

Dataset

Page
Date

Train

Test

1 2 3 4 5 6 7

$k=5$

k fold

1 2 3	4 5 6	7	1 2

odds → The odds that India wins. The match against Pakistan is 4:1 with

India → 4 matches
Pakistan → 1 match

Total number of matches = 5

$$P(\text{India} \text{ wins}) = 4/5 \rightarrow P(\text{India} \text{ loses}) = 1/5$$

$$P(\text{PAH} \text{ wins}) = 1/5$$

$$P(\text{India} \text{ wins}) = 4/5$$

$$P(\text{India} \text{ loses}) = 1/5$$

$$\text{odds} = \left(\frac{P}{1-P} \right) \rightarrow 4:1$$

$$\frac{P(\text{India} \text{ win})}{P(\text{India} \text{ lose})} = \frac{(4/5)}{(1/5)} = (4)$$

Odds → $4 : 1 = (4/1)$

Shot on AWESOME A05s

$$\text{odds} = \left(\frac{P}{1-P} \right) \rightarrow \boxed{\text{Date: } \quad \quad}$$

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

$$\log(\text{odds}) = \log\left(\frac{P}{1-P}\right) \rightarrow \text{eq } ②$$

✓ Log of odds = Logit function

$$\text{log(odds)} = \log(P/(1-P))$$

P → Range → 0 ↔ 1

$$\log\left(\frac{P}{1-P}\right) = \log(P) - \log(1-P) \rightarrow \begin{cases} \log(\%) \\ \downarrow \\ \log(a) - \log(b) \end{cases}$$

$$= \log(0) - \log(1-0) = \log(0) - \log(1)$$

$$= (-\infty) - 0 = -\infty$$

when

$$P=0$$

$$\log(0) = -\infty$$

$$\log(1) = 0$$

$$\textcircled{1} \quad \text{log odds} = \log(P/(1-P))$$

$P \rightarrow \text{Range} \rightarrow 0 \rightarrow 1$

$$\log(P/(1-P)) = \log(P) - \log(1-P)$$

$$= \log(1) - \log(1-1) = \log(1) - \log(0)$$

$$= 0 - (-\infty) = (+\infty)$$

$\log(a/b)$
\downarrow
$\log(a) - \log(b)$

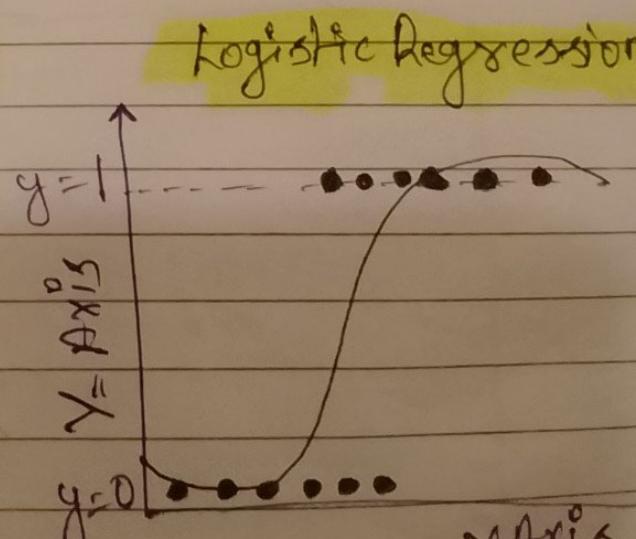
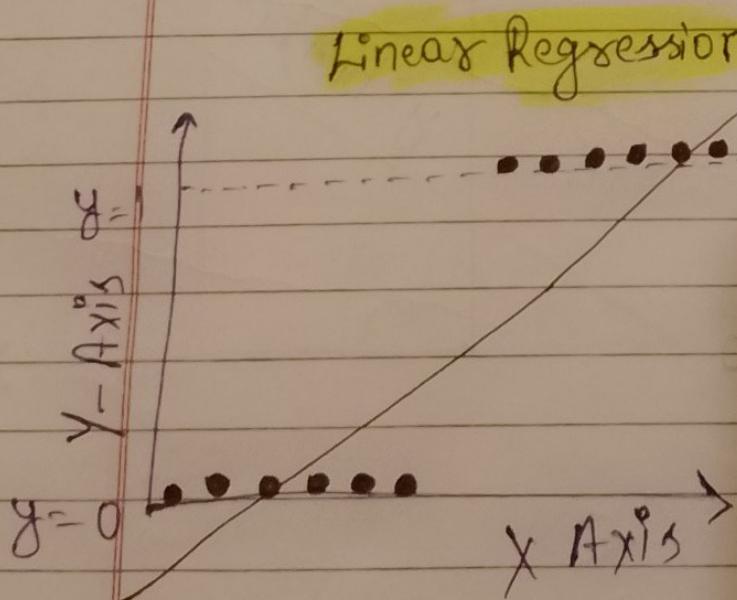
when
 $P=1$
 $\log(0)$
 $= -\infty$
 $\log(1)$
 $= 0$

$$\textcircled{2} \quad P=0 \rightarrow \log(P/(1-P)) \rightarrow -\infty$$

$$P=1 \rightarrow \log(P/(1-P)) \rightarrow +\infty$$

$y = b_0 + b_1 x_1 \rightarrow$ Linear Regression

(13) Logistic Regression



Logistic Model

$$P = \frac{1}{1 + e^{-(b_0 + b_1 x)}}$$

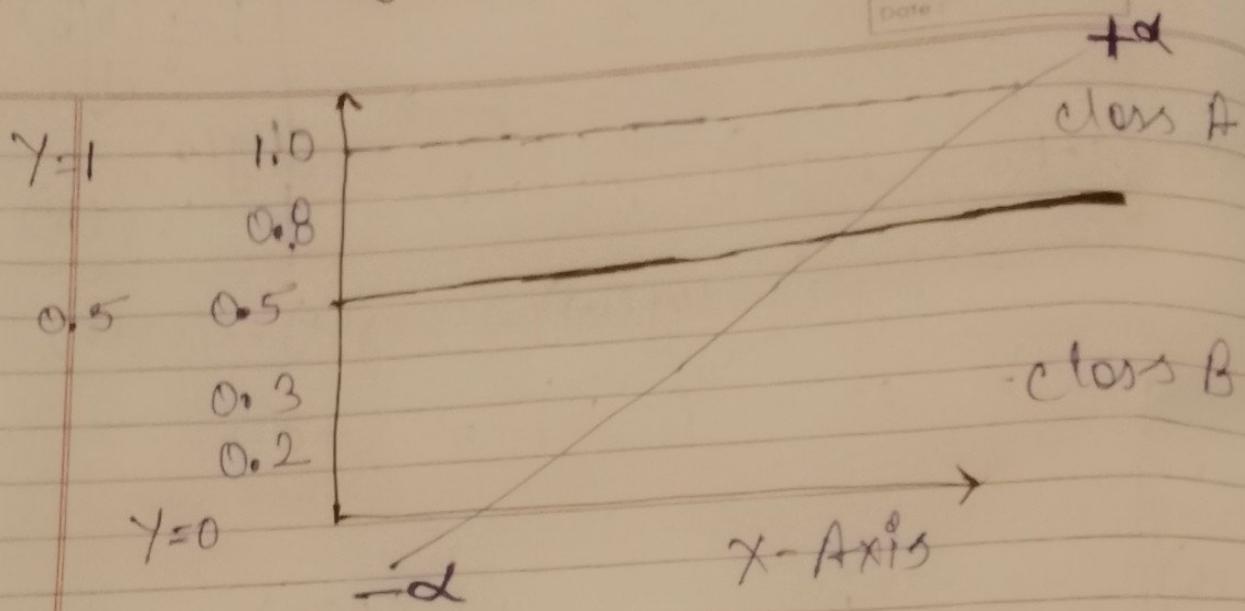
(14)

logistic Regression

Page	
Date	

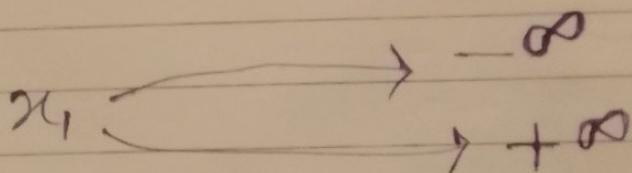
(16)

[Fin]

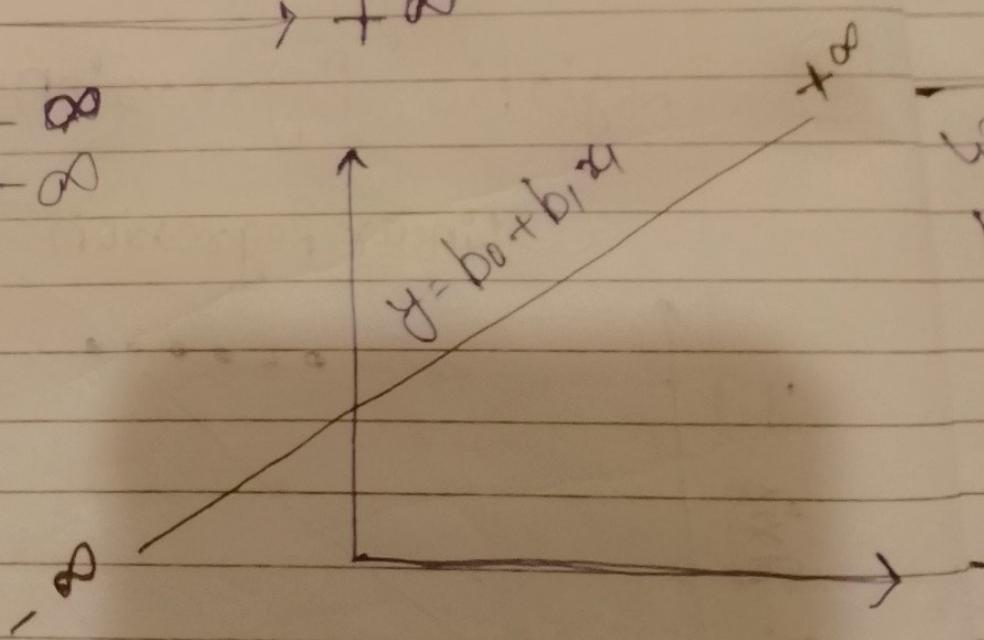


(15)

$$y = b_0 + b_1 x_1$$



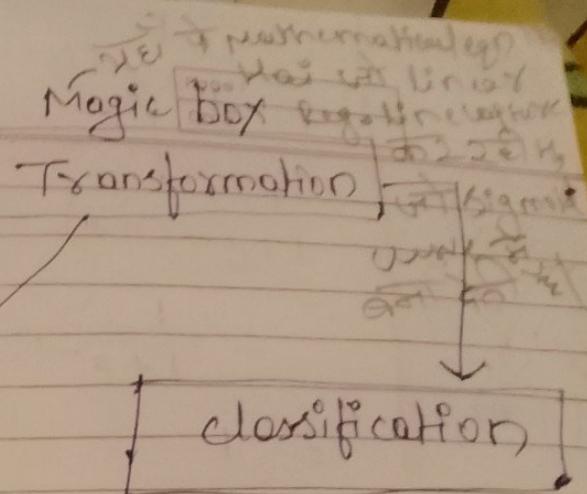
$$y \rightarrow -\infty \quad +\infty$$



(16)

Linear Reg

$$y = b_0 + b_1 x_1$$



Mathematical equation

$$\left(\frac{1}{1+e^{-z}} \right) \rightarrow \text{going to capture linear regression}$$

This is exponential
value = $e = 2.713$.

(17)

$$\frac{1}{1+e^{-(b_0+b_1 x_1)}}$$

when

$$b_0 + b_1 x_1 = y = \frac{1}{1+e^{-(\infty)}}$$

 $(-\infty)$

$$= \frac{1}{1} = \text{large number}$$

$$= 0$$

$$z = b_0 + b_1 x_1$$

$$x_1 \rightarrow -\infty$$

$$y = b_0 + b_1 x_1 \rightarrow +\infty$$

$e^\infty = \text{large number}$

$(2.713)^\infty = \text{large number}$

$1 + \text{large no.} = \text{large number}$

$$(18) \quad \frac{1}{1+e^{-(b_0+b_1x_1)}} = \frac{1}{1+e^{-c+\alpha}}$$

when both b_0 & b_1 \downarrow
 $(+\infty)$

$$= \frac{1}{1+e^{-\infty}} = \frac{1}{1+1} = \frac{1}{2}$$

$$= \frac{1}{1+\frac{1}{e^{\infty}}} = \frac{1}{1+\frac{1}{\text{large no.}}} = \frac{1}{1+0} = 1 \quad \checkmark$$

x_1 \rightarrow food \downarrow
 Donet \uparrow

$$y = b_0 + b_1 x_1$$

$\downarrow \rightarrow -\infty$
 $\uparrow \rightarrow +\infty$

$e^\alpha = \text{large number}$

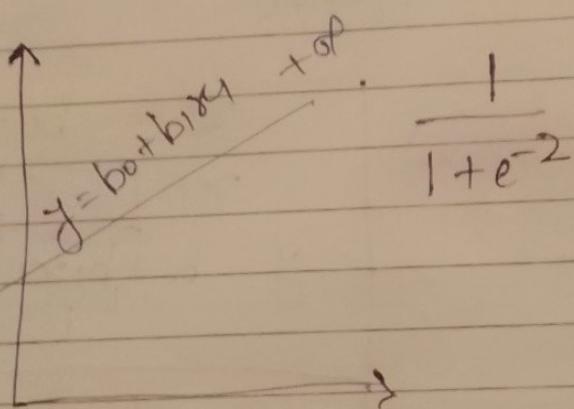
$(2.713)^\alpha = \text{large number}$

$1 + \text{large} = \text{large number}$

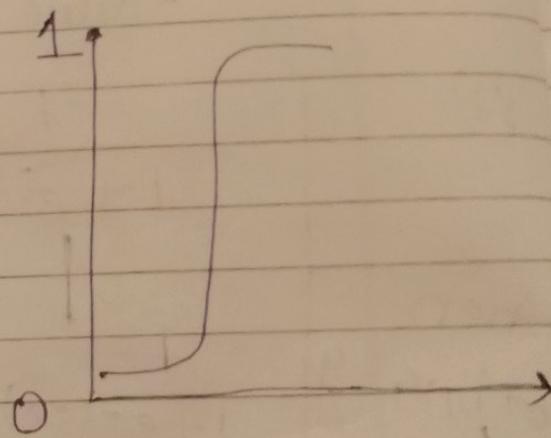
(21)

