

### Data Mining Technologies:

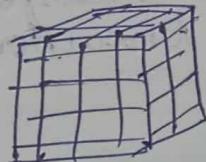
- i) DBMS
- ii) ML
- iii) OLAP (Online Analytical Processing)

### OLAP (Warehouse)

> OLAP multidimensional data i.e) cubic nature (Data cubes)  
operations of OLAP → Roll up, Drill down, slice, Dice, Pivot, cross

Roll up → low to higher level  
eg: city → country      Drill down = higher to lower  
eg: Year → month

Data cube:



ML

→ supervised → Labels classification  
→ unsupervised → clustering  
                    ↳ similarities

### STAGES OF DATA MINING PROCESS:

1. Understand the problem & data
2. Prepare the data
3. Build the model
4. Evaluate the result
5. Implement changes & Monitor

## Data Mining Techniques :

- Classification • Regression • Clustering • Association rule
- Anomaly Detection • Time Series Analysis • Neural networks
- Decision Tree • Ensemble methods • Text Mining

### Classification :

- Σ For categorical data it is used
- > It is a type of supervised

Association rule:

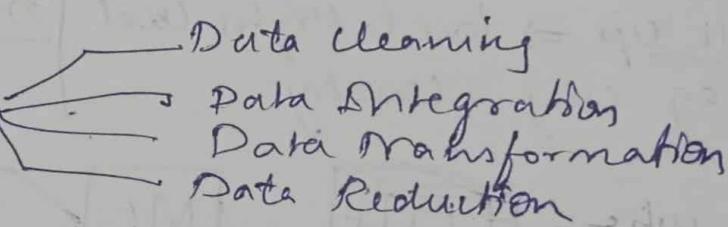
?  $A \times B$

Regression : > Continuous data | Ensemble : combine multiple models.

Neural nets : ➤ Backtracking

- ### Applications of DM :
- Research • Education sector
  - Transportation • Market Basket Analysis
  - Business Transactions • Intrusion Detection
  - Scientific Analysis • Finance & Banking sector
  - Insurance & Healthcare

## Data Preprocessing :



## Applications of DM :

### i) Education sector :

Educational institutions use DM to analyze student performance, identify at-risk students & improve programs.

> The performance Analysis → Decision Trees or Random Forest Alg

> Personalized learning → collaborative filtering techniques

### ii) Research:

- > Economics & Bioinformatics → Hierarchical clustering
  - groups genes with similar patterns
- > Climate modeling → Time series Analysis (ARIMA)
  - Historical climate data to predict future trends

### iii) Transportation:

- > Traffic flow prediction → Time series Analysis (ARIMA, LSTM)
  - predict future conditions based on historical data
- > Incident Detection → Anomaly detection (Isolation Forest)
  - detect unusual patterns indicates accidents
  - enabling faster response time

### iv) Market Basket Analysis:

- > Retail sales optimization → Apriori Algorithm
  - identifying frequent item sets
- > Product Recommendations - Association Rule Mining
  - understanding customer purchase patterns

### v) Business Transactions:

- > customer segmentation → K-Means clustering
  - grouping customers with similar purchasing patterns
- > credit scoring → Logistic Regression
  - Assessing the creditworthiness of customers to make informed lending decisions
- > Risk Mgmt → Decision Trees
  - assessing the risk level of transactions

## vii) Intrusion Detection :

- > Deep packet Inspection - Deep learning (convolutional Neural N/Ws)
  - To analyze packet content and identify hidden malicious patterns.
- > Hybrid Intrusion Detection systems - Ensemble learning (eg Bagging, Boosting)
  - Analyze both host & n/w data

## viii) Scientific Analysis :

- > Astronomy → k-means clustering
  - classifying stars based on their characteristics
- > Physics → Principal Component Analysis
  - Reducing the dimensionality of particle collision

## viii) Finance & Banking sector :

- > Fraud Detection - Anomaly Detection (e.g. Autoencoders)
  - detecting unusual transaction patterns
- > Anti Money Laundering - Association rule mining
  - Identifying patterns of transactions that are indicative of money laundering, such as series of small deposits followed by large withdrawals.

## ix) Insurance & Healthcare :

- > Medical Image Analysis → CNN
  - to analyze X-rays, MRIs, CTs
- > Customer retention in Insurance - support vector machines (SVM)
  - Identifying factors that influence policy holder retention to improve customer loyalty.

Attribute : Field, feature, dimension

Types : Nominal, binary, ordinal, discrete, continuous  
numeric ratio scale

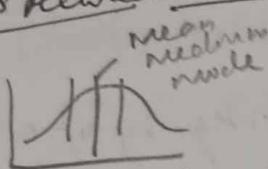
Discrete  $\rightarrow$  classification  $\hookrightarrow$  supervised

Continuous  $\rightarrow$  Regression

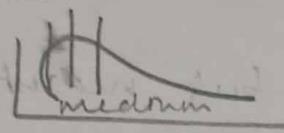
3 areas of statistical description :

> central tendency  $\rightarrow$  Measuring dispersion of data

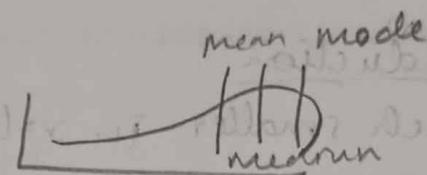
Skewed data :



mean mode



+ very skewed



- very skewed

Major Tasks in Data Preprocessing :

- i) Data cleaning  $\rightarrow$  incomplete (missing), incorrect, inconsistent (noisy)
- ii) Data Integration
- iii) Data Reduction

IV) Data Transformation & data discretization

cleaning :

Noisy  $\rightarrow$  Binning, Regression  $\begin{cases} \text{Linear} \\ \text{Multiple} \\ \text{Clustering} \end{cases}$

Discrepancy detection  $\rightarrow$  Tools: Data scrubbing, Data Auditing

Integration:

- > Schema identification problem, Entity identification problems
- > integrating data from diff sources
- > Removing redundancy  $\rightarrow$  correlation, covariance analysis

Correlation Analysis :

>  $\chi^2$  (chi-square test)  $\chi^2 = \sum \frac{(\text{Observed} - \text{Expected})^2}{\text{Expected}}$

$\hookrightarrow$  For nominal Data

> For numeric data  $\rightarrow$  Pearson's product moment coefficient of correlation  
coefficient  $r_{A,B} = \frac{\sum_{i=1}^n (a_i - \bar{A})(b_i - \bar{B})}{(n-1)\sigma_A \sigma_B} = \frac{\sum_{i=1}^n (a_i b_i) - n \bar{A} \bar{B}}{(n-1)\sigma_A \sigma_B}$

$r_{A,B} > 0$  - very correlated |  $< 0$  - very correlated |  $= 0$  independent

Correlation (linear relationship):

$$a'_k = \frac{(a_k - \text{mean}(A))}{\text{std}(A)} \quad b'_k = \frac{(b_k - \text{mean}(B))}{\text{std}(B)}$$

$$\Rightarrow \text{correlation}(A, B) = A' \cdot B'$$

Correlation coefficient  $\Rightarrow r_{A,B} = \frac{\text{cov}(A, B)}{\sigma_A \sigma_B}$

$$\text{cov}(A, B) = \frac{\sum_{i=1}^n (a_i - \bar{A})(b_i - \bar{B})}{n}$$

Data Reduction:

- > much smaller in volume but maintains integrity of original data
- > Data reduction strategies:

- ↳ Dimensionality reduction  $\rightarrow$  PCA, wavelet transforms
- ↳ Numerosity reduction  $\rightarrow$  Attribute subset selection
- ↳ Data compression  $\rightarrow$  lossy, lossless

Data Transformation:

Methods  $\rightarrow$  smoothing, Attribute / feature construction, Aggregation

Normalization techniques  $\rightarrow$  min-max, Z-score, decimal scaling

↳ range  $-1 \dots 1$  or  $0 \dots 1$

> Min - Max  $v'_i = \frac{v_i - \min}{\max - \min}$  ( $\text{new\_max} - \text{new\_min}$ ) +  $\text{new\_min}$

> Z-score  $v'_i = \frac{v_i - \mu}{\sigma}$   $\rightarrow$  Decimal-scaling  $v'_i = \frac{v_i}{T_{ij}}$

## (AOI) Data Generalization by Attribute Oriented Induction:

- > Data generalization → replacing relatively low level values with higher level concepts
  - > Concept description → form of data generalization
  - > concept description → generate description for data characterization & comparison : sometimes called class description (specifically while objects of class usage)
  - > characterization → provide a concise summarization of given collection of data.
  - > discrimination (class comparison) → provide description by comparing two (or) more data collection.
- Why AOI? → complex data types and aggregation  
→ user control vs automation

> These are the reasons for using AOI than OLAP

### AOI for Data characterization:

- proposed in 1989
- data cube based on materialized view of data (off-line aggregation)
- AOI based on query oriented, generalization based online data analysis technique.
- It first called task relevant data and perform generalization by attribute removal (or) attribute generalization.
- Aggregation is performed by merging identical tuples.

Query : use Big-University-DB

mine characteristics as "Science Students"

in relevance to name, gender, major, birthplace,  
birth-date, residence, phone#, spa  
from Student

General  
specific (or)  
where status in "Graduate"  
where status in {"MSc", "MA", "MBA", "PhD"}

## Two essential operations of AOE

op 1. attribute removal

case 1 : There is large distinct values for an attribute and there is no generalization

case 2 : Large distinct values and its higher level concepts are expressed in terms of other attribute.

op 2. attribute generalization :

large set of distinct values with generalization operator  
approaches to control generalization

- i) attribute generalization threshold control
- ii) generalized relation threshold control

eg : name - op1, gender - op2, birthplace op1 (case 2),  
DOB - op2, major\_subj - op2, residence - op2,  
phone\_no - op1, GPA - op2

## Efficient Implementation of AOE

efficiency is analyzed as follows :

1. Relational query collect task-relevant data into working relation, efficiency depends on query processing
2. collect statistics on working relation
3. drives the prime relation

### Algorithm :

1. W. (DMA query)  $\rightarrow$  I/p (Data Mining query)

2. Prepare for generalization

\* collect distinct values

\* removed (or) compute its level

Consider Big University DB

I/p consists  
like names, age, residence, etc..

