

INDEX

Department of : Artificial Intelligence And Data Science

CLASS T.E

SR. NO.	NAME OF EXPERIMENT	Expt. Conducted on	Expt. Checked on	PAGE NO.	SIGN	REMARK
1.	Data-Wrangling - 1	2/2/23	5/2/23	1		
2.	Data-Wrangling - 2	9/2/23	12/2/23	3		
3.	Descriptive Statistics	16/2/23	17/2/23	6		
4.	Data Analytics - 1	23/2/23	25/2/23	12		
5.	Data Analytics - 2	21/3/23	23/3/23	15		
6.	Data Analytics - 3	9/3/23	11/3/23	19		
7.	Text Analytics	16/3/23	21/3/23	22		Revised
8.	Data Visualization - 1	23/3/23	27/3/23	28		
9.	Data Visualization - 2	13/4/23	15/4/23	31		

CERTIFICATE

This is to certify that Mr./Miss Derrat Yashraj Deepak of
Class T.E (A.I.D.S) Roll No. 223703 has satisfactorily completed the term work of the subject,
Software Laboratory-3 for VI Semester of 2022-23.

Date : / /

Staff Member

In-charge

H.O.D.

PRINCIPAL

Rahul

Shmp
P

B. S. Zade



INDEX

Department of : Artificial Intelligence And Data Science

CLASS T.E

SR. NO.	NAME OF EXPERIMENT	Expt. Conducted on	Expt. Checked on	PAGE NO.	SIGN	REMARK
10.	Data Visualization - 3	20/4/23	22/4/23	36		
11.	Simple Program Scala	27/4/23	30/4/23	40		<i>Folvo</i>
12.	GINA.	30/4/23	1/5/23	42		
13.	Mini Projects	2/5/23	8/5/23	46		

CERTIFICATE

This is to certify that Mr./Miss Devrat Yashraj Deepak _____ of
Class T.B (A.I.D.S) Roll No. 2237031 has Satisfactorily Completed the term work of the subject,
Software Laboratory 3 for VI Semester of 2022-2023

Date : / /

Folvo
Staff Member
In-charge

S
f H.O.D.

B-386
PRINCIPAL



Assignment No. 1

6/2/23

Title: Data Wrangling - 1

Objective: Perform the following operations using Python on any open source dataset (e.g. data.csv)

1. Import all the required Python libraries.
2. Locate open source data from web (e.g. <https://www.kaggle.com>) Provide clear description of the data and its source.
3. Load the dataset into pandas datframe.
4. Data Preprocessing : Check the missing values in the data using pandas isnull(), describe() function to get some statistics . Check dimension of data.
5. Data Formatting and Data Normalization : Summarize the types of variables by checking the data type (i.e character , numeric , integer , factor and logical)
6. Turn Categorical variables into quantitative variables using python.

Theory :

Data Wrangling : It is the process of gathering and collecting and transforming raw data into another format for better understanding , decision making , accessing and analysis in less time . Data Wrangling is also known as Data Munging .



Pandas: Pandas is a Python library used for working with datasets. It has function for analysing, cleaning, exploring and manipulating data.

Data Preprocessing:

1. Check for missing values in the data using pandas isnull(): The isnull() method returns a DataFrame object where all the values are replaced with a Boolean value True or NULL values, and otherwise False.

dataframe.isnull()

The notnull() method returns a DataFrame object where all the values are replaced with Boolean value True for NOT NULL values, otherwise False.

dataframe.notnull()

2. describe(): The describe() method returns description of the data in the DataFrame. If the DataFrame contain numerical data, the description contain these information for each column:

Counts - The number of non-empty values.

mean - average.

Std - Standard deviation



min - the minimum value

25% - The 25% percentile.

50% - The 50% percentile.

75% - The 75% percentile.

max - maximum value.

Q. Pandas.size, .shape & n.dim are used to return size, shape and dimension of data frames & series.

Syntax : `dataframe.size`

Returns : Return size of dataframe / series which is equivalent to total number of elements
This is ~~rows X columns~~.

Syntax : ~~`dataframe.shape`~~

Returns : Returns ~~Shape~~ (Rows, columns) of dataframe

Syntax : `dataframe.ndim`

Returns : Returns dimension of dataframe / series.
1 for one dimension (series), 2 for two dimension (dataframe)



Data Formatting & Normalization:

`type()` function is used to determine the type of data type.

`dtype()` function returns the data type of the array.

Turn Categorical variables into quantitative variables in Python:

Using `replace()` method

Replacing is one of the methods to convert categorical terms into numeric. For example, we will take a dataset of peoples salaries based on their level of education. This is an ordinal type of categorical variable. We will convert their education level into numeric terms.

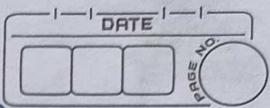
`replace()` replace = None, value = None, inplace = False, limit = None, regex = False, method = pad'

Method 2: using `get_dummies()`

Replacing the values is not the most efficient way to convert them. Pandas provide a method called `get_dummies` which will return the dummy variable columns.

(3)

16/2/23



Assignment No. 2

Title : Data Wrangling - 2

Objective : Create an 'Academic Performance' dataset of students and perform the following operations using python.

1. Scan all variables for missing value and inconsistencies. If there are missing values and/or inconsistencies, use any of the suitable techniques to deal with them.
2. Scan all numeric variables for outliers. If there are outliers, use any of the suitable techniques to deal with them.
3. Apply data transformation on at least one of the variables. The purpose of this transformation should be one of the following reasons:
to change the scale for better understanding of the variables, to convert non-linear relation into linear one, or to decrease the skewness and convert the distribution into a normal distribution.

Theory :

Missing Values (Data Cleaning) :



Data Cleaning can be applied to remove data noise and correct inconsistencies in data. Data cleaning routine work to "clean" the data by filling the missing values, smoothing noisy data, correct inconsistencies in data. If the data is dirty then it could lead to inaccurate results.

In pandas missing values are represented by `NaN`: Python singleton object used for missing data in Python.

`NaN`: Not a Number Special floating point value recognized by all systems that use the Standard IEEE floating point representation. Function for deleting, removing, replacing until values in Pandas DataFrame.

`isnull()`: returns True for NaN values

`notnull()`: returns False for NaN values

`dropna()`: delete all null values

`fillna()`: replaces NaN values with some value of their own.

`replace()`: Same as `fillna()`.

Outliers: It is a data object that deviates significantly from the rest from the data objects & behaves in a different manner an outlier is an object that deviates significantly from the rest of the objects.



* Handling of Outliers:

① Trimming / removing the outliers:
we ~~cannot~~ remove the outliers from the dataset.

② Quartile based floating & capping:

Outliers is capped at certain value above the 90th percentile value/floored at a factor below the 10th percentile value.

③ Mean / Median imputation:

As the mean value is highly influenced by the outlier it is advised to replace the outliers with the median value.

④ Data Transformation:

In data transformation, the data is transformed in consolidated into appropriate for mining.

Data transformation, the data involve the following

a. Smoothing:

It removes noise from the data such as techniques including regression & clustering. It helps in the predicting the patterns and collecting data.

b. Aggregation:

Data collection of aggregation is the method of storing and presenting data in summary format. The data may be obtained from multiple data sources.



to integrate data sources into a data analytics description.

c. Discretization:

It is a process of transforming continuous data into set of small variables, most data mining activities in the real world require continuous attributes.

d. Attribute Construction:

Where new attributes are created and applied to assist mining process from the given set of attributes. This simplifies the original data & makes mining more efficient.

e. Generalization:

To convert low level of converting all data variable to high level data attributes using concept hierarchy.

f. Normalization:

Data normalization involves converting all data variable into a given range. Techniques of normalization.

1. Min-Max

$$v' = \frac{v - \min_A}{\max_A - \min_A} (new.\max_A - new.\min_A) + new.\min_A$$

2. Z-Score:

This value /score helps to understand how far the data point from mean.

$$Z\text{-Score} = (\text{data point} - \text{mean}) / (\text{std. deviation})$$

3. Decimal Scaling:

In this technique, we move the decimal point of values of the attribute.

A value v_i of attribute A is normalized by the following formula.

Normalized value : $(v_i / 10^i)$
of attribute

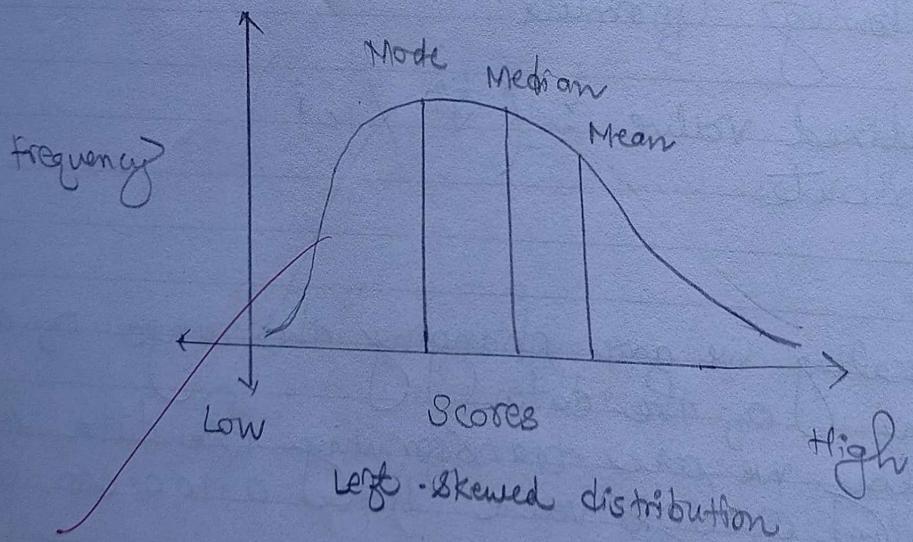
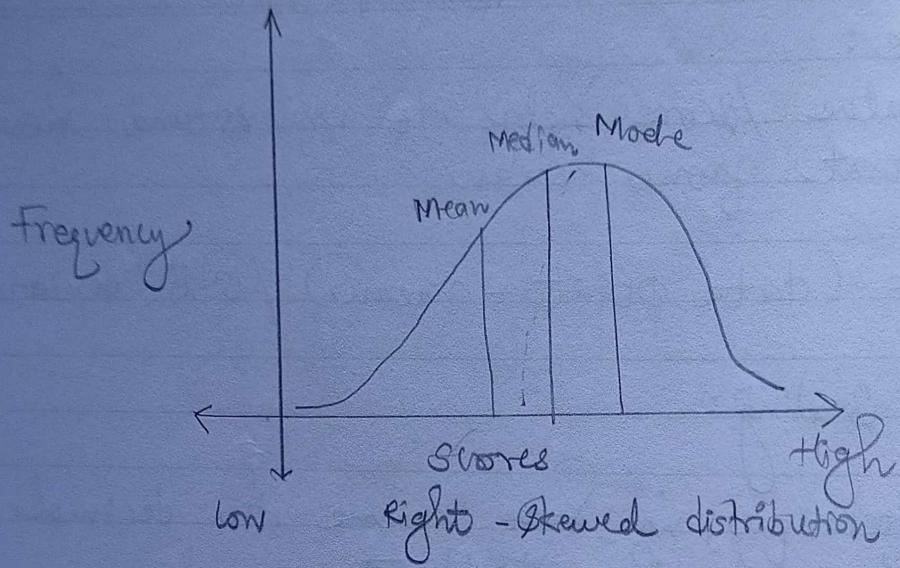
4. Scaling:

In scaling we are changing the range of the distribution of the data.

