

## Logistic Regression

- It is also called "Classification model" and help to make predictions in cases where the output is a Categorical variable.
- Logistic Regression is the most easily interpretable of all classification models.

### Logarithmic Functions

→ Log with base 10 has following property

$$G = \log_{10} G$$

$$20 = \log_{10} 20$$

$$6^n = 20$$

$$(10^{\log 6})^n = 10^{\log 20}$$

$$10^{n \log 6} = 10^{\log 20}$$

$$\Rightarrow n \log 6 = \log 20$$

$$n = \frac{\log 20}{\log 6} \approx 1.67$$

#### 1. Product Property

$$\boxed{\log_b (a \cdot c) = \log_b a + \log_b c} \quad (\text{a,b,c positive no., } b \neq 1)$$

#### 2. Second the Quotient property

$$\boxed{\log_b \frac{a}{c} = \log_b a - \log_b c} \quad (\text{a,b,c positive, } b \neq 1)$$

#### 3. Third the power property

$$\boxed{\log_b (a^c) = c \log_b a}$$

$c$  = real no.,  $a$  &  $b$  positive no.,  $b \neq 1$

# Univariate Logistic Regression

- Binary Classification
- Sigmoid Function
- Likelihood Function
- Building a logistic regression model in python
- Odds and log odds

## Binary Classification

→ It is a most common use of logistic regression model.

Eg:-

1. Customer will default or not

2. Spam/ham example

→ In this classification we have only 2 categories

Note:- Simple boundary decision method would not work as some of the datapoint will become mis-classifying few data points.

→ Target variable has only 2 possible values.

### # Sigmoid Curve

→ Now the issue of Simple boundary decision method can be somehow solved by the use of probability.

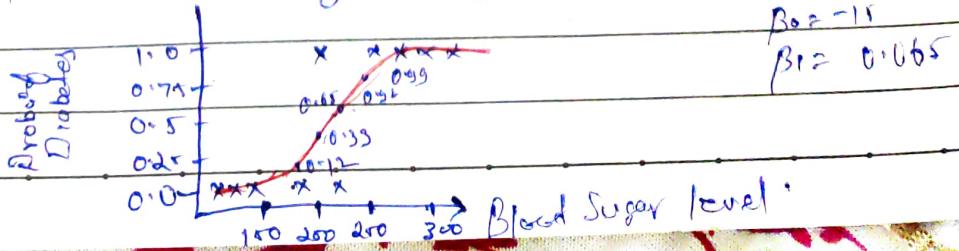
→ In method we assign low probability for the point below 180 and high probability for point above 250.

→ We draw a curve. One possible curve to draw is called "Sigmoid Curve".

Eg:-

$$Y(\text{Probability of Diabetes}) = \frac{1}{1+e^{-(B_0+B_1x)}}$$

By changing  $B_0$  &  $B_1$  we can get different Sigmoid Curve.

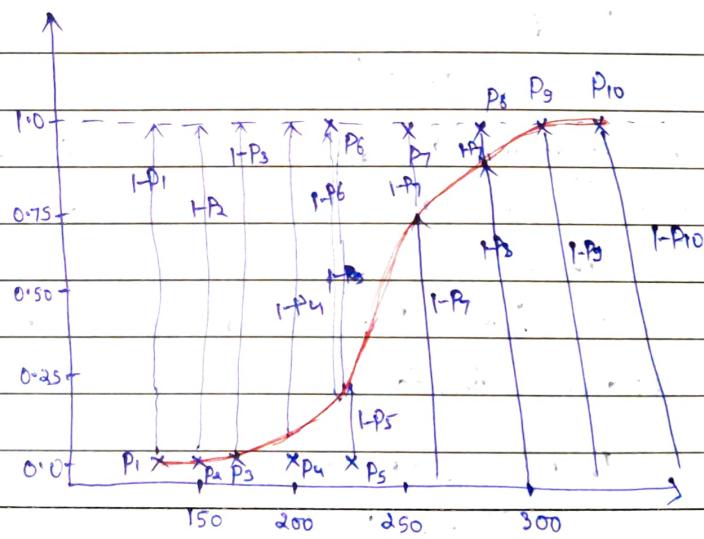


Note:- The main reason of not using straight line is that it is not steep enough. In the sigmoid curve at starting we have low probability value which rises suddenly and after that we have a lot of high values.

In the straight line the value rise from low to high ~~is~~ <sup>very</sup> uniform and hence "boundary" region where the probability transition from high to low is not present.

### # Finding the Best Fit Sigmoid Curve - I

→ It is a process of finding the **best fit Sigmoid curve** means finding  $B_0$  &  $B_1$  which best fits the dataset.



We need to maximize the product of these probability to find the best fitted line.

$$\text{Product} = (1-p_1)(1-p_2)(1-p_3)(1-p_4)(1-p_5)(p_5)(p_7)(p_8)(p_9)(p_{10}) \dots$$

$\downarrow$

This product is called **Likelihood function**.

$\left[ \text{At } [(1-p_i)(1-p_i)] \text{ for all non-diabetics } \right] * \left[ (p_i)(p_i) \text{ for diabetics} \right]$

This process where we vary the beta until we find the best fit curve for the probability of diabetes is called **Logistic regression**.

Point No.	1	2	3	4	5	6	7	8	9	10
Diabetes	No	No	No	Yes	No	Yes	No	Yes	Yes	Yes

$$(1-P_1)(1-P_2)(1-P_3)(1-P_4)(1-P_5)(1-P_6)(P_7)(P_8)(P_9)(P_{10})$$

### # Finding the Best Fit Sigmoid Curve - II

- Now we can randomly choose the value of  $\beta_0$  &  $\beta_1$  and find the respective Probability value.
- Now calculate the likelihood.
- Keep on keep on doing the same until the best fit line is not found. (max value of likelihood)

### \* Odds and Log Odds

#### → Equation of Logistic Regression

$$P = \frac{1}{1 + e^{-(\beta_0 + \beta_1 x)}}$$

- But it is difficult to understand the trend beta parameters.
- As we can't predict P when the value of x changes by few interval.
- The relationship beta P and x is too complex to see any apparent trends.

$$1-P = \frac{e^{-(\beta_0 + \beta_1 x)}}{1 + e^{-(\beta_0 + \beta_1 x)}}$$

$$\frac{P}{1-P} = \text{odd}$$

⇒ Ratio of Probability of an event occurring to the probability of the event not occurring.

$$\frac{P}{1-P} = e^{(\beta_0 + \beta_1 x)}$$

$$\log\left(\frac{P}{1-P}\right) = \beta_0 + \beta_1 x$$

$$\log\left(\frac{P}{1-P}\right) = \text{Log odds.}$$

$$\text{Ex:- } \frac{P}{1-P} = 4$$

( $P$  = probability of diabetic)  
 $(1-P)$  = prob. of not diabetic)

$$P = 4 - 4P$$

$$P = 4/5 = 0.80$$

$$P(\text{Diabetes}) = 4 * P(\text{No Diabetes})$$

Note:-

When the value of  $\alpha$  changes linearly (by const value)  
 then odd becomes multiple of previous odd value.

### Other form of the Logistic Regression Equation

#### 1. Sigmoid Curve (logit equation)

$$P = \frac{1}{1 + e^{-(\beta_0 + \beta_1 x)}}$$

Extremely low at start, extremely high at end & intermediate at middle

#### 2. Probit form of logistic regression

$$P = \Phi^{-1}(\beta_0 + \beta_1 x)$$

#### 3. Cloglog of logistic Regression

$$P = \log(-\log(1 - (e^{\beta_0 + \beta_1 x})))$$

→ Last 2 egn can also be used in logistic regression as their graph also gives the same trend.

Note:-

links = smf.families.links.logit = smf.GLM(y~train, (smf.odd-constant(x~train)).

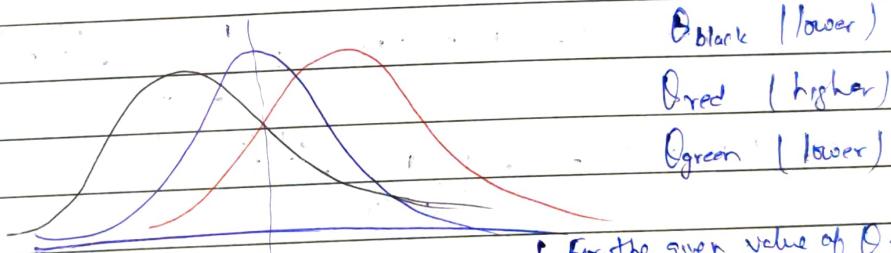
logit → family = smf.families.Binomial(link=links.logit)) logm.fit().summary()

Probit → family = smf.families.Binomial(link=links.probit)) logm.fit().summary()

cloglog → family = smf.families.Binomial(link=links.cloglog)) logm.fit().summary()

## Maximum Likelihood Cost Function

- In case of logistic regression we use Sigmoid function
- ↓ where  $P = \frac{1}{1+e^{-(B_0+B_1x)}}$
- We need to find the optimal value of  $B_0$  and  $B_1$  such that the likelihood function is maximised.
- The optimisation method used for to do so is called Maximum likelihood estimation or MLE
- A random variable is a variable whose possible values are numerical outcomes of a phenomenon or event
- There are 2 types of random variable
  - Discrete  $\rightarrow$  head-tail, dice
  - Continuous  $\rightarrow$  length, time, weight



$$p(x; \theta_0, \sigma) = P(x; \theta)$$

$$J(\theta) = p(x; \theta) = P(x_1, x_2, \dots, x_N; \theta)$$

For the given value of  $\theta$  the product  $P(x_1), P(x_2), \dots, P(x_N)$  will have to maximum value

Assumption:-

$$\text{Independence} = \prod_{i=1}^N P(x_i; \theta) = L(\theta) \quad (\text{Product of } P(x_1) \cdot P(x_2) \cdots P(x_N))$$

$$L(\theta) \log L(\theta) = \log L(\theta) = \sum_{i=1}^N \log P(x_i; \theta)$$

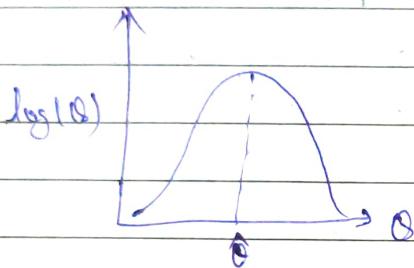
Note:- Taking log is allowed because log is a monotonic function  
 If we have a maxima or minima at  $L(\theta)$ , then if we will take log then  $\log L(\theta)$  will go hold at the same value of maxima & minima at along derivatives also.

