

C (4)	P (4)	A (2)	Total (10)	Sign

<b>Assignment No.</b>	5
<b>Title:</b>	1.Implement logistic regression using Python/R to perform classification on Social_Network_Ads.csv dataset. <u>2.</u> Compute Confusion matrix to find TP, FP, TN, FN, Accuracy, Error rate, Precision, Recall on the given dataset
<b>Roll No:</b>	
<b>Class:</b>	T.E.
<b>Date:</b>	
<b>Subject:</b>	DSBDA Lab



## **Group A**

### **Assignment No: 05**

#### **Aim:**

1. Implement logistic regression using Python/R 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.

#### **Objective:**

Students should be able to data analysis using logistic regression using Python for any open source dataset.

#### **Prerequisite:**

1. Basic of Python Programming
2. Concept of Regression.

#### **Theory:**

1. Logistic Regression
2. Differentiate between Linear and Logistic Regression
3. Sigmoid Function
4. Types of Logistic Regression
5. Confusion Matrix Evaluation Metrics

## 1. Logistic Regression

Classification techniques are an essential part of machine learning and data mining applications. Approximately 70% of problems in Data Science are classification problems. There are lots of classification problems that are available, but logistic regression is common and is a useful regression method for solving the binary classification problem. Another category of classification is Multinomial classification, which handles the issues where multiple classes are present in the target variable. For example, the IRIS dataset is a very famous example of multi-class classification. Other examples are classifying article/blog/document categories.

Logistic Regression can be used for various classification problems such as spam detection. Diabetes prediction, if a given customer will purchase a particular product or will they churn another competitor, whether the user will click on a given advertisement link or not, and many more examples are in the bucket.

Logistic Regression is one of the most simple and commonly used Machine Learning algorithms for two-class classification. It is easy to implement and can be used as the baseline for any binary classification

problem. Its basic fundamental concepts are also constructive in deep learning. Logistic regression describes and estimates the relationship between one dependent binary variable and independent variables.

Logistic regression is a statistical method for predicting binary classes. The outcome or target variable is dichotomous in nature. Dichotomous means there are only two possible classes. For

example, it can be used for cancer detection problems. It computes the probability of an event occurring.

It is a special case of linear regression where the target variable is categorical in nature. It uses a log of odds as the dependent variable. Logistic Regression predicts the probability of occurrence of a binary event utilising a logit function.

**Linear Regression Equation:** Linear regression shows the linear relationship between two variables. The equation of linear regression is

similar to the slope formula what we have learned before in earlier classes such as linear equations in two variables. It is given by;

$$Y = a + bX$$

$$a = \frac{[(\sum y)(\sum x^2) - (\sum x)(\sum xy)]}{[n(\sum x^2) - (\sum x)^2]}$$

