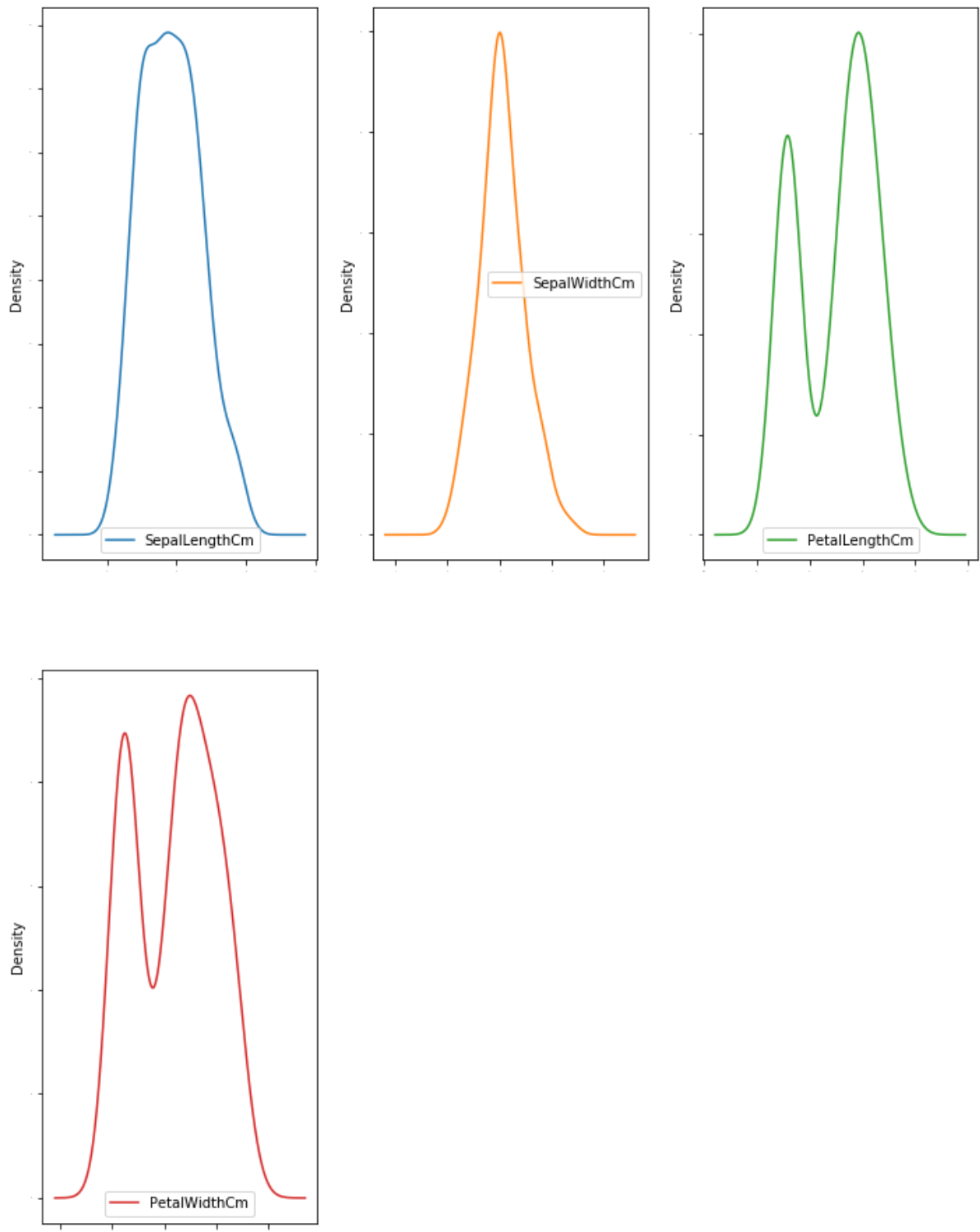
```
Out[32]: <function matplotlib.pyplot.show(*args, **kw)>
```



