

Regression Analysis (predictive Data Analysis)
predict real valued output for given input.

i/p
Amount of spend for Ad
Amount of carbon-14
Rainfall in previous month
previous stock prices
No of cigarettes
size of house, area, no.
of rooms location

o/p
sales (profit)
Age of fossil
Rainfall in current
Future stock prices
lung capacity
house price

Types of Regression

Linear Regression :-

- (1) simple linear Regression (Univariate)
- (2) Multiple linear Regression (Multivariate)

Polynomial Regression

(1) Simple linear Regression

Data:- pair of variables (one i/p, one o/p)

$$D = \{ (x_i, y_i) / i = 1 \text{ to } n \}$$

x_i = Today's temperature

y_i = Tomorrow's Rainfall

x_i is independent Variable

y_i is dependent Variable.

No of cigarettes (x_i)

lung capacity (y_i)

0

45

5

42

10

33

15

31

20

29

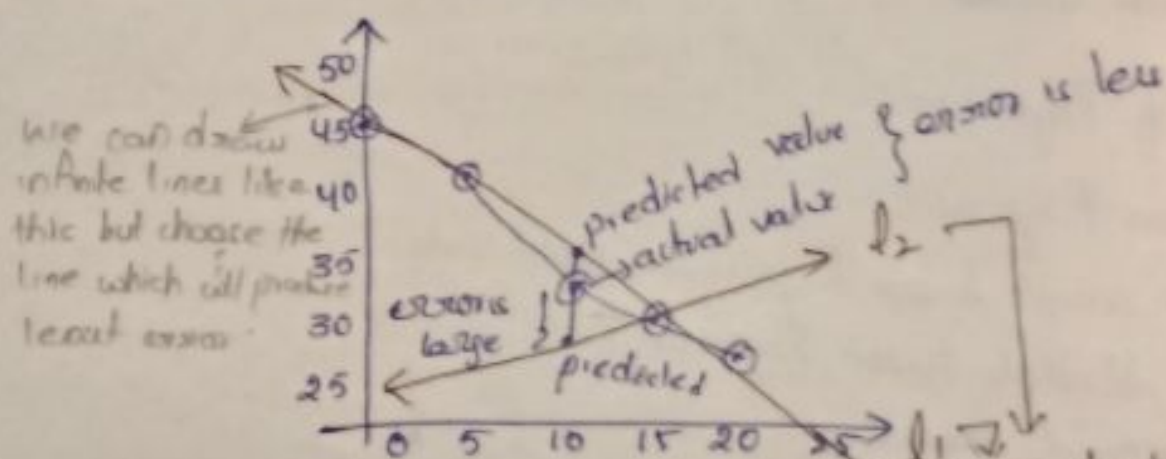
25

?

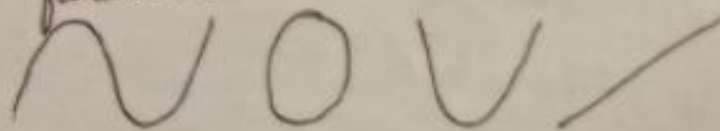
Task:

predict real valued y for given real valued x using a regression model f .

Model Structure :- $f(x) = w_0 + w_1 x$



we can use line plot (or) scatter plot function :-



some function is slightly varied from original.

if line slope $\Rightarrow y = mx + c$

m = slope

c = intercept.

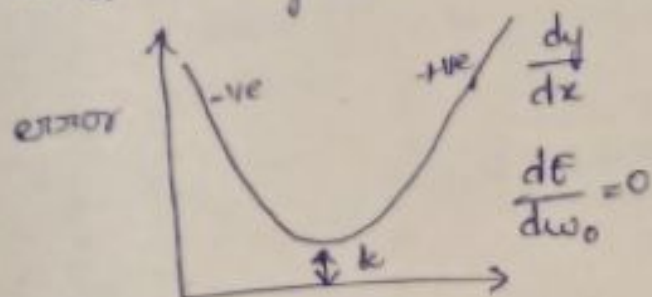
if m or c change line may change but x is changing it will not go beyond line.

Model Parameter :- $\Theta = \{w_0, w_1\}$.

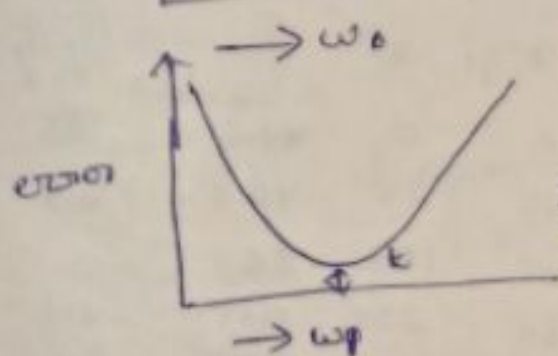
Error function :- Sum Absolute Error = $\sum_{i=1}^n |y_i - f(x_i)|$

$$\text{Sum Squared Error} = \sum_{i=1}^n (y_i - f(x_i))^2$$

Training:- given a training set find the value of regression parameters such that model ~~best~~ fits the for training data.



From there a choose the one which has least error (k).
Ex: w_1



$\frac{dE}{dw_1} = 0 \rightarrow$ select if zero.
or else select w_1 .

$$E = \text{Sum squared error} = \sum_{i=1}^n (y_i - f(x_i))^2 = \sum_{i=1}^n (y_i - w_0 - w_1 x_i)^2$$

[derivative]

$$\frac{\partial E}{\partial w_0} = -2 \sum_{i=1}^n (y_i - w_0 - w_1 x_i) = 0$$

$$\Rightarrow \sum y_i - \sum w_0 - w_1 \sum x_i = 0$$

$$\Rightarrow \sum w_0 = \sum y_i - w_1 \sum x_i$$

$$= n w_0 = \sum y_i - w_1 \sum x_i$$

$$\therefore w_0 = \frac{\sum y_i - w_1 \sum x_i}{n}$$

$$w_0 = \bar{y} - w_1 \bar{x}$$

$$\frac{\partial E}{\partial \omega_1} = -2 \sum_{i=1}^n (y_i - \omega_0 - \omega_1 x_i) \cdot x_i = 0$$

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

$$\Rightarrow \sum x_i y_i - \omega_0 \sum x_i - \omega_1 \sum x_i^2 = 0$$

$$\Rightarrow \sum x_i y_i - (\bar{y} - \omega_1 \bar{x}) n \bar{x} - \omega_1 \sum x_i^2 = 0$$

$$\Rightarrow \sum x_i y_i - n \bar{x} \bar{y} - n \bar{x}^2 \omega_1 - \omega_1 \sum x_i^2 = 0$$

$$\Rightarrow \sum x_i y_i - n \bar{x} \bar{y} - n \bar{x}^2 \omega_1 - \omega_1 \sum x_i^2 = 0$$

$$\boxed{\omega_1 = \frac{\sum x_i y_i - n \bar{x} \bar{y}}{\sum x_i^2 - n \bar{x}^2}}$$

25/1/19

No of cigars	x_i	y_i Wing capacity	$x_i y_i$	x_i^2	\hat{y}_i $= \omega_0 + \omega_1 x$	$(\hat{y}_i - \bar{y})$	$(\hat{y}_i - \bar{y})^2$
0	0	45	0	0	44.6	8.6	73.96
5	5	42	210	25	40.3	4.3	18.49
10	10	33	330	100	36	0	0
15	15	31	465	225	31.7	-4.3	18.49
20	20	29	580	400	27.4	-8.6	73.96
			1585	750			

$$\bar{x} \Rightarrow \frac{0+5+10+15+20}{5}$$

$$\Rightarrow \frac{50}{5} = 10$$

$$\bar{y} = \frac{45+42+33+31+29}{5}$$

$$= \frac{180}{5}$$

$$= 36$$

$$\frac{\partial E}{\partial \omega_1} = -2 \sum_{i=1}^n (y_i - \omega_0 - \omega_1 x_i) x_i = 0$$

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

$$\Rightarrow \sum x_i y_i - \omega_0 \sum x_i - \omega_1 \sum x_i^2 = 0$$

$$\Rightarrow \sum x_i y_i - (\bar{y} - \omega_1 \bar{x}) n \bar{x} - \omega_1 \sum x_i^2 = 0$$

$$\Rightarrow \sum x_i y_i - n \bar{x} \bar{y} - n \bar{x}^2 \omega_1 - \omega_1 \sum x_i^2 = 0$$

$$\Rightarrow \sum x_i y_i - n \bar{x} \bar{y} - n \bar{x}^2 \omega_1 - \omega_1 \sum x_i^2 = 0$$

$$\boxed{\omega_1 = \frac{\sum x_i y_i - n \bar{x} \bar{y}}{\sum x_i^2 - n \bar{x}^2}}$$

26/1/19

No of cigars	x_i	y_i Wing capacity	$x_i y_i$	x_i^2	\hat{y}_i $= \omega_0 + \omega_1 x_i$	$(\hat{y}_i - \bar{y})$	$(\hat{y}_i - \bar{y})^2$
0	0	45	0	0	44.6	8.6	73.96
5	5	42	210	25	40.3	4.3	18.49
10	10	33	330	100	36	0	0
15	15	31	465	225	31.7	-4.3	18.49
20	20	29	580	400	27.4	-8.6	73.96
			1585	750			

$$\bar{x} \Rightarrow \frac{0+5+10+15+20}{5}$$

$$\Rightarrow \frac{50}{5} = 10$$

$$\bar{y} = \frac{45+42+33+31+29}{5}$$

$$= \frac{180}{5}$$

$$= 36$$

$$w_0 = 36 - w_1(10)$$

$$w_1 = \frac{1585 - 1800}{750 - 500}$$

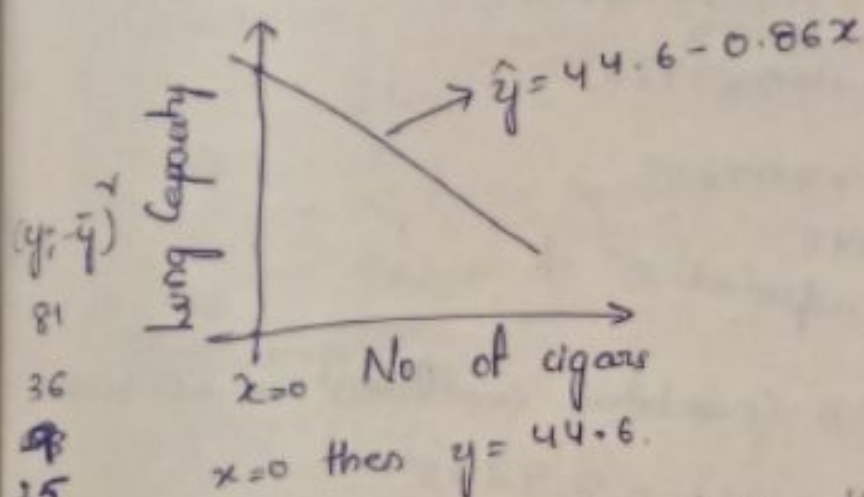
$$\Rightarrow \frac{-215}{250}$$

$$= -0.86$$

$$w_0 = 36 - (-0.86)(10)$$

$$= 36 + 8.6$$

$$= 44.6$$



Now we have to look whether the model is good enough or not can be done by Coefficient of Determination (R^2)

R^2 measure % of variation in y explained by the model.

$$R^2 = \frac{\text{Explained Variation}}{\text{Total Variation}} = \frac{SSM}{SST}$$

$$= \frac{\sum_{i=1}^n (\hat{y}_i - \bar{y})^2}{\sum_{i=1}^n (y_i - \bar{y})^2} = \frac{184.9}{200} \Rightarrow 0.9245$$

92% variation is explained by the model.

100% variation will not be explained as there will be noise for all the time.

$$\text{Adjusted } R^2 = 1 - \left(\frac{N-1}{N-d} \right) (1-R^2)$$

N = no. of data points = 5

d = no of features = 1 [(x_i) input]

$$1 - \left(\frac{5-1}{5-1} \right) (1 - (0.9245)^2)$$

$$1 - (1 - 0.85470025)$$

$$1 - (0.14529975)$$

$$= 0.85470025$$

$$\approx 0.9245$$

features \uparrow , adjusted $R^2 \downarrow$

r = pearson correlation coefficient = -0.96151

$$r^2 = (-0.96151)^2 = 0.9245$$

This will work for 2 variables only.