O/p  $\rightarrow$  GPA,

residence, DOB,

3. P  $\leftarrow$  generalization (W) & consider prime relation  $\rightarrow$  final O/p

## Attribute oriented induction for class comparison:

» discrimination (or) comparison  $\Rightarrow$  distinguish target class from its contrasting class

### How class comparison performed?

1. Data collection

2. Dimension relevance analysis

3. synchronous generalization

4. Presentation of derived comparison.

1. Data collected by query processing

2. select only higher relevant dimensions

3. Generalization performed on target class controlled by user (or) expert

4. Visualized in the form of tables, graphs & rules.

### Example:

use Big-University-DB to mine comparison as "grad. vs UG"

In relevance to name, gender, major, birth-place, birth-date, residence, phone #, gpa

for "graduate-students" where status is "graduate"  
versus "UG"

where status is "UG"

analyze count %

from student

For  $\Leftarrow$  Target class (prime relation)

Graduate	Major	Age	gpa	Count %
SC	21-25	good	5.53	
SC	26-30	good	5.02	
Business	over 30	Beneath	9.68	

Contrasting class  $\rightarrow$  For UG

Major	age	gpa	Count %
SC	16-20	fair	5.53
SC	16-20	good	9.53
SC	26-30	good	2.32
Business	over 30	Ex	0.63

Statistical Measures : understand the distribution of data

- ↳ Measure of central tendency  $\leftarrow$  mean, median, mode
- ↳ Measuring dispersion  $\leftarrow$  variance, std deviation

Association Rules : If then stmnt

- > Market Basket Analysis motivation eg
- > Frequent pattern

Representation  $a \Rightarrow b$  support - together

confidence - If that then

> support(A) =  $\frac{\text{no. of transactions in which A appears}}{\text{total no. of transactions}}$

> Confidence ( $A \Rightarrow B$ ) =  $\frac{S(A \cup B)}{S(A)}$

Closed Freq Itemset :

Maximal Frequent Itemset

Immediate supersets have  
same support

For Identifying frequent Itemset, remove the items that are less than min support count.

Algs for Frequent Itemsets :

1) Apriori : (with candidate generation)  $\xrightarrow{\text{Join Step}} \xrightarrow{\text{Prune Step}}$

2) FP Growth (without candidate generation)

Apriori eg :

Tid	Contents
T1	I1, I2, I3
T2	I2, I4
T3	I2, I3
T4	I1, I2, I4
T5	I1, I3
T6	I2, I3
T7	I1, I3
T8	I1, I2, I3, I5
T9	I1, I2, I3

$\Rightarrow$

At last

Items	Freq
I1, I2, I3, I5	1
I2, I3	2
I1, I3	2

so we foot

before Freq set

(min support count)

where min support count  $\rightarrow 2$

I <sub>1d</sub>	Items
I <sub>1</sub>	A, C, D
I <sub>2</sub>	B, C, E
I <sub>3</sub>	A, B, C, E
I <sub>4</sub>	B, E

min support count }  $\rightarrow 2$

A - 2  
B - 3  
C - 3  
 $\times$  D - 1  
E - 3

AB - 1 X  
AC - 2  
AE - 1 X  
BC - 2  
BE - 2  
CE - 2

ABC - 1 X  
BCE - 2

ACE - 1 X  
ABE - 1 X

$S(A \cup B) = \frac{2}{9}$

i) For I<sub>1</sub>, I<sub>2</sub>, I<sub>3</sub>:  
From 1<sup>st</sup> eg take freq itemset.

I<sub>1</sub>, I<sub>2</sub>, I<sub>3</sub> - {2}

Generate subsets

I<sub>1</sub>, I<sub>2</sub>, I<sub>3</sub>, I<sub>1</sub>I<sub>2</sub>

- Association rules:
- 1) I<sub>1</sub>  $\rightarrow$  I<sub>2</sub>, I<sub>3</sub>
  - 2) I<sub>2</sub>  $\rightarrow$  I<sub>1</sub>, I<sub>3</sub>
  - 3) I<sub>3</sub>  $\rightarrow$  I<sub>1</sub>, I<sub>2</sub>

- 4) I<sub>1</sub>, I<sub>2</sub>  $\rightarrow$  I<sub>3</sub>
- 5) I<sub>1</sub>, I<sub>3</sub>  $\rightarrow$  I<sub>2</sub>
- 6) I<sub>2</sub>, I<sub>3</sub>  $\rightarrow$  I<sub>1</sub>

I<sub>1</sub>, I<sub>2</sub>, I<sub>3</sub>

I<sub>2</sub>  $\rightarrow$  I<sub>1</sub>, I<sub>3</sub>

$S = \frac{2}{9}$

$I_1 \rightarrow I_2, I_3$

$S = \frac{2}{9}$

If A  $\rightarrow$  B is association rule

(C) Confidence level =  $\frac{\text{support}(A \cup B)}{\text{support } A}$

For rule 1: C =  $\frac{S(I_1, I_2, I_3)}{S(I_1)}$

$$= \frac{2/9}{6/9} = \frac{1}{3}$$

For rule 2: C =  $\frac{S(I_1, I_2, I_3)}{S(I_2)}$

and other confidence

> Similar process for  $I_1, I_2, I_3$  which is in ppt

Picking min support count [ 1, 2, 3, ... select 2  
 [ 1, 3, 4, ... select 3  
 [ 4, 5, 6, ... select 4 ] ] ]

### > FP tree (Frequent pattern tree)

> To overcome the disadvantage of candidate gen in Apriori, FP growth is used

eg 1:	TID	Items bought	f - 4	l - 2
	1	f a c d g i m r	a - 3	0 - 2
	2	a b c f l m o	c - 4	l - 1
	3	b f h j o	d - 1	h - 1
	4	b c k s p	g - 1	b - 3
	5	a f c e l p m n	i - 1	s - 1
			m - 3	u - 1
			p - 3	n - 1
			j - 1	

Given min support count = 3

Find desc ord  $\rightarrow$  f c a b m p l o e d ...

only items to consider since count  $\geq 3$

Now change transaction table (based on desc ord)

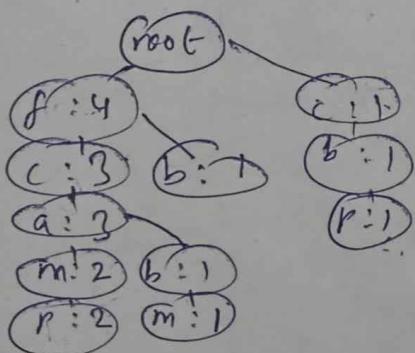
1	f c a m p
2	f. c a b m
3	f b
4	c b p
5	f c a i m p

(construct tree) based on these transactions from (root) of 3, make node with child nodes & have a count which should be same initial count

> desc order since

tree size will be lower

Final tree



Conditional Pattern Base

$\Rightarrow$  Write the ascending order of considered items & fill the pattern base from the path to reach the item

Items	Conditional pattern base	count of item "P" through the path	condition FP tree is 3rd col
P	{f, c, a : 2}, {c, b : 1}		
m	{f, c, a : 2}, {f, c, a, b : 1}		
b	{f, c, a : 1}, {f : 1}, {c : 1}		
a	{f, c : 3}		
c	f{f : 3}		
f	∅		

for item P  $\Rightarrow c : 2 + 1 \Rightarrow c : 3$ ,  $f : 2 \times m : 2 \times b : 1 \times$

$m : f, c, a : 2 + 1 \Rightarrow f, c, a : 3$  ✓ if  $f : 3 - c : 3, a : 1$  ?

> after Pruning for all items

> draw the final table for frequent pattern generated

- draw the final table for frequent pattern generated
- Steps**: Count  $\rightarrow$  Descending Order  $\rightarrow$  Eliminate which  $<$  min support count
- construct new table  $\rightarrow$  draw FP tree
  - find conditional pattern base  $\rightarrow$  Then condition FP tree  $\rightarrow$  final table with frequent pattern

Prblm:

(\*) Min support = 60%. min confidence = 80%. Apply

Apriori & FD growth.

If asked to print Preq  
Itemset then find sets alone  
no need for Pruning confidence,  
support

TID Items bought

T<sub>1</sub> M O N K E Y

T<sub>2</sub> D O N K E X

T<sub>3</sub> M A K E

T<sub>4</sub> M U C K Y

T<sub>5</sub> C O O K I E

If asked for  
strong association  
rule Prule confidence  
& support f.s

Workshop

# i) Apriori Alg

M - 3 ✓  
 O - 4 ✓  
 N - 2  
 K - 5 ✓  
 E - 4 ✓  
 Y - 3 ✓  
 D - 1  
 A - 1  
 U - 1  
 C - 2  
 T - 1

MO - 1  
 MK - 3 ✓  
 ME - 2  
 MY - 2 ✗  
~~OE~~  
 OK - 3 ✓  
 OE - 3 ✓  
 OY - 2  
 KE - 4 ✓  
 KY - 3 ✓  
 EY - 2

MKY - 2  
 MOK - 1  
 MOE - 1  
 KYO - 2  
 MKE - 2  
 MEY - 1  
 MOY - 1  
 OKE - 3 ✓  
 EOY - 2

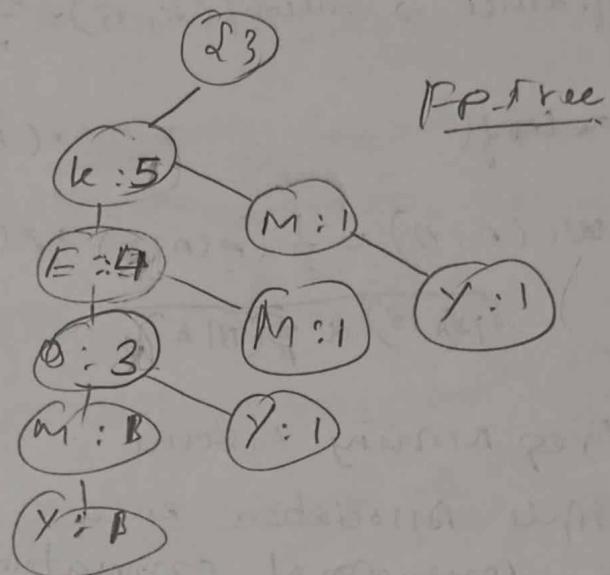
$OKE \Rightarrow$  1)  $O \rightarrow KE$       4)  $KE \rightarrow O$   
 2)  $K \rightarrow OE$       5)  $OE \rightarrow K$   
 3)  $E \rightarrow OK$       6)  $OK \rightarrow E$

1)

### iii) FP Growth :

$K, E, O, M, Y$  :  $k=5, E=4, O=4, M=3, Y=3$

$T_1$	KEOMY
$T_2$	KEOY
$T_3$	KE M
$T_4$	KMY
$T_5$	KEO



FP-tree

- > All confidence  $\rightarrow \text{allconf}(A, B) = \frac{\text{sup}(A \cup B)}{\max(\text{sup}(A), \text{sup}(B))}$
- > ~~markconf~~  $= \max(P(A|B), P(B|A))$
- >  $\text{kull}(A, B) = \frac{1}{2} (P(A|B) + P(B|A))$
- > cosine:  $\sqrt{P(A|B)} \times \sqrt{P(B|A)}$

