

NAME: Nilkanth Jadhav - 8087089261

STD: \_\_\_\_\_

DIV: \_\_\_\_\_

ROLL NO.: \_\_\_\_\_

SUBJECT: R with Data Science

## INDEX

SR. NO.	DATE	TITLE	PAGE NO.
	30/7/18	Exploratory Data Analysis & Visualization	
	02/8/18	Central limit theorem	
	23/8/18	Machine Learning	
	24/8/18	Regression	
	28/8/18	Ridge & Lasso Regression	
	30/8/18	Time series, Stepwise regression	
	03/9/18	Classification <ul style="list-style-type: none"><li>- KNN</li><li>- Logistic Regression</li><li>- KNN for regression</li><li>- Naive Bayes classification</li></ul>	
	05/9/18	Decision Trees	
	06/9/18	Unsupervised Learning <ul style="list-style-type: none"><li>Clustering<ul style="list-style-type: none"><li>- K-means</li><li>- Hierarchical</li><li>- DB Scan</li></ul></li></ul>	
	19/9/18	SYM	
	20/9/18	ANN	
	26/9/18	Market Basket Analysis.	

# Exploratory Data Analysis and Visualization (EDA)

- Understand the problem (Project charter) optional.
- collect data related to problem.  
what is current status about problem?
- Brain-storming session with management team.  
discuss about problem and find why problem?  
define metadata, take surveys about sub-factors.
- Collect data from warehouse.
- Process data - extraction, integration, clean-up.
- Transforming data - making data ready for analysis, filtering, sorting, numeric  $\rightarrow$  category (discretization), column combination  $\rightarrow$  new column
  - Box-Cox Transformation & Johnson transformation  
un-normalized data to normalized data.
- Understand data - mathematical summary + visualization summary
  - column wise analysis.
  - pairwise comparison,  $y \rightarrow x_1, y \rightarrow x_2, \dots, y \rightarrow x_{25}$

↑  
univariate & bivariate

Univariate analysis - column by column.

- Categorical
- Quantitative

## Univariate

## Categorical

## Mathematical Summary

Count, proportion, percent  
ratio, mode

## Visualization

Pi chart, bar chart,  
histogram, box plot.  
dot plot

## Quantitative

mean, median, mode,

std. deviation, range,

(min , max, percentile)

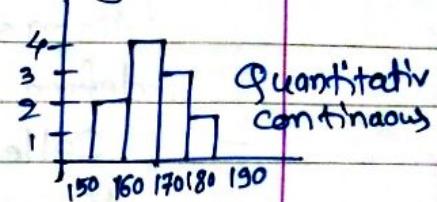
8, 8, 8- } 5 no. summer

IQR. ( $\theta_1 = -\theta_2$ ) distribution

of 50% middle elements.

## Box & Whiskers plot

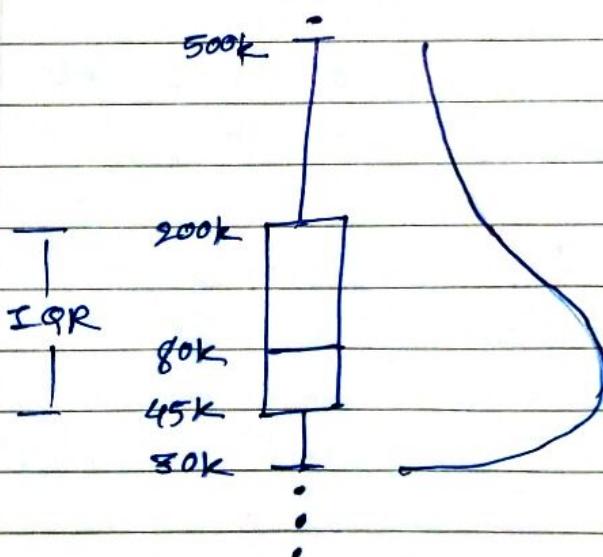
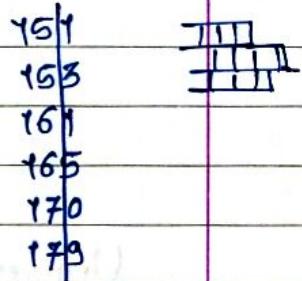
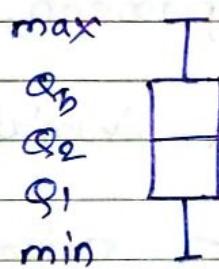
## Histogram



Class interval  
how data is distributed  
order

- Information loss  
in summarization  
(stem-leaf plot)

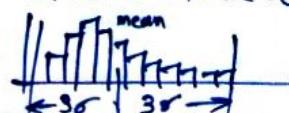
Sort data &  
partition into  
stem - leaf.



## Outlier detection

$1.5 * IQR$  will limit bar graph shooting out of paper, outliers of higher values are represented by dots.

$\mu \pm 3\sigma$  are outliers.



Bivariate

## Mathematical Summary

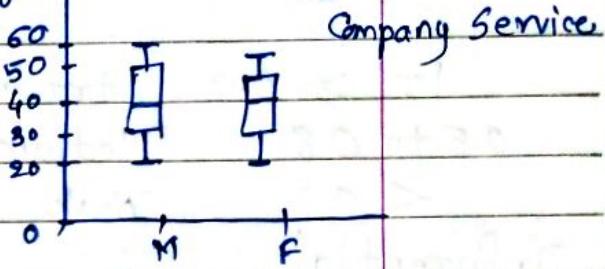
## Visualization Summary

- $\checkmark C \rightarrow Q \rightarrow \text{location} \rightarrow \text{Salary}$ , "whether cond"  $\rightarrow$  "Acad", Gender vs. Hemoglobin  
 $\times Q \rightarrow C \rightarrow \text{smoking} \rightarrow \text{Cancer}$ , location  $\rightarrow$  Gender, season  $\rightarrow$  dieses.  
 $\checkmark C \rightarrow C$   
 $\checkmark Q \rightarrow Q \rightarrow \text{experience} \rightarrow \text{salary}$ ,  $\therefore \leftarrow \frac{\text{no. of hr. study}}{\text{salary}} \rightarrow \text{saving}$ , petrol  $\rightarrow$  avg. time  $\rightarrow$  product,  $\frac{\text{distance}}{\text{Time}}$   
 $\frac{\text{no. of product sold}}{\text{Time}}$

$C \rightarrow Q$  Comparative five number summary

	Male	Female
Max	60	50
Q <sub>3</sub>	45	35
Q <sub>2</sub>	35	30
Q <sub>1</sub>	27	25
Min	20	20

side by side Boxplot



$C \rightarrow C$  Smoking Cancer

Y	Y
Y	N
N	N
N	Y
Y	Y
:	:

Gender

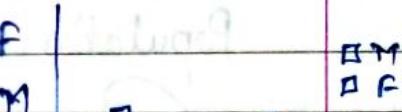


Table  
Two-way  
contingency

		B	P
		0.47	0.53
		0.45	0.55

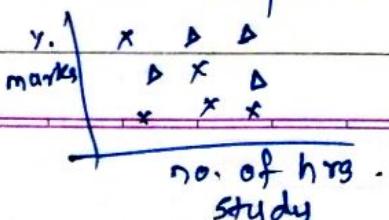
Comparative proportion.

- Gender equality independent of city
- Salary distribution for medium band is dependent.

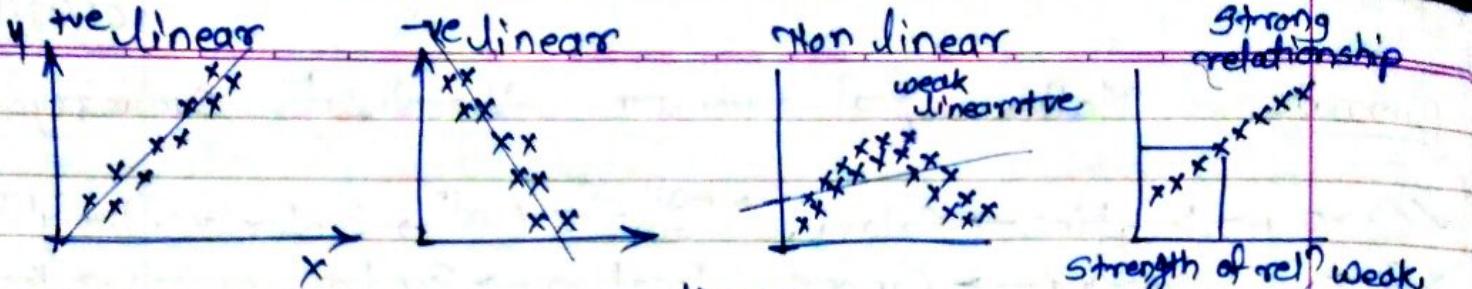
	B	P	L	M	H
B	334	471			= 705
P	762	916			= 1680

$Q \rightarrow Q$  Experience vs salary

Scatterplot ( $Q \rightarrow Q$ )



Stratified Scatterplot.



$X \& Y$

- Direction of reln.
- Linearity of rel?
- Strength of rel?

has

linear  
Co-relation coefficient  
( $r$ )  $-1 < r < +1$

Strength of rel? weak



$r = 0$  No relationship.



$|r| \geq 0.8$  Strong rel?

0.5 to 0.8 Medium - - -

< 0.5 Weak - - -

## Inferential

Census every 10 years. ~~every 10 years~~

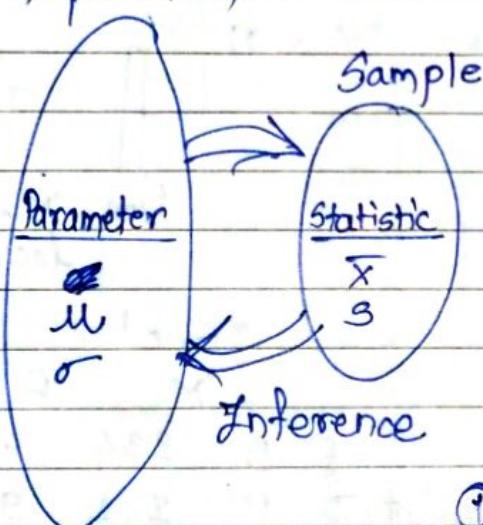
Tiger Census

## Population

Same,

$\mu \rightarrow \bar{x}$

$\sigma \rightarrow s$



Inference is  
 $\bar{x}$  estimates  $\mu$   
 $s$  estimates  $\sigma$

- We estimate population parameters using sample statistic

① - Point estimate  $\mu = \bar{x}$   
there will be error

② - Interval estimate

confidence interval =  $(\bar{x} \pm \text{margin of error})$   
 $\mu$  varies in interval

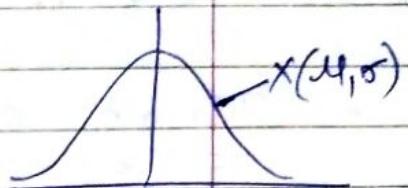
③ - Test of hypothesis.

- Population parameters remains same as population does not change. Sample statistic may change as sample changes due to random items.

Study  
well

## Central Limit Theorem -

- Centering, spread, shape.



- mean of  $X$  and mean of  $\bar{X}$  are centered at same point

$$\text{mean } X \rightarrow \mu$$

$$\text{mean } \bar{X} \rightarrow \bar{\mu}$$

when no. of samples are high

$$\bar{X} = \frac{\sum X}{N}$$

no. of  
sample  
taken

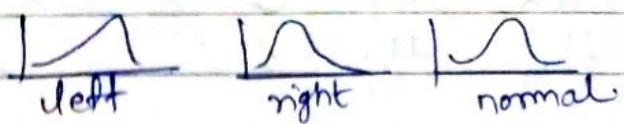
$$\text{sd. } \bar{X} = \frac{\sigma}{\sqrt{n}}$$

no. of pts  
sample

- For large enough sample size of  $\bar{X}$  will tend to be normal irrespective of distribution of  $X$ .

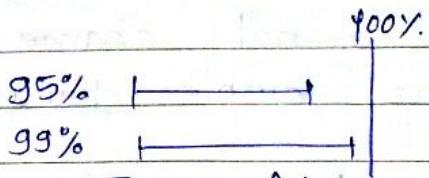
Sample size,  $n=30 >$  its normal

~~faster to get normal~~



## Interval Estimate -

- Confidence level/interval



$$= \bar{x} \pm \text{margin of error.} = \bar{x} \pm \text{confidence multiplier} \times \frac{\sigma}{\sqrt{n}}$$

% Area will be covered

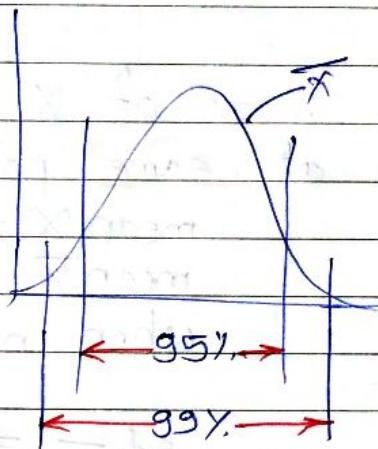
- Margin of error.

$$= 2 \times \frac{\sigma}{\sqrt{n}}$$

$$= 3 \times \frac{\sigma}{\sqrt{n}}$$

2.64

How many std. deviation times ~~center~~ range is away from std. deviation



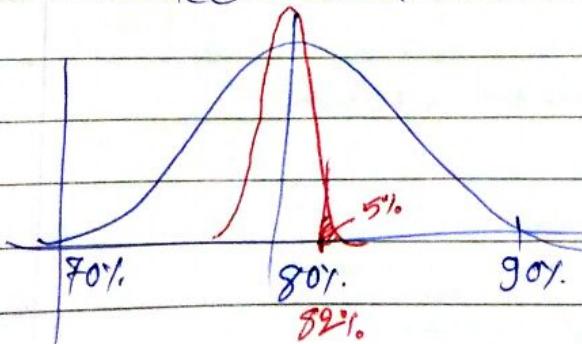
## Test of hypothesis -

- Domain claims. & management claim  
alternative hypothesis. null hypothesis ( $H_0$ )  
change No change

$$H_0 : \mu_p = 80\%$$

$$H_a : \mu_p > 80\%$$

- Collect data and test data for hypothesis.



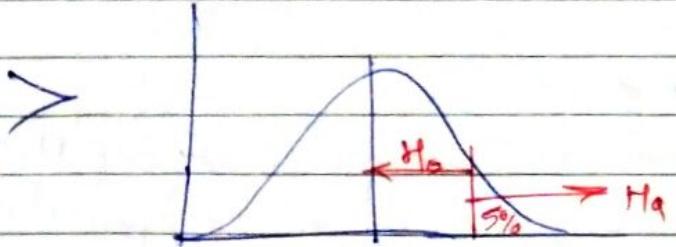
25 samples.

$$\bar{x} = 80\%$$

$$6d\bar{x} = \frac{5}{\sqrt{25}} = 1$$

$$\sigma = 5$$

your eg.  $H_0$  = coin fair  $H_a$  = unfair  $\alpha$  = threshold.  
level of significance



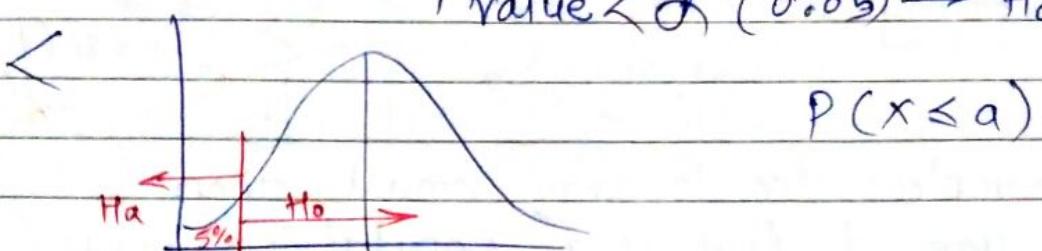
$$P(X > a)$$

$$5\% = \alpha = 0.05$$

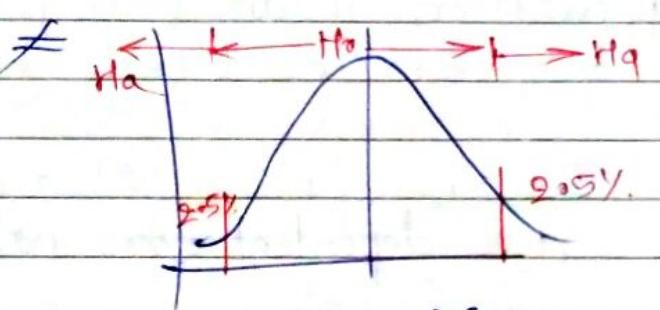
level of significance

$$P(X > b) = p\text{-value}$$

$P\text{value} < \alpha (0.05) \rightarrow H_a$  is accepted



$$P(X \leq a)$$



## Types of hypothesis tests

### Univariate

#### Categorical.

#### proportion

$$H_0: p = 0.20$$

$$H_a: p < 0.20$$

①  $\chi^2$ -test for population proportion.