11.9 Multiple linear Regression

i/p
TV Ad, Internet Ad, Newspaper Ad o/p
sales

Data $\{x_{11}, x_{12}, \dots, x_{1d}, y_1\}$

$\{x_{21}, x_{22}, \dots, x_{2d}, y_2\}$

$\{x_{n1}, x_{n2}, \dots, x_{nd}, y_n\}$

$D = \frac{1}{2} (x_{11}, x_{12}, \dots, x_{1d}, y_1^*) / i=1, 2, 3, \dots, n$
 Task:- predicting the values of y
 Model Structure

$$\hat{y} = f_y = f(x_1, x_2, \dots, x_d)$$

$$= w_0 + w_1 x_1 + w_2 x_2 + \dots + w_d x_d$$

↳ not polynomials (power 1), so Variables

$$= w_0 + \sum_{i=1}^d w_i x_i$$

Parameter:- $d+1$ parameters

$$\theta = \frac{1}{2} w_0, w_1, \dots, w_d$$

we have to find these parameters such that error function is minimum

Error function:-

$$E = SSE = \sum_{i=1}^n (y_i - f(x_{i1}, x_{i2}, \dots, x_{id}))^2$$

Training:-

$d+1$ parameters $\Rightarrow d+1$ equations

Matrix Notation is useful to find the Parameters w_0, w_1, \dots, w_d .

$$3x + 4y = 6$$

$$9x + 22y = 26$$

$$\begin{bmatrix} 3 & 4 \\ 9 & 22 \end{bmatrix} \begin{bmatrix} x \\ y \end{bmatrix} = \begin{bmatrix} 6 \\ 26 \end{bmatrix}$$

$$\begin{bmatrix} x \\ y \end{bmatrix} = \begin{bmatrix} 3 & 4 \\ 9 & 22 \end{bmatrix}^{-1} \begin{bmatrix} 6 \\ 26 \end{bmatrix}$$

Normal Equation (for Simple Linear Regression)

$$\frac{\partial E}{\partial \omega_0} = -2 \sum_{i=1}^n (y_i - \omega_0 - \omega_1 x_i) = 0$$

$$\frac{\partial E}{\partial \omega_1} = -2 \sum_{i=1}^n (y_i - \omega_0 - \omega_1 x_i) x_i = 0$$

$$X = \begin{bmatrix} 1 & x_1 \\ 1 & x_2 \\ \vdots & \vdots \\ 1 & x_n \end{bmatrix} \quad y = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{bmatrix}$$

$$\begin{bmatrix} n & \sum x_i \\ \sum x_i & \sum x_i^2 \end{bmatrix} \begin{bmatrix} \omega_0 \\ \omega_1 \end{bmatrix} = \begin{bmatrix} \sum y_i \\ \sum x_i y_i \end{bmatrix}$$

↓ Above equation and matrix format are equal.

$$\begin{bmatrix} \omega_0 \\ \omega_1 \end{bmatrix} = \begin{bmatrix} n & \sum x_i \\ \sum x_i & \sum x_i^2 \end{bmatrix}^{-1} \begin{bmatrix} \sum y_i \\ \sum x_i y_i \end{bmatrix}$$

$$X^T X = \begin{bmatrix} 1 & 1 & 1 & \dots & 1 \\ x_1 & x_2 & x_3 & \dots & x_n \end{bmatrix} \begin{bmatrix} 1 & x_1 \\ 1 & x_2 \\ 1 & x_3 \\ \vdots & \vdots \\ 1 & x_n \end{bmatrix}$$
$$= \begin{bmatrix} n & \sum x_i \\ \sum x_i & \sum x_i^2 \end{bmatrix}$$

$$X^T y = \begin{bmatrix} 1 & 1 & 1 & \dots & 1 \\ x_1 & x_2 & x_3 & \dots & x_n \end{bmatrix} \begin{bmatrix} y_1 \\ y_2 \\ y_3 \\ \vdots \\ y_n \end{bmatrix}$$
$$= \begin{bmatrix} \sum y_i \\ \sum x_i y_i \end{bmatrix}$$

$$\boxed{W = [X^T X]^{-1} [X^T y]}$$

Normal Equation (for Multiple Linear Regression)

$$X = \begin{bmatrix} 1 & x_{11} & x_{12} & \dots & x_{1d} \\ 1 & x_{21} & x_{22} & \dots & x_{2d} \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 1 & x_{n1} & x_{n2} & \dots & x_{nd} \end{bmatrix} \quad y = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{bmatrix}$$

$n \times (d+1)$ $n \times 1$

$$W = \begin{bmatrix} w_0 \\ w_1 \\ \vdots \\ w_d \end{bmatrix}$$

$(d+1) \times 1$

$$W = (X^T X)^{-1} (X^T y) \Rightarrow \underline{(d+1) \times (d+1)}$$

square matrix and symmetric

Every square symmetric matrix is invertible

Limitations of Linear Regression

The relationship

* True Relationship of x and y might be non linear

* complexity of the Algorithm $O(Nd^2 + d^3)$

N = no. of datapoints.

d = no. of features.

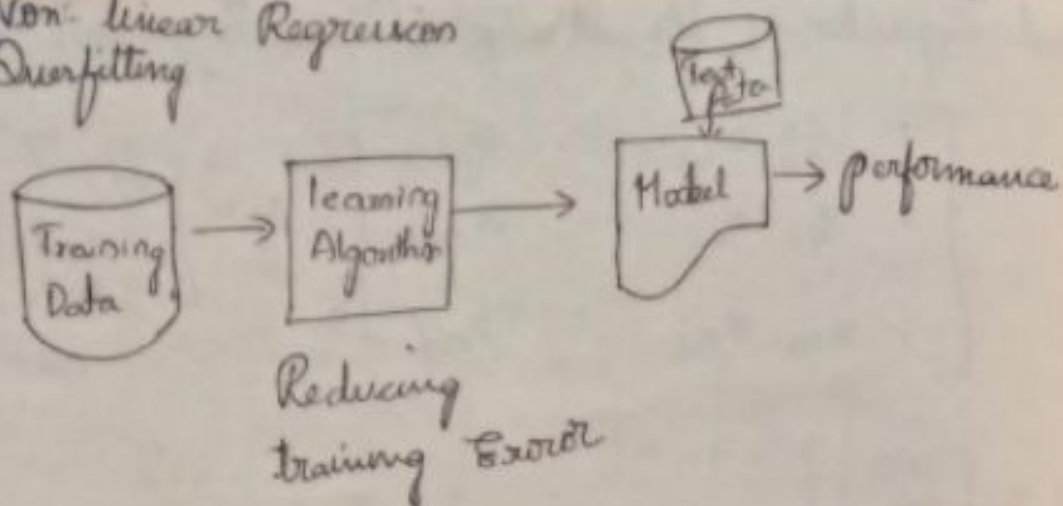
* It includes all variables in the model. but all variables may not be related to y .

* This approach may fail if two variables are correlated namely multicollinearity

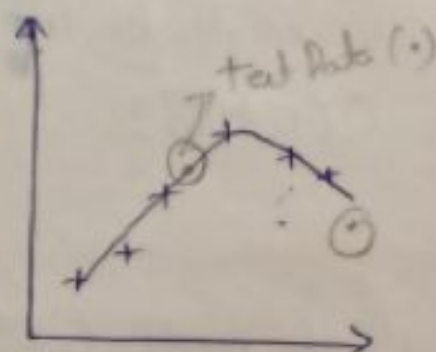
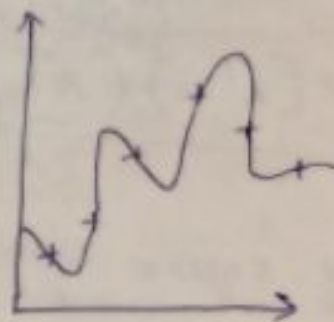
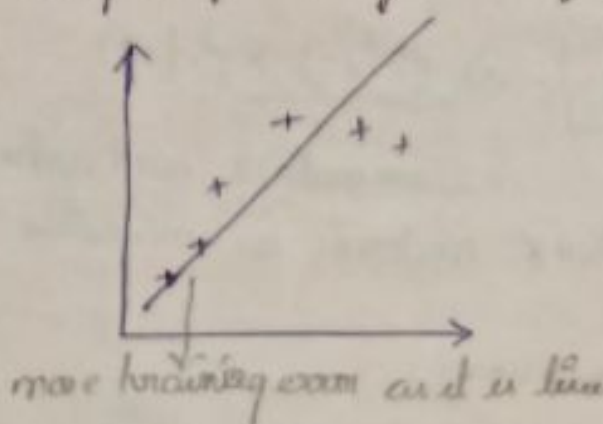
* correlation / collinearity among x variables can cause numerical instability.

Non-linear Regression

30/11/20 Overfitting

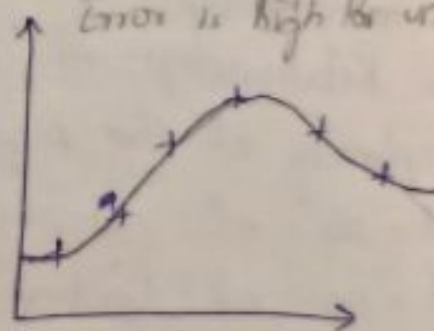


Complexity vs Goodness of fit



$$y = w_0 + w_1 x + w_2 x^2$$

- + some training error
- + works well for unknown data



$$y = w_0 + w_1 x + w_2 x^2$$

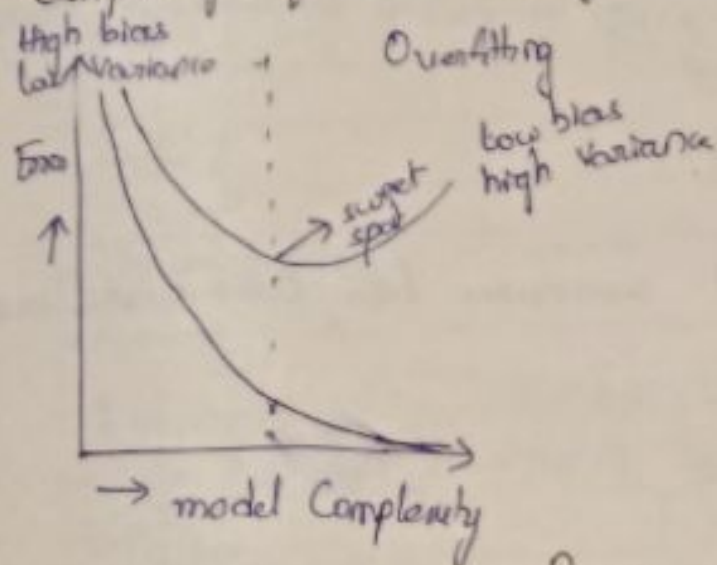
Occam's Razor Principle:-

simpler solutions are more likely to be correct than the complex ones. ~~use the~~

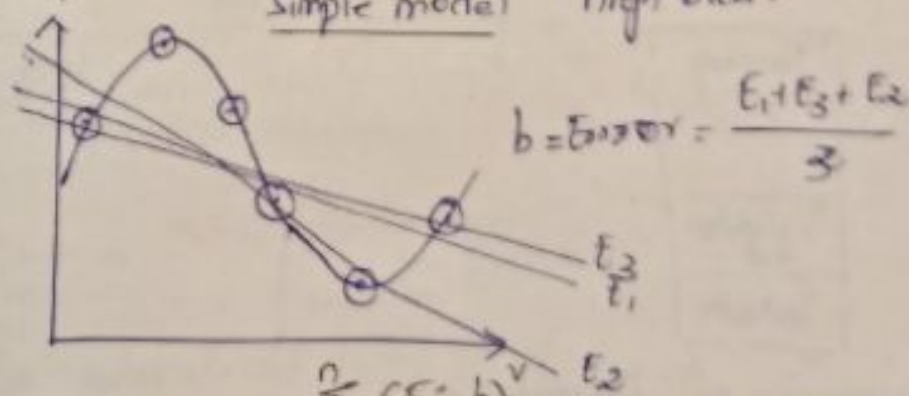
Use the simplest model which gives acceptable accuracy on training set

5/1/19

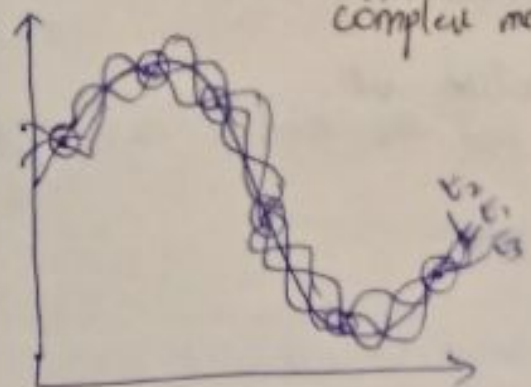
Complexity of model vs generalization



① Bias of a model = $\frac{\sum_{i=1}^n E_i^2}{n-1}$ = b Low Variance high bias.
Simple model



② Variance model = $\frac{\sum_{i=1}^n (E_i - b)^2}{n-1}$ complex model



$$b = \frac{E_1 + E_2 + E_3}{3}$$

$$V = \frac{\sum_{i=1}^n (E_i - b)^2}{n-1}$$

Low bias high Variance

Bias :-

Bias of a model is the average error it makes on various training sets

$$B = \left(\sum_{i=1}^n E_i \right) / n$$

Variance :-

Variation of a model error it makes on different training sets

$$V = \frac{\sum_{i=1}^n (E_i - b)^2}{n}$$

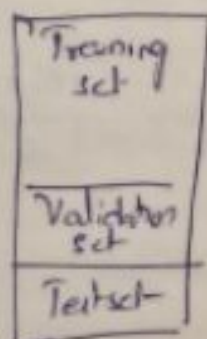
We need to find a

Find a model that minimizes both bias and Variance

Regularization :- $\begin{cases} L_1 \\ L_2 \end{cases}$

Model Validation :-

① Holdout method :-



Purpose of Training set

use this data to train each model

Purpose of Validation set

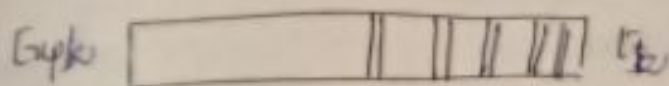
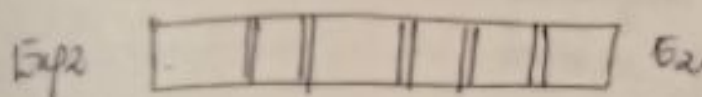
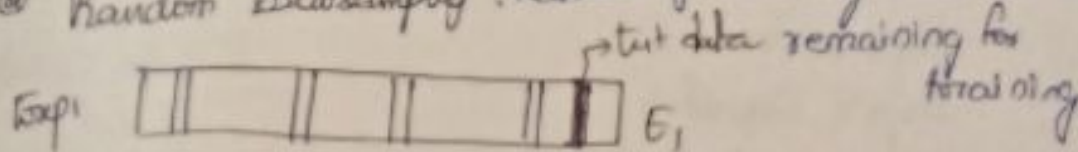
use this data to test the models and select the model with minimum error.

