* Assignment DA1 *

* Title:-

Summary statistics, data visualization and boxplot for the features on the Iris dataset or any other dataset.

* Problem statements:-

Download the iris flower dataset into Dataframe. Use Python|R for perform following:

① How many features are there and what are their types (e.g. numeric, nominal etc.)

② Compute and display summary statistics for each feature available in dataset (eg min value, max value, mean, range, variance, standard deviation and percentiles.

③ Data visualization:- Create a histogram for each feature in the dataset to illustrate feature distributions. plot each histogram

④ Create a boxplot from each feature in dataset. All of the boxplots should be Combined single plot. Compare distributions and identify Outliers.

* Learning Objective:-

① Learn to use dataset, dataframes, features of dataset in an application.

② Learn to compute summary statistics for the features.

③ Learn to use visualization technique.

\* Learning Outcome:-
                    Students will be able to compute
the statistics on the features of dataset.
use histograms and boxplot on the
features of dataset.


\* Related Mathematics:-

Mathematical model
    Let s be the system set.
    $S = \{ s; e; x; y; fme; DD; NDP; FC; SC\}$
    where dataset is loaded into dataframe.
    S = start state
    e = end state   i.e. Summary statistics of
                        each feature is computed.
    x = set of inputs.
    $x = \{x1\}$
    where
    $X1$ = IRIS or any other dataset
    where,
    $Y$ = set of outputs.
    ① Number of features and their types.
    ② Summary statistics of each feature
    ③ Data visualization (histogram, boxplot)

    Fme is set of main functions.
    $fme = \{ f1, f2, f3\}$.
    where,
            $f1$ = function to load dataset into dataframe.
        $f2$ :- function to get number of features
        $f3$ :- function to draw histogram for each
                                            feature.
        $f4$ :- function to get feature type.
        $f5$ :- function to draw boxplot for each
                                            feature.

DD:- Deterministic Data
NDD:- Non-Deterministic data.
FC:- failure Case
No failure Case identified for application.

Theory:-

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

A dataset is Collection of data. Most Commonly dataset Correspond to the Contents of single database table, on single statistical data matrix where each Column of table represent Perticular variable.

Mean, standard deviation, variance size, min and max are the fundamentals of data analytics process.

* Mean = $\dfrac{\text{sum of data enteries.}}{\text{No. of data enteries.}}$

Population mean: $\mu = \dfrac{\sum x}{n}$

Sample mean: $\bar{x} = \dfrac{\sum x}{n}$

* **Range:-** Difference between max. and min. data entries in set.

$$Range = \begin{pmatrix} Max. \ data \\ entry \end{pmatrix} - \begin{pmatrix} Min \ data \\ entry \end{pmatrix}$$

* **Standard deviation:-** It measures variability and consistency of the sample on population in most real word application consistency is great advantage.

$$Population \ S.D. = \sigma = \sqrt{\frac{\sum(x-\mu)^2}{N}}$$

$$Sample \ S.D. = S = \sqrt{\frac{\sum(x-\bar{x})^2}{n-1}}$$

* **Variance:-** The average squared deviation from mean is variance.

* **Percentile:-** Let $p$ be any integer between 0 to 100. The $p^{th}$ percentile of dataset is data value at which $p$ percent of the value in dataset are less than or equal to this value.

Python Commands/Algorithm Steps:-

① Importing numpy for numerical calculation if required or pandas can be used. Import pandas for data access and manipulation. Import seaborn for plotting boxplot.

② Reading csv content in dataset named dataframe.

⇒ dataset = pd.read-csv(r, "iris.csv")

③ Printing shape i.e. Rows & Columns.

⇒ print (dataset.shape)

④ printing columns and their datatype.

⇒ dataset.dtypes.

⑤ To check column is numeric or nominal.

⇒ is_numeric-dtype(dataset['Column'])

⑥ Describe mean, median, min, max, percentile.

⇒ print (dataset.describe())

⑦ printing histograms.

⇒ dataset['Column'].plot.hist()

⑧ printing boxplots:-

⇒ sns.boxplot (data=dataset)

⑨ for comparing outliers.

⇒ sns. distplot (dataset['Column'], label="...",
                   color ='blue')

Test Cases:-

| | Expected Output | Actual Output | Result |
|---|---|---|---|
| ① | Import dataset | Iris. csv imported in dataframe | Pass. |
| ② | Printing no.of Rows & Column | using shape in python. No.of rows & Column printed (150,5) | Pass. |

| Expected output | Actual output | Result. |
|---|---|---|
| ③ Printing Columns & datatypes. | 5 Columns printed with datatype | Pass |
| ④ Check Column is numeric or not | All 4 Column for histogram are numeric | Pass. |
| ⑤ Printing histogram | 4 histogram for sepal length, sepal width, petal length, petal width printed | Pass. |
| ⑥ Printing boxplot | Box plot with all 4 Column printed showing min, max & percentile | Pass |
| ⑦ Comparing outlier | outliers are Compared by using dist plot | Pass. |

Conclusion:- Analysis of iris flower dataset is performed and shown using visualization technique such as histogram and boxplot.