What Histogram Tells About Data

The histograms are displayed for each column of the data. Sepalength and sepalwidth shows almost normal distribution, petal length and width seems to be bimodal

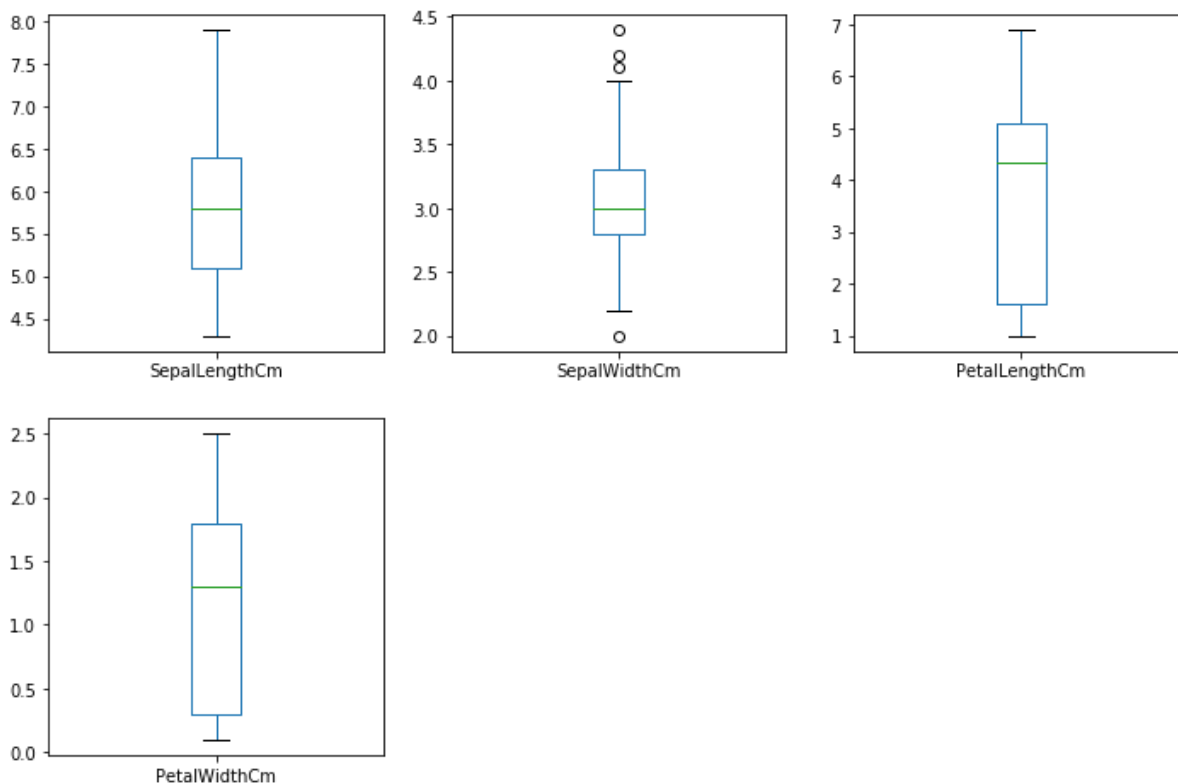
```
In [33]: #Density plot
dfIris.plot(kind="density", subplots = True, layout = (2,3), sharex = False,
legend = True, fontsize=1, figsize=(12,16) )
plt.show()
```



What Densityplots Tells About Data

The density plot is consistent with histogram and are displayed for each column of the data. Sepallength and sepalwidth shows almost normal distribution, petal length and width seems to be bimodal

```
In [34]: #Boxplot
dfIris.plot(kind="box", subplots = True, layout = (2,3), sharex = False,
sharey = False, figsize=(12,8) )
plt.show()
```