Purpose of Test set

use this data to calculate unbiased error for the selected model.

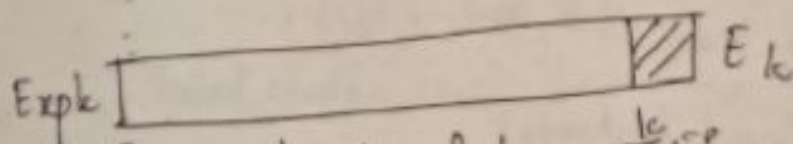
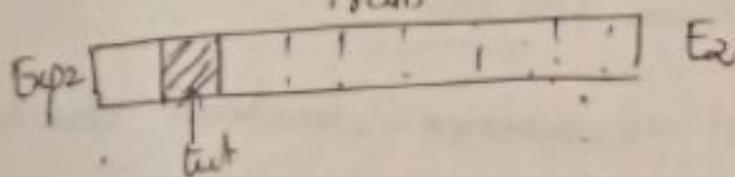
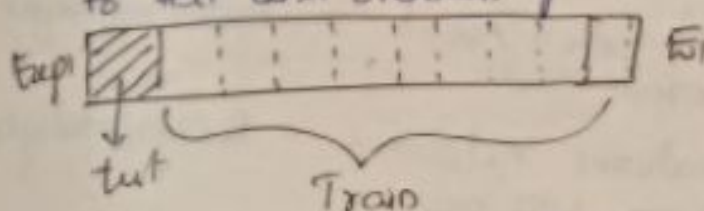
③ Cross Validation:

① Random Subsampling: Randomly selecting a test set and



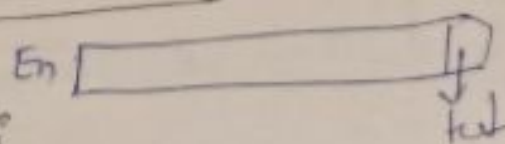
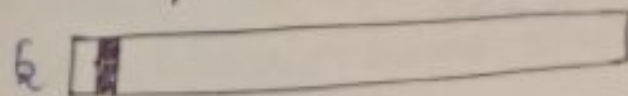
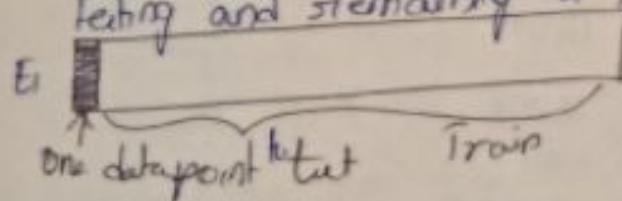
True estimate of error $E = \frac{\sum_{i=1}^k E_i}{k}$

② k-fold Cross validation: use a fold of dataset to test and remaining for training.



True estimate of $k = \frac{\sum_{i=1}^k E_i}{k}$

③ Leave one out: Each time selecting a datapoint for testing and remaining for training



$E_o = \frac{\sum_{i=1}^n E_i}{n}$

11/2/19

Classification - supervised learning

- Decision Trees
- Baye's Classifier
- k Nearest Neighbours
- support vector machine
- ANN
- Ensemble method

Task	i/p	o/p
① Categorizing Email	Features Extracted from email headers and Content	category ↓ o/p spam/not spam
② Loan Approval	Information About person	Approve/Reject
③ Identifying Tumour cell	Features extracted from MRI scan	Benign/Malignant

Decision Tree:- It is a tree structured classifier.

Data:- $\{ \underbrace{x_{11}, x_{12}, \dots, x_{1d_1}}_{\text{input features}}, \underbrace{y_1}_{\text{class label}} \}$

ID	Home owner	Hospital status	Annual Income	Defaulted
1	Y	M	125k	N
2	N	N	220k	Y
	↑ cat	↑ cat	↑ number	

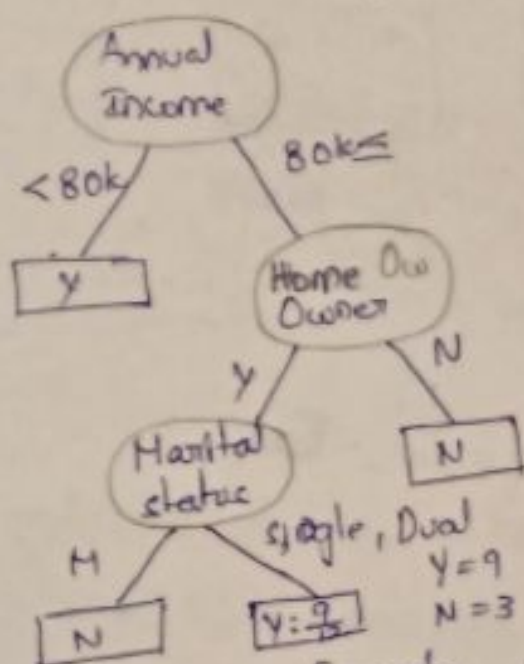
Task: The task is to learn a model that maps each attribute set x into one of the predefined class label y

N	S	180k	?
---	---	------	---

Model structure :- - A tree

- Two types of nodes
 - Decision Node - Internal Node - Attribute
 - Consequence Node - Leaf Node - class

Eg:-



* given some training Examples we have to generate Decision Tree

• many Decision Trees are possible

• we prefer a smaller Tree

means less no. of nodes and low depth

* finding the smallest Tree that fits the data is computationally hard problem.

* Greedy Algorithm - will be used (local optimization)

Parameters of model

- which Attribute to select for split?
- when to stop?

② All attributes belongs to the ^{same} class

- ① All Attributes are already used (no more features)
- ③ Too few examples to make informative split

Error Function: Measure of Node impurity.

+	9
-	1

more pure

+	: 5
-	: 5

less pure.

① Gini index :-

② Entropy

③ Misclassification Error

above 3 are used for calculating informative splits

Training will done by below Algorithms.

- Hunt's Algorithm
- CART
- ID3, C4.5
- SLIQ, SPRINT

Basic outline of ID3 Algorithm: Top down Approach

$A = \{A_1, A_2, A_3, \dots, A_n\}$

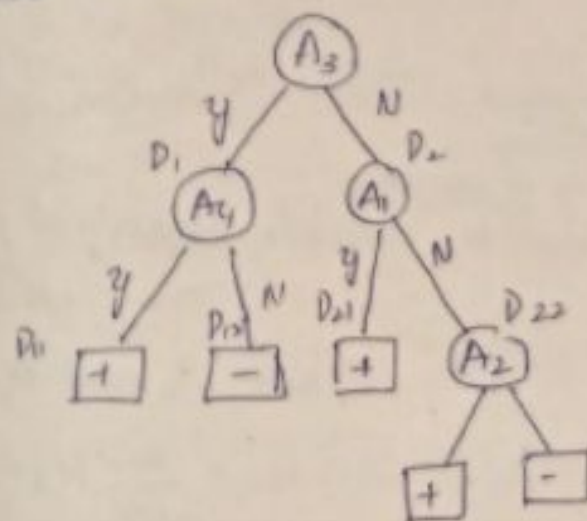
① select best Attribute among all called for Next node

② Assign the selected Attribute as Decision Attribute for the Node

③ For each value of the Attribute create new descendant

④ sort the training Examples ~~into~~ to leaf nodes according attribute value

⑤ If all training Examples are perfectly classified then stop else iterate over new leaf Nodes.



4/2/19

which attribute to select?

- which attribute gives smallest errors

- greedy Approach.

- Measure of impurity

- Entropy

- Gini

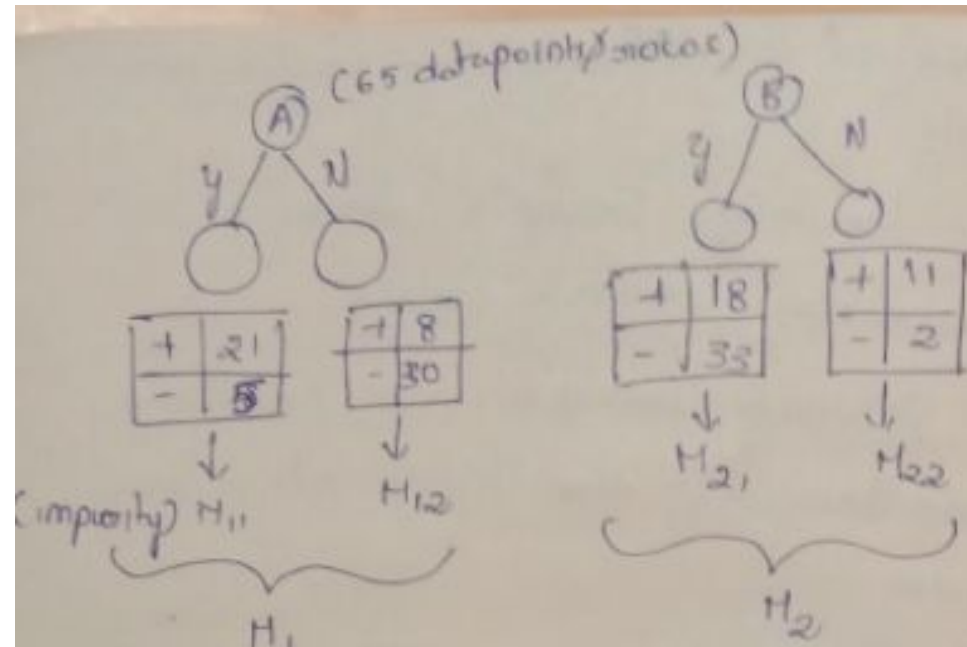
- Misclassification Error

} Information gain

+	29
-	35

Attribute = {A, B}

P = impurity



Gain = $P - H_1$ (impurity of children) \rightarrow impurity of parent

Gain = $P - H_2$

We don't know how to calculate impurity of parent (children)

" " " " " " " " combine H_{11} and H_{12}

combining H_{11} and $H_{12} \therefore H_1 = \frac{26}{64} H_{11} + \frac{38}{64} H_{12}$

① finding Best Attribute :- $H_2 = \frac{n_1}{n} H_{21} + \frac{n_2}{n} H_{22}$

- 1 Compute the impurity measure p before splitting
- 2 compute the impurity measure after splitting
- step-A: compute the impurity measure for each ^{child} node

Measure of Impurity

Entropy :-

Expectation : It is the fancy name for Average

Team (x) : A B C D

Probability of win: 0.3 0.4 0.2 0.1

Reward (v(x)) : 10k 5k 10k -30k

What is the Expected Reward?

$$\text{Expected Reward} = \frac{3 \times 10 + 4 \times 5 + 2 \times 10 + 1 \times (-30)}{10}$$

$$= \frac{3}{10} \times 10 + \frac{4}{10} \times 5 + \frac{2}{10} \times 10 + \frac{1}{10} \times (-30)$$

$$= 0.3 \times 10 + 0.4 \times 5 + 0.2 \times 10 + 0.1 \times (-30)$$

For any Random Variable :- $E(x) = \sum p(x) \cdot v(x)$

$$= \sum \text{probability of } (v(x)) \times \text{value of } x$$

Information content :-

$$I(A) \propto \frac{1}{P(A)}$$

$$I(A) = f(P(A)) = \log\left(\frac{1}{P(A)}\right) = -\log(P(A))$$

$$I(A) = -\log(P(A))$$

Information content of Independent Events:

X = which team is going to win

Y = lunch is going to be good in mess

$$I(X \cap Y) = I(X) + I(Y)$$

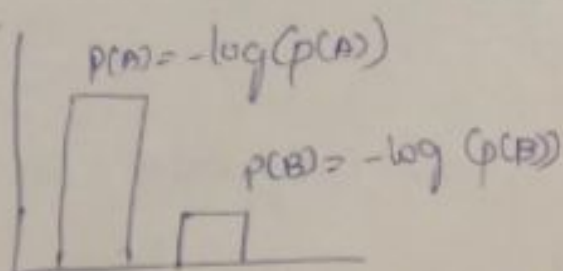
$$f(P(X \cap Y)) = f(P(X)) + f(P(Y))$$

$$f(P(X)) + f(P(Y)) = f(P(X)) + f(P(Y))$$

$$\log(P(X) \cdot P(Y)) = \log(P(X)) + \log(P(Y))$$

X, Y are independent

Entropy:



Expected $IC = p(A) I(A) + p(B) I(B)$
 $= -p(A) \log_2(p(A)) - p(B) \log_2(p(B))$
 $= \text{Entropy}$