Adv Programming Patterns :

no repeated  
↑ predicates

- > Multiple Association rules
- > multi-dimensional association mining. (Interdimensional)
  - single-dimensional " " also known as ~~extra dim~~
- > Repeated predicates - hybrid dimensional assoc rules.
- ↳ static → Discretization → converting numeric to nominal values
- ↳ dynamic →

single-dimensional association rules :

e.g.: buys(A, 'car')  $\Rightarrow$  buys(B, 'petrol')

Multi-dimensional association rules :

- > with no repeated predicates  $\Rightarrow$  inter-dimensional assoc rule
- > with repeated predicates  $\Rightarrow$  hybrid-dimensional rules

2 approaches for handling quantitative data

- ↳ discretization using concept hierarchy  $\Rightarrow$  static (values)
  - predefined or known
- ↳ discretization w/ clustering  $\Rightarrow$  dynamic (unknown values)
  - dynamic quantitative assoc rules

## UNIT - I

$5 \times 10 = 50$

- > Problems on Data Preprocessing Cleaning & transformation
  - > Attributes, Distance measures, Removing noise
  - > Normalization  $\rightarrow$  Mean Median Mode, Chebyshev min-max
  - > Binning, PCA, wavelet  $\rightarrow$  Numerosity, dimensionality reduction
  - > AOI
- CONT - 2 > Association Rule mining  $\leftarrow$  FP growth  $\leftarrow$  Apriori
- > For FP growth draw tree for each transaction Frequent & Strong rules

Classification vs Prediction

- > For continuous data  $\rightarrow$  regression  $\rightarrow$  numeric prediction
- classification  $\Rightarrow$  1. learning step / 2. classification step

Distance Measures b/w numeric data points :

Minkowski distance

$$d(i, j) = \sqrt[h]{|x_{i1} - x_{j1}|^h + |x_{i2} - x_{j2}|^h + \dots + |x_{ip} - x_{jp}|^h}$$

p-dimension & h-order (user-defined)

if  $h=1$  means it becomes  $L_1$  norm (or) Manhattan distance

$$d(i, j) = |x_{i1} - x_{j1}| + |x_{i2} - x_{j2}| + \dots + |x_{ip} - x_{jp}|$$

if  $h=2$  means Euclidean distance

$$d(i, j) = \sqrt{|x_{i1} - x_{j1}|^2 + |x_{i2} - x_{j2}|^2 + \dots + |x_{ip} - x_{jp}|^2}$$

if  $h=\infty$  supremum distance /  $L_{\max}$  norm /  $L_\infty$  norm

$$d(i, j) =$$

e.g.: data set given calculate euclidean distance

$x_1$	12	30	15	30	19
$x_2$	25	40	23	45	56
$x_3$	38	11	35	63	42

$$d(1,2) = \sqrt{|12-25|^2 + |30-40|^2 + |15-23|^2 + |30-45|^2 + |19-56|^2}$$

$$= \sqrt{13^2 + 10^2 + 8^2 + 15^2 + 37^2}$$

$$d(1,2) = \sqrt{169 + 100 + 64 + 225 + 1369} = 43.898$$

$$d(1,3) = \sqrt{|12-38|^2 + |30-11|^2 + |15-35|^2 + |30-63|^2 + |19-42|^2}$$

$$= \sqrt{26^2 + 19^2 + 20^2 + 33^2 + 23^2}$$

$$= 55.272$$

$$d(2,3) = \sqrt{13^2 + 29^2 + 12^2 + 18^2 + 14^2} = 40.915$$

∴ Distance b/w 2,3 is lesser with 40.915

	$x_1$	$x_2$	$x_3$
$x_1$	0	43.89	55.27
$x_2$	43.89	0	40.9
$x_3$	55.27	40.9	0

Min-Max sum:

- 1) 13, 15, 16, 16, 19, 20, 23, 29, 35, 45, 44, 53, 62, 69, 72. Use min-max alg to solve for age 42. On the range [0.0, 1.0]

Soln:

Given  $v = 42$ ,  $\min = 13$ ,  $\text{new-min} = 0.0$   
 $\text{new-max} = 1.0$

For age = 42  $v' = \frac{v - \min}{\max - \min} \times (1.0 - 0.0) + 0.0$

$$\therefore \frac{42 - 13}{72 - 13} \times 1 = \frac{29}{59} = 0.492$$

$$\text{Age} = 45 = \frac{45 - 13}{72 - 13} \times 1 = 0.5423$$

Note : If  $v$  is not given in ques find min-max  
for all datas given

$$2) 200, 300, 400, 600, 1000 \quad [0, 1]$$

$$\frac{\text{Min-Max}}{v = 200} , \quad v' = \frac{200 - 200}{1000 - 200} \times (1 - 0) + 0 = 0$$

$$v = 300 , \quad v' = \frac{300 - 200}{1000 - 200} \times (1 - 0) = \frac{1}{8} = 0.125$$

$$v = 400 , \quad v' = \frac{400 - 200}{1000 - 200} = \frac{200}{800} = 0.25$$

$$v = 600 , \quad v' = \frac{600 - 200}{1000 - 200} = \frac{400}{800} = 0.5$$

$$v = 1000 , \quad v' = \frac{1000 - 200}{1000 - 200} = 1$$

$$v' = \frac{v - \mu}{\sigma}$$

Z-score :

$$v' = \frac{v - \mu}{\sigma} \quad \sigma = \sqrt{\frac{\sum (x - \mu)^2}{N}}$$

$$\mu = \frac{200 + 300 + 400 + 600 + 1000}{5} = \frac{2500}{5} = 500$$

$$\sigma = \sqrt{\frac{(200 - 500)^2 + (300 - 500)^2 + (400 - 500)^2 + (600 - 500)^2}{5}}$$

$$\sigma = \sqrt{\frac{(200 - 500)^2 + (300 - 500)^2 + (400 - 500)^2 + (600 - 500)^2 + (1000 - 500)^2}{5}}$$

$$= \sqrt{\frac{90000 + 40000 + 10000 + 10000 + 250000}{5}}$$

$$\sigma = 282.842$$

Let  $v = 550$

$$v' = \frac{550 - 500}{282.84} = 0.1767$$

For sum 1) Z-score :

$$\bar{x} = 35.133 \quad \sigma = 19.94$$

For 42,  $v' = \frac{42 - 35.133}{19.94}$

For sum 2) Decimal-scaling

$$= \frac{v}{10^j}$$

for  $v = 200$ ,  $j = 3$   
then  $v = 200 / 10^3 = 0.2$

sums  $\leftarrow$  UNIT-1  $\rightarrow$  theory

- > Bonning method
- > Chi-square test
- > wavelet transform
- > PCA
- > 3 normalization (Min Max, Z-score, Decimal scaling)
- > Distance of data points (Distance b/w each & every point)  
(Default : Euclidean or else specified in ques)

UNIT-2  $\Rightarrow$  Assoc rule mining  $\rightarrow$  A priori

FP growth sums  
 $\hookrightarrow$  Pseudocode, Adv & Disadv  
 $\hookrightarrow$  & sums (Mentioned in John Phane)

> Solve sums using full explanation & detail

> Preprocessing how we remove Noisy data

> Sampling  $\begin{cases} \text{with replacement} \\ \text{without replacement} \end{cases}$

## Classification

1. One R classifier  $\rightarrow$  rule:

$\rightarrow$  simple, accurate

$\rightarrow$  select the rule with smallest total error as its one rule to construct rule

$\rightarrow$  create frequency table

eg: dataset in ppt (play tennis)

eg: frequent table from the dataset

outlook	play tennis	
	yes	no
sunny	2	3
rainy	3	2
overcast	4	0
	$\sum \text{ve} = 9$	$\sum \text{ne} = 5$

Based on this eg:

$\Rightarrow$  If sunny THEN play tennis = no

2. Decision Tree Algorithm: (DT) or ID3 Alg:

Decision Tree Induction:

3 parameters : i)  $D \rightarrow$  Data ii) Attribute list

iii) Attribute selection measures (to find which node should be root / leaf (internal))

$\hookrightarrow$  Information gain  $\xrightarrow{\text{or entropy}}$   $\text{Info}(D) = - \sum_{i=1}^M p_i \log_2(p_i)$

particular attribute  $\leftarrow \text{Info}_A(D) = \sum_{j=1}^k \frac{|D_j|}{|D|} \times \text{Info}(D_j)$

$$\boxed{\text{Gain}(A) = \text{Info}(D) - \text{Info}_A(D)} \quad \text{or} \quad \boxed{\text{Entropy}(A)}$$

Dataset before (playing tennis) set is taken

$$\text{Yes} = 9 \quad \text{No} = 5 \quad P_{\text{Yes}} = 9/14 \quad P_{\text{No}} = 5/14$$

$$\text{Info}(D) = - \frac{9}{14} \log \frac{9}{14} - \frac{5}{14} \log \frac{5}{14} = 0.9403$$

For outlook attribute,

sunny	rainy	Overcast
$y \text{ No}$ 2 3	$y \text{ No}$ 3 2	$y \text{ No}$ 4 0

$$\text{Info}_{\text{outlook}}(D) = \frac{5}{14} \left( -\frac{2}{5} \log_2 \frac{2}{5} - \frac{3}{5} \log_2 \frac{3}{5} \right)$$

$$+ \frac{5}{14} \left( -\frac{3}{5} \log_2 \frac{3}{5} - \frac{2}{5} \log_2 \frac{2}{5} \right)$$

$$\text{for rainy } + \frac{4}{14} \left( -\frac{4}{4} \log_2 \frac{4}{4} \right) = 0.693$$

for overcast

$$\text{rain cover}(D) = 0.9403 - 0.693 = 0.247$$

For Temperature :

	$y$	$N$	$y$	$N$	$C$	$y$	$N$
H	2	2	M	4	2	3	1

$$\text{tot } y = 9 \quad \text{tot } N = 5$$

$$\text{Info}(D) = 0.9403$$

$$\begin{aligned} \text{Info}_{\text{temp}}(D) &= \frac{4}{14} \left( -\frac{2}{4} \log \frac{2}{4} - \frac{2}{4} \log \frac{2}{4} \right) + \\ &\quad \frac{6}{14} \left( -\frac{4}{6} \log \frac{4}{6} - \frac{2}{6} \log \frac{2}{6} \right) + \frac{4}{14} \left( -\frac{3}{4} \log \frac{3}{4} - \frac{1}{4} \log \frac{1}{4} \right) \\ &= \frac{4}{14} \left( \frac{1}{2} + \frac{1}{2} \right) + \frac{6}{14} (0.39 + 0.52 + 3) + \frac{4}{14} (0.3113 + 0.5) \\ &= 0.9111 \end{aligned}$$

$$\text{rain(temp)} = 0.9403 - 0.9111 = 0.0292$$

Humidity :

