# LP1 Assignment DA R1

Summary statistics, data visualization and boxplot for the features on the Iris dataset

## Date - 9th August, 2020.

## Assignment Number - DA R1

## Title

Summary statistics, data visualization and boxplot for the features on the Iris dataset

## Problem Definition

Download the Iris flower dataset or any other dataset into a DataFrame. (eg https://archive. ics.uci.edu/ml/datasets/Iris) Use Python and Perform following:

- How many features are there and what are their types (e.g., numeric, nominal)?
- Compute and display summary statistics for each feature available in the dataset. (eg. minimum value, maximum value, mean, range, standard deviation, variance and percentiles
- Data Visualization-Create a histogram for each feature in the dataset to illustrate the feature distributions. Plot each histogram.
- Create a boxplot for each feature in the dataset. All of the boxplots should be combined into a single plot. Compare distributions and identify outliers.

## Learning Objectives

- Learn to use dataset, dataframes, features of dataset in an application
- Learn to compute summary statistics for the features.
- Learn to use visualization techniques.

## Learning Outcomes

I will be able to compute statistics on the features of the dataset, use histograms and boxplot on the features of the dataset.

# Software Packages and Hardware Apparatus Used

- Operating System : 64-bit Ubuntu 18.04
- Programming Language : Python 3
- Jupyter Notebook Environment : Google Colaboratory
- Python Libraries : Numpy, Pandas, MatPlotLib

# Related Mathematics

## Mathematical Model

Let S be the system set:

$S$ = {s; e;X; Y; Fme; Ff;DD;NDD; Fc; Sc}
       where Dataset is loaded into the dataframe

       s=start state
              Iris Dataset

       e=end state
              Summary statistics for each feature is computed.

       X=set of inputs
       X = {X1}
              Where X1 = IRIS Dataset
                     ● 5 Features
                            ○ 4 Numerical Feature
                            ○ 1 Nominal Feature
                     ● Data Count - 150

       Y=set of outputs
              1) Number of features
              2) Types of features
              3) Minimum value for each feature  in the dataset
              4) Maximum value for each feature  in the dataset
              5) Mean for each feature  in the dataset
              6) Range for each feature  in the dataset
              7) Standard deviation  for each feature  in the dataset
              8) Variance for each feature  in the dataset
              9) Percentiles for each feature  in the dataset

10) Histogram for each feature  in the dataset
11) Boxplot for each feature  in the dataset

Fme is the main function
It calls friend functions

Ff is the set of friend functions
Ff = {f1,f2,f3,f4,f5,f6}
where
f1 = function to load dataset into dataframe
f2 = function to to get number of features
f3 = function to get feature type
f4 = function to get minimum,maximum,mean,range,standard deviation,variance and percentile for each feature
f5 = function to draw histogram for each feature
f6 = function to draw boxplot for each feature

DD= Deterministic Data
IRIS dataset
- 5 Features
  - 4 Numerical Feature
  - 1 Nominal Feature
- Data Count - 150

NDD=Non-deterministic data
No non deterministic data

Fc =failure case:
No failure case identified for this application

# Concepts related Theory

Data analysis is a process of inspecting, cleansing, transforming, and modelling data with the goal of discovering useful information, informing conclusions, and supporting decision-making.Data analysis has multiple facets and approaches, encompassing diverse techniques under a variety of names, while being used in different business, science, and social science domains.

A data set (or dataset) is a collection of data. Most commonly a data set corresponds to the contents of a single database table, or a single statistical data matrix, where every column of the table represents a particular variable, and each row corresponds to a given member of the data set in question.

Mean, standard deviation, regression, sample size determination and hypothesis testing

are the fundamental data analytics methods.

## Mean

The sum of all the data entries divided by the number of entries.

$$\text{Population Mean: } \mu = \frac{\Sigma x}{N}$$

$$\text{Sample Mean: } \overline{x} = \frac{\Sigma x}{n}$$

## Range

The difference between the maximum and minimum data entries in the set.