## MLE For Continuous Distributions

→ MLE is basically a technique to find the parameters that maximize the likelihood of observing the data point assuming they were generated through a given distribution.

Ex:- For normal (Gaussian) distribution the parameters are the mean  $\mu$  and the standard deviation  $\sigma$ .

Probability Density Function of normal distribution

$$f(x_i) = \frac{1}{\sigma\sqrt{2\pi}} e^{-\frac{(x_i-\mu)^2}{2\sigma^2}}$$


$$L(\theta) = \sum_{i=1}^N \log(P(x_i; \theta))$$

Suppose this is a Gaussian distribution

$$P(x_i; \theta) = P(x_i; \mu, \sigma) = \frac{1}{\sigma\sqrt{2\pi}} e^{-\frac{(x_i-\mu)^2}{2\sigma^2}}$$

$$\log(P(x_i; \theta)) = \log\left(\frac{1}{\sigma\sqrt{2\pi}} e^{-\frac{(x_i-\mu)^2}{2\sigma^2}}\right)$$

$$\log(P(x_i; \theta)) = \log\left(\frac{1}{\sigma\sqrt{2\pi}}\right) + \log\left(e^{-\frac{(x_i-\mu)^2}{2\sigma^2}}\right)$$

$$J(\theta) = J(\mu, \sigma) = \sum_{i=1}^N \left[ \log\left(\frac{1}{\sigma\sqrt{2\pi}}\right) + -\frac{(x_i-\mu)^2}{2\sigma^2} \right]$$

Now we need to find  $\hat{\mu}$  &  $\hat{\sigma}$

- ① For iterative approach we use Gradient Descent
- ② Close form

Q)

$$J(\mu_0) = -\sum_{i=1}^N \log(\sigma \sqrt{2\pi}) - \frac{(x_i - \mu)^2}{2\sigma^2}$$

constant

$$\frac{dJ}{d\mu} = 0 \Rightarrow \frac{d}{d\mu} \left( -\sum_{i=1}^N \frac{(x_i - \mu)^2}{2\sigma^2} \right)$$

$$= \frac{-1}{2\sigma^2} \sum_{i=1}^N 2(x_i - \mu)(-1) = 0$$

$$= \sum_{i=1}^N x_i = \sum_{i=1}^N \bar{x} \Rightarrow (\bar{x} = \frac{\sum x_i}{N})$$

$\bar{x} = \frac{\sum x_i}{N}$

→ Sample Mean

$\hat{\sigma} = \sqrt{\frac{\sum (x_i - \bar{x})^2}{N}}$

### (MLE) Maximum Likelihood Estimation for Discrete Distributions

- Let us Bernoulli Distribution for discrete distribution.
- Implement MLE by finding the parameters of a Bernoulli Distribution.

Note:- Binary Logistic Regression is also a part Bernoulli Distribution.

N observation (Discrete)

Bernoulli Distribution ( $p$ ) {H, T}

$$J(p) = p(Y; p) = p(y_1, y_2, \dots, y_N; p) = \prod_{i=1}^N p(y_i; p)$$

$\downarrow$

$y_i$

Maximise the probability of certain 'Y' vector, for the given parameter ' $p$ '.

$p(y_i)$  = probability of  $y_1$  and so on.

All events are mutually independent hence product is used.

$J(p) = \log(L(p)) = \sum_{i=1}^N \log(p(y_i; p))$

$$y_i = \text{Head} = 1 \text{ or } 0$$

$$P(y_i; p) = p^{y_i} (1-p)^{1-y_i}$$

$y_i=1 \text{ Head } \Rightarrow (P)$   
 $y_i=0 \text{ Tail } \Rightarrow (1-P)$

$$J(p) = \sum_{i=1}^N \log [p^{y_i} (1-p)^{1-y_i}]$$

$$= \sum_{i=1}^N [\log(p^{y_i}) + \log((1-p)^{1-y_i})]$$

$$J(p) = \sum_{i=1}^N [y_i \log(p) + (1-y_i) \log(1-p)]$$

$$\frac{dJ(p)}{dp} = 0 \Rightarrow \sum_{i=1}^N \left[ \frac{y_i}{p} + \frac{(1-y_i)}{1-p} \right] = 0$$

$$0 = \sum_{i=1}^N \frac{y_i}{p} - \sum_{i=1}^N \frac{(1-y_i)}{(1-p)}$$

$$\Rightarrow \sum_{i=1}^N \frac{y_i}{p} = \sum_{i=1}^N \frac{(1-y_i)}{(1-p)}$$

$$\Rightarrow \frac{\sum_{i=1}^N y_i}{p} = \frac{\sum_{i=1}^N (1-y_i)}{1-p}$$

$$\text{Now; } \sum_{i=1}^N y_i = N_H \quad \& \quad \sum_{i=1}^N (1-y_i) = N_T$$

$$\Rightarrow \frac{N_H}{p} = \frac{N_T}{1-p}$$

$$\Rightarrow \frac{p}{1-p} = \frac{N_H}{N_T}$$

$$\Rightarrow p = \frac{N_H}{N}$$

$$\hat{p} = \frac{\sum_{i=1}^N x_i}{N}$$

Maximum Likelihood Estimation for Logistic Function  
 → Logistic Classification model is also known as Logit model or logistic regression.

$$P(Y|\beta) = \prod_{i=1}^N P(y_i|\beta)$$

$y_i$  = Male/Female  
 $x_i$  = length of hair

$$l(\beta) = \sum_{i=1}^N [y_i \ln(\beta) + (1-y_i) \ln(1-\beta)]$$

$$= \sum_{i=1}^N y_i \ln[P(\text{Female}/x_i)] + (1-y_i) \sum_{i=1}^N \ln[1 - P(\text{Female}/x_i)]$$

$$P_{\text{Female}/x_i}$$

$$P(\text{Male}/B_{10}=5)$$

$$P(\text{Female}/B_{10}=5)$$

$$P(\text{Female}/x_i) = \sigma(\beta_0 + \beta_1 x_i)$$

$$= \sigma(\beta^T x_i)$$

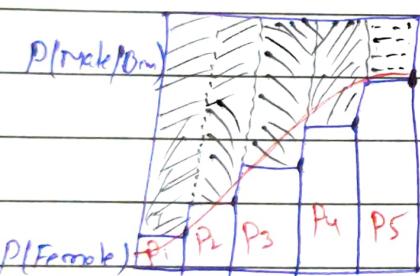
$$= \frac{1}{1 + e^{-\beta^T x_i}}$$

$$P(\text{Female}/x_i) = \frac{e^{\beta^T x_i}}{1 + e^{\beta^T x_i}}$$

$$P(\text{Male}/x_i) = \frac{1 - e^{\beta^T x_i}}{1 + e^{\beta^T x_i}}$$

$$= \frac{1}{1 + e^{\beta^T x_i}}$$

$$P(\text{Male}/x_i) = \frac{e^{-\beta^T x_i}}{1 + e^{-\beta^T x_i}}$$



$$l(\beta) = \sum_{i=1}^N y_i \log P(\text{Female}/x_i) + (1-y_i) \log (1-P(\text{Female}/x_i))$$

$$= \sum_{i=1}^N \left[ \log (1-P(\text{Female}/x_i)) + \sum_{i=1}^N y_i \log \frac{P(\text{Female}/x_i)}{(1-P(\text{Female}/x_i))} \right]$$

$$= \sum_{i=1}^N \log \frac{1}{1 + e^{\beta^T x_i}} + \sum_{i=1}^N y_i \beta^T x_i$$

$$J(\beta) = - \sum_{i=1}^N \log(1 + e^{\beta^T x_i}) + \sum_{i=1}^N y_i \beta^T x_i \rightarrow \text{Cost Function}$$

Using Gradient Descent Method to find parameter values.