Transformation (standardized form of normal distn)  
 $Z : X \rightarrow Z = \frac{X - \mu}{\sigma}$   
 where  $\mu = 0, \sigma = 1$

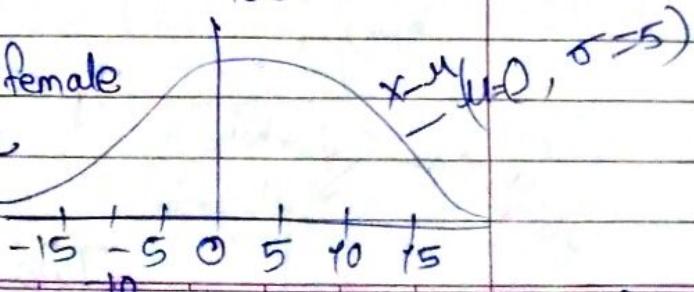
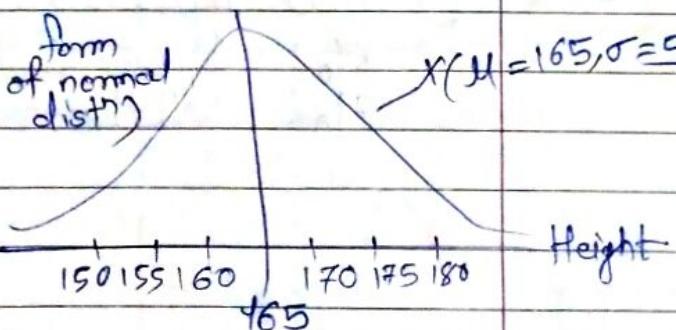
$$X(\mu = 165, \sigma = 5)$$

#### Application of $z$ -tron<sup>s</sup>.

- Comparing positions.

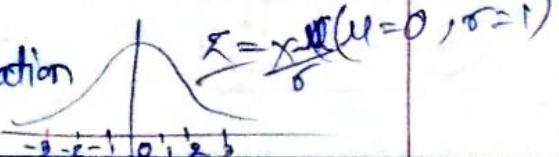
eg. board marks male vs. female  
height male vs. female

- Two normal distn can be compared with help of  $z$  transform.



unstd → std

Standardization



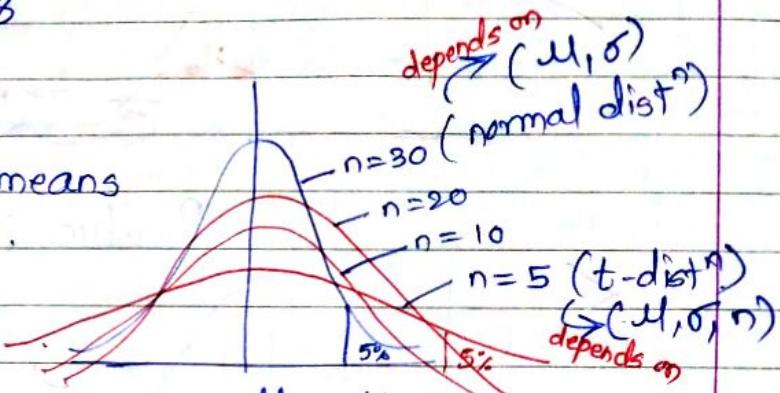
## Quantitative

$$H_0: \mu_p = 0.8$$

$$H_a: \mu_p > 0.8$$

t-distribution

- less no. of sample means  
uncertainty is high.



- If sample size is very small then we use t test for population mean ( $n < 30$ )
- $\Sigma$  test of population mean. (only t test is considered., one sample t-test known as)

## Bivariate

$X \times Y$   
 $c \rightarrow \alpha$   
 $c \rightarrow c$   
 $\alpha \rightarrow \alpha$

Domain claim -  $X$  and  $Y$  are dependent or not.

$$H_0: X \text{ and } Y \text{ are independent}$$

$$H_a: X \text{ and } Y \text{ are dependent}$$

Smoking and Cancer. { dependent ideally.  
Hr. of study and marks. }

$c \rightarrow \alpha$  ② Two independent samples

$$H_0: \text{Pune \& Mumbai salaries are different.}$$

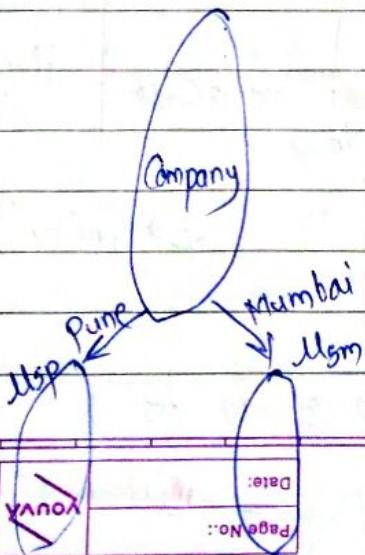
$$H_a: \text{Weather \& Car breakdown are dependent}$$

$$H_0: \mu_{\text{SalaryPune}} = \mu_{\text{SalaryMumbai}}$$

$$\mu_{sp} - \mu_{sm} = 0$$

$$H_a: \mu_{sp} < \mu_{sm}$$

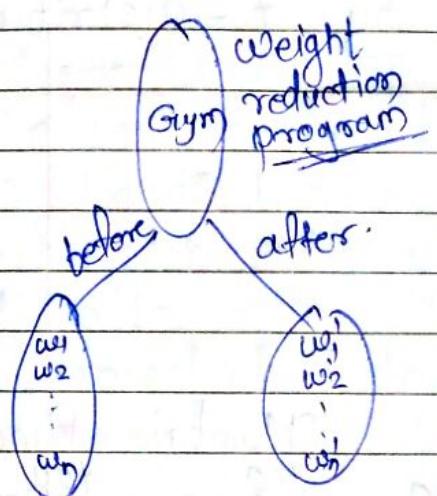
$$\mu_{sp} - \mu_{sm} < 0$$



Two sample t-test.

e.g. p value = 0.063 no difference  $> 0.05$

- 2) Two dependent samples - Two elements are connected with each other
- eg. fitness club joining for weight loss.
- weight should be measured before and after experimentation.
  - same group of people / machines / things should be under consideration before and after experimentation.
  - sales before & after training.



$$H_0: \mu_d = 0$$

$$H_a: \mu_d > 0$$

paired t-test

$$\underline{\mu_d}$$

$$d_1 = w_1 - w_1'$$

$$d_2 = w_2 - w_2'$$

$$d_n = w_n - w_n'$$

$$p\text{-value} = 0.0963$$

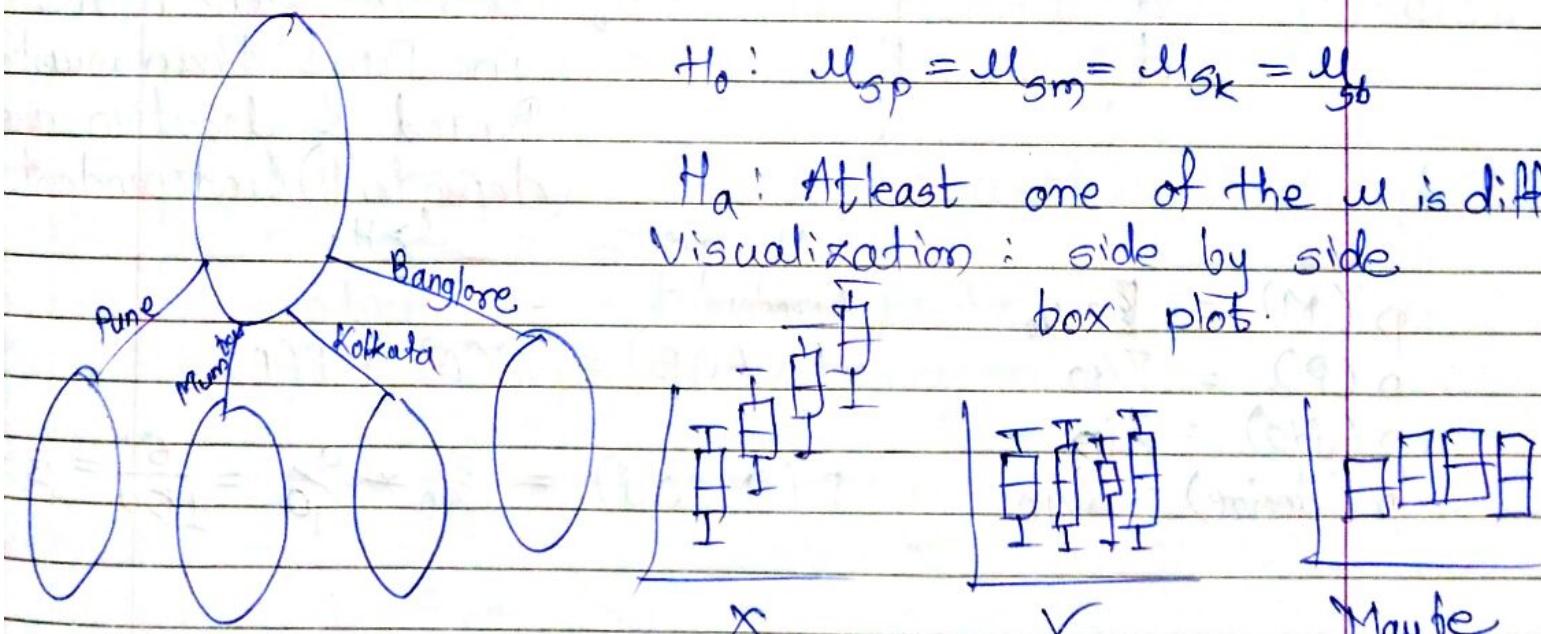
$H_0$  is true.

- 3) More than two independent sample.

$$H_0: \mu_{sp} = \mu_{sm} = \mu_{sk} = \mu_{sb}$$

$H_a$ : Atleast one of the  $\mu$  is diff.

Visualization: side by side  
box plot.



To compare means in box plot, f test.

$$F = \frac{\text{Variability between samples}}{\text{Variability within samples}} > 1$$

F statistics  $> 1$   
should be

### ANOVA (Analysis of variance) F test

$< 0.05$  good  
Overall model significance given by f-distribution  
Individual significance given by t-distribution.

C → C

Priority vs severity  
weather vs delay  
brand vs location

Brand      Location

jyo      Pune

desire      Mumbai

M      M

P      P

d      P

d      P

:

:

:

:

$$P(M) = \frac{3}{10}$$

$$P(P) = \frac{7}{10}$$

$$P(jyo) = \frac{9}{10}$$

$$P(\text{desire}) = \frac{8}{10}$$

variables are  
independent

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

$$P(M \cap j) = \frac{3}{10} * \frac{9}{10} = \frac{6}{100} = \frac{3}{50}$$

$H_0$ : desire sale more  
in Pune than mumbai  
Brand & location are  
(dependent)/ independent

$\hookrightarrow H_a$

	Mumbai	Pune
jyo	100 $E=60$	100 $E=140$
desire	200 $E=240$	600 $E=560$
M	300	700

## Chi Square Statistics —

$$\chi^2 = \sum \frac{[\text{Actual} - \text{Expected}]^2}{\text{Expected}}$$

$$\chi^2 = \sum \frac{[A-E]^2}{E}$$

$$= \frac{[100-60]^2}{60} + \frac{[100-140]^2}{140} + \\ \frac{[200-240]^2}{240} + \frac{[600-560]^2}{560}$$

## Chi Square Test —

Chi square distribution will have distribution which 5% value, if  $\chi^2$  falls then they are dependent otherwise independent.

Assumptions in ppt check.

S → S salary - saving , salary - expenditure , experience - salary , age - salary , stock market of day , fund value , (time) stress - strain

$$y = mx + c \quad c - \text{intercept} , m = \text{slope}$$

$$\text{Salary} = m * \text{experience} + c$$

$$\text{Salary}' = m * (\text{expt} + 1) + c$$

$$\text{salary}' - \text{salary} = m = \text{increment}$$

$$\text{Sales} = m * \text{advertising expences} + c$$

$$y = B_0 + B_1 x$$

$$y = \beta_0 + \beta_1 x$$

$H_0$ : Linear relationship bet' y and x does not exist  $\beta_1 = 0$

$H_a$ :  $\beta_1 \neq 0$

$\beta_1$ -t distribution should be considered.

↳ t-test for population slope

Judgement eg.

		Truth		evidence
		Innocent	Guilty	
Judge I.	G	✓	Type 2 ( $\beta$ )	$\alpha \downarrow \beta \uparrow$
		Type 1 ( $\alpha$ )	✓	$\alpha \uparrow \beta \downarrow$

→ More sample size  
more accurate  
judgement

Supplier eg.

Pc.	defective
500	Acceptable 25
100	Acceptable 5

1) Producers' risk ( $\alpha$ )	500	Innocent $\rightarrow$	20
Type 1	100	Guilty $\rightarrow$	6

Power of test -

$$1 - \beta$$

$\beta$  should be as less as

2) Consumers' risk ( $\beta$ )	500	Guilty $\rightarrow$	30
Type 2	100	Innocent $\rightarrow$	4

AB testing ( $\alpha \beta$  testing)

# Machine Learning

23/08/18

SDLC  $\Rightarrow$  Requirement - Analysis - Design - Dev (Coding) - Testing - Implementation - maintenance & support

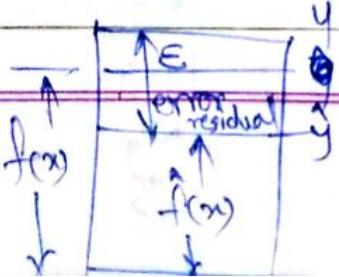
- Business needs are identified and conveyed to software company
- Analysis and design - break tasks into tasks modules
- Code modules done - testing - Implementation -
- Machine learns on own based on experience without being explicitly programmed.
- If we provide more experience to machine, the performance of machine will improve.

e.g. Attrition (y)

- Understand (y)
- define scope of y.
- Baseline y - what is current situation
- Define causes of y.  $\rightarrow X$
- define metadata of X
- Collect data of X
- Process data - cleaning transformation
- EDA & visualization
- Inference which are important X (remove unwanted X)
- 10<sup>th</sup> ML

$y = f(x) + \epsilon$  - relationship bet' x and y.  
predict  $\hat{y}$  (signal noise  $\epsilon$  is error factor)

estimate machine tries to learn best value of  $f \Rightarrow \hat{f}$ ;  $\hat{y} = \hat{f}(x)$  (hat 1)



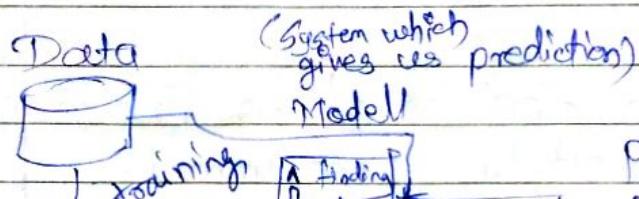
error  
residual

$$y - \hat{y} = f(x) + \epsilon - \hat{f}(x)$$

$$= [f(x) - \hat{f}(x)] + \epsilon$$

Reducable  
error

Irreducible  
error.