In scaling we are transforming the data so that it gets value specific value like 0-100 or 0-1.

Common types of scaling:

X _{new} :	$\frac{X_{old}}{X_{max}}$
--------------------	---------------------------





b. Min-Max Scaling :

$$X_{\text{new}} = \frac{X_{\text{old}} - X_{\text{min}}}{X_{\text{max}} - X_{\text{min}}}$$

This is the most popular they simple. Feature Scaling, it takes each value & subtracts the minimum & then divides by the range (max-min)

Skewness :

If the value of a specific independent variable are skewed depending on the model, Skewness may violate model assumption or may reduce the interpretation of feature importance.

"Skewness is a degree of asymmetry observed in a probability distribution that deviates from the symmetrical normal distribution in given set of data".

$$\text{Mean} = \text{Median} = \text{Mode} = 0$$

Types of Skewness :

1. Positive Skewed or right Skewed (Positive Skewness)

A positively skewed or right skewed distribution



Has a long right tail, it is a sort of distribution where unlike symmetrically distributed data where all measures of the central tendency (mean, mode, median) equal each other, with positively skewed data, the measures of the central tendency

Positively Skewed Distribution is a type of distribution where mean, mode, median of distribution are positive rather than negative or zero.

$$\text{Mean} > \text{Median} > \text{Mode}$$

2. Negative Skewed or left Skewed (Negative skewness)

A negatively skewed or left skewed distribution has long left tail. It is in the straight comes reverse of positively skewed distribution. In statistics, negatively skewed distribution. Plotting right side of graph and the tail of the distribution is spreading on the left side.

$$\text{Mode} > \text{Median} > \text{Mean}$$

Rel 3
16/2/23

Conclusion: Thus we implemented data wrangling 2 by applying data transformation on datasets also checking for missing values and handled it. Identifying outliers and handled with it using mentioned techniques.

(B)



Assignment No. 3

19/2/2023

Title : Descriptive Statistics - Measures of Central Tendency and Variability.

dim : Study the various measures of central tendency and variability with respect to the dataset.

Objectives : Perform the following operations on any dataset (e.g. data.csv).

① Provide summary statistics (mean, median, mode, minimum, maximum, Standard deviation) for a dataset (age, income etc) numeric variables grouped by one of the qualitative (categorical) variable.

e.g. Categorical Variable → age.

quantitative variable → income.

Provide summary statistics for income grouped by age. Create list of categorical variable.

② Write a python program to display some basic statistical details like percentile, mean, standard deviation etc. ~~Specified of Iris-Setosa, Iris-Versicolor & of Iris.csv dataset.~~

Theory: Central Tendency and Variation are two measures of descriptive summary used in Statistics to summarize data.

Measures of Central Tendency Shows the centre or middle of the dataset is located, whereas measures of variation Shows the dispersion among data values.

Measures of Central Tendency:

"The central tendency is used to measure ~~if~~ represent whole dataset using single value.

1. Mean : The mean is average of given observations is defined as the sum of all the observations divided by the total number of observations.

$$\bar{x} = \frac{\sum_{i=1}^n f_i x_i}{\sum_{i=1}^n f_i}$$

function to calculate mean : numpy .mean()
e.g. df.mean()

2. Median: The median is central or middlemost value of dataset.

~~Odd number of observations = $\left(\frac{n+1}{2}\right)^{th}$ observation.~~

~~Even number of observations = average of $\frac{n}{2}^{th}$~~

$\left(\frac{n}{2} + 1\right)^{\text{th}}$ observation.

$$\text{Median} = L + \left(\frac{\frac{n}{2} - C.F.}{f} \right) \times h$$

L - Lower limit of median class, h = class size
 C.F. - Cumulative frequency, f - frequency

e.g. df.median()

Q. Mode : The mode is the most repeated value of among the given set of data. Among the given set of observations, mode is the value which has the maximum frequency.

$$\text{Mode} = L + \left(\frac{f_1 - f_0}{2f_1 - f_0 - f_2} \right) \times h$$

f_1 is frequency of model class, h is class size.

f_0 is frequency of the class preceding model class.

f_2 is frequency of the class succeeding model class.

e.g. df.mode()



Measures of Variability / Dispersion :

A measure of dispersion is a summary statistic that represents the amount of dispersion in a dataset. How spread out are the values? High dispersion signifies that they tend to fall further away.

① Standard Deviation:

The standard deviation is the average amount of dispersion in your dataset.

$$\textcircled{1} \quad \sigma = \sqrt{\frac{\sum (X - \mu)^2}{N}} \quad \begin{array}{l} \sigma = \text{population} \\ \text{Standard deviation.} \end{array}$$

$$\sum = \text{Sum of}$$

μ = population mean , N = number of values in population.

Formula Explanation

$$\textcircled{2} \quad s = \sqrt{\frac{\sum (X - \bar{x})^2}{n-1}} \quad \begin{array}{l} s = \text{sample Standard} \\ \text{deviation.} \end{array}$$

$\sum = \text{Sum}$

X = each value

\bar{x} = sample mean.



e.g. df. Std()

2. Variance : The variance is the average of squared deviations from the mean. Variance is the square of Standard deviation.

Formula

Explanation

$$\textcircled{1} \quad \sigma^2 = \frac{\sum (x - \mu)^2}{N} \quad \sigma^2 = \text{Variance}$$

μ = mean

n = number of values in sample.

$$\textcircled{2} \quad s^2 = \frac{\sum (x - \bar{x})^2}{n-1}$$

e.g df.Var()

Percentile : In statistics, a percentile is a term that describes how a score compares to other scores from the same distribution.

$$P_x = \frac{x(n+1)}{100}$$



Function to calculate percentile: percentile()

e.g. df.percentile()

Groupby() Function in Python:

Pandas dataframe. groupby() function is used to split the data into groups, based on some criteria. Pandas objects can be split on any of their axes.

e.g. df.groupby(["car"]).mean()

Conclusion: Central Tendency and Variation are two measures used in Statistics to summarize data. Measures of central tendency shows where the centre or middle of the data is located, whereas measure of variations shows the dispersion among data values.

Feb 3
2/2/23

Assignment No. 4

2/3/2023

Title: Data Analytics- 1

Objective : To predict the value of prices of the houses using given features. Creates a linear regression model using Python to predict home prices using Boston Housing Dataset ! {
<https://www.kaggle.com/c/boston-housing>}.

The Boston Housing dataset contains information about various houses in Boston through different parameters. There are 506 samples & 14 feature variables in this dataset.

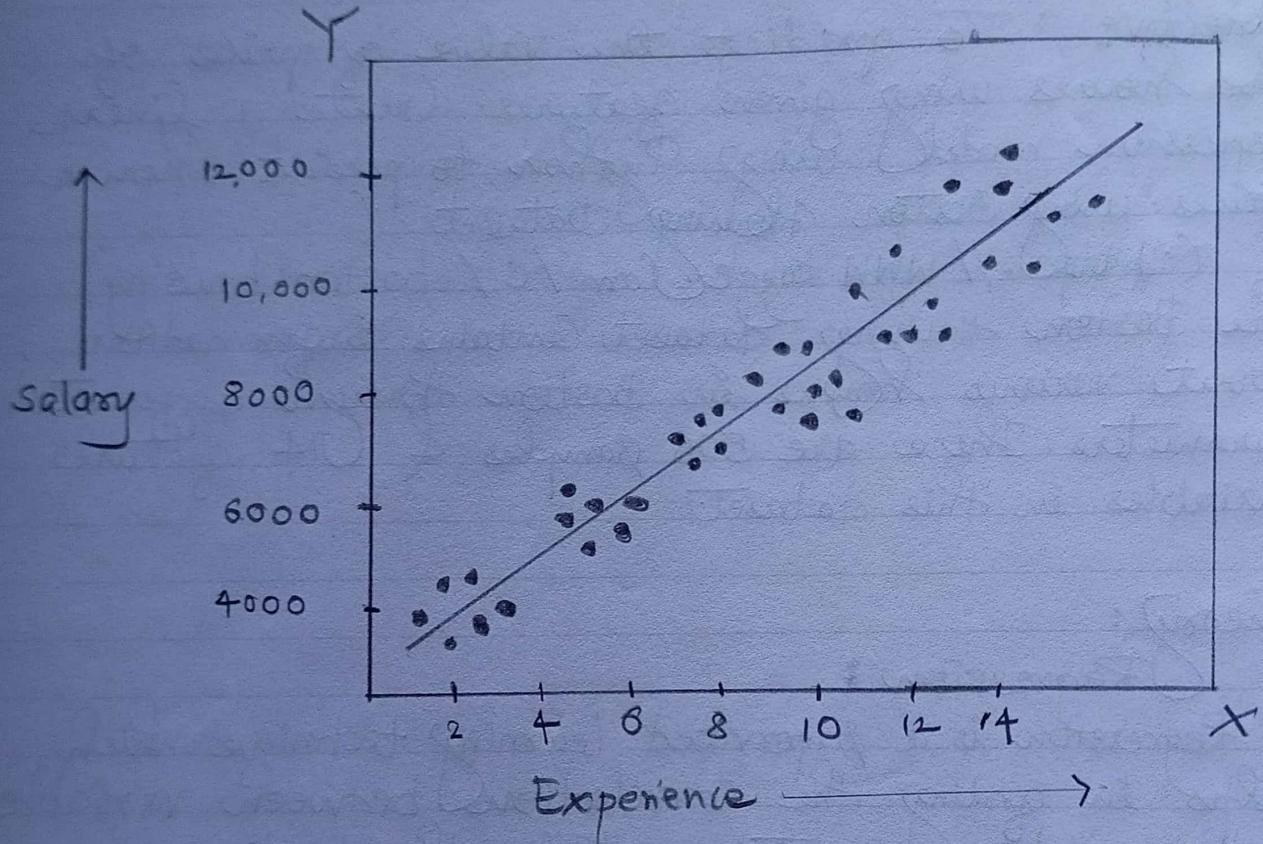
Theory :

Regression :

Regression is a supervised learning technique which helps in finding the correlations between variables and enables us to predict the continuous output variable based on the one or more predictor variables. It is mainly used for prediction, forecasting, time series modeling, and determining the causal-effect relationship between variables.