$$\frac{dL(B)}{dB_i} = - \sum_{j=1}^N \frac{e^{B_m i}}{1 + e^{B_m i}} m_{i,j} + \sum_{j=1}^N y_{i,j} \pi_{i,j}$$

$$\hat{B}_m i = \beta_0 + \beta_1 \pi_{i,j}$$

$\pi_{i,j}$   
 $i = \text{Sample}$   
 $j = \text{feature}$

$$\frac{dL(B)}{dB_j} = \sum_{i=1}^N (y_i - P(\text{Female}|x_{i,j})) \pi_{i,j}$$

↑                    ↑  
 actual            estimated

$$e_i = y_i - P(\text{Female}|x_{i,j})$$

$$= [e_1, e_2, \dots, e_N] \begin{bmatrix} \pi_{1,j} \\ \pi_{2,j} \\ \vdots \\ \pi_{N,j} \end{bmatrix} = [\pi_{1,j}, \pi_{2,j}, \dots, \pi_{N,j}] \begin{bmatrix} e_1 \\ e_2 \\ \vdots \\ e_N \end{bmatrix}$$

$$\nabla L(B) = \begin{bmatrix} \frac{dL(B)}{dB_0} \\ \frac{dL(B)}{dB_1} \end{bmatrix} = \begin{bmatrix} \pi_{1,0}, \pi_{2,0}, \dots, \pi_{N,0} \\ \pi_{1,1}, \pi_{2,1}, \dots, \pi_{N,1} \end{bmatrix} \begin{bmatrix} e_1 \\ e_2 \\ \vdots \\ e_N \end{bmatrix} = X^T e$$

$$X \rightarrow N \times (d+1)$$

Read  $\Rightarrow$  Obtaining the parameter : Newton's Method  
 (Addition before  $\log$  odds & odds )

## Obtaining the Parameters: Newton's Method

- Newton-Raphson method is a second order technique.
- By using this method we will compute the value of betas.

$$\frac{dL(\beta)}{d\beta_j} \quad \text{Good Descent: } \beta_j^{\text{new}} = \beta_j^{\text{old}} - \eta \frac{dL(\beta)}{d\beta_j}$$

$$\text{Newton: } \beta_j^{\text{new}} = \beta_j^{\text{old}} - H^{-1} \frac{dL(\beta)}{d\beta_j} \quad (H^{-1} = \text{Hessian matrix})$$

it is second order derivative terms.

$$\frac{dL(\beta)}{d\beta_j} = \sum_{i=1}^N (y_i - P(\text{Female}|\alpha_i)) m_{ij}$$

$$\frac{dL(\beta)}{d\beta_k d\beta_j} = - \sum_{i=1}^N \frac{(1 + e^{\beta_{\alpha_i}}) e^{\beta_{\alpha_i}} m_{ik} - e^{\beta_{\alpha_i}} e^{\beta_{\alpha_i}} m_{ik} m_{ij}}{(1 + e^{\beta_{\alpha_i}})^2}$$

$$= \sum_{i=1}^N P(\text{Female}|\alpha_i) + (P(\text{Female}|\alpha_i))^2$$

$$\boxed{\frac{dL(\beta)}{d\beta_k d\beta_j} = \sum_{i=1}^N (P(\text{Female}|\alpha_i)) (1 - P(\text{Female}|\alpha_i))}$$

## Multivariate Logistic Regression (Model Building)

→ When only one independent variable is not enough to get a better result then we need to add more feature.

### Univariate Logistic Regression

$$P = \frac{1}{1 + e^{-(B_0 + B_1 m)}}$$

### Multivariate Logistic Regression

$$P = \frac{1}{1 + e^{-(B_0 + B_1 m_1 + B_2 m_2 + B_3 m_3 + \dots + B_n m_n)}}$$

\* Build a multivariate logistic regression model in Python

\* Conduct feature selection for logistic regression

\* Automated methods: RFE - Recursive Feature Elimination

\* Manual Method: VIF & p-value check

# Steps to Build a Logistic Regression Model

Two most commonly used metrics to evaluate a Classification Model

1. Sensitivity
2. Specificity

\* Sensitivity:- (True Positive Rate) (Recall)

It is defined as the ratio of correctly predicted by total no. of actual Predicted positive output of a model.

Ex.

Actual / Predicted		Predicted	
		Not Churn	Churn
Not Churn	Not Churn	3269	966
	Churn	595	692

From the above table we can say that the model has actually predicted 595+692 no. of Churn but out of which the correctly predicted is 692.

Actual / Predicted		Not Churn	Churn
Not Churn	Not Churn	True Negative	False Positive
	Churn	False Negative	True Positive

$$\text{Sensitivity} = \frac{\text{No. of actual Yeses Correctly Predicted}}{\text{Total No. of actual Yeses}}$$

$$\text{Sensitivity} = 100 \times \frac{692}{(692 + 595)} = 53.768\%$$

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

Probability of an actual 'yes' cases → predicted correctly.

\* Specificity

$$\text{Specificity} = \frac{\text{No. of actual Nos Correctly Predicted}}{\text{Total Number of actual Nos}}$$

$$\text{Specificity} = \frac{TN}{TN + FP}$$

$$\text{Specificity} = \frac{3269}{3269 + 366} \times 100 = 89.931\%$$

\* Accuracy

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

\* Positive Rate :- Predicted +ve when it was -ve

$$\text{Positive Rate} = \frac{FP}{TN + FP}$$
 (False positive Rate)  $(1 - \text{Specificity})$

\* Positive Predictive Value (Precision)

$$\text{Positive Prediction Rate} = \frac{TP}{TP + FP}$$

Probability that predicted 'yes' is actually 'yes'

True +ve  
Total Predicted +ve

\* Negative Predictive value

$$\text{Negative Prediction Rate} = \frac{TN}{TN + FN}$$

Note:- A Good model is the one in which TPR is high and FPR is low.

(Harmonic Mean)

$$F = 2 \cdot \frac{\text{precision} \cdot \text{recall}}{\text{precision} + \text{recall}}$$

High recall, low precision  $\rightarrow$  optimistic Model

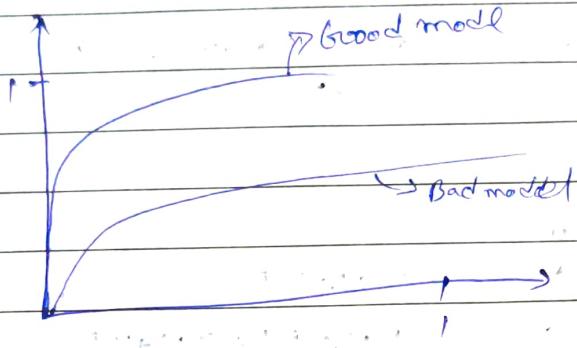
M	T	W	T	F	S	S
Page No. ....						
Date						YOUVA

## ROC Curve (

### Receiver Operating Characteristics

- Used in Radar
- Developed by Electricals Engineer in 1950's.
- It is used in all the classification model.

→



- Good model has higher area under curve and vice versa.
- Good ROC touches the upper left corner of the graph.
- There is a trade-off b/w the True Positive Rate and False Positive Rate, or simply a tradeoff b/w

### Sensitivity and Specificity

- When we plot TPR against FPR we get a graph which shows the trade-off b/w them and this curve is known as ROC.

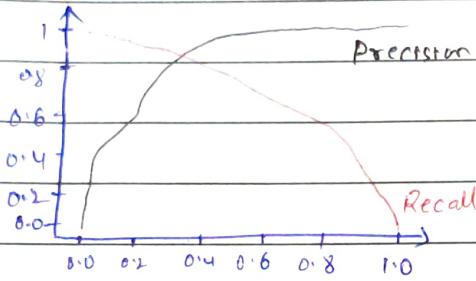
Note:- 1. If TPR increases then FPR also increases.

2. For a completely random model, the ROC curve will pass through the upper 45° line.
3. When Sensitivity is increasing,  $(1 - \text{Specificity})$  and since  $(1 - \text{Specificity})$  is increasing, it simply means that Specificity is decreasing.
4. Cutoff point is taken from X-axis.

Read (FP)

M	T	W	T	F	S	S
Page No.:						
Date:						YOUVA

- When we draw a curve of both accuracy, Sensitivity, Specificity from different value of cut-off (prob) we choose the point where all the metrics (accuracy, Sensitivity & Specificity) are fairly descent and almost equal.



- AUC is near to 1 then model is working satisfactorily  
→ AUC is low (around 0.5) the the model is not working properly and just guessing randomly.  
→ If we want to create a balance model, without business need, then we can use an ROC curve to see the tradeoff between sensitivity and specificity and accordingly choose an optimal trade off point where both these values along with accuracy are descent.

## Logistic Regression

1. Types of Logistic Regression
  - a. Multinomial Logit
  - b. Binary Logit

## Classification Techniques

1. Support Vector Machine
2. Neural Network
3. Random forest
4. Gradient Boosting
5. Deep Learning

Note:- Logistic regression is most common classification technique instead of other methods because

1. Easy to understand and intuitive variable explanation  
In other classification techniques it becomes very complex to make understand the features of the model to the manager.

2. Linear Relationship With log of odds.

As probability increases the log of odds linearly increases and linearity is easy to understand in comparison to other complex relationships.