	$y$	$N$	total
High	3	4	7
Normal	6	1	7

$$\text{Info}_{\text{humidity}}(D) = \frac{7}{14} \left( -\frac{3}{7} \log \frac{3}{7} - \frac{4}{7} \log \frac{4}{7} \right) + \frac{7}{14} \left( -\frac{6}{7} \log \frac{6}{7} - \frac{1}{7} \log \frac{1}{7} \right)$$

$$= 0.4927 + 0.2959$$

$$= 0.7886$$

$$\text{rain} = 0.9043 - 0.7886 = 0.1157$$

windy:

		Y	N	
strong	3	3	6	
weak	6	2	8	
	9	5		

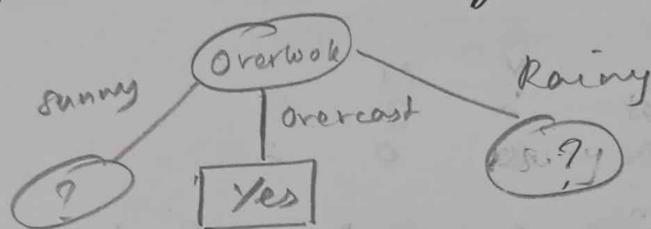
$$\text{Info}_{\text{windy}}(D) = \frac{6}{14} \left( -\frac{3}{6} \log \frac{3}{6} - \frac{3}{6} \log \frac{3}{6} \right) + \frac{8}{14} \left( -\frac{6}{8} \log \frac{6}{8} - \frac{2}{8} \log \frac{2}{8} \right)$$

$$= \frac{6}{14} (1) + \frac{8}{14} (0.3112 + 0.5)$$

$$= 0.8920$$

$$\text{gain} = 0.9403 - 0.8920 = 0.0483$$

Among all gain overcast is high so the root is overcast



Here from the given dataset overcast has all yes so yes is leaf node of overcast but for sunny, rainy we have to make separate tables & calculate info & gain.

outlook	Temp	Humidity	windy	PlayFenses
sunny	Hot	High	weak	No
sunny	Hot	High	strong	No
sunny	Mild	High	weak	No
sunny	Cool	Normal	weak	Yes
sunny	Mild	Normal	strong	Yes
Rainy	Mild	High	weak	Yes
Rainy	Cool	Normal	weak	Yes
Rainy	Cool	Normal	strong	No
Rainy	Mild	Normal	weak	Yes
Rainy	Mild	High	strong	No

For sunny  $N = 5$   $Y_0 = 2$   $N_0 = 3$

$$\text{Info}_{\text{sunny}}(D) = -\frac{2}{5} \log \frac{2}{5} - \frac{3}{5} \log \frac{3}{5} = 0.159 + 0.133$$

$$\text{Info}_{D=\text{sunny}} = 0.292$$

For temp :

	$Y$	$N$
H	0	2
M	1	1
C	1	0

$$\begin{aligned} \text{Info}_{\text{temp}}(D) &= \frac{2}{5} \left( -\frac{2}{2} \log \frac{2}{2} \right) + \frac{2}{5} \left( -\frac{1}{2} \log \frac{1}{2} - \frac{1}{2} \log \frac{1}{2} \right) \\ &\quad + \frac{1}{5} \left( -\frac{1}{1} \log \frac{1}{1} \right) \end{aligned}$$

For Humidity :

	$Y$	$N$
H	0	3
N	2	0

$$\text{Info}_{\text{Humid}}(D) = \frac{3}{5} \left( -\frac{2}{3} \log \frac{2}{3} \right) + \frac{2}{5} \left( -\frac{2}{2} \log \frac{2}{2} \right)$$

Gini-Index :

$$\text{Gini-Index}(D) = (1 - \sum (p_i)^2)$$

Pruning Tree: What, why, categories:  
 ↗ prepruning  $\Rightarrow$  halting its construction early  
 ↗ postpruning  $\Rightarrow$  removing subtrees from a fully grown tree

Decision tree issues  $\Rightarrow$  Repetition & replication

Model Evaluation:

> important for all models / Algs

Confusion Matrix:

		Predicted class		
		Yes	No	
Actual class	Yes	TP	FN	P.
	No	FP	TN	N
		<u>P'</u>	<u>N'</u>	

Evaluation Measures:

$$\text{Accuracy} = \frac{TP + TN}{P + N}$$

% of test set tuples that are correctly classified

$$\text{Error rate} = \frac{FP + FN}{P + N}$$

opposite to accuracy

Class imbalance problem:

$$\begin{array}{l} \text{specificity, sensitivity} \\ \text{TP}/\text{TN} \\ \text{TP}/P \end{array}$$

$$\begin{array}{l} \text{specificity} = TN/N \\ \text{sensitivity} = TP/P = \text{recall} \end{array}$$

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

(Measure of exactness)

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

(measure of completeness)

$$F_1 \text{ score} = \frac{2}{\frac{1}{\text{Precision}} + \frac{1}{\text{Recall}}} = \frac{2 * (\text{Precision} * \text{Recall})}{\text{Precision} + \text{Recall}}$$

Weighted F-score  $F_B = (1 + \beta^2) * \frac{(\text{Precision} * \text{Recall})}{(\beta^2 * \text{Precision}) + \text{Recall}}$

Assessing accuracy:

> Holdout > Random subsampling > Cross validation > Bootstrap

Holdout  $\rightarrow$  training & test data | Cross validation  $\rightarrow$  k-fold, N-fold  
 Random subsampling  $\rightarrow$  Continuous holdout | Bootstrap  $\rightarrow$  Sampling with replacement

Bayesian classifier:

> statistical classifier > Bayes theorem  $P(A|B) = \frac{P(B|A) \cdot P_A}{P(B)}$   
 $B \rightarrow \text{evidence } A \rightarrow \text{hypothesis}$

For buys-computer dataset:

$$P(\text{buys yes}) = \frac{9}{14} \quad P(\text{buys no}) = \frac{5}{14}$$

$$P(\text{age} > 40) = \frac{3}{9} \quad \frac{9}{5} \rightarrow \begin{cases} \text{If age} > 40 \\ \text{No is zero/s} \end{cases}$$

$$P(\text{income - medium}) = \frac{9}{9} \quad \frac{2}{5}$$

$$P(\text{student - no}) = \frac{3}{9} \quad \frac{4}{5} \quad \begin{cases} \text{if all other below} \\ \text{No is like } \frac{1}{6} \end{cases}$$

$$P(\text{credit - fair}) = \frac{6}{9} \quad \frac{2}{5} \quad \begin{cases} \text{No is like } \frac{2}{6} \\ \text{(income - medium)} \end{cases}$$

$$\frac{9}{14} \times 0.032 \quad 0.051 \times \frac{5}{14} = \frac{4}{6}, \frac{2}{6} \dots$$

$$\frac{0.021}{\boxed{\text{Yes}}} > \frac{0.018}{\text{No}}$$

& also total

$$\begin{aligned} \text{buys no} &= \frac{6}{15} \\ \text{Yes} &= \frac{9}{15} \end{aligned}$$

$$\therefore x = \{\text{age} > 40, \text{income - medium}, \text{student - no}, \text{credit - fair}\} - \text{Yes}$$

[ If probability is zero, add +1 for every attribute count ] & also add the total count in samplespace value ]

→ Problem : use Euclidean | Theory : Manhattan for eg

### KNN Algorithm : (k Nearest Neighbor)

→ Non-parametric Algorithm

→ both classification & regression

→ lazy learner Algorithm

[ k → no. of neighbours considered ]

eg :	x	y	class	
	7	7	Bad	
	7	4	Bad	
	3	4	Good	
	1	4	Good	

Nearest instance (3,7)  
Predict the class

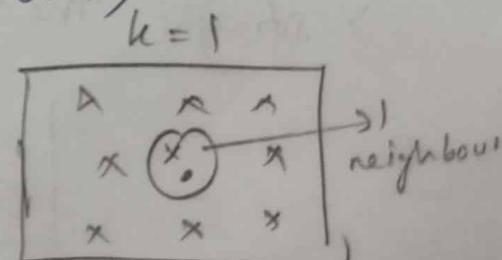
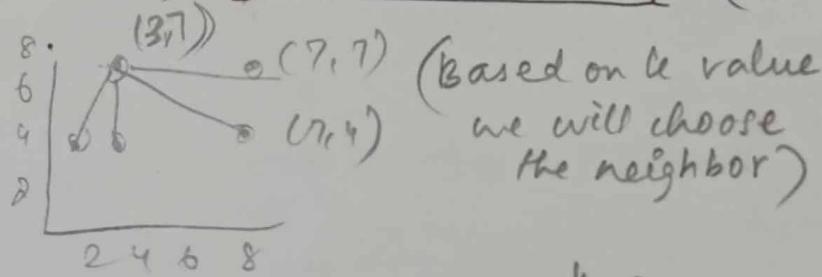
distance (Manhattan b/w new instance & given x, y)

Bad	$4+0 = 4$ ②	Rank based on distance if $k=1$ good $k=2$ bad $k=3$ good	∴ if $k=3$ given then (3,7) is majority good
Bad	$4+3 = 7$ ④		
Good	$0+3 = 3$ ①		
Good	$2+3 = 5$ ③		

Limitations :> what k value, distance metrics  
> computational complexity

$$\text{For KNN} \Rightarrow k = \sqrt[3]{N}$$

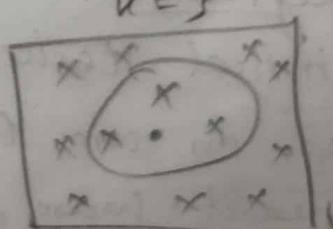
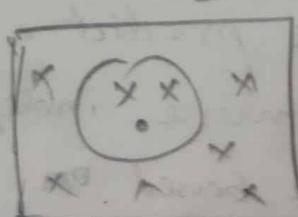
$N \rightarrow$  size of the dataset  
(k should be odd)



( k → only odd bcoz

there may be same number of classes

if even like 2 Good, 2 Bad )



$k=5$  (given)

Height (cm)	Weight (kg)	Class	Distance
167	51	Underweight	$3+4 = 7 \textcircled{5}$
182	62	Normal	$12+7 = 19 \textcircled{9}$
176	69	N	$6+14 = 20 \textcircled{6}$
173	64	N	$3+9 = 12 \textcircled{7}$
172	65	N	$2+10 = 12 \textcircled{8}$
169	56	UW	$1+1 = 2 \textcircled{1}$
173	58	N	$3+3 = 6 \textcircled{4}$
170	57	N	$0+2 = 2 \textcircled{2}$
174	56	N	$4+1 = 5 \textcircled{3}$
170	55	?	Distance

If  $k=1$  UW  
 $k=2$  N      Even though for  $k=5$  UW  
 $k=3$  N      Majority is Normal  
 $k=4$  N      So for  $(170, 55)$  is Normal  
 $k=5$  UW

### Numeric Prediction:

#### → Linear Regression

- > Linear models (Regression Model) → Under supervised learning
- > straight line through a set of data points
- > straight line represents best fit prediction equation



