

Logistic Regression :-

- It comes under supervised learning [both x & y].
- It is used to solve classification problems [target is categorical]

eg: Diabetic data, cancerous data etc.,

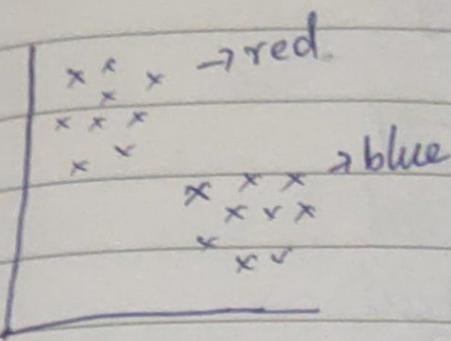
- It is used to predict the probability of categorical dependent variable.

Assumption of Logistic Regression.

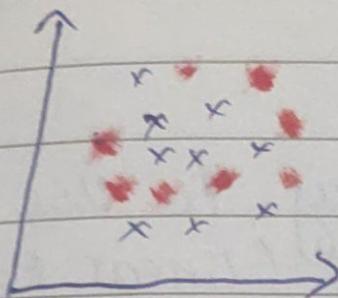
- The target variable is binary i.e. Yes/No, True/False, 0/1..
- No or Little multicollinearity
- Outliers should be handled.
- Sample size should be sufficiently large.

→ Data is linearly separable.

Linearly Separable → There should be clear separation among the classes.



Non-Linear. → overlapping of classes.



How Logistic Regression works?

Prob. Statement :-

I want to predict whether student will pass or not in exam.

DATA:

No. of hours studied (x)	Result (y)
1	Fail
4	Fail
3	Fail
8	Pass
10	Pass
12	Pass.
5	?

What happens if we try to use LR?
(Linear Regression)

→ Linear regression will fit a straight line which separates two classes.

$$y = mx + c$$

→ If we use Linear regression to make predictions it will make use straight line and give prediction either 0 & 1.

→ If data has outliers then line gets shifted, it results in misclassification.

Sigmoid Function :-

It is one of the mathematical functions that we use to transform continuous values into a probability values which lies within a range 0 to 1.

$$h_0(x) = g(z) = \frac{1}{1+e^{-z}} \rightarrow \text{O/P will be } [0,1]$$

Where $z = mx + c$ \rightarrow straight line
 $z = \beta_0 + \beta_1 x$ eqn.

Set a threshold (For Predictions).

$$h_0(x) > 0.5 \rightarrow y=1 \mid \text{Pass}$$

$$h_0(x) < 0.5 \rightarrow y=0 \mid \text{Fail}$$

<u>x</u>	<u>Result</u>	\hat{y}	$h_0(x)$	<u>Prediction</u>
2	Fail(0)	-50	0.2	Fail
3	Fail(0)	40	0.4	Fail
5	Fail(0)	-20	0.3	Fail
8	Pass(1)	70	0.6	Pass
10	Pass(1)	80	0.8	Pass
12	Pass(1)	100	0.9	Pass.

* How it works?

→ It will fit lines

→ Make Predictions

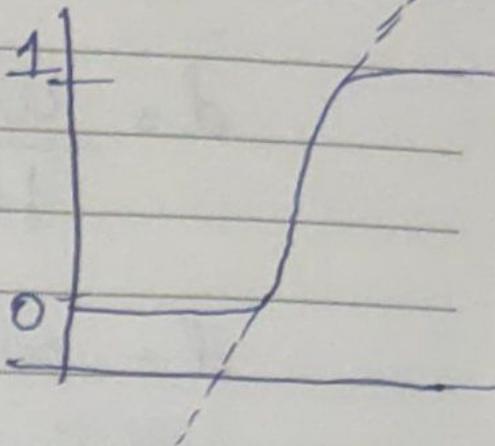
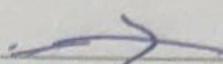
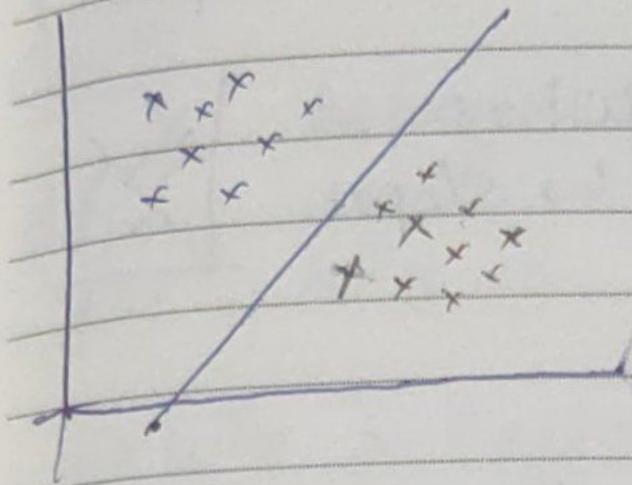
→ Transform continuous values using Sigmoid function