+	29	29/64
-	35	35/64

calculable Entropy

$$\text{Entropy} = -\frac{29}{64} \log_2\left(\frac{29}{64}\right) - \frac{35}{64} \log_2\left(\frac{35}{64}\right)$$

$$= -\frac{29}{64} \log_2\left(\frac{29}{64}\right) - \frac{35}{64} \log_2\left(\frac{35}{64}\right)$$

$$\Rightarrow \frac{+0.155776209}{-0.143342461}$$

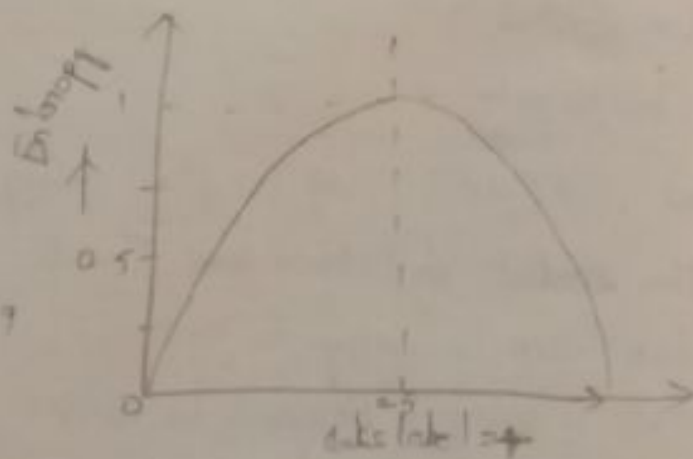
+	10
-	0

+	0
-	10

+	5
-	5

+	9
-	1

0.469

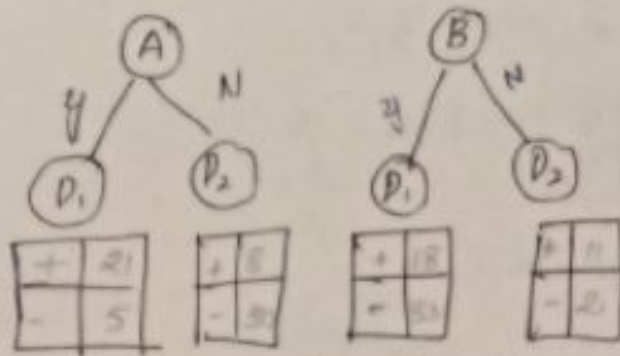


$$E(0) = 0.99$$

5/2/19

$$D = \begin{bmatrix} + & 29 \\ - & 35 \end{bmatrix}$$

$$E(D) = 0.99$$



$$H_{11} = E(D_1) = 0.71$$

$$H_{12} = E(D_2) = 0.74$$

$$H_{21} = E(D_1) = 0.94$$

$$H_{22} = E(D_2) = 0.62$$

play or not?

Day	outlook	Temp	Humidity	wind	play
1	sunny	Hot	High	weak	no
2	sunny	Hot	High	strong	no
3	overcast	Hot	High	w	yes
4	Rainy	Mild	High	w	yes
5	Rainy	cool	Normal	w	yes
6	Rainy	cool	Normal	s	no
7	overcast	cool	Normal	s	yes
8	sunny	Mild	High	w	no
9	sunny	cool	Normal	w	yes
10	Rainy	mild	Normal	w	yes
11	sunny	mild	Normal	s	yes
12	overcast	mild	High	s	yes
13	overcast	Hot	Normal	w	yes
14	Rainy	Mild	High	s	No

$$\begin{aligned} (A) \text{ gain} &= 0.99 - \left(\frac{26}{64} \times 0.71 + \frac{38}{64} \times 0.74 \right) \\ &= 0.99 - 0.73 \\ &= 0.26 \end{aligned}$$

$$\begin{aligned} (B) \text{ gain} &= 0.99 - \left(\frac{51}{64} \times 0.94 + \frac{13}{64} \times 0.62 \right) \\ &= 0.99 - 0.62 \end{aligned}$$

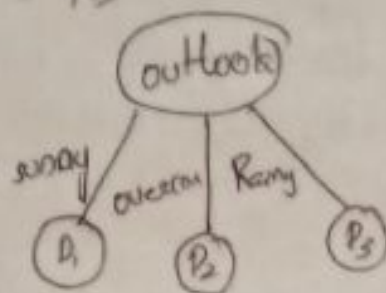
$$\text{Reduction in impurity} = 0.37, 0.115$$

select A for split since high gain

Entropy of root node

$D = \begin{array}{|c|c|} \hline \text{yes} & 9 \\ \hline \text{no} & 5 \\ \hline \end{array} \quad E(D) = 0.94$

Gain (outlook, D)



class	count		
	S	O	R
yes	2	4	3
no	3	0	2

$$E(D_1) = 0.97 \Rightarrow -\frac{2}{5} \log_2\left(\frac{2}{5}\right) - \frac{3}{5} \log_2\left(\frac{3}{5}\right)$$

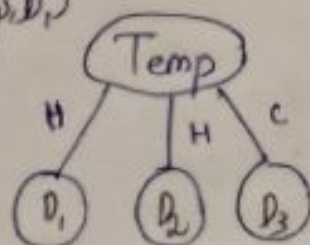
$$E(D_2) = 0$$

$$E(D_3) = 0.97$$

$$\begin{aligned} \text{gain}(\text{outlook}, D) &= 0.94 - \left(0.97 \times \frac{5}{14} + 0 \times \frac{4}{14} + 0.97 \times \frac{5}{14}\right) \\ &\Rightarrow 0.25 \end{aligned}$$

Gain (Temp, D)

(-)-(-)



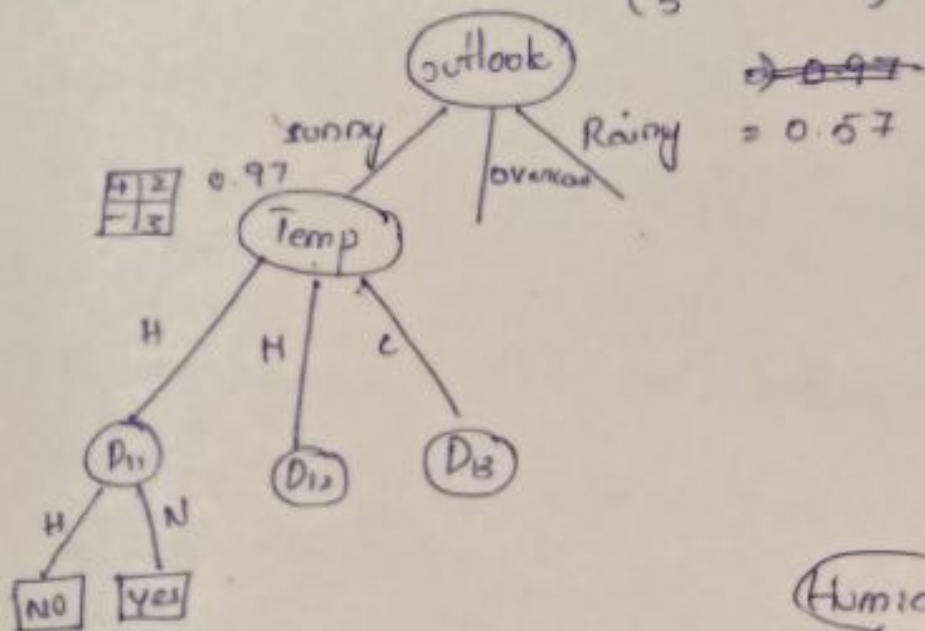
class	count		
	Hot	mild	cool
yes	2	4	3
no	2	2	1

$$E(D_1) = 0.1$$

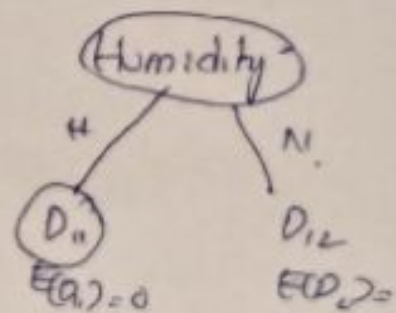
$$E(D_3) = 0$$

$$E(D_2) = 1$$

$$\text{Gain}(\text{Temp}, D_1) = 0.97 - \left(\frac{2}{5} \times 0 + \frac{2}{5} \times 1 + \frac{1}{5} \times \right.$$



$$\text{Gain}(\text{Humidity}, D_1) = 0.97 - 0 = 0.97$$



$$\begin{aligned} \text{Gain}(\text{wind}, D_1) &= 0.97 - \left(\frac{3}{5} \times 0.918 + \frac{2}{5} \times 1 \right) \\ &= 0.02 \end{aligned}$$

Gini index

$$Gini(t) = 1 - \sum_{i=0}^{c-1} [p(i/t)]^2$$

+	5	$\rightarrow p = \frac{1}{2}$
-	5	$\rightarrow p = \frac{1}{2}$

$$1 - \left[\left(\frac{1}{2}\right)^2 + \left(\frac{1}{2}\right)^2\right] = 0.5$$

+	10	1
-	0	0

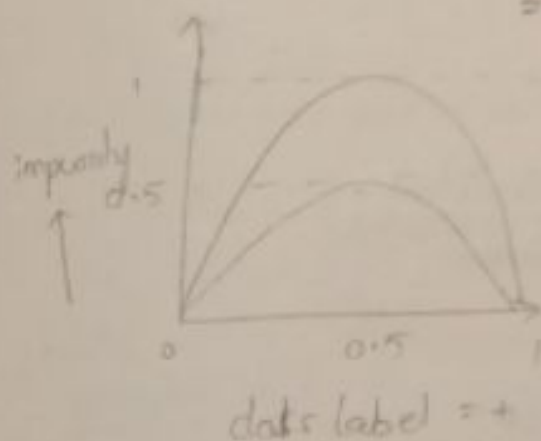
$$Gini(t) = 1 - [(1)^2 + (0)^2] = 0$$

+	10	0
-	10	1

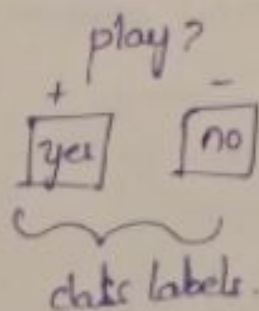
$$Gini(t) = 1 - [0^2 + (1)^2] = 0$$

+	8	0.8
-	2	0.2

$$\begin{aligned} Gini(t) &= 1 - [(0.8)^2 + (0.2)^2] \\ &= 1 - [0.64 + 0.04] \\ &= 0.32 \end{aligned}$$

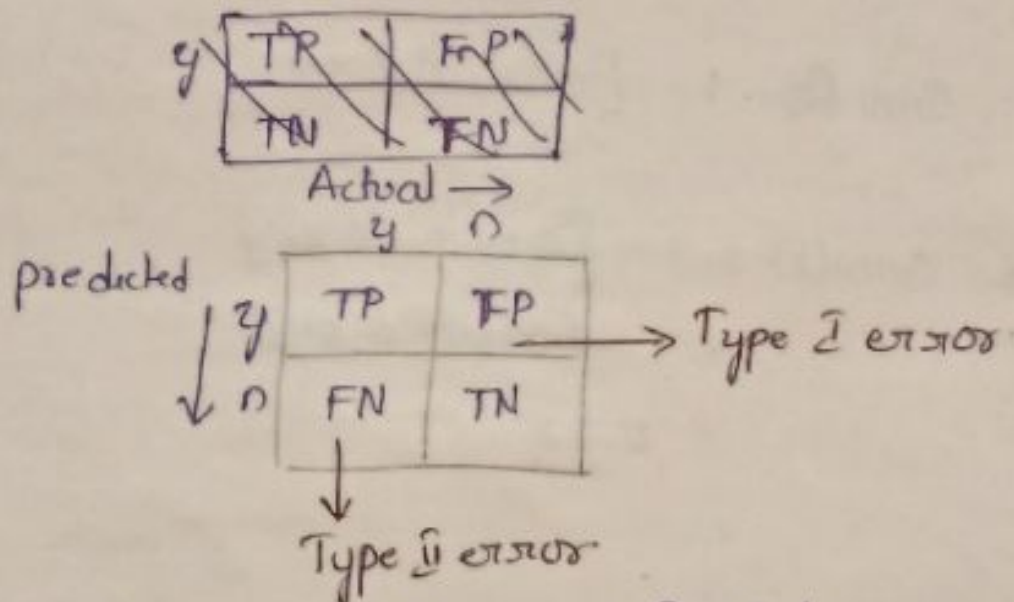


Metric to evaluate model performance (classification) :-



True label	predicted label	Result Type
yes	yes	True +ve (TP)
yes	no	True -ve (TN)
no	yes	False +ve (FP)
no	no	False -ve (FN)

Confusion Matrix :-



$$\begin{aligned}
 \textcircled{1} \text{ Accuracy} &= \frac{\# \text{ correct Prediction}}{\# \text{ total prediction}} \\
 &= \frac{TP + TN}{(TP + TN + FN + FP)}
 \end{aligned}$$

$$\begin{aligned}
 \textcircled{2} \text{ Error Rate} &= \frac{\# \text{ Error inconsistent prediction}}{\# \text{ total prediction}} \\
 &= \frac{FP + FN}{TP + TN + FN + FP}
 \end{aligned}$$

Fig:-

True	predicted
y	N
N	N
N	N
y	y
y	y
n	n
y	n
y	y
n	n
N	y

Actual

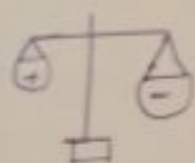
	y	N
predicted y	3	1
N	2	4

$$\text{Accuracy} = \frac{3+4}{10} = 0.7$$

$$\text{Error Rate} = \frac{2+1}{10} = 0.3$$

Limitation with Accuracy:-

- class imbalance problem .