### Linear Regression:

- statistical technique predict a continuous values based on one (or) more independent attributes
- eg: Predict house price based on living area
- Predict future income of student based on college, major & GPA

## Least square Method

$x$	$y$	$x - \bar{x}$	$y - \bar{y}$	$(x - \bar{x})(y - \bar{y})$	$(x - \bar{x})^2$
1	2	-2	-2	4	4
2	4	-1	0	0	1
3	5	0	1	0	0
4	4	1	0	0	1
$\bar{x} = \frac{5}{8}$	$\bar{y} = \frac{5}{4}$	2	1	$\frac{2}{6}$	$\frac{4}{10}$

$$m(\text{slope}) = \frac{\sum (x - \bar{x})(y - \bar{y})}{\sum (x - \bar{x})^2} = \frac{6}{10}$$

$$y = mx + c \rightarrow m \text{ intercept}$$

$$\text{let } y = 4 \quad (\text{so } \bar{y} = 4) \quad \bar{x} = 3$$

$$4 = \frac{6}{10} \times 3 + c \Rightarrow 4 = \frac{3}{5} \times 3 + c$$

$$\boxed{c = 2.2}$$

2) Predict for  $x = 4$   $y = ?$

$x$	$y$
2	50
3	60
5	80
7	90
9	95
mean	$\frac{5.2}{75}$

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

$x - \bar{x}$	$y - \bar{y}$	$(x - \bar{x})(y - \bar{y})$	$(x - \bar{x})^2$	$m = \frac{221}{34.4}$	$\bar{y} = m\bar{x} + c$
-3.2	-25	80	10.24	$m = 6.424$	$75 = 6.424 \times 5.2 + c$
-2.2	-15	33	4.84		
-0.2	5	-1	0.04		
2.2	15	33	4.84		
4.2	20	86	14.44		
		$\frac{221}{34.4}$	$\frac{34.4}{34.4}$	$c = 41.5930$	$\boxed{c = 41.6}$

$$\text{Now if } x = 4, \quad y = 6.424 \times 4 + 41.6 = 67.29$$

$$\boxed{y = 67.3}$$

3). Find the Regression line for

$x_i$	$y_i$	$x - \bar{x}$	$y - \bar{y}$	$(x - \bar{x})^2$	$(x - \bar{x})(y - \bar{y})$
1	4	-3.0	-7.0	9.0	-21
3	10	-1.0	-1.0	1.0	+1
5	14	1.0	3.0	1.0	3
7	16	3.0	5.0	9.0	15
		<u>3.2</u>	<u>8.8</u>	<u>20</u>	<u>40</u>
<u>4</u>	<u>10.8</u>				

$$m = \frac{40}{20} = 2.0 \Rightarrow \boxed{y = 2x + c}$$

$$\bar{y} = 2\bar{x} + c \Rightarrow 11 = 2 \times 4 + c \Rightarrow c = 3$$

$$\therefore \boxed{y = 2x + 3}$$

Linear Regression :

$$Y = \beta_0 + \beta_1 x + E$$

$$E(Y) = \beta_0 + \beta_1 x$$

Best line for Regression  
shd cross through  
Mean value

I) least square error:  $\bar{y} = b_0 + b_1 \bar{x}$

$$b_1 = \frac{\sum (x_i - \bar{x})(y_i - \bar{y})}{\sum (x_i - \bar{x})^2}$$

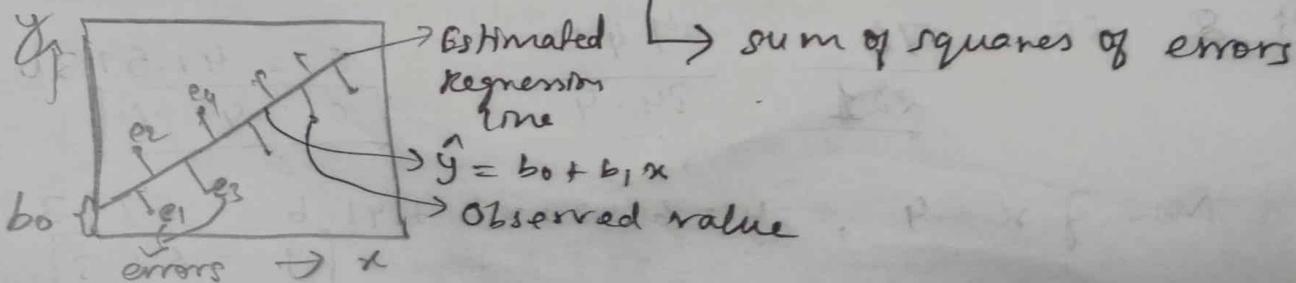
↓  
minimum least square error

sum of squares

$$S_{xx} = \sum_{j=1}^n (x_j - \bar{x})^2 \quad S_{yy} = \sum_{f=1}^n (y_f - \bar{y})^2$$

$$S_{xy} = \sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})$$

$$\text{Slope } (m) = \frac{S_{xy}}{S_{xx}} \quad SSE = S_{yy} - \frac{S_{xy}}{S_{xx}}$$



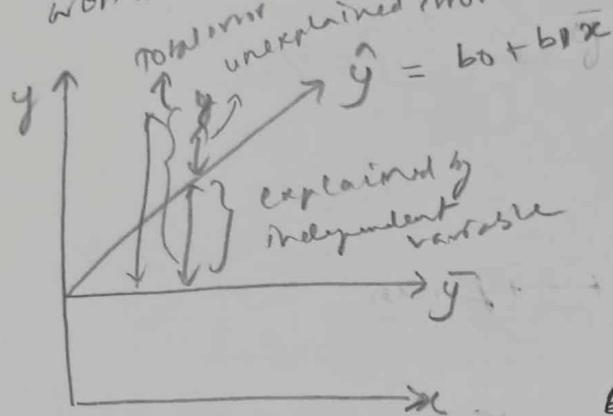
## II) $R^2$ - coefficient of Determination

$$R^2 = \frac{SSR}{SST}$$

$SSR \rightarrow$  sum of square of Regression  
 $SST \rightarrow$  total sum of squares

$$SST = SSR + SSE$$

This is used when  $x$  value (Independent variables) won't have relationship with dependent variables.



$R^2$  - Explained error

Total error

Interval b/w 0 to 1

Total error

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

Explained error,  $\sum (\hat{y} - \bar{y})^2$

Total error,  $\sum (y - \bar{y})^2$

sum:

	x	y	$x - \bar{x}$	$y - \bar{y}$	$(x - \bar{x})(y - \bar{y})$	$(x - \bar{x})^2$
1	1	14	-1	-1	1	1
9	3	24	1	4	4	9
4	2	18	0	-2	-2	4
1	1	17	-1	-3	3	1
9	3	27	1	7	7	9
<u>24</u>	<u>10</u>	<u>20</u>	<u>0</u>	<u>0</u>	<u>0</u>	<u>40</u>

$$m = 20/4 = 5$$

$$\Rightarrow \bar{y} = m\bar{x} + c \Rightarrow 20 = 5 \times 2 + c \Rightarrow c = 10$$

$$\hat{y} = mx + c \Rightarrow \boxed{\hat{y} = 5x + 10}$$

$$\frac{5 \times 220 - (20 \times 10)}{5 \times 24 - 100}$$

For  $R^2$

	$(y - \bar{y})$	$\hat{y}$	$\hat{y} - \bar{y}$	$(\hat{y} - \bar{y})^2$
36	15	5	8	64
16	25	5	0	25
4	20	0	0	0
9	15	5	5	25
49	25	5	5	25
				<u>100</u>
		<u>114</u>		

$$R^2 = \frac{\sum (\hat{y} - \bar{y})^2}{\sum (y - \bar{y})^2}$$

$$\therefore R^2 = \frac{100}{114} = 0.877 = 88\%$$

### Multiple Linear Regression :

$$Y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \dots + \beta_n x_n$$

sums:

house size	$x_1$	$x_2$	$y$
1	1500	3	800
2	1600	3	220
3	1700	3	340
4	2100	4	980
5	2300	4	500

$\Sigma x_1$	$\Sigma x_2$	$\Sigma y$
8500	15	3000
$\Sigma x_1^2$	$\Sigma x_2^2$	$\Sigma xy$
182500	45	12500
$\Sigma x_1^3$	$\Sigma x_2^3$	$\Sigma x_1 x_2$
20475000	3375	101250
$\Sigma x_1^4$	$\Sigma x_2^4$	$\Sigma x_1^2 x_2$
256250000	15125	1012500

$$\boxed{0.1 \times 2 = 0.1 \times 0 + 3000 = 3000}$$

Find regression line & test the significance :

17) (hours studied)

	$x$	$y$	$x - \bar{x}$	$y - \bar{y}$	$(x - \bar{x})(y - \bar{y})$	$(x - \bar{x})^2$
1		50	-2	-13	26	4
2		55	-1	-8	8	1
3		65	0	2	0	0
4		70	1	7	7	1
5		75	2	12	24	4
	<u>3</u>	<u>63</u>			<u>65</u>	<u>10</u>

$$m = \frac{65}{10} = 6.5$$

$$Y = 6.5x + c \quad \text{For } c \Rightarrow 63 = 6.5(3) + c \\ \Rightarrow c = 43.5$$

$$\Rightarrow \boxed{Y = 6.5x + 43.5}$$

$$\text{if } m = 7 \text{ is taken} \Rightarrow \boxed{\hat{Y} = 42 + 7x} \Rightarrow (\text{Take this})$$

Test the significance

Null hypothesis  $H_0: \beta_1 = 0$  (No relationship b/w  $x$  &  $y$ )

Alternate hypothesis  $H_1: \beta_1 \neq 0$  (Same value of  $x$  used to predict  $y$ )

$$\text{calculate std error of slope } (\beta_1) = \frac{s}{\sqrt{\sum(x_i - \bar{x})^2}}$$

$s \rightarrow \text{std dev of residuals}$

Sum of square of residuals,  $\text{SSR} = \sum(y - \hat{y})^2$

$$\text{SSR} = 1^2 + 1^2 + 2^2 + 0 + 2^2 \quad (y - \hat{y} = 1 - 1.2 = 0 - 2) \\ = 10$$

$$\text{residual variance } s^2 = \frac{\text{SSR}}{n-2} = \frac{10}{5-2} = 3.33$$

$$\text{std dev, } s = \sqrt{3.33} = 1.825$$

$$SE(\beta_1) = \frac{1.825}{\sqrt{10}} \left( \frac{s}{\sqrt{\sum(x-\bar{x})^2}} \right) = 0.577$$

$$t - \text{statistic} : t = \frac{\beta_1}{SE(\beta_1)} = \frac{7}{0.577} = 12.13$$

$t \leq \text{table value} \rightarrow \text{accept } H_0$

$t > \text{table value} \rightarrow \text{reject } H_0$

$\therefore \text{table value for } \alpha = 0.05 \approx 3.182$

$$12.13 > 3.182$$

$\therefore \text{Reject } H_0$

x	y	$x - \bar{x}$	$y - \bar{y}$	$(x - \bar{x})(y - \bar{y})$	$(x - \bar{x})^2$	$\hat{y}$	$y - \hat{y}$
2	81	-9	-11.7	+105.3	81	85.09	
4	93	-7	0.3	-2.1	49	86.78	
6	91	-5	-1.7	8.5	25	88.47	
8	97	-3	4.3	-12.9	9	90.16	
10	85	-1	-7.7	7.7	1	91.85	
12	87	1	-5.7	-5.7	1	93.54	
14	89	3	-3.7	-11.1	9	95.23	
16	99	5	6.3	31.5	25	96.92	
18	102	7	9.3	65.1	49	98.61	
20	103	9	10.3	92.7	81	100.308	
<u>II</u>	<u>92.7</u>			<u>279</u>	<u>330</u>		

$$(\beta_1) = \frac{279}{330} = 0.8454$$

$$\hat{Y} = \beta_1 \bar{x} + c \Rightarrow 92.7 = 0.8454 \times 11 + c$$

$$c = 83.4$$

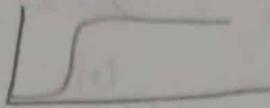
$$\boxed{\hat{Y} = 0.8454x + 83.4}$$

## Logistic Regression:

statistical model / method used to model the relationship b/w one or more independent variable and binary dependent variable.