What Boxplots Tells About Data

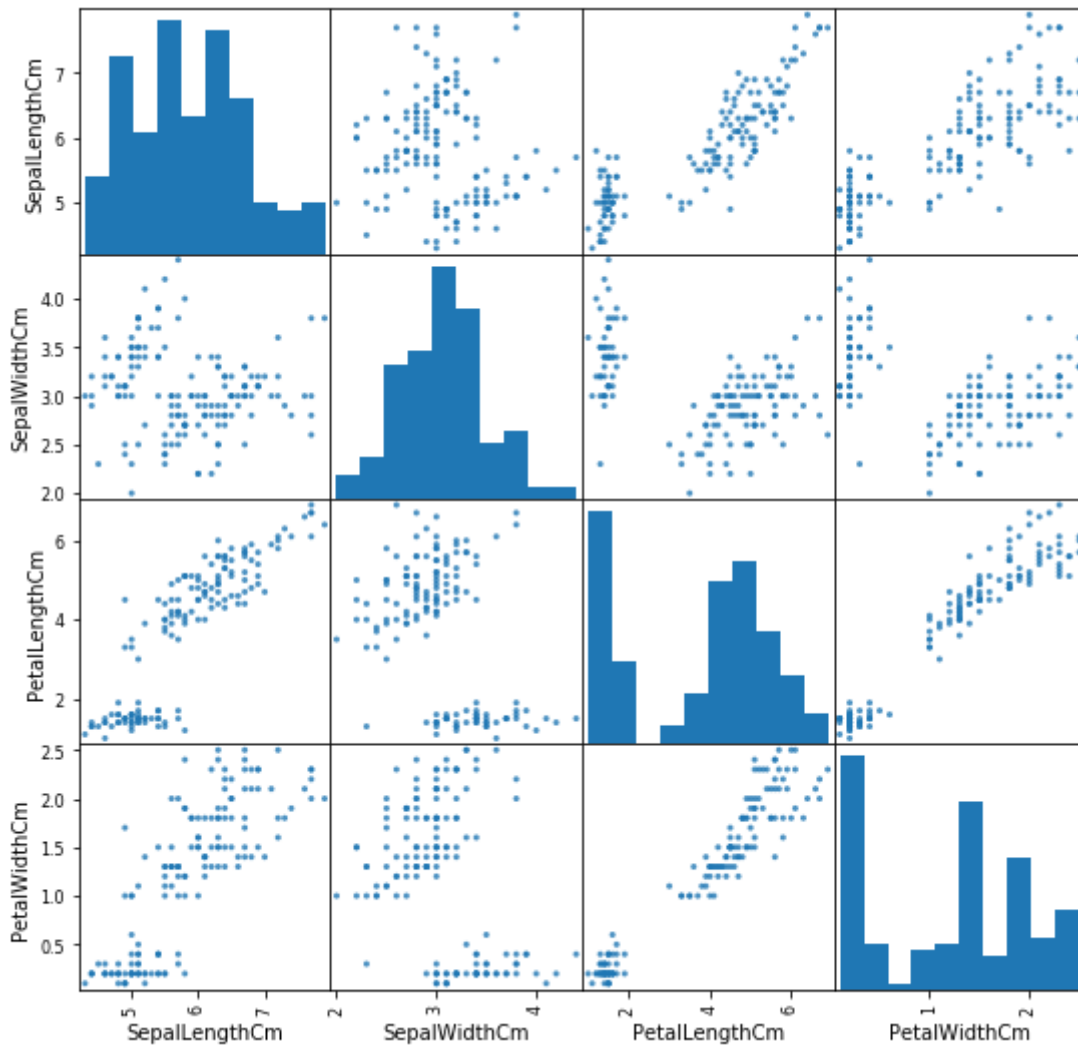
- For sepal length the values range from 4.5cm to 8cm
- For sepal width the values range from 2cm to 4.5cm
- For petal length the values range from 1cm to 7cm
- For petal width the values range from 0 to 2.5cm
- Sepalwidth has evident outliers

Multivariate Analysis

The above plots shows characteristics of each column like distribution of data, range, dispersion or independent features of each column. Now we would like to see how each feature is affected by other features and we use multivariate analysis for this. Scatter matrix is the best technique for this analysis.