What is Linear Regression ?

The objective of a linear regression model to find a relationship between a one or more features (independent variables) and a continuous target variables (dependent variables). When there is only feature it is called Uni-variate / Simple



$$\text{Equation: } Y = aX + b$$

Here, Y = dependent variables (target variables)
 X = Independent Variables (predictor variables)
 a, b are the linear coefficients.



Linear Regression and if there are multiple features it is called Multiple Linear Regression.

Linear Regression is a statistical regression method which is used for predictive analysis.

It is one of the very simple and easy algorithms which is used for predictive analysis and shows a relationship between the continuous variables.

- The relationship between variables in the linear regression model can be explained using the below image. Here we are predicting the salary of an employee on the basis of the year of experience.

Where can be Linear Regression be used?

It is very powerful technique and can be used on to understand the factors the influence profitability. It can be used to forecast sales in the coming months by analyzing the sales data for previous months. It can also be used to gain various insights about customer behaviour. By the end of the day we will build a model which looks like the below picture i.e determine a line which best fits the data.

Some popular applications of linear regression are:

- Analyzing trends and sales estimates.
- Sales forecasting



- Real estate prediction.
- Arriving at ETA's in traffic.

Evaluation:

1. RMSE : It measures the average difference between values predicted by a model and actual values. It provides an estimation of how well the model is able to predict the target value (accuracy). The lower the value of the Root Mean Squared Error, the better the model is.

2. MSE : The Mean Squared Error measures how close regression line is to a set of data points. It is risk function corresponding to the expected value of the squared error loss. Mean Squared Error is calculated by taking the average, specifically the mean of errors squared from data as it relates to a function. Lesser the MSE \Rightarrow Smaller is the error \Rightarrow Better the estimator.

3. AE = Absolute error refers to the magnitude of difference the prediction of an observation and true value that of observation.

Conclusion: Thus, we performed, data analytics, by creating linear regression model to predict home prices. We have also evaluated performance of the model.

Feb
6/3/23

Assignment No.5

Title : Data Analytics 2.

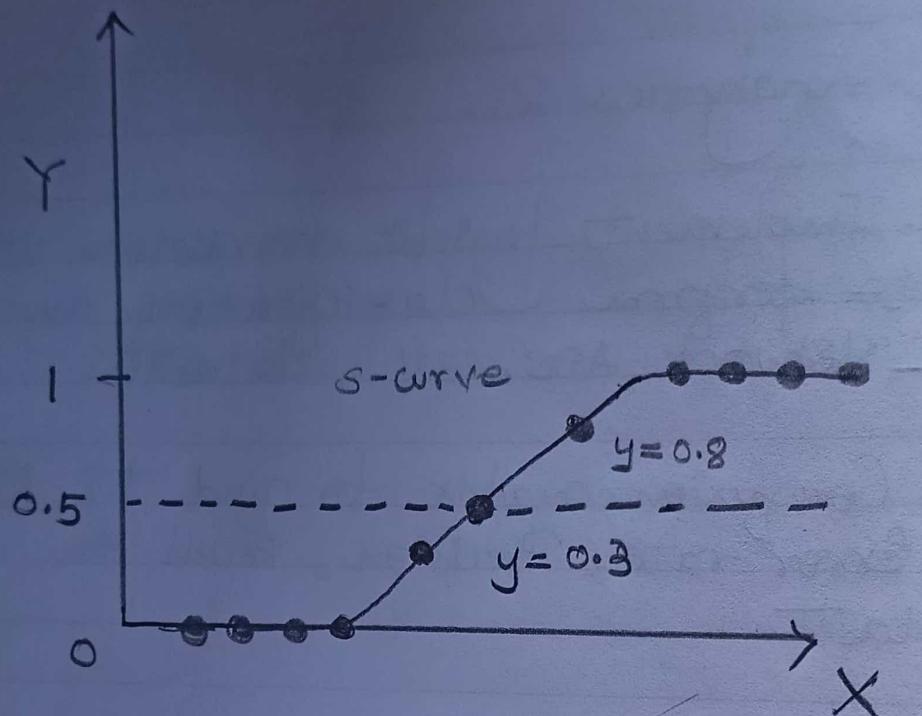
Objective : 1. Implement logistic regression using Python to perform classification on Social Network Ads.csv dataset.

2. Compute Confusion matrix to find TP, FP, TN, FN Accuracy, Error rate, Precision, Recall on the given dataset.

Theory:

What is Logistic Regression?

- It is used for predicting the categorical dependent variable using a given set of independent variables.
- Logistic Regression predicts the output of a categorical dependent variable. Therefore the outcome must be a categorical or discrete value. It can be either Yes or No, 0 or 1, true or false. But instead of giving the exact value as 0 or 1, it gives the probabilistic values which lies between 0 & 1.
- In Logistic Regression, instead of getting a straight regression line, we fit an "S" shaped logistic function which predicts two maximum values (0 or 1).
- In Logistic Regression, indicates curve the



logistic Function



likelihood of something such as whether the cells are cancerous or not, a mouse is obese or not based on its weight, etc.

- Logistic Regression is significant machine learning algorithm because it has the ability to provide probabilities and classifying the observation new data using continuous and discrete datasets.
- Logistic Regression can be used to classify the observation using different types and can easily determine the most effective variables used for the classification. The below image is showing the logistic function.

The Logistic Regression equation can be obtained from the Linear Regression equation. The mathematical steps to get Logistic Regression equations are given below:

We know the equation of the straight line can be written as:

$$y = b_0 + b_1 x_1 + b_2 x_2 + b_3 x_3 + \dots + b_n x_n$$

In Logistic Regression y can be between 0 & 1 only so for this let's divide the above equation by $(1-y)$:

$$\frac{y}{1-y}; 0 \text{ for } y=0; \text{ and } \infty \text{ for } y=1$$

① Accuracy : $\frac{TN + TP}{TN + TP + FN + FP}$

② Error rate : $1 - \text{Accuracy} = \frac{FN + FP}{TN + TP + FN + FP}$

③ Precision : $\frac{TP}{FP + TP}$

④ Recall : $\frac{TP}{TP + FN}$

⑤ F₁ score : $\frac{2 \times (\text{Recall} \times \text{Precision})}{(\text{Recall} + \text{precision})}$



But we need range between $(-\infty \text{ to } +\infty)$, then take algorithm of the equation. It will become

$$\log \left[\frac{y}{1-y} \right] = b_0 + b_1 x_1 + b_2 x_2 + b_3 x_3 + \dots + b_n x_n$$

The above equation is the final equation of Logistic Regression.

- Accuracy: It defines the ratio of correctly predicted with all parameters.
- TN : It stands for True Negative. The case is negative & predicted positive.
- FN : (False negative) the case was true predicted negative and positive.
- FP : (False positive) the case was negative but predicted positive.
- TP : (True Positive) Case was positive & predicted positive.
- Precision : Accuracy of positive predictions.



• Recall : It is the ability of a classifier to find all positive instances . It is defined as ratio of True positives to the sum of true positives & false negatives.

• F1-score : It is harmonic mean of precision & recall such that best score is 1 & worst is 0.0 . F1 scores are less than accuracy measure as they embed precision & recall into their computation.



Conclusion : Thus we performed data analytics by creating logistic regression model and computed confusion matrix to find TP, FP, TN, FN, accuracy, error rate, Precision, Recall.

F1 score

$$\frac{2}{3}$$
$$6 \sqrt{3} / 23$$



Assignment No. 5:

Title: Data Analytics 5.

Objective : 1. Implement simple Naive Bayes classification algorithm using Python on Iris dataset

2. Compute Confusion matrix to find TP, FP, TN, FN, Accuracy, Error rate, Precision, Recall on the given dataset.

Theory :

Naive Bayes Classifier Algorithm :

- Naive Bayes algorithm is a supervised learning algorithm, which is based on Bayes theorem and used for solving classification problems.
- It is mainly used in text classification that include a high-dimensional training set.
- It is a probabilistic classifier, which means it predicts on the basis of the probability of object.
- Some popular examples of Naive Bayes Algorithm are spam filtrations, sentimental analysis, and classifying articles.

Bayes Theorem :

Bayes theorem is also known as Bayes law; which is used to determine the probability of a hypothesis with prior knowledge. It depends on the conditional probability probability.

The formula for Bayes theorem :

$$P(A|B) = \frac{P(B|A) P(A)}{P(B)}$$

where,

$P(A|B)$ is posterior probability : Probability of Hypothesis A on the observed event B.

$P(B|A)$ is likelihood probability : Probability of the evidence given that the probability of a hypothesis is true.

$P(A)$ is Prior Probability Probability of Hypothesis before observing the evidence.

$P(B)$ is Marginal Probability : Probability of the Evidence.

Types of Naive Bayes Model:

There are types of Naive Bayes model, which is given below:

- Gaussian: The Gaussian model assumes that the features follow a normal distribution. This means if predictions take continuous values instead of discrete, then the model assumes that these values are sampled from the Gaussian Distribution.
- Multinomial: The multinomial Naive Bayes classifier is used when data is multinomial distributed. It is primarily used for the document classification problems, it means a particular document belongs to which category such as sports, politics, education.