- very few + examples

- most are -ve Examples

Eg: 3% are cancer, 97% are not cancer

	y	n
y	0	0
n	3	97

12/2/19

Precision And Recall

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

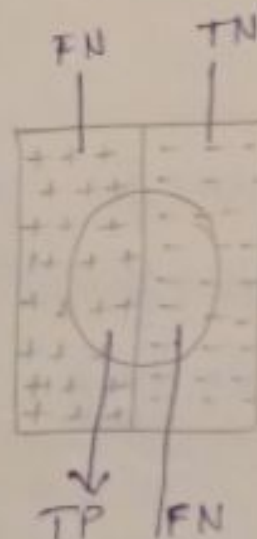
$$\text{Recall} = \frac{TP}{TP + FN}$$

$$\text{Precision} = \frac{1}{\text{Venn Diagram}}$$

→ measure of exactness

$$\text{Recall} = \frac{1}{\text{Rectangle}}$$

→ measure of completeness



F-Measure:- combines precision & Recall

$$F_1 = 2 * \frac{\text{Precision} \times \text{Recall}}{\text{Precision} + \text{Recall}}$$

$$0 \leq F_1 \leq 1$$

1 = perfect classification

F_1 = evenly Weighted

F_2 = Weights recall more

F_3 = Weights precision more

0	0
3	97

Accuracy = 0.97 precision = 0, Recall = 0
 $F_1 = 0$

		Actual		
		cat	Dog	Rabbit
Predicted	cat	5	2	0
	Dog	3	3	2
	Rabbit	0	1	11

cat?

	+	-
+	5	2
-	3	

Dog?

	+	-
+	3	5
-	3	

Rabbit?

	+	-
+	1	1
-	2	

$$\text{Precision} = \frac{5}{7}$$

$$\text{Recall} = \frac{5}{8}$$

$$F_1 = 2 \times \frac{\frac{5}{7} \times \frac{5}{8}}{\frac{5}{7} + \frac{5}{8}}$$

$$\Rightarrow 0.67$$

$$\text{precision} = \frac{3}{8}$$

$$\text{Recall} = \frac{3}{6}$$

$$F_1 = 2 \times \frac{\frac{3}{8} \times \frac{3}{6}}{\frac{3}{8} + \frac{3}{6}}$$

$$= 0.42$$

$$\text{Precision} = \frac{11}{12}$$

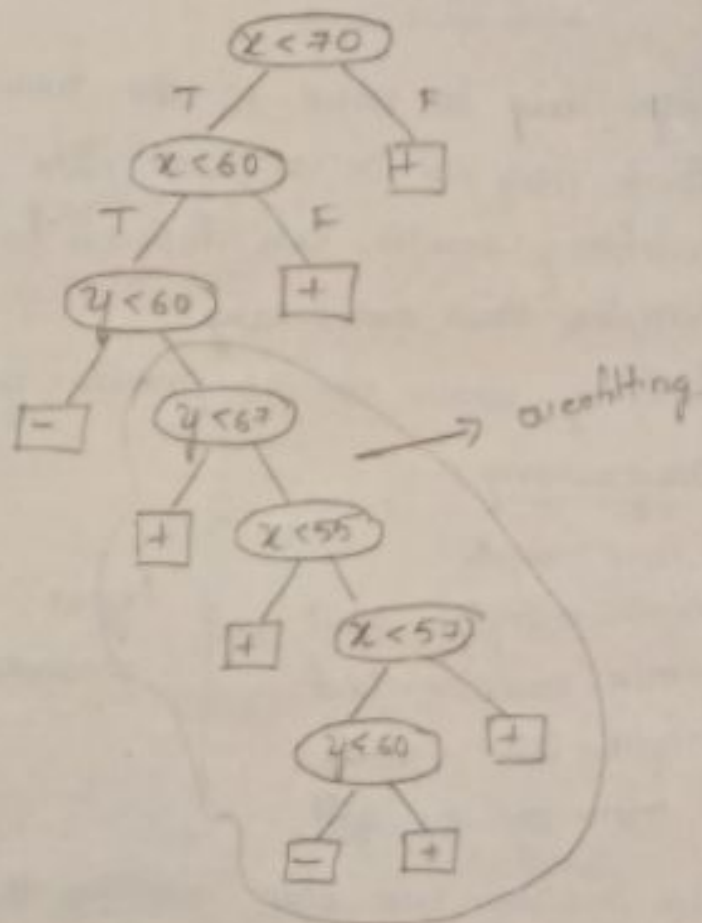
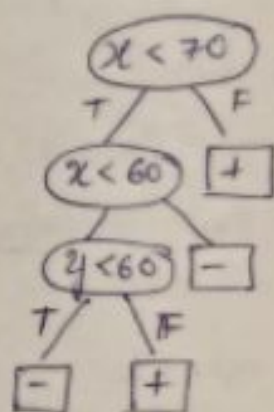
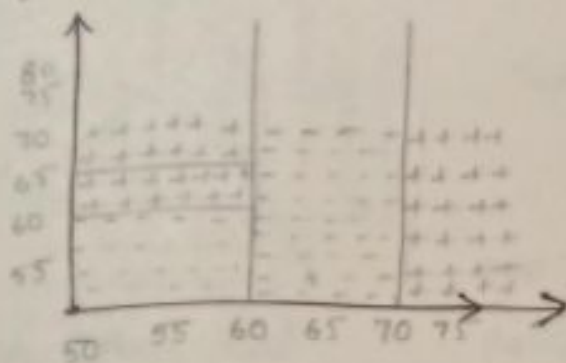
$$\text{Recall} = \frac{11}{13}$$

$$F_1 = 2 \times \frac{\frac{11}{12} \times \frac{11}{13}}{\frac{11}{12} + \frac{11}{13}}$$

$$= 0.88$$

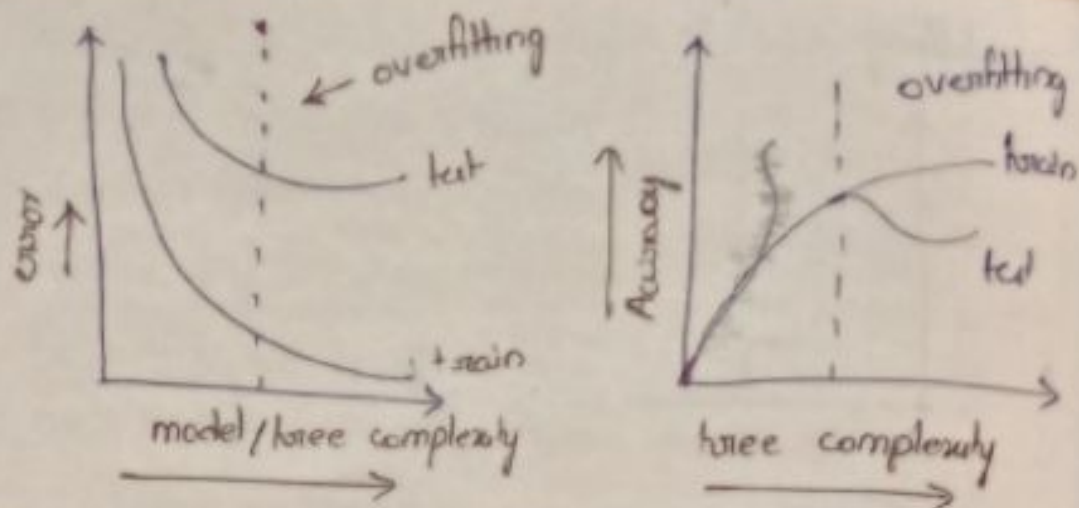
$$F = \frac{0.67 + 0.42 + 0.88}{3} = 0.65$$

overfitting



11/12/19

A model m is said to overfit the training data if there is another model m' , such that m' has more errors than m on training data but m' has less errors than m on test data.



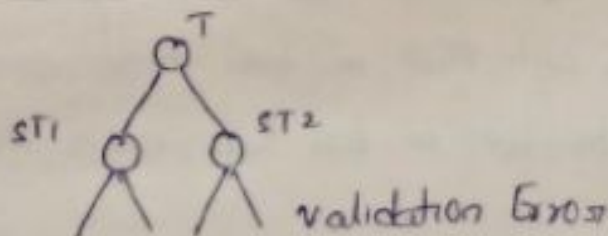
- There may be noise in the training data
- There may not be sufficient data to take decision
- overfitting results in a decision tree that is more complex than necessary
- Training errors do not provide good estimate

Pruning :-

- max-depth
 - min-sample-split
 - min-sample-leaf
 - max-leaf-nodes
 - min-impurity-split
- } hyper parameters.

In pruning, we stop growing the tree

Cross Validation:-



$E(T)$

$E(T - ST1), E(T - ST2)$

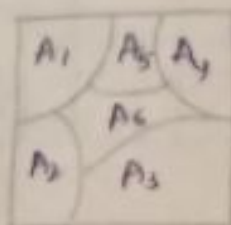
select the node whose removal gives least validation error

15th Feb 19

Baye's Classifier

probability crash course

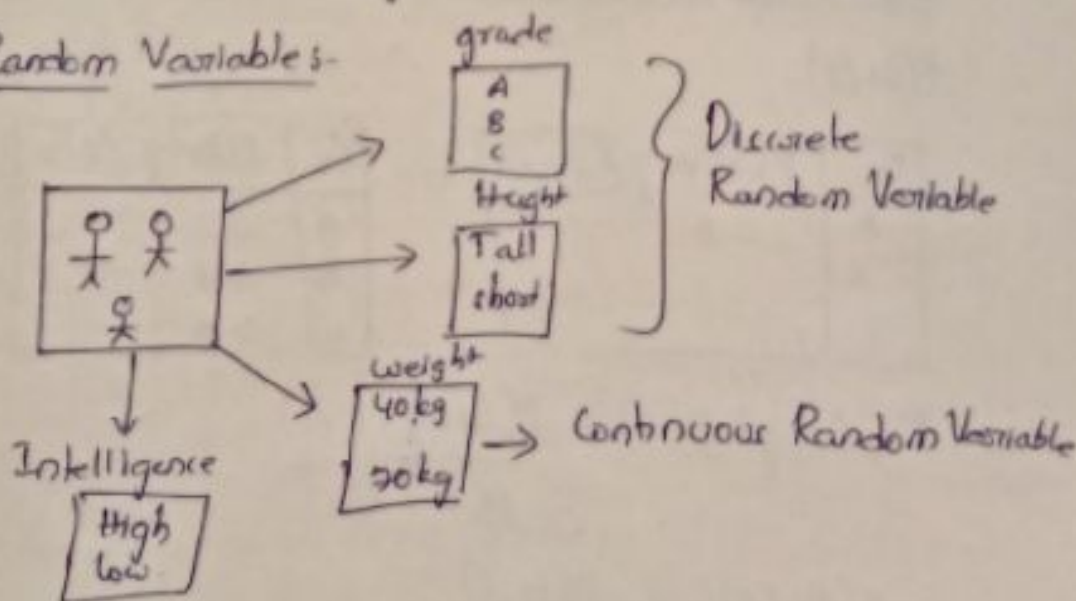
Axioms:-



- For any event A , $p(A) \geq 0$
- $p(\Omega) = \sum_i p(A_i) = p(\Omega)$

Disjoint Events:- $A_i \cap A_j = \emptyset$
 $\forall i \neq j$

Random Variables:-



Marginal Distribution:

A	40
B	30
C	30

G_i	$p(G_i = g)$
A	0.4
B	0.3
C	0.3

$$p(G_i = g) \forall g \in A, B, C$$

$$p(\Omega)$$

joint Distribution

G and I

G	I	$P(G=g, I=i)$
A	H	0.3
A	L	0.1
B	H	0.15
B	L	0.15
C	H	0.05
C	L	0.25

	A	B	C	
High	30	15	5	
low	10	15	25	
				100

$$P(G=g, I=i) \quad \forall (g,i) \in \frac{1}{2} A, B, C \times \frac{1}{2} H, L$$

$$P(G, I)$$

Conditional Distribution

$P(G/I)$

G	$P(G=g / I=H)$
A	0.6
B	
C	

G	$P(G=g / I=L)$
A	0.2
B	0.3
C	0.5

$$P(G=g / I=i) = \frac{P(G=g, I=i)}{P(I=i)}$$

$$P(G/I) = \frac{P(G, I)}{P(I)}$$

$$P(G, I) = P(G/I) \times P(I)$$

\downarrow joint Distribution \rightarrow conditional \rightarrow marginal

Joint Distribution of Random Variable

x_1	x_2	x_3	...	x_n	$p(x_1, x_2, x_3, x_4 \dots x_n)$

2^n : Assuming each random variable take 2 values

$$\begin{aligned}
 p(x_1, x_2, \dots, x_n) &= p(x_2, x_3, x_4, \dots, x_n | x_1) \cdot p(x_1) \\
 &= p(x_3, x_4, \dots, x_n | x_1, x_2) \cdot p(x_2 | x_1) \cdot p(x_1) \\
 &= p(x_4, \dots, x_n | x_1, x_2, x_3) \cdot p(x_3 | x_1, x_2) \cdot p(x_2 | x_1) \cdot p(x_1)
 \end{aligned}$$