Parametric

- form of model

e.g.  $y = mx + c$

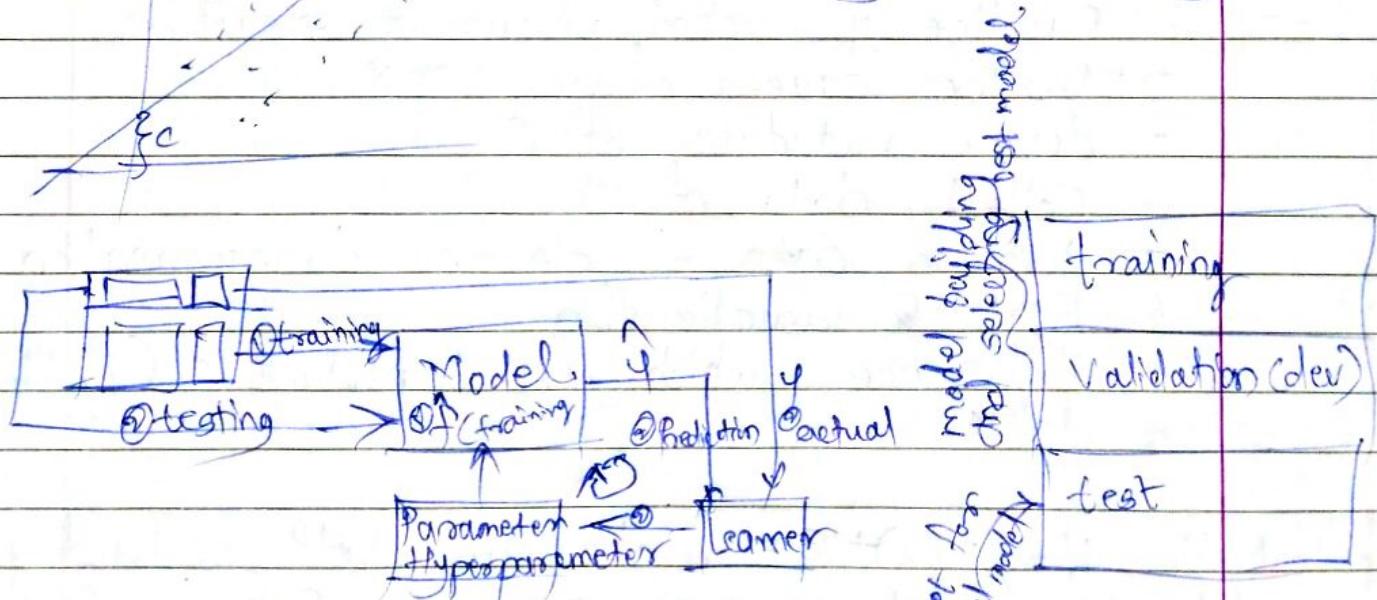
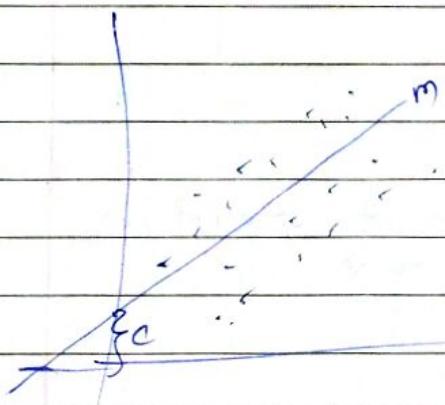
parameter.  
 $m$  &  $c$

nonparametric

- form decided by  
algorithm itself.

(hyperparameters)

Existing data divided in two  
parts randomly into  
training and testing data



$M_1$

$M_2$

$M_3$

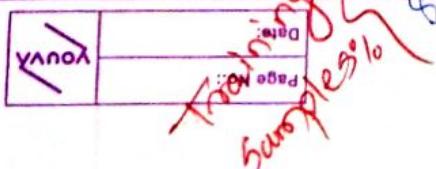
$M_4$

$M_5$

$$x \rightarrow y, y - \hat{y}$$

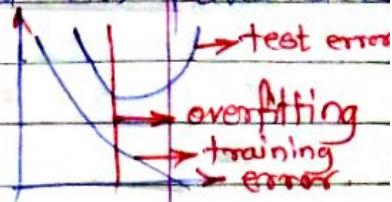
$$e_1, e_2, \dots$$

train  
validation



## Machine learning example -

Predicting weight of person based on parameters like height, fat %, etc.



- $f(x) = \text{height}$  → ①
- $f(x) = \text{height} + \% \text{ fat}$  → ②
- eq ① gives less information about weight as it only considers height.
- eq ② will be more optimised as it also considers % fat in function equation for predicting weight.
- $f(x) = \text{height} + \% \text{ fat} + \text{bone-density}$  → ③
- eq ③ will be more appropriate to consider but introducing too many parameters may lead to introduction of error in prediction (problem of overfitting)
- Cross validation approach uses complete dataset to train model and test model.
- Validation set approach uses % of data.

### LOOCV - Leave one out cross validation ( $k=n$ )

Tr	$e_1^1$	$e_1^2$
$\vdots$	$\vdots$	$\vdots$
$\frac{n-2}{n-1}$	$e_{n-2}^1$	$e_{n-2}^2$
$\vdots$	$e_{n-1}^1$	$e_{n-1}^2$
mean = $e_m^1$	$e_m^2$	$\vdots$
$e_m^1$	$M_n \rightarrow e_1^1$	$e_m^1$

$$M_1 \rightarrow \begin{cases} e_1^1 \\ e_1^2 \\ \vdots \\ e_n^1 \end{cases} \quad \{e_m^1\}$$

$$M_2 \rightarrow \begin{cases} e_1^2 \\ e_2^2 \\ \vdots \\ e_n^2 \end{cases} \quad \{e_m^2\}$$

$$M_n \rightarrow \begin{cases} e_1^n \\ e_2^n \\ \vdots \\ e_n^n \end{cases} \quad \{e_m^n\}$$

- has significant cost over computation

- any two successive models will have almost same as majority rows are same during training

In correlated model predictions & errors are much similar.

such models are known as correlated models.

## K-fold cross validation.

1	
2	
3	
4	
5	

	Test	
M1	1,2,3,4	5
M2	1,2,3,5	4
M3	1,2,4,5	3
M4	1,3,4,5	2
M5	2,3,4,5	1

Suggested value  
of k = 10.

$$M_1 e_{\text{mean}} = \frac{e_1 + e_2 + \dots + e_5}{5}$$

Cross validation -

Sampling without replacement

$M_2 e_{\text{mean}}$

$M_3 e_{\text{mean}}$ .

whichever model has  $e_{\text{mean}}$  least is selected.

Bootstrap - Sampling with replacement

100 sample rows  $\rightarrow$  5 ~~sample~~ test data of 100 sample row  
 $5 \times 100 = 500$  rows required.  
 duplication is allowed here

## Weather Prediction - Categorical data.

Sunny (S), Cloudy (C), Rainy (R)

(Mean squared  
Error)

Predictions	Actual
S	S
S	C
R	C
R	R
C	C
S	S
C	R
R	R
S	S
C	C

30% error

Misclassification =  $\hat{g}_0$

$$\text{Quantitative} = \frac{\sum (y - \hat{y})^2}{n}$$

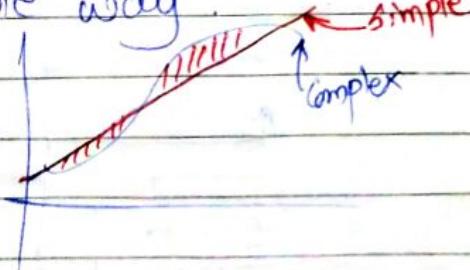
Categorical has misclassification

$$\text{Residual sum of squares} = \sum (y - \hat{y})^2$$

$$y - \hat{y} = \text{Residual} = \text{Reducable} + \text{Irreducible}$$

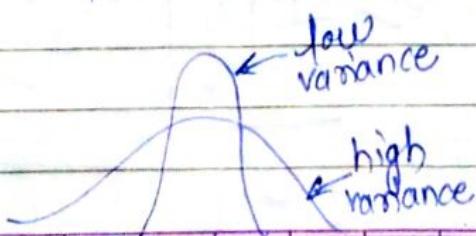
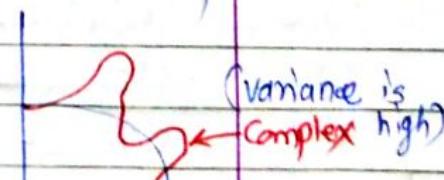
Bias error:

- fitting function in favour of some parameters in simple way.



Variance ( $\sigma^2$ ) goes high

- Complementing function arises due to introduction of excessive parameters



- Variance should be minimum so that function will not turn over around in space.

## Model Selection - Performance

$$\text{Error} = f - \text{accuracy}$$

- Accuracy
- Time to train & test
- Computation cost
- Scalable
- Simple
- Interpretable (One should understand model what's going on inside.)

## Regression

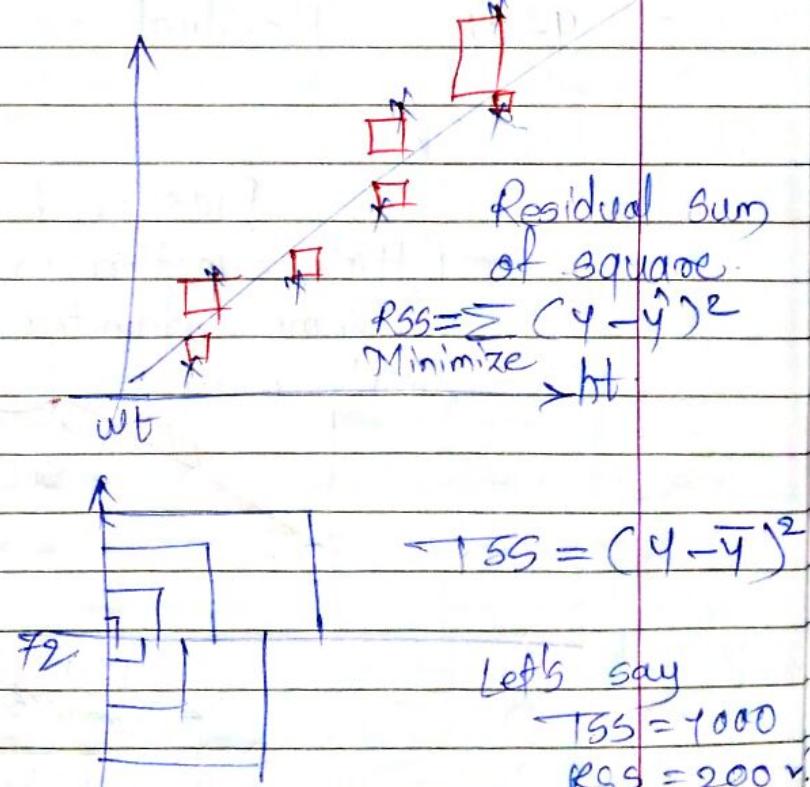
### ① SLR - Simple Linear Regression

$$TSS = \sum (y - \bar{y})^2$$

Total sum of square (More error)

Residual should be minimum; zero mean & constant variance

HL	y	(Mean)
	$\bar{y} = 72$	
452	55	17
102	63	9
170	71	1
185	84	-12
163	59	13
171	68	4
179	75	-3
190	89	-17
180	77	-5
181	80	-8



### Reduction in Error

Total Error:

$R^2 = \frac{\text{Reduction in error}}{\text{Total Error}}$

Total Error

$$R^2 = \frac{TSS - RSS}{TSS}$$

wt.

$- R^2$  gives idea how model is good.

Greater  $R^2 \rightarrow$  Good model

④ MLR - Multiple Linear Equation (X power is 1 and all X are additive)

$$y = \beta_0 + \beta_1 x$$

$$\frac{\partial \text{RSS}}{\partial \beta_0} = 0 \quad &$$

RSS

1000

↓

$$\beta_0 = \bar{y} - \beta_1 \bar{x}$$

$$\frac{\partial \text{RSS}}{\partial \beta_1} = 0$$

2000 - 1 predictor (ht)

$$\beta_1 = \frac{\sum (x - \bar{x}) * \sum (y - \bar{y})}{\sum (x - \bar{x})^2}$$

100 - 2

↓

50 - 3

↓

45

43

↓

0

$$\beta_0 = \bar{y} - \beta_1 \bar{x}$$

$$wt = \beta_0 + \beta_1 * ht + \beta_2 fat$$

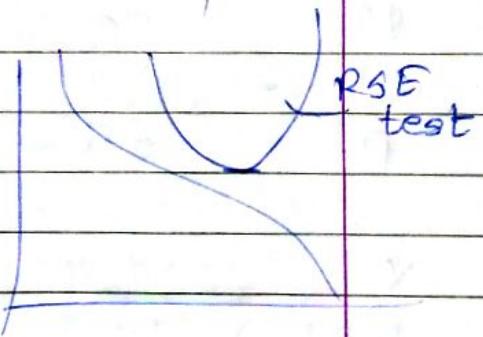
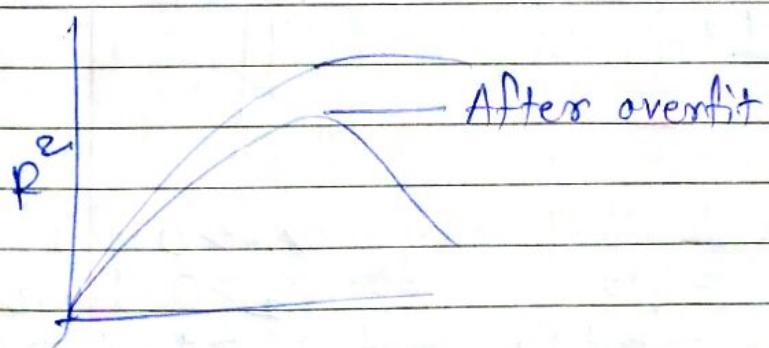
$$wt = \beta_0 + \beta_1 * ht + \beta_2 fat + \beta_3 boned?$$

RSE - Residual Standard Error.

$$RSE = \sqrt{\frac{\text{RSS}}{n-p-1}}$$

n = no. of rows

p = no. of predictors



Residual Std. Error.

RSE calculation

$$\frac{100}{10-1-1}$$

$$\frac{50}{10-2-1}$$

$$\frac{45}{10-3-1}$$

Yours	Date _____
Page No. _____	_____

$$\frac{43}{10-4-1}$$

- ① - null hypothesis independent
- alternative - u dependent

② Normality test -

Curve is normal - null hyp.

Curve is not normal - alternat. hyp.

③ Col: Summary, visualization, test of hypothesis  
P < 0.05

for project -

① 4-X Univariate analysis for all X and Y  
Search for outliers.