## 2. Logistic Regression Nuances

### a. Sample Selection (Most important)

1. Seasonal or Cyclical fluctuations populations (Big Billion Day, Diwali)
2. Representative population (used sample base on recent Business Strategy)
3. Rare incidence population (Fraud detection)

### b. Segmentation.

We are building different models based on different nature

M	T	W	T	F	S	S
Page No.:						YOUVA
Date:						

of the 8 imputations are behaving differently; which is made into 8 segment of WPs  
 → We combine all these model to make it a single model

Ques

Note:- Building Segmentation for logistic regression is different from other Segmentation as it is because the Segmentation of logistic regression is build in order to ensure that the overall predicted power of Segmented System is higher than the predictive power of a Single model.  
 → We need to make sure that the Segmentation variable is done in such a way that the predictors are slightly different or has different power across the different Segment.

### c. Variable Transformation

→ How we transform a variable before making a model.

#### 1. Dummy Variable Transformation

Adv:- [ Dummy Variable becomes stable in case of continuous variables makes model stable.]

Disadv:- Changing Continuous variables to dummies, all the data will be compressed into few categories and that might result in Water Clumping

#### 2. Weight of Evidence (WoE) transformation

→ It is generally used in finance sector.

$$\text{WOE} = \log \left[ \frac{\text{good in bucket}}{\text{Total Good}} \right] - \log \left[ \frac{\text{bad in bucket}}{\text{Total Bad}} \right]$$

$$- \quad \text{WOE} = \log \left( \frac{\text{Y. of Good}}{\text{Y. of Bad}} \right)$$

→ Calculating woe values for fine binning and coarse binning

→ The importance of woe for fine binning and coarse binning

→ The usage of woe transformation

### Advantage of WOE

1. It reflects group identity (trend of distribution)
2. It helps in treating missing values logically for both types of variables.
3. Variable remains stable over time.

Note:- The graph of WOE should follow increasing or decreasing trend across bins.

If trends are not monotonic, then we need to combine the bucket/bins of that variable and then calculate the WOE again.

→ Disadvantage → small or dummy variable.

→ We may end doing some score clumping which decrease the model efficiency in small band.

→ Decrease in predictive power.

### 3. Information Value

→ It is an important indicator of predictive power

$$IV = WOE \times \left( \frac{\text{Good in bucket}}{\text{Total Good}} - \frac{\text{Bad in Bucket}}{\text{Total Bad}} \right)$$

## 3. Continuous Variable Transformation

### Technique

1. Spline transformation

2. Interaction Variable

3. Mathematical transformation

4. Principal component transformation

M	T	W	T	F	S	S
Page No.						
Date:	18/01/22				YOUNA	

## Advantage of WOE

1. It reflects group identity (trend of distribution)
2. It helps in treating missing values logically for both types of variables.
3. Variable remains stable overtime.

Note - The graph of WOE should follow increasing or decreasing trend across bins.

→ If trends are not monotonic, then we need to combine the bucket/bins of that variable and then calculate the WOE again.

→ Disadvantage → small as dummy variable.

→ We may end doing some score clumping which decrease the model efficiency in small band.

→ Decrease in predictive power.

## 3. Information Value

→ It is an important indicator of predictive power

$$IV = WOE \times \left( \frac{\text{Good in bucket}}{\text{Total Good}} - \frac{\text{Bad in Bucket}}{\text{Total Bad}} \right)$$

## 3. Continuous Variable Transformation

### Technique

1. Spline transformation

2. Interaction Variable

3. Mathematical transformation

4. Principal component transformation

# Logistic Regression

## 1. Types of Logistic Regression

## 2. Logistic Regression Nuances

### a. Sample Selection

### b. Segmentation

### c. Variable Transformation

#### 1. Dummy Variable transformation

#### 2. Weight of evidence (WOE) transformation

#### 3. Continuous Variable transformation

#### 4. Interaction Variables

#### 5. Splines

#### 6. Mathematical transformation

#### 7. Principal Component transformation

## 3. Model Evaluation (performance) measures

### (a) Discriminatory Power

### (b) Accuracy

### (c) Stability

→ KS Statistics

→ Gini (Receiver operating characteristics)

→ Rank ordering

→ Sensitivity

→ Specificity

→ Sensitivity

→ Specificity

→ Compare actual versus predicted by model

## 4. Model Validation

## 5. Model Tracking or model governance

Read about → KS Statistics

→ Rank ordering

## Challenges in Logistic Regression.

### 1. Low event rate

If it is a scenario in which we have very few true condition in the dataset.

Ex:- Fraud in credit card use

Method to Fix it

- (a) Use of Sampling where we create a balanced sampling and increase the sampling rate and perform cross validation on it.
- (b) Collect a data over wider time frame and get more of the 1's.
- (c) Change the definition of 1's. Suppose we assume failure do increase machine failure chances.  
Ex:- Credit card default payment Then in this case we can use multiple times credit card failure & late payment as an early sign of default.

### 2. Missing values

We need to find that missing value is random or it's due to any pattern.

Method to be used

- (a) Imputation using ~~recent~~ WOE
- (b) Imputation using Mean
- (c) Imputation using Median
- (d) Imputation using predictive pattern
- (e) Markov Chain Monte Carlo
- (f) Expectation Maximisation

### 3. Truncated Data

Mean in this process we don't analyse how it would have behaved if accepted.

Ex: Rejecting a credit card application based on a predefined parameter then in this case we don't know how that customer have performed after giving the credit card.

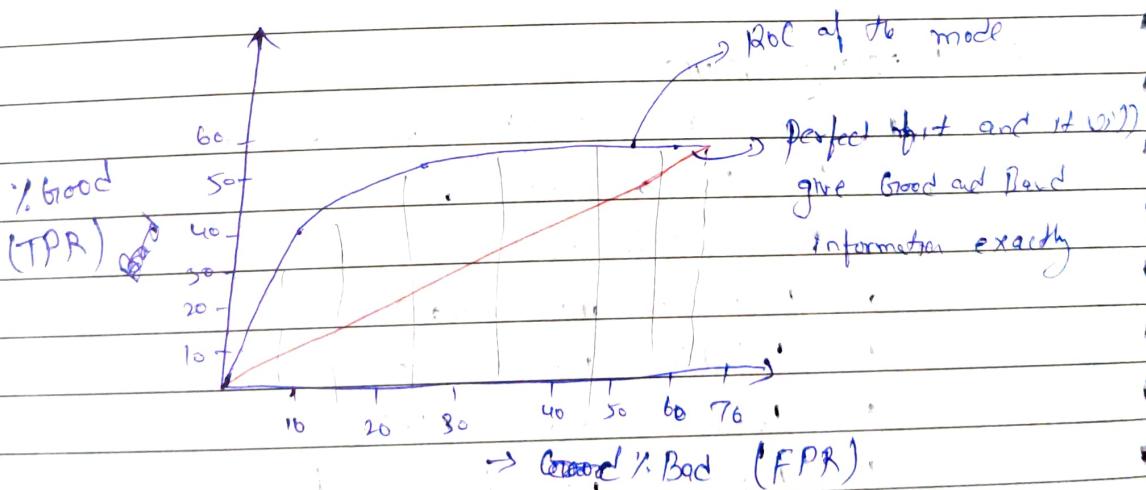
Methods:

- Reject Inference
- Hedging Correction Method

### Model Evaluation Method.

1. Gini (ROC) (Receiver Operating Characteristics)

→ GI is used to find Good and Bad present in the model dataset.



→ Point on the ROC curve of a model gives us no. of Good and Bad at that particular point.

→ The area under the curve is called Gini of the model.

→ More steeper the ROC curve higher will be the Gini of the model.

$$\text{Gini} = 2 \times \text{Area Under ROC Curve} - 1$$

## 4. Model Validation

- Every sample dataset has its own predictive power.
- We split our dataset into Train and Test whenever we build any model.
- The "real" test of model is when we test our model on completely different set of data.

### Method of Model Validation

#### 1. In-Sample Validation

- In this method we use 70% dataset for training and 30% dataset for testing the model.

- We check whether our model gives the same performance on 30% of test dataset or not.

#### 2. Out-of-time Validation

- In this process suppose we took a data from year 2015 to 2019 to train our model and we test our model on 2013, 2014, 2020, 2021 dataset and check how our model is working on it.

#### 3. K-cross Validation (Bagging)

- It is based on the concept of **Bagging and Boosting**.
- Suppose we don't have enough dataset and in this case the model can overfit by memorizing the data and its performance will fall on test dataset.
- To fix this problem we can take various <sup>(many times)</sup> 30% data from the overall sample and testing the model.
- In this process it's not the same 30% dataset used for validation but instead different set of 30% dataset and check model is working well on these different set of 30% dataset.
- It helps us to understand the stability of the model <sup>across</sup> on different set of sample.
- K-iteration means k time different set of data will be used for training and test the model. (**K-fold cross validation**)

## 4. Stability

- Model stability is most important part of a model.
- If models are not stable then we will not be able to make right decision over-time.
- Stability is as important than predictive power and in some cases stability is more important than predictive power.
- Even if the predictive power of the model is not that good but need to be stable in order to get work over-time to take decision.

### Factor on which Stability Depends

#### (a) Performance Stability

- All the factor like <sup>below</sup> should have similar value on dev/test sample for a model.
  - In-Sample validation
  - Gini
  - Capture Rate
  - Sensitivity
  - Specificity
- Result of in-sample validation approximately match those of out-of-time validation.