General Equation:-

$$p(x_i^*)_{i=2}^n = \prod_{i=2}^n p(x_i^* | x_1, \dots, x_{i-1}^*)$$

Calculating Marginal distribution from joint Distribution

G	I	$p(G, I)$	
A	H	0.3	} $\rightarrow 0.4$
A	L	0.1	
B	H	0.15	} 0.3
B	L	0.15	
C	H	0.05	} 0.3
C	L	0.05	

$$P(G=g) = \sum_{I=i} p(G=g, I=i)$$

$$P(G=g) = \sum_{I=i} p(G=g, I=i) \quad P(I=i) = \sum_{G=g} p(G=g, I=i)$$

$$P(G) = \sum_I p(G, I) \quad P(I) = \sum_G p(G, I)$$

n random Variables :-

$$p(x_1) = \sum_{x_2, x_3, \dots, x_n} p(x_1, x_2, \dots, x_n)$$

x_1	x_2	x_3	x_4	y
outlook	Temp	Humidly	wind	play
3	2	2	2	2

$$p(y, x_1, x_2, x_3, x_4)$$

$$p(y/x_1, x_2, x_3, x_4) = \frac{p(y, x_1, x_2, x_3, x_4)}{\sum_{x_1, x_2, x_3, x_4} p(y, x_1, x_2, x_3, x_4)}$$

Feb 14

Baye's classification :-

$p(x), p(y) \rightarrow$ marginal

$p(x, y) \rightarrow$ joint

$p(y/x) \rightarrow$ conditional

$$p(x, y) = p(y/x) p(x)$$

$$\frac{p(x, y)}{p(x)} = p(y/x)$$

$$p(x, y) = p(y/x) \cdot p(x)$$

$$\Rightarrow \frac{p(x, y)}{p(x)} = \frac{p(y/x) \cdot p(x)}{p(x)}$$

$$p(y/x) = \frac{p(y/x) \cdot p(x)}{p(x)} \rightarrow \begin{array}{l} \text{likelihood probability} \\ \text{prior probability} \end{array}$$

→ posterior probability

\Rightarrow 3 i/p variables x_1, x_2, x_3 and one o/p variable y

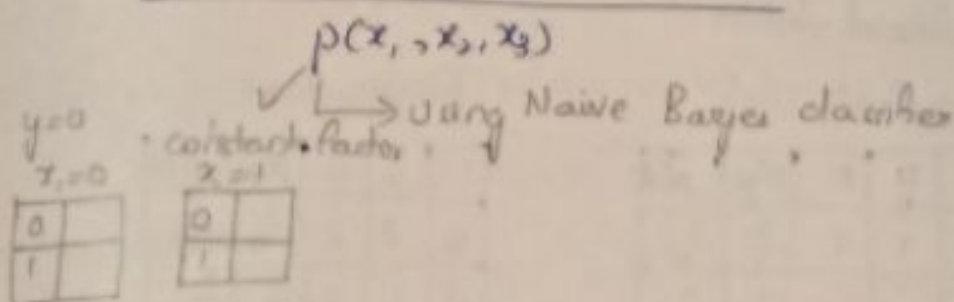
$$\begin{aligned} p(y/x_1, x_2, x_3) &= \frac{p(y, x_1, x_2, x_3)}{p(x_1, x_2, x_3)} \\ &= \frac{p(x_1, x_2, x_3/y) p(y)}{p(x_1, x_2, x_3)} \end{aligned}$$

$$2^0 - 1 = 2^4 - 1 = 15$$

$$= \frac{p(x_2, x_3 | y, x_1) p(x_1 | y) p(y)}{p(x_1, x_2, x_3)}$$

$$= \frac{p(x_3 | y, x_1, x_2) p(x_2 | y, x_1) p(x_1 | y) p(y)}{p(x_1, x_2, x_3)}$$

$$= \frac{p(x_3 | y) p(x_2 | y) p(x_1 | y) p(y)}{p(x_1, x_2, x_3)}$$



Two Approaches:-

- ① Naive Bayes
- ② Bayes Belief N/w.

① Naive Bayes classifier

→ Conditional Independence

Intelligence

Grade

Grade

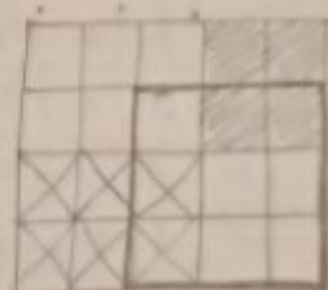
$$P(\text{Gr}/I, G) = P(\text{Gr}/I)$$

→ Naive Assumption

$$p(x_2 | y, x_1) = p(x_2 | y)$$

→ x_2 and x_1 are independent from each other in the class variable y (assumption)

The conditional independent Assumption states that the features are independent of each other given class label



$$\boxed{X} = \frac{6}{20} \quad \boxed{\text{shaded}} = \frac{4}{20}$$

$$P(\boxed{X} / \text{inner square}) = \frac{2}{9}$$

$$P(\boxed{\text{shaded}} / \text{inner square}) = \frac{2}{9}$$

Ex: day = {R, C, H, S} play = ?

y=yes; n=no

$$p(y/x) = \frac{p(\text{outlook}/y) p(\text{Temp}/y) p(\text{Hum}/y) p(\text{wind}/y) p(y)}{p(\text{outlook}, \text{Temp}, \text{Hum}, \text{wind})}$$

$$p(N/x) = \frac{p(\text{out}/n) p(T/n) p(H/n) p(w/n) p(n)}{p(\text{Co}, T, H, w)}$$

likelihood tables:

outlook:

	y	N	$p(y)$	$p(n)$
S	2	3	$\frac{2}{9}$	$\frac{3}{5}$
O	4	0	$\frac{4}{9}$	$\frac{0}{5}$
R	3	2	$\frac{3}{9}$	$\frac{2}{5}$
Total	9	5	1	1

Temperature:

	y	N	$p(y)$	$p(n)$
H	2	2	$\frac{2}{9}$	$\frac{2}{5}$
M	4	2	$\frac{4}{9}$	$\frac{2}{5}$
C	3	1	$\frac{3}{9}$	$\frac{1}{5}$
Total	9	5	1	1

Humidity:

	y	N	$p(y)$	$p(n)$
H	3	4	$\frac{3}{9}$	$\frac{4}{5}$
N	6	1	$\frac{6}{9}$	$\frac{1}{5}$
Total	9	5	1	1

wind

	y	N	$p(y)$	$p(n)$
S	3	3	$\frac{3}{9}$	$\frac{3}{5}$
w	6	2	$\frac{6}{9}$	$\frac{2}{5}$
Total	9	5	1	1

prior table:

class	Freq	prob
y	9	$9/14$
N	5	$5/14$
Total	14	1

$$p(y/w) = \frac{p(R/y) p(C/y) p(H/y) p(S/y) p(y)}{p(x)}$$

$$= \frac{(\frac{3}{9}) (\frac{3}{9}) (\frac{3}{9}) (\frac{3}{9}) (\frac{9}{14})}{1} = 0.00793$$

$$p(w/x) = \frac{p(r/n) p(c/n) p(h/n) p(s/n) p(n)}{\alpha}$$

$$= \frac{\left(\frac{2}{5}\right) \left(\frac{1}{5}\right) \left(\frac{4}{5}\right) \left(\frac{3}{5}\right) \left(\frac{5}{14}\right)}{\alpha} = \frac{0.013714285}{\alpha}$$

Eg day = {s, H, N, w}

$$p(y/x) = \frac{p(\text{outlook}/y) p(T/y) p(H/y) p(w/y) p(y)}{p(O, T, H, w)}$$

$$p(N/x) = \frac{p(O/N) p(T/N) p(H/N) p(w/N) p(N)}{p(O, T, H, w)}$$

likelihood tables:-

outlook:

	y	n	p(y)	p(n)
s	2	3	$\frac{2}{9}$	$\frac{3}{5}$
O	4	0	$\frac{4}{9}$	$0/5$
R	3	2	$\frac{3}{9}$	$\frac{2}{5}$
Total	9	5	1	1

Temp:

	y	n	p(y)	p(n)
H	2	2	$\frac{2}{9}$	$\frac{2}{5}$
M	4	2	$\frac{4}{9}$	$\frac{2}{5}$
C	3	1	$\frac{3}{9}$	$\frac{1}{5}$
Total	9	5	1	1

Humidity:

	y	n	p(y)	p(n)
H	3	4	$\frac{3}{9}$	$\frac{4}{5}$
N	6	1	$\frac{6}{9}$	$\frac{1}{5}$
Total	9	5	1	1

Wind:

	y	n	p(y)	p(n)
s	3	3	$\frac{3}{9}$	$\frac{3}{5}$
w	6	2	$\frac{6}{9}$	$\frac{2}{5}$
Total	9	5	1	1

prior table:-

class	F	P
y	9	$\frac{9}{14}$
n	5	$\frac{5}{14}$
total	14	1

$$p(y/x) = \frac{p(e/y) p(h/y) p(n/y) p(w/y) p(y)}{p(e, h, n, w)}$$

$$= \frac{(\frac{2}{9})(\frac{2}{9})(\frac{6}{9})(\frac{6}{9})}{\alpha} = \frac{0.0219}{\alpha}$$

$$p(n/x) = \frac{p(e/n) p(h/n) p(n/n) p(w/n) p(n)}{p(e, h, n, w)}$$

$$= \frac{(\frac{3}{5})(\frac{2}{5})(\frac{1}{5})(\frac{1}{5})}{\alpha} = \frac{0.0192}{\alpha}$$

19th Feb 2019

$$p(y/x) = \frac{p(x/y) p(y)}{p(x)}$$

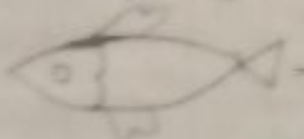

Likelihood / class condition
prior / Apriori
Bayes' Theorem.

Naive Bayes Theorem

posterior / posteriori

$$p(y/x, x_1) = \frac{p(x_2/y) p(x_1/y) p(y)}{p(x)}$$

Data: $\{x_1, x_2, \dots, x_n, y_i\}$
 $\underbrace{\hspace{10em}}_{\text{all features}} \quad \downarrow \text{class label}$

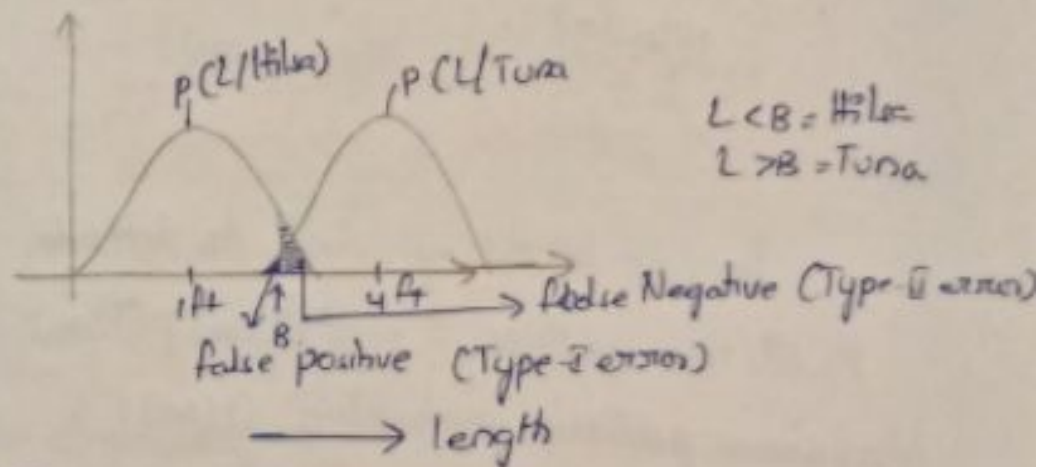
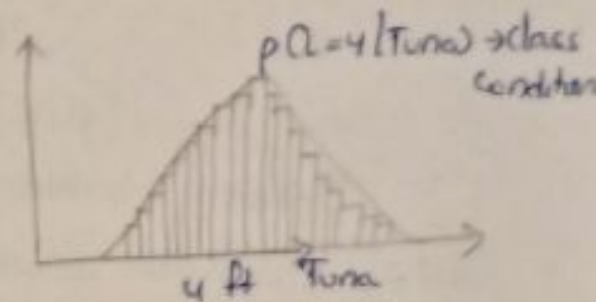
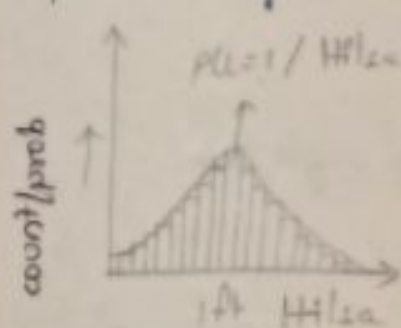
Eg.  - Hila  - Tuna

length $x_{(1)}$	class
1.5	H
3	T
1	H
4.5	T

Task:- classify the fishes into 2 classes based on length
find the decision boundary - B^*

Model:- 1000 - Hilsa 1000 - Tuna

plot Histogram



Hilsa?