② 4-X Summary and visualization will give  
apparent conclusion

③ Test of hypothesis -  $p < 0.05$  dependent  
 $p \geq 0.05$  independent

Variable  
will be

28/02/18

X is categorical

Location Salary

$$\textcircled{1} \quad \begin{array}{c|c} X & Y \\ \hline P & 30 \end{array} \quad y = \beta_0 + \beta_1 X \quad X \text{ is dummy variable}$$

$$M \quad 35 \quad \text{Baseline} \quad y = \beta_0 \Rightarrow X \rightarrow 0 \text{ Pune}$$

$$P \quad 70 \quad y = \beta_0 + \beta_1 \Rightarrow X \rightarrow 1 \text{ Mumbai}$$

$$M \quad 65$$

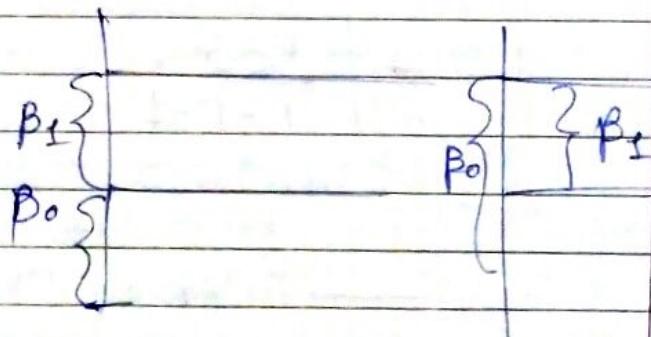
$$P \quad 85$$

$$M \quad 90$$

$$M \quad 70$$

$$M \quad 50$$

$$P \quad 65$$



$$\textcircled{2} \quad \begin{array}{c|c} X & Y \\ \hline P & 30 \end{array} \quad y = \beta_0 \Rightarrow X_1 \rightarrow 0, X_2 \rightarrow 0 \text{ Pune}$$

$$M \quad 35 \quad y = \beta_0 + \beta_1 \Rightarrow X_1 \rightarrow 1, X_2 \rightarrow 0$$

$$B \quad 70 \quad y = \beta_0 + \beta_2 \Rightarrow X_1 \rightarrow 0, X_2 \rightarrow 1$$

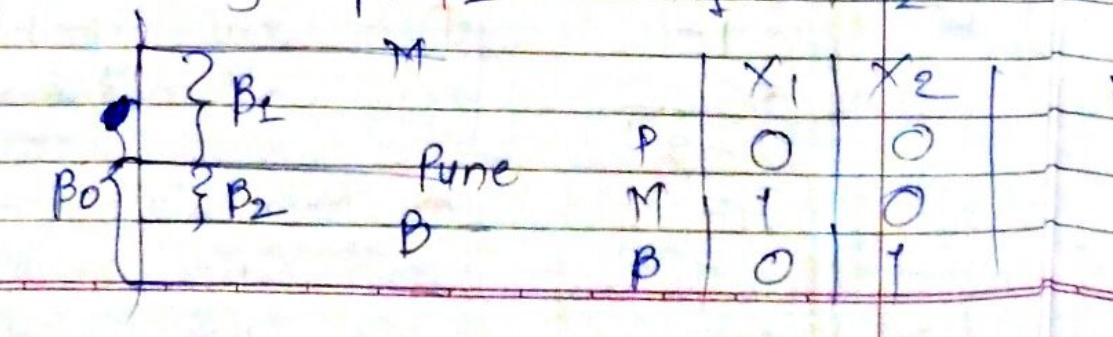
$$B \quad 65$$

$$M \quad 85$$

$$M \quad 90$$

$$P \quad 70$$

$$B \quad 50$$



$\beta_0, \beta_1, \beta_2$  eq?

$\beta_0 \rightarrow$  Pune, female.

$\beta_1 \rightarrow M, F$ .

$\beta_0 + \beta_1 \rightarrow B, F$

$\beta_0 + \beta_2 \rightarrow P, M$

$\beta_0 + \beta_1 + \beta_2 \rightarrow$  Mumbai, M

$\beta_0 + \beta_1 + \beta_2 \rightarrow B, M$

salary

B.F

P,F

M,F

$\beta_0$

$\sum \frac{2}{3} \beta_2$

$\beta_1$

Exp<sup>r</sup>.

#### ④ Polynomial Regression.

In MLR we have eq<sup>n</sup>.

$$y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \beta_3 x_3 + \dots$$

where  $x$  have power of 1 and all terms are additive.

$$y = \beta_0 + \beta_1 x_1 + \beta_2 x_1^2 + \beta_3 x_1^3 \dots$$

polynomial regression. powers are  $> 1$

⑤

$$y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \beta_3 \underline{x_1 x_2}$$

Interaction - eg. assembly line example where  $x_1$  &  $x_2$  together value add.

no changes new assembly line without workers.

in production new workers without assembly line

- new workers on new assembly will increase production of company

first ⑤ models minimizes RSS.

RSS - Residual Sum of squares.

$\downarrow$

Irreducible	$\epsilon$
Reducible	$(f - \hat{f})$
- Bias	
- Variance	

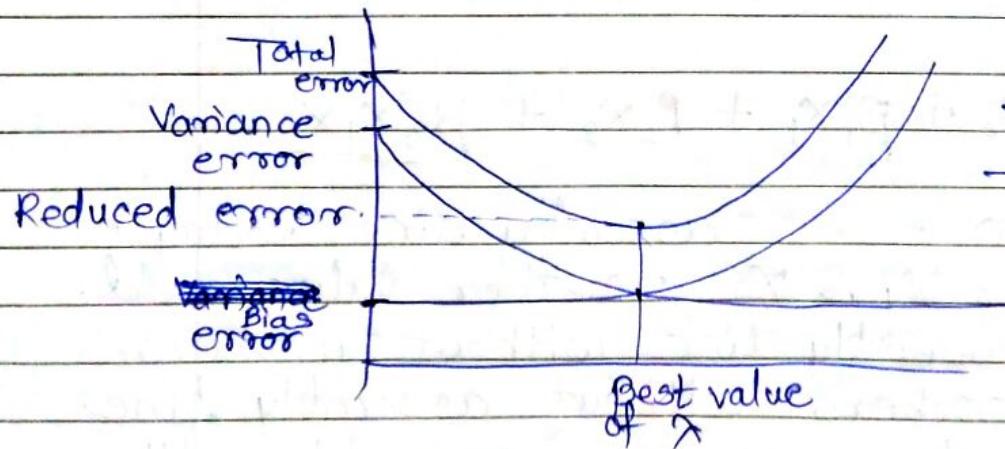
⑥ Ridge Regression } Reduce Variance from  
⑦ Lasso model, such models  
are known as regularization methods.

- RR  $\rightarrow \min_{\beta} [RSS + \lambda (\beta_1^2 + \beta_2^2 + \beta_3^2 + \dots + \beta_p^2)]$   
Tuning parameter  $L_2$  norm

- Lasso  $\rightarrow [RSS + \lambda (|\beta_1| + |\beta_2| + \dots + |\beta_p|)]$   
 $L_1$  norm

-  $L_p$  norm  $\Rightarrow |\beta_1|^p + |\beta_2|^p + \dots + |\beta_p|^p$

-  $\lambda$  is a tuning parameter.



-  $\epsilon$  is negligible  
- When variance error is more in any model this techniques are useful.

- If no. of parameters are more variance errors is more.

is not connected with  $\lambda$

- Empirical model. (multiple readings are taken total & error)

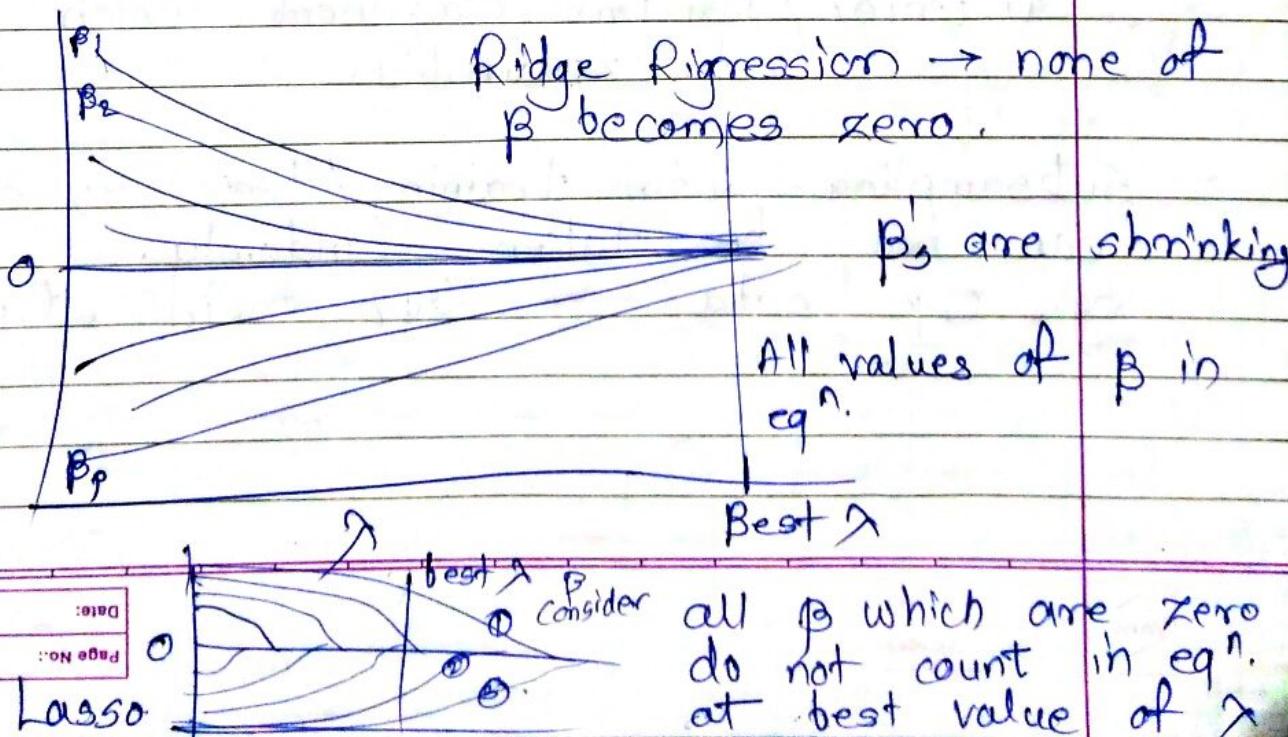
$G\beta_0 \rightarrow \uparrow$

$\beta_1 \downarrow \beta_2 \downarrow \beta_3 \downarrow \beta_4 \downarrow \dots \beta_p \downarrow$

2  
0.9  
0.2  
 $10^{-1}$   
 $10^{-2}$   
 $10^{-4}$   
 $10^{-6}$

- $\beta$ 's are shrinking.
- Ridge Regression & Lasso are shrinkage method
- Read RR & Lasso Theory.

- 1) Create a vector of  $\lambda$  values, very small to very large
- 2) for each  $\lambda$ 
  - a) Form the model. (we obtain  $\beta_0, \beta_1, \dots, \beta_p$  by minimizing eq<sup>n</sup>)
  - b) Test the model.  $\rightarrow \hat{y} - y \rightarrow$  error  
continue for all  $\lambda$
- 3) Plot error vs.  $\lambda$
- 4) find best  $\lambda$  (minimum error)
- 5) find  $\beta$  corresponding to best  $\lambda \rightarrow$  (Best model)



Model	Error	# of predictors
-------	-------	-----------------

MLR	10126	20
Ridge	9018	20
Lasso	9216	7

Choosing of model depends on Business domain  
 eg. card fraud - error must be minimum  
 Insurance, policy - error can be tolerated

Ridge & Lasso gives improvement in error.

⑧ Time series.  $y$  vs  $T$  for forecasting

Finding one or more value in future.

Time series has four parts.

1) Trend - upward / downward.

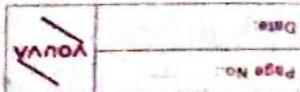
2) Season - regular pattern. eg. electricity use monthly, weekly, quarterly, < 1yr

3) Cycle - > 1yr trend, period not constant  
 needs lots of data around 25-30yr

4) Error / Random component which is not explainable.

- Subsampling into training / test is complicated,  
 can not be taken randomly.

eg. 5yr data  $\rightarrow$  4yr train, 1 yr test.



Do from ppt. - Time series.

- static & dynamic time series.

- Trend analysis. - should not contain season

- Decomposition.

Multiplicative

TS.

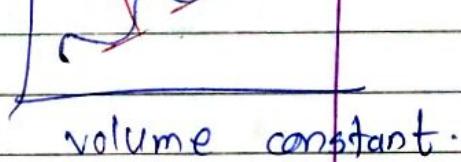


volume increasing

$$y = \text{trend} * \text{season} * \text{error}$$

season & trend are  
interactive - dependent

Additive  
timeseries.

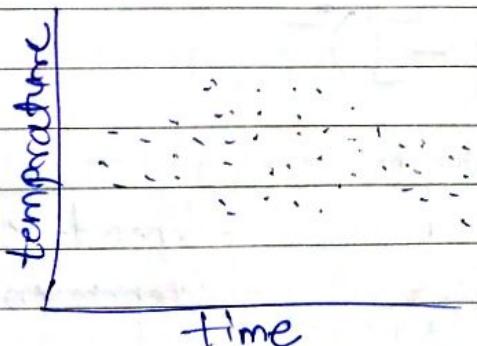


$$y = \text{trend} + \text{season} + \text{error}$$

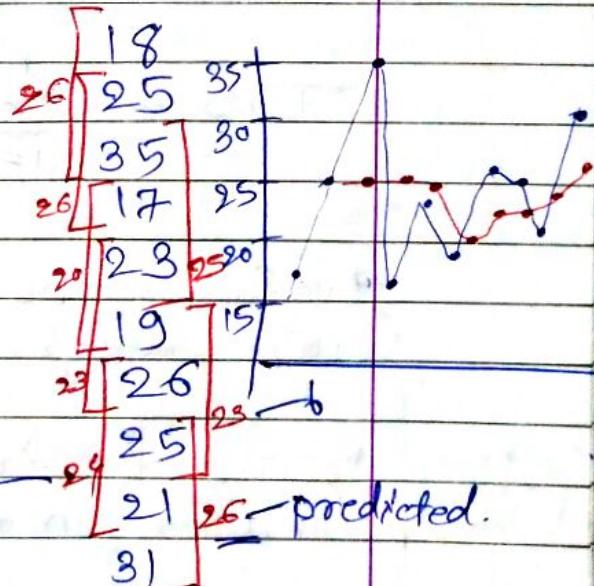
trend & season  
are independent

- Moving Average.

eg. Temperature data.



Smoothing of curve -  
define span, eg. 3.



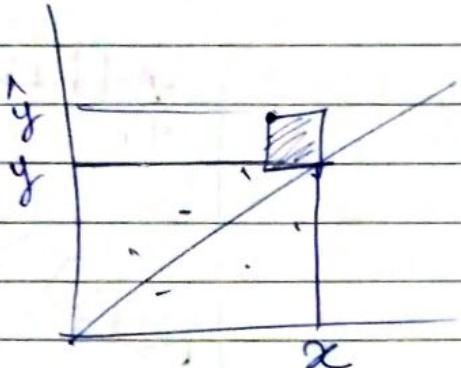
26 predicted

Exponential smoothing - coefficient reduces exponentially

$$RSS = \sum_{i=1}^n (y_i - \hat{y}_i)^2$$

Actual      Predicted.

$(y_i - \hat{y}_i)$  is residual



Least square Estimate

$$\beta_0 \quad \beta_1$$

$$37.28 \quad -5.3$$

$$\begin{array}{ccc}
 41.1 & -6.4 & - \\
 37.28 & -5.3 & - \\
 33.45 & -4.2 & -
 \end{array}
 \left. \begin{array}{c} \\ \\ \end{array} \right\} \text{Confidence interval.}$$

plot  $(x, y)$

$y \sim X$  convention

$R^2 = \frac{\text{Reduction in error}}{\text{Total Error}}$

$$TSS = \sum_{i=1}^n (y_i - \bar{y})^2$$

$\varnothing$  vs  $\mathcal{G}$  slope zero independent null  
slope not zero dependent alternative

Degree of freedom decreases as known variables increases

$$RSE = \sqrt{\frac{RSS}{n-p-1}}$$

Young	Date
Older	No

↑ denominator is degree of freedom

RSS

Summary

plot(modelname, 1)

Jn (mpg ~ wt + hp, cars)

Model	RSS	RSE	R <sup>2</sup>	R <sup>2</sup> adj.
wt				
wt + hp				

1. ✓ SLR  $\leftarrow \text{Jm}(\text{mpg} \sim \text{wt})$

2. ✓ MLR  $\leftarrow \text{Jm}(\text{mpg} \sim \text{wt} + \text{hp}, \text{cars})$

3. Categorical.  $C \rightarrow Q$

cars \$ cyl  $\leftarrow \text{as.factor(cars$cyl)}$   
model2  $\leftarrow \text{Jm}(\text{mpg} \sim \text{cyl}, \text{cars})$

4. (2+3) ✓ Jm (mpg ~ wt + hp + cyl, cars)

5. + Interaction ✓ Jm (mpg ~ wt + hp + cyl + wt\*hp, cars)  
term\*

{ 6. Jm (mpg ~ hp, cars)

Identity  
terms only  
allow power

7. Jm (mpg ~ hp + I(hp^2), cars)

one power  
term (I)

8. Jm (mpg ~ poly(hp, 3), cars)  $y = \beta_0 + \beta_1 h p + \beta_2 h p^2 + \beta_3 h p^3$

(5+ poly) 9. Jm (mpg ~ wt + hp + cyl + wt\*hp + I(hp^2), cars)

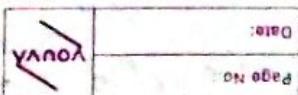
RSS } Training error.  
R<sup>2</sup>

- min - RSE - error in model.
- max - R<sup>2</sup> adj. - Reduction in error
- anova(carsm1, carsm2, ...)
- predict(carslm2, data.frame(wt = c(2.5, 3.5),  
ht = c(100, 200))  
)
- Stepwise model.