→ Set threshold.

→ If threshold is ≥ 0.5 classify as pass and if threshold is less than (< 0.5) classify as fail.

Linear Regression

Sigmoid



How classification works?

$$y = m_1 x_1 + m_2 x_2 + m_3 x_3 + c$$

$$y = \omega^T x + c$$

where $\omega^T = \begin{bmatrix} m_1 \\ m_2 \\ m_3 \end{bmatrix}$, $x = [x_1, x_2, x_3]$

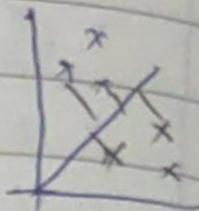
Let us say line passes through origin,

$$c = 0$$

$$y = \omega^T x$$

How to find distance b/w point and plane?

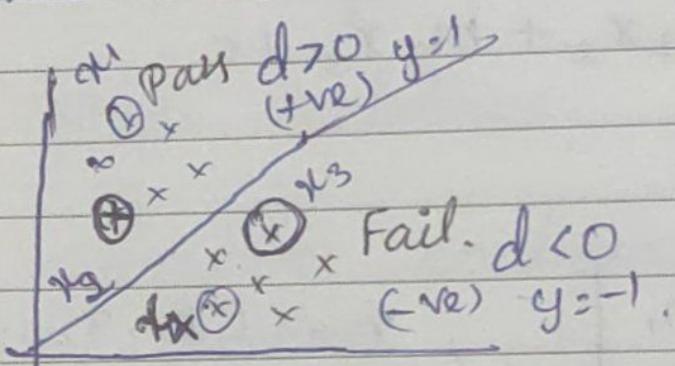
$$d = \frac{w^T x}{\|w\|}, \|w\|=1 \hookrightarrow \text{slopes.}$$



$$d = w^T x$$

$y \times d = y \times w^T x \geq 0$ then we say

the model has done correct classification



→ Any point above the line / plane will have positive distance

→ Any point below the line / plane will have negative distance.

→ All (+ve) classes are denoted by 1.

\rightarrow all (-ve) classes are denoted by -1.

Case 1
 x_1 belongs to +ve class

then,

$$y=1, w^T x > 0$$

$$y \times w^T x \geq (+ve) \times (+ve) = +ve$$

$y \times w^T x > 0, \therefore x_1$ is correctly classified

case 2

x_2 belongs -ve class

then

$$y=-1, w^T x > 0$$

$$y \times w^T x = (-ve) (+ve) = -ve$$

$y \times w^T x < 0$ so x_2 is misclassified.

case 3

x_3 belongs -ve class

$$y=-1, w^T x < 0$$

$$y \times w^T x = (-ve) (-ve) = +ve$$

$$y \times w^T x \geq 0$$

so, x_3 is correct classification.