$$b = \frac{[n(\sum xy) - (\sum x)(\sum y)]}{[n(\sum x^2) - (\sum x)^2]}$$

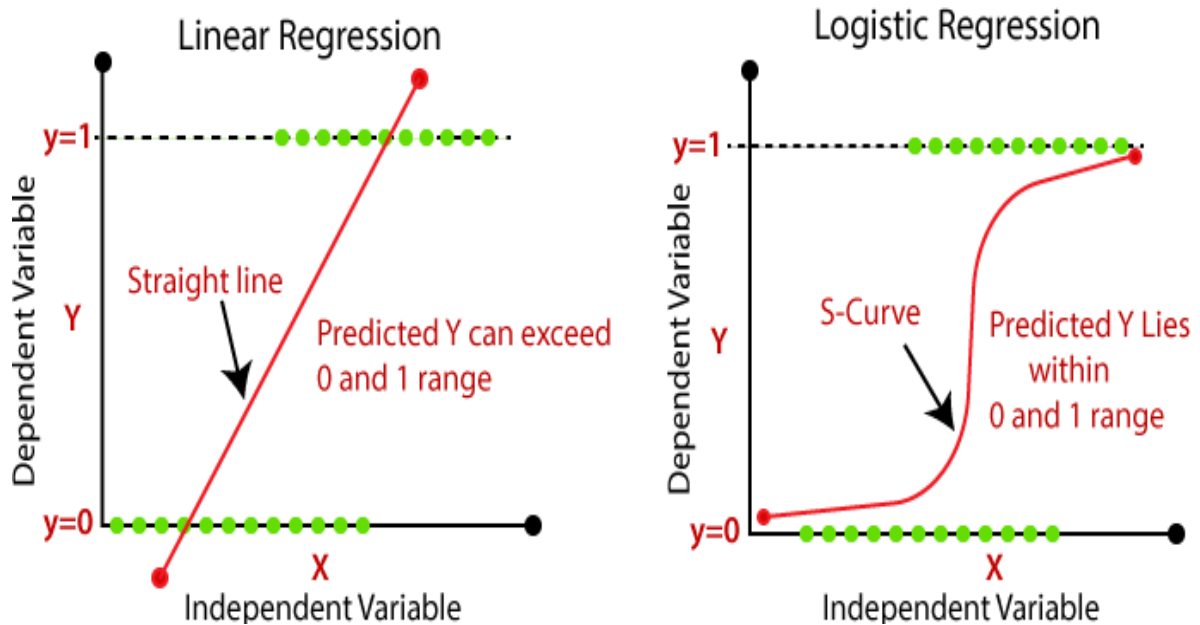
Now, here we need to find the value of the slope of the line, b, plotted in scatter plot and the intercept, a.

**Sigmoid Function:**

$$S(x) = \frac{1}{1 + e^{-x}} = \frac{e^x}{e^x + 1} = 1 - S(-x).$$

## 2.Differentiate between Linear and Logistic Regression

Linear Regression	Logistic Regression
Linear regression is used to predict the continuous dependent variable using a given set of independent variables.	Logistic Regression is used to predict the categorical dependent variable using a given set of independent variables.
Linear Regression is used for solving Regression problem.	Logistic regression is used for solving Classification problems.
In Linear regression, we predict the value of continuous variables.	In logistic Regression, we predict the values of categorical variables.
In linear regression, we find the best fit line, by which we can easily predict the output.	In Logistic Regression, we find the S-curve by which we can classify the samples.
Least square estimation method is used for estimation of accuracy.	Maximum likelihood estimation method is used for estimation of accuracy.
The output for Linear Regression must be a continuous value, such as price, age, etc.	The output of Logistic Regression must be a Categorical value such as 0 or 1, Yes or No, etc.
In Linear regression, it is required that relationship between dependent variable and independent variable must be linear.	In Logistic regression, it is not required to have the linear relationship between the dependent and independent variable.
In linear regression, there may be collinearity between the independent variables.	In logistic regression, there should not be collinearity between the independent variable.



### 3.Sigmoid Function

A **sigmoid function** is a mathematical function having a characteristic "S"-shaped curve or **sigmoid curve**.

A common example of a sigmoid function is the logistic function shown in the first figure and defined by the formula:

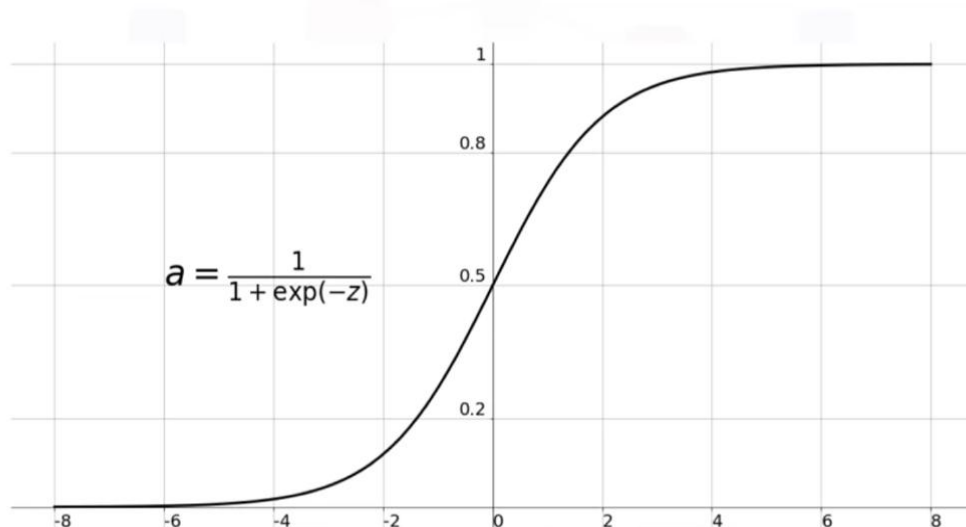
$$S(x) = \frac{1}{1 + e^{-x}} = \frac{e^x}{e^x + 1} = 1 - S(-x).$$