The classifier uses for frequency of words the predictor

- Bernoulli: The Bernoulli classifier works similar to the multinomial classifier, but the predictor variables are the independent Boolean variables. Such as if a particular word is present or not in a document. This model is also famous for document classification tasks.

Steps to implement:

- Data Pre-processing step.
- Fitting Naive Bayes to the Training set.



- predicting the test result
- Test accuracy of the result (Creation of confusion matrix)
- Visualizing the test set result.

Conclusion: Thus, we performed, data analytics by creating Naive Bayes classification on iris dataset and computed confusion matrix to find TP, FP, TN, FN, Accuracy, Error rate, Per Precision, Recall.

~~Robin~~
29/3/23

Assignment No. 7

Title: Text Analytics

Objective : 1. Extract, Sample, document and apply following document pre-processing methods:

Tokenization, POS Tagging, Stop words removal, Stemming, Lemmatization.

2. Create representation of documents by calculating Term Frequency and Inverse Document Frequency.

Theory : Tokenization : In this step, the text is split into smaller units.

Output : Sentences are tokenized into words.

POS Tagging : It is also called grammatical tagging is the process of marking up a word in text as corresponding to a particular part of speech, based on both its definition and its context.

~~Stop words removal~~ : Stopwords are the commonly used words and are removed from the text as they do not add any value to the analysis.
Eg. [i, me, my, myself, we, our, etc.]



Stemming: It is also known as the text standard step where the words are stemmed or diminished to their root / base form. For example, words like 'programmer', 'programming', 'program' will be stemmed to 'program'.

Lemmatization: It stems the word but makes sure that it does not lose its meaning. It has a predefined dictionary that stores the context of words in the dictionary while diminishing.

Introduction TF-IDF

TF-IDF stands for 'Term Frequency - Inverse Document Frequency'. This is a technique to quantify words in a set of documents. We generally compute a score for each word to signify its importance in the document and corpus.

Terminology t - term (word) d - document (set of words) N - Count of Corpus - the total documents Set

Term Frequency: This measures the importance of the documents in a whole set of corpus. This is very similar to TF but the only difference is that TF is the frequency counter for a term t in document d, whereas DF is the count of term t in the document set N.

$$df(t) = \text{Occurrence of } t \text{ in } N \text{ documents.}$$



Inverse Document Frequency: IDF is the inverse of the document frequency which measures the informativeness of term t . When we calculate IDF, it will be very low for the most occurring words such as stop words.

$$idf(t) = N / df$$

Finally, by taking a multiplicative value of TF and IDF we, get the TF-IDF score. There are many different variations of TF-IDF but for now, let us concentrate on this basic version.

$$tf-idf(t, d) = tf(t, d) * \log(N / (df + 1))$$

Conclusion: Thus, we perform text analytics by initially performing pre-processing on the text and then applying TF-IDF on it.

Feb 2
27/4/23.



Assignment No. 8

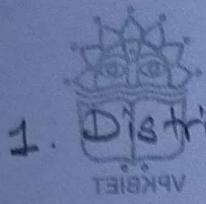
Title: Data Visualization I.

Objective : 1. Use the inbuilt dataset 'titanic'. The dataset contains 891 rows and contains information about the passengers who boarded the unfortunate Titanic ship. Use the Seaborn library to see if we can find any patterns in data.

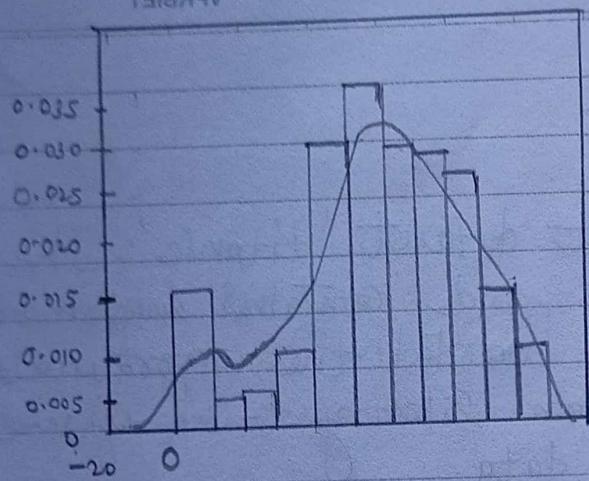
2. Write a code to check how the price of the ticket (column name: 'Fare') for each passenger is distributed by plotting a histogram.

Theory : Data visualization is the graphical representation of information and data. By using visual elements like charts, graphs and maps, data visualization tools provide an accessible way to see and understand trends, outliers, and patterns in data. In the world of big data, and data visualization tools and technologies are essential to analyze massive amounts of information and make data driven decisions.

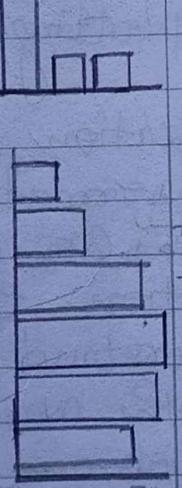
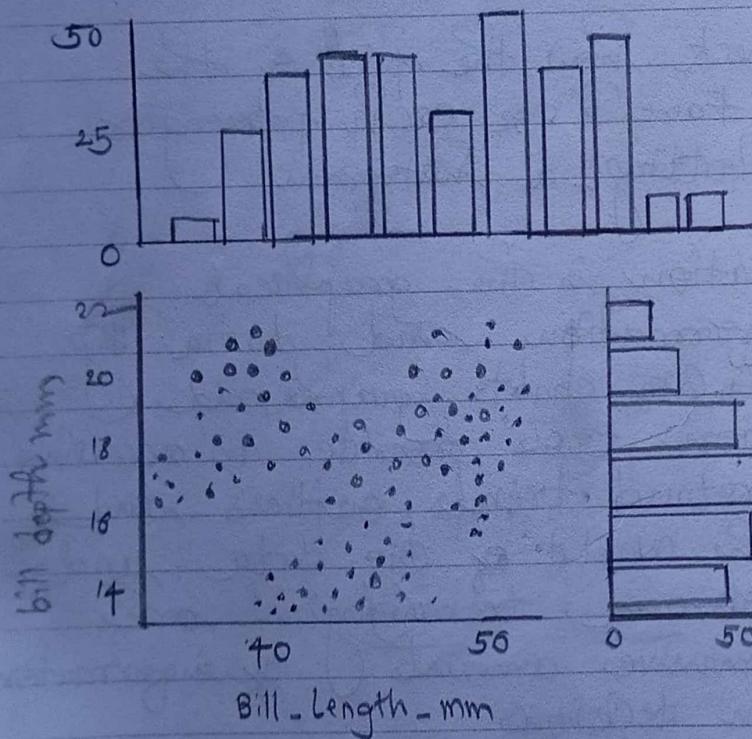
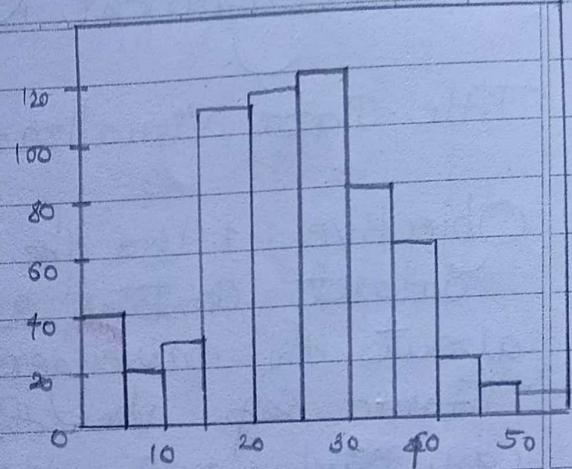
~~The Seaborn library is built on top of Matplotlib and offers many advanced data visualization capabilities. Through, the Seaborn library can be used to draw a variety of charts, and matrix plots, grid plots, regression plots etc. it can be used to draw distribution and categorical plots as well.~~



1. Distribution plot

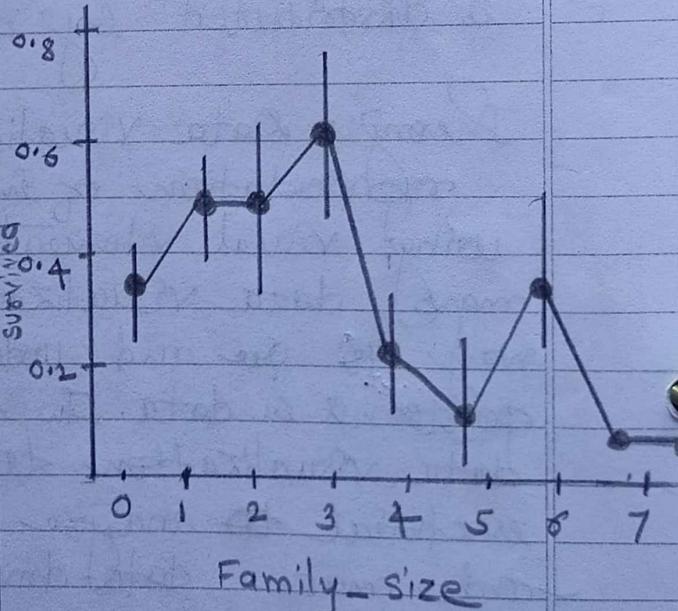


2. Histogram plot



3. Joint plot

4. Pair plot



- Distribution plot

Distribution plots visually assess the distribution of sample data by comparing the empirical distribution of the data with the theoretical values expected from a specified distribution.

`sns.distplot (dataset ['fare'])`

- Histogram plot

In Statistics, a histogram is a graphical representation of the distribution of data. The histogram is represented by a set of rectangles, adjacent to each other, where each bar represents a kind of data.

`sns.distplot (dataset ['Fare'], kde = False, bins = 10)`

- Joint plot

Joint plot is the way of understanding the relationship between two variables and the distribution of the individual each variables.

~~`sns.jointplot (x = 'age', y = 'fare', data = dataset)`~~

`sns.jointplot (x = 'age', y = 'fare', data = dataset, kind = 'hex')`



• pair plot :

Pair plot is used to understand the best set of the features to explain a relationship between two variables onto form the most separated clusters.

Conclusion: Thus we successfully implemented Simple Data Visualization techniques using Python on Titanic Data Set.

Refer
24/3/23



Assignment No.