Case 4:

x_4 belongs to +ve class

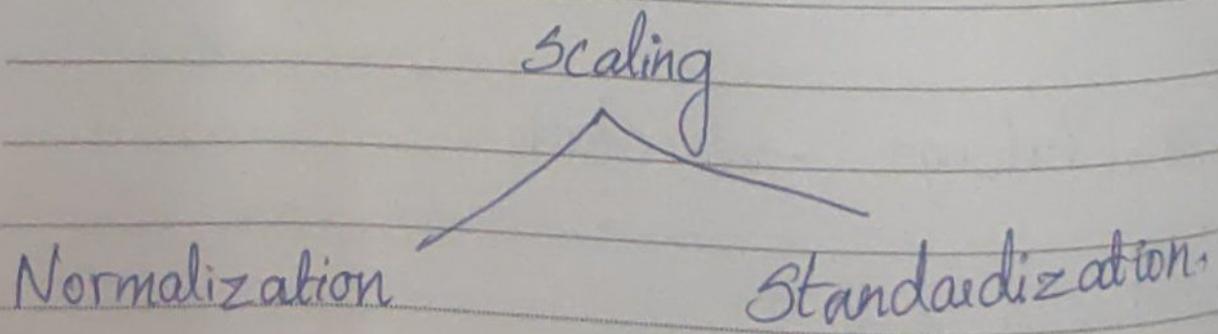
$$y=1, w^T x < 0$$

$$y \times w^T x = (+ve) (-ve) = -ve$$

$y \times w^T x < 0$, x_4 is misclassified.

Scaling:

scaling is always applied on continuous numerical column.



Normalization :-

→ Min max scalar

It transforms all the values range into a scale of 0 to 1 range

$$\text{Minmax} = \frac{x - x_{\min}}{x_{\max} - x_{\min}}$$

Standardization :-

→ Standard Scalar

$$z = \frac{x - \mu}{\sigma}$$

$\mu \rightarrow$ mean

$\sigma \rightarrow$ std.

It transforms all the values into a scale of z-value.

Qmp:-

Scaling is important to transform characteristic continuous data into certain scale

* - Scaling is very important in distance based algorithm.

Metric to Evaluate classification problems

Different metrics are .

- ① confusion Matrix
- ② Accuracy
- ③ Precision
- ④ Recall / sensitivity
- ⑤ F1 score .

Note:

classification algs are evaluated based on misclassification that model has done .

Confusion Matrix:-

		Predicted	
		0	1
Actual	0	TN	FP
	1	FN	TP

0-0 → TN

0-1 → FP

0-1 → FP

1-1 → FN

		Predicted	
		Non diabetic	diabetic
Actual	Non diabetic	Correct	wrong
	Diabetic	wrong	Correct

Date

Ex:

$TN \rightarrow$ Actual is Non-diabetic and predicted as dia non-diabetic

$TP \rightarrow$ Actual is diabetic and predicted as diabetic.

$FP \rightarrow$ Actual is Non-diabetic and predicted as diabetic
↓(Type-1 error)

$FN \rightarrow$ Actual is diabetic and predicted is Non-diabetic.

↓ Type-2 error

→ Confusion Matrix gives information about no. of correct classification & misclassification.

Correct classification = $TP + TN$

Mis classification = $FN + FP$.

③ Accuracy Score:-

Correct Predictions

$$\text{Accuracy} = \frac{\text{Total Predictions}}{\text{Total Predictions}}$$

$$\text{Accuracy} = \frac{TP + TN}{TP + TN + FP + FN}$$

→ Accuracy measures strength of model.

→ It talks about the percentage of correct predictions.

→ Accuracy score is best metric when data is balanced.

→ If data is not balanced then accuracy is not best metric; In that case, we use recall, precision and F1 score.

④ Precision

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

Out of positive predictions, what percentage are truly positive:

ex: $[0 \text{ } -1] \rightarrow$ how many are $-1]$

④ Recall

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

Out of all actual positive, what percentage is truly positive.

$[1 \text{ } -1] \rightarrow$ how many are $-1]$.

⑤ F1 score :-

harmonic mean

\rightarrow F1 score is (average) of Precision and Recall.

$$\text{F1 score} = 2 \times \frac{\text{Precision} \times \text{Recall}}{\text{Precision} + \text{Recall}}$$

eg:

F1-score gives total positive predictions out of error.

actual predicted

eg. $y \quad \hat{y}$

① confusion Matrix.

0 0
1 1
0 0
1 1
0 1

		Pred	
		Actual	0 1
		0	2 3
		1	2 2
		0	1

1 0
0 1
0 1
1 0

TP \rightarrow 2

TN \rightarrow 2

FP \rightarrow 3

FN \rightarrow 2

$$\textcircled{2} \text{ Accuracy} = \frac{TP + TN}{TP + TN + FP + FN} = \frac{2+2}{2+2+3+2} = 4/9 \\ = 0.44$$

$$\textcircled{3} \text{ Precision} = \frac{TP}{TP + FP} = \frac{2}{2+3} = 2/5 = 0.4$$

$$\textcircled{4} \text{ Recall} = \frac{TP}{TP + FN} = \frac{2}{2+2} = 2/4 = 0.5$$

$$5) F1\text{-Score} = \frac{2 \times \text{Precision} \times \text{Recall}}{\text{Precision} + \text{Recall}}$$

$$= \frac{2 \times 0.4 \times 0.5}{0.4 + 0.5} = 0.40$$

Classification

Binary classification

Multiclassification

→ If target columns has two categories / output

Diabetic

eg: Diabetic

Non diabetic

→ If target column has three or more categories / output.

Setosa

Virginica

Virginica

Logistic Regression works well with binary classification.

① Ratings

1

2

3

→ It's Not work well with Logistic regression.

ROC - AUC curve:

→ It is one of the metric used to evaluate Binary classification problems such as logistic regression.

ROC → [Receiver Operator characteristic]

Using this ROC, we evaluate our model at different thresholds.

$$h_0(x) < 0.5 \rightarrow y = 0$$

$$h_0(x) \geq 0.5 \rightarrow y = 1$$

<u>y</u>	<u>$h_0(x)$</u>	<u>$\hat{y}(\cos)$</u>	<u>$\hat{y}(0.3)$</u>	<u>threshold</u>
0	0.2	0	0	
1	0.4	0	1	
0	0.3	0	1	
1	0.7	1	1	
1	0.8	1	1	
0	0.5	1	1	

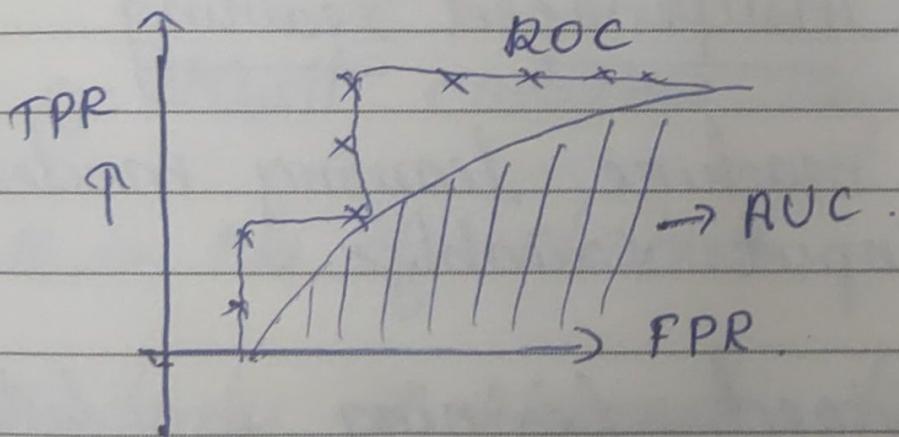
ROC will plot a graph of TPR vs FPR

TPR \rightarrow True Positive Rate.

$$\text{TPR} \Rightarrow \frac{\text{TP}}{\text{TP} + \text{FN}} = \frac{2}{2+1} = \frac{2}{3} \rightarrow \text{at } 0.5 \text{ threshold}$$

FPR \rightarrow False Positive Rate.

$$\text{FPR} = \frac{\text{FP}}{\text{FP} + \text{TN}} = \frac{1}{1+2} = \frac{1}{3} \text{ at } 0.5 \text{ threshold}$$



AUC [Area Under the curve]

\rightarrow AUC tells about the strength of the model by measuring area below the ROC curve.

Date

→ More the area below the ROC curve better is the model.

Range of AUC = [0,1]

AUC value near to '0' we say bad model

AUC Value near to '1' we say good model