Special cases of the sigmoid function include the Gompertz curve (used in modeling systems that saturate at large values of x) and the ogive curve (used in the spillway of some dams). Sigmoid functions have domain of all real numbers, with return (response) value commonly monotonically increasing but could be decreasing. Sigmoid functions most often show a return value (y axis) in the range 0 to 1. Another commonly used range is from -1 to 1.

A wide variety of sigmoid functions including the logistic and hyperbolic tangent functions have been used as the activation function of artificial neurons. Sigmoid curves are also common in

statistics as cumulative distribution functions (which go from 0 to 1), such as the integrals of the logistic density, the normal density, and Student's  $t$  probability density functions. The logistic sigmoid function is invertible, and its inverse is the logit function. The sigmoid function could generate some non-zero values, resulting in a dense representation.

## Sigmoid Function



### 4. Types of Logistic Regression

There are three main types of logistic regression: binary, multinomial and ordinal. They differ in execution and theory. Binary regression deals with two possible values, essentially: yes or no. Multinomial logistic regression deals with three or more values. And ordinal logistic regression deals with three or more classes in a predetermined order.

**1] Binary logistic regression:** Binary logistic regression was mentioned earlier in the case of classifying an object as an animal or not an animal—it's an either/or solution. There are just two possible outcome answers. This concept is typically represented as a 0 or a 1 in coding.

**2]Multinomial logistic regression:** Multinomial logistic regression is a model where there are multiple classes that an item can be classified as. There is a set of three or more predefined classes set up prior to running the model.

**3]Ordinal logistic regression:** Ordinal logistic regression is also a model where there are multiple classes that an item can be classified as; however, in this case an ordering of classes is required. Classes do not need to be proportionate. The distance between each class can vary.

## 5.Confusion Matrix Evaluation Metrics

### What is a confusion matrix?

It is a matrix of size  $2 \times 2$  for binary classification with actual values on one axis and predicted on another.

		ACTUAL	
		Negative	Positive
PREDICTION	Negative	TRUE NEGATIVE	FALSE NEGATIVE
	Positive	FALSE POSITIVE	TRUE POSITIVE

- **Confusion Matrix**

Let's understand the confusing terms in the confusion matrix: **true positive**, **true negative**, **false negative**, and **false positive** with an example.

#### EXAMPLE

A machine learning model is trained to predict tumor in patients. The test dataset consists of 100 people.



		ACTUAL	
		Negative	Positive
PREDICTION	Negative	60	8
	Positive	22	10

- **Confusion Matrix for tumor detection**

- **True Positive (TP)** — model correctly predicts the positive class (prediction and actual both are positive). In the above example, **10 people** who have tumors are predicted positively by the model
- **True Negative (TN)** — model correctly predicts the negative class (prediction and actual both are negative). In the above example, **60 people** who don't have tumors are predicted negatively by the model.
- **False Positive (FP)** — model gives the wrong prediction of the negative class (predicted-positive, actual-negative). In the above example, **22 people** are predicted as positive of having a tumor, although they don't have a tumor. FP is also called a **TYPE I** error.
- **False Negative (FN)** — model wrongly predicts the positive class (predicted-negative, actual-positive). In the above example, **8 people** who have tumors are predicted as negative. FN is also called a **TYPE II** error.

