# Assignment : C1

- **Problem Statement :** Download the iris flower dataset into a Dataframe using Python /R and perform following operations :-
  - How many features are there and what are there types (eg. numeric, nominal )?
  - Compute and display summary statistics for each feature available in the dataset (eg. Minimum value, maximum value, mean, range, standard deviation)
  - Data Visualization - Create a histogram for each feature in the dataset to illustrate the feature distribution Plot each histogram.
  - Create a boxplot for each feature in the dataset. All of the boxplots should be combined into a single plot. Compare distribution and identify outliers.

- **Objective :**
  - To learn the concept and terminologies in data analytics.
  - To learn how to display summary statistics and data visualization.

- **Outcomes :** we will be able to :-
  - Learn concepts of Data Analytics.
  - Learn the concepts of statistics and data Visualization.

- Software / : OSS windows /ubuntu distribution.
  Hardware : Python libraries, Python framework,
  Requirements : R studio, Jupyter notebook, Anaconda
  Navigator.

- Theory : IRIS Dataset :
  The dataset is multivariate dataset introduced by
  the British statistician and biochemist
  Renold Finsher in 1936.
  Dataset consist of 50 sample from each of 3
  species of IRIS which are sentosa, virginia and
  versicolor. Four features were measured from each
  sample. The length and width of the sepal and
  petal in centimeters. Based on the combination of
  these four features, Fisher developed a linear
  discriminant model to distinguish the species from
  each other. The dataset contains a set of 150
  records under five attributes - sepal length,
  sepal width, Petal length, Petal width and species.
  The data is loaded in python as follows :

```
from sklearn datasets import load-iris
data = load-iris ()
df = Pd.DataFrame (data= data ['data'], columns=
                   data ['feature-names'])
```

To Display the features and more types :

```
print (list (df. columns))
x= df. drop (['target'], axis=1)
x. dtypes.
```

```
y = df ['target']
    y.value-counts ()
```

- Summary Statistics :
   i] Mean : The avarage value of set of values.

   $$\bar{x} = \frac{\sum x_i}{n}$$   $x_i$ = value of attributes, $n$ = total no. of items.

   df [" feature - name "] . mean ()

   ii] Range : The lowest and highest value in dataset.
       range = max - min.
       df ["feature-name"] . min )) , df ["feature-name"] . max ()

   iii] Standard Deviation : $\sigma = \sqrt{\dfrac{\sum (x_i - \bar{x})^2}{N}}$

   iv] Variance : $\sigma^2$.

- Data Visualization :
   Histogram : It is suitable for visualization of numeric data over a continuous interval or a certain time period. The histogram organize large amount of data and provides a visualization quickly using a single dimension. Histgram is a graphical display of data using bars of different heights. It is similar ta a Bar charts but a histogram groups number of ranges. The height of each bar shows how many falls into each range. df. hist.()
                          plt. show ()

2. Box plot : It allows quick graphical examination of one or more datasets. It may seems primitive than a histogram but they do have some advantages. They are usefull for comparing distribution between groups of data.

Box Plot using Jupyter Notebook :

Combined Box Plot : x.boxplot()

seperate for each feature : sns. boxplot (x=df ["target"], y=df ["feature-name"])

- Test Cases :

| | Description | Expected output | Result |
|---|---|---|---|
| i] | Feature name | sepal-length, sepal-width, petal-length, petal-width, dtype: object<br>target, dtype: int64 | Pass |
| ii] | Summary Statistics | count : 150 (all features)<br>sepal-length : mean : 5.8<br>Min : 4.3, Max : 7.9, std : 0.82<br>sepal-width : mean : 3.05, Min : 2.0<br>Max : 4.4, Std : 0.43,<br>Petal-length : Mean : 3.75, Min : 1.0<br>Max : 6.9, Std : 1.76<br>Petal width : Mean : 1.19,<br>Min : 0.10, Max : 2.50, Std : 0.76 | Pass |

| Summary Statistics | target : Mean : 1.0 Min : 0.0 , Max : 2.0 Std : 0.81 | Pass |
|---|---|---|

- Conclusion : Thus we successfully computed the given operation like Summary Statistics and Data Visualization on IRIS flower Dataset.