he

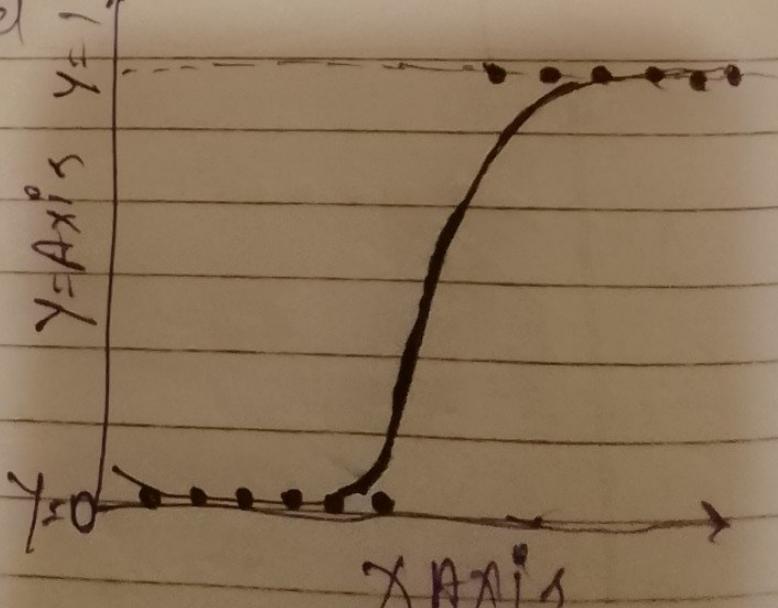
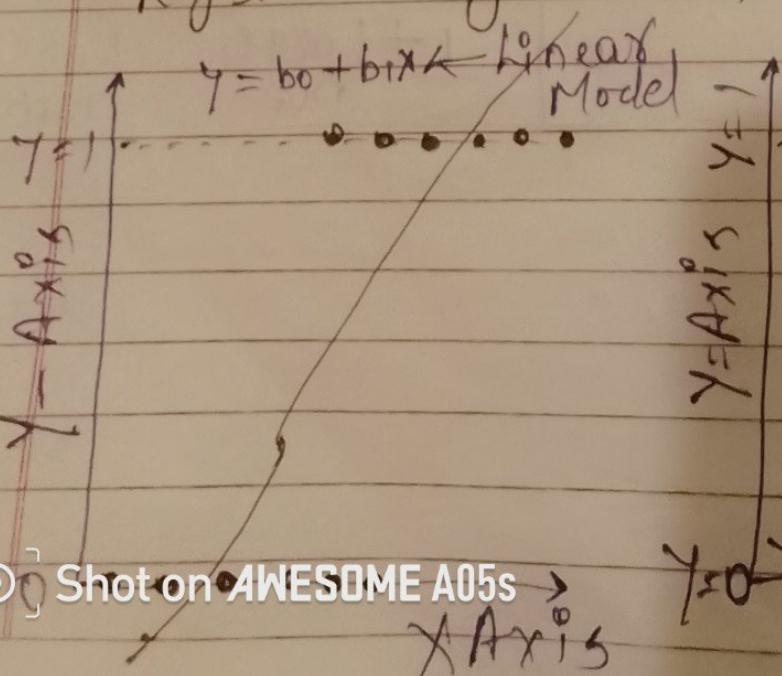
(22)



$$\frac{1}{1+e^{-2}}$$



(20) Logistic Regression



logistic Model

$$P = \frac{1}{1 + e^{-(b_0 + b_1 x)}}$$

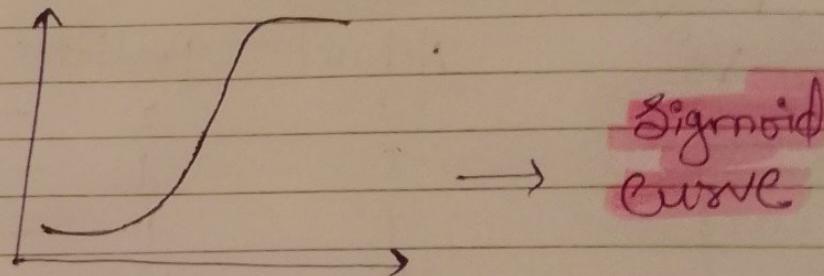
Page

Date

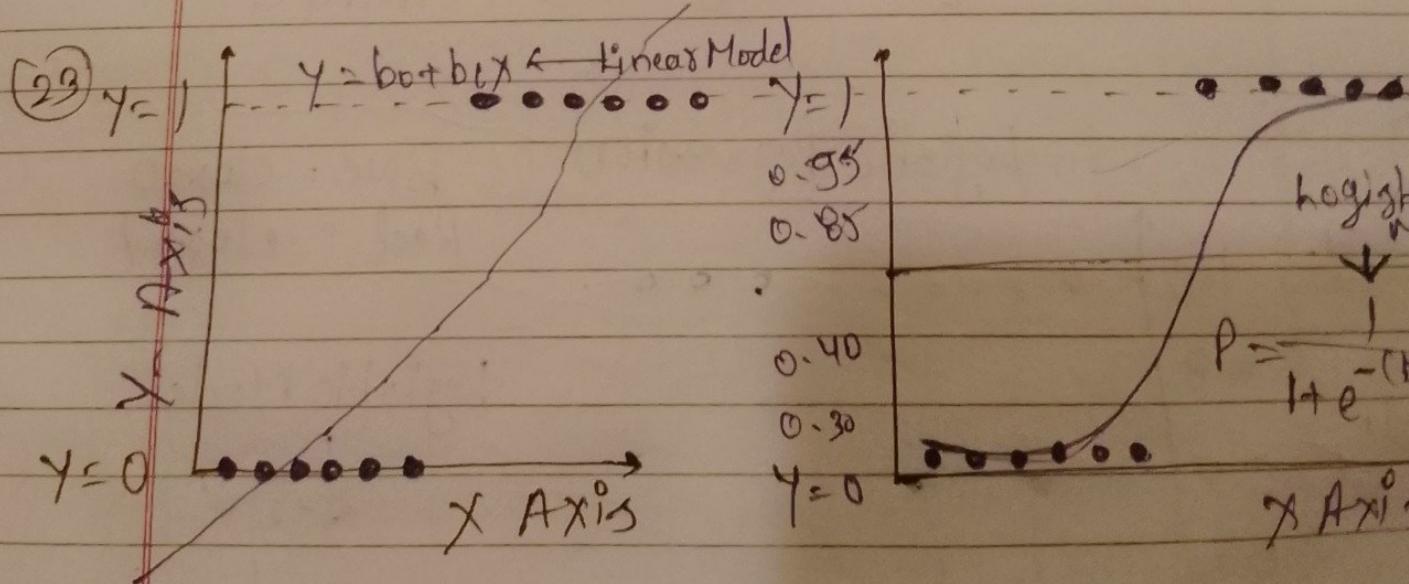
(21) $P=0 \rightarrow \log(P/(1-P)) \rightarrow -\infty \rightarrow b_0 + b_1 x_1(x_1)$
 $P=1 \rightarrow \log(P/(1-P)) \rightarrow +\infty \rightarrow b_0 + b_1 x_1(x_1)$

$\log(P/(1-P)) = b_0 + b_1 x_1$ ↗ logit function
logistic Regression

(22) $\left(\frac{1}{1+e^{-x}} \right) \rightarrow$ Sigmoid equation



Sigmoid
curve



(24)

$$P = \frac{1}{1+e^{-2}} = \frac{1}{1+e^{-(b_0+b_1x_1)}}$$

$$x_1 = 1.1015 \\ P = 0.85$$

$$x_1 = 1.000 \\ P = 1$$

(27)

(28)

(25)

$$\text{Regression} \rightarrow \text{Residual} = A - P$$

classification

→ correctly classified

→ incorrectly classified

 $\begin{cases} 0 \rightarrow N \text{ Diabetic} \\ 1 \rightarrow \text{Diabetic} \end{cases}$

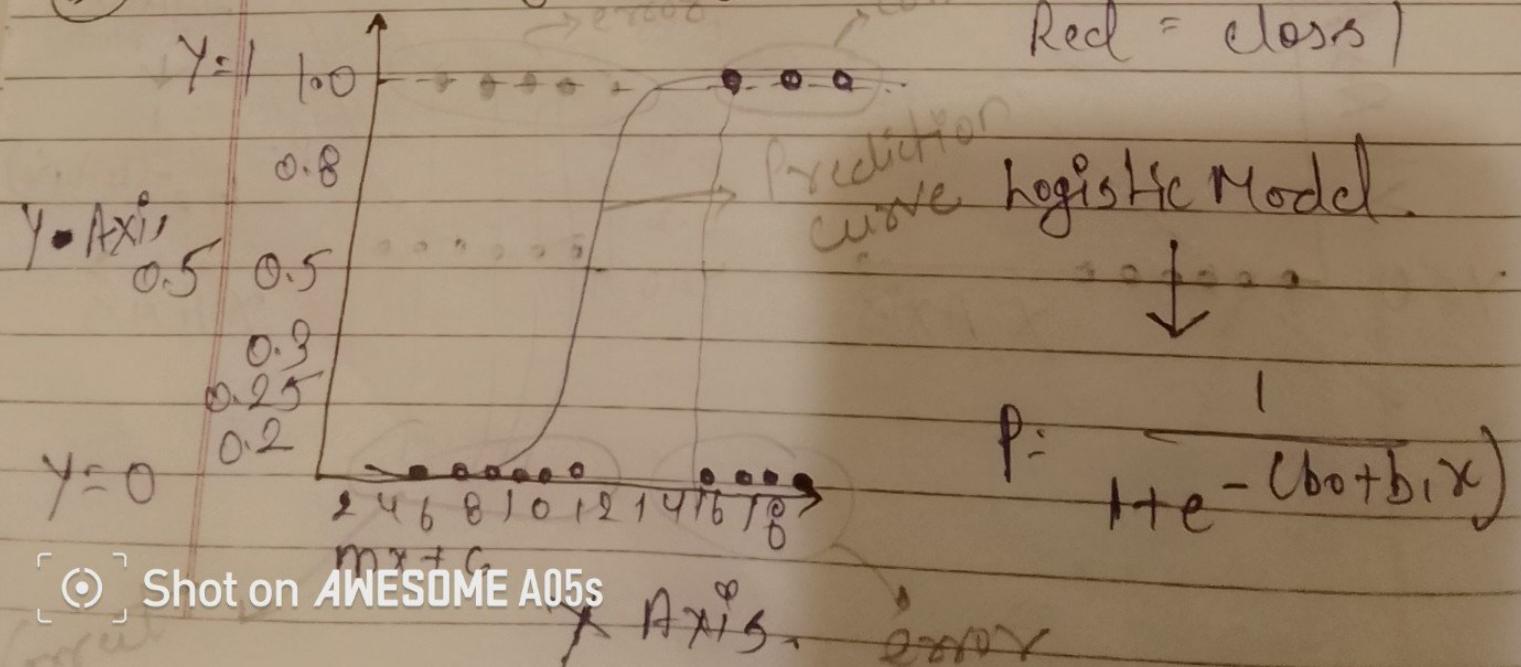
Actual	Predict	
1	0	→ Misclassified
0	1	→ Misclassified
1	1	→ correctly classified
0	0	→ correctly classified

(26)

Logistic Regression

Blue: class 0

Red: class 1



Shot on AWESOME A05s

X	0	$1-P$	\downarrow	P
$P(X)$				

when
 $y=0$

$$= P^y \times (1-P)^{1-y}$$

$$= P^0 \times (1-P)^{1-0}$$

$$= 1 \times (1-P)$$

P = Probability of success

$1-P$ = Probability of failure

$y \rightarrow$ class

$y=0$ class

$y=1$ class

(3)

Bernoulli distribution

y	0	1	\downarrow	P
$P(y)$	$1-P$			

when
 $y=1$

$$= P^y \times (1-P)^{1-y}$$

$$= P^1 \times (1-P)^{1-1}$$

$$= P^1 \times (1-P)^0$$

$$= P \times 1 = P$$

P = Probability of success

$1-P$ = Prob. of failure

$y \rightarrow$ class

$y=0$ class

$y=1$ class

(3)

Bernoulli distribution

y	0	1	\downarrow	P
$P(y)$	$1-P$			

when
 $y=0$

$$= P^y \times (1-P)^{1-y}$$

$$= P^0 \times (1-P)^{1-0}$$

$$= 1 \times (1-P) = 1-P$$

P = Probability of success

$1-P$ = Prob. of failure

$y \rightarrow$ class

$y=0$ class

$y=1$ class

Shot on AWESOME A05s

$$(36) \text{ when } y=0, P^y \cdot (1-P)^{1-y} \rightarrow P$$

when $y=1, P^y \cdot (1-P)^{1-y} \rightarrow 0$

x	y	0	1
$P(y)$		$1-P$	P

$$f(x) = P^y \cdot (1-P)^{1-y}$$

(37)

$$L = (P^y)(1-P)^{1-y}$$

↳ Probability of success
Probability of failure

Likelihood is defined as the Probability of success \times Probability of failure

$$L = (P^y)(1-P)^{1-y}$$

(38)

$$L = (P^y)(1-P)^{1-y}$$

$$\log(L) = \log[(P^y)(1-P)^{1-y}]$$

$$\log(L) = \log(P^y) + \log(1-P)^{1-y}$$

$$\log(L) = y \log(P) + (1-y) \log(1-P)$$

Log Likelihood

$\log(a \cdot b)$

$\log a + \log b$

$\log(a^b)$

$b \log a$

(39)

Null Model	Fitted Model	Full Model	Saturated Model
↓ When you build Model without any Predictors	↓ If you build Model with at least one Predictor	↓ If you build Model with all predictors	↓ If you build Model which Predictors everything right (Perfect Model)
Intercept Model			

(40)

 $\log(L) \rightarrow$ Saturated Model

$\log L = 0$

saturated

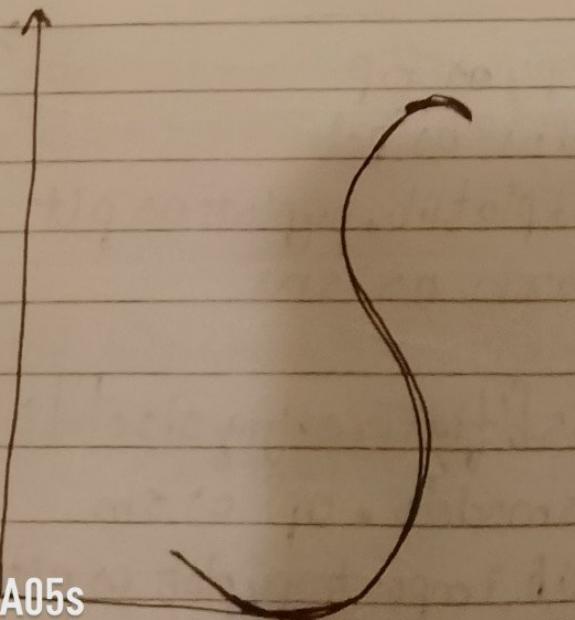
(-ve)

null model
 $(\log L) - 300$

Good model
should have
log likelihood
closer to 0

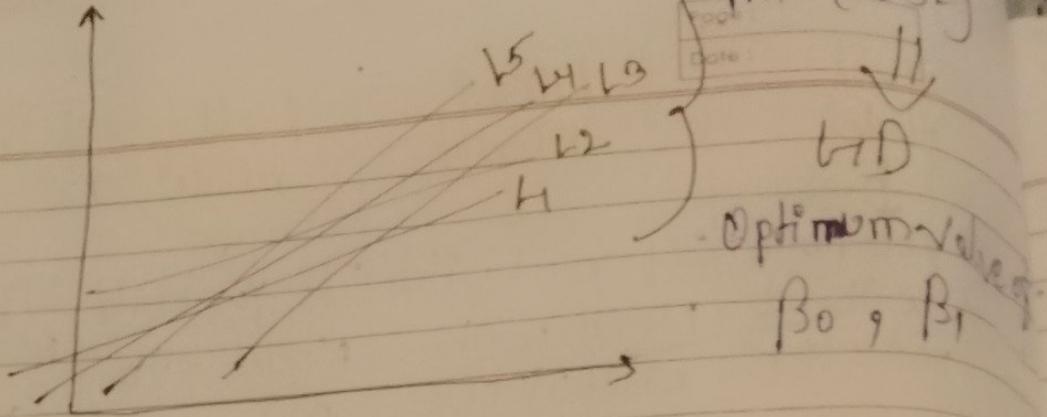
far from
null model

(41)