With the help of these four values, we can calculate True Positive Rate (TPR), False Negative Rate (FPR), True Negative Rate (TNR), and False Negative Rate (FNR).

$$TPR = \frac{TP}{Actual\ Positive} = \frac{TP}{TP + FN}$$

$$FNR = \frac{FN}{Actual\ Positive} = \frac{FN}{TP + FN}$$

$$TNR = \frac{TN}{Actual\ Negative} = \frac{TN}{TN + FP}$$

$$FPR = \frac{FP}{Actual\ Negative} = \frac{FP}{TN + FP}$$

Even if data is imbalanced, we can figure out that our model is working well or not. For that, **the values of TPR and TNR should be high, and FPR and FNR should be as low as possible.**

With the help of TP, TN, FN, and FP, other performance metrics can be calculated.

- **Accuracy:-** Accuracy (ACC) is calculated as the number of all correct predictions divided by the total number of the dataset. The best accuracy is 1.0, whereas the worst is 0.0. It can also be calculated by  $1 - \text{ERR}$ . Accuracy is calculated as the total number of two correct predictions (TP + TN) divided by the total number of a dataset (P + N).

$$ACC = \frac{TP + TN}{TP + TN + FN + FP} = \frac{TP + TN}{P + N}$$

- **Error rate:-** Error rate (ERR) is calculated as the number of all incorrect predictions divided by the total number of the dataset. The best error rate is 0.0, whereas the worst is 1.0. Error rate is calculated as the total number of two incorrect predictions (FN + FP) divided by the total number of a dataset (P + N).

$$ERR = \frac{FP + FN}{TP + TN + FN + FP} = \frac{FP + FN}{P + N}$$

- **Precision:-** Out of all the positive predicted, what percentage is truly positive.

$$Precision = \frac{TP}{TP + FP}$$

The precision value lies between 0 and 1.

- **Recall:-** Out of the total positive, what percentage are predicted positive. It is the same as TPR (true positive rate).

$$Recall = \frac{TP}{TP + FN}$$

### Algorithm:

**Step 1: Import libraries and create alias for Pandas, Numpy and Matplotlib**

**Step 2: Import the Social\_Media\_Adv Dataset**

**Step 3: Initialize the data frame**

**Step 4: Perform Data Preprocessing**

- Convert Categorical to Numerical Values if applicable.
- Check for Null Value.
- Covariance Matrix to select the most promising features.
- Divide the dataset into Independent(X) and Dependent(Y) variables.
- Split the dataset into training and testing datasets.
- Scale the Features if necessary.

**Step 5: Use Logistic regression( Train the Machine ) to Create Model**

```
# import the class
from sklearn.linear_model import LogisticRegression
```

```
# instantiate the model (using the default
parameters)
logreg = LogisticRegression()
# fit the model with data
logreg.fit(xtrain,ytrain)
# y_pred=logreg.predict(xtest)
```

**Step 6: Predict the y\_pred for all values of train\_x and test\_x**

**Step 7: Evaluate the performance of Model for train\_y and test\_y**

**Step 8: Calculate the required evaluation parameters**

```
from sklearn.metrics import
precision_score, confusion_matrix, accuracy_score, recall_score
cm= confusion_matrix(ytest, y_pred)
```

## **Conclusion:**

In this way we have done data analysis using logistic regression for Social Media Adv. and evaluate the performance of model.

## **Oral Question:**

1) Consider the binary classification task with two classes positive and negative. Find out TP, TP, FP, TN, FN, Accuracy, Error rate, Precision, Recall

N = 165	Predicted YES	Predicted NO
Actual YES	TP = 150	FN = 10
Actual NO	FP = 20	TN = 100

2) Comment on whether the model is best fit or not based on the calculated values.

3) Write python code for the preprocessing mentioned in step 4. and Explain every step in detail.