### Anova

F stat p value will tell how good model  
 $p < 0.05$ , model value should be less than 0.05.

In anova, models should be progressively increasing, then only it can give comparative evaluation.



# Stepwise Regression -

30/08/18

$y \quad x_1 \quad x_2 \quad x_3 \quad x_4 \quad x_5$

- First take SLR of  $y \quad x_1$
- min p value gives the most dependent variable  $x$   $y \quad x_2$   $\boxed{y \quad x_3}$  — 1<sup>st</sup> min. found.
- Now pair  $(y, x_3)$  and remaining.

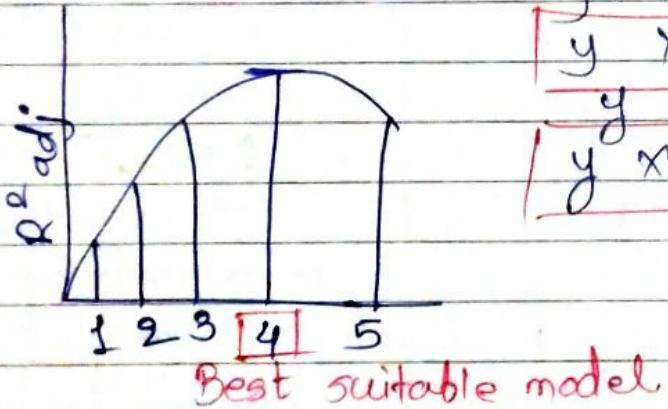
## Stepwise Regression - Forward

Backward

Mixed.

$y \quad x_3 \quad x_1$   
 $y \quad x_3 + x_2$  — 2<sup>nd</sup> min.  
 $y \quad x_3 \quad x_4$   
 $y \quad x_3 \quad x_5$

Model having R adj - maximum is good.



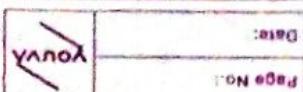
$y \quad x_3 \quad x_2 \quad x_1$   
 $y \quad x_3 \quad x_2 \quad x_4$   
 $y \quad x_3 + x_2 + x_5$  — 3<sup>rd</sup> min  
 $y \quad x_3 \quad x_2 \quad x_5 \quad x_1$  — 4<sup>th</sup> min  
 $y \quad x_3 \quad x_2 \quad x_5 \quad x_1 \quad x_4$

$$5c_1 \quad 5c_2 \quad 5c_3 \quad 5c_4 \quad 5c_5 \\ 5 \quad 10 \quad 10 \quad 5 \quad 1 = 31$$

R library for stepwise regression — Search

library (leaps)  
regsubsets

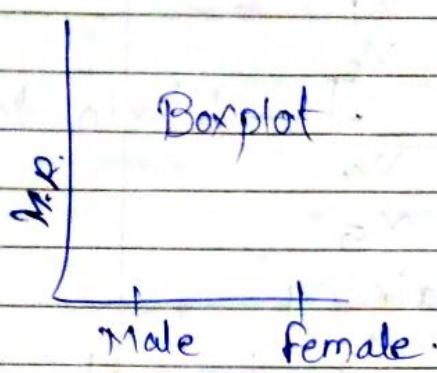
Variable Reduction and Model Selection.R



# Project -

## Regression

$y = \text{monthly rate}$  (-Attrition)



test - <sup>t-test</sup> Independent sample,  
more than two Independent samples

paired t-test:

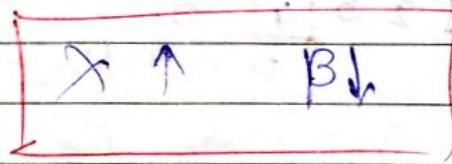
training before - after  
weight loss before - after

## Data preparation

- Decide Q and C
- Datatype check, if required to change to factors.

$$\begin{matrix} \gamma_1 & \gamma_2 & \gamma_3 \\ \beta_0 & \beta_0' & \beta_0'' \end{matrix} \quad \begin{matrix} \dots & \dots \\ \beta_{100} & \beta_{100}' \end{matrix}$$

$$\beta_{10}$$



$$(x_1, y_1)$$

$$(x, y)$$

$$\frac{y - y_2}{y_1 - y_2} = \frac{x - x_2}{x_1 - x_2}$$

$$(x_2, y_2)$$

To find  $\beta$  value at certain  $x_a$ ,  
predict function is used.

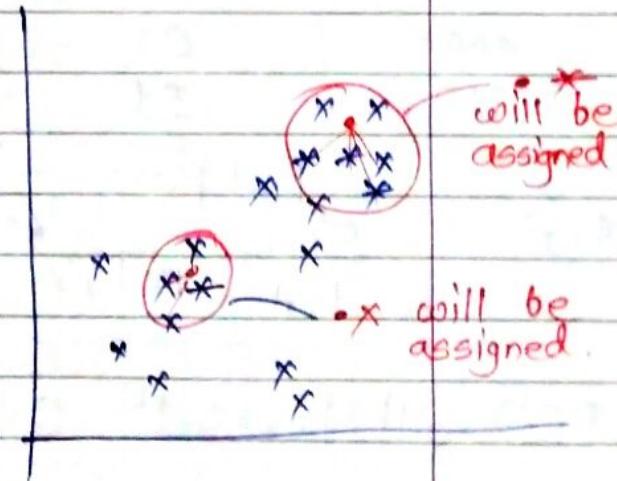
## Classification -

- Prediction of quantitative data is regression.
- " qualitative " classification
- Classification -
  - Characteristics of animal. - will classify.
  - Mail spam or not spam
  - Employee rating.
- Based on training data classifier will classify the new data.

① KNN - K nearest neighbour.

- K is user given input
- K decides of neighbouring sample points
- Mode - highest frequency.

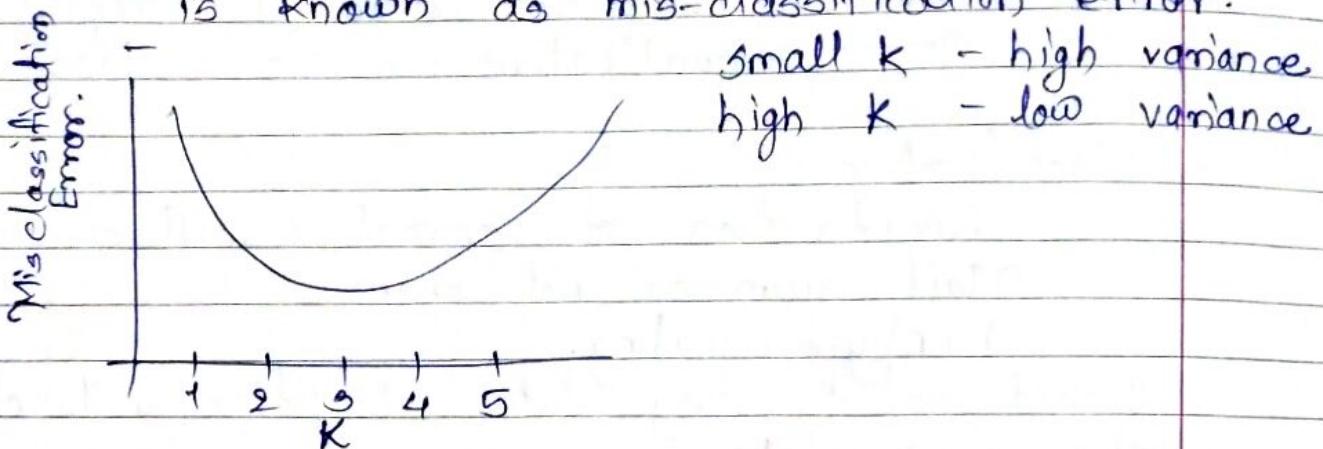
Age	Exp.	Salary
H		
L		
H		
L		
M		
H		
H		



- Here, we can not make training & testing separately. They go in parallel.

- Misclassification error

- Any point which is not predicted correctly is known as mis-classification error.



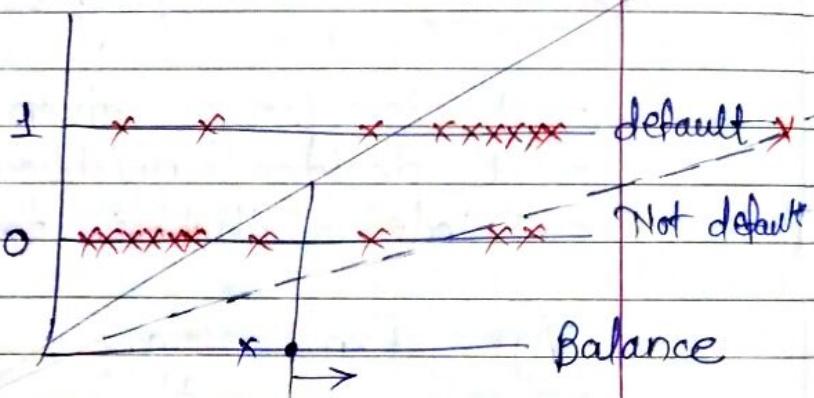
small  $K$  - high variance

high  $K$  - low variance

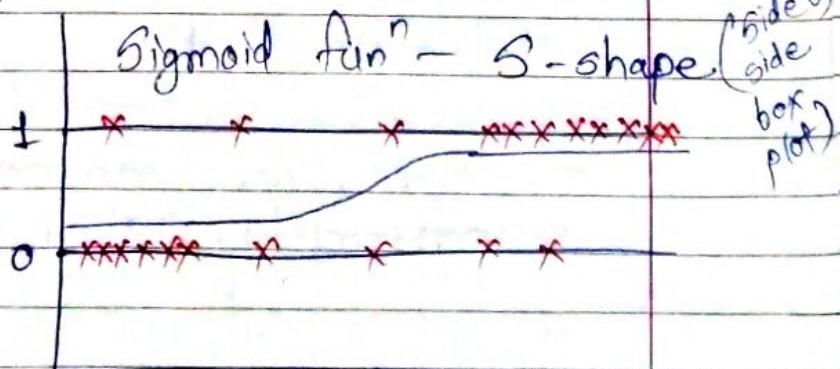
## ② Logistic Regression - (Only for two classes)

$X$	$y$
Credit Card Bill Amt.	Person is default/not
7,000	0
40,000	1
3,000	0
50,000	1

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



- Outliers impact the regression line hence linear regression is not useful in this scenario.



Large  $X - 1$  den=0 - Sigmoid function.

Small  $X - 0$  den= $\infty$

$$P(\text{Balance} = \text{Default}) = P(y/x = \text{Default}) = \frac{e^{-(\beta_0 + \beta_1 x)}}{e^{-(\beta_0 + \beta_1 x)} + 1} \quad \text{OR} \quad \frac{e^{\beta_0 + \beta_1 x}}{1 + e^{\beta_0 + \beta_1 x}}$$

X Test	$\frac{y}{x}$	< 0.5 cutoff	< 0.8	< 0.9
15,000	0.6	1	0	1
11,000	0.4	0	0	1
7,000	0.35	0	0	0
20,000	0.63	1	0	1
25,000	0.71	1	0	1
12,000	0.45	0	0	1

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

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

$$\frac{P(\text{Success})}{P(\text{Failure})} = \frac{P(x)}{1 - P(x)} \cdot e^{\beta_0 + \beta_1 x} = \text{odds}.$$

$$\text{odds} = \frac{9}{4} = 2.25 = \frac{9/13}{4/13}$$

$$P(\text{Success}) = \frac{P(x)}{1 - P(x)}$$

$$\frac{9}{4} = \frac{P(x)}{1 - P(x)} \therefore P(x) = \frac{9}{13}$$

If we increment  $x$  by 1, odds will increase

$$\text{odds}' = \text{odds} * e^{\beta_1}$$

$$\begin{aligned} \text{odds}' &= e^{\beta_0 + \beta_1(x+1)} \\ &= e^{\beta_0 + x\beta_1 + \beta_1} \\ &= e^{\beta_0 + x} * e^{\beta_1} \end{aligned}$$

$$\text{odds} = e^{\beta_0 + \beta_1 x}$$

ranges between 0 to 1

$$\log(\text{odds}) = \beta_0 + \beta_1 x$$

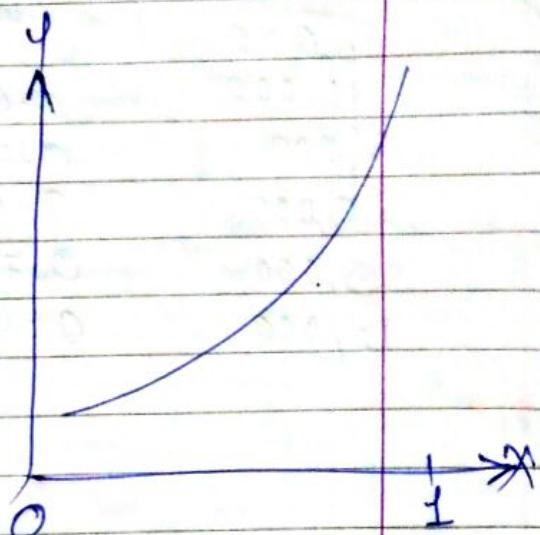
$$\log\left(\frac{P(x)}{1-P(x)}\right) = \beta_0 + \beta_1 x$$

diff? Names

- logit

- logistic Regression

$$y = \beta_0 + \beta_1 x$$



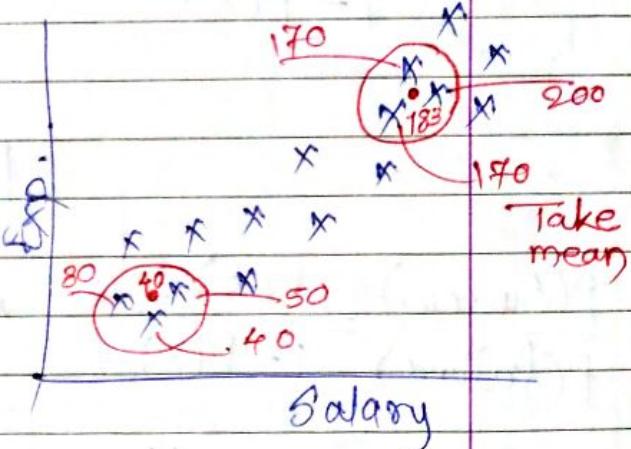
### ③ KNN for regression —

Exp Age Sal(k) - Take k=3

32 55 200

4 25 60

80 55 ?



c1pg — KNN

c1 — Logistic

c1/pg — Decision Tree

— Classification Tree

— Regression Tree

— Bagging

— Random Forest

c1 — Naive Bayes Classification

SVN ANN

c1

c1/Rg

## ④ Naive Bayes Classification

### Conditional Probability

Statistical Conditional Probability

$$P(A \cap B) = P(A) * P(A|B) \quad \text{--- (1) Bayesian} \\ = P(B) * P(B|A) \quad \text{--- (2).}$$

e.g. Effect of smoking on Cancer.

Smoke	Non Smoke	Cancer	
100	100	200	Yes
900	600	800	No.

Group formed already

Prospective study  $P(\text{Cancer}/\text{Smoke})$   $P(\text{Person will have cancer when smoke})$

Retrospective study  $P(\text{Smoke}/\text{cancer})$   $P(\text{Person smokes and has got cancer})$   
 Study continues for years.