**Range = (Max. data entry) – (Min. data entry)**

## Standard Deviation

The standard deviation measures variability and consistency of the sample or population. In most real-world applications, consistency is a great advantage.

$$\text{Population Standard Deviation} = \sigma = \sqrt{\frac{\Sigma(x-\mu)^2}{N}}$$

$$\text{Sample Standard Deviation} = s = \sqrt{\frac{\Sigma(x-\overline{x})^2}{n-1}}$$

## Variance

Variance is the average squared deviation from the mean.

## Percentile

Let p be any integer between 0 and 100. The pth percentile of the data set is the data value at which p percent of the value in the data set is less than or equal to this value.

# Steps for Execution

1. Download Iris Dataset
2. Open Google Colaboratory
3. Upload Iris Dataset to Google Colaboratory
4. Import Python Packages like Numpy, Pandas, MatPlotLib

5.  Get features from the Dataset into Pandas Dataframe
6.  Give Feature Names to Pandas Dataframe
7.  Get Feature Count and Type of Feature
8.  Compute Statistics in the Problem Definition
9.  Generate Histograms and Boxplots for features of the dataset

# Useful Python Functions

## Pandas Functions

### read_csv

Read Data from Iris Dataset downloaded and uploaded to Google Colab

### shape

Get (number of samples,number of features)

### iteritems

Iterate through features/columns of the dataset

### describe

Get Useful Statistics on the dataset like mean, max value, min value, standard deviation and percentiles

### var

Get Variance of features of the data set

### hist

Plot Histogram of features of the dataset

### max

Get the max value from all the features

### min

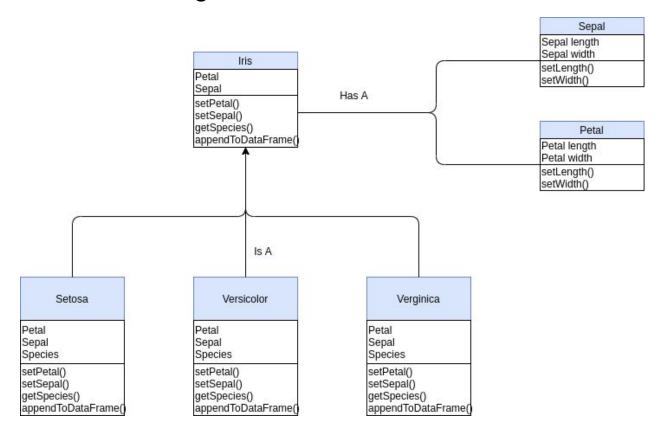Get the min value from all the features

plot

> Plot types of graph including box plot

## Matplot Library Functions

show

> Show the plotted graph using Matplot Library
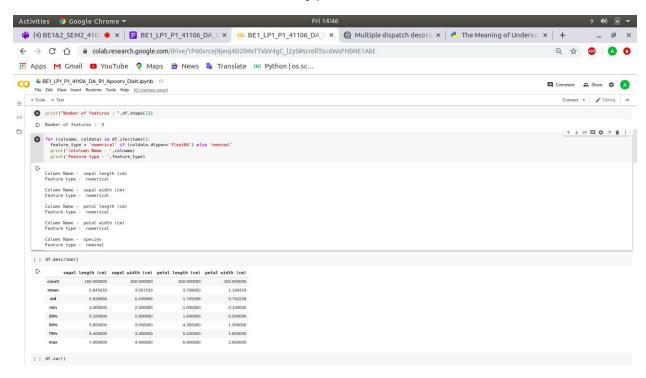
# Flowchart Design