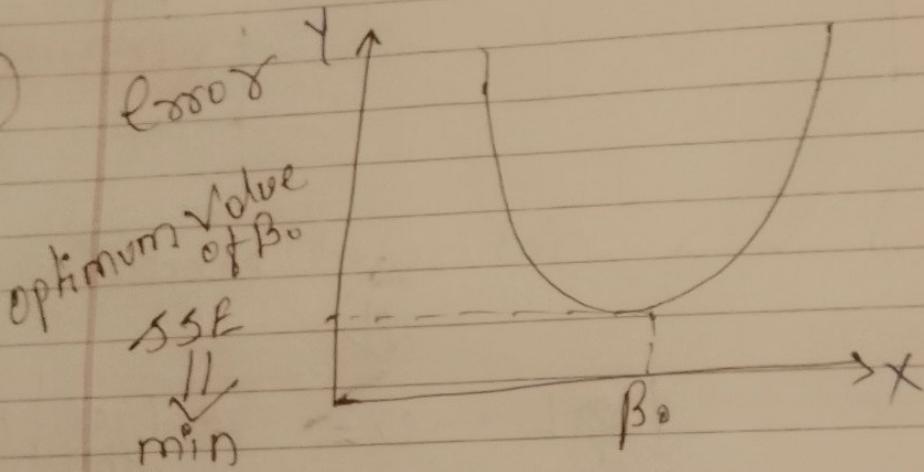


Shot on AWESOME A05s

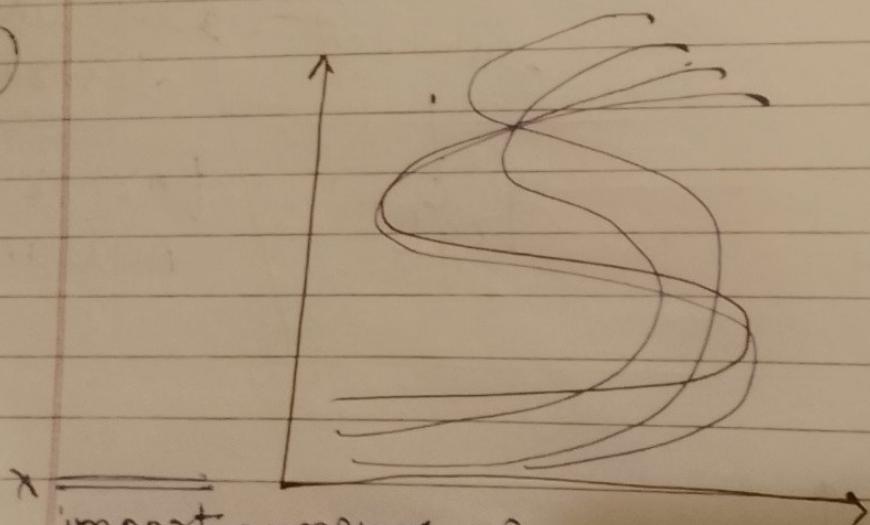
(42)



(43)



(45)



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

```
plt.rcParams['figure.figsize'] = [15,8]
```

```
import statsmodels.api as sm  
from matplotlib import pyplot as plt
```



Shot on AWESOME A05s

point (logreg.summ)
optimization terminated
current function value: 0.24436
Iterations: 8

① End

Optimum
Value of
B0 w/
55

②

logistic
KELTHOG

Linear regression based on logit
Logistic regression based on probability

Stage ①

Metrics

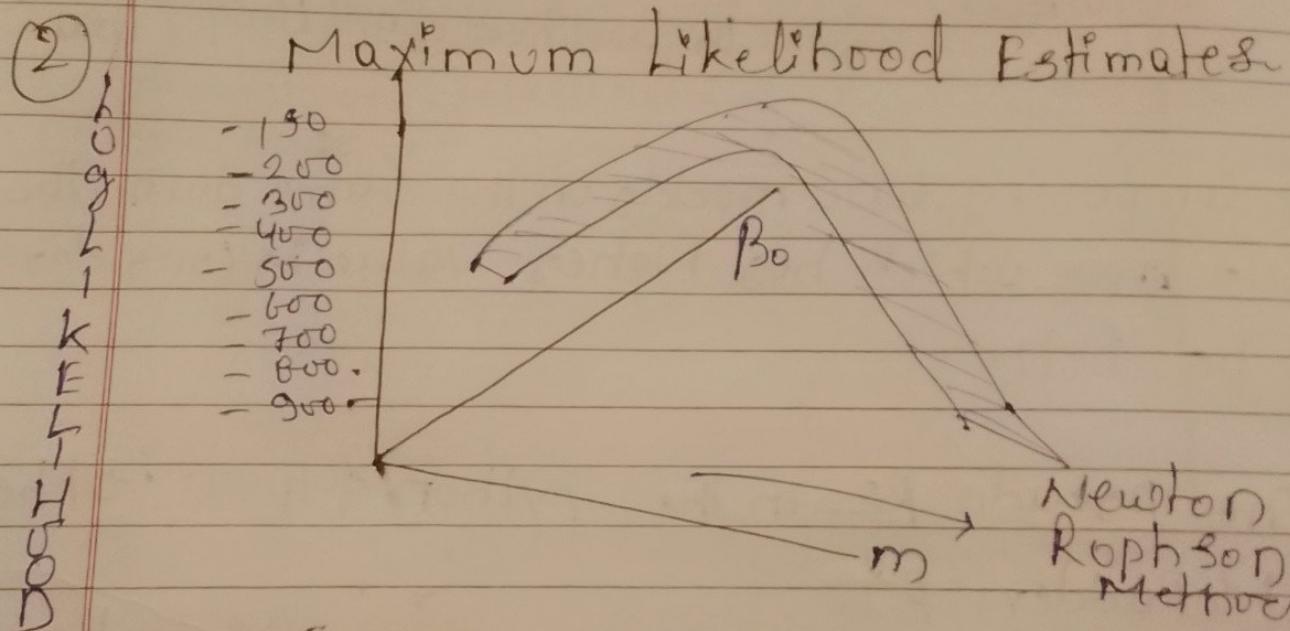
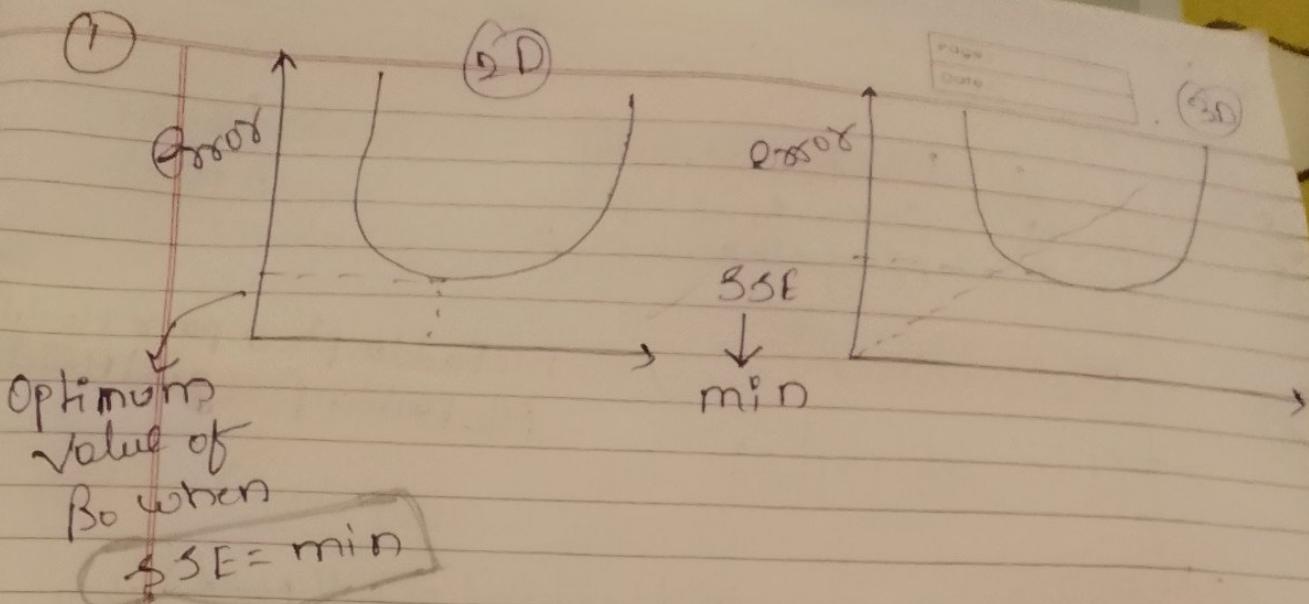
Model Evaluation
Metrics

- Deviance
- AIC
- Pseudo R²

Stage ②

Model Performance
Metrics

- Confusion Matrix
- Cross Entropy
- ROC



Stage ① Metrics	Stage ② Metrics
Model Evaluation Metrics	Model Performance Metrics
Deviance ✓	Confusion Matrix ✓
AIC ✓	Cross Entropy ✓
Pseudo R ² ✓	R ² ✓

Stage ① Metrics: Deviance, AIC, Pseudo R².

Stage ② Metrics: Confusion Matrix, Cross Entropy, R².

Annotations:

- Stage ① Metrics:** Handwritten note: "This tells you about how good the model was built".
- Stage ② Metrics:** Handwritten note: "How good were predictions done".

(4)

Pseudo R^2

The Pseudo R^2 are on the Probabilistic

This called as Pseudo R^2 because it based on the Probabilistic

range 0 to 1

(5)

McFadden R^2

Linear regression

(6)

Cox-Snell R^2

R^2 based on the variation

(7)

Nagelkerke R^2

McFadden R^2

McFadden $R^2 = 0$ poor Model

McFadden $R^2 = 1$ good Model

(8)

This is defined as

$$R^2_{\text{McFadden}} = 1 - \frac{\text{In likelihood of full Model}}{\text{In likelihood of null Model}}$$

$\frac{\text{Num}}{\text{Den}}$

→ closer to 1 = saturated

If comparing two models on the same data, the model which has higher value is considered to be better

The pseudo R^2 in the python output is the McFadden R^2

This is done so that the values are bounded.

Pseudo

$R^2 = 1$

Overfit
(model)

(poor model)

⑥

null
model

Deviance

$\log L$

full model \rightarrow all features

Row
Sum

Saturated
Model
 $\log(1) = 0$

closer to zero

⑦

⑧

$$D = -2 \ln \left[\frac{\text{Likelihood of fitted Model}}{\text{Likelihood of Saturated Model}} \right]$$

$$D = -2 \log \left[\frac{\text{Likelihood of fitted model}}{\text{Likelihood of Saturated model}} \right]$$

\downarrow

$$\log(1) = 0$$

fitted
model

\hookrightarrow likelihood of fitted model
likelihood of saturated model

saturated
model

① Deviance is a measure

② It measures the difference between the fitted model and saturated model.

③ Deviance should tends towards zero

④ If Deviance is low it means that the fitted model is closer to the Saturated model

⑤ If the deviance is high it means the fitted model is far off from the Saturated model

(10)

AIC (Akaike Information Criteria)

(11)

- ① Lower the AIC better is the model
- ② AIC is measuring the loss of information

(11) The Akaike Information Criteria (AIC) is a relative measure of model

- ✓ evaluation for a given dataset

✓ It is given by

L : log likelihood

k : parameters to be estimated

False
False
True
True

The AIC gives a trade off between the model accuracy and model complexity
i.e it prevents us from overfitting

$$AIC = -2 \log L + 2k$$

k is the no. of Parameters.

$\log L \rightarrow$ Log Likelihood of the model

$$\textcircled{12} \quad \text{AIC} = -2 \log L + 2k$$

lower the AIC better is the model

→ AIC is positive value

$$\rightarrow \left. \begin{array}{l} \text{AIC}_{M_1} = 1000 \\ \text{AIC}_{M_2} = 2000 \end{array} \right\} M_1 \text{ is better than } M_2$$

→ AIC is Relative Measure

$$\textcircled{13} \quad \left| \begin{array}{l} \text{Errors} \\ (\text{Reg}) \\ (\text{class}) \end{array} \right| \quad \begin{array}{l} A - P \\ \rightarrow \text{Missclassified} \end{array}$$

	Actual	Predicted	
False Positive (FP)	0	1	Misclassified
False Negative (FN)	1	0	
True Negative (TN)	0	0	Correctly classified
True Positive (TP)	1	1	

$$\textcircled{14} \quad \left| \begin{array}{l} \text{Errors} \\ (\text{Reg}) \\ (\text{class}) \end{array} \right| \quad \begin{array}{l} A - P \\ \rightarrow \text{Missclassified} \end{array}$$

	Actual	Predicted	Model	0 → Negative 1 → Positive
FP	0	1		Mismatch → False
FN	1	0		Correctly → True
TN	0	0		
TP	1	1		

Precision - out of TP + FP, how many are true

Recall -

, how many are actually

(24)

Predict

Page _____
Date _____

Actual

	O(N)	O(P)
O	TN	FP
I	FN	TP

Recall (positive) = $\frac{TP}{TP+FN}$ Recall (negative) = $\frac{TN}{TN+FP}$

(25)

Predictive.