Using Bay's Theorem.

$$\therefore \text{from eq' (2)} \quad P(A|B) = \frac{P(A) * P(A|B)}{P(B)}$$

$$P(\text{Cancer}/\text{Smoke}) = \frac{P(\text{Smoke})}{P(\text{Cancer})} * P(\text{Smoking})$$

e.g. Heart attack prediction.

Ex: BMI	High	BP	high	Cholesterol	Y
Y	N	H/M/L	<del>Y/N</del>	<del>Y/N</del>	Heart attack
Y	H	Y	Y	Y	1
Y	M	N	N	N	0
N	L	Y	Y	Y	1
Y	L	N	N	Y	0

$x_4 \quad x_3 \quad x_2 \quad x_1 \quad Y$

$$P(B/A) = \frac{P(B) * P(A/B)}{P(A)}$$

$$P(HA/x) = \frac{P(HA) * P(x/Ha)}{P(x)}$$

Prior probability

Posterior Probability.

$$P(X)_{HA=y} = P(Ex=N \cap BMI=H \cap BP=y \cap CH=N) / HA=y$$

Assumption - Naive (simplistic)

- All  $x_i$ 's are independent of each other.

$$\therefore P(X) = P(Ex=N) * P(BMI=H) * P(BP=y) * P(CH=N)$$

- All probabilities can be calculated from respective columns.

$$\therefore P(X|HA=y) = P(Ex=N|HA=y) * P(BMI=H|HA=y) * P(BP=y|HA=y) * P(CH=N|HA=y)$$

Laplace Correction

$$H \quad M \quad L \\ 160 \quad 40 \quad 0 = 200$$

$$161 \quad 41 \quad 1 = 203$$

$$\frac{161}{203} \quad \frac{41}{203} \quad \frac{1}{203}$$

- In Naive Bayes it is recommended to convert features to qualitative data.

e.g. If it is not possible to convert data to qualitative

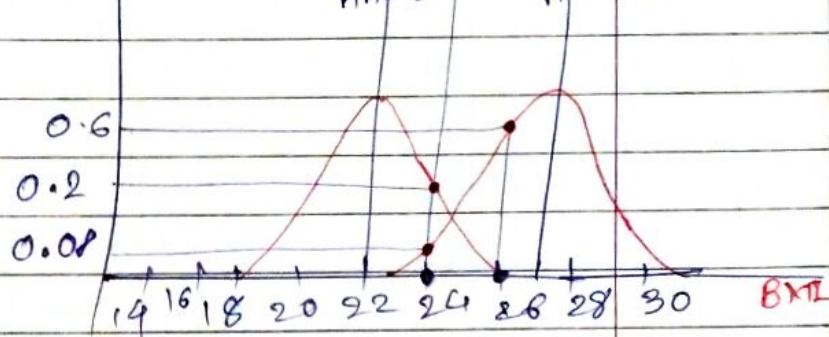
BMI	HA	Segmented 0's & 1's.	
22	0	0	5
23	0	2	5
25	1	2	3
27	0	2	7
30	1	2	1
32	1	mean = 2	$\sigma = 2$
21	0	$\sigma = 2$	

Assume they are normally distributed.

$$P(24 \mid HA=0) = 0.08$$

$$P(24 \mid HA=1) = 0.2$$

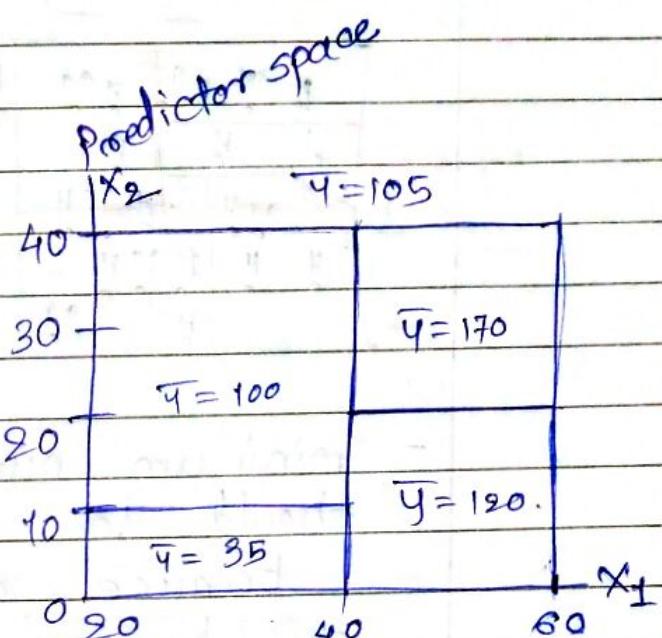
$$P(26 \mid HA=0) = 0.6$$



## Decision Trees

~~Classification Tree~~

	$X_2$	$X_1$	$y$
	Age	Exp <sup>n</sup> .	Salary(k)
22	1	30	
26	3	40	
55	30	180	
50	27	160	
38	15	100	
42	19	120	
35	14	? (100)	

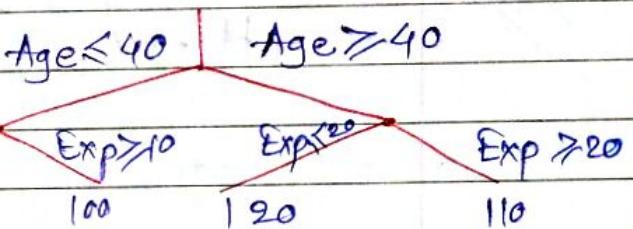


- If there's no function formulation to find relationship we take mean of data.

$$\bar{y} = 105$$

- Binary split

- Branches, <sup>internal</sup> nodes, terminal node (leaf)
- Easy to understand and implement



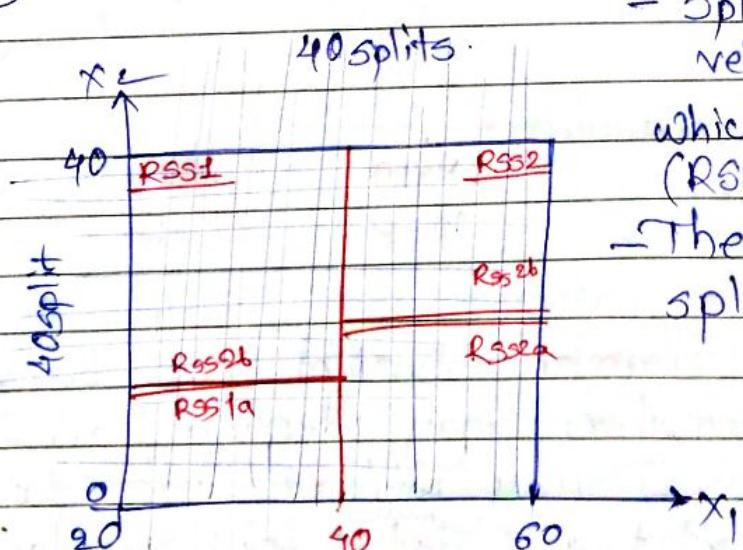
- Split can come hr. or ve<sup>n</sup>. where (RSS)<sub>sum</sub> is min.

which will give  $(RSS_1 + RSS_2)$

- The less error split will be selected

Hierarchy  
nested if-else

Split can take any order vertical horizontal



- Stopping criteria - atleast minimum number of points should be included (eg. 5 pt. in box)

$$\frac{\max}{\min} \text{ ratio } RSS = 5$$

- When salary is categorical ( $H, M, L$ )

<del>H H H M H M M H H H H H H H H H H M H H H H M H</del>	$H=6$ $M=4$	Mode $H=14$ $M=7$	Misclassification = 7
<del>H H H H H H H H H H H H H H H</del>	$H=8$ $M=3$	- Node purity more important	

- minimum number of point in a region should be slightly high eg. 10
- because a single or couple of pts. can be ambiguously classified.

~~Regression  
Tree~~

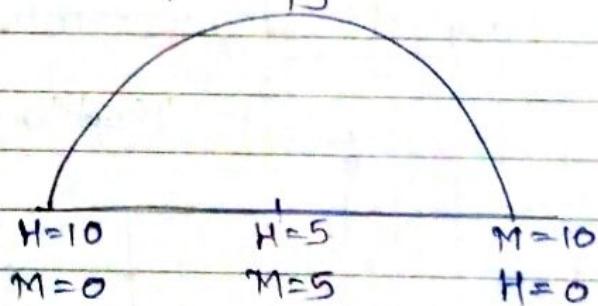
$$\begin{matrix} M & - & 4 \\ H & - & 10 \end{matrix} \rightarrow \begin{matrix} M & - & 2 & 2 \\ H & - & 5 & 5 \end{matrix}$$

or

$$\begin{matrix} M & - & 4 & 0 \\ H & - & 7 & 3 \end{matrix}$$

This split is preferable as it gives pure node.

Entropy Curve



- Cross Entropy
- Gini Index

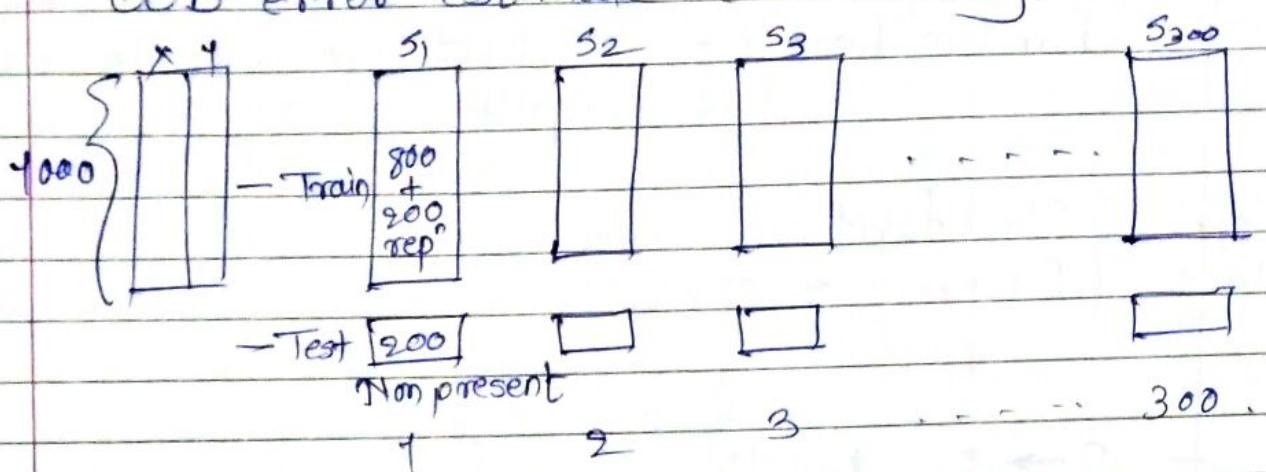
- Entropy is measure of disorder.

- ~~Bagging~~
- Tree performance might not be good
  - Tree formation may depend on ~~train~~ train data, which may mislead test data result
  - Using bootstrap (multiple sample with replacement) is used to improve accuracy
  - Centre limit theorem - take mean of multiple predicted values.

Bootstrap  
aggregation  
(Bag of  
trees)

Bagging

## OOB error estimate (Out of Bag)



25 predictors - { strong }  
 3 medium } Predictors  
 21 Weak }

- Because of 1 strong & 3 medium predictors all 300 trees would look alike co-related trees.
- Central limit theorem is applicable for spreaded data.
- Bagging will not good result in such scenarios. Hence predictors should be equally important.

Random Forest

To solve above problem, take predictor subset randomly.

- Out of 25 predictor split subset of 5 randomly
- Each tree is looking different hence known as Random forest

P predictors - classification  $\sqrt{P}$   
 (Recommendation) - regression  $P/3$

Strong Predictor =  $\sum \text{RGS}$   
 no. of occurrence

Multiple Tree { Bagging - All predictors in each split  
 Random Forest - Subset of p. in each split

Single Tree { Classification Tree  
 Regression Tree

- $y \rightarrow c$  logistic regression
- $\mu \pm 3\sigma$  } Outlier  
 $1.5 \text{ IQR}$  }
- Categorical data has no outlier.
- Data Prep. Inference
- EDA + Visualization + ML.
- Matrix plot will give good relationship between X's

Income & and other X  
 Attrition & and other X

## Unsupervised Learning - There's no 'Y'

Employee data - each employee is object.

F.Name L.Name Skill Exp Salary Loc Gender

Emp 1

Emp 2

⋮

Grouping can be done on attributes.

- Clustering is done on objects.
  - Similar type of objects form clusters.
  - eg. similar employee, customer, product
  - eg. Library is example of having clusters of subjects.
  - Clusters are formed based on similarities.
  - Similarities needs to be done based on all attributes.
  - dissimilarity varies between 0 & 1
- 0 → similar.  
1 → dissimilar.

Dissimilarity co-efficient calculation -

	1	2	3	
1	0	$d_{12}$	$d_{13}$	$d_{12} = d_{21}$
2	$d_{21}$	0	$d_{23}$	$d_{1j} = d_{ji}$
3	$d_{31}$	$d_{23}$	0	

↑ Useful portion.

- cont'd

Nominal data - There is no order and it is a categorical data.

	Tech.	Loc?	Responsibility
Emp 1	Oracle	Mumbai	Coding
2	Java	Pune	Testing
3	Oracle	UK	Coding
4	R	Pune	Coding

$$\text{Emp 1 \& Emp 2 } D_{12} = \frac{3}{3} = 1$$

	1	2	3	4
1	0	1		
2	1	0		
3	0.33	1	0	
4	0.66	0.66	0.66	0

Ordinal data -

	Ordinal Rating	Rank	Normalized Rank (N.R.)
Emp 1	0.66 G	1 Outst.	0
2	0.66 Outst.	2 VG	0.33
3	0.33 VG	3 G	0.66
4	0.66 G	4 A	1

	1	2	3	4
1	0			
2	0.66	0		
3	0.33	0.33	0	
4	0	0.66	0.33	0

Binary Data

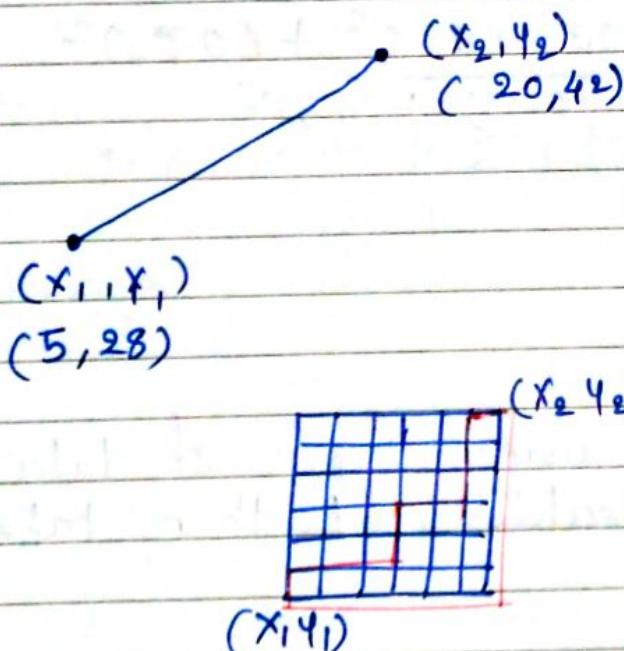
Test Case 1

	Result	High Priority	Delivered
P	1	Y	N
F	0	N	Y
F	0	Y	Y
P	1	N	N

$$T_{23} \quad \begin{matrix} & 0 & 1 \\ 0 & | & | \\ 1 & | & | \end{matrix} = \frac{1}{3}$$

Numeric Dissimilarity -

- Distance is calculated to find dissimilarity



$$d_E = \sqrt{(y_2 - y_1)^2 + (x_2 - x_1)^2}$$

$$= \sqrt{14^2 + 15^2}$$

$$(x_2, y_2) = 14\sqrt{2}$$
 Euclidian dist?

$$d_m = |y_2 - y_1| + |x_2 - x_1|$$

$$= 14 + 15 = 29$$

Manhattan dist.