# Footnotes

- As per the Python naming conventions, Use one leading underscore only for non-public methods, instance variables and Classes. In this assignment, class _Petal and _Sepal are treated as Non Public Classes. Private instance variables include _species of Classes Setosa, Virginica and Versicolor; and _length and _width of classes _Petal and _Sepal.
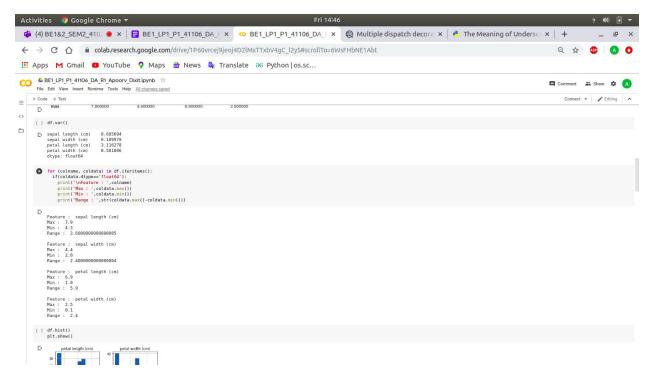
- Python does not support Method Overloading by default, To emulate method Overloading in Python, Programmer can do the following things (https://codippa.com/how-to-perform-method-overloading-in-python/)
  - Use Default Parameters
  - Modify program according to different number of parameters
- Python does not have the concept of Private Instance of Variables. To emulate the same, programmers use the concept of name mangling. In name mangling, a double underscore prefix is placed before variable names. It causes the Python interpreter to rewrite the attribute name in order to avoid naming conflict.
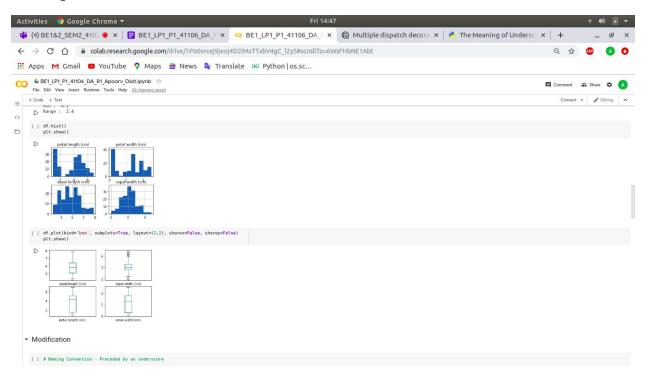
# Output Screenshots

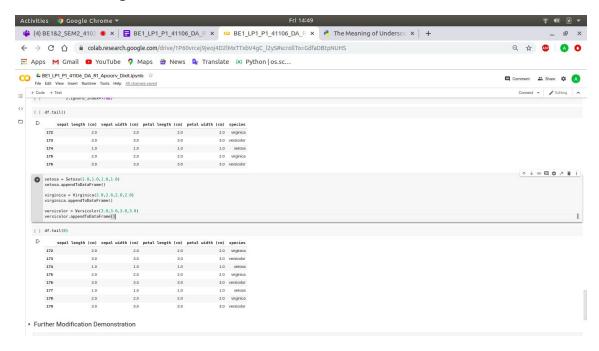## Number of Features, Feature Type and Statistics

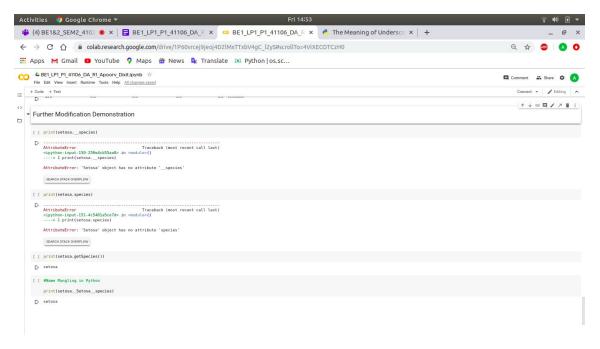# Variance and Range



# Histogram Plots and Box Plots

# Modification - AppendToDataFrame() function demonstration of Setosa, Virginica and Versicolor classes



# Modification - Demonstrating Name Mangling in Python

# Conclusion

We have successfully computed statistics on the features of the Iris dataset, and used histograms and boxplot on the features of the dataset.