In [35]: `#Scatter matrix`

```
scatter_matrix(dfIris, alpha=0.8, figsize=(9,9))
plt.show()
```



What scatter matrix tells us

As sepal length increase petal length and width increase Sepal width exhibits almost null correlation with other features and hence has data points scattered all over Petal length and petal width also shows a positive linear correlation

In [36]: `#class distribution i.e. how many records are in each class`
`print(dfIris.groupby('Species').size())`

```
Species
Iris-setosa      50
Iris-versicolor  50
Iris-virginica   50
dtype: int64
```

Each of the groups are equally represented in the data.

Separate Dataset into Input & Output NumPy arrays

Training data Set (X) consist of the independent variables or predictors Desired Output (Y) consist of the dependent variable or that which we are trying to predict

```
In [37]: # store dataframe values into a numpy array
array = dfIris.values
# separate array into input and output by slicing
# for X(input)[:, 1:4] --> all the rows, columns from 1 - 5 (6 - 1)
# these are the independent variables or predictors
X = array[:,0:4]
# for Y(input)[:, 4] --> all the rows, column 5
# this is the value we are trying to predict
Y = array[:,4]
```

Split Input/Output Arrays into Training/Testing Datasets

```
In [38]: # split the dataset --> training sub-dataset: 67%; test sub-dataset: 33%
test_size = 0.33
#selection of records to include in each data sub-dataset must be done r
randomly
seed = 7
#split the dataset (input and output) into training / test datasets
X_train, X_test, Y_train, Y_test = train_test_split(X, Y, test_size=test
_size,
random_state=seed)
```

- 33% of the data is test data
- 67% of the data is train data
- Split is based on random sampling

Build and Train the Model

```
In [39]: #build the model
model = LogisticRegression()
# train the model using the training sub-dataset
model.fit(X_train, Y_train)
#print the classification report
predicted = model.predict(X_test)
report = classification_report(Y_test, predicted)
print(report)
```

	precision	recall	f1-score	support
Iris-setosa	1.00	1.00	1.00	14
Iris-versicolor	0.93	0.78	0.85	18
Iris-virginica	0.81	0.94	0.87	18
accuracy			0.90	50
macro avg	0.91	0.91	0.91	50
weighted avg	0.91	0.90	0.90	50

```
/Users/jacobdenny/anaconda3/lib/python3.7/site-packages/sklearn/linear_model/logistic.py:432: FutureWarning: Default solver will be changed to 'lbfgs' in 0.22. Specify a solver to silence this warning.
```

```
FutureWarning)
```

```
/Users/jacobdenny/anaconda3/lib/python3.7/site-packages/sklearn/linear_model/logistic.py:469: FutureWarning: Default multi_class will be changed to 'auto' in 0.22. Specify the multi_class option to silence this warning.
```

```
"this warning.", FutureWarning)
```

A note on classification report - precision, recall, F1-score and support

Precision

The precision is the ratio $tp / (tp + fp)$ --> where tp is the number of true positives and fp the number of false positives The precision represents the ability of the classifier not to label a positive sample as negative

Recall

The recall is the ratio $tp / (tp + fn)$ --> where tp is the number of true positives and fn the number of false negatives The recall represents the ability of the classifier to find all the positive samples

F-beta

The F-beta score can be interpreted as a weighted harmonic mean of the precision and recall --> where an F-beta score reaches its best value at 1 and worst score at 0 The F-beta score weights recall more than precision by a factor of β $\beta == 1.0$ means recall and precision are equally important

Support

The support is the number of occurrences of each class in y_true

We have perfect prediction of the Iris-setosa

For Iris-versicolor, 93% Iris-versicolor were classified correctly but 7% were classified as Iris-versicolor wrongly