Actual

	O(N)	O(P)
O	TN	FP
I	FN	TP

$$\text{Precision}_{\text{(positive)}} = \frac{TP}{TP+FP} = \frac{1}{T+I} = \frac{1}{2} = 0.5$$

$$\text{Precision}_{\text{(negative)}} = \frac{TN}{TN+FN}$$

Precision = Precision is the ratio of correctly predicted positive classes to the total positive. It is defined as $(TP / (TP + FP))$

Recall = is also a ~~Sensitivity~~ & True Positive rate.
It is the ratio of correctly predicted positive observations to all the observations $(TP / (TP + FN))$

Note \Rightarrow Precision and Recall are inversely associated with each other, it means that if the precision increases, the recall will go down and vice-versa.

Shot on AWESOME A05s

(26) $F_{1\text{ score}}(\text{Pos}) = 2 \times \frac{\text{Recall}(\text{Pos}) \times \text{Precision}(\text{Pos})}{\text{Recall}(\text{Pos}) + \text{Precision}(\text{Pos})}$

$F_2\text{ score}(\text{Neg}) = 2 \times \frac{\text{Recall}(\text{Neg}) \times \text{Precision}(\text{Neg})}{\text{Recall}(\text{Neg}) + \text{Precision}(\text{Neg})}$

f_1 score is harmonic mean of Precision and Recall.

(27)

		Prediction	
		O(N)	I(P)
Actual	0	TN	FP
	1	FN	TP

Precision(Pos) = $\frac{TP}{TP+FP}$

Precision(Neg) = $\frac{TN}{TN+FN}$

$F_{1\text{ score}}(P) = 2 \times \frac{R \times P}{R+P}$

$F_{1\text{ score}}(N) = 2 \times \frac{R \times P}{R+P}$

Recall(Pos) = $\frac{TP}{TP+FN}$

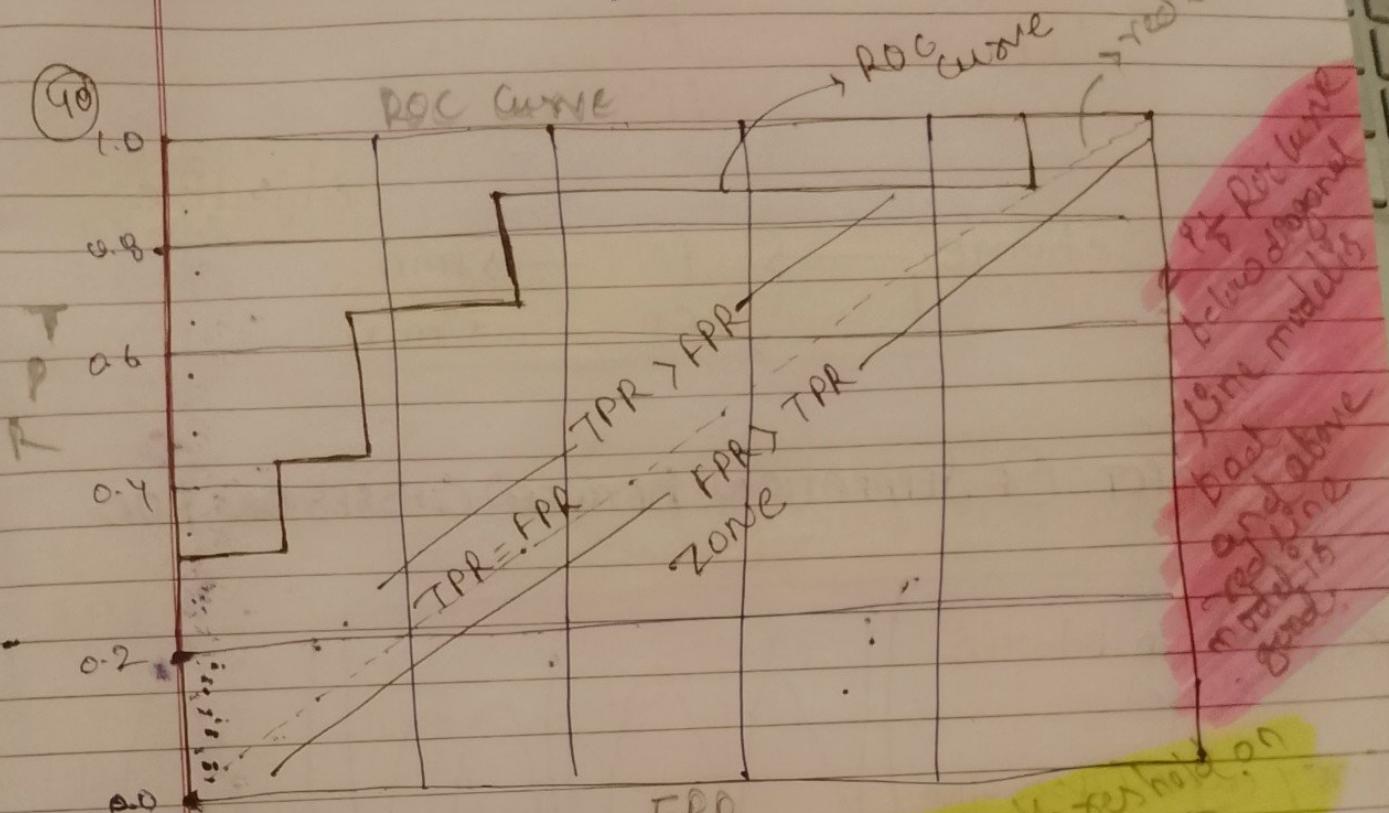
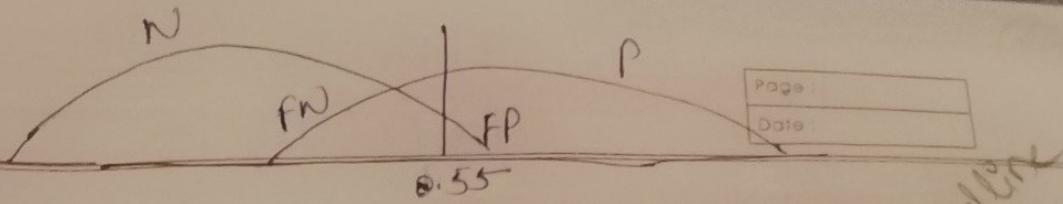
Recall(Neg) = $\frac{TN}{TN+FP}$

(28) Interpretation COEFFICIENTS OF LOGISTIC REGRESSION

$$\text{Log(Odds)} = b_0 + b_1 x_1$$

$y = b_0 + b_1 x_1$ 1 unit of change in x_1 , brings b_1 units of change in y

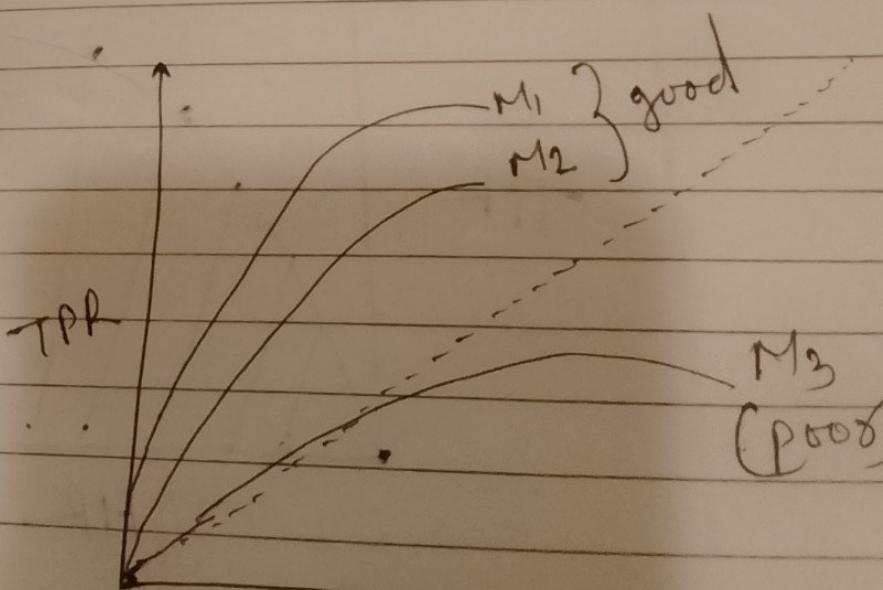
1 unit of change in $x_1 \rightarrow$ Brings b_1 units of change in Log(Odds)



TPR	FPR	Threshold
0.33	0	0.78
0.33	0.09	0.73
0.44	0.091	0.62
0.44	0.18	0.61
0.66	0.18	0.55
0.67	0.36	0.46
0.88	0.36	0.35
0.89	0.91	0.159
1	0.91	0.157

Q9

Do you plot Y-axis threshold? No
ROC curve basis NO
we plot ROC curve for TPR & FPR



(42)

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

Page _____
Date _____

(45)

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

change \rightarrow FN \rightarrow min
 \rightarrow FP \rightarrow min

objective

(43)

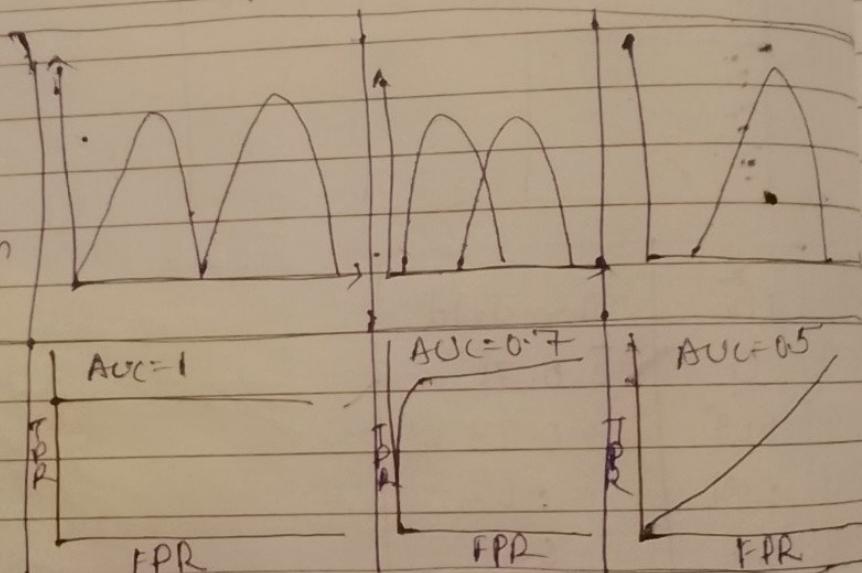
EFFECT OF SEPARATION BETWEEN CLASSES ON ROC

separation between classes

Grey curve Observation

ROC

curve

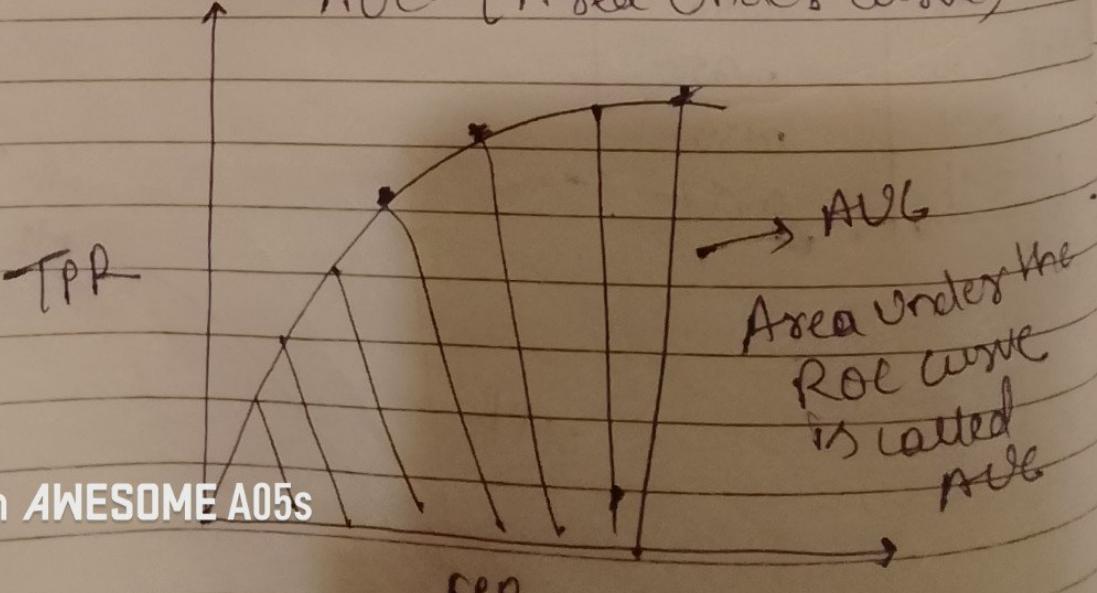


(46)

R₁
R₂
R₃

R₄
R₅
R₁₀₀

AUC = [Area Under Curve]



Shot on AWESOME A05s

(47)

(45)

Kappa	Interpretation	Page Date
< 0	No consistency	
0 - 0.2	Slight consistency	
0.2 - 0.4	Fair consistency	
0.4 - 0.6	Moderate consistency	
0.6 - 0.8	Substantial consistency	
0.8 - 1	Almost perfect consistency	

Reliability
 → consistency
 of performance
 which is
 measured
 or when
 kappa =
 score

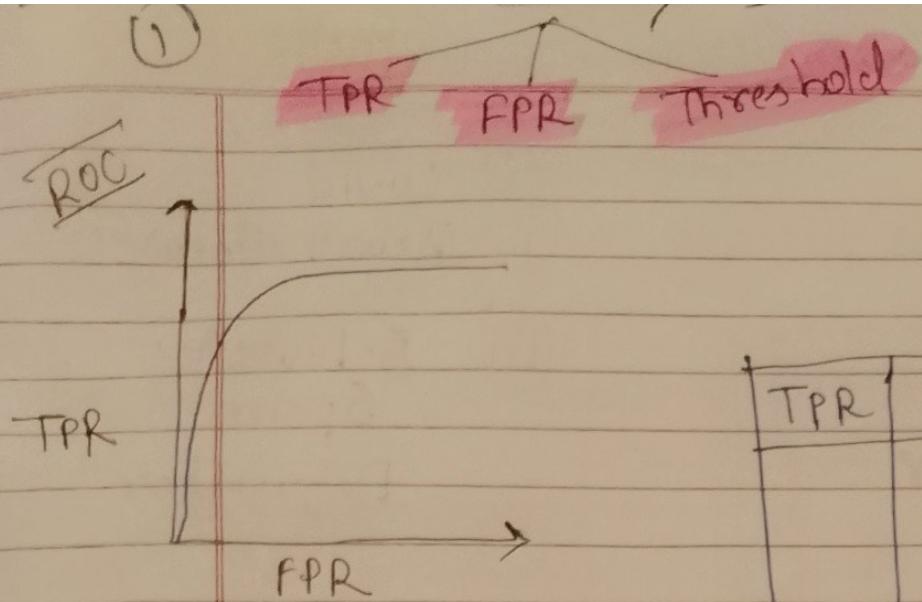
≤ 0.4] poor consistency
 > 0.4] good consistency

(46)

	C ₁	C ₂	C ₃	C ₄	C ₅	Y	P
R ₁					0	0	
R ₂					1	0	
R ₃					1	1	
R ₄					0	1	
R ₅					0	0	
					0	1	
R ₁₀₀					1	0	

(47)

(1)

Page _____
Date _____

FPR \rightarrow min.
TPR \rightarrow max.