#### (b) Variable stability

##### → Variable distribution Stability

- If a dataset has 8 feature then all the feature should be present on train and test dataset.

##### → Population stability index (PSI)

- Percentage of distribution of feature in the train and test dataset should be nearby.

- The distribution of amount of feature should be stable between train and test dataset.

## Good Model

- The sample used for model building hasn't changed too much and has the same general characteristics.

### (C) Predictive Pattern

→ While calculating the WOE for different group we need to check the WOE of different groups in test and train set should be same.

Note:- 1. Sgn of WOE should not change on test dataset

2. We calculate the PSI of WOE and if change in PSI of <sup>WOE of</sup> a group is less than 0.1 on test set then <sup>we say</sup> ~~it is accepted~~ <sup>it is stable</sup> and if it is more than 0.25 it is not ~~accepted~~ stable.

Work:- A Good model will be stable when

i. Performance Stability:-

Result of in-sample validation approximately match those of out-of-time validation.

ii. Variable Stability:-

The sample used for model building hasn't changed too much and has the same general characteristics.

## 5. Model Tracking or Model Governance

- We track the performance of a model over the time with new dataset to check whether our model works fine or not.
- If the performance decreases by slight then we can do "recalibration" else we will have to build a model.

## 6. Model Recalibration

- In this method we don't build the model again we just update the coefficient of the variables present in the model.
- It is used which the Gini of the model is decrease by 0.1 or below and
- If over the time Gini decreases by 0.25 then we need to rebuild the model.

SGTL - ML - C34 - CC - G104

\* Closed Form Soln.

→ It is a equation where we put value of  $B_0, B_1$  and find value of  $y$ .

$$\omega^* = (X^T X)^{-1} X^T y$$

→ Inverse  $(X^T)$  is not always possible

→ N-Cube operation it takes time.

\* Why Approximation

→ Because it is not possible to determine  $\omega$  using closed form soln.

$$E(\omega) = \frac{1}{2} \sum_{i=1}^N [g(x_i; \omega) - y_i]^2$$

Find  $\min E(\omega)$

\* Logistic regression brings  $y$  value from  $-\infty$  to  $+\infty$  between 0 to 1 and for that we used Sigmoid function.

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

\* Accuracy =  $\frac{\# \text{Correct Prediction}}{\# \text{Total Predictions}}$

Is accuracy a good metric? No.

(Q1. Why can't linear regression is used in place of logistic regression for binary classification? (Bc)

→ 1. Distribution of Error term

- The distribution of data in the case of linear and logistic regression is different

- Linear regression assumes that error terms are normally distributed.

→ In case of binary classification, this assumption does not hold true.

### 2. Model Output

- In linear regression the output is a continuous.

→ In BC, linear regression may predict the value that can go beyond 0 and 1.

- If we want our output as probability then its output range is restricted to 0 and 1 and based on it the model output will be 0 and 1.

### 3. Variance of Residuals Errors

- Linear regression assume that the variance of random error is constant and this assumption is also violated in case of logistic regression.

(Q2. What are the outputs of logistic model and logistic function?

→ Logistic model output the logits. i.e log odds; and the logistic function output the probability.

$$\text{Logistic Model} = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \beta_3 x_3 + \dots + \beta_n x_n$$

$$\text{Logistic Function} = f(x) = \frac{1}{1 + e^{-(\beta_0 + \beta_1 x_1 + \beta_2 x_2 + \beta_3 x_3 + \dots + \beta_n x_n)}}$$

Q3. How to interpret the result of a logistic regression model?  
Or, what are the meaning of the different betas in a logistic regression model.

- $\beta_0$  is the baseline in a Logistic regression Model.
- It is log odds for an instance when all the attributes ( $X_1, X_2, X_3, \dots, X_n$ ) are zero.
- $\beta_0$  is the log odds for an instance when none of the attributes is taken into consideration.
- All the other Betas are the values by which the log odds change by a unit change in a particular attribute by keeping all other attributes fixed or unchanged (control variables).

Q4. What is odd ratio?

- $OR_{X_i; X_0} = e^{\sum_{i=1}^k \beta_i (X_{ii} - X_{0i})}$
- $X_i$  and  $X_0$  stand for two different groups for which the odds ratio needs to be calculated.
- $X_{ii}$  stands for the instance 'i' in the group  $X_i$ .
- $X_{0i}$  stands for the instance 'i' in the group  $X_0$ .
- $\beta_0$  stands for the coefficient of the logistic regression model.
- Baseline is not included in the formula.

Q5. What is the Maximum Likelihood Estimator (MLE)?

- It chooses those set of unknown parameters (estimator) that maximise the likelihood function.
- To find the MLE is to use calculus and setting the derivative of the logistic function with respect to an unknown parameter to zero and solving gives the MLE.

Note:- MLE and ordinary square estimation gives the same results for linear regression if the dependent variable is assumed to be normally distributed. MLE doesn't assume anything about independent variable.

Q.6. What are the different method of MLE and When each method preferred.

→ In case of logistic regression There are 2 approaches to MLE:

1. Unconditional Method
2. Conditional Method.

These are algorithms that we different likelihood funts.

#### Conditional Method

1. gt don't estimate unwanted
2. gt can't be done with of joint probability.
3. gt gives unbiased in such case.
4. Statisticians suggest that Conditional MLE is used When in doubt.

#### Unconditional Method

1. gt estimates the values of unwanted parameter also.
2. gt can directly be developed with joint probabilities.
3. If the no. of parameter is higher relative to the no. of instances, then it gives biased result.

Q.7. What is the Output of MLE program.

→ Maximised likelihood value.

→ Estimated likelihood value.

Q8. Why can't we use Mean Square Error (MSE) as a cost function for logistic regression.

- In logistic regression "we" use Sigmoid function and perform a non-linear transformation to obtain the probabilities.
- Squaring this non-linear transformation will lead to non-convexity with local minima.
- Finding the global minimum in such cases using gradient descent is not possible.
- Due to this reason MSE is not suitable for logistic regression.
- Cross-entropy or log loss is used as a cost function for logistic regression.
- In cost function for logistic regression, the confident wrong predictions are penalized heavily and confident right predictions are rewarded less.
- By optimising this cost function, convergence is achieved.

Note:- Accuracy gives equal importance to both False positive and False negatives. Which is not required in many cases.

## Cost Function

- \* Cost Function
- \* Min/Max of a Cost Function
- \* Differentiation of a function
- \* Types of Minimisation
  - Constrained
  - Unconstrained
- \* Solution of Unconstrained minimisation
  - Closed Form
  - Gradient Descent

## # Cost Functions

→ It helps us to reach the optimal solution.

### Rule of Differentiation

$$y = am^b$$

$$\frac{dy}{dm} = b a m^{b-1}$$

$$y = \sin mx$$

$$\frac{dy}{dm} = \cos mx$$

$$y = \cos mx$$

$$\frac{dy}{dm} = -\sin mx$$

$$y = e^{mx}$$

$$\frac{dy}{dm} = e^{mx}$$

$$y = \ln mx$$

$$\frac{dy}{dm} = \frac{1}{m}; \text{ for } m > 0$$

$$\text{Ex: } ① y = 10m^2 \Rightarrow \frac{dy}{dm} 2 \times 10m^{2-1} = 20m$$

$$\text{Slope at } m=4 \Rightarrow \frac{dy}{dm} = 20m = 80$$

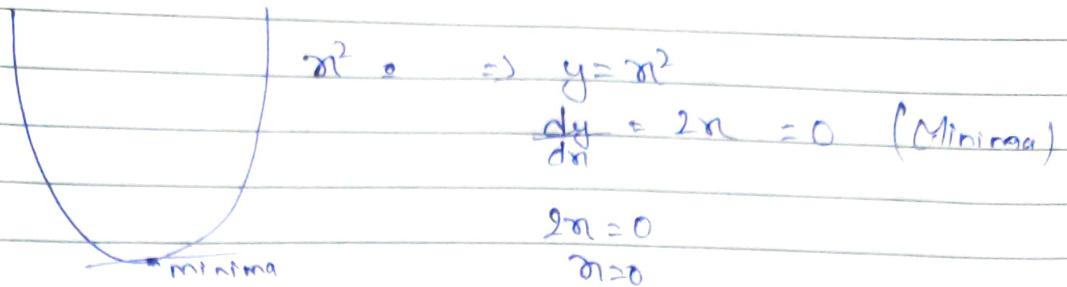
$$② y = 10 = 10m^0 \quad (m^0 = 1)$$

$$\frac{dy}{dm} = 0 \times 10m^{0-1} = 0$$

$$③ y = 2m^2 + 3m + 4$$

$$\frac{dy}{dm} = 2 \times 2m^{2-1} + 1 \times 3m^{1-1} + 0$$

$$\frac{dy}{dm} = 4m + 3$$



$$y = \theta_1^2 + \theta_2^2$$

$$\frac{dy}{d\theta_1} = 2\theta_1 = 0$$

$$2\theta_1 = 0$$

$$\theta_1 = 0$$

$$y = \theta_1^2 + \theta_2^2$$

$$\frac{dy}{d\theta_2} = 2\theta_2 = 0$$

$$2\theta_2 = 0$$

$$\theta_2 = 0$$

$$* J(m, c) = [y_1 - (m\theta_1 + c)]^2 + \dots$$

$$\frac{\partial J}{\partial m} = 0$$

$$\frac{\partial [y_1 - (m\theta_1 + c)]^2 + \dots}{\partial m}$$

$$2[y_1 - (m\theta_1 + c)](-\theta_1)$$

$$= 0$$

$$\frac{\partial J}{\partial c} = 0$$

$$\begin{aligned} & \partial(y_1 - (m\theta_1 + c))^2 (-1) + \dots / \partial c \\ & 2(y_1 - (m\theta_1 + c))(-1) \end{aligned}$$

$$= 0$$

Note :- To minimise a func

$$[y' = 0 \text{ and } y'' > 0]$$

To maximise a func

$$[y' = 0 \text{ and } y'' < 0]$$

→ For Straight line we find the sum of squared Error and find the optimal value of 'm' & 'c' by differentiating the eqn of straight line w.r.t to 'm' & 'c'