45 = 0.8 Min. Max normalization - Converts value bet' 0 & 1 proportionally

$$= \frac{\text{value} - \text{min}}{\text{Max} - \text{min}}$$

$$= \frac{\text{value} - 5}{50}$$

$$36 = 0.62$$

## Ordinal data cont? - more than 2 category

As per previous

Rating	Normalised Rank	Rank	SAL	Normalized Rank
Outstn	0	1	H	$\frac{1-1}{3-1} = 0$
VG	0.33	2	M	$\frac{2-1}{3-1} = 0.5$
G	0.66	3	L	$\frac{3-1}{3-1} = 1$
A	1			

	$x_1$	$x_2$
Emp 1	$0 = 0$	$H = 0$
2	$VG = 0.33$	$M = 0.5$
3	$G = 0.66$	$H = 0$
4	$V = 0.33$	$L = 1$

$$d_{12}(\text{Emp 1, 2}) = \sqrt{(0.33 - 0)^2 + (0.5 - 0)^2}$$

$$\begin{aligned} d_{23}(\text{Emp 2, 3}) &= \sqrt{(0.33 - 0.66)^2 + (0.5 - 0)^2} \\ &= \sqrt{(0.33)^2 + (0.5)^2} \end{aligned}$$

Strategy<sup>2</sup> dist' =  $\frac{\sum \text{val}}{\text{Count}}$

Strategy<sup>3</sup> If there are many types of data first calculate separately and then take mean.

## Clustering

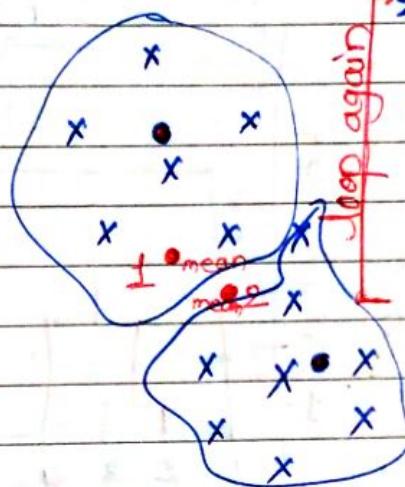
- K-means
- Hierarchical
- DB Scan

## Clusture properties -

- similar pts. in same cluster
- two clusters should be dissimilar

### ① K-means

- first mean
- Second mean

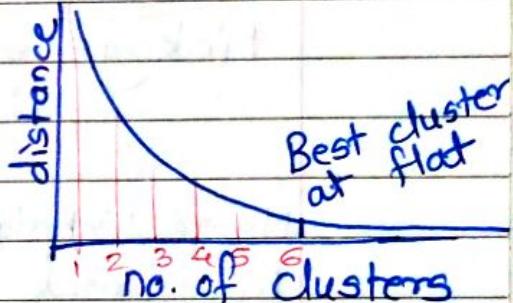


1. Randomly assign cluster number to each point.
2. Calculate mean for each cluster.
3. ReAssign cluster number to each pt. based on minimum distance from calculated mean.
4. When there is no change in assignment of pts. - stop.

Within cluster

$$\sum (\text{dist}_1)^2 + \sum (\text{dist}_2)^2$$

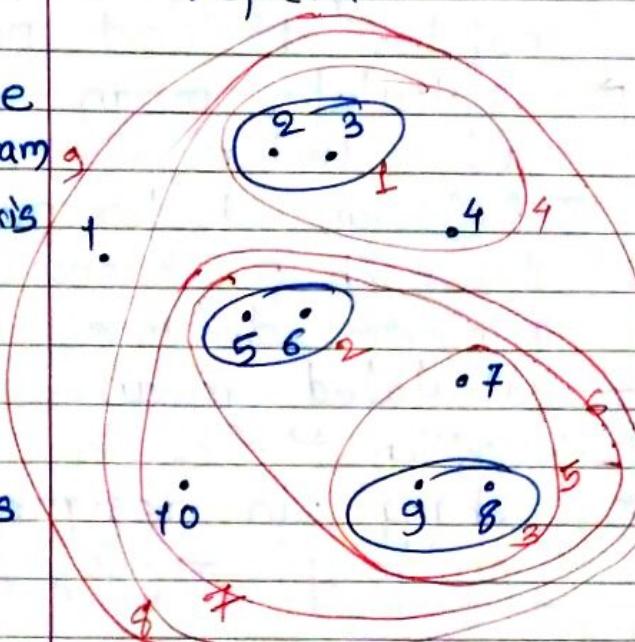
Between =  $d^2$



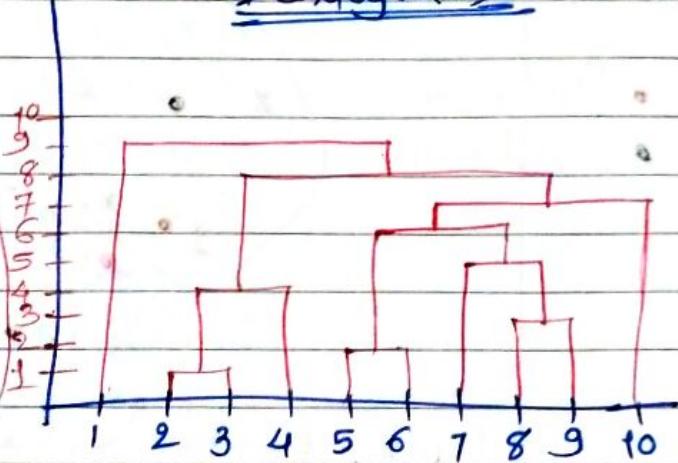
## Hierarchical -

- All pts. are clusters individually at the begining
- Close ones form one cluster and repeat.

Cut the dendrogram on Y axis for that many no. of clusters



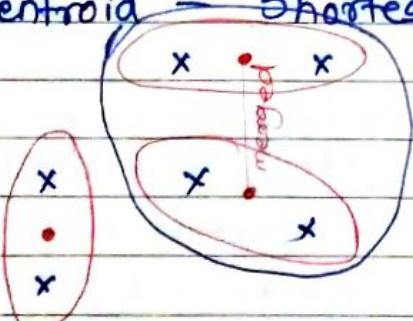
## Dendrogram



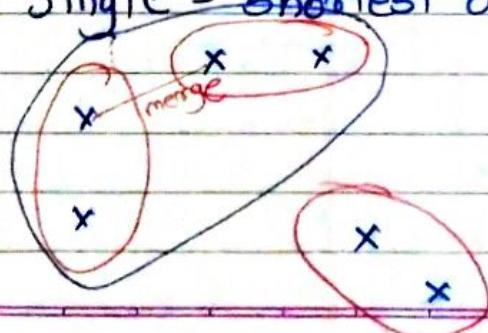
Linkage - Dist<sup>n</sup> bet<sup>n</sup> pt. & cluster  
— u — two clusters

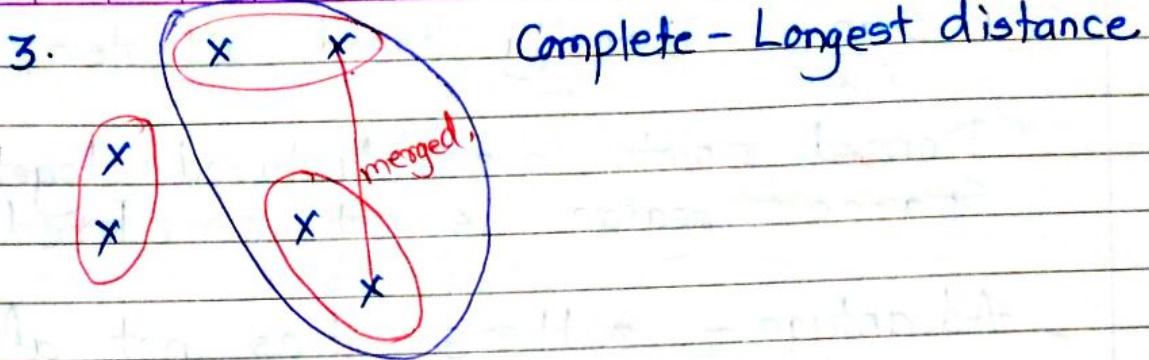
Merge Techniques 1. Centroid - Shortest centroid dist<sup>n</sup>

1. Centroid
2. Single
3. Complete
4. Average

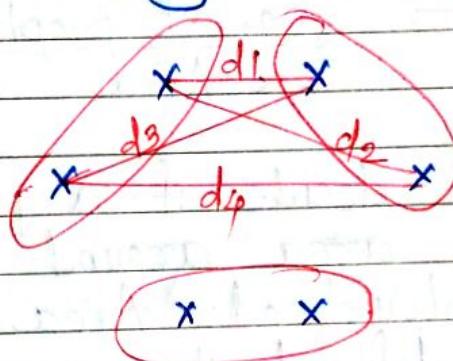


2. Single - shortest distance bet<sup>n</sup> any two pts of diff. clusters



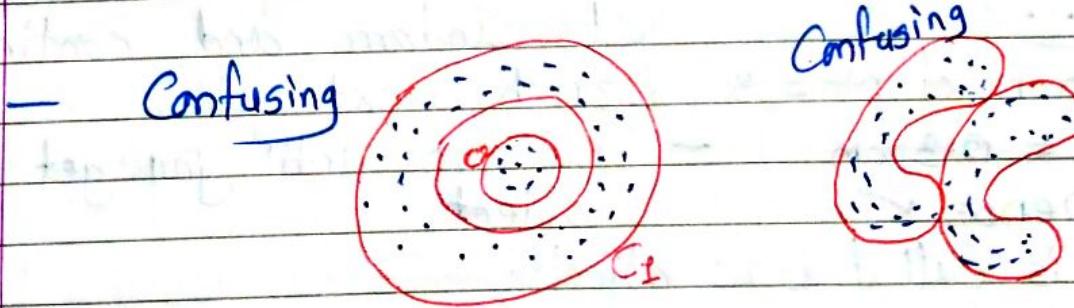


4. Average -  $(d_1 + d_2 + d_3 + d_4) / 4$  minimum



### Disadvantages —

- Spherical cluster - K means
- Hierarchical



- Convex shape clusters are formed by K means & Hierarchical

- Outliers affects these two as it forces to include every pt.

## DB Scan - Density Based Cluster.

Densed points are clustured together  
Sparse region is not considered.

Advantage - outliers does not affect DB scan. Thoes are left behind.

Population Density =  $\frac{\text{no. of people}}{\text{Area}}$

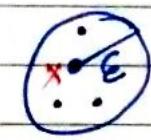
- First we need to decide if a pt. is dense.
- We consider an area around that pt.



$E$  neighbourhood Area

$E$  is defined by user

min. pts. are defined by user



min. no. of pts = 3

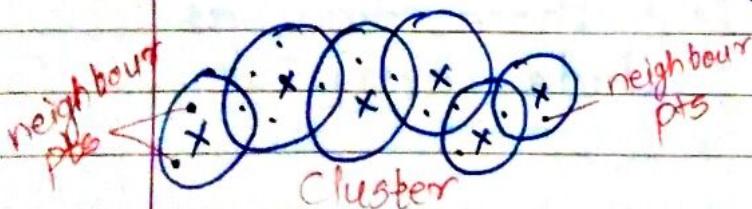
- Starting core pt. is selected at random and continue to search next pt.

$$E = 0.3 \text{ cm}$$

Dense ✓

X is called core object.

- Go on until you get non-core object
- Group together all neighbouring core points.
- Core objects and their neighbours form cluster.



## Columnwise Grouping - Assosiation

### Market Basket Analysis.

Basket 1 bread butter milk egg cheese

	bread	butter	milk	egg	cheese
2	✓	✓	✓	-	-
3	✓	✓	-	-	✓
4	✓	✓	-	✓	-
5	(-)	✓	✓	✓	-
6	✓	-	✓	✓	✓

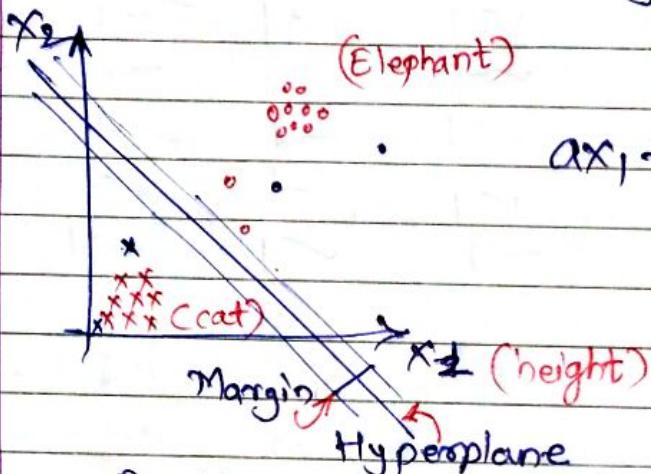
Hyperplane is used to segregate classes.

Two dimensions - line.

Three - " - plane.

$> 3$  - " - Hyperplane

(Weight)



- Farther points can be said confidently of some class.
- Pts. close to lines are less confident about class
- Try to fit thickest band inbetween edge pts. (outermost points)
- nearest pts to line should be farthest.
- SVM is mostly used for binary classification

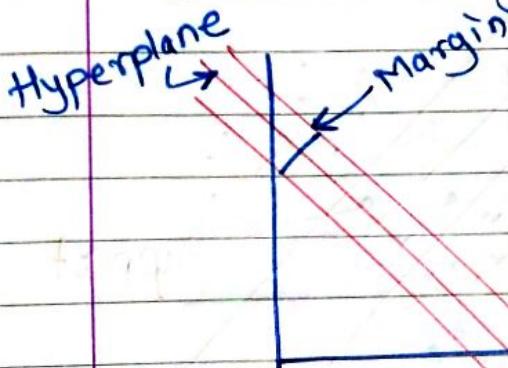
No overlap  
of data pts

Maximum margin classifier is used for perfectly linearly separable classes.

- Due to few points newly introduced margin can be varied.
- points could be messed up, there will be overlap b/w two classes eg. cat - dog
- At the expense of some misclassification majority should be classified.

$\epsilon$  - Penalty / slack penalty  
 crossed margin but not hyperplane  
 crossed margin also hyperplane

$$\begin{array}{l} 0 \leq \epsilon \leq 1 \\ \epsilon > 1 \end{array}$$



Tuning Parameter  $c$  (budget)  
 $\sum_i \epsilon \leq c$

- Budget will control how many pts. will continue to cross the borders.
- Margin decision is optimization problem.

- As value of  $c$  increases margin increases.
- Model is build using training data
- Cost  $\propto \frac{1}{\text{budget}}$

cost is less, more pts. will cross.

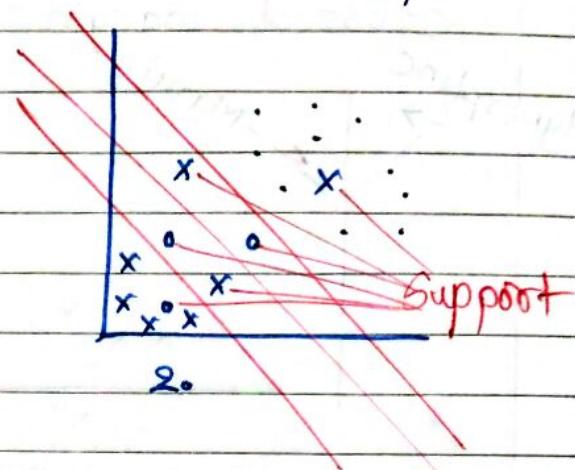
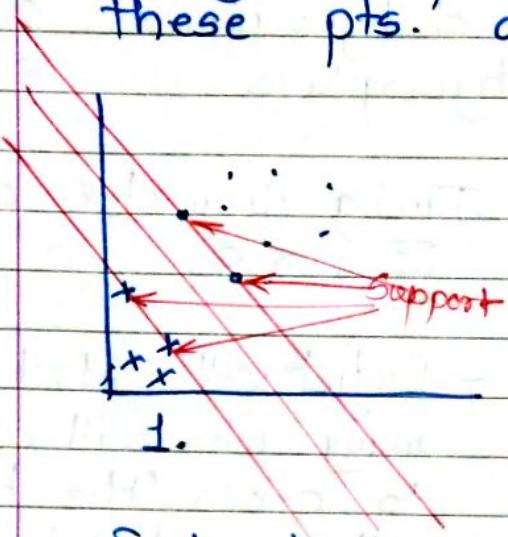
cost is more, less pts will cross.

