

Name : Vikas Mane
Roll No: B1921152
PRN: 72000291F

Assignment1

Title: Download the Iris flower dataset or any other dataset into a Data Frame. (e.g., <https://archive.ics.uci.edu/ml/datasets/Iris>) Use Python/R 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. (e.g. 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.

Aim: Implement a dataset into a data frame. Implement the following operations:

1. Display data set details.
2. Calculate min, max, mean, range, standard deviation, variance.
3. Create histogram using hist function.
4. Create boxplot using boxplot function.

System Requirements: Python, Data Visualization libraries, Jupyter Notebook.

Theory:

Data analytics (DA): Data analytics (DA) is the *process of examining data sets in order to draw conclusions about the information they contain, increasingly with the aid of specialized systems and software*. Data analytics technologies and techniques are widely used in commercial industries to enable organizations to make more-informed business decisions and by scientists and researchers to verify or disprove scientific models, theories and hypotheses.

Data Science: Dealing with unstructured and structured data, Data Science is a field that comprises of everything that related to *data cleansing, preparation, and analysis*.

Data Science is the combination of statistics, mathematics, programming, problem-solving, capturing data in ingenious ways, the ability to look at things differently, and the activity of cleansing, preparing and aligning the data.

In simple terms, it is the umbrella of techniques used when trying to extract insights and information from data.

Big Data: Big Data refers to humongous volumes of data that cannot be processed effectively with the traditional applications that exist. The processing of Big Data begins with the raw data that isn't aggregated and is most often impossible to store in the memory of a single computer.

A buzzword that is used to describe immense volumes of data, both unstructured and structured, Big Data inundates a business on a day-to-day basis. Big Data is something that can be used to analyze insights which can lead to better decisions and strategic business moves.

The definition of Big Data, given by Gartner is, “Big data is high-volume, and high-velocity and/or high-variety information assets that demand cost-effective, innovative forms of information processing

Data Analytics: Data Analytics the science of examining raw data with the purpose of drawing conclusions about that information.

Data Analytics involves applying an algorithmic or mechanical process to derive insights. For example, running through a number of data sets to look for meaningful correlations between each other.

It is used in a number of industries to allow the organizations and companies to make better decisions as well as verify and disprove existing theories or models.

The focus of Data Analytics lies in inference, which is the process of deriving conclusions that are solely based on what the researcher already knows.

Applications of Data Analysis:

- **Healthcare:** The main challenge for hospitals with cost pressures tightens is to treat as many patients as they can efficiently, keeping in mind the improvement of the quality of care. Instrument and machine data is being used increasingly to track as well as optimize patient flow, treatment, and equipment used in the hospitals. It is estimated that there will be a 1% efficiency gain that could yield more than \$63 billion in the global healthcare savings.
- **Travel:** Data analytics is able to optimize the buying experience through the mobile/ weblog and the social media data analysis. Travel sights can gain insights into the customer's desires and preferences. Products can be up-sold by correlating the current sales to the subsequent browsing increase browse-to-buy conversions via customized packages and offers. Personalized travel recommendations can also be delivered by data analytics based on social media data.
- **Gaming:** Data Analytics helps in collecting data to optimize and spend within as well as across games. Game companies gain insight into the dislikes, the relationships, and the likes of the users.
- **Energy Management:** Most firms are using data analytics for energy management, including smart-grid management, energy optimization, energy distribution, and building automation in utility companies. The application here is centered on the controlling and monitoring of network devices, dispatch crews, and manage service outages. Utilities are given the ability to integrate millions of data points in the network performance and lets the engineers use the analytics to monitor the network.

Objective: To learn the concept of how to display summary statistics for each feature available in the dataset

Outcome: After completion of this assignment, it is expected to implement various statistics.

Input: Structured Dataset: Iris Dataset

Output:

1. DatasetDetails.
2. Min, Max, Mean, Variance value and Percentiles of probabilities
- 3 Histogram using HistFunction.
4. Boxplot using Boxplot Function.

Result / Conclusion: We learned to analyse data distribution, and data statistics for given dataset.

Program:

```
# Load libraries
import pandas
import matplotlib.pyplot as plt

# Load dataset
url = "https://archive.ics.uci.edu/ml/machine-learning-
databases/iris/iris.data"
names = ['sepal-length', 'sepal-width', 'petal-length', 'petal-width',
'class']
dataset = pandas.read_csv(url, names=names)

# shape
print(dataset.shape)

# head
print(dataset.head(20))

# descriptions
print(dataset.describe())

# class distribution
print(dataset.groupby('class').size())

# histograms
dataset.hist()
plt.show()

# box and whisker plots
dataset.plot(kind='box', subplots=True, layout=(2,2), sharex=False,
sharey=False)
plt.show()
```

Output:

