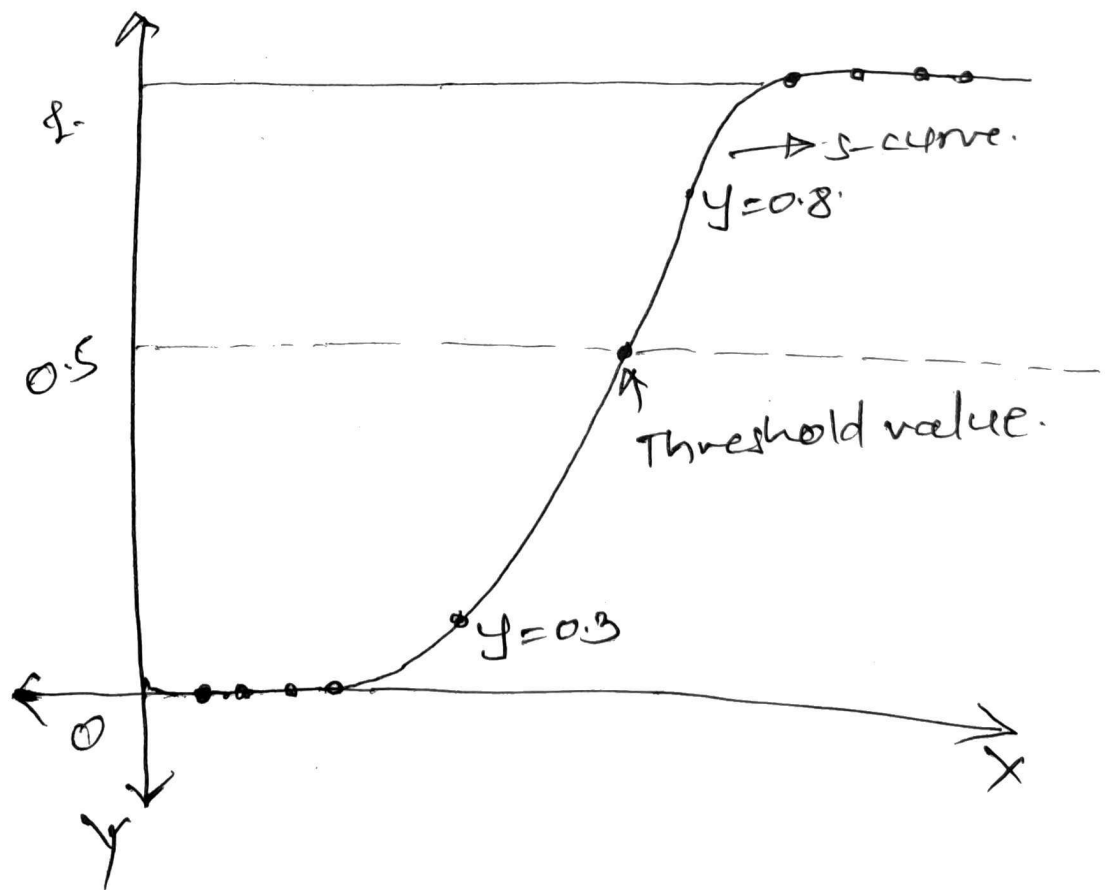


* Logistic Regression.

- Supervised learning technique.
- Used for predicting categorical dependent variables using a given set of independent variables.
- Predicts o/p of categorical dependent variable.
- It is used for solving classification problem.
- It gives probabilistic values which lie between 0 and 1.
- We can fit 'S' shaped logistic function which predicts maximum values (0 & 1).



* Logistic funⁿ: (sigmoid funⁿ):

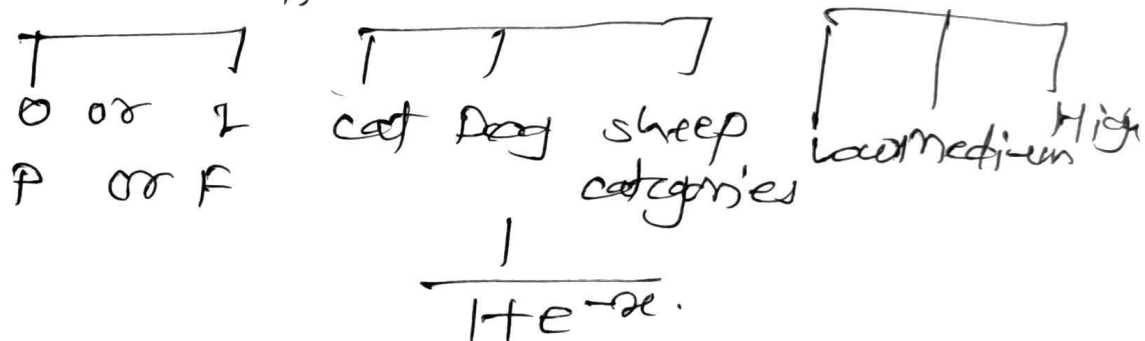
- used to map predicted values to probabilities
- maps real value into another value within a range 0 & 1.
- value of logistic regression must be bet 0 & 1 which cannot go beyond the limit so it forms a curve of 's' shape. & this curve is also called as sigmoid funⁿ or logistic funⁿ.
- threshold value is probability of either 0 or 1. values above ~~threshold~~ ^{tends to 1} & below threshold values tends to 0.

* Assumptions for logistic Regression:

- Dependent variables must be categorical in nature.
- independent variables should not have multi collinearity

$$y = b_0 + b_1x_1 + b_2x_2 + \dots + b_nx_n$$

Types: Binomial, multinomial, ordinal.



Applⁿ:

credit scoring., Hotel booking,
text editing, medicine, Gaming

Ad & dis ad: It constructs linear boundaries

- easier to setup & train - underfitting.
- When outcomes are linearly separable then it is most efficient
- should not be used if no. of observed are larger than no. of features
- nonlinear problems cannot be solved.

Parameters

Correlation

Statistical measure that determines the association or correlation between two variables.

Usage

Represent linear relationship between variables.

Dependent & independent var.

No difference.

Indicates

extent to which two variables move together.

objective

To find a numerical value expressing relationship between variables.

Regression

It describes how to numerically relate an independent variable to dependent variable.

To fit best line & to estimate one variable based on another

Both are different

Indicates impact of change of unit on the estimated variable in known variable (X).

To estimate values of random variables on the basis of values of fixed variables

1) SST (sum of squares of total):

$$SST = \sum (y_i - \bar{y})^2.$$

sum of squared differences betⁿ individual data pts (y_i) & mean of response variable (\bar{y}).

2) SSR (sum of squares of regression):

sum of squared differences betⁿ predicted data pts (\hat{y}_i) & the mean of response variable (\bar{y}).

$$SSR = \sum (\hat{y}_i - \bar{y})^2.$$

3) sum of squares error (SSE): sum of squared differences betⁿ predicted data pts (\hat{y}_i) & observed data pts (y_i)

$$SSE = \sum (\hat{y}_i - y_i)^2.$$

$$SST = SSR + SSE$$

R squared :- referred to as coefficient of determination, measure of how well a linear regression model fits a dataset.

- it represents prop. of variance in the response variable that can be explained by predictor var.
- it ranges from 0 to 1.
0 - response var. cannot be explained by predictor var. at all
1 - —||— perfectly —||—

$$R \text{ squared} = \frac{SSR}{SST}$$

e.g. $SSR = 137.5$

$SST = 156$

$R^2 = ?$

Adjusted R squared :-

- ~~measures~~ It adjusts statistic based on the no. of independent variables in model.
- It tells how well data fit a curve/line

$$R^2_{adj} = 1 - \left[\frac{(1-R^2)(n-1)}{n-k-1} \right]$$

- It adjusts ^{no. of} terms in line/curve.

n - no of pts in sample.

k - no of independent regressors.

- You need R^2 when working with samples.

- Always less than or equal to R^2

- not necessary when you have data from entire population

eg. A fund has a sample R squared

value close to 0.5 & it is doubtlessly offering higher risk adjusted returns with sample size of 50 & I predict find adjusted R square value.

$$n=50 ; k=5. R\text{-square} = 0.5.$$

- If you add useless variables ad. R^2 will decrease & if you add useful var. it will increase.