- models are tested using multiple  $C$  values on test data and chosen best

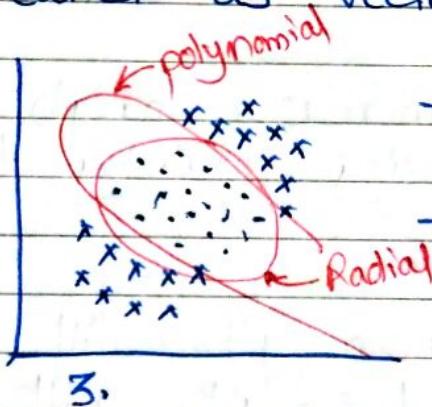
- 1. Support Vector Classifier. (eg. cat-dog)
  - overlap of data points
  - linear classifier
- 2. Maximum Margin classifier. (eg. cat-elephant)
  - Perfectly separable classes
  - No overlap pts.

- 3. Radial Classifier (PTO) / Polynomial (SVM)
    - Non linear hyperplane
    - overlapping data pts.
- Support Vector Machine also known as.

→ Margin depends on nearest pts.  
these pts. are called as support vector



- Each pt. has multiple dimensions hence called as vector.



- Such data will require non linear classifier.
- eg: radial classifier

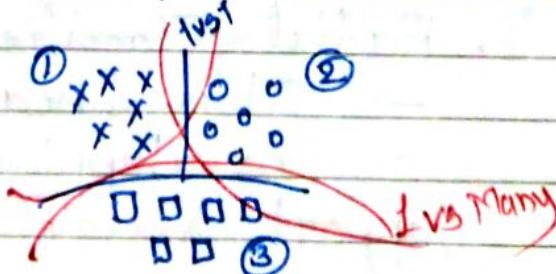
- Kernel is used to simplify solution of polynomial classifiers.

- In 1. & 2. Kernel is used linear.

3. Kernel is nonlinear  
polynomial - degree  
radial - gamma  $\gamma$

- In case of more than two classes (eg. three class)

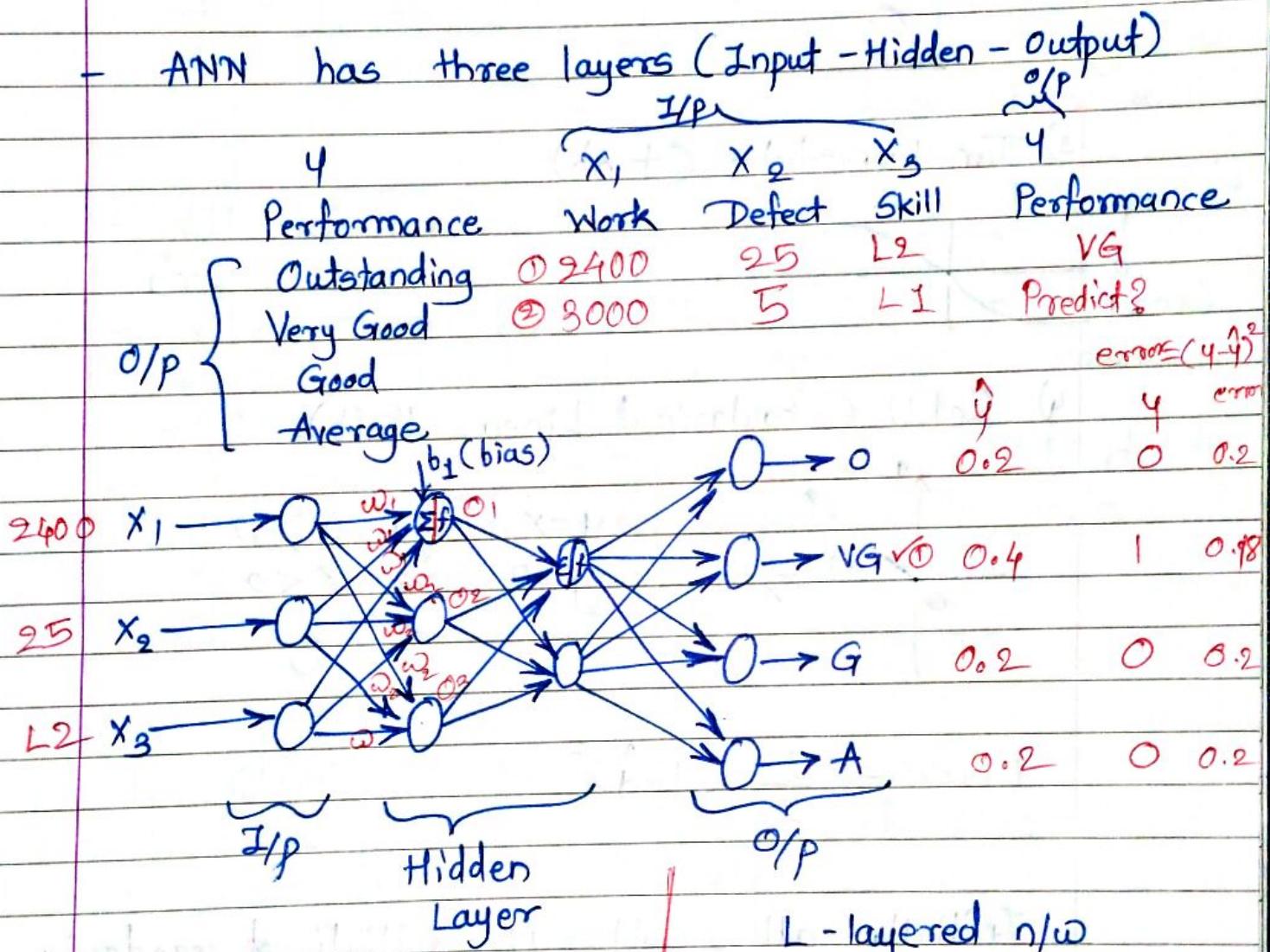
① 1 vs 1 - take two classes at once



② 1 vs many

## ANN

- Each neuron - processes signal independently
  - passes signal to each connected neuron
  - can enhance / surpass signal
- ANN has three layers (Input - Hidden - Output)



similar  
MLR  
equation

$$\sum = w_1 x_1 + w_2 x_2 + w_3 x_3$$

$f$  is a function which makes signal to non-linearity

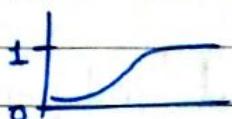
$$f(b_a + w_4 O_1 + w_5 O_2 + w_6 O_3)$$

- Outputs will be probabilities.
- $f$  is known as activation function

L - layered n/w  
 $L = \text{O/P layers} + \text{Hidden layers}$   
 $= 1 + 2 = 3 \text{ layers}$

default function

1) Sigmoid function

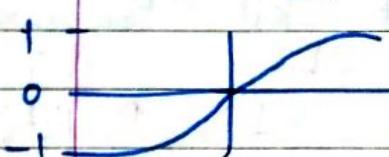


$$\frac{e^{b_0 + \omega_1 x_1 + \omega_2 x_2 + \omega_3 x_3}}{1 + e^{b_0 + \omega_1 x_1 + \omega_2 x_2 + \omega_3 x_3}}$$

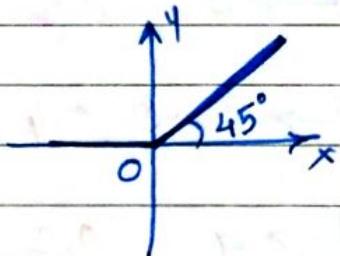
2) Gaussian function (Normal)



3) Tan-hyperbolic (tanh)



4) ReLU (Regularized Linear Unit)



$$y = x \quad \text{for } x \geq 0$$

$$y = 0 \quad \text{for } x < 0$$

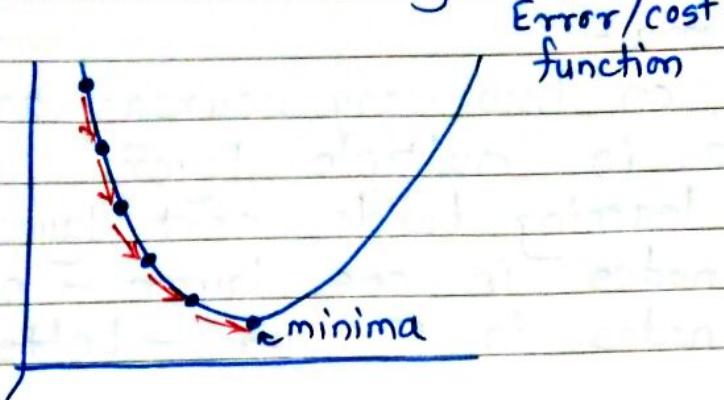
$$\text{Error} = \frac{1}{2} (\hat{y} - y)^2$$

Initially all weights are initialised randomly in between -1 to 1

Forward propagation - signal

Backward propagation - Error

**Stop**      **Go down**  
**Gradient Descent Algorithm**



- Maximum slope is calculated by partial differentiation
- it will give  $\Delta w$  to adjust current weight so that error will minimize.

- All  $w$ 's and all biases will try to minimize error.
- When there's no further change in network it is converged. It can take 1000 rows to millions rows to get train.
- Learning rate is a parameter which decides convergence steps. As large learning rate might cause overshooting and might not reach minima.



1. Classification, more than 2 classes  
no. of output nodes = no. of classes.
2. Classification, 2 classes  
no. of o/p nodes = 2
3. Regression (only summation fun', no activation)  
no. of o/p nodes = 1

### Design Criteria -

- Depending on how many neurons are required setup n/w in multiple layers.  
Stepwise learning builds next layer.
- eg. 25 nodes in one layer - not good  
5 nodes in 5 layer - better

### TensorFlow by Google

Recurrent network - Input is taken from  
(feedback n/w) one of the output.

## Market basket Analysis

- What are group of products that are bought together by customers.
- Business can take decisions about products to keep in store.
- Recommender systems can be built using market basket analysis.

cheese      Bread      Butter      Milk      Egg      Potato      Onion

Basket 1	1	1	1	1	1	0	0	0
2	1	1	0	0	1	0	0	0
3	0	1	1	1	1	0	0	0
4	0	1	0	0	0	1	1	0
5	0	1	1	1	0	1	0	0
6	1	1	1	1	0	0	0	0
7	0	0	1	0	1	0	0	0
8	1	0	1	1	1	1	0	0
9	0	0	1	0	1	0	0	0
10	0	1	1	1	1	0	0	0

Suppose data,  
not accurate

Support

Bread + Butter =  $\frac{6}{10} = 60\%$  frequent combinations

Cheese + Bread. =  $\frac{4}{10} = 40\%$  are calculated in percentage

Bread + Milk. =  $\frac{5}{10} = 50\%$ .

Cheese + Bread + Butter =  $\frac{3}{10} = 30\%$ .

$T_{C_1} + T_{C_2} + T_{C_3} + \dots + T_{C_7}$  combinations possible.

Good  
for  
Business

{ support  $\geq$  min. support &  
confidence  $\geq$  min. confidence &  
Lift  $> 1$

Page No.:  
Date: youva

- Lets assume a customer buys only at max. of 150 products in basket.
  - First find combinations and frequency occurred in basket.
  - Bread + Butter is commonly bought together, whereas Potato - onion combination bought less together.
  - With help of domain expert, we decide minimum support for combinations.
  - For supermarkets cutoff will be smaller, for shops cutoff can be higher.
- Confidence
- Further we evaluate the impact of products on each other.  
eg. A+B, when A sold how many times B was sold. ( $A \Rightarrow B$ )

$$\text{Bread} + \text{Butter} = \frac{6}{10}$$

$$(\text{Bread} \Rightarrow \frac{7}{8} \text{ Butter})$$

Bread

- Further based on confidence level cutoff can be decided eg. 67% or 2/3

Correlation

Positively correlated	$> 1$
Not correlated	$= 1$
Negatively correlated.	$< 1$

*Lift value*

$$\text{Lift} = \frac{\text{Confidence} (\text{in presence of other})}{\text{self} (\text{on own})}$$

(B in presence of A)

$$\text{Lift}(B) = \frac{A \Rightarrow B}{B} = \frac{70}{80} < 1 \quad -ve$$

$$(B \text{ on own}) \quad \frac{90}{80} > 1 \quad +ve$$

$$\text{Support}_{(A \Rightarrow B)} = \frac{n_{(A \cap B)}}{N} = P(A \cap B)$$

$$\text{Confidence}_{(A \Rightarrow B)} = \frac{n_{(A \cap B)}}{n_A} = P(B/A)$$

$$\text{Lift} = \frac{\text{Confidence}}{n_B/N} = \frac{\text{conf}(A \Rightarrow B)}{\text{sup}(B)}$$

$$P(A \cap B) = P(A) * P(B) \quad \text{provided } A \& B \text{ independent}$$

$$\text{Lift} = \frac{P(A \cap B)}{P(A) * P(B)} \rightarrow \frac{\text{conf}(A \Rightarrow B)}{\text{sup}(B)}$$

- Group of items known as "item set"
- Support  $\geq$  min. support. known as "frequent item set"
- frequent item set mining

## A-priority Algorithm (R-shrikant R-agarwal)

$$A+B+C \geq 50\%$$

then  $A+B \geq 50\%$ .  
 $B+C \geq 50\%$ .  
 $A+C \geq 50\%$ .

TID	Transaction	Item	Count
1	I1, I2, I5	I1	- 6
2	I2, I4	I2	- 7
3	I2, I3	I3	- 6
4	I1, I2, I4	I4	- 2
5	I1, I3	I5	- 2
6	I2, I3		
7	I1, I3		min-supp = 2
8	I1, I2, I3, I5		min-conf = 50%
9	I1, I2, I3		Lift > 1

~~I1 I2 - 4~~

~~I1 I2 I3 - 2~~

~~I1 I3 - 4~~

~~I1 I2 I5 - 2~~

~~minsup~~

~~condition~~

~~not meet~~

~~I1 I4 - 1~~

~~I1 I2 I5~~

~~I1 I5 - 2~~

~~I2 I3 I4~~

~~I2 I3 - 4~~

~~I2 I3 I5~~

~~I2 I4 - 2~~

~~I2 I4 I5~~

~~I2 I5 - 2~~

~~I2 I4 I5~~

~~I3 I4 - 0~~

~~I1 I2 I3 I5~~

~~I3 I5 - 1~~

~~I1 I2 I3 I5~~

~~I4 I5 - 0~~

~~I1 I2 I3 I5~~

Check sub combinations

~~Confidence = 50%~~

~~I1  $\Rightarrow$  I2, I3~~

~~I2, I3  $\Rightarrow$  I1~~

~~2/4~~

~~I2  $\Rightarrow$  I1, I3~~

~~I1, I3  $\Rightarrow$  I2~~

~~2/4~~

~~I3  $\Rightarrow$  I1, I2~~

~~I1, I2  $\Rightarrow$  I3~~

~~2/4~~

$$\begin{array}{l}
 \text{2/6} \quad - I_1 \Rightarrow I_2 I_5 \quad | \frac{1}{(6/3)} = 1.5 \\
 \text{2/7} \quad I_2 \Rightarrow I_1 I_5 \quad | \frac{1}{(7/9)} = 9/2 \\
 \text{2/8} \quad \checkmark I_5 \Rightarrow I_1 I_2 \quad | \frac{1}{(4/9)} = 2.25 \quad \checkmark I_1 I_2 \Rightarrow I_5 \frac{2}{4}
 \end{array}$$

$$\text{Lift} = \frac{\text{confidence}}{(\text{sup.}/\text{Total})}$$

Therefore  $I_5 \not\Rightarrow I_1 I_2$  or  
 $I_1 I_2 \Rightarrow I_5$  is considerable.