TPR	FPR	Threshold

Good model

Diffr is high
can't get a
good model.

(2)

Youden's Index

FPR	TPR	Threshold	Difference
0	0.89	0.78	0.89
0.01	0.89	0.73	0.88
0.02	0.91	0.62	0.89
0.02	0.95	0.61	0.93
0.03	0.95	0.55	0.92
0.04	0.96	0.46	0.92
0.06	0.97	0.35	0.91
0.09	0.97	0.189	0.88
0.09	1	0.167	0.91

① The objective is FPR should be as Min. as possible And TPR should be as max as possible

② Hence the difference between TPR - FPR will be max

③ Take that Threshold comes pending to the max. diff. and develop the model

(3)

Imbalanced Dataset

Shot on AWESOME A05s

④

2 classes

Negative(0)

- (i) No covid
- (ii) No heart disease
- (iii) No Diabetics
HAM
Non Defaulters

Positive(1)

- i covid
- ii heart disease
- iii d. Diabetic
Spam
Defaulter

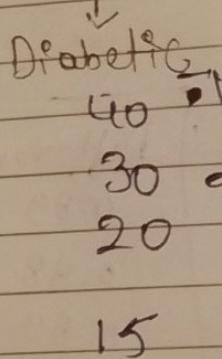
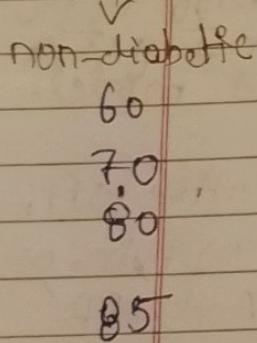
⑤

Balanced Dataset

Target Variable(y)

100

y



fairly balanced

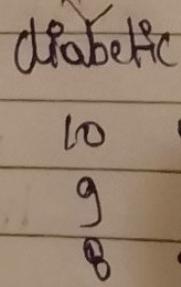
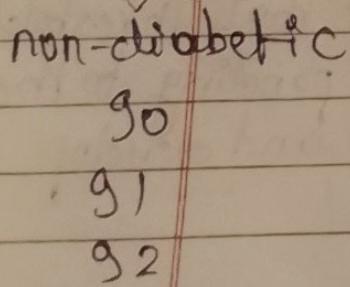
⑥

Balanced Dataset

Target Variable(y)

100

y



Imbalanced dataset.

This is not a rule
but a guideline.

90-10 is considered a
imbalanced dataset



Shot on AWESOME A05s

FLAM

(1)

⑦

Handling imbalanced Data.

⑧ Up sample minority class ✓

⑨ Down sample majority class ✓

change the performance metric ✓

try synthetic sampling approach ✓

Use different algorithm ✓

⑩ non-disease
90
majority

Page
Date

Up sample minority class

→ closer to go vs 90

⇒ problem in
redundancy
of data

Duplicacy of
data!

⑪

SMOTE

Imbalance Dataset

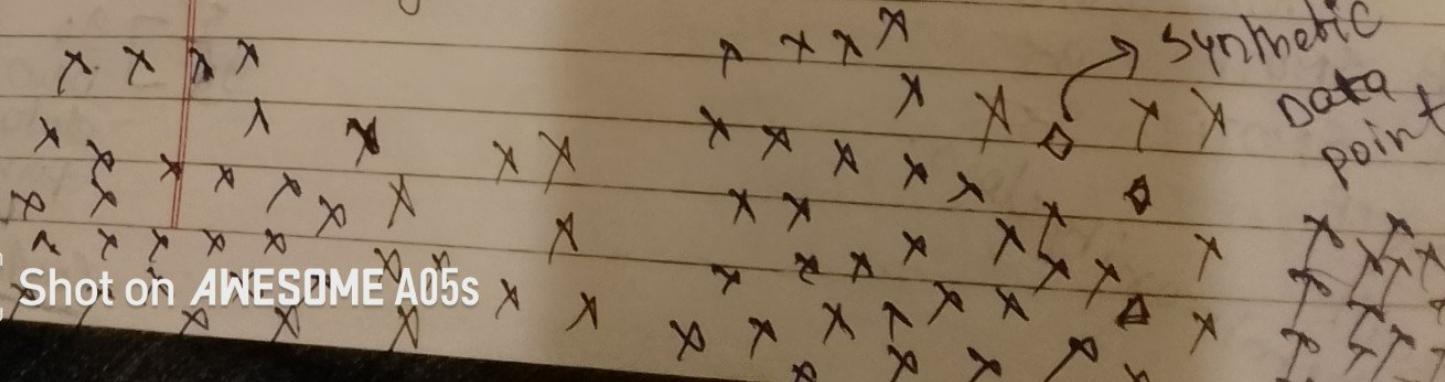
} Synthetic minority
oversampling
Technique

① Synthesizing data points for the minority class.

Synthesizing means generating new data points
based on existing data points.

② Function randomly picks up a data point from
minority class and applies kNN technique.

③ New datapoint which is called Synthetic data
point is added between the point picked up
and its neighbours



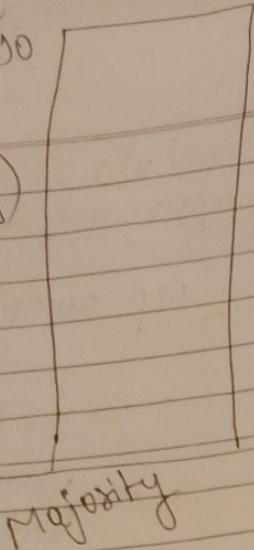
©

Shot on AWESOME A05s

⑩ non diabetic go

Majority

⑪ Normal
go



Diabetic Patient

Minority)

10

Minority

Closer to go

⇒ problem is
Redundancy
of data

⇒ Duplicacy of data

⑫

⑫ So Using SMOTE you will see the better performance of the model

⑬ SMOTE is balancing the number of observation in negative and positive class

⑭ SMOTE is applied only on Training data

⑮ So when SMOTE is applied on Training data Model starts learning the pattern well and hence gives better Result.

⑯

classes,

No covid

B SMOTE 95

A SMOTE 95

SMOTE

Minority class

⑰ Shot on AWESOME A05s

Covid

95]

new
data

pa
co

[] Shot on AWESOME A05s

(22)

Variable

	Numerical	Categorical
-	Euclidean	Cosine
-	Squared Euclidean	Hamming
-	Manhattan	Manhattanobis
-	Chebychev	
-	Minkowski	

(23)

DISTANCE MEASURES

Depending on the type of data we have different distance measures

Numerical Data

- Euclidean distance
- Manhattan distance
- Minkowski distance
- Chebychev's distance

String Data

- Cosine distance
- Edit distance
- Longest Common Sequence
- Hamming distance

(24)

Edit Euclidean distance -

numerical data

Consider two points $(5, 6)$ and $(1, 3)$. Obtain the Euclidean distance

$$\sqrt{\sum_{i=1}^n (x_i - y_i)^2} = \sqrt{(5-1)^2 + (6-3)^2}$$

$$\sqrt{16+9} = \sqrt{25} = 5$$

(25)

 β ($\beta^{(1)}$)

$$\sqrt{(x_1 - y_1)^2 + (x_2 - y_2)^2 + (x_3 - y_3)^2}$$

$$\sqrt{(1-3)^2 + (2-4)^2 + (3-5)^2}$$

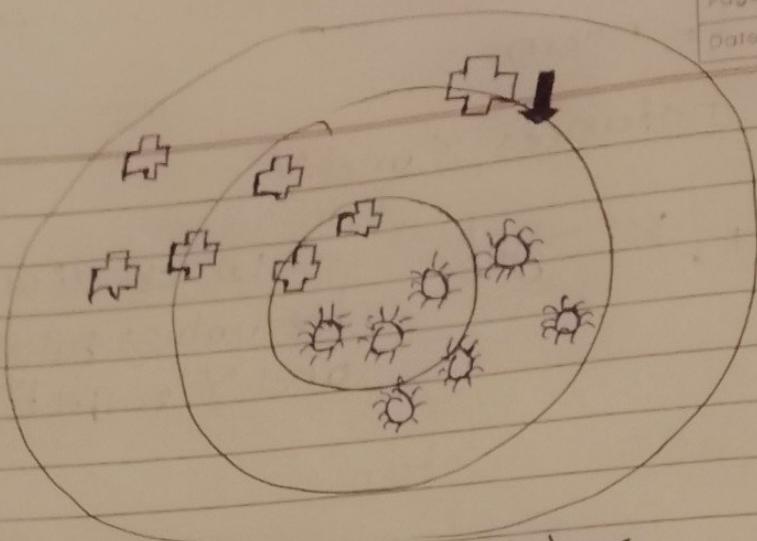
$$\sqrt{(-1)^2 + (-2)^2 + (-2)^2}$$

$$\sqrt{4+4+4} = \sqrt{12} = 2\sqrt{3}$$

Euclidean
method

Shot on AWESOME A05s

(34)

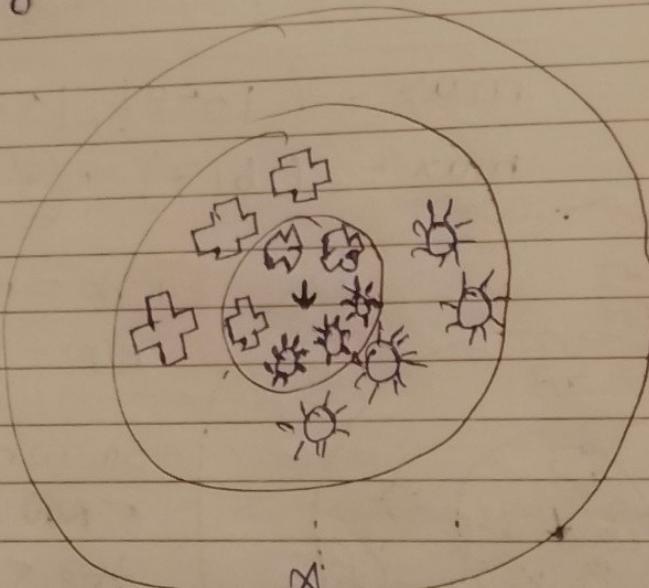


Page:

Date:

- ① KNN is impacted by outliers
- ② $k=1$ is very small to decide the nearest neighbour

(35)



$k \rightarrow$ Avoid having
K value of 1
as there is a
possibility of
tie

(36)

	Observation	X	Humidity	Y	Temperature	Rain
1		58	19	0	0	0
2		62	26	0	0	0
3		40	30	0	0	0
4		36	35	0	0	0
5		87	19	1	0	0
6		93	18	1	0	0
7		79	16	1	0	0
8		69	17	1	0	0
9		62	33	0	0	0
10		71	15	1	0	0

9	55	33
12	78	19
13	60	20
14	58	35
15	35	39

A → 1 → D₁
 A → 2 → D₂
 A → 3 → D₃
 A → 4
 A → 5

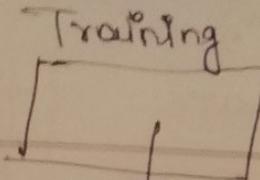
$$\begin{aligned}
 & \sqrt{(84-58)^2 + (37-14)^2} \\
 & \sqrt{(84-62)^2 + (37-26)^2} \\
 & \sqrt{(84-40)^2 + (37-30)^2}
 \end{aligned}$$

(H) (T)
 (84, 37)

Observation	Euclidean distance (sorted)	Class Label (rainfall)
5	18.25	1
12	18.97	1
6	21.02	1
7	21.59	0
9	22.96	0
2	24.6	1
8	25.00	1
10	25.53	0
14	26.08	0
11	29.27	1
13	29.41	0
1	31.62	0
3	34.55	0
4	48.04	0
15	49.04	0

It uses Voting
 principal to decide
 the class.

✓ KNN
Hence this is called lazy learner



It is only memorizing the coordinates of data points

Testing
All calculations are done

It is calculating the euclidean distance

(15)

calculations are done in testing phase
Hence it acts hazy to do all calculation at the end.

(16)

KNN

Advantages

- ① Easy to implement
- ② No training required
- ③ New data can be added at any time
- ④ Effective if training data is large

Disadvantages

- ① To choose apt value for k
- ② Computational expensive
- ③ Can not tell which features gives the best result

(17)

KNN Applications

→ Image Classification

→ Handwriting recognition

→ Predict credit rating of customers

→ Replace missing values } → Just KNN Imputer

Shot on AWESOME A05s

We have one application

Benson
Day 6

Naive bayes

① Probability based

It is totally based
Probability based

- ① Probability is how likely an event is to occur
- $P = \frac{\text{No. of ways an event can occur}}{\text{Total possible events.}}$

- ② The probability of an event always lies in between 0 and 1

- ③ 0 indicates impossibility of the event and 1 indicates a certain event.

Q: There are 40 candidates in a team with equal calibre. Out of which 25 are men and 15 are women. A person is randomly chosen to be team leader. What is the probability that the person is a woman?

$$\text{No of ways event can occur} = 15$$

$$\text{Total no. of outcomes} = 40$$

$$\text{So the probability: } 15/40 = 0.375$$

- ② Bayes Theorem \rightarrow Conditional Probability

$$\text{Cond Prob} = P(A/B)$$



Finding the probability of event A given event B has already happened.

$$\text{Cond Prob} = P(B/A)$$

Finding the probability of B (event) given A has already happened.



Note

$$P(A/B) = \frac{P(A \cap B)}{P(B)} \quad \checkmark \quad - \text{eq } \textcircled{1}$$

$$P(B/A) = \frac{P(A \cap B)}{P(A)} \quad \checkmark \quad - \text{eq } \textcircled{2}$$

By eq $\textcircled{1}$ & eq $\textcircled{2}$ \checkmark

$$P(A \cap B) = P(A/B) \cdot P(B) \quad \leftarrow \textcircled{3}$$

$$P(A \cap B) = P(B/A) \cdot P(A) \quad \leftarrow \textcircled{4}$$

$$P(A/B) \cdot P(B) = P(B/A) \cdot P(A) \quad \text{By eq } \textcircled{3} \text{ & } \textcircled{4} \quad \checkmark$$

\downarrow
P
(given information)

$$\boxed{P(A/B) = \frac{P(B/A) \cdot P(A)}{P(B)}} \quad \text{Bayes Theorem}$$

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

given information

$$\boxed{P(B/A) = \frac{P(A/B) \cdot P(B)}{P(A)}} \quad \text{Bayes Theorem}$$

Note Bayes theorem uses conditional probability

conditional prob. ques

- A pair of fair dice is rolled. If the sum of numbers that appear is 6, find the probability that one of the dice shows 2?

$$P\left(\frac{\text{loan} = \text{NA}}{\text{CS} = \text{High}}\right) = \left(\frac{1}{2}\right) \cdot \left(\frac{2}{3}\right) / \left(\frac{2}{3}\right) = \underline{\underline{\frac{1}{3}}}$$

$$P\left(\frac{\text{loan} = \text{Approved}}{\text{CS} = \text{High}}\right) = \left(\frac{2}{3}\right)$$

$$P\left(\frac{\text{loan} = \text{NA}}{\text{CS} = \text{High}}\right) = \left(\frac{1}{3}\right)$$

$$P\left(\frac{\text{loan} = \text{A}}{\text{CS} = \text{High}}\right) > P\left(\frac{\text{loan} = \text{NA}}{\text{CS} = \text{High}}\right)$$

→ Hence the person will be given a loan.