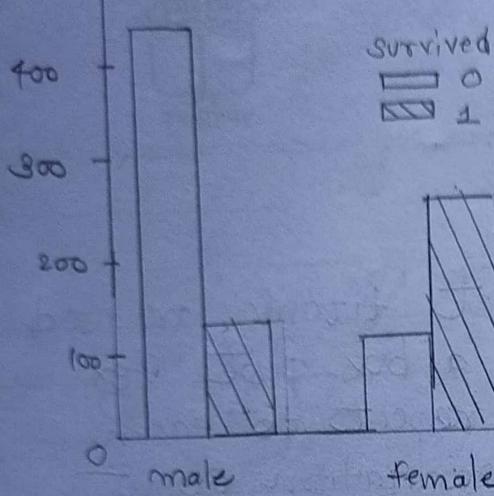
Title : Data Visualization 2

Objective : 1. Use the inbuilt dataset 'titanic' and use it in above problem. Plot a box plot for distribution of age with respect to each gender along with the information about whether they survived or not. (column : 'sex' and 'age').

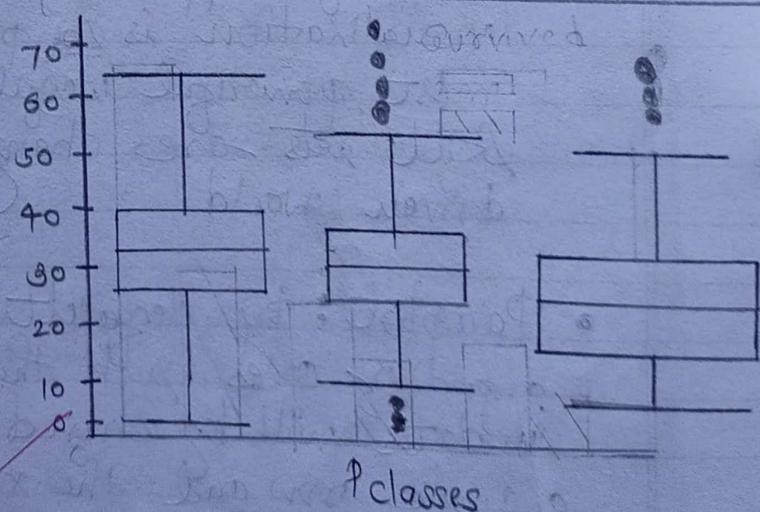
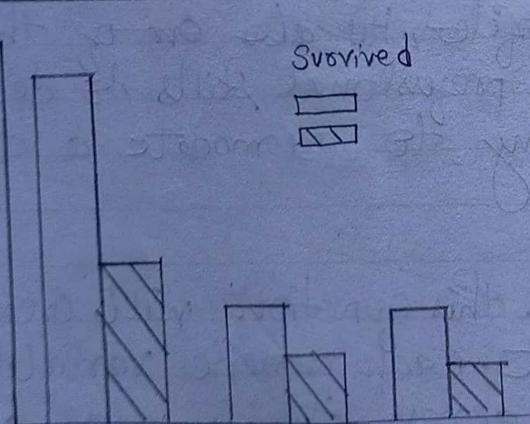
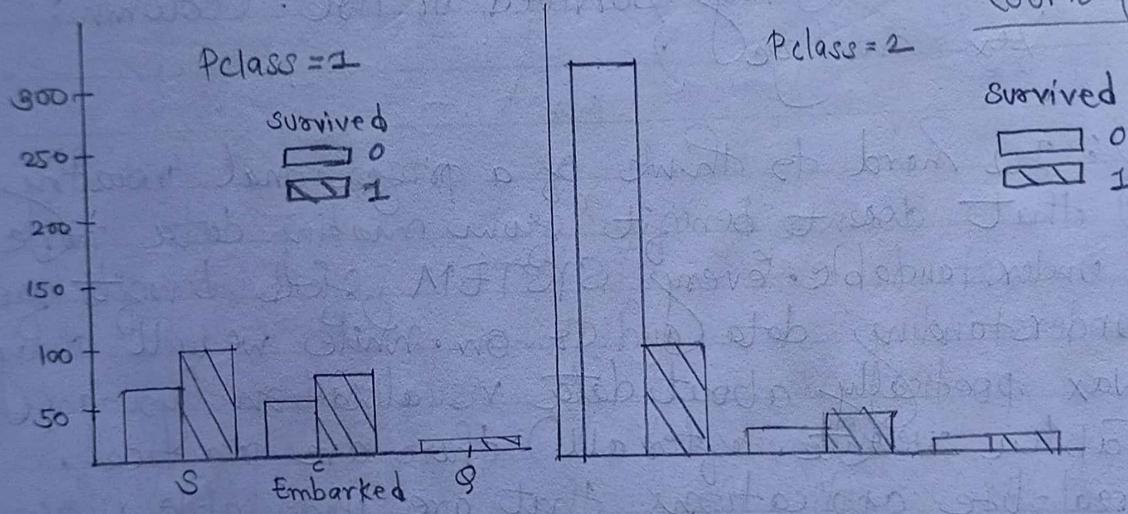
Theory: It's hard to think of a professional industry that doesn't benefit from making data more understandable. Every system yields benefits from understanding data (and so on). While we will always wax poetically about data visualization (you are (in Tableau website, after all)) there are practical, real-life applications that are undeniable. Since visualization is so prolific, it's also one of the most powerful useful professional skills to develop. Skill sets are changing to accommodate a data-driven world.

- **Pairplot:** By default, this function will create a grid of axes such that each numeric variable in data will be shared across the y-axis across a single row and the x-axis across a single column.
- **Barplot:** A barplot (or barchart) is one of the most common type of graphic. It shows the relationship between a numeric variable and a categorical variable.

Bar plot



Count plot



Box plot



- **Count plot:** A count plot be thought of as its a histogram across a categorical, instead of quantitative variable. The basic API and options are identical to those of barplot(), so you can compare counts across nested variables.

sns.countplot(x=df[["class"]])

- **Box plot:** A box and whisker plot also called box plot - displays the five-number summary of set of data. The five-number summary is the minimum, first quartile median, third quartile and maximum. In a box plot we draw a box from the first quartile to the third quartile. A vertical line goes through the box at the median.

sns.boxplot(x=df[["age"]])

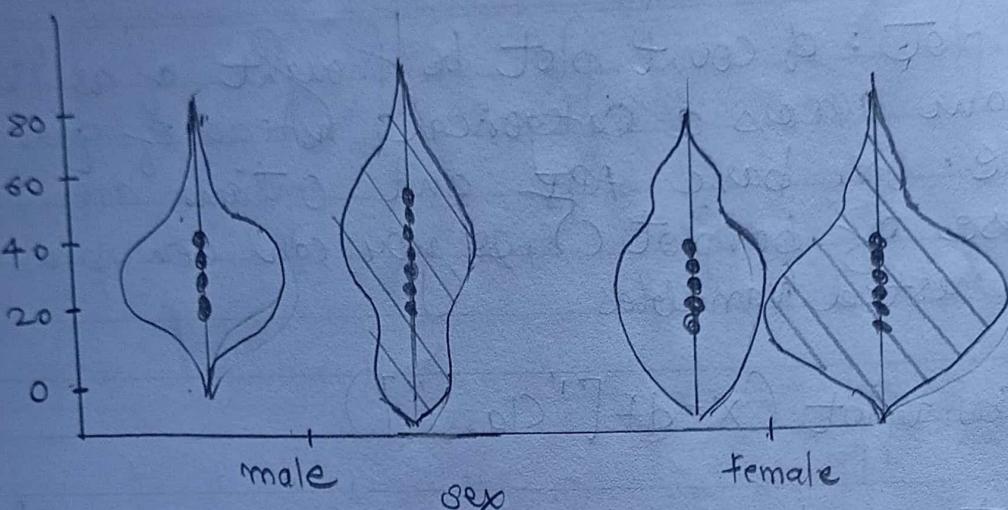
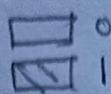
- **Violin plot:** A violin plot is hybrid of box plot and a kernel density plot, which shows peaks in the data. It is used to visualize the distribution of the numerical data.

sns.violinplot(x=df[["age"]])

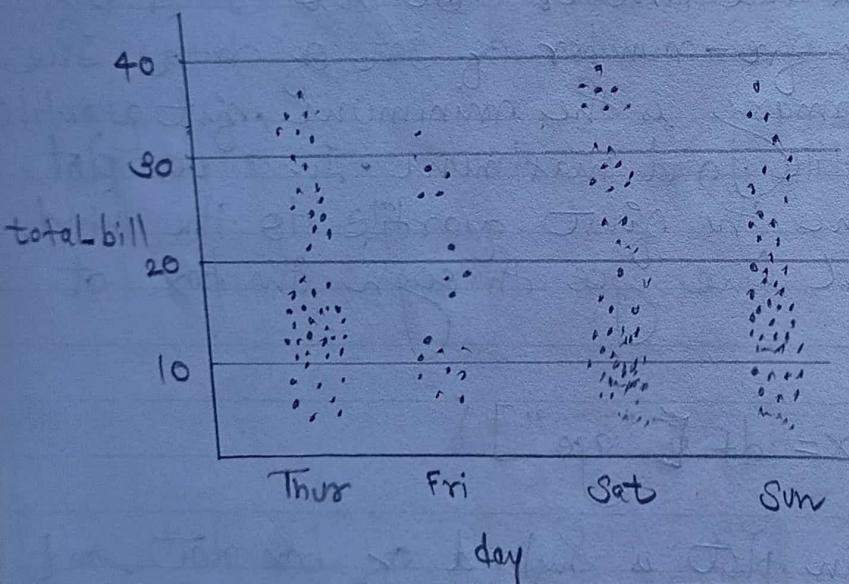
- **Strip plot:** Strip plot can be drawn on its own, but it is also good complement to a box or violin plot in cases where you want to show all observations along with some representation of the underlying distribution.