## Naive Bayes

- $g_f$  is another type of Supervised Classification
- $g_f$  is a probabilistic classifier which returns the probability of a test point belonging to a class rather than the label of the test point.
- $g_f$  is a type of "Name Based" classifier.  
Ex: Spam or Ham
- $g_f$  is based on Bayes' Theorem

## Conditional Probability

$$P(\text{Head}) = \frac{\# \text{Favourable Outcomes}}{\# \text{Total Outcomes}}$$

$$\left\{ P(\text{Play}) = \frac{20}{30} = 66\% \right.$$

$$\left. P(\text{Play}|\text{Rain}) = \frac{1}{10} = 10\% \right.$$

$$P(\text{Spam}) = \frac{1}{5} = 20\%$$

$$P(\text{Ham}) = \frac{4}{5} = 80\%$$

$$P(\text{Spam}|\text{Word} = \text{VIAGRA}) = 70\%$$

### Prior

$$P(\text{Play}) = 66\%$$

$$P(\text{Spam}) = 20\%$$

### Posterior

$$P(\text{Play}|\text{Rain}) = 10\%$$

$$P(\text{Spam}|\text{Word} = \text{Viagra}) = 70\%$$

## Conditional Probability

A and B are events

$P(A)$  = Probability of an email being spam

$P(B)$  = Probability of the word being VIAGRA

$$P(A) = 20\%$$

$$P(A|B) = 70\%$$

To understand the relative probabilities.

## Naive Bayes

- It is another type of Supervised Classification.
- It is a probabilistic classifier which returns the probability of a test point belonging to a class rather than the label of the test point.
- It is a type of "Name Based" classifier.  
Ex: Spam or Ham
- It is based on Bayes' Theorem.

## Conditional Probability

$$P(\text{Head}) = \frac{\# \text{ Favourable Outcomes}}{\# \text{ Total Outcomes}}$$

$$\left\{ \begin{array}{l} P(\text{Play}) = \frac{20}{30} = 66\% \\ P(\text{Play}|\text{Rain}) = \frac{1}{10} = 10\% \end{array} \right.$$

$$P(\text{Spam}) = \frac{1}{5} = 20\%$$

$$P(\text{Ham}) = \frac{4}{5} = 80\%$$

$$P(\text{Spam}|\text{Word} = \text{VIAGRA}) = 70\%$$

### Prior

$$P(\text{Play}) = 66\%$$

$$P(\text{Spam}) = 20\%$$

### Posterior

$$P(\text{Play}|\text{Rain}) = 10\%$$

$$P(\text{Spam}|\text{Word} = \text{Viagra}) = 70\%$$

## Conditional Probability

A and B are events

$P(A)$  = Probability of an email being spam

$P(B)$  = Probability of the word being VIAGRA

$$P(A) = 20\%$$

$$P(A|B) = 70\%$$

→ It is need to understand the relative probabilities.

# Bayesian estimate Vs Maximum Likelihood estimate

M	T	W	T	F	S	S
Page No.:						
Date:	26/1/22				YOUVA	

## Bayes' Theorem

(B) Sachin Scores a Century  
Sachin Doesn't Score a Century (B)  
Total

	Indra Win (A)	Indra Lose (A)	Total
(B) Sachin Scores a Century	10	2	12
Sachin Doesn't Score a Century (B)	50	38	88
Total	60	40	100

Prior

$$P(A) = \frac{60}{100} \quad P(B) = \frac{12}{100}$$

Joint Probability

$$P(A \cap B) = P(\text{Indra Win and Sachin's Century})$$

$$P(A \cap B) = \frac{10}{100} = P(B|A)$$

Conditional Probability

$$P(A|B) = \frac{10}{12} = \frac{\frac{10}{100}}{\frac{12}{100}} = \frac{P(A \cap B)}{P(B)} \quad \left. \begin{array}{l} \cap = \text{intersection} \\ \cap = \text{and} \end{array} \right.$$

$$\boxed{P(A \cap B) = P(A|B) \cdot P(B)} \\ = \frac{\frac{10}{12}}{\frac{10}{100}} = \frac{10}{100}$$

$$\boxed{P(A \cap B) = P(B|A) \cdot P(A)} = \frac{10}{60} \times \frac{60}{100} = \frac{10}{100}$$

$$P(A|B) \cdot P(B) = P(B|A) \cdot P(A)$$

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

→ Bayes' Theorem

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

A and B are independent

$$P(A \text{ and } B) = P(A|B) \cdot P(B) = P(A) \cdot P(B)$$

## Naive Bayes - With One Feature

- It is a probability classifier.
- It uses probability as a core criteria to classify.
- Common use is "Spam & Ham filter".

Ex:-

Types of Mushroom	Cap-shape	Cap Surface	Cap Color	Bruises	Odor
Edible	Convex	Smooth	White	Brunies	Almond
"	"	"	"	"	Anise
"	Bell	Scaly	"	"	Almond
"	"	"	Yellow	"	Anise
Poisonous	Convex	Smooth	Brown	"	Dingy
Edible	"	"	Gray	No	None
Poisonous	"	Fibrous	"	"	Foul
"	"	Scaly	"	"	"
"	Flat	Smooth	White	Brunies	"
"	Flat	Scaly	Pink	"	None
Edible	Convex	Fibrous	Brown	No	"
Poisonous	Bell	Scaly	Yellow	No	"
Edible	Bell	Fibrous	White	No	"
Poisonous	Knobbed	Scaly	Cinnamon	No	Musty

$$P(C=\text{Edible} | m=\text{Convex}) = \frac{P(m=\text{Convex} | C=\text{Edible}) * P(C=\text{Edible})}{P(m)}$$

$$P(C=\text{Poisonous} | m=\text{Convex}) = \frac{P(m=\text{Convex} | C=\text{Poisonous}) * P(C=\text{Poisonous})}{P(m)}$$

$P(m)$  is a common factor so need not be considered

$$\therefore P(\text{Edible}) = 7/14 = 0.5$$

$$\therefore P(\text{Poisonous}) = 7/14 = 0.5$$

$$P(\text{Cap-shape} = \text{Convex} | \text{edible} = \text{Yes}) = \frac{4}{7}$$

$$P(\text{Cap-shape} = \text{Convex} | \text{edible} = \text{No}) = \frac{3}{7}$$

$$P(\text{edible} = \text{Yes} | \text{x} = \text{Convex}) = \frac{P(\text{x} = \text{Convex} | \text{edible} = \text{Yes}) * P(\text{edible} = \text{Yes})}{P(\text{x})}$$

$$= \frac{\frac{4}{7} * \frac{1}{2}}{P(\text{x})} = \frac{4}{14 * P(\text{x})}$$

$$P(\text{edible} = \text{No} | \text{x} = \text{Convex}) = \frac{P(\text{x} = \text{Convex} | \text{edible} = \text{No}) * P(\text{edible} = \text{No})}{P(\text{x})}$$

$$= \frac{\frac{3}{7} * \frac{1}{2}}{P(\text{x})} = \frac{3}{14 * P(\text{x})}$$

Note:-  $P(\text{x})$  is common and acts as a scaling factor and hence can be removed by comparing the two classes

$$\frac{4}{14 * P(\text{x})} > \frac{3}{14 * P(\text{x})} \Rightarrow \frac{4}{14} > \frac{3}{14}$$

$$\Rightarrow P(C_i | x) = \frac{P(x|C_i) P(C_i)}{P(x)}, C_i = \text{Classes}$$

$x = \text{Features of the data point}$

→ Probability is simply calculated by counting the no. of instances/occurrences for categorical data.

→ The class assigned to the new test point is the class for which  $P(C_i|x)$  is greater.

→ Naive Bayes would properly in text classification

## Conditional Independence in Naive Bayes

→ When we have multiple features then we multiply each feature with class of different features.

$$P(\text{edible} | \text{x} = (\text{convex}, \text{smooth}))$$

$$\propto P(x = (\text{convex}, \text{smooth}) | \text{edible}) * P(\text{edible})$$

$$= P(\text{smooth} | \text{edible}) * P(\text{convex} | \text{edible}) * P(\text{edible})$$

Note:- Convex and Smooth are ~~not~~ independent events.

P

→ When we divide the equation on separate feature like below is called Naive assumption and hence called Naive Bayes.

$$P(C | x_1, x_2, \dots) = P(x_1 | C) * P(x_2 | C) \dots * P(x_n | C)$$

→ In other word Naive Bayes follows an assumption that the variables are conditionally independent.