### Binary classification:

e.g.: Predicting whether patient has a disease or not whether loan will be approved.



$$\text{Sigmoid func, } \hat{P} = \frac{e^x}{1+e^{-x}}$$

$$Z = \beta_0 + \beta_1 x_1 + \dots + \beta_n x_n$$

calculate the probability if a student studied 34 hrs he is pass or fail?

hrs X	Pass/Fail Y
29	0
15	0
33	1
28	1
39	1

$$\text{Probability prediction } \hat{P} = \frac{1}{1+e^{-Z}}$$

Soln 1 (Logistic Regression)  
Given  $x = 34 \text{ hrs}$

$$\text{odd}(z) = -64 + 2 \times \text{hrs.} + 0.3$$

$$\text{odd}(z) = -64 + 2 \times 34$$

$$\beta_0 = 0.5$$

$$\therefore \hat{P} = \frac{1}{1+e^{-4}} \Rightarrow \frac{1}{1+0.0183} = 0.98 \quad \beta_1 = 28$$

$$\therefore \boxed{\hat{P} = 0.98} \quad (98\% \text{ pass percentage})$$

### Soln 2 (Linear Regression):

X	Y	$X - \bar{X}$	$Y - \bar{Y}$	$(X - \bar{X})(Y - \bar{Y})$
29	0	0.2	-0.6	-0.12
15	0	-13.8	-0.6	8.28
33	1	4.2	0.4	1.68
28	1	-0.8	0.4	-0.32
39	1	10.2	0.4	4.08
$\bar{X}$	$\bar{Y}$			
28.8	0.6			13.6

Type 2 sum:  $\hat{P} = 95\% = 0.95$  (Chosen) Ans?

$$\frac{1}{1+e^{-z}} = 0.95 \Rightarrow \frac{1}{0.95} = 1+e^{-z}$$

$$\Rightarrow 0.052 = e^{-z} \Rightarrow \ln(0.052) = \ln(e^{-z})$$

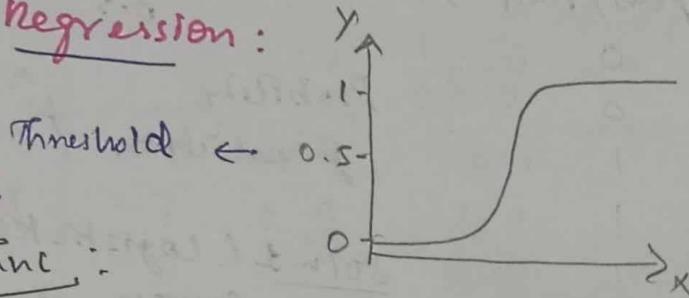
$$\Rightarrow -z = \ln(0.052) \Rightarrow z = \underline{\ln(e^{-z})}$$

$$\Rightarrow \text{Arrhen} \quad z = -64 + 2^* \text{hrs}$$

$$\Rightarrow \frac{2 \cdot 95 + 64}{2} = \text{hrs} \Rightarrow \boxed{\text{hrs} = 33.475 \text{ hrs}}$$

33.475 hrs required to get minimum 95% probability

### Logistic Regression:



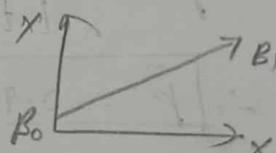
[For theory]

Derivation of sigmoid func:

logit (odds) =  $\frac{\text{Probability of an event happening}}{\text{Probability of an event not happening}}$

↓  
odds of success

$$\text{logit function } \in \text{Odds } (\theta) = \frac{p}{1-p}$$



$$Y = B_0 + B_1 x$$

$$\log\left(\frac{p(x)}{1-p(x)}\right) = B_0 + B_1 x$$

Take log on odds

$$\log\frac{p(x)}{1-p(x)} = B_0 + B_1 x$$

exponentiating both sides

$$e^{\ln \frac{p(x)}{1-p(x)}} = e^{B_0 + B_1 x}$$

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

Al\_Nav\_91\_9707771111  
Al\_Nav\_91\_9707771111@imauveronic.com

Let  $y = e^{\beta_0 + \beta_1 x}$

then  $\frac{P(x)}{1-P(x)} = y$

i.e)  $P(x) = y / (1 - P(x))$

$$P(x) + y(P(x)) = y \Rightarrow P(x) = \frac{y}{1+y}$$

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

( $\div$  by  $e^{\beta_0 + \beta_1 x}$  on numerator & denom)

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

Sigmoid func i.e  $P(x) = \frac{1}{1 + e^{-x}}$

### Multiple Logistic Regression :

Consider Dataset

	Age ( $x_1$ )	BMI ( $x_2$ )	Blood Pressure ( $x_3$ )	Diabetes ( $y$ )
women,				
$\beta_0 = 0.67$	45	28	140	1
$\beta_1 = 0.58$	50	25	135	0
$\beta_2 = 0.47$	40	30	150	1
$\beta_3 = 0.36$	35	24	125	0

Predict for Age 25 BMI 20 BP 120 Diabetes = ?

Given  $x_1 = 25$   $x_2 = 20$ ,  $x_3 = 120$

$$Z = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \dots + \beta_n x_n$$

1.  $Z = 0.67 + 0.58 \times 25 + 0.47 \times 20 + 0.36 \times 120$

$$Z = 67.77$$

$$\hat{P} = \frac{1}{1+e^{-2}} = \frac{1}{1+0} = 1$$

$\therefore$  He is diabetic for given values

### Logistic Regression

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

Use maximum likelihood estimates for estimating the coefficient

likelihood function:

Product of individual probabilities

$$L(\beta_0, \beta_1) = \prod_{i=0}^n P(Y_i=1/X_i)^{Y_i} \times (1-P(Y_i=1/X_i))^{1-Y_i}$$

$Y_i$  - actual outcome

log likelihood

$$\ell(\beta_0, \beta_1) = \sum_{i=1}^n [Y_i \log P(Y_i=1/X_i) + (1-Y_i) \log (1-P(Y_i=1/X_i))]$$

start with initial coefficients let  $\beta_0 = 0$   $\beta_1 = 0.5$

calculate predicted probability

i.e.) for  $X = 2$

$$P(Y=1/X=2) = \frac{1}{1+e^{-(0+0.5 \times 2)}} = 0.731$$

$$P(Y=1/X=4) = \frac{1}{1+e^{-(0+0.5 \times 4)}} = 0.88$$

$$P(Y=1/X=6) = \frac{1}{1+e^{-(0+0.5 \times 6)}} = 0.95$$

$$P(Y=1/X=8) = \frac{1}{1+e^{-(0+0.5 \times 8)}} = 0.92$$

$$P(Y=1/X=10) = \frac{1}{1+e^{-(0+0.5 \times 10)}} = 0.94$$

compare loglikelihood & actual  $y = [0, 0, 1, 1]$

given Predict whether accident P/F based on  
no. of hrs shedded

Hours X	2	4	6	8	10
P/F Y	0	0	1	1	1

for  $X = 2, Y = 0$

$$l(\beta_0 + \beta_1) = \log(1 - 0.73) =$$

$$\text{for } X = 4, Y = 0 \\ l(\beta_0 + \beta_1) = \log(1 - 0.88) =$$

for  $X = 6, Y = 1$   $\rightarrow$  (For  $X=1$  (Don't take as 1 - value))  
 $l(\beta_0 + \beta_1) = \log(0.95)$

update the coefficient:

$$\beta_0 \text{ new} = \beta_0^{\text{old}} + \alpha * \frac{\partial l(\beta_0, \beta_1)}{\partial \beta_0}$$

$$\beta_1 \text{ new} = \beta_1^{\text{old}} + \alpha * \frac{\partial l(\beta_0, \beta_1)}{\partial \beta_1}$$

$\alpha \rightarrow$  learning rate  $\partial \beta_0, \partial \beta_1 \rightarrow$  partial derivatives of  $\beta_0 \& \beta_1$

### Wald Test:

→ determine if the set of independent variables are significant in a logistic regression model

$$W = \frac{\hat{\beta}_j}{\text{SE}(\hat{\beta}_j)} \quad \hat{\beta} \rightarrow \text{estimated coefficient}$$

$\text{SE} \rightarrow \text{std err}$

→ can be applied for every coeff

### LR (Log Likelihood Ratio Test)

$$LR = -2 \times (\text{lo} - \text{lr})$$

↓  
null model

↓  
without the predictor

↓  
full model

↓  
with predictor

Score test :

$$S = U(\hat{\beta}_0)^T \cdot I(\hat{\beta}_0) \cdot U(\hat{\beta}_0)$$

Logistic Regression model to predict whether a customer will purchase a product  $y=1$  and NO  $y=0$  on their income  $x$ .  $\hat{\beta}_1 = 0.8$   $SE(\hat{\beta}) = 0.2$

Test the null hypothesis

$$H_0 : \beta_1 = 0 \quad H_1 : \beta_1 \neq 0$$

Wald test :

$$w = \frac{\hat{\beta}}{SE(\hat{\beta})} = \frac{0.8}{0.2} = 4$$

Compare with normal distribution under  $H_0$ , two tailed test. Calculate  $\phi$ -value using  $z$  distribution

$$P = 2 \times (P(Z > 4)) \approx 2 \times 0.00003 = 0.00006 < 0.5$$

$\therefore$  P value is very small, reject  $H_0$  (which means  $\beta_1$  significantly affects  $y$ ).