violin plot

survived

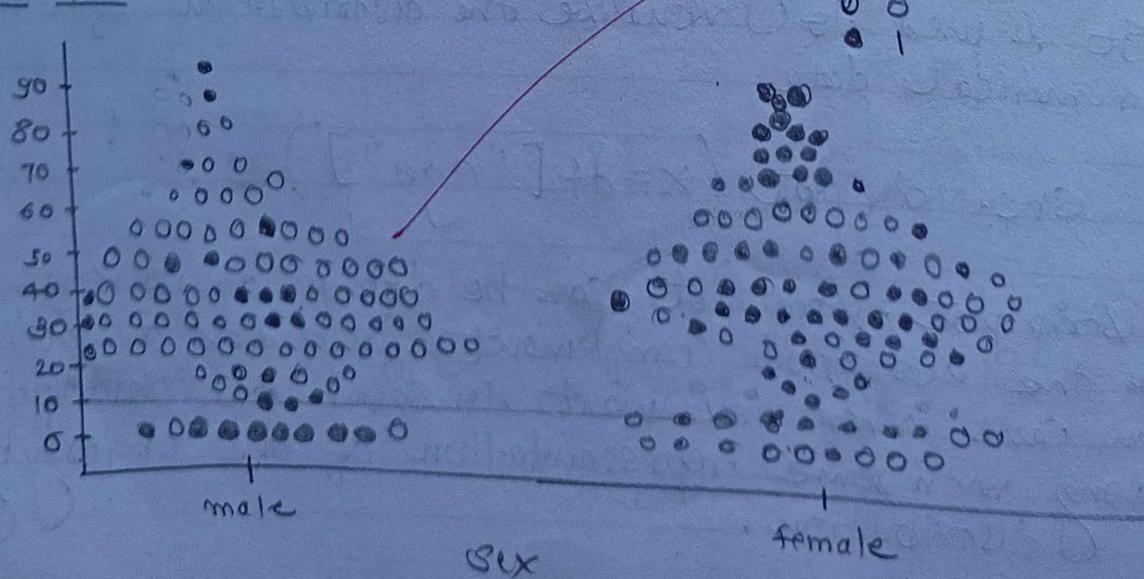


strip plot



swarm plot

survived





sns.stripplot (data = tips, x = "total_bill")

- **Swarm plot:** This function is similar to stripplot() but the points are adjusted only along the categorical axis so that they don't overlap. This gives a better representation of the distribution of values, but it doesn't scale well to large numbers of the distribution of values, but it does not scale well to large numbers of observations. This type of plot is sometimes called a "beeswarm".

Conclusion: Thus we have successfully implemented simple data visualization techniques using Python on Titanic dataset.

RDBMS
24/3/23

Assignment 10

Data Visualization III

Download the Iris flower dataset or any other dataset into a DataFrame. (e.g., <https://archive.ics.uci.edu/ml/datasets/Iris> (<https://archive.ics.uci.edu/ml/datasets/Iris>)). Scan the dataset and give the inference as:

1. List down the features and their types (e.g., numeric, nominal) available in the dataset.
2. Create a histogram for each feature in the dataset to illustrate the feature distributions.
3. Create a boxplot for each feature in the dataset.
4. Compare distributions and identify outliers

In [1]:

```
import seaborn as sns  
import matplotlib.pyplot as plt  
import pandas as pd
```

In [2]:

```
df=pd.read_csv('IRIS.csv')
```

In [3]:

```
df
```

Out[3]:

	sepal_length	sepal_width	petal_length	petal_width	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
...
145	6.7	3.0	5.2	2.3	Iris-virginica
146	6.3	2.5	5.0	1.9	Iris-virginica
147	6.5	3.0	5.2	2.0	Iris-virginica
148	6.2	3.4	5.4	2.3	Iris-virginica
149	5.9	3.0	5.1	1.8	Iris-virginica