BAYES THEOREM - FORMULA

$$P(t/x) = \frac{P(t) \cdot P(x/t)}{P(x)}$$

conditional prob. of t given that predictor X i.e. posterior probability

Probability of class it is prior probability

conditional Prob. of X given that its class label is t, i.e. likelihood

Probability of value taken by the predictor variable.
i.e. evidence.

$$P(A/B) = P(B/A) \cdot P(A) \rightarrow \text{(prior probability)}$$

Posterior probability.

P(B)

likelihood

Evidence:

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

$$P(t|x) = \frac{P(x|t) \cdot P(t)}{P(x)}$$

where

t = target variable

x = independent variable

$$P\left(\frac{\text{loan}=1}{\text{CS=high}}\right) = \frac{P\left(\frac{\text{CS=high}}{\text{loan}=1}\right) \cdot P(\text{loan}=1)}{P(\text{CS=high})}$$

$$P(t|x) = \frac{P(x|t) \cdot P(t)}{P(x)} \quad \checkmark$$

$$P(t/x_1, x_2) = \frac{P(x_1, x_2 | t) \cdot P(t)}{P(x_1, x_2)} \quad \checkmark$$

Salary
Experience
etc.

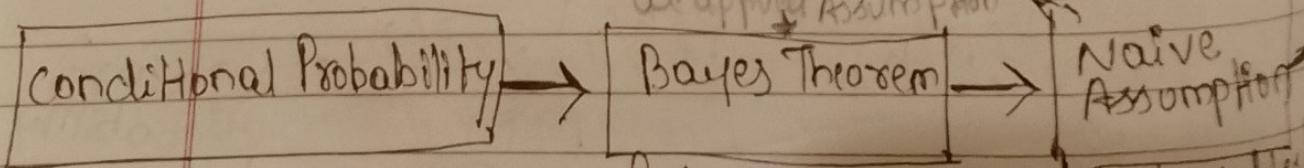
This where
Bayes theorem
talks.

for classification

Uses Bayes Theorem

Assumes Independence between the variables - Naive

Becoz in real life
we apply Assumption



$$P(A|B) = P(ANB)/P(B)$$

Prob. of A given B

$$P(A|B) = P(ANB)/P(B)$$

Prob. of B given A

$$P(B|A) = P(ANB)/P(A)$$

Variables
features
are
indep.
of each
other

when
indep
multi

P(Diabetic
A)

- Age
- Blood
- Ins
- BM

If B rep
real val
param
diabetic

P(Diabetic)

P



Shot on AWESOME A05s

$$P(A \cap B) = P(B|A) * P(A)$$

lets say that parameters are.

- Glucose
- Blood Pressure
- Insulin
- BMI

If B represents condition leading to diabetic. In real world, there might be several conditions or parameters that might lead a patient to be diabetic.

$$P(\text{Diabetic} | \text{Glucose}, \text{BP}, \text{Insulin}, \text{BMI}) =$$

$$\frac{P(\text{Glucose}, \text{BP}, \text{Insulin}, \text{BMI} | \text{Diabetic})}{P(\text{Diabetic})}$$

$$= \frac{P(\text{Glucose}, \text{BP}, \text{Insulin}, \text{BMI})}{P(\text{Glucose}) * P(\text{BP}) * P(\text{Insulin}) * P(\text{BMI})}$$

naive

e
option
variables
atures
re
dependent
and
independent
variables

Naive Bayes

when events are
independent
multiplied

Naive means basic
assumption that variables
are independent of each
other.

$$P(t/x_1 x_2 x_3) = \frac{P(x_1 x_2 x_3/t) \cdot P(t)}{P(x_1 x_2 x_3)}$$

$$P(t/x_1 x_2 x_3) = \frac{P(x_1 n x_2 n x_3/t) \cdot P(t)}{P(x_1 n x_2 n x_3)}$$

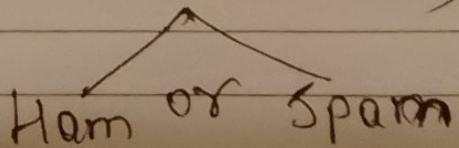
$$P(A \cap B \cap C) = P(A) \cdot P(B) \cdot P(C)$$

↳ if A, B, C are independent of each other then

$$= \frac{P(x_1/t) \cdot P(x_2/t) \cdot P(x_3/t) \cdot P(t)}{P(x_1) \cdot P(x_2) \cdot P(x_3)}$$

Session ①

- Naive Bayes is used only for classification
- It is probability based. This is one reason it can only be used for classification
- Wide application of NB → Mail Box
 - ↳ Predict(Mail)
- High priority tickets
- All classification examples



Lonely	2	1
Horoscope	20	5
work	5	12
Snacks	0	5
Money	21	7

Email (Good, work) → Spam
Ham

$$P\left(\frac{t}{x_1 x_2}\right) = \frac{P(x_1 \cap x_2 / t) \cdot P(t)}{P(x_1 \cap x_2)}$$

$t = \text{Ham} / \text{Spam}$

Good and work
 \downarrow \downarrow
 x_1 x_2

Assumption → Good is independent of Work

$$P\left(\frac{\text{SPAM}}{\text{Good}, \text{work}}\right) \rightarrow \textcircled{1}$$

$$P\left(\frac{\text{HAM}}{\text{Good}, \text{work}}\right) \rightarrow \textcircled{2}$$

$$P(\text{SPAM}/\text{Good}, \text{work}) >$$

$$P(\text{HAM}/\text{Good}, \text{work}) \text{ then}$$

it will be classified as SPAM
else HAM

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

```
from sklearn.preprocessing import LabelEncoder  
from sklearn.tree import DecisionTreeClassifier, plot_tree
```

Decision Tree

INFORMATION THEORY

Information theory is based on the intuition that.

Event	Information Gain	Example
Most likely event	No Information	The sun rose this morning
Likely event	Little Information	The sunrise at 6:30 AM this morning
Unlikely event	Maximum Information	There was a solar eclipse this morning

① Let $I(x)$ denote the information of an event X

② It is the self-information of an event X at x

$$I(x) = -\ln P(x)$$

③ Since we have considered natural log its unit is not

④ For log with base 2, we use units called bits or

⑤ Shot on **AWESOME A05s**

Shannons Entropy

- Entropy is the measure of information for classification problem, i.e. it measures heterogeneity of a feature

The entropy of a feature is calculated as

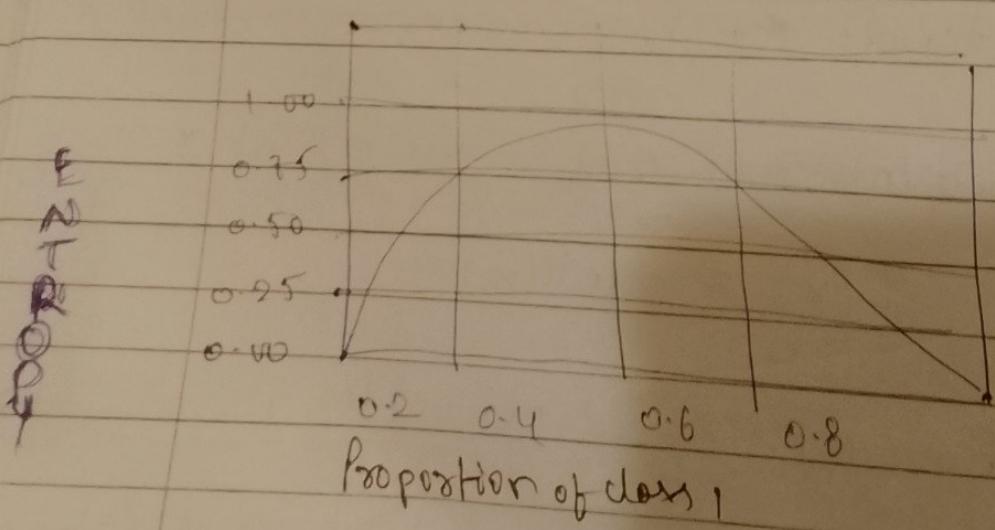
$$E = -\sum_{i=1}^C p_i \log_2 p_i \quad \text{where } p_i \text{ is the probability of occurrence of the class.}$$

A lower entropy is always preferred

Entropy is always non-negative

Consider a feature with two class the entropy for various proportion of classes is given below

	class 1	0	0.1	0.3	0.5	0.7	0.9	1
class 2	1	0.9	0.7	0.5	0.3	0.1	0	
Entropy	0	0.46	0.88	1	0.88	0.46	0	



$$\textcircled{1} \quad p\text{-default} = 6/20 \\ p\text{-non-default} = 14/20$$

Date _____
Page No. _____

$$\text{Entropy} = -(p\text{-default} * np.\log_2(p\text{-default})) + p\text{-non-default} * np.\log_2(p\text{-non-default})$$

$$\text{Entropy} \\ 0.8812908992306927$$

$$\textcircled{2} \quad p\text{-default} = 10/20 \\ p\text{-non-default} = 10/20$$

$$\text{Entropy} = -(p\text{-default} * np.\log_2(p\text{-default})) + p\text{-non-default} * np.\log_2(p\text{-non-default})$$

$$\text{Entropy} \\ 1.0$$

Conditional Entropy

The conditional entropy of one feature given others is calculated as from the contingency table of the two features

$$E(C|X) = \sum_{c \in C} P(c) E(c)$$

It is the sum of the product of probability of occurrence of the each class and the entropy of it

Q Obtain the condition entropy of loan given the rest of the person

loan

	reject	Approved.
Income Low	8	1
High	2	4

Entropy of loan? $E(\text{loan})$

Q Shot on **AWESOME A05s**

ent_loan
0.9122050340544896

Entropy $E(\text{loan} / \text{income})$

$$\text{ent_loan_given_income} = \left(\frac{2}{15}\right) * \text{entropy} \left(\frac{8}{15}, \frac{1}{3}\right) + \left(\frac{6}{15}\right) * \text{entropy} \left(\frac{2}{6}, \frac{4}{6}\right)$$

ent_loan_given_income
0.6692733344871833

INFORMATION GAIN

- ① Information gain is the decrease in entropy at a node
- ② Information Gain (T/x) = Entropy(T) - Entropy(T/x)
- ③ To construct the decision Tree, the feature with highest information gain is chosen
- ④ Information gain is always positive

Q Obtain the information gain in the feature loan due to income

$$IG_{\text{loan_income}} = \text{ent_loan} - \text{ent_loan_given_income}$$

IG_loan_income
0.24902249956730627

Decision Tree Algorithms II

Hunts Algorithm

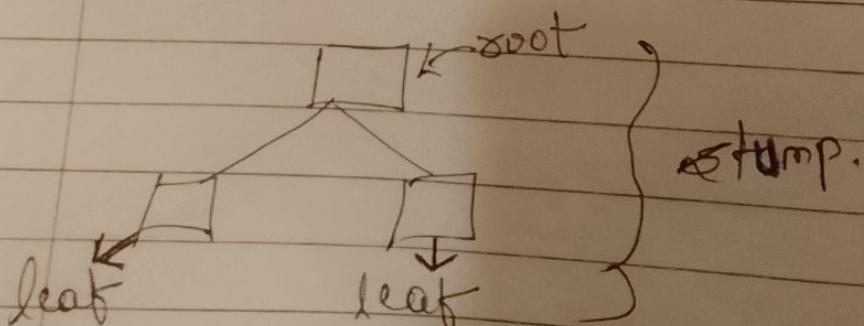
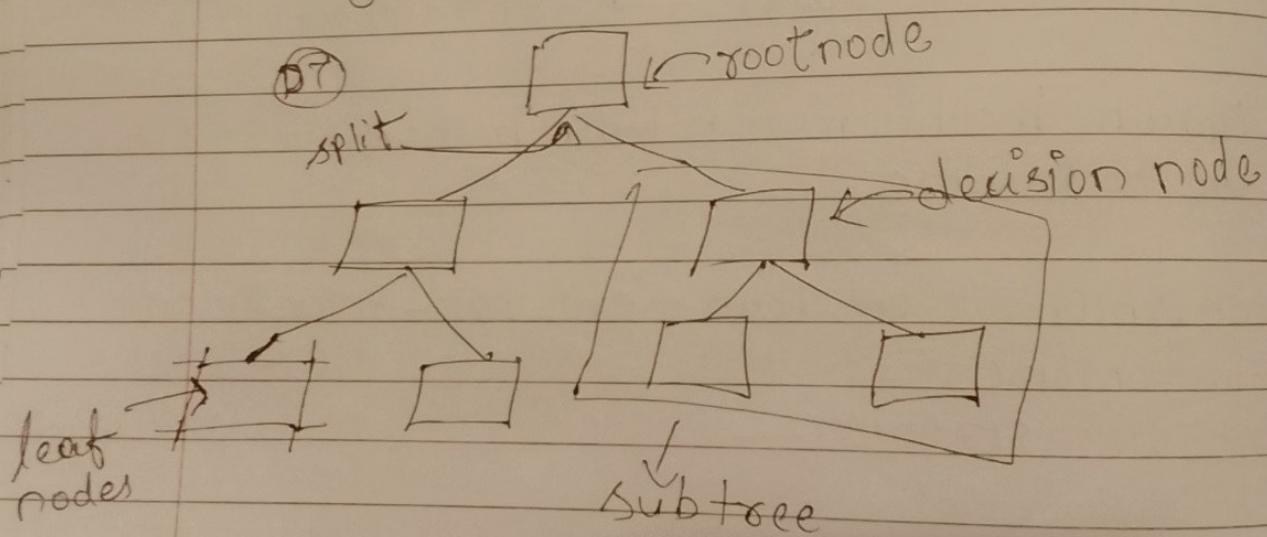
Date /
Page No.

let S_n to be the training samples associated with node n and y_c be the class labels. The algorithm is follows

- ① If all samples belong to the same class y_c , then n is a leaf node with label y_c
- ② If S_n has samples with more than one class, an attribute value is selected to partition the samples into smaller subsets such that the samples in the subsets belong to same class

↓
where information gain
that feature will be your root node

Hunts Algorithm



what's the purity of windy = False (Gini Index)