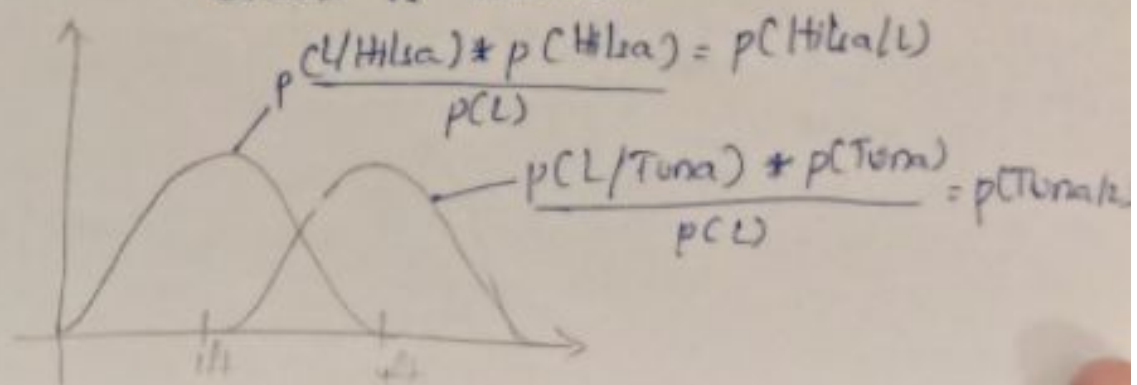
Error

Type-I error : actually hilsa but classified as Tuna

Type-II error : actually Tuna but classified as Hilsa

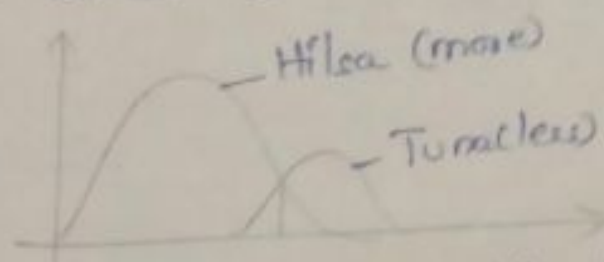
Error : Type I + Type II

Training :- find optimal decision boundary (B^*) such that error is minimum



$$P(B^*/Hilsa) = P(B^*/Tuna)$$

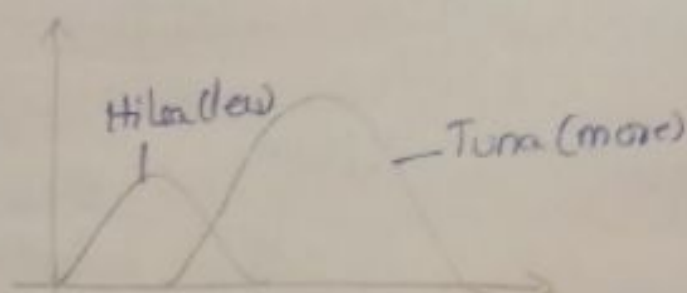
Case 1: Kolkata at Monsoon:



$$P(H) > P(T)$$

Decision boundary shifts to right side

Case 2: California in winter

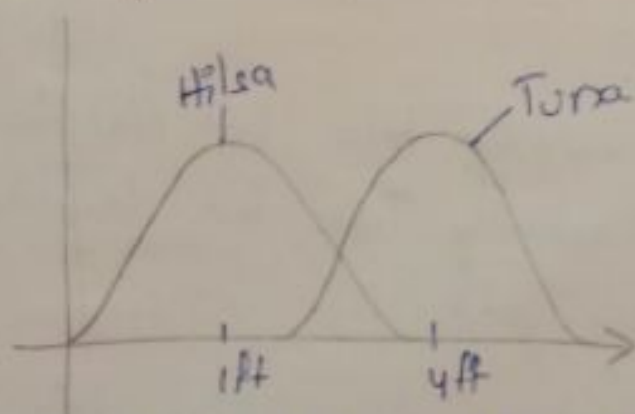


$$P(H) < P(T)$$

Decision boundary shifts to left side

$$P(B^*/Hilsa)P(Hilsa) = P(B^*/Tuna)P(Tuna)$$

Maximum posterior probability (MAP):



$$\text{if } P(Hilsa/L) < P(Tuna/L)$$

Tuna

else

Hilsa

21st Feb 19

day = { overcast, mild, normal, sunny }

play = ?

$$p(y/x) = \frac{p(o/y)p(m/y)p(n/y)p(s/y)p(y)}{p(o, m, n, s)}$$

$$= \frac{(\frac{4}{9})(\frac{4}{9})(\frac{6}{9})(\frac{3}{9})(\frac{9}{14})}{\alpha} = \frac{0.282}{\alpha}$$

$$p(n/x) = \frac{p(o/n)p(m/n)p(n/n)p(s/n)p(n)}{p(o, m, n, s)}$$

$$= \frac{(\frac{0}{5})(\frac{2}{5})(\frac{1}{5})(\frac{3}{5})(\frac{5}{14})}{\alpha} = \frac{0}{\alpha}$$

if one of the conditional probability is '0', then entire exp becomes '0' which we don't want.

original : $p(x_i/y) = \frac{N_{x_i y}}{N_y}$

Laplace Smoothing :-

$$p(x_i/y) = \frac{N_{x_i y} + 1}{N_y + c}$$

c = no. of classes

m = no. of features

$$\therefore p(o/n) = \frac{0+1}{5+2} = \frac{1}{7} ; p(m/n) = \frac{2+1}{5+2} = \frac{3}{7}$$

$$p(n/w) = \frac{1+1}{5+2} = \frac{2}{7} ; p(s/n) = \frac{3+1}{5+2} = \frac{4}{7}$$

$$p(n o/x) = \frac{(\frac{1}{7})(\frac{3}{7})(\frac{2}{7})(\frac{4}{7})(\frac{5}{14})}{\alpha} = \frac{0.00356}{\alpha}$$

$$p(yes/x) = \frac{(\frac{5}{11})(\frac{5}{11})(\frac{7}{11})(\frac{4}{11})(\frac{3}{11})}{\alpha} = \frac{0.031}{\alpha}$$

Text Classification:-

	Doc	Words	class
Train	0	India Delhi India	1
	1	India India Hyderabad	1
	2	India Mumbai	1
	3	Beijing china India	0
Test	4	India, India, India, Beijing china	?

mail - spam / not spam

Reviews - +ve / -ve

News - sports / politics / Entertainment

Bag of words Representation

		0	1	2	3	4	5	
Documents	0	2	1	0	0	0	0	Train
	1	2	0	1	0	0	0	
	2	1	0	0	1	0	0	
	3	1	0	0	0	1	1	
	4	3	0	0	0	1	1	
		words						

Multinomial
Naive Bayes

Vocabulary = { India, Delhi, Hyd, Mumbai, china, Beijing }
 0 1 2 3 4 5

	0	1	2	3	4	5
0	1	1	0	0	0	0
1	1	0	1	0	0	0
2	1	0	0	1	0	0
3	1	0	0	0	1	1
4	1	0	0	0	1	1

Bernoulli Naive Bayes

$$P(C|D) = \frac{P(w_1|C)P(w_2|C)P(w_3|C)P(C)}{P(w_1, w_2, w_3)}$$

$$P(w|c) = \frac{N_{w,c} + 1}{N_c + |V|}$$

N_c = total no. of words in the documents of class 'c'

$N_{w,c}$ = total no. of times word w appears in the documents of class 'c'

$|V|$ = size of the vocabulary

$$P(\text{India} | I) = \frac{5+1}{8+6} = \frac{6}{14} = \frac{3}{7}$$

$$P(I|D) = \frac{P(\text{India} | I)^3 P(\text{Beijing} | I) P(\text{China} | I) P(I)}{P(D)}$$

$$= \frac{\left(\frac{3}{7}\right)^3 \times \frac{1}{4} \times \frac{1}{14} \times \frac{3}{4}}{P(D)} = \frac{0.000301}{P(D)}$$

$$P(C|D) = \frac{P(\text{China} | C) P(\text{India} | C) P(\text{Beijing} | C) P(C)}{P(D)}$$

$$= \frac{\left(\frac{2}{9}\right) \left(\frac{2}{9}\right) \left(\frac{2}{9}\right) \left(\frac{1}{4}\right)}{P(D)} = \frac{0.0027}{P(D)}$$

22nd Feb 19

Advantages:-

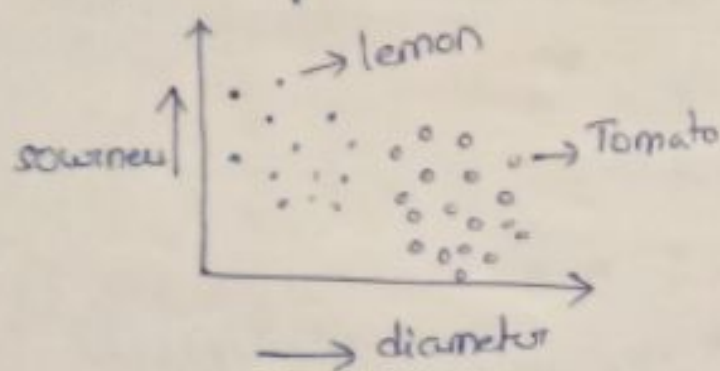
- (1) Very fast, requires less storage
- (2) It is robust to noise and irrelevant to features
- (3) It is very good / cool in a domain with many imp features
- (4) It is optimal when the independent assumptions hold.

Disadvantages:-

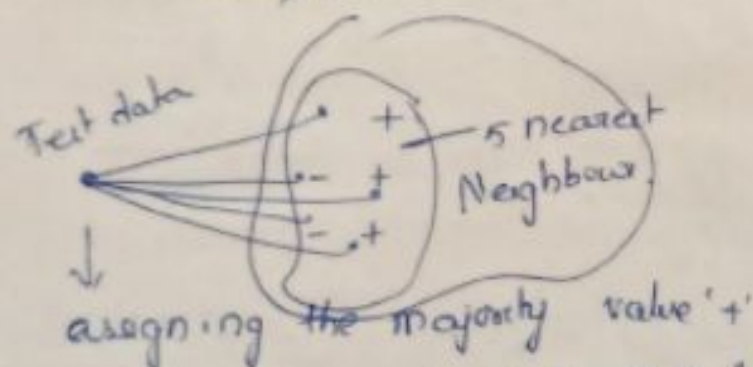
- (1) Independence assumption may not hold in reality
- (2) (The distribution) Because of finite dataset, we may not be able to estimate the distribution accurately.

k-Nearest Neighbour Classifier:-

- If it walks like a duck and quacks like a duck then it is probably a duck



- 1-NN \Rightarrow lemon
- 3-NN \Rightarrow lemon
- 5-NN \Rightarrow Tomato



- Requires :-
- (a) set of labelled data
 - (b) distance metric (to compute distance b/w)
 - (c) value of 'k'

To classify an unknown record:

- compute distances to other records
- identify its nearest neighbors
- assign majority class label of the nearest neighbors to the test record

Issues:-

- The value of k (deciding)
- choice of distance metric
- computation of complex unit



* If k is too small, then the classifier is sensitive to noise

* If k is very / too large, the neighborhood may include points from other classes

* ' k ' is a hyper parameter

Rule of thumb $k = \sqrt{N}$; N is no. of datapoints.

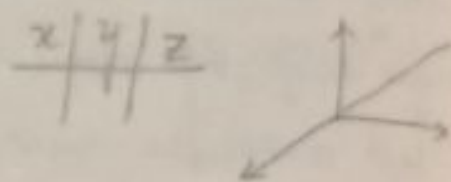
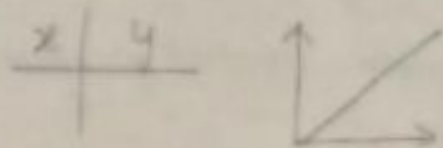
Distance Metrics :-

- Euclidean Distance
 - Manhattan Distance / City Block
 - correlation
 - Canberra
 - Mahalanobis
 - Quadratic
 - Rank correlation
 - chi-square
 - value difference metric
 - HDVM \rightarrow Heterogeneous
- } Numerical
- } \rightarrow categorical

25/2/2019

Distance Metrics :-

- simple attribute (single attribute)
- Vectors (multiple attributes)



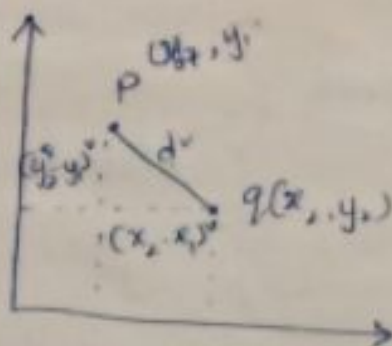
Simple Attribute :-

Attribute Type	Dissimilarity	similarity
Normal	$d = \begin{cases} 0 & \text{if } p=q \\ 1 & \text{if } p \neq q \end{cases}$	$s = \begin{cases} 1 & \text{if } p=q \\ 0 & \text{if } p \neq q \end{cases}$
Ordinal	$d = \frac{ x-y }{n-1}$ values are mapped to 0 to n-1	$s = 1-d$
Interval/Ratio	$d = x-y $	$s = -d$ $s = \frac{1}{1+d}$ $s = e^{-d}$

A - 0
B - 1
C - 2
D - 3
E - 4
 $d(B, E) = \frac{|4-1|}{4-1} = \frac{3}{3} = 1$

Vectors :-

① Euclidean Distance



$$d(p, q) = \sqrt{(x_2 - x_1)^2 + (y_2 - y_1)^2}$$

$$d(p, q, r) = \sqrt{(x_2 - x_1)^2 + (y_2 - y_1)^2 + (z_2 - z_1)^2}$$

$$d(p, q) = \sqrt{\sum_{k=1}^n (p_k - q_k)^2}$$

n = no. of dimensions (Attributes)

p_k and q_k are kth attribute

③ Minkowski Distance

Generalization of Euclidean Distance

$$d(p, q) = \left(\sum_{k=1}^n |p_k - q_k|^r \right)^{1/r}$$

r = parameter chosen by user

$r=2 \Rightarrow$ Euclidean Distance

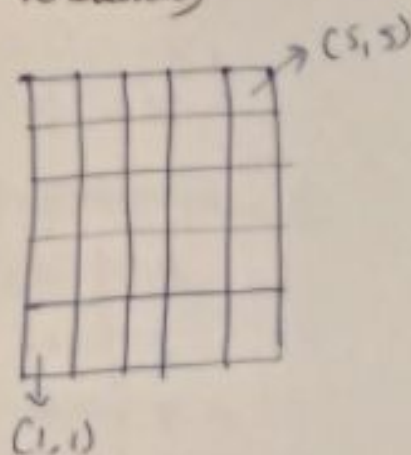
$r=1 \Rightarrow$ CityBlock / Manhattan Distance