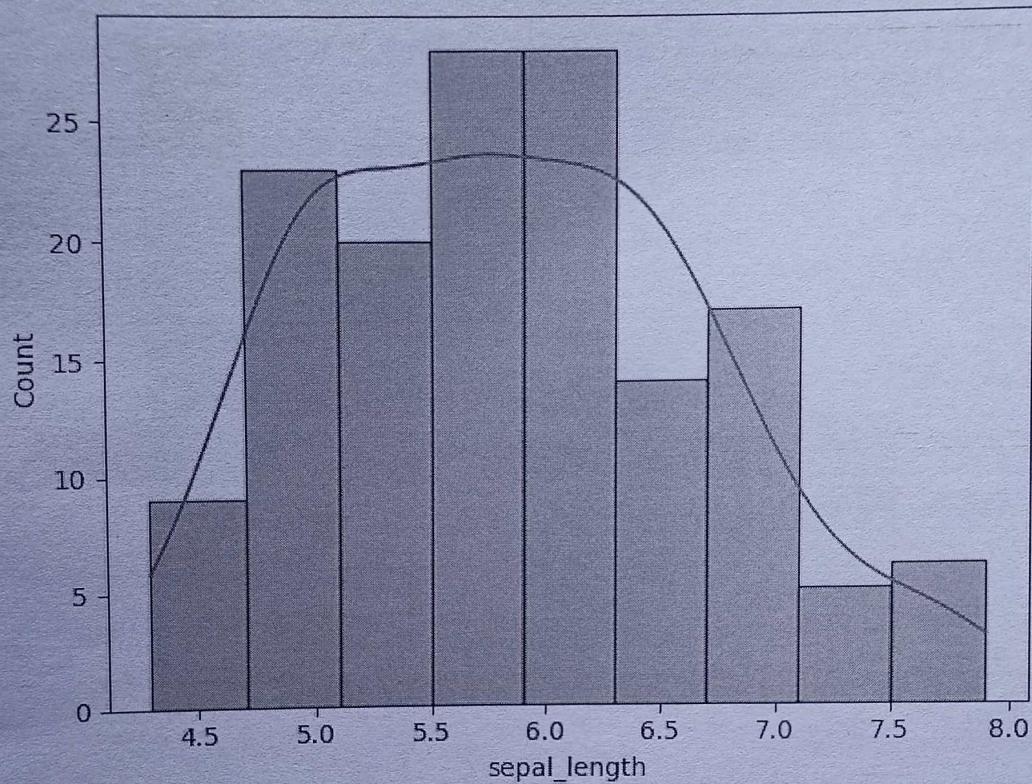
150 rows × 5 columns

In [4]:

```
sns.histplot(df['sepal_length'], kde=True)
```

Out[4]:

```
<AxesSubplot: xlabel='sepal_length', ylabel='Count'>
```

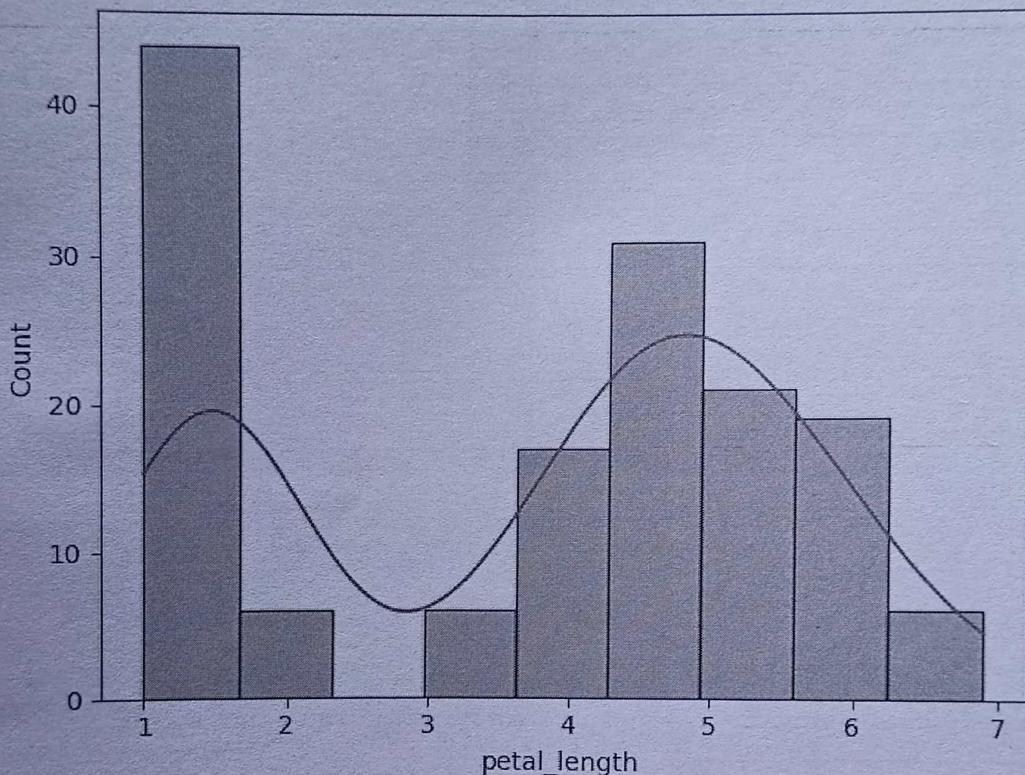


In [5]:

```
sns.histplot(df['petal_length'], kde=True)
```

Out[5]:

```
<AxesSubplot: xlabel='petal_length', ylabel='Count'>
```

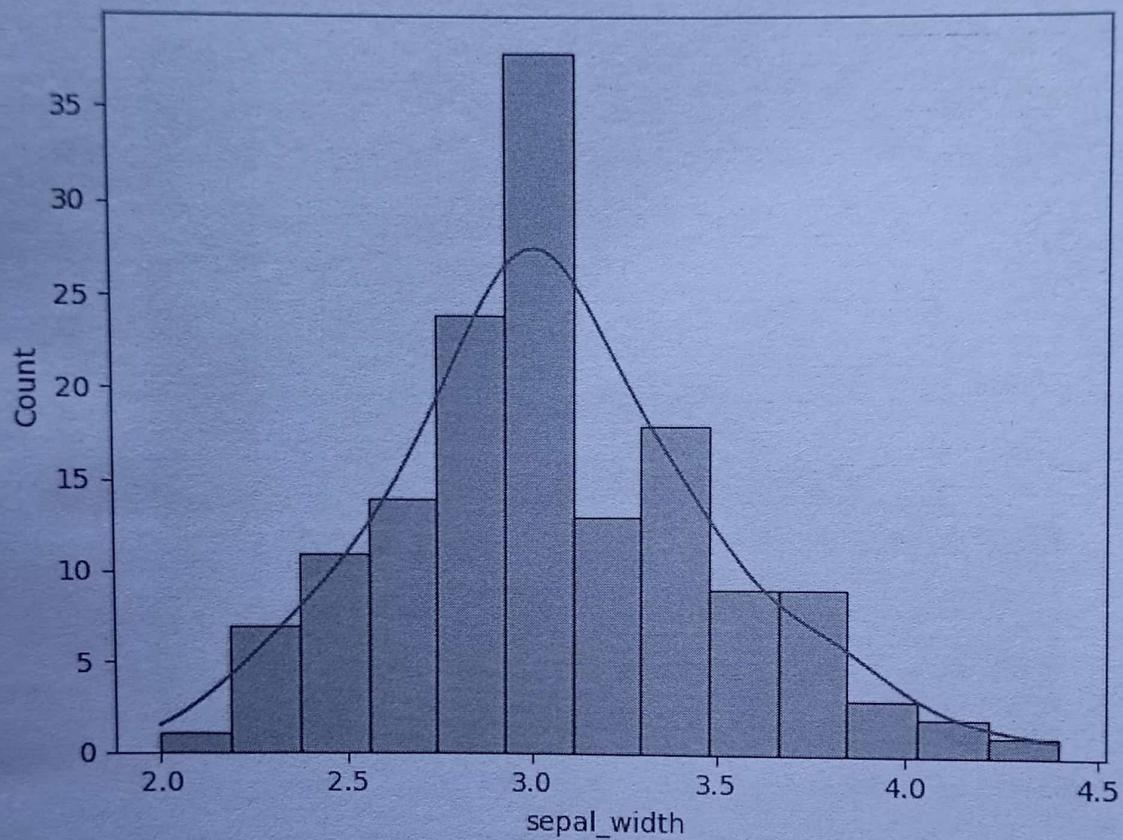


In [6]:

```
sns.histplot(df['sepal_width'],kde=True)
```

Out[6]:

```
<AxesSubplot: xlabel='sepal_width', ylabel='Count'>
```

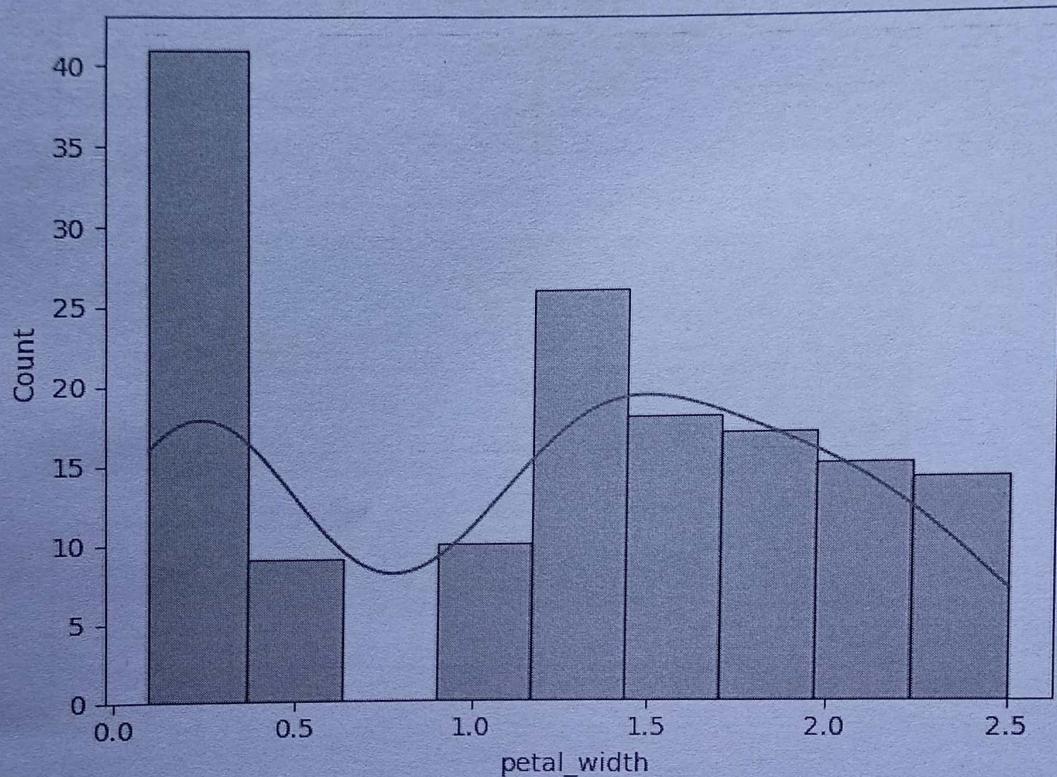


In [7]:

```
sns.histplot(df['petal_width'], kde=True)
```

Out[7]:

```
<AxesSubplot: xlabel='petal_width', ylabel='Count'>
```

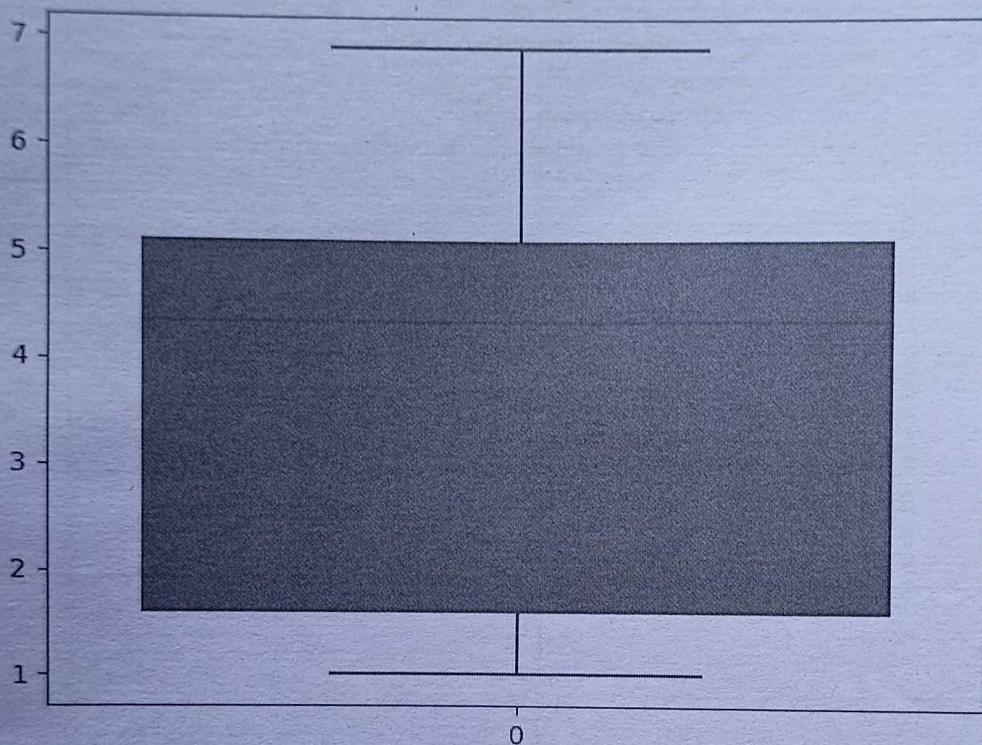


In [8]:

```
sns.boxplot(df['petal_length'],color='Red')
```

Out[8]:

<AxesSubplot: >

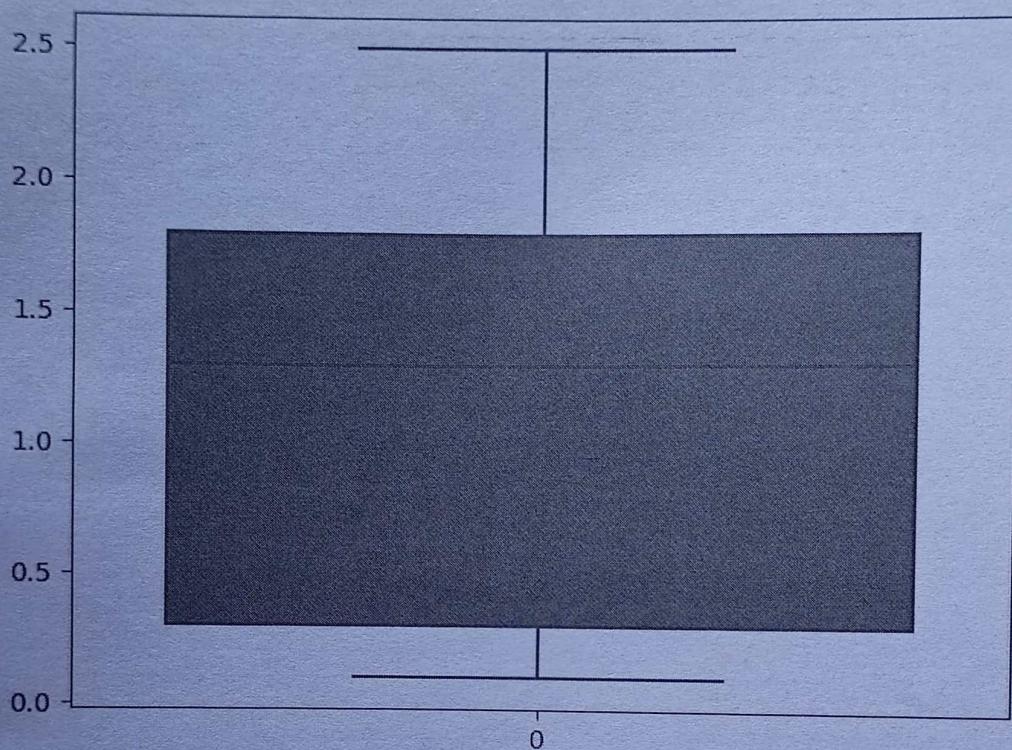


In [9]:

```
sns.boxplot(df['petal_width'],color='Red')
```

Out[9]:

<AxesSubplot: >

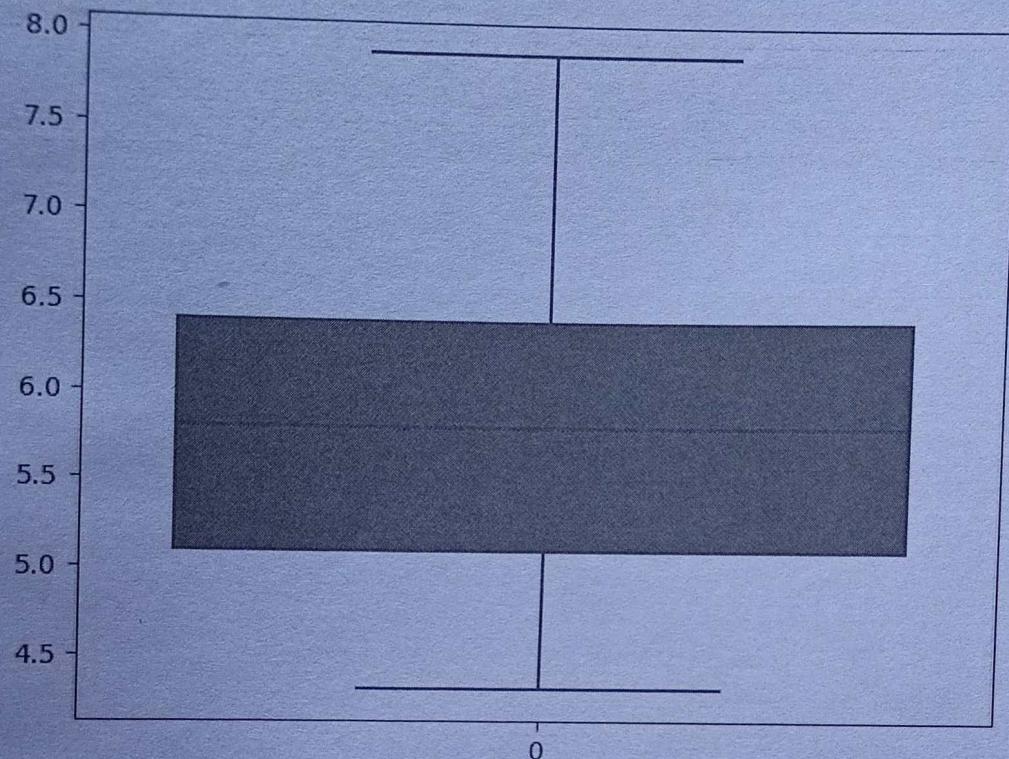


In [10]:

```
sns.boxplot(df['sepal_length'], color='Red')
```

Out[10]:

<AxesSubplot: >

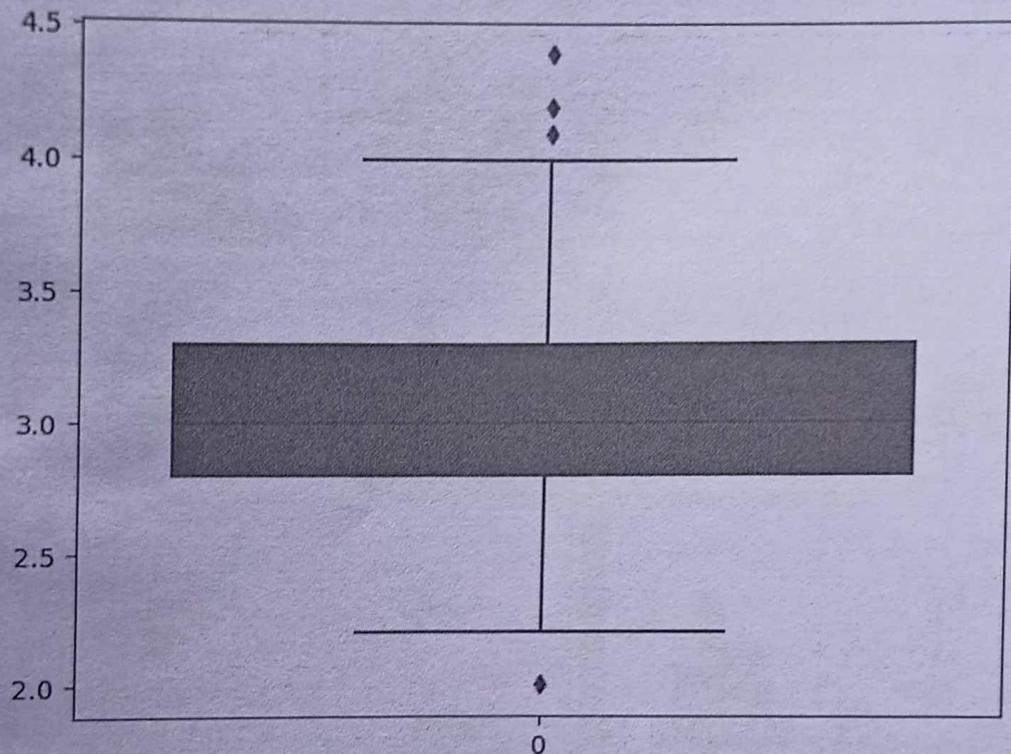


In [11]:

```
sns.boxplot(df['sepal_width'],color='Red')
```

Out[11]:

<AxesSubplot: >

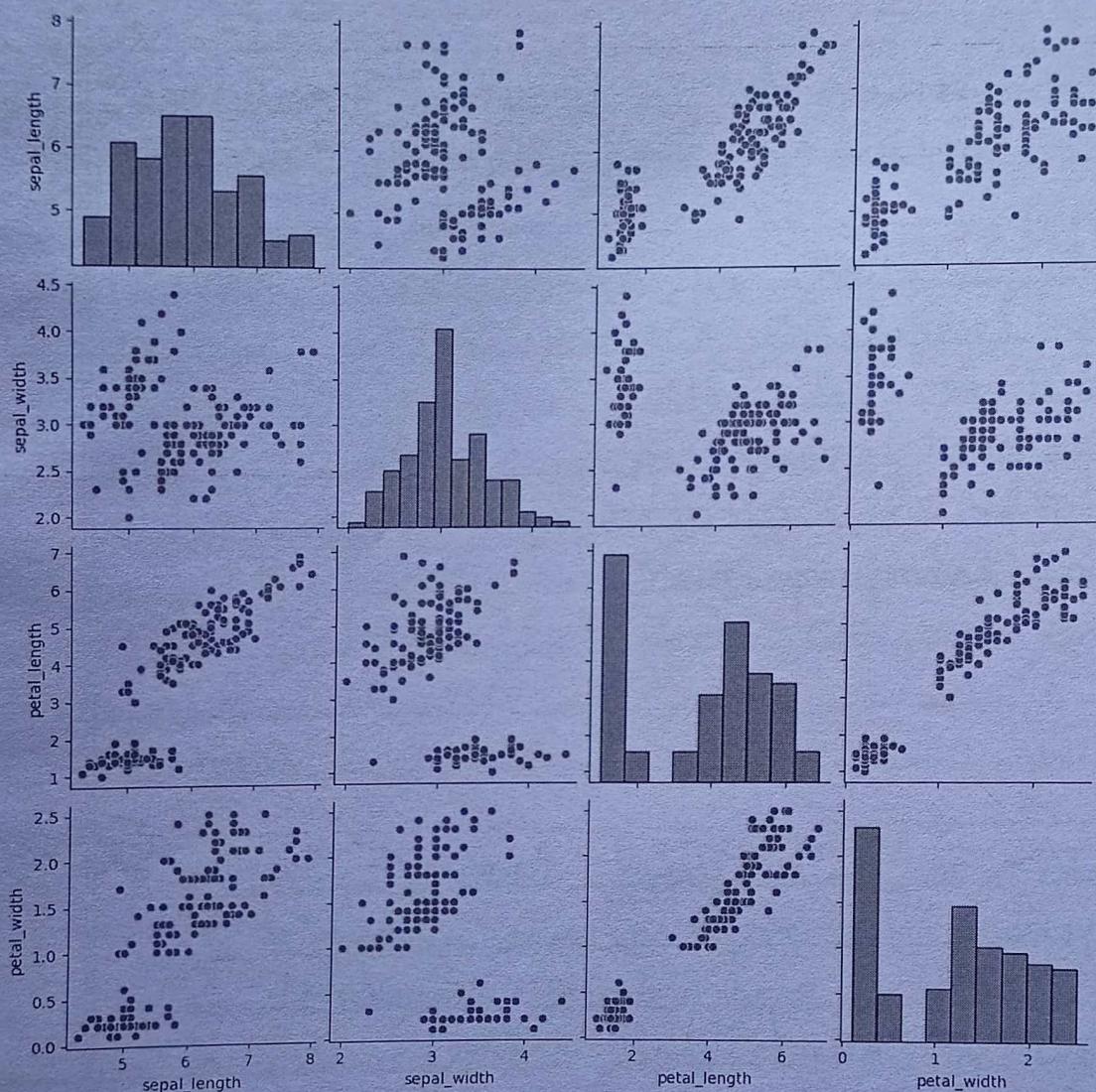


In [12]:

```
sns.pairplot(df)
```

Out[12]:

```
<seaborn.axisgrid.PairGrid at 0x20f245f47d0>
```





Assignment No. 11

Title : Scala

Objective : To write a simple program in SCALA using Apache Spark framework.

Theory : Scala stands for Scalable language. Scala is general-purpose, high level, multi-paradigm programming language. It is a pure object-oriented programming language which also provides support to the functional programming approach. Scala programs can convert to byte codes and can run on the JVM.

Features of SCALA :

- Object Oriented : Every value in SCALA is an object. It is purely object-oriented programming language. The behaviour and type of objects are depicted by the classes and traits in SCALA.
- Functional : It is also a functional programming language as every function is a value and every value is an object. It provides the support for the higher-order functions, nested functions, anonymous functions etc.
- Statically Typed : The process of verifying and then enforcing the constraints of types is done at compile time in Scala. In most cases, the user has no need to specify a type.



- Extensible: New language constructs can be added to scale in the form of libraries. Scala is designed to interpolate with the JRE (Java Runtime Environment).
- Concurrent & Synchronous Processing: Scala allows the user to write codes in an immutable manner that makes it easy to apply the parallelism (Synchronization) and concurrency.

Syntax of Scala:

① Objects in scala:

A typical scala program creates many objects which as you know, interact by invoking methods.

```
var obj = new Dog();
```

② Functions in scala:

One can divide up the code into separate function

```
def function_name ([parameter-list])  
[ : return_type ]
```

{

}

// function body,



Variable in Scala:

Variables are nothing but reserved memory to store values

Var my-var: String = "foo"

Conditional Statements:

: If statement - Nested if else

: If-else statement - If else-if

if (condition):
{ }

// Statement to be executed

}

if (condition)
{
// if block statements
}

else {

// else block statements

}

Conclusion: Thus we performed a single string computation program in Scala.

Feb 15/15