	Play	No	Yes
windy			
False	2	6	
True	3	3	

$$\text{gini-index} = 1 - ((2/8)^2 + (6/8)^2)$$

$$\text{gini-index}$$

$$0.375$$

$$\text{gini-index} = 1 - ((3/6)^2 + (3/6)^2)$$

$$\text{gini-index}$$

$$0.5$$

Yes overcast class in outlook.

$$\text{gini-index} = 1 - ((0/4)^2 + (4/4)^2)$$

$$\text{gini-index}$$

$$0.0$$

Classification Error

① The classification error of a variable is calculated

①

$$\boxed{\text{Error} = 1 - \max p_c^2}$$

where p_c : probability of occurrence of the class

② For samples belonging to one class, the classification error is 0 and for equally distributed samples,

the classification error is 0.5

③ Shot on **AWESOME A05s**

Decision Tree in sklearn

Date / /

Page No.

data

```
x = data.drop(['columns': 'Play'])
```

```
y = data['Play']
```

```
x.dtypes
```

```
dt = DecisionTreeClassifier(criterion='entropy')
```

```
le = LabelEncoder()
```

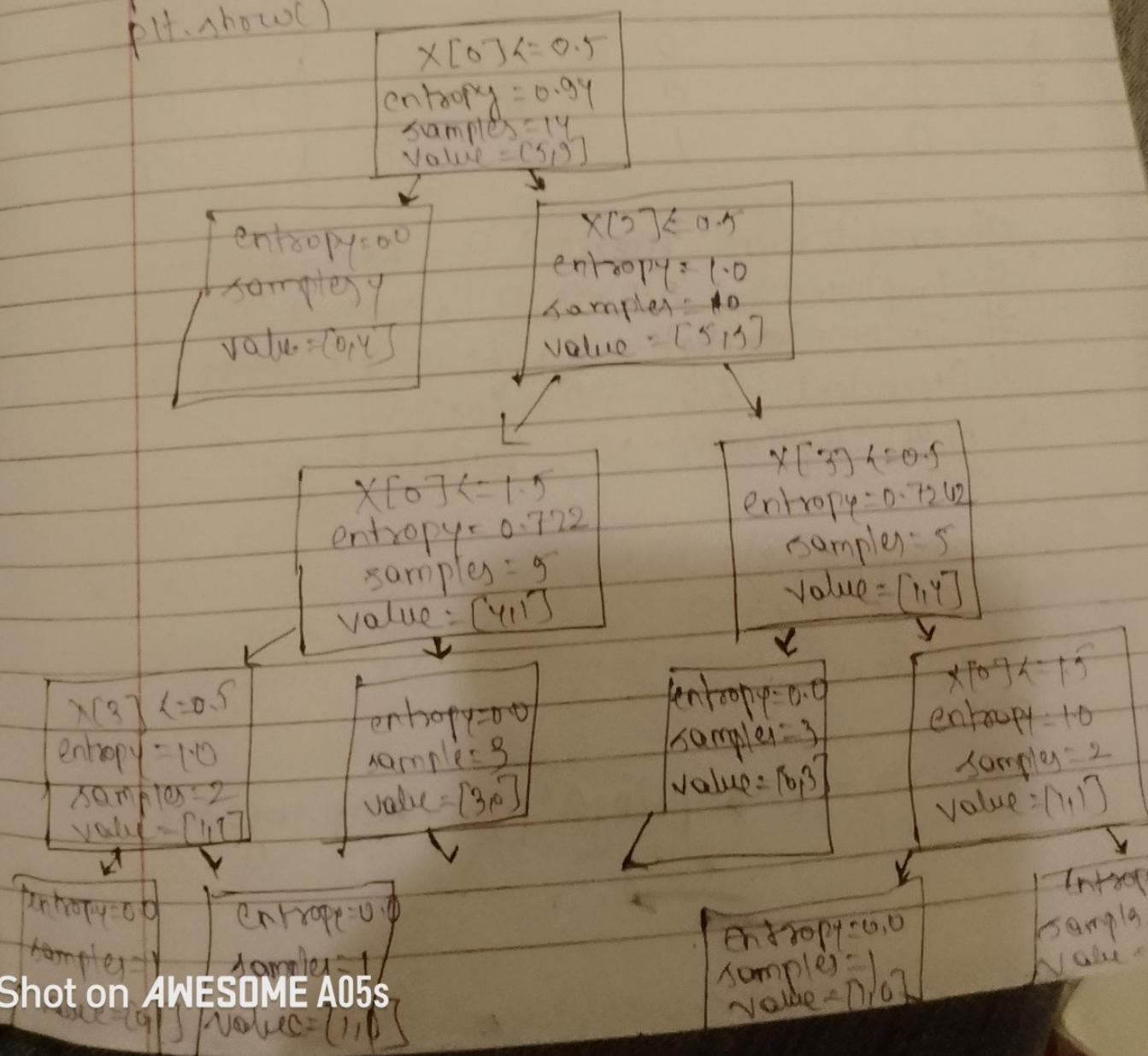
```
for i in data.columns:
```

```
    data[i] = le.fit_transform(data[i])
```

```
plt.figure(figsize=(7,5))
```

```
plot_tree(dt)
```

```
plt.show()
```



ENSEMBLE LEARNING

Date / /
Page No.

- Ensemble learning algorithms combine multiple models into one predictive model.
- Decisions from several weak learners are combined to increase the model performance.

ENSEMBLE LEARNING MODELS

Bagging

Homogeneous models can be built — independently and their outputs are aggregated at the end.

Example
Random Forest

Boosting

Homogeneous models can be built sequentially, previous model dictates the features the succeeding model will focus on.

Example
Ada Boost

Stacking

Heterogeneous base models can be built, outputs from former model are used as inputs to the meta model.

Example
Voting Classifier

RANDOM FOREST CLASSIFIER

The subsets of variables are selected at random to build a decision tree.
Forest of "Decision" trees.

- Random forest consists of several independent decision trees that operates as an ensemble.

Shot on AWESOME A05s

It is an ensemble learning algorithm based on bagging.

→ Train decision tree models on bootstrap samples where variables are selected at random. The aggregate output from these trees is considered as the final output.

BOOTSTRAP

- ① Bagging is composed of two parts
 - ① Bootstrapping
 - ② Aggregation
- ② Random sampling with replacement.
- ③ For a data with i observations, a random sample with replacement of size i , is a bootstrap sample.
- ④ The observations n_3 and n_5 are not included in bootstrap sample; they are the out-of-bag (OOB) samples.

Original Data

n_1
 n_2
 n_3
 (n_4)
 n_5
 n_6

Bootstrap Sample

n_1
 n_5
 n_1
 n_5
 (n_4)
 n_2

... (with replacement)

Plots
did you know?

Out-of-bag sample

Date / /
Page No.

Almost 36.8% of the training data has the potential to be out-of-bag sample

How?

Since, we consider the bootstrap sample i.e. a random sample with replacement. The probability of not selecting a sample is $(1 - 1/N)^N$. For large N , the probability $(1 - 1/N)^N$ is approximately $e^{-1} \approx 0.368$.

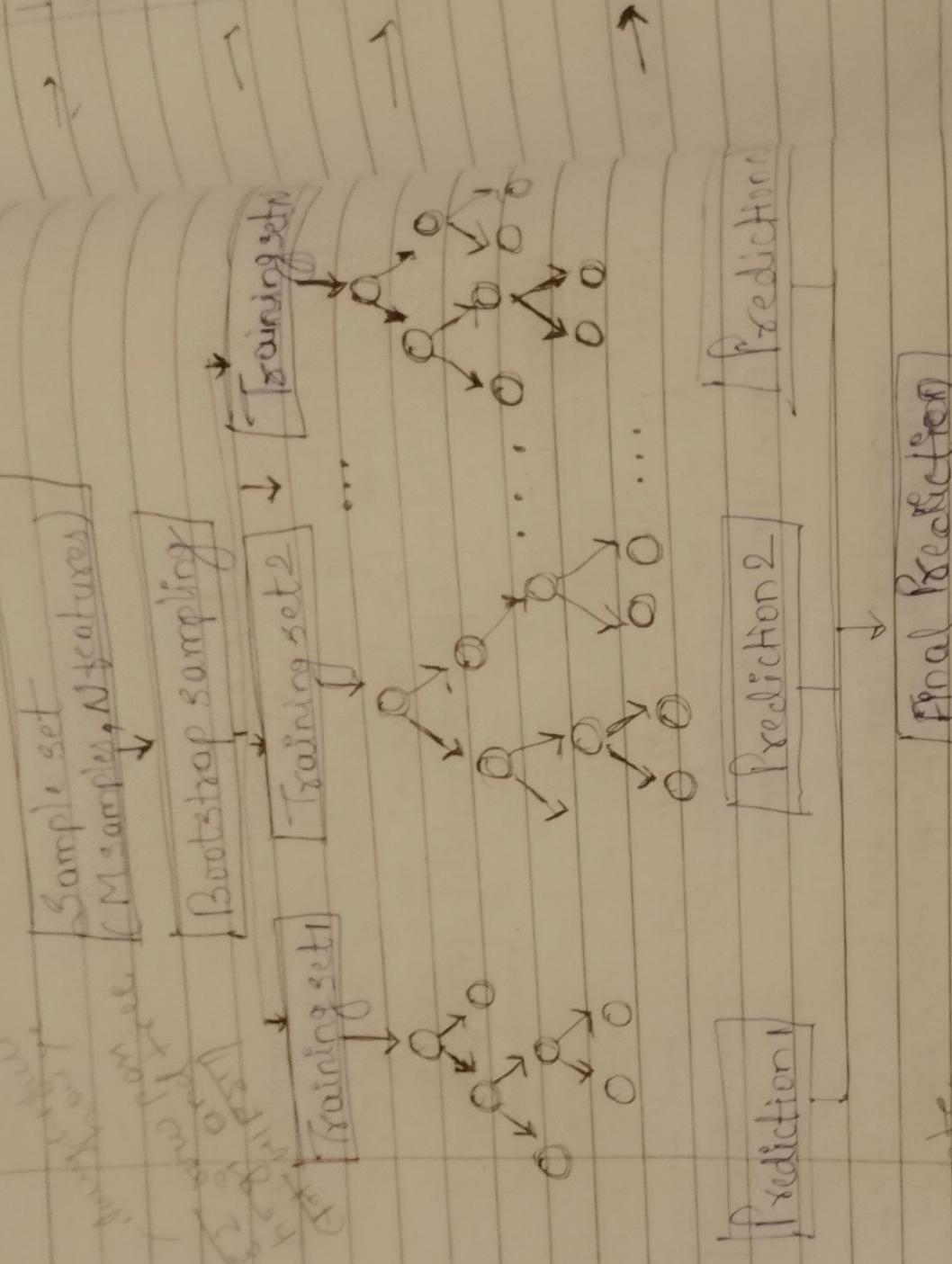
Aggregation

Model predictions undergo aggregation to combine them for the final prediction, to consider all the possible outcomes.

The aggregation can be done based on the total number of outcomes, or the probability of predictions derived from the bootstrapping of every model in the procedure.



Steps for prediction using Random forest



To note

For classification, the final prediction considers the mode of the predicted labels and for regression, it considers the average of the predicted values.

Advantages of Random Forest

- Improved Accuracy — Random Forest often achieves higher accuracy compared to individual models
 - Increased Robustness — Random Forest more robust to outliers and noisy data
 - Reduced Variance — Random forests reduce variance by taking many trees from n bootstrapped samples and averaging their prediction leading to lower variance.
 - Reduced Bias
 - Random forests also reduce bias in the dataset by decorrelating the trees which is a problem with imbalanced datasets. Random forests overcome this problem by forcing each split to consider only a subset of predictors.
 - Therefore, on average $(p-m)/p$ of the splits will not even consider the strong predictor and so other predictors will have more of a chance.
 - We can think of this process as decorrelating the trees, thereby making the average of resulting trees less variable and hence more reliable.
 - Reduced Overfitting — It helps in reducing overfitting especially when individual models tend to overfit the training data.
- ◎ Shot on **AWESOME A05s**

Disadvantages of Random Forest

→ Increased Complexity

Ensembles introduce additional complexity to the modeling process

→ Computational Overhead

Ensemble methods require retraining and maintaining multiple models, which can be computationally expensive especially when dealing with large datasets or complex models.

→ Interpretability

Ensemble models are typically less interpretable compared to individual models.

Ensemble Technique (Bagging)

Random Forest

```
from sklearn.ensemble import RandomForestClassifier
```

```
rfc = RandomForestClassifier(n_estimators=5, max  
⑤ and  
features=None)
```

```
rfc.fit(x_train, y_train)
```

RandomForestClassifier

```
RandomForestClassifier(max_features=None, n_estimators=5)
```

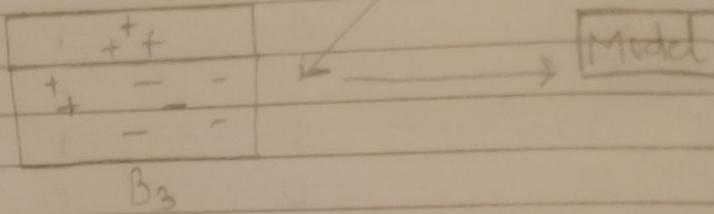
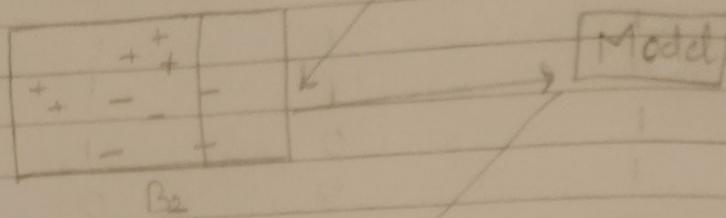
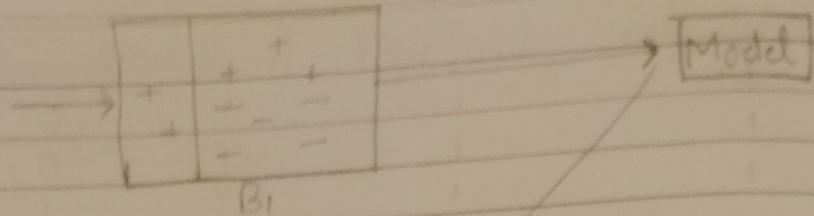
Ada Boost

Boosting

Step	Credit-Score	Employed	Income	Dependent	Loan
1.	1	1	1	0	1
1.	1	1	1	0	1
1.	1	1	1	1	1
1.	1	1	1	1	0
1.	1	1	1	1	0
Excel data	1	1	1	1	0
1.	1	0	1	1	0
1.	1	0	0	0	1
Excel pe 2015/11/21	0	1	1	1	1
Ada	0	0	1	0	1
Boost	0	1	0	0	1
Boosting	0	0	0	0	1
example	0	1	0	0	1
Hal	0	0	1	0	1

What is Boosting

- The ensemble of weak learners that learn sequentially
 - In each iteration weights of the samples are adjusted, such that the misclassified samples have a higher weight, therefore higher chance of getting selected to train the next classifier
 - Boosting reduces bias and variance
 - Here B_1, B_2 and B_3 are base learners and B_4 is an ensemble of weak learners



A 3x3 grid representing the ensemble model B_4 . The columns are labeled with '+' and '-' signs. The rows are labeled with '+' and '-' signs. The entries in the grid are: Row 1, Col 1: '+'; Row 1, Col 2: '+'; Row 1, Col 3: '+'; Row 2, Col 1: '+'; Row 2, Col 2: '-'; Row 2, Col 3: '-'; Row 3, Col 1: '-'; Row 3, Col 2: '-'; Row 3, Col 3: '-'. To the right of the grid is the equation $B_4 = B_1 + B_2 + B_3$.

→ The ensemble of weak learners that learn sequentially

→ In each iteration weights of the samples are adjusted, such that misclassified samples have a higher weight, therefore higher chance of getting selected to train the next classifier

→ Boosting reduces bias and variance

→ Here B_1 , B_2 and B_3 are base learners and B_4 is the boosting ensemble of weak learners

Boosting = advantages

- Enhances the efficiency of weak classifiers
- Both precision and recall can be enhanced through boosting algorithms

Boosting = Disadvantages

- Loss of simplicity and explanatory power
- Increased computational complexity

How boosting differs from bagging?

BAGGING

Base learners learn in parallel

Random Sampling

Reduces variance

BOOSTING

Base learners learn sequentially

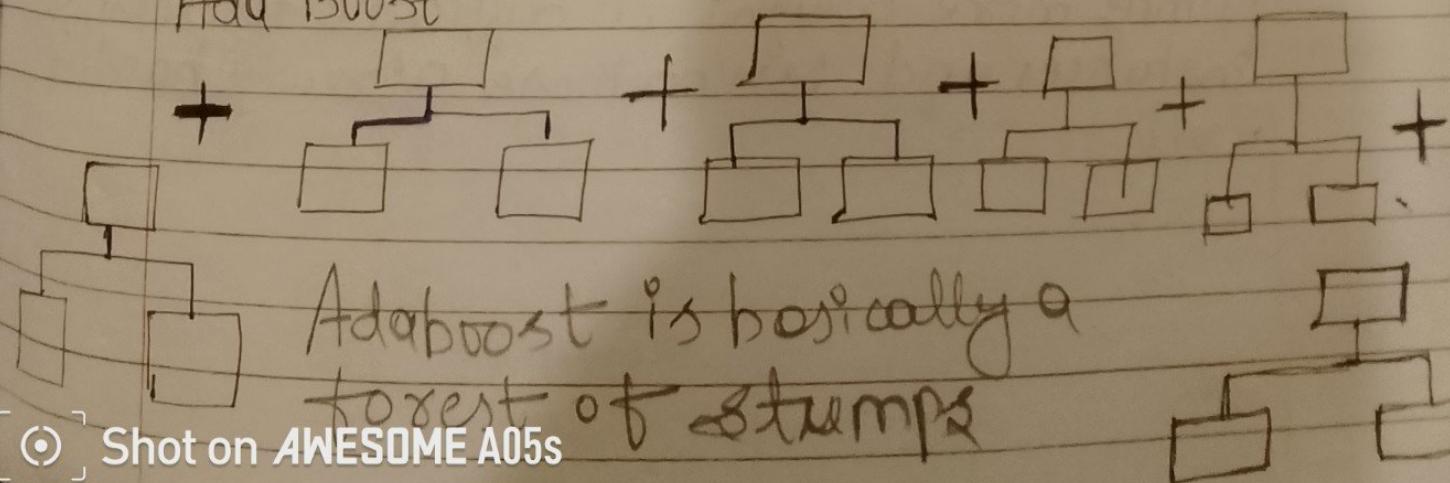
Non-Random Sampling

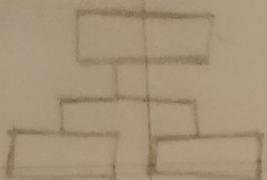
Reduces bias and variance

BOOSTING ALGORITHMS ARE :-

- Ada Boost
- Gradient Boosting
- XGBoost

Ada Boost





Stump is a tree with just one parent
Date _____
node and two leaves Page No. _____

Step 9

Credit

- 1
- 2
- 3
- 4
- 5
- 6
- 7
- 8
- 9
- 10
- 11
- 12
- 13
- 14
- 15

- A single algorithm or a stump may classify the objects poorly
- But if we combine multiple classifiers with a selection of training set at every iteration and assigning the right amount of weight in the final voting, we can have good accuracy score for overall classifier

Ada Boost

- Ada Boost is short for Adaptive Boosting
- It was formulated by Yoav Freund and Robert Schapire
- It is sensitive to noisy data and outliers
- Usually, decision trees are used for modelling
- Multiple sequential models are created, each correcting the errors from last model
- Ada Boost works by weighing the observations, putting more weight on difficult to classify instances and less on those already handled well

Step 10

Co
Ro
t

#

0.887943
(3)

Step 2 Initialize 0.067 to all rows
 ↓
 Initial weight Date / /
 Page No.

	Credit-Score	Employed	Income	Dependent	Loan	Initial Weight
1	1	1	1	0	1	0.067
2	1	1	1	0	1	0.067
3	1	1	1	0	0	0.067
4	1	1	1	1	0	0.067
5	1	1	1	1	0	0.067
6	1	1	1	1	0	0.067
7	1	0	1	1	0	0.067
8	1	0	0	0	1	0.067
9	0	1	1	1	1	0.067
10	0	0	1	0	1	0.067
11	0	1	0	1	1	0.067
12	0	1	1	0	1	0.067
13	0	0	0	0	1	0.067
14	0	1	0	0	1	0.067
15	0	0	1	0	1	0.067

Step 3 calculate entropy

Row Labels	Count of Credit-Score		Column Labels	Grand Total
	0	1		
0	0	7	7	
1	5	3	8	
Grand Total	5	10	15	

Credit-Score

0 7

5 3

$$(7/15) * \text{entropy}(0/7, 7/7) + (8/15) * \text{entropy}(5/8, 3/8)$$

$$= 0.509031468226648$$

1 4

Employed

4 6

$$(5/15) * \text{entropy}(5/5, 4/5) + (10/15) * \text{entropy}(4/10, 6/10)$$



0 4 Income

5 6

$$(4/15) * \text{entropy}(0/4, 4/4) + (11/15) * \text{entropy}(5/11, 6/11)$$

$$0.7289534884164346$$

1 8 Dependent

9 2

$$(9/15) * \text{entropy}(1/9, 8/9) + (6/15) * \text{entropy}(4/6, 2/6)$$

$$0.66$$

Step 1 Stump① लंगाईँ and stump② को क्या करें
 (निम्नी entropy शब्द का आहु है, उसमें Info.
 ब्रॉड सबसे उचित है, और वही से stump①
 लंगाईँ।

1

2

8

Stump①

i

credit < 0.5
sample 15
(5/10)

(iv) Total

Total

No

v

class①
COP
Credit-Score

Prediction (5/3)

class②

1	1	0	अवधि predictor
2	1	0	column अवधि
3	1	0	जाहिरा and
4	1	0	प्रति credit score
5	1	0	प्रति 1 रुपये 0
6	1	0	भविष्यातीर्ण अवधि
7	1	0	1 रुपये 0
8	1	0	अवधि अवधि का
9	0	0	Credit Score का
10	0	0	अवधि 1 अवधि का

(iii)

3rd prediction (Out column which has
worse loan Main compare and
misclassification data find out)

Credit-Score	Loan	Initial-Weight	Prediction
1	1	0.067	0
2	1	0.067	0
3	1	0.067	0

{ 2) Misclassification
data has ↑

(iv) Total error find out & like no. of misclassification data \div by
Total no. of data.

$$\text{Total error} = \frac{3}{15} = \underline{\underline{0.2}}$$

(v) Amount of say = $(1/2) * \log[(1-TE)/TE]$ formula
 $= 0.69315$

Amount of say and Total Error

→ when stump does a good job then Total Error is small
Amount of Say will be a large positive value

→ When stump classifies only half of the samples correctly and half incorrectly, then Total Error is 0.5. Amount of say will be zero

→ When stump does a very poor job i.e. the stump gives an opposite classification, then the total error would be close to 1 and the Amount of say will be a large neg. Value

Decision Tree

- ① It deals with categorical entities of target variable, where relationship is not linear.
- ② Can be used for both regression & classification.
- ③ It is like a flowchart that makes decision based on a series of questions and answers. It starts of a "root Node" with question of "leaf node" (represents decision).
- ④ It divides data into tree-like structure (root node of leaf node).
- ⑤ This distribution is based on Gini Index & Entropy.
- ⑥ This process combines until we reach the stopping criteria, such as max-depth or a min no. of samples at node.

Ex → Bank wants to decide if person is eligible for loan based on credit score.

← entropy → Measure of Information

- Is always +ve
- lower entropy is preferred.

∴ $[-\sum P_c \log_2(P_c)]$:- P_c = Prob of occurrence
Here \log represents no. of classes in target of class.

← Conditional entropy

helps to measure how much uncertainty is there in one thing when we know the information about another thing.

$$\{ E(T|X) = \sum_{x \in X} P(x) E(x) \}$$

↳ Shot on AWESOME A05s

Decrease in Entropy at a node:

$$E(T|X) \{ 1 - \text{entropy} \}$$

- ① To make decision tree highest information gain is used
- ② Always +ve

→ Measure of Impurity ① Gini Index ② Classification error
 ③ Entropy $(1 - p_1^2 - p_0^2)$

- ④ Gini Index (1 - $p_1^2 - p_0^2$)
- ⑤ Samples belonging to same class, $G.I. = 0$
- ⑥ For equally distributed sample, $G.I. = 0$
- $G.I.(T/X) = \text{entropy}(T) - \text{entropy}(T/X)$
- $\text{Info. gain}(T/X) = G.I.(T) - G.I.(T/X)$

$T = \text{target}$
 $X = \text{features}$

- ⑦ Classification error
- ⑧ Samples of same class, $CE = 0$
- ⑨ Samples of equally dist. $CE = 0.5$

⑩ Pruning Reduce complexity of tree.

- ⑪ Pre Pruning It stops tree before it fully grows
- ⑫ Post Pruning Tree is allowed to grow completely then prune

To reduce Overfitting, we use pruning by passing Hyperparameters

- max_depth
- max_samples_split
- max_sample_leaf
- max_leaf_nodes
- min_sample_leaf
- min_sample_split

⑬ Importance of Node

- [% of data in parent * Ent. of Parent - % of data in left child * Ent. of left child - % of data in right child * Ent. of Right child]

⑭ Shot on AWESOME A05s

$$\text{Importance} = \frac{\text{Imp. of Node}}{\text{Imp. of all Nodes}}$$

XGBoost

- ① XGBoost combines the predictors of multiple weak learners to create a single, strong predictive model.
- ② It involves building models sequentially, where each new model focuses on correcting the errors made by previous ones.
- ③ The no. of boosting rounds (trees) is a tunable hyperparameter.
- ④ "XGBoost Assign weights" to data points based on errors made by previous trees (Data points that are more difficult to predict require higher weights).
- ⑤ XGBoost provides a measure of "feature importance" that helps us to identify which features have the most influence on model's predictions.
- ⑥ XGBoost has built-in handling of missing values.
- ⑦ XGBoost has 3 components :-

① Objective function

This defines loss function to be minimized during training

$$\left\{ \begin{array}{l} \text{Mean Squared error} \rightarrow \text{Regression} \\ \text{log loss} \rightarrow \text{classification} \end{array} \right.$$

② Weak learners (Base learners)

XGBoost typically uses decision tree as weak learners (stumps)

③ Regularization

XGBoost includes "L1 (Lasso)" & "L2 (Ridge)" regularization terms to control the complexity of model & reduce overfitting.



Shot on AWESOME A05s

Gradient Boosting

- Gradient Boosting can use a wide range of base learners, such as decision trees and linear models.
- It updates weights based on the gradients, which are less sensitive to outliers.

The boosting process begins with weak learners (base learners)

Step 1 This base learner is trained on the dataset and predictions are made.

Step 2 The difference b/w predicted values and the actual values is calculated and this difference is also k/s residual error, becomes target for next weak learners

Step 3 Weak learners are then trained to predict the residual error of previous learners. Each new learner focuses on minimizing the error made by the ensemble upto that point

Step 4 By combining predictions from multiple weak learners gradient Boosting is able to capture complex relationship in data and make accurate predictions

Gradient Boosting works well on numerical & categorical data and can handle large datasets

If the problem has a pattern of wise classifications or errors that can be fixed iteratively, Gradient Boosting can be interesting approach.



ADA Boost

- It is basically a forest of stumps
- Ada boost assigns weights to the predictions made by each weak learners by iterations
- Predictions of more accurate weak learners hold great weight in the final ensemble

Ada Boost learns from its past mistakes & adopts its learning process accordingly

weights:

Ada boost progresses through its iterations, it adopts & adjusts these weights based on the performance of the weak learners.

Error minimization: Once the weights are appropriate updated, Ada Boost proceeds to train a weak learner on the modified training data.

The aim is to minimize the error made by models during iteration 3

Step 1: Assign weights to each sample

In beginning each sample will have same weights

$$\text{Sample weight} = 1 / \text{No. of samples}$$

Step 2: Build the stumps

Build stumps with each variable, using Gini Index or entropy for each variable.

Smaller Gini Index = Base Learning Model

Shot on AWESOME A05s

Naive Bayes

Based on Probability
It will be helpful if an algorithm can label required received emails as important () emails or junk (spam) emails for a user

Probability: It is how likely an event is to occur

$P \Rightarrow \frac{\text{No. of ways an event can occur}}{\text{Total possible events}}$

Prob of an event always lies in b/w 0 & 1.

0 indicates impossibility of event and 1 indicates a certain event

Bayes Theorem: \rightarrow Conditional Probability

$P(A|B) \rightarrow$ finding the prob. of event A given event B has already happened

$P(B|A) \rightarrow$ finding the prob. of B given event A has already happened.

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

BAYES THEOREM

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

posterior prob.

Naive Assumptions

Variables / features are independent of each other

Note When we apply Bayes theorem with a Naive assumption that events are independent then they can be multiplied.

- ① Naive Bayes is used only for classification
- ② It is prob-based. This is the reason it can only be used for classification
- ③ Application of NB → Mail Box
 - Ham
 - Spam

Note
④ Naive Bayes fails when frequency of occurrence of word = 0
⑤ To overcome this problem we introduce Laplace theorem
Increase the frequency to avoid 0
Binary Cross Entropy

OOB

Entire Node