For Iris-virginica , 81% Iris-virginica were classified correctly but 19% were classified as Iris-virginica wrongly
94% Iris-virginica were classified correctly and 6% were misclassified

```
In [40]: #score the accuracy leve
result = model.score(X_test, Y_test)
#print out the results
print("Accuracy: %.3f%%" % (result*100.0))
```

Accuracy: 90.000%

The model correctly classifies 90% of the times

Classify/Predict Model

Now use the model to predict the group of a given flower with given measurements for different features. The new record has the following predictors:

- sepal length in cm = 5.3
- sepal Width in cm = 3.0
- petal length in cm = 4.5
- petal width in cm 1.5

```
In [41]: predicted=model.predict([[5.3, 3.0, 4.5, 1.5]])  
print("The model prediction for the given features is :",predicted[0])
```

The model prediction for the given features is : Iris-virginica

Based on the model's accuracy score: there is 90% chance that this new record is a Iris-virginica.

Evaluate the model using the 10-fold cross-validation technique.

- Specify the number of times of repeated splitting,
- Fix the random seed, Must use the same seed value so that the same subset can be obtained for each time the process is repeated,
- Split the whole data set into folds, for logistic regression
- Use accuracy level to evaluate the model/algorithm, train the model
- Run k-fold cross-validation to validate/evaluate the model,
- Print out the evaluation results


```
In [42]: # evaluate the algorithm
# specify the number of time of repeated splitting, in this case 10 folds
s
n_splits = 10
# fix the random seed
# must use the same seed value so that the same subsets can be obtained
# for each time the process is repeated
seed = 7
# split the whole dataset into folds
'''In k-fold cross-validation, the original sample is randomly partitioned into k equal sized subsamples. Of the k subsamples, a single subsample is retained as the validation data for testing the model, and the remaining k - 1 subsamples are used as training data. The crossvalidation process is then repeated k times, with each of the k subsamples used exactly once as the validation data. The k results can then be averaged to produce a single estimation. The advantage of this method over repeated random sub-sampling is that all observations are used for both training and validation, and each observation is used for validation exactly once
'''
kfold = KFold(n_splits, random_state=seed)
# for logistic regression, we can use the accuracy level to evaluate the model / algorithm
scoring = 'accuracy'
# train the model and run K-fold cross validation to validate / evaluate the model
results = cross_val_score(model, X, Y, cv=kfold, scoring=scoring)
# print the evaluation results
# result: the average of all the results obtained from the K-fold cross validation
print("Accuracy: %.3f (%.3f)" % (results.mean(), results.std()))
```