LR test:

Predict whether a patient will develop disease based on age

$$\text{Full model: } \log \frac{P(y=1)}{P(y=0)} = \beta_0 + \beta_1 x$$

$$\text{Null model: } \log \frac{P(y=1)}{P(y=0)} = \beta_0$$

Consider full & null models are -120 & -140

$$\begin{aligned} \therefore LR &= -2 \times (-140 - (-120)) \\ &= -2 \times -20 \\ &= 40 \end{aligned}$$

test statistic follows a chi-square distribution with degree of freedom equal to the difference in the no. of parameters between two models.

P-value calculated using chi-square distribution with ~~df~~  $df = 1$

$$p = P(X^2 > 4.0) \approx 0.0001$$

$\therefore$  p-value is very small. Reject  $H_0$  age significantly contributes to the model.

### Score test:

- $X_1 \rightarrow$  education level ;  $X_2 \rightarrow$  income of owning a car ( $Y$ )
- Score statistics for adding  $X_1 = 5.8$ , for adding  $X_2 = 12.1$ .
- Follow chi-square distribution critical value of  $X^2$  at a 5% significance is 3.84
- $5.8 > 3.84 \Rightarrow$  reject  $H_0$
- $12.1 > 3.84 \Rightarrow$  reject  $H_0 \therefore$  Both  $X_1$  &  $X_2$  should be included in the model.

### Problem:

Age	Monthly Premium (MP)	Year as customer	Renewed
20	150	2	Yes
45	200	5	No
50	180	4	Yes
60	220	6	No

$$\hat{B}_0 (\text{Intercept}) = 3.5 \quad \hat{B}_1 = -0.03 \quad \hat{B}_2 = -0.02 \quad \hat{B}_3 = 0.3$$

- i) Fit the logistic reg.
- ii) Evaluate the model performance.
- iii) Predict whether renewed for Age = 35, MP = 160, Years = 3

i) Logistic Regr Model fit:

$X_1 \rightarrow \text{Age}$        $Y \rightarrow \text{Renewed}$   
 $X_2 \rightarrow \text{APP}$   
 $X_3 \rightarrow \text{Year}$

$$Y = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \beta_3 X_3$$

$$Y = 3.5 + 0.05 X_1 + 0.02 X_2 + 0.3 X_3$$

ii) For evaluating model performance:

iii) Calculating predicting probability for each  $x_1, x_2, x_3$   
 to compare with already given Renewed (Yes/No).

$X_1$	$X_2$	$X_3$	Given Y	$\hat{P}$
30	150	2	Yes	$3.5 - 0.05 \times 30 - 0.02 \times 150 = Z \Rightarrow \frac{1}{1+e^{-Z}} = 0.4$
45	200	5	No	$0.22 + 0.3 \times 2 = 0.4 = Z \Rightarrow \frac{1}{1+e^{-Z}} = 0.4$
50	180	4	Yes	0.198
60	220	6	No	0.109

Prediction Probability,  $\hat{P} = \frac{1}{1+e^{-Z}}$  where  $Z = \beta_0 + \beta_1 X_1 + \beta_2 X_2$

Here every  $\hat{P}$  is  $< 0.5$  so No but in given data Yes is there

III) For  $X_1 = 35, X_2 = 160, \text{Year } X_3 = 3$

$$\hat{P} = \frac{1}{1+e^{-(3.5 - 0.05 \times 35 - 0.02 \times 160 + 0.3 \times 3)}} = 0.366$$

.. Prediction for given input as Not Renewed.

## Multiple Logistic Regression:

### Forward selection & Backward selection methods:

#### Forward selection:

- Start with null model i.e. only intercept
- add variables based on their significance
- at each step variable improve the model
- continue adding until no more significant variable is found.

## Backward Selection (Elimination) :

- start with all possible predictor variables in the model
- gradually remove predictor that do not contribute for prediction.
- continue removing until all remaining variables are statistically significant.

Example :

Predict whether student will pass / fail based on hours of study, attendance, GPA

where hours of study, attendance are more significant than GPA.

Interpretation of Parameters :

Log odds & Odds Ratio :

coefficients  $\beta_0, \beta_1, \dots, \beta_k$  change outcome of one unit change in prediction by holding other variables constant.

Odds ratio is calculated as

$$e^{\beta_i}$$

e.g.  $\beta = 0.5$  for hrs of study.

$$\therefore \text{odds ratio} \Rightarrow e^{0.5} = 1.65$$

$\Rightarrow$  each additional hrs of study, odds of passing the exams are 1.65 times higher, while holding others constant.

Multi nominal Logistic Regr : where the category is more than 2 like for eg : High, Medium, Low. This model holds 1 category for reference & performs LR for other 2,

Significance of coefficients :

$\Rightarrow$  P-value is used to access the significance and uncertainty of the coefficient estimates.

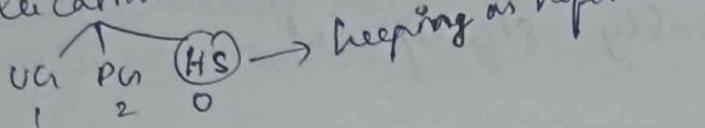
P-value < 0.05 reject hypothesis  
 ≥ 0.05 accept "



## Categorical Data Analysis :

Multinomial Logistic Regression for more than 2 categories.

e.g.: Education level



1 2 0

Sum: Perform Forward & Backward selection:

Age	Monthly income	Customer satisfaction	Retained
25	3000	7	0
35	4000	9	1
45	5000	8	1
30	3500	6	0
40	4500	7	1

→ P.Md odds ratio  
 ↘ Curve some write-ups (theory)

$$\text{Given } \beta_0 = -5.5, \beta_1 = 0.002, \beta_2 = 0.8, \beta_3 = 0.03$$

(Based on odds ratio give some explanation of eqns in each step )

## General Linear Models:

② Distribution & Loss function

Theory

GLM

log reg, Prod, Bwd

(PPT) Mult log reg,

↳ what is GLM  
 ↳ expo family distributions & eqns

↳ basic descri for each member

↳ Link functions

↳ Derivation for ref (not for exam)

Sum & Algo

Decision tree

KNN (PPT)

Name Bayes (PPT)

Linear Regression

$$\$ \times 10 = 40$$

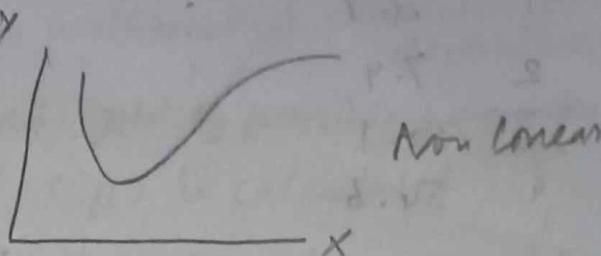
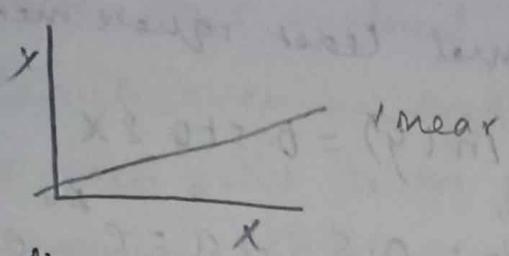
other 5

$$1 \times 10 = 10$$

$$\frac{50}{50}$$

↳ Derivation of sigmoid function

## Non Linear Regression :



Equation :

$$Y = f(X, \beta) + \epsilon$$

$f(X, \beta) \rightarrow$  Regression function  $X \rightarrow$  independent variable

$\beta \rightarrow$  parameter that model aims to estimate

$\epsilon \rightarrow$  error term

### i) Polynomial Regression :

uses higher order polynomial functions

$$Y = \beta_0 + \beta_1 x + \beta_2 x^2 + \dots + \beta_n x^n$$

### ii) Logistic Regression :

→ Binary classification > probability estimation > non-linear relationships

### iii) Exponential Regression :

follow exponential growth or delay pattern

$$Y = a e^{(\beta x)}$$

Linearization (Techniques) Transform : → to transform non linear eqn to linear

#### i) log :

$$\text{eg: } Y = a e^{bx}$$

$$\ln(Y) = \ln(a) + bx$$

$$\therefore Y = \beta_0 + \beta_1 x$$

$$\Rightarrow \beta_0 = \ln(a), \beta_1 = b$$

x	y	$\ln(y)$
1	2.7	
2	7.4	
3	20.1	
4	54.6	

Find the  $\beta_0$  &  $\beta_1$  using  
normal least square method.

$$\Rightarrow \ln(y) = \beta_0 + \beta_1 x$$

$$\Rightarrow \beta_0 = 0.5 \quad \therefore a = e^{\beta_0} = e^{0.5}$$

$$\beta_1 = 0.8 \quad b = 0.8$$

∴ Final Non linear eqn 
$$y = 1.648 e^{0.8x}$$

### II) reciprocal :

$$\frac{1}{y} = \frac{q}{P} \frac{1}{x} + \frac{1}{P}$$

### III) Polynomials :

$$y = \beta_0 + \beta_1 x + \beta_2 x^2 + \dots + \beta_n x^n$$

### (Theory)

#### Iterative Procedure for NLS (Non Linear Least square)

> Q is used to find best fit parameters.

> Aims to minimize sum of squared residuals (Res)

i.e.)  $S(\beta)$

$$S(\beta) = \sum_{i=1}^n (y_i - f(x_i, \beta))^2$$

$y_i \rightarrow$  observed value  
 $f(x_i, \beta) \rightarrow$  model predicted

#### Common Derivative Procedures :

i) Grid search

value using  $\beta$   
 $S(\beta) \rightarrow$  sum of squared  
residuals

ii) Newton Raphson

iii) steepest descent

iv) Marquardt's Method

## i) Grid search :

- > brute force method & No mathematical calculation for derivatives
- > Search over a predefined grid of possible parameter values
- > At each point in the grid  $s(\beta)$  is calculated.
- > Minimum  $s(\beta)$  is chose & less accurate compared to others.

## ii) Newton Raphson :

- > Iterative Procedure & involves mathematical derivatives
- & Use Second Order Info (Hessian Matrix of second derivatives)
- for parameter estimation
- > fast & converges rapidly near the solution.

### Steps :

1. Start with initial guess of  $\beta_0$

2. Update 
$$\beta_{\text{new}} = \beta_{\text{old}} - H^{-1}(\beta_{\text{old}}) \cdot \nabla s(\beta_{\text{old}})$$

$H(\beta_{\text{old}})$  - Hessian Matrix  $\nabla s(\beta_{\text{old}})$  - gradient (1<sup>st</sup> derivative)

3. Repeat until converges

## iii) steepest Descent (Gradient) :

- > First Order optimization technique

\* Start with initial guess for  $\beta_0$

\* Update parameter by moving in direction opposite to gradient 
$$\beta_{\text{new}} = \beta_{\text{old}} - \alpha \cdot \nabla s(\beta_{\text{old}})$$

$\alpha \rightarrow$  learning rate ,  $\nabla s(\beta_{\text{old}}) \rightarrow$  gradient of objective function  $s(\beta)$

\* continue updating until no change in the objective function.

## N) Margquardt's method :

- > hybrid technique that combines both Newton Raphson and steepest descent.
- > start with initial guess  $\beta_0$
- > update  $\beta_{\text{new}} = \beta_{\text{old}} - (\mathbf{J}^T \mathbf{J} + \lambda \mathbf{I})^{-1} \mathbf{J}^T \mathbf{y}$
- $\mathbf{J} \rightarrow$  Jacobian Matrix       $\mathbf{y} \rightarrow$  Vector of residuals
- $\lambda \rightarrow$  damping parameter       $\mathbf{I} \rightarrow$  Identity matrix
- $\lambda$  is large  $\rightarrow$  steepest descent
- $\lambda$  is small  $\rightarrow$  Newton Raphson
- > Adjust  $\lambda$ .

## Semi Parametric Regression :

- > combine both parametric & non parametric models.
- > Predicting income ( $Y$ ) based on age ( $x_1$ ) and education level ( $x_2$ )

$$Y = \beta_0 + \beta_1 x_1 + f(x_2) + \epsilon$$

## Additive Regression :

- > Response variable depend on each predictor through an independent smooth function.

$$Y = \beta_0 + f_1(x_1) + f_2(x_2) + \epsilon$$

e.g.: House price      Square footage      Distance from city

## Non-Parametric:

- > It do not assume any functional form
- > Fully data driven

non-parametric :

- kernel regression - weighted avg of nearby data
- smoothing splines - smooth curve through data
- Local Regression - Perform Regression locally using nearby data.

eg: bike rental demand vs temperature.

UNIT-4Covariance :

$$\text{cov}(x, y) = \frac{\sum (x - \bar{x})(y - \bar{y})}{n} \quad | \quad \text{Correlation: } \text{cor} = \frac{\text{cov}(x, y)}{\sigma_x \cdot \sigma_y}$$

Problem 1:

$$\text{Prod 1} \quad 120 \quad 125 \quad 130 \quad 135 \quad 140 \quad 145 \quad 150 \quad \sum x = 823 \quad \bar{x} = 117.57$$

$$\text{Prod 2} \quad 115 \quad 120 \quad 125 \quad 130 \quad 135 \quad 140 \quad 145 \quad \sum y = 910 \quad \bar{y} = 130$$

~~Since  $E(x \cdot y) = \bar{x} \cdot \bar{y}$ ,  $\text{cov}(x, y) = 0$~~ ,  $\text{cor} = 1$

	day	1	2	3	4	5
x	A	320	350	410	490	510
y	B	380	440	520	590	610
$x - \bar{x}$		-96	-66	-6	74	94
$y - \bar{y}$		-128	-68	12	82	102

$$\bar{x} = 416 \quad \bar{y} = 508$$

$$x - \bar{x} \quad -96 \quad -66 \quad -6 \quad 74 \quad 94$$

$$y - \bar{y} \quad -128 \quad -68 \quad 12 \quad 82 \quad 102$$

$$3) \quad \begin{array}{ccccccc} x & 7 & 8 & 9 & 10 & 11 & 12 & 9.5 \\ y & 1 & 2 & 3 & 4 & 5 & 6 & 3.5 \end{array}$$

$$\sigma_x = \frac{17.5}{6}$$

$$\therefore \text{cov}(x, y) = \frac{17.5}{6}$$

$$x - \bar{x} \quad -2.5 \quad -1.5 \quad -0.5 \quad 0.5 \quad 1.5 \quad 2.5$$

$$y - \bar{y} \quad -1.5 \quad -0.5 \quad 0.5 \quad 1.5 \quad 2.5$$

$$= 2.92$$

$$(x - \bar{x})(y - \bar{y}) \quad 8.25 \quad 2.25 \quad 0.25 \quad 0.25 \quad 1.25 \quad 0.25 \quad 17.5$$

$$\text{cor} = \frac{2.92}{\sqrt{17.5} \times \sqrt{17.5}}$$

$$(x - \bar{x})^2 \quad 49 \quad 36 \quad 9 \quad 4 \quad 1 \quad 0 \quad 17.5$$

$$(y - \bar{y})^2 \quad 1 \quad 4 \quad 9 \quad 16 \quad 25 \quad 36 \quad 17.5$$

$$\text{cor} = 1$$

## Features of covariance :

$$\& \text{cov}(x, a) = 0 \quad \Rightarrow \text{cov}(x, x) = \text{Var}(x)$$

$$\& \text{cov}(x, y) = \text{cov}(y, x) \quad \text{cov}(x+y, z) = \text{cov}(x, z) + \text{cov}(y, z)$$

$$\& \text{cov}(ax, by) = ab \text{cov}(x, y)$$

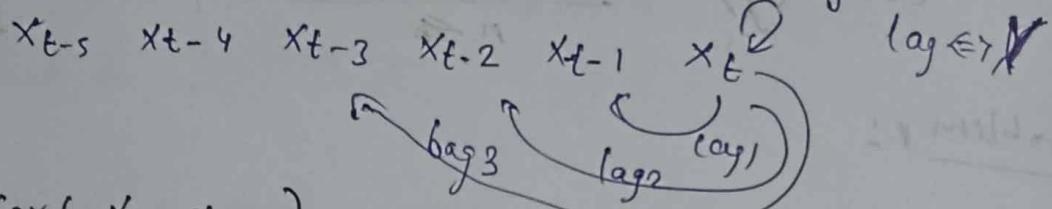
& covariance of independent variable is zero

## Auto covariance :

For single variable different times relation,

e.g.: 2017 2018 2019 2020 2021 2022

7 8 9 10 11 12 lag<sup>0</sup>



$$Y_1 = \text{cov}(x_t, x_{t-1}) = E(x_t, x_{t-1}) - E(x_t) \cdot E(x_{t-1})$$

$$Y_2 = \text{cov}(x_t, x_{t-2})$$

$$Y_3 = \text{cov}(x_t, x_t) = \text{Var}(x_t)$$

Time series :  $\rightarrow$  Auto cov, cor, ARMA, ARIMA

Components : 1. Trend 2. Seasonality 3. Cyclic 4. Irregular

stationarity of stochastic process

stationary  $\rightarrow$  mean, var, autocorr are all constant over time

all constant  $\rightarrow$  strictly stationary

only mean, var constant  $\rightarrow$  weakly stationary

## Yule-Walker Estimation:

$\Gamma \rightarrow$  covariance matrix  
 In order to find correlation, 1st have to find covariance  
 consider AR eqn of  $p=2$

$$y_t = \phi_1 y_{t-1} + \phi_2 y_{t-2} + \varepsilon_t$$

(X)  $y_{t-k}$  on both sides

$$y_t y_{t-k} = \phi_1 y_{t-k} y_{t-1} + \phi_2 y_{t-2} y_{t-k} + \varepsilon_t y_{t-k}$$

Expectation on both sides

$$E(y_t y_{t-k}) = \phi_1 E(y_{t-k} y_{t-1}) + \phi_2 E(y_{t-2} y_{t-k}) + E(\varepsilon_t y_{t-k})$$

$$Y_k = \phi_1 Y_{k-1} + \phi_2 Y_{k-2}$$

$\therefore$  by  $r_0$  (variance)

$$\frac{Y_k}{r_0} = \phi_1 \frac{Y_{k-1}}{r_0} + \phi_2 \frac{Y_{k-2}}{r_0}$$

$$\Rightarrow P_k = \phi_1 P_{k-1} + \phi_2 P_{k-2} \quad P \rightarrow \text{correlation}$$

$$\text{Let } k=1 \Rightarrow P_1 = \phi_1 P_0 + \phi_2 P_{-1}$$

$$= \phi_1 P_0 + \phi_2 P_1 \quad (P_{-1}) = (P_1)$$

$$\Rightarrow P_1 = \phi_1 + \phi_2 P_1 \rightarrow ① \quad P_0 \text{ (no contribution)} = 1$$

$$P_1 - \phi_2 P_1 = \phi_1$$

$$P_1(1 - \phi_2) = \phi_1 \Rightarrow P_1 = \frac{\phi_1}{1 - \phi_2}$$

$k=2$

$$\Rightarrow P_2 = \phi_1 P_1 + \phi_2 P_0$$

$$P_2 = \phi_1 P_1 + \phi_2$$

$$P_2 = \phi_1 \left( \frac{\phi_1}{1 - \phi_2} \right) + \phi_2$$

$$= \frac{\phi_1^2 - \phi_2^2 + \phi_2}{1 - \phi_2}$$

If & only if  $k \neq 0$

This whole term  
can be neglected  
if not it is a standard deviation  
(e.g.) variance

## Problems for Unit - 2

### 1. Normalisation

Min Max, Z-score, decimal scaling,

- > Brinney method → Data cleaning
- > Dissimilarity (proximity matrix) for all data
- > Chi-square test
- > Wavelet transform → Data preprocessing



### Unit - 2 Problems

- > Apriori classification metrics ↴ Performance evaluation ↴
- > FPGrowth
- > Decision tree
- > NB
- > KNN
- > Linear Regr.

### Unit - 3 Problems

- > Testing of hypothesis
- > Linear Regr.
- > (Rare case) Logistic Regr.

### Unit - 4 Problem

- > Covariance & Correlation

## UNIT - 1

- > OLAP, OLTP - operation, details, diff b/w
- > Data Mining Architectures
  - ↳ Techniques
  - ↳ Task / Functionalities
  - ↳ Applications
- > Data Preprocessing (X) (X)
  - ↳ Chi-square, z-score (Write Problem  $\rightarrow$  binning)
  - ↳ Theory every problems

## > Attribute oriented Analysis (AOE)

> Dissimilarity matrix, Distance Matrix

## Association rule mining UNIT - 2 Problem (X)

- > APriori, FP Growth (X) (Steps shd be clear)
- > Classification (Write Procedure by our own)
  - Decision tree • Naive Bayes • KNN • Linear Models

## > Test of Hypothesis - UNIT - 3 (Theory X)

> Logistic Regr - Derivation, Diff b/w linear & logistic

> Fwd & Backwd  $\rightarrow$  Multi Logistic Regr

> GLM (Based on marks write the no. of distributions)

> Non-linear Regression (what is what)

> Semiparametric Regress.

> Additive

> Nonparametric ..

} Handwritten notes  
in GCR

## UNIT - 4

> Problem  $\rightarrow$  Covariance & Correlation (For univariate series Data)

> Distr of Time Series Data

> Estimation of ARMA  $\rightarrow$  Yule Walker, Max Likelihood (X)  
Least Square

Prescriptive analysis : with optimization  
Decision & Risk Analysis To decision tree.