Ex:-  $P(A \text{ and } B | C)$

If  $P(A|C)$  is same for all values of B,

$P(B|C)$  is same for all values of A

then there is Conditional independence betw A and B given C.

This is when  $P(A \text{ and } B | C) = P(A|C) * P(B|C)$  imply that A is not conditioned on B or vice versa.

When we divide the numerator with denominator in Bayes theorem then we get **Base probability**

M	T	W	T	F	S
Page No..	YOUVA				
Date:					

## Deciphering Naive Bayes

- \*  $P(\text{edible}|\alpha) = P(\alpha|\text{edible}) * P(\text{edible}) \rightarrow \text{Bayes Theorem}$
- \*  $P(\text{poisonous}|\alpha) = P(\alpha|\text{poisonous}) * P(\text{poisonous})$

$$P(\text{edible}|\alpha) > P(\text{poisonous}|\alpha)$$

$$\therefore \rightarrow \text{if } P(C_i|\alpha) > P(C_j|\alpha)$$

$\rightarrow \alpha$  is classified as  $C_i$

$\rightarrow$  This is called "**Maximum A posteriori Classification Rule**" (MAP)

$$\star P(C_i|\alpha) = \frac{P(\alpha|C_i) * P(C_i)}{P(\alpha)}$$

$P(C_i)$  = Prior probability  $[P(C_i=0), P(C_i=1)]$

$\rightarrow$  It is defined as the probability of occurrence of an event before collection of new data or feature is called as Prior

$\rightarrow$  Prior probability highly influences the class of new test point

$P(\alpha|C_i) [P(\alpha|C_i=1), P(\alpha|C_i=0)]$

Likelihood Function

$\rightarrow$  It maximises probability of observing data

Ex:-

$$P(\alpha=\text{convex}|\text{edible}) = 0.85$$

If edible is yes, then 85% chance mushroom is convex.

$\rightarrow$  It tells the likelihood of a point occurring in a class.

$\rightarrow$  The conditional independence assumption is leveraged while computing the likelihood probability.

Note:- 1. If prior is neutral (50%) then the likelihood may largely decide the outcome.

2. If prior is too powerful then likelihood often barely affects the result.

$\delta_1 = \pi_1, \delta_2 = \dots$

$$P(\delta_i | C_j) = P(\delta_i | C_i=0) \times P(\delta_i | C_i=1) \dots \times P(\delta_i | C_i=j)$$

↑  
Likelihood Calculation

M	T	W	T	F	S	S
Page No.:					YOUVA	

### Posterior

$$P(C_i | \delta_i) [P(C_i=0|\delta_i), P(C_i=1|\delta_i)]$$

Posterior Probability

→  $g_t$  is called calculated after incorporating the feature of the data point.

→  $g_t$  combines prior belief and case-specific information.

→  $g_t$  is balanced outcome of the prior and the likelihood.

### Naïve Bayes Assumption

- Variables are independent given class.
- Or variables are conditionally independent.
- Decreases Computational time of algorithm.

Q: How point 3 happens?

## Introduction - Naive Bayes for Text Classification

- It is widely used for
  - Document Classification
  - Natural language processing (NLP)
  - Sentiment Analysis
  - Spam, Ham classification

### Document Classifier - Pre Processing Steps

#### Test Dataset

Document → Feature Representation → Class

1. upgrad is a education institute		Education
1. Educational greatness depends on ethics		Education
2. A story of great ethics and educational greatness		Education
3. Sholey is a great cinema		Cinema
4. Good movie depends on good story		Cinema

Dictionary before Stop word removal	
0: and	
1: UpGrad	
2: is	
3: a	
4: Institute	
5: of	
6: Educational	
7: Cinema	
8: Movie	
9: Story	
10: good	

Dictionary After Stop word removal	
0: UpGrad	
1: Institute	
2: Educational	
3: Cinema	
4: Movie	
5: Story	
6: good	

Stop Words	
0: and	
2: is	
3: a	
5: of	

## Dictionary / vocabulary

Document	Cinema	defends	educational	ethics	good	uploaded	movie	/
	0	0	0	1	0	0	1	0
	1	0	1	1	0	0	1	0
	2	0	0	1	1	0	1	1
	3	1	0	0	0	0	1	0
	4	0	1	0	0	2	0	0

### ② Bag of Word Representation

The sentences are broken down into words and the ordering doesn't matter anymore as if it is put in a bag and shuffled.

Cinema, ... → uploaded, movie

$$\text{Deducation} = \left[ \begin{array}{cccccc} 0, 0, 1, 0, 0, 1, 0, \\ 0, 1, 1, 1, 0, 0, 1, 0, 0 \\ 0, 0, 1, 1, 0, 1, 1, 0, 0 \end{array} \right] \quad \text{Poster}$$

$$P(E) = 3/5$$

$$P(\text{Cinema}) = 2/5$$

$$\text{DCinema} = \left( \begin{array}{cccccc} 1, 0, 0, 0, 1, 0, 0 \\ 0, 1, 0, 0, 2, 0, 0 \end{array} \right) \quad 87$$

$$\frac{P(\text{Ed} | w_1, w_2, \dots, w_n)}{P(\text{Cin} | w_1, w_2, \dots, w_n)} \quad \text{Posterior}$$

$$P(\text{Cin} | w_1, w_2, \dots, w_n) = \frac{P(w_1, w_2, \dots, w_n | \text{Cin})}{P(w_1, w_2, \dots, w_n)}$$

$$= P(w_1 | c_i) P(w_2 | c_i) \dots P(w_n | c_i)$$

$$P(w_1, w_2, \dots, w_n) = P(w_1 | c_1) P(w_2 | c_2) \dots P(w_n | c_n)$$

$$P(\text{Ed} | \text{Text Doc}) = P(\text{Ed} | \text{Grand}) P(\text{Ed} | \text{Story}) P(\text{Ed} | \text{Edu})$$

$$P(\text{Grand}) P(\text{Story}) \rightarrow ignore$$

$$= (3/5) \times (1/1) \times 3/5$$

$$P(\text{Cin} | \text{Text Doc}) = 0 \times 1/7 \times 2/5$$

$$\text{Text Doc} = "Grand Story"$$

$$P(\text{Ed} | \text{TD}) = (3/5) \times (1/1) \times (3/5)$$

$$P(\text{Cin} | \text{TD}) = 1/7 \times (1/2) \times (2/5)$$

Document	Cinema	defends	educational	ethics	good	uploaded	movie
1	0	0	0	1	0	0	1
2	1	0	1	1	0	0	1
3	2	0	0	1	1	0	1
4	3	1	0	0	0	0	1
5	4	0	1	0	0	2	0

## Laplace Smoothing

→ Some time we see "zero probability problem" - probability of a word which has never appeared in a class. (it may appear in the dataset of another class) is 0.

→ This issue can be solved using "Laplace Smoothing"  
→ We need to ignore the word which doesn't belong to any class as "feature".

Ex:-

Text Document = "very good great Educational"

$$P(Edu | Text Doc) = \frac{0}{11} \times \frac{3}{11} \times \frac{3}{11}$$

$$P(Cinema | Text Doc) = \frac{3}{11} \times \frac{0}{11} \times \frac{0}{11}$$

→ Word "good" doesn't appear for education. we to add 1 to all the feature of both classes so it will appear 1 time.

	Reduced	P(w C)	cinema	P(c C) = Cm
Cinema	0+1	$\frac{1}{(11+1)} = \frac{1}{12}$	$1+1=2$	$\frac{2}{12+8} = \frac{2}{20} = \frac{1}{10}$
Story	1+1	$\frac{1}{12}$	$1+1=2$	$\frac{2}{20} = \frac{1}{10}$
Education	3+1	$\frac{3}{12}$	$0+1=1$	$\frac{1}{10}$
Ethnic	2+1	$\frac{3}{12}$	$0+1=1$	$\frac{1}{10}$
good	0+1	$\frac{1}{12}$	$3+1=4$	$\frac{4}{20} = \frac{1}{5}$
great	3+1	$\frac{4}{12}$	$0+1=1$	$\frac{1}{5}$
upgrade	0+1	$\frac{1}{12}$	$0+1=1$	$\frac{1}{5}$
Movie	1+1	$\frac{2}{12}$	$2+1=3$	$\frac{3}{20} = \frac{3}{20}$

Similarly we can do it again as we don't have "g"

$$P(e | Text Doc) = \frac{1}{12} \times \frac{4}{12} \times \frac{4}{12}$$

$$P(gm | Text Doc) = \frac{4}{20} \times \frac{1}{10} \times \frac{1}{10}$$

$$\frac{1}{12} \times \frac{4}{12} \times \frac{4}{12} \times \frac{1}{10} \times \frac{1}{10}$$

## Quick Introduction to Bernoulli Naive Bayes

- Prior to this we had gone through "multinomial Way of classifying documents".
- Now we will go through "Bernoulli Naive Bayes"

### Bernoulli Theorem: Introduction

One document =  $D = \begin{matrix} 0 & 0 & 1 & 0 & 0 & 1 & 0 & 1 & 0 & 0 & 0 & 1 \\ 0 & 1 & 1 & 1 & 0 & 0 & 1 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 1 & 1 & 0 & 1 & 1 & 0 & 0 & 0 & 1 & 0 \\ 1 & 0 & 0 & 0 & 0 & 1 & 0 & 0 & 0 & 1 & 0 & 0 \\ 0 & 1 & 0 & 0 & 2 & 0 & 0 & 0 & 1 & 0 & 1 & 0 \end{matrix}$  5th Document.

↑ This word appear 2 times in the 5th document

→ In "Bernoulli's Naive Bayes" we don't care about how many times a word occurs, but we care about whether a word occurs or not.

→ Now the 5th Document Will be replaced by following value and other field remains same:

5th Doc =  $\begin{matrix} 0 & 1 & 0 & 0 & 1 & 0 & 0 & 0 & 1 & 0 & 1 & 0 \end{matrix}$  ↗ Becomes 1 from 2.

Note:- Bernoulli Naive Bayes is concerned only with whether the word is present or not in a document, whereas the Multinomial Naive Bayes count the no. of occurrence of the words as well.

ML-C034-BG12-004

## Linear Regression

$x^{(i)}$  → input features

$y^{(i)}$  → o/p target variable

Train =  $\{ (x^{(i)}, y^{(i)}) \mid i=1 \dots N \}$

Machine Learning Model

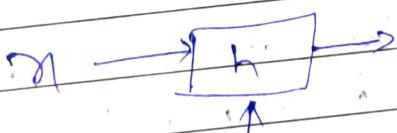
(distribution of data)  $X$  = Space of o/p sample  
 $y$  = Space of o/p value

Goal of Supervised learning → Given

$h: X \rightarrow Y$

hypothesis function

Such a way function  $h(x)$  is a good



learning algorithm

linear reg

$$h_0(x) = \theta_0 + \theta_1 x_1 + \theta_2 x_2$$

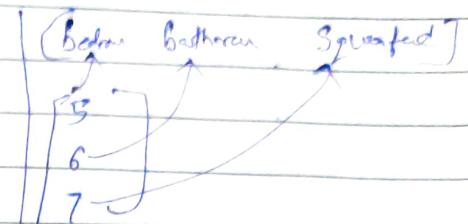
$\theta_i$  → parameter or weights

→ that control how you map function  $h(x)$  on to  $y$

$$h_{\theta}(x) = \sum_{i=0}^d \theta_i x_i$$

Least  
(~~approx~~) ( $d = \text{no. of I/P variable}$ )

$$h_{\theta}(x) = \sum_{i=0}^d \theta_i x_i$$



Choose  $\theta$  such that minimize  $J(\theta)$

$$J(\theta) = \frac{1}{2} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)})^2$$

$m = \text{no. of Sample}$

Least Mean Square (LMS)  $\rightarrow$  gt

$\rightarrow$  Choose  $\theta$  such that minimize  $J(\theta)$

$$\theta_j := \theta_j - \alpha \frac{\partial J(\theta)}{\partial \theta_j}$$

$$\frac{\partial}{\partial \theta} = \frac{\partial}{\partial \theta_0}, \dots, \frac{\partial}{\partial \theta_d}$$

For Single term

$$\begin{aligned} \frac{\partial J(\theta)}{\partial \theta} &= \frac{1}{2} \left[ (h_{\theta}(x^{(i)}) - y^{(i)}) x^{(i)} + (h_{\theta}(x^{(i)}) - y^{(i)}) x^{(i)} \right] \\ &= (h_{\theta}(x^{(i)}) - y^{(i)}) \alpha_j \end{aligned}$$

$$\theta_j := \theta_j + \alpha (y^{(i)} - h_{\theta}(x^{(i)})) x_j^{(i)}$$

Two Way

1. Batch gradient descent
2. Stochastic "

Batch Gradient descent  
Repeat until Converges

$$\theta_j := \theta_j + \alpha \sum_{i=1}^m (y^{(i)} - h_{\theta}(x^{(i)})) x_j^{(i)}$$

for every j

$\rightarrow$  Assuming learning rate is not too large, the above eqn converges to global minima.

$\rightarrow$  For every data point based on

See all Sample

See 1 Sample \$  
Update

Batch GGD  
(Slow but reliable)

Stochastic GGD  
(fast)

SGD get 0 clean to minima but never reaches  
reachable to global minima.

→ Iterative approach (BGD & SGD)

⇒ Non-iterative approach

find the optimal set of 'weight' in a Single - Setup

$$y = mn + c$$

$$x = \begin{bmatrix} 1 & m^T \\ 1 & n^{(1)T} \\ 1 & n^{(2)T} \\ 1 & n^{(m)T} \end{bmatrix} \quad y = \begin{bmatrix} y^{(1)} \\ y^{(2)} \\ \vdots \\ y^{(m)} \end{bmatrix}$$

$$x_0 - y = \begin{bmatrix} h(m^n - y^{(1)}) \\ h(m^n - y^{(2)}) \\ \vdots \\ h(m^n - y^{(m)}) \end{bmatrix}$$

$$J(\theta) = \frac{1}{2} \sum_{i=1}^m (h(m^n - y^{(i)}))^2$$

$$Z^T Z = \sum Z_i^2$$

$$J(\theta) = \frac{1}{2} (x_0 - y)^T (x_0 - y)$$

## Matrix derivation

$$f: \alpha \begin{bmatrix} A_1 & \dots & A_n \end{bmatrix}$$

$$\nabla_\theta f(\theta) = \begin{bmatrix} \partial A_1, \dots, \partial A_n \end{bmatrix}$$

grad in  $n \times n$  matrix

$$J(\theta) = \frac{1}{2} (x_0 - y)^T (x_0 - y)$$

$$\nabla_\theta J(\theta) = \frac{1}{2} \nabla_\theta [(x_0)^T x_0 - (x_0)^T y - y^T x_0 - y^T y]$$

$$= \frac{1}{2} \nabla_\theta [2(x^T x)\theta - 2(x^T y)^T \theta]$$

$$= \frac{1}{2} \nabla_\theta [2(x^T x)\theta - 2(x^T y)^T \theta]$$

$$= (x^T x)\theta - x^T y$$

$$\begin{pmatrix} a & b \end{pmatrix}^T = c^T B^T$$

Equate the first derivative to zero

$$x^T x \theta = x^T y$$

linear regression has closed form loop

$$\theta = (x^T x)^{-1} x^T y$$

$x^T$  is not possible to find anywhere

Matrix need to be Singular to find  
X inverse.

M	T	W	T	F	S	S
Page No.:	YOUVA					
Date:						

$$X_{(n \times d)} \rightarrow X^T_{(d \times n)}$$

## Logistic Regression

$h_0(x) \rightarrow \theta^T x \rightarrow$  Sigmoid funt

$$h_0(x) = g(\theta^T x) = \frac{1}{1 + e^{-\theta^T x}}$$

$$g(z) = \frac{1}{1 + e^{-z}}$$

$$\begin{aligned} \text{As } z &\rightarrow \infty & g(z) &\rightarrow 1 \\ z &\rightarrow -\infty & g(z) &\rightarrow 0 \end{aligned}$$

$$\theta^T x = \theta_0 + \sum_{j=1}^d \theta_j x_j$$

$$= \sum$$

$$\begin{aligned} \frac{dg(z)}{dz} &= \frac{-1}{(1 + e^{-z})} \cdot \frac{d}{dz}(1 + e^{-z}) \\ &= \frac{e^{-z}}{(1 + e^{-z})^2} \end{aligned}$$

$$\begin{aligned} \frac{df(x)}{dx} &= \frac{1}{n} \\ \frac{df(x)}{\partial x^2} &= \frac{1}{n^2} \cdot \frac{d^2}{dx^2} \end{aligned}$$

derivative of Sigmoid

$$[g'(z) = g(z)(1-g(z))]$$

$$h_{\theta}(x)$$

$$y=1, y=0$$

likelihood ←

$$\begin{cases} P(y=1|x; \theta) = h_{\theta}(x) \\ P(y=0|x; \theta) = 1 - h_{\theta}(x) \end{cases}$$

$$P(y|x; \theta) = (h_{\theta}(x))^y (1-h_{\theta}(x))^{1-y}$$

Assume in  $y^{(i)}$  example are generated independent  
we can write down the likelihood function

$$L(\theta) = P(y|x; \theta)$$

$$= \prod_{i=1}^n P(y^{(i)}|x^{(i)}; \theta)$$

$$= \prod_{i=1}^n (h_{\theta}(x^{(i)}))^{y^{(i)}} (1-h_{\theta}(x^{(i)}))^{1-y^{(i)}}$$

$\log L(\theta)$  is monotonically increasing

$$\log L(\theta) = \sum_{i=1}^n y^{(i)}$$

Maximizes likelihood using gradient descent

$$\theta := \theta + \alpha \nabla \log L(\theta)$$

# Tensorflow Pytorch $\rightarrow$ Autodiff

M	T	W	T	F	S	S
Page No.						
Date					YOUVA	

$$\frac{\partial J(\theta)}{\partial \theta} = \left( \frac{y}{g(\theta^m)} - \frac{(1-y)}{1-g(\theta^m)} \right) g(\theta^m)(1-g(\theta^m))$$

$$\left( \frac{y}{g(\theta^m)} - \frac{(1-y)}{1-g(\theta^m)} \right) g(\theta^m)(1-g(\theta^m))$$

find Cost function then find derivative