③ CityBlock Distance (Manhattan Distance)

$$d(p, q) = \sum_{k=1}^n |p_k - q_k|$$

$$= |5-1| + |5-1|$$

$$\Rightarrow 4+4=8$$



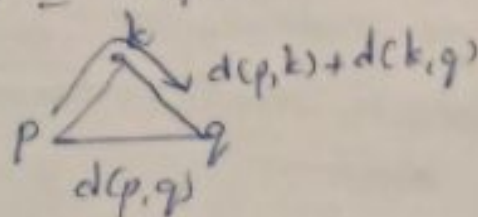
Euclidean properties:-

① $d(p, q) \geq 0$ (distance is always +ve)

② $d(p, p) = 0$ (distance at the particular point is 0)

③ $d(p, q) = d(q, p)$ (distance from $p \rightarrow q$ is equal to $q \rightarrow p$)

④ $d(p, q) \leq d(p, k) + d(k, q)$ [Triangle Inequality] property



Issues with Euclidean Distance:-

① scale Effect

Income (1000)	House size (100 sq.ft)
75	19
52	17
$d = 23.07$	

Income Rupees	House size (100 sq.ft)
75000	19
52000	17
$d = 23000$	

Income 1000	House size (sq.ft)
75	1900
52	1700
$d = 201.2$	

* Different features may have different measurement scales.

* Features with larger value will have greater impact on Distance - It may bias the performance of classifier.

② sparsity in the High Dimensional Data.

$$\begin{matrix} 1 & 1 & 1 & 1 & 1 & 1 & 1 & 1 & 1 & 1 & 1 & 0 \\ 0 & 1 & 1 & 1 & 1 & 1 & 1 & 1 & 1 & 1 & 1 & 1 \end{matrix} \} d = 1.4142 \rightarrow \text{more similar}$$

$$\begin{matrix} 1 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 1 \end{matrix} \} d = 1.4142 \rightarrow \text{more similar}$$

③ Collinearity Issue:-

x_1	x_2	x_3
2	1	1
4	6	3
6	8	5

If variables are collinear then co-related variable will have more impact on distance measure than independent variable.

Scale effects - Handling:-

Feature Scaling:- ① Normalization - Minmax scaling
② Standardization

① Normalization:- Data is scaled in the range of 0

$$X_{\text{norm}} = \frac{X - X_{\min}}{X_{\max} - X_{\min}}$$

- problem is Normalization suppresses the outliers

② Standardization:-

$$Z = \frac{X - \mu}{\sigma} = \frac{X - \bar{x}}{s}$$

\bar{x} = mean of the data

s = standard deviation of data

Eq:-

Hanks	Z-score	Normalization
90	-0.0823	-0.08
95	0.7438	0.74
92	0.2479	0.24
93	0.4132	0.41
98	1.2397	1.24
95	0.7438	0.74
98	1.2397	1.24
90	-0.0826	-0.08

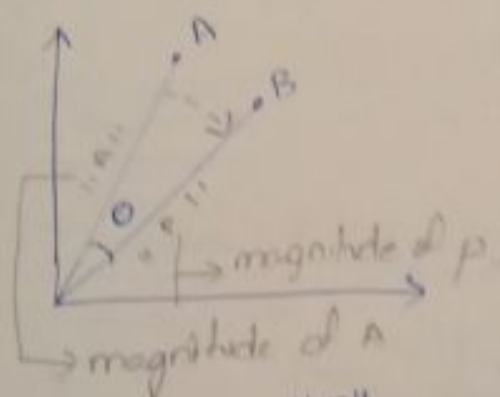
let $\bar{x} = 90.5$

$s = 6.05$

$Z \sim N(0, 1)$ Normalization
 \hookrightarrow Mean

Sparsity in High Dimensional Data:— Handling

(i) Cosine Similarity :-



$\cos 0^\circ = 1$ (most similar)

$\cos 90^\circ = 0$

$\cos 180^\circ = -1$ (least similar)

$$\cos \theta = \frac{\|P\|}{\|A\|} = \frac{\vec{A} \cdot \vec{B}}{\|\vec{A}\| \cdot \|\vec{B}\|}$$

$$\|P\| = \frac{\vec{A} \cdot \vec{B}}{\|\vec{B}\|} = \frac{\sum_{i=1}^n A_i \cdot B_i}{\sqrt{\sum_{i=1}^n A_i^2} \sqrt{\sum_{i=1}^n B_i^2}} = \frac{5}{\sqrt{42} \sqrt{6}}$$

$\vec{A} = [3, 2, 0, 5, 0, 0, 0, 2, 0, 0]$

$\vec{B} = [1, 0, 0, 0, 0, 0, 0, 1, 0, 2]$

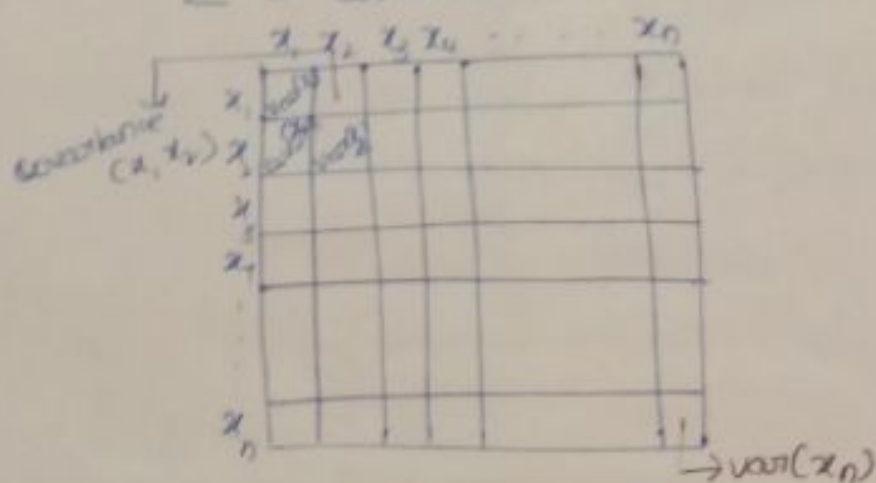
$= \frac{5}{15.874}$

$= 0.31497$

26-1-19
(5) Mahalanobis Distance :-

$$d(p, q) = \sqrt{(p-q)^T \Sigma^{-1} (p-q)}$$

Σ = covariance matrix of Data



$$p = \begin{bmatrix} \end{bmatrix}_{n \times 1} \quad \text{--- } n \times 1$$

$$q = \begin{bmatrix} \end{bmatrix}_{n \times 1}$$

$$p - q = n \times 1$$

$$(p-q)^T = 1 \times n$$

$$\Sigma^{-1} = n \times n$$

$$= \sqrt{\underbrace{(p-q)^T}_{1 \times n} \underbrace{\Sigma^{-1}}_{n \times n} \underbrace{(p-q)}_{n \times 1}}$$

$$\Sigma_{j,k} = \frac{1}{n-1} \sum_{i=1}^n (x_{ij} - \bar{x}_j)(x_{ik} - \bar{x}_k)$$

$$\boxed{d(p, q) = \sqrt{(p-q)^T (p-q)}}$$

Σ = Identity Matrix, $HD \in ED$

Binary Attribute :-

$$p = 0 \ 1 \ 0 \ 0 \ 1 \ 1 \ 0$$

$$q = 1 \ 0 \ 1 \ 0 \ 1 \ 0 \ 0$$

Jaccard Coefficient :

	1	0
1	1	2
0	2	2

$$M_{00} = \# \text{ of } p=0 \ q=0$$

$$M_{01} = \# \text{ of } p=0 \ q=1$$

$$M_{10} = \# \text{ of } p=1 \ q=0$$

$$M_{11} = \# \text{ of } p=1 \ q=1$$

→ measure of similarity in presence

$$J = \frac{M_{11}}{M_{01} + M_{10} + M_{11}} \Rightarrow \frac{2}{2+2+1} = \frac{1}{5}$$

0 - absent

1 - present

So M_{00} doesn't have importance.

Asymmetric Binary Attribute \Rightarrow gives more importance to presence (1)

Symmetric Binary Attribute \Rightarrow Both outcomes are equally important

→ measure of dissimilarity -

$$\text{Simple Binary Coefficient} = \frac{M_{01} + M_{10}}{M_{00} + M_{01} + M_{10} + M_{11}}$$
$$= \frac{4}{7}$$

Nominal Attribute :-

$$\text{simple matching coefficient} = \frac{mm}{n}$$

n = total no of attributes

mm = # of mismatches

Text Classification

0 1 2 3 4 5 class cosine Distance

0	2	1	0	0	0	0	1	0.81
1	2	0	1	0	0	0	1	0.81
2	1	0	0	1	0	0	1	0.64
3	1	0	0	0	1	1	0	0.87
4	3	0	0	0	1	1	2	

words

S.NO	Name	Gave Birth	can Fly?	lives in water	Have legs	class	Distance
1	Human	y	n	n	n	m	2/4
2	salmon	n	n	y	n	nm	1/4
3	python	n	n	n	n	nm	2/4
4	whale	y	n	y	n	m	0
5	Frog	n	n	y	y	nm	2/4
6	komodo	n	n	n	y	nm	3/4
7	Bat	y	y	n	y	m	3/4
8	pigeon	n	y	n	y	nm	1
9	Cat	y	n	n	y	m	2/4
10	Dolphin	y	n	y	n	m	0
11	—	y	n	y	n	?	

m

$$k = \sqrt{10} = 3 \quad 2m$$

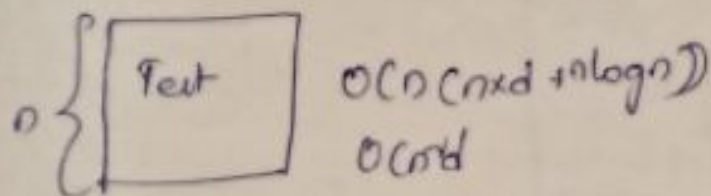
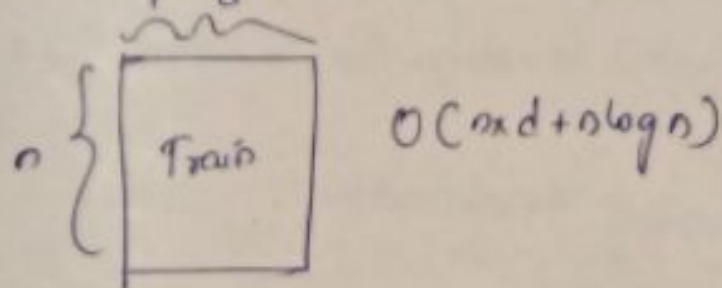
$$1nm$$

2-mammals are there (majority)
so mammal

Performance of k-NN

1 Time Complexity - Expensive

To determine the Nearest Neighbour on test data we must compute distances to all training data. It is expensive.



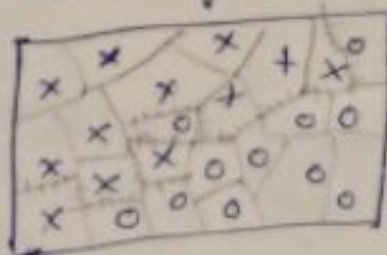
- Remove redundant data - condensation
- pre sort data into fast data structures
 - \rightarrow kd tree
- Compute only approximate distance (LHF)

2 Space complexity - must store all training data

Training - Condensation

Condensation

Decision bounding in k-NN -



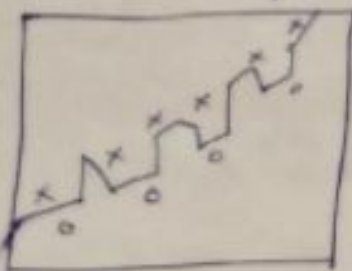
A Voronoi diagram divides the space into cells, such that each ^{cell} contains one sample and every location within cell is closer to that sample than any other sample.

The Decision Boundary separates the class Region into
Voronoi NN Decision Rule

knowledge of the Decision boundary is sufficient to
classify new data. Retain only those points necessary to
generate an identical boundary. This is called Condensation

Two types:-

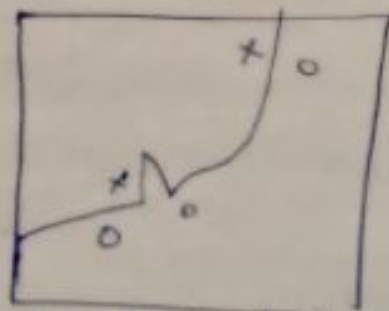
① Decision Boundary Condensation



This condensation is subset with to Nearest Neighbour
Decision boundary

- A subset whose nearest neighbour decision boundary
is identical to the boundary of the entire training set

Minimum consistent set condensation



The smallest subset of training data that correctly
classifies all the original data