Accuracy: 0.880 (0.148)

```
/Users/jacobdenny/anaconda3/lib/python3.7/site-packages/sklearn/linear_model/logistic.py:432: FutureWarning: Default solver will be changed to 'lbfgs' in 0.22. Specify a solver to silence this warning.
FutureWarning)
/Users/jacobdenny/anaconda3/lib/python3.7/site-packages/sklearn/linear_model/logistic.py:469: FutureWarning: Default multi_class will be changed to 'auto' in 0.22. Specify the multi_class option to silence this warning.
"this warning.", FutureWarning)
/Users/jacobdenny/anaconda3/lib/python3.7/site-packages/sklearn/linear_model/logistic.py:432: FutureWarning: Default solver will be changed to 'lbfgs' in 0.22. Specify a solver to silence this warning.
FutureWarning)
/Users/jacobdenny/anaconda3/lib/python3.7/site-packages/sklearn/linear_model/logistic.py:469: FutureWarning: Default multi_class will be changed to 'auto' in 0.22. Specify the multi_class option to silence this warning.
"this warning.", FutureWarning)
/Users/jacobdenny/anaconda3/lib/python3.7/site-packages/sklearn/linear_model/logistic.py:432: FutureWarning: Default solver will be changed to 'lbfgs' in 0.22. Specify a solver to silence this warning.
FutureWarning)
/Users/jacobdenny/anaconda3/lib/python3.7/site-packages/sklearn/linear_model/logistic.py:469: FutureWarning: Default multi_class will be changed to 'auto' in 0.22. Specify the multi_class option to silence this warning.
"this warning.", FutureWarning)
/Users/jacobdenny/anaconda3/lib/python3.7/site-packages/sklearn/linear_model/logistic.py:432: FutureWarning: Default solver will be changed to 'lbfgs' in 0.22. Specify a solver to silence this warning.
FutureWarning)
/Users/jacobdenny/anaconda3/lib/python3.7/site-packages/sklearn/linear_model/logistic.py:469: FutureWarning: Default multi_class will be changed to 'auto' in 0.22. Specify the multi_class option to silence this warning.
"this warning.", FutureWarning)
/Users/jacobdenny/anaconda3/lib/python3.7/site-packages/sklearn/linear_model/logistic.py:432: FutureWarning: Default solver will be changed to 'lbfgs' in 0.22. Specify a solver to silence this warning.
FutureWarning)
/Users/jacobdenny/anaconda3/lib/python3.7/site-packages/sklearn/linear_model/logistic.py:469: FutureWarning: Default multi_class will be changed to 'auto' in 0.22. Specify the multi_class option to silence this warning.
"this warning.", FutureWarning)
/Users/jacobdenny/anaconda3/lib/python3.7/site-packages/sklearn/linear_model/logistic.py:432: FutureWarning: Default solver will be changed to 'lbfgs' in 0.22. Specify a solver to silence this warning.
FutureWarning)
/Users/jacobdenny/anaconda3/lib/python3.7/site-packages/sklearn/linear_model/logistic.py:469: FutureWarning: Default multi_class will be changed to 'auto' in 0.22. Specify the multi_class option to silence this warning.
"this warning.", FutureWarning)
/Users/jacobdenny/anaconda3/lib/python3.7/site-packages/sklearn/linear_model/logistic.py:432: FutureWarning: Default solver will be changed to 'lbfgs' in 0.22. Specify a solver to silence this warning.
```

```

FutureWarning)
/Users/jacobdenny/anaconda3/lib/python3.7/site-packages/sklearn/linear_
model/logistic.py:469: FutureWarning: Default multi_class will be chang
ed to 'auto' in 0.22. Specify the multi_class option to silence this wa
rning.
    "this warning.", FutureWarning)
/Users/jacobdenny/anaconda3/lib/python3.7/site-packages/sklearn/linear_
model/logistic.py:432: FutureWarning: Default solver will be changed to
'lbfgs' in 0.22. Specify a solver to silence this warning.
    FutureWarning)
/Users/jacobdenny/anaconda3/lib/python3.7/site-packages/sklearn/linear_
model/logistic.py:469: FutureWarning: Default multi_class will be chang
ed to 'auto' in 0.22. Specify the multi_class option to silence this wa
rning.
    "this warning.", FutureWarning)
/Users/jacobdenny/anaconda3/lib/python3.7/site-packages/sklearn/linear_
model/logistic.py:432: FutureWarning: Default solver will be changed to
'lbfgs' in 0.22. Specify a solver to silence this warning.
    FutureWarning)
/Users/jacobdenny/anaconda3/lib/python3.7/site-packages/sklearn/linear_
model/logistic.py:469: FutureWarning: Default multi_class will be chang
ed to 'auto' in 0.22. Specify the multi_class option to silence this wa
rning.
    "this warning.", FutureWarning)

```

- using the 10-fold cross-validation to evaluate the model / algorithm, the accuracy of this logistic regression model is 88%

Above, the model predicted the flower type of the new record as Iris-virginica. According to the model's accuracy score obtained from the model evaluation using 10-fold cross-validation there is only 88% chance that this new record is an Iris-virginica as opposed to 90% obtained by single sampling done earlier.

In []: