

→ flow →  
→ ( ) see →

## Data Mining & Warehousing

Mining: Process of extracting something you are interested in from a collection

Warehousing: How you organise your data so that data mining can better utilise that organized data in an efficient manner

After 1990s, we can better share & utilise the data

Stream data → real time data (Eg- stock market)

### \* DBMS vs Data Mining

In data mining some predictions are involved.

3, 15, 3, 19, 27 → marks in course A

Using DBMS we can write SQL queries for performance, total marks etc. but not predict his ~~permo~~ & performance in another course next sem. But we can do that using Data Mining

"interesting" (slide 7)

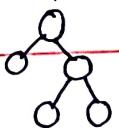
That satisfies your requirement

5, 15, 25, 10, 20 → data

$$x \rightarrow f \rightarrow y = f(x)$$

5, 10, 15, 20, 25 marks in ↑ order → Inform'

↓ create some model

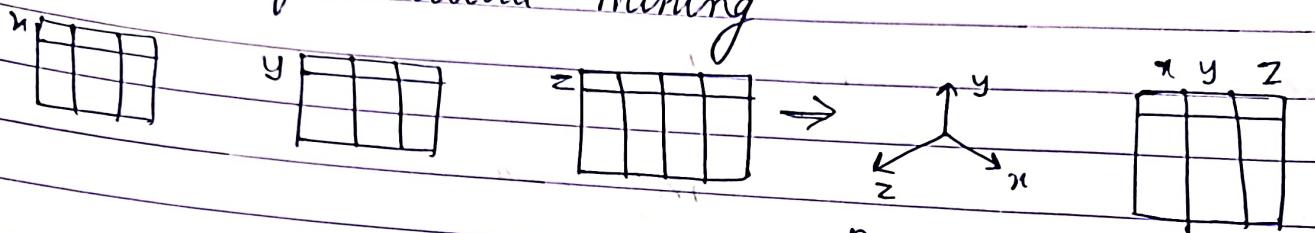


Knowledge → using accumulated inform'

## \* KDD Process

Data available on diff servers may have a heterog. org". So we need to integrate that into homog. org" (data integ")

In data warehouse data is stored in a format suitable for data mining



Web Mining → extracting data from web

Cube → Multidimensional data organisation  
(Homogeneous org", stored in warehouse)

"Present" of results: Better way to represent the data  
so some conclusion can be drawn

I.\* Preprocessing: Activity performed before execution of your algorithms, performed on ip data

## Lab

- 1. Preprocessing
  - < cleaning
  - missing value
- 2. Feature selection
- 3. Algo
- 4. Result → Acc
  - Preci
  - Recall
  - F1 score

.arff (attribute relation file format)

- 3 section
  - 1. @relation
  - 2. @attribute
  - 3. @ data

- Visualization → ROC
  - AUC

-x-

- Data is available in diff. db's in heterog. form.
- so we need to transform it into a standard form, as per the requirement. Depends on algo, diff types of Data Integration is there

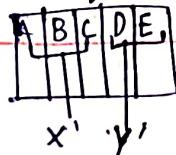
- Normalization: 1-100 marks range → 1-10 range req. by algorithm. How you represent data values in a different range, which is required by the algorithm

- Feature Selection: Attribute in a table → Feature/variable, dimension (linear algebra). ↑ features ↑ complexity.

Redundancy in data is common. Ways to ↓ size of data:

→ Remove redundant / irrelevant features (Feature selec")

→ Feature extraction: Using initial features, design some new features more relevant for your application



→ Issue:

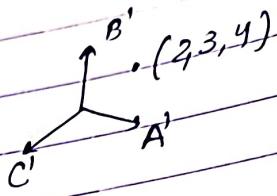
F.S has the probability of information loss, unless you remove only irrelevant/redundant ones. But in reduction, inform<sup>n</sup> loss should not be there → performance of model affected

Dimensionality reduction:

Dimension → Attribute

Due to complexity of algo., we need to ↓ dimension of dataset

A'	B'	C'
-	-	-
C'	A'	B'



A	B	C	D
1	1	1	1
x'			

3D "inform" transformed into 1D, → D.R  
After feature selection → transform must be w/o inform<sup>n</sup> loss

### Mechanisms

- 1) PCA
- 2) LDA

## II \* Data Mining

① Pattern discovery: Based on patterns we take business decisions (e.g. customer purchasing behaviour)

### ② Association & Correlation:

Transactional dataset

Transac <sup>n</sup>	item1	item2	item3
T <sub>1</sub>	Bread	Milk	Soap
T <sub>2</sub>	Bread	Milk	
T <sub>3</sub>	Soap	Bread	
T <sub>4</sub>	Bread	Milk	

Analysis → the customers buys bread along with milk, and in regular intervals. ∵ Bread & Milk are associated to each other



Transac <sup>n</sup>	Bread	Milk	Soap
T <sub>1</sub>	1	1	1
T <sub>2</sub>	1	1	0
T <sub>3</sub>	1	0	1
T <sub>4</sub>	1	1	0

75% of time he buys bread + milk  
business prepare their inventories  
in diff. places based on this data.

- Actual dataset → considering every purchased item as an attribute

Correlation - Milk → Bread

$$-1 \leq \rho \leq 1$$

↑ no corr. ↑ +ve correlation

Mechanism: PCC (Pearson correlation coefficient)

Numerical quantification of associa<sup>n</sup> b/w 2 attributes/variables

- ③ Classification: Classify samples in different categories

2 classes → Binary classification

> 2 " → Multiclass "

T	B	M	S	class	
				+	valid transac <sup>n</sup>
				+	
				-	invalid
				+	

Binary

#### ④ Clustering:

"Classific" works on labelled data and  
"Clustering" "unlabelled data

labelled  $\rightarrow$  class attribute present in data

class label not known  $\rightarrow$  clustering

$\Rightarrow$  Classific algos  $\xrightarrow{\text{called}}$  supervised learning algorithm

You learn in presence of supervisor (check & rectify)  
Learning w/ labelled dataset

$\Rightarrow$  Clustering  $\rightarrow$  Unsupervised ML algo

$\Rightarrow$  Hybrid  $\rightarrow$  Semi supervised ML algo

#### ⑤ Outlier analysis:

T	B	M	height	class
3.00				
5.3				
5.2				
6.0				
4.7				
6.5				
5.00				

height permissible limit 0-7  
on analysis you find 99%.

students have height b/w 0-6

Then 6.5 is treated as an outlier

Outlier: related to your data but are rare samples

Outlier analysis is useful in Credit Card system.

A B C D class

-	-	-	-	valid	Eg1 - 1-2% customers X paying bill on time
-	-	-	-	valid	→ invalid customers
-	-	-	-	valid	
-	-	-	-	valid	relevant customer ✓ (satisfied cond' for
-	-	-	-	invalid	becoming bank's customer) but X paid bill on time → rare sample

\*Outlier Analysis: How you identify (criteria) the rare samples ("identific" of rare by formulating some criteria)

Make business schemes → EMI offer, offers to pay bill on time

Eg2- Identifying unsuccessful packet transmissions, what can be the causes etc.

### III. Post processing:

If algo executed but result is unordered → need to represent in proper way

#### ① Pattern evaluation:

Transac"	B	M	S	3 Bread, Milk	2 Bread, Milk, Soap	patterns
T <sub>1</sub>	1	1	1			
T <sub>2</sub>	1	1	0			
T <sub>3</sub>	1	0	0			
T <sub>4</sub>	1	1	1			

All patterns may not be useful

Unuseful patterns  $\rightarrow$  space complexity  $\uparrow$

P. Evaluation: Mechanism with which only useful patterns are shown

Eg - pattern useful if frequency  $> 75\% \rightarrow$  threshold criteria set to select relevant and useful patterns

② Pattern selec"

③ " Interpretat"

④ " Visualization

### \* Dimensions of Data Mining

① Data to be mined: Many sources / types

stream  $\rightarrow$  video data  
(pos/coordinates)  
spatiotemporal  $\rightarrow$  temperature, GPS  
time-series  $\rightarrow$  stock price  
sequence  $\rightarrow$  first computer then speakers then pen drive  
 $\rightarrow$  DNA structure  
graph  $\rightarrow$  social networking  
 $\rightarrow$  lattice structure

Diff. DBs as per type of data

② Knowledge to be mined:

Characterization: "summariz" of your data related to the target class. Eg - invalid in credit card system

Eg - Those who got A+ in Data Mining, did they also get A+ in DBMS  $\rightarrow$  good students for data science

Discrimination: Using the summarisation, how we discriminate the class of with other classes

### Descriptive vs Predictive Data Mining:

- 1) Descriptive: With your data you are formulating the descriptions (Eg- assoc<sup>n</sup>/correl<sup>n</sup>/existence of patterns)  
What method we use to describe our data (Eg- Pattern anal., cluster ana, asso, correl<sup>n</sup>)
- 2) Predictive: On the basis of certain features we are predicting some other features. Eg- if one class label for a row is missing. Using some model ~~on~~ on the data we can predict the unknown class. so all ML algor → predictive)

## Ch-2: Getting to know your data

Matrix: Arrangement of data in rows or columns

Data matrix: Data in numeric form in matrix

$A_1$	$A_2$	$\dots$	$A_N$
:	:		
			$M \times N$

Document matrix

TF-IDF: Term Frequency - Inverse Document Frequency

Data Structured

unstructured

semi structured

Algos we study are mainly based on structured data

### \* Important Characteristics of Structured Data

1) Dimensionality: No. of attributes in the database. Issue  $\rightarrow$  How to tackle  $\uparrow$  dimen.

Curse of dimen. problem: If you  $\uparrow$  dimension of your database, complexity of your algo.  $\uparrow \rightarrow$  performance  $\downarrow$  memory  $\uparrow$  computation  $\uparrow$  need for structure  $\uparrow$

2) Sparsity:

Many terms can have freq. 1 or 0

$T_1 \quad T_2 \quad T_3 \quad T_4 \quad T_5 \quad T_6$  In sparse data many useful analysis  
 $\begin{matrix} 1 & 0 & 0 & 0 & 0 & 0 \\ 0 & \cancel{X} & \cancel{X} & \cancel{X} & \cancel{X} & 1 \end{matrix}$  cannot be done Eg- checking if 2 docu.  
 $0 \quad 0 \quad 0 \quad 0 \quad 0 \quad 0$  are similar, here you will think its same since most words have 0 frequency

Sol  $\rightarrow$  remove ~~terms~~ having frequency 0, then only you can get good result for similarity/dissimilarity

3) Resolution: Size of data  
 $M \times N$  (No. of tuples)

For good learning, sufficient samples should be there.

↓ samples → underfit

↑ " → overfit (also ↑ attributes)

Issue → optimum no. of samples

4) Distribution: Values corresponding to attributes

Eg- Gaussian dist" func"

Tendency of data, in which direc. its tilts to

Data Objects:

In database ~~not~~ all tuples are data objects

Attributes:

qualitative Nominal ( $=, \neq$ ) that can distinguish with posse

(categorical) ordinal ( $<, >$ ) tells us about arrangement

quantitative Interval (+, -) values & continuous → discrete in categories

(Numeric) Ratio ( $\times, /$ ) ratios well defined

No. of prop. of Bsn prop

Binary → 1 or 0 → 2 val

Gender → Binary

Habits of smoking

Cancer + only few ppl have issue of cancer

Autism being att occurrence & same as b.

Interval:

X True zero point

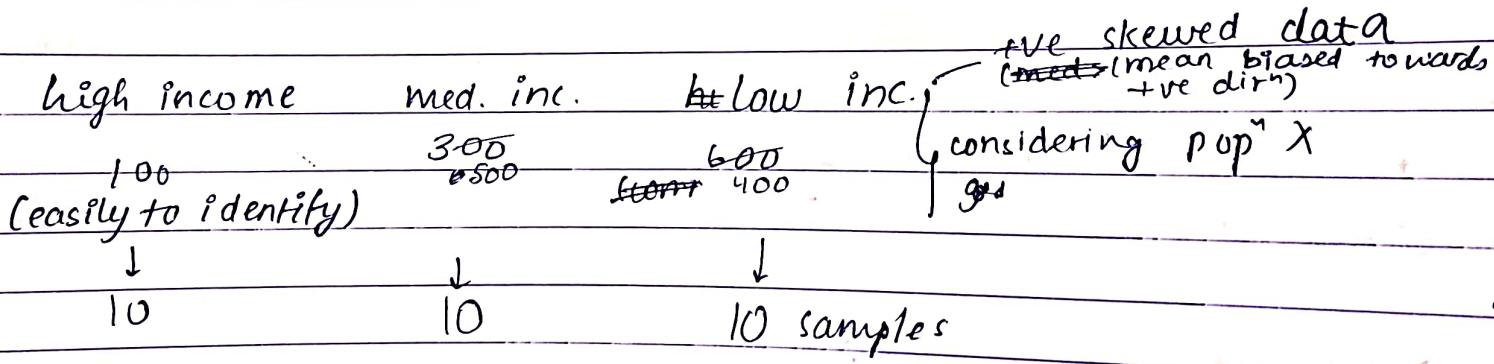
Multiplication factor exactly the same reference pt,  
well defined  $\rightarrow$  zero point

Discrete: Shirt size

Continuous: real values within in certain range

### \* Statistical QM

1) Mean (sample vs population)



Go w/ the ~~any~~ sample ~~not~~ ~~not~~ not the pop<sup>n</sup>

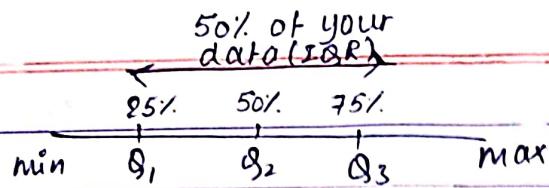
Tendency biased towards the ~~not~~ outlier values.

10, 12, 15, 20, 21, 180

outlier

most freq

-ve skewed data (mean biased towards -ve dirn)  
Ques how to vel-ve skew?



$Q_2 \rightarrow \text{median}$

IQR  $\rightarrow$  denote behaviour of 50% of your data  
 Outlier =  $1.5 \times \text{IQR}$  (standard way to identify outliers)

### \* Similarity & Dissimilarity (Proximity)

In your database if you want to know the similarity/dissim. of data objects in numerical form

	hair	weight	gender
Obj 1	black	32.5	1
Obj 2	red	16.2	0

Data Matrix: p attributes n data objects  
 Resolution  $\rightarrow n \times p$

$d(2,1)$  dissimilarity of 2nd object w/ 1st object  
 $[\text{Sim}(2,1) = 1 - d(2,1)]$

#### ① Nominal attributes:

$$\underline{\text{MI}}: d(i,j) = \frac{p-m}{p}$$

p  $\rightarrow$  no. of attributes

m  $\rightarrow$  matches across those attributes.

	(Nominal) EmployeeId	(Nominal) Hard. of material	
1		V-H	
2		S	$d = \frac{2-0}{2} = 1$
3		V-S	$d \text{ lies b/w } [0,1]$
4		S	$d=1 \rightarrow \text{highly dissimilar}$
5		V-S	
similar	{ 1, 2, 3, 4 }		
	5		
	3		

$$d(6,3) = \frac{2-2}{2} = 0$$

Symmetric binary attribute (Gender)  
Asymm. (Cancer detec<sup>n</sup>)

SMC (Simple Matching Coefficient)  $\text{Sim}(i,j) = 1 - d(i,j)$  (0,0) or (1,1) matching!

	A	B	C	D	E
1	0	0	0	0	0
2	0	1	0	0	0
	x	x			
	1	0			

1 0 → r frequency of 10

1 2 3 4  
0 8 4 0

1 q r  
0 s t

Symmetric binary attribute

(Eg-students giving exam in similar way → 1,1 0,0)

Symmetric binary variables  $\Rightarrow d(i,j) = \frac{r+s}{q+r+s+t}$  (1,0 ones)

Asymmetric .. ..  $d(i,j) = \frac{r+s}{q+r+s}$  (0,0 removed → negative result → x help in dissimilarity)

Similarity in asym. " "  $\text{sim}_{\text{Jaccard}}(i,j) = \frac{q}{q+r+s}$  (1,1) out of those doing some right (1,1) (1,0) (0,1) (0,0) → wrong answer

Normalizing:

1) Z-score normalisation

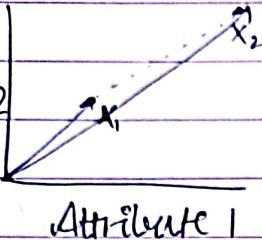
$$\begin{array}{|c|c|} \hline X & X' \\ \hline 1 & 0 \\ 2 & : \\ 3 & \\ \vdots & \\ 100 & 1 \\ \hline \end{array} \quad z = \frac{x - \mu}{\sigma}$$

$$\begin{array}{|c|c|} \hline X & X' \\ \hline 1 & -1 \\ 2 & -0.5 \\ 3 & 0 \\ 4 & 0.5 \\ 5 & 1 \\ \hline \end{array} \quad 1 \rightarrow \frac{1-3}{2} = -1$$

$$2 \rightarrow \frac{2-3}{2} = -0.5$$

2) Absolute Z-score normalisation (Mean abs. deviation)

## ② Numeric attributes :



2 dimensions → cartesian system

point	att. 1	att. 2	
$x_1$	1	2	point in vector space
$x_2$	3	5	
$x_3$	2	0	
$x_4$	4	5	

vectors  $4\mathbf{i} + 5\mathbf{j}$

For 2D → Euclidean distance  $\sqrt{(x_1 - x_2)^2 + (y_1 - y_2)^2}$

	$x_1$	$x_2$	$x_3$	$x_4$
$x_1$	0			
$x_2$	3.61	0		
$x_3$	5.1	5.1	0	
$x_4$	4.24	1	5.39	0

Minkowski distance: absolute diff. b/w

2 data obj. compared across p data objects

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

$h=1$  Manhattan distance

$h=2$  Euclidean ..

$h=\infty$  Supremum ..

$$\max_f |x_{if} - x_{jf}|$$

(max distance for a particular distance f)

$$d(x_1, x_2) = \sqrt[2]{(1-3)^2 + (2-5)^2}$$

③ max → supremum dist = 3

(3) Ordinal attributes  $\rightarrow$  Discrete  
 Qualitative, Categorical attribute  
 Order is imp.

$(0,1)$  map T shirt size

$\frac{1-1}{4-1} = 0$	$1 - S$
$2 - M$	
$\frac{3-1}{3-1} = 0.6 \leftarrow 3 - L$	
$\frac{4-1}{3-1} = 1 \leftarrow 4 - XL$	
$0.333 \leftarrow 2 - M$	
$0 \leftarrow 1 - S$	

can be converted to continuous (value taken within a range like 1-10)

\* To transform discrete to category continuous:  
 Every category has a rank as per their order  
rank  $m_f \rightarrow$  no. of category values against that attribute

$$\begin{aligned} S &\rightarrow 1 \\ M &\rightarrow 2 \\ L &\rightarrow 3 \\ XL &\rightarrow 4 \end{aligned}$$

$$z_{if} = \frac{\text{rank} - 1}{\text{no. of cat. vals} - 1}$$

Using the numeric measure we can calculate dissimilarity

#### (4) Mixed Attributes

Weighted Average  $d(i,j) = \frac{\sum_{f=1}^P \delta^{(f)} d_{ij}^{(f)} \text{weight}}{\sum_{f=1}^P \delta^{(f)}}$

0.2 0.5 0.2 0.1 weights  
 ordinal nominal Binary Numeric

	$X_1$	$X_2$	$X_3$	$X_4$
1	[			]
2				
3				
4				
5				
6				

weights assigned to diff. attributes as per <sup>their</sup> importance in the app

Dif. func<sup>n</sup>s/wt. comput<sup>n</sup> mechanism can be defined

numerical  
Normalize  $\rightarrow \frac{x_{if} - x_{if}^{\text{min}}}{x_{\text{max}} - x_{\text{min}}}$   
 $x_{if} = 21$  (or many value in final matrix)  
 $x_{if}^{\text{min}} = 0$  (or many value in matrix)

Nominal	$\rightarrow$	prepare dissimilarity matrix	$\begin{bmatrix} 0 \\ d(2,1) \\ d(3,1) & d(3,2) \end{bmatrix}$
Ordinal	$\rightarrow$	"	"
Binary	$\rightarrow$	"	$\begin{bmatrix} 0 \\ d(2,1) \\ d(3,1) & d(3,2) \end{bmatrix}$
Numeric	$\rightarrow$	"	"

$$d(2,1) = 0.2 \times d(2,1)_{\text{ordinal}} + 0.5 \times d(2,1)_{\text{nom}} + 0.2 \times d(2,1)_{\text{bin}} + 0.1 \times d(2,1)_{\text{num}}$$

$$0.2 + 0.5 + 0.2 + 0.1$$

Q.	SNO	0.1	0.5	0.4	final
		Test 1 (Nominal)	Test 2 (Ordinal)	Test 3 (Numeric)	
	1	code A	Fair	22	
	2	code B	Excellent	45	
	3	code A	good	64	
	4	code C	excellent	21	

Ans- ~~1234~~

Nominal

	1	2	3	4
1	0			
2	$\frac{45-0}{45} = 1$	0		
3	0	1	0	
4	1	1	1	0

$$\begin{bmatrix} 0 \\ d(2,1) \\ d(3,1) & d(3,2) \\ d(4,1) & d(4,2) & d(4,3) \end{bmatrix}$$

Ordinal Numeric

~~Ans~~

	1	2	3	4
1	0			
2	$45-22-23$	0		
3	42	19	0	
4	1	24	43	0

max  $\rightarrow 43$   
divide all by 43  $\rightarrow [0,1]$

Ordinal

Fair	1	$\frac{1-1}{3-1} = 0$
Good	2	$\frac{2-1}{2} = 0.5$
Excellent	3	$\frac{3-1}{2} = 1$

Ordinal,

		3	4
1	0		
2	$1-0=1$	0	
3	0.5	0.5	0
4	1	0	0.5

$$d(2,1) = 0.1 \times 1 + 0.4 \times 23 + 0.5 \times 1$$

$$= 10.8$$

$$d(3,1) = 0.1 \times 0 + 0.4 \times 42 + 0.5 \times 0.5$$

$$= 17.05$$

$$d(3,2) = 0.1 \times 1 + 0.4 \times 19 + 0.5 \times 0.5$$

$$= 7.95$$

$$d(4,1) = 0.1 \times 1 + 0.4 \times 1 + 0.5 \times 1$$

$$= 1$$

$$d(4,2) = 0.1 \times 1 + 0.4 \times 24 + 0.5 \times 0$$

$$= 9.7$$

$$d(4,3) = 0.1 \times 1 + 0.4 \times 43 + 0.5 \times 0.5$$

$$= 17.55$$

Write these 2 Qns in assignment

	1	2	3	4
1	0			
2	9.8	0		
3	17.05	7.95	0	
4	1	9.7	17.55	0

n.w Q.	S.NO	(Binary) Test 4	(Binary) Test 5	(Nominal) Test 6	Symmetric binary
1	1	1	0	S	
2	2	0	0	L	
3	3	0	1	XL	
4	4	1	1	L	

Ans- Test 4 (Binary) & 5

	1	2	3	4
1	0			
2	0.5	0		
3	1	0.5	0	
4	0.5	0.1	0.5	0

$$d(2,1) = \frac{1}{\sqrt{2}} \begin{bmatrix} 0 & 0 \\ 1 & 1 \end{bmatrix} = \frac{1}{2}$$

$$d(3,1) = \frac{1}{\sqrt{2}} \begin{bmatrix} 1 & 0 \\ 0 & 0 \end{bmatrix} = \frac{\sqrt{2}}{2} = 1$$

$$d(3,2) = \frac{1}{\sqrt{2}} \begin{bmatrix} 1 & 1 \\ 0 & 0 \end{bmatrix} = \frac{1}{2}$$

Test 6

0			
1	0		
1	1	0	
1	1	1	0

## Cosine Similarity

Terms in all the documents  $\rightarrow$  attributes

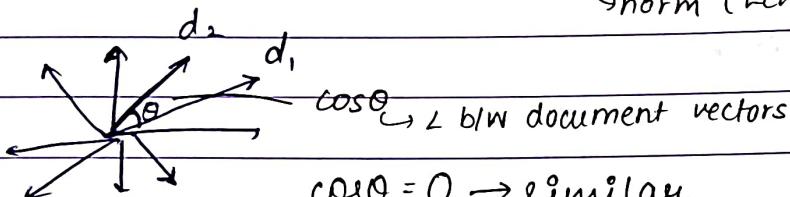
Document	team	coach	hockey
document 1	5	0	3
document 2	3	0	2

Compare text documents by finding the cosine similarity

$$\begin{array}{l} x = (2, 0, 3) \\ y = (1, 2, 3) \end{array} \quad \begin{array}{c} x \\ + \\ y \\ \hline 3 \\ 2 \\ 5 \\ 6 \end{array}$$

dot product

Similarity  $\rightarrow \cos(d_1, d_2) = \frac{\vec{d}_1 \cdot \vec{d}_2}{\|\vec{d}_1\| \|\vec{d}_2\|}$   $d_1, d_2 \rightarrow$  vectors  
 $\|\vec{d}\|$  norm (length of vector)



$$\cos \theta = 0 \rightarrow \text{similar}$$

$$\cos \theta = 90^\circ \rightarrow 90^\circ / \text{orthogonal}$$

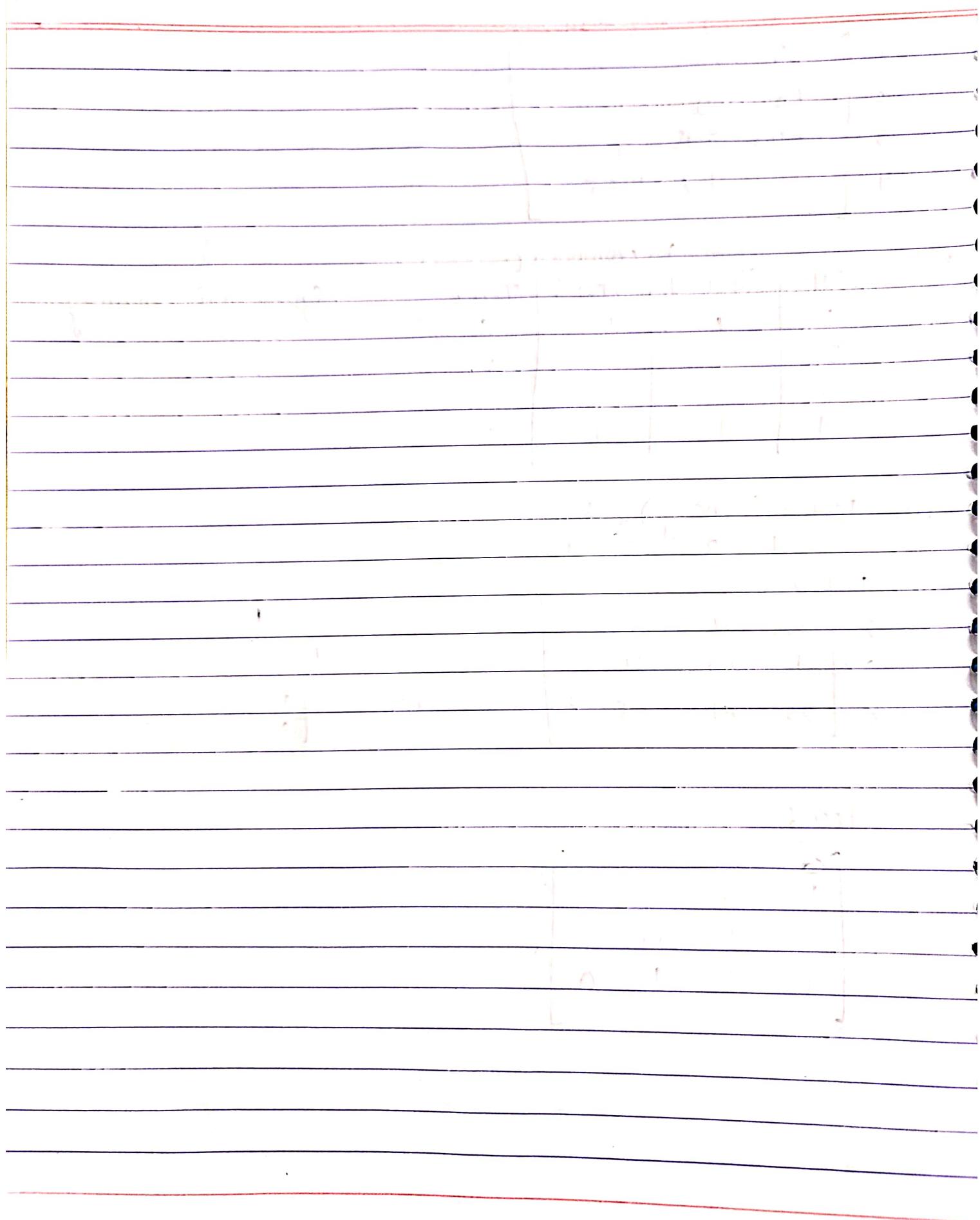
$90^\circ \rightarrow \text{dissimilar}$

$$\cos \theta = 0.42$$

$$\cos \theta = 0.30 \quad (\text{more similar}) ?$$

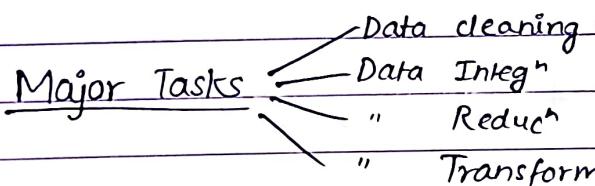
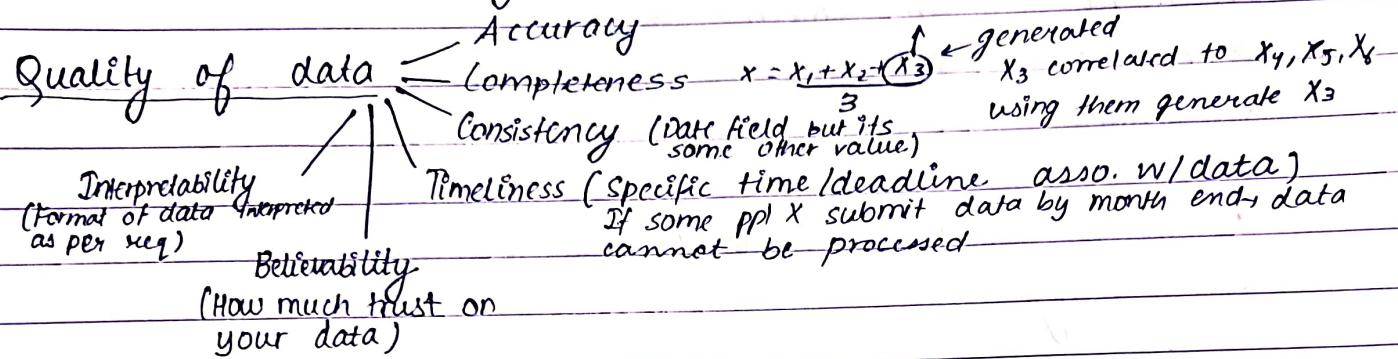
As  $\theta \uparrow$ , dissimilarity  $\uparrow$

$$d(i, j) = 1 - \cos(i, j)$$



## Chapter 3: Data Preprocessing

Objective: Prepare your data for requirement for data mining algorithm



### ① Data cleaning

a)

$X_1$	$X_2$	$X_3$	$X_4$	Handling missing values
✓	✓	✓	✓	Ignore missing values (good if dataset v. large & only few missing values)
✓	✓	✓	-	Eg - 1L total, 4 missing (X advised if moderate size)
✓	-	✓	-	
✓	-	✓	-	

4 dimensional space  
 (Multiple regression → linear for only 2)

1) Ignore missing values (good

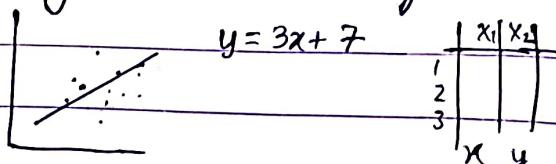
if dataset v. large & only few missing values)

Eg - 1L total, 4 missing

(X advised if moderate size)

2) Statistical method (mean...)

3) Regression analysis



4) Predict" mechanism (decision tree)  
 Eg - 1000 total 3 missing using 997 → train model  
 $3 \rightarrow$  Predict



4, 8, 15, 21, 21, 24, 25, 28, 34

① Equal Partitioned Mean

Bin1: 4, 8, 15

Bin2: 21, 21, 24

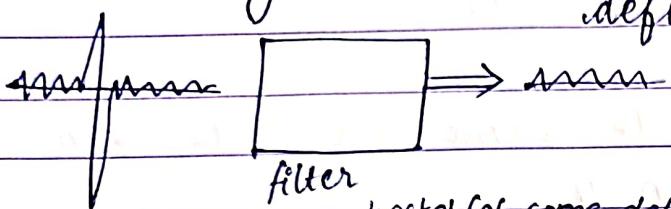
Bin3: 25, 28, 34

② Bin Mean (replace all by mean)  
Bin1: 9, 9, 9  
Bin2: 22, 22, 22

### b) smoothing noisy data

some random values/variance in your data  
How to handle?

Smoothing method: smooth your data within the defined limit



→ basket (of some defined size)

Eg - Binning method

4, 8, 15, 20, 21, 27, 32, 35, 91, 43, 47, 102

outlier/noise

→ sort

→  $w=4$  (basket width)

Bin 1 = 4, 8, 15, 20

Binning by mean → calc.

Bin 2 = 21, 27, 32, 35

mean of each bin &

Bin 3 = 91, 43, 47, 102

replace all values in the bin by that

$$\text{Mean(Bin 1)} = 12$$

$$\text{Bin 1} = 12, 12, 12, 12$$

$$\text{Bin 2} = 28, 28, 28, 28$$

Another is binning by max value (Take max in a bin)

$$\text{Bin 3} = 59, 59, 59, 59$$

(24, 23, 27, 32, 34, 36)

$$2+6=8$$

?, ?

$$64+2=6$$

Binning can be used as discretization process & data reduction (No. of bits req. is less)

Reason - Data integ<sup>r</sup> (causes redundancy as well)

### d) Resolve inconsistencies

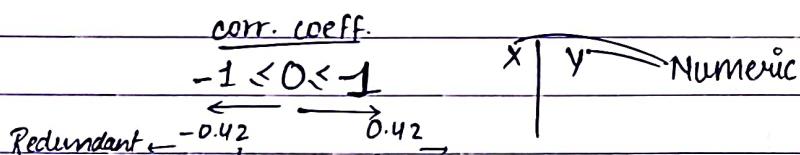
In a field of certain dtype, there is value of another dtype  
How to handle?

1) Numeric attribute  $\rightarrow$  correlation coefficient  $[-1, 1]$  (Pearson)

2) Nominal attribute  $\rightarrow$  Chi-squared test

A.cust-id = B.cust-# both same attributes  $\rightarrow$  on integ<sup>r</sup>  
there is redundancy

Chi-square test  $\chi^2 = \sum \frac{(\text{observed} - \text{expected})^2}{\text{expected}}$  Used to find redund. or inconsis. in data.



If no correlation ( $\rho=0 \rightarrow$  then  $X$  &  $Y$  are independent)

$\rho=1 \rightarrow$  similar  $\rightarrow$  redundant

LST  $\rightarrow$  assumes attributes are independent hypothesis  $\rightarrow$  if rejected  $\rightarrow$

$X$  (has 2 nominal values)

		Play chess	Not play chess
Y	Like sci fiction	250	200
	Don't like sci fiction	50	1000

	X	Y
PC	PC	LSF
	NPC	LSF
	PC	NLSF

Observed freq. of diff. nominal values

$\rightarrow$  Recognize freq. of PC & NPC  
 $\rightarrow$  Make contingency table

$$\text{Expected frequency } E_{ij} = \frac{A(A=a_i) \times B(B=b_j)}{n}$$

check

$$e_{11} = A(A=PC) \times B(B=PS)$$

$$= \frac{300}{n} \times \frac{450}{150} = 90$$

Prob. of occurrence of each  $\rightarrow$  multiply both  $\rightarrow$  consider independent

~~(m-1) x (n-1)~~ no. of indep. var. wrt ~~not mut~~

Significance level 0.001  $\chi^2$  values are but this much.

507.93 (highly correlated  $\rightarrow \chi^2 > 15$  upper limit)

Deg. of freedom	significance	$\chi^2$
15	0.001	15

Within 15  $\rightarrow$  independent

> 15  $\rightarrow$  highly correlated

	X	Y
1	2	9
2	3	28
3	4	65
4	5	?
5	6	217

Using linear Reg. method, find missing value?

$$\cos(x, y) = ? \quad 0.926$$

$$\text{corr}(x, y) = ? \quad 0.964$$

Write your observations about cos & corr & what they indicate

Ans-

$$28 - 9 = y - 65$$

$$3 - 2 \quad 5 - 4$$

$$19 = y - 65$$

$$y = 84$$

$$y = mx + c$$

$$y = 19x + c$$

$$65 = 19 \times 4 + c$$

$$c = -11$$

$$y = 19x - 11$$

$$y =$$

$$\cos(x, y) = ?$$

$$= \frac{2x9 + 3x28 + 4x65 + 5x126 + 6x217}{\sqrt{2^2 + 3^2 + 4^2 + 5^2 + 6^2} \sqrt{9^2 + 28^2 + 65^2 + 126^2 + 217^2}}$$

$$= \frac{2294}{9.48 \times 260.87} = 0.927$$

$$\text{corr}(x, y) = \frac{\text{Covariance}(x, y)}{\text{std}(x) \cdot \text{std}(y)}$$

$$\text{Covariance}(x, y) = \frac{\sum_{k=1}^n (x_k - \bar{x})(y_k - \bar{y})}{n-1}$$

$$\text{std}(x) = \sqrt{\frac{\sum_{k=1}^n (x_k - \bar{x})^2}{n-1}}$$

$$\text{corr}(x, y) = ? \quad \bar{x} = 4 \quad \bar{y} = 89$$

$$\text{Covariance}(x, y) = \frac{(-4+2)(9-89) + (3-4)(28-89) + (4-4)(65-89) +}{4}$$

$$\frac{(5-4)(126-89) + (6-4)(217-89)}{4}$$

$$= \frac{(-2)(-80) + (-1)(-61) + 37 + 2 \times 128}{4} = \frac{160 + 61 + 37 + 256}{4} = 128.5$$

$$\text{Std}(x) = \sqrt{\frac{(2-4)^2 + (3-4)^2 + (5-4)^2 + (6-4)^2}{4}}$$

$$= \sqrt{\frac{10}{4}} = 1.58$$

$$\text{Std}(y) = \sqrt{\frac{80^2 + 61^2 + 37^2 + 128^2 + 24^2}{4}}$$

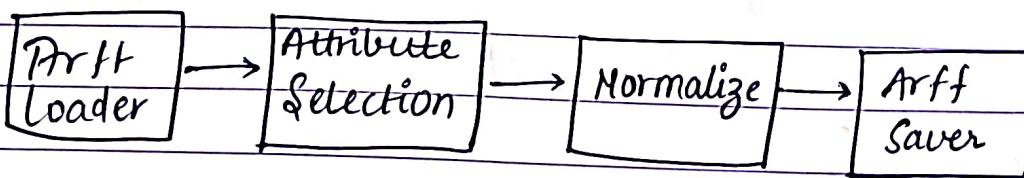
$$= \sqrt{\frac{6400 + 3721 + 1369 + 16384 + 576}{4}}$$

$$= 84.3$$

$$\text{corr}(x,y) = \frac{128.5}{1.58 \times 84.3} = 0.964$$

## Lab

Data Sources > Arff Loader (Drag & drop) > Configure > add file name  
 Filters > supervised > attribute > attribute Selection



Arff loader icon > Dataset > Join w/ Attribute selec<sup>n</sup>  
 unsupervised > normalize > Drag & drop > join

Data Sinks > Arff Saver > Drag, drop, join > Right click icon > configure >  
 give directory to save file  
 ► click > Run this Now (below pgm tab) > see log to see save  
 successful

Q. Is  $\vec{v}$  highly non correlated?

## \* Data Reduction

Attribute or sample level reduction

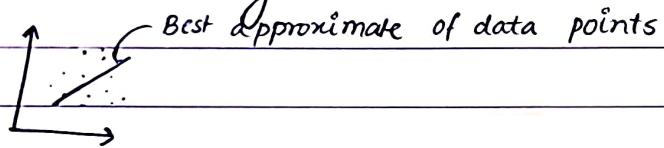
1) Dimensionality reduction (Attribute level reduc<sup>n</sup>)

→ Wavelet transforms

→ PCA

→ Feature subset selection, feature creation  
(10 features combined to create 1)

2) Numerosity reduction (Model level sol<sup>n</sup>)



3) Data Compression

→ Binary Compression

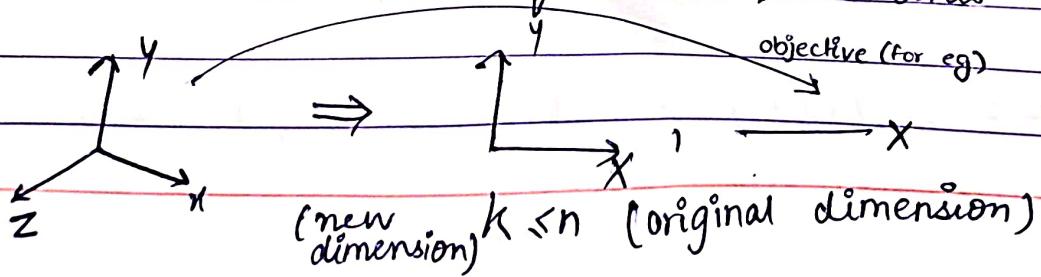
## \* PCA

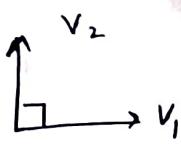
With the given data, we have to find principal components  
(10 attributes  $\rightarrow$  10 PC we can find)

	$X_1$	$X_2$	
1	-	1	$\Rightarrow$
-	-		
100	-	10	

$\leq 2$  dimension

2 variables can be transformed like this in terms of PCs





orthogonal  $\rightarrow$  b/w vectors  $90^\circ$   
 orthonormal  $\rightarrow$  if length of these vectors is unit  
 length &  $L = 90^\circ$

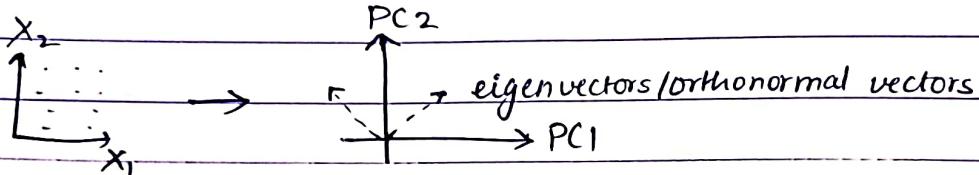
Ho

### \* Variance method

$$\text{Variance}(x) + \text{Var}(y) + \text{var}(z) = \text{var}(x)$$

total variance is same in original & transformed data  $\rightarrow$  no data loss

- ① Normalize your data (to avoid bias) (let 1-10 range)
- ② Compute k orthonormal vectors



③

$$X_i = a \cdot PC_1 + b \cdot PC_2$$

- ④ Sort PCs in  $\downarrow$  order of strength

$[PC_1, PC_2]$  if you want to choose 1  $\rightarrow$  select PC1

$[PC_1, PC_2, PC_3]$  data reduc<sup>n</sup> to 2  $\rightarrow$  select first two

Eg -  $X \quad y$

132 80

130 82

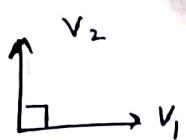
136 78

140 80

134 84

130 82

$$A: \bar{X} = 134 \quad \bar{y} = 81$$



Orthogonal  $\rightarrow$  b/w vectors  $90^\circ$   
 orthonormal  $\rightarrow$  if length of these vectors is unit  
 length &  $L=90^\circ$

Ho

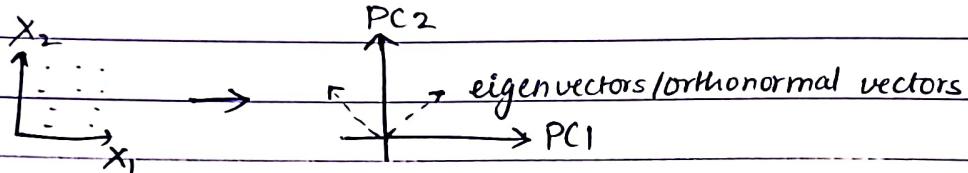
### \* Variance method

$$\text{Variance}(x) + \text{Var}(y) + \text{var}(z) = \text{var}(x)$$

Total variance is same in original & transformed data  $\rightarrow$  no data loss

① Normalize your data (to avoid bias) (Let 1-10 range)

② Compute k orthonormal vectors



③

$$X_i = a \cdot PC_1 + b \cdot PC_2$$

④ Sort PCs in  $\downarrow$  order of strength

PC1, PC2 if you want to choose 1  $\rightarrow$  select PC1

PC1, PC2, PC3 data reduc<sup>n</sup> to 2  $\rightarrow$  select first two

Eg -	X	y
132	80	
130	82	
136	78	
140	80	
134	84	
132	82	

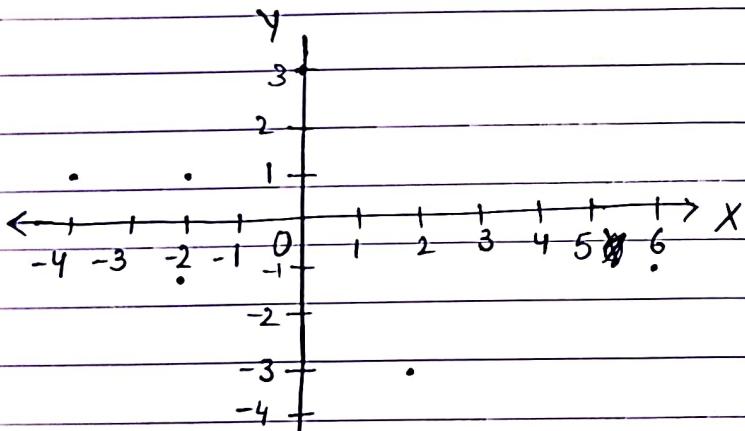
$$A: \bar{X} = 134 \quad \bar{y} = 81$$

Variance = How much variation in your data wrt mean  
 Covariance = How 1 variable vary wrt another (if you vary x, what will be var in y)

### ① Normalize by mean

$X \text{ of } X - \bar{X}$	$Y - \bar{Y} \text{ by } Y$
-2	-1
-4	1
2	-3
6	-1
0	3
-2	1

Centre data / normalised data



### ② Covariance matrix (∴ variance method)

$$\text{cov}(x, x)$$

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

	X	Y
X	10.6	-2.6
Y	-2.6	3.67

$$\begin{aligned} \text{cov}(x, x) &= \frac{(-2)x(-2) + (-4)^2 + (2)^2 + 6^2 + (-2)^2}{6} \\ &= \frac{4 + 16 + 4 + 36 + 4}{6} = 10.66 \end{aligned}$$

$$\begin{aligned} \text{cov}(x, y) &= \frac{(-2)x(-1) + (-4) + (2)(-3) + 6(-1) + (-2)(1)}{6} \\ &= \frac{2 - 4 - 6 - 6 - 2}{6} = -\frac{16}{6} = -2.66 \end{aligned}$$

$$\begin{aligned}\text{cov}(Y, Y) &= \frac{(-1)^2 + 1 + (-3)^2 + (-1)^2 + 3^2 + 1}{6} \\ &= \frac{1+1+9+1+9+1}{6} \\ &= \frac{22}{6} = 3.66\end{aligned}$$

$$A = \begin{bmatrix} 10.6 & -2.6 \\ -2.6 & 3.67 \end{bmatrix}$$

orthogonal + unit vector orthonormal

③ Find eigen values & eigenvectors  $= a \begin{bmatrix} 1 \\ 0 \end{bmatrix} + b \begin{bmatrix} 0 \\ 1 \end{bmatrix}$

$$\det(A - \lambda I) = 0$$

large

$$(0,1) \xrightarrow{\quad} \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix}$$

$$\det = 1$$

A matrix can be decomposed

into submatrices, which are det of matrix  $\rightarrow$  gives volume formulated by that object formed by the matrix data

$$\underbrace{\begin{bmatrix} 1 & 1 & 1 & 1 & 1 \end{bmatrix}}_{\text{Eigen vectors}} = \underbrace{\begin{bmatrix} 1 \\ 1 \\ 1 \\ 1 \\ 1 \end{bmatrix}}_{\text{Eigen vectors}} + \underbrace{\begin{bmatrix} 1 \\ 1 \\ 1 \\ 1 \\ 1 \end{bmatrix}}_{\text{Eigen vectors}} + \underbrace{\begin{bmatrix} 1 \\ 1 \\ 1 \\ 1 \\ 1 \end{bmatrix}}_{\text{Eigen vectors}}$$

Cond' is that these eigen vectors should be orthonormal to each other / independent to each other

$$AI = \lambda I$$

$$AI = 0$$

$$A \cdot V = 0$$

$$A - \lambda I = \begin{bmatrix} 10.6 - \lambda & -2.6 \\ -2.6 & 3.67 - \lambda \end{bmatrix} = (10.6 - \lambda)(3.67 - \lambda) - (-2.6)^2$$

$$= 38.9 - 14.27\lambda + \lambda^2 - 6.76 = 0$$

$$\lambda^2 - 14.27\lambda + 32.14 = 0$$

$$\lambda = \frac{14.27 \pm \sqrt{14.27^2 - 4 \times 32.14}}{2}$$

$$= \frac{14.27 \pm \sqrt{75.01}}{2}$$

$$= \frac{22.93}{2}, \frac{5.61}{2}$$

$$= \underbrace{11.468}_{\lambda_1}, \underbrace{2.8}_{\lambda_2}$$

Eigenvector

$$| A \cdot v = \lambda v |$$

$$| (A - \lambda I) \cdot v = 0 | \quad (\text{orthogonal})$$

Let  $\lambda = \lambda_1 (11.46)$

$$\begin{bmatrix} 10.6 - \lambda & -2.6 \\ -2.6 & 3.67 - \lambda \end{bmatrix} \begin{bmatrix} x \\ y \end{bmatrix} = 11.46 \begin{bmatrix} x \\ y \end{bmatrix}$$

$$\begin{bmatrix} 9.6 & -2.6 \\ -2.6 & 2.67 \end{bmatrix} \begin{bmatrix} x \\ y \end{bmatrix} = 11.46 \begin{bmatrix} x \\ y \end{bmatrix}$$

should be 0?

$$-0.865x = 2.6y$$

$$-0.332x = y$$

$$-2.6x = 7.795y$$

$$\begin{bmatrix} 9.6y & -2.6x \\ -2.6x + 2.67y \end{bmatrix} \begin{bmatrix} 11.46x \\ 11.46y \end{bmatrix} = \begin{bmatrix} 1 \\ -0.33 \end{bmatrix}$$

$$9.6y - 14.06x = 0 \Rightarrow y = 1.46x$$

$$-2.6x - 8.79y = 0$$

$$\begin{bmatrix} -0.865x & -2.6y \\ -2.6x & -7.795y \end{bmatrix} \begin{bmatrix} x \\ y \end{bmatrix} = 11.46 \begin{bmatrix} x \\ y \end{bmatrix} \rightarrow \text{should be } 0?$$

$$\begin{bmatrix} -0.865x - 2.6y \\ -2.6x - 7.795y \end{bmatrix} = \begin{bmatrix} 11.46x \\ 11.46y \end{bmatrix}$$

$$-0.865x - 2.6y = 11.46x$$

$$y = -0.135x$$

$$y = -4.73x$$

$$x = t \text{ or } x = 1$$

$$y = -4.73$$

$$\begin{bmatrix} t \\ -4.73t \end{bmatrix} = t \begin{bmatrix} 1 \\ -4.73 \end{bmatrix}$$

$$1st \text{ eigenvector} = \begin{bmatrix} 1 \\ -4.73 \end{bmatrix}$$

$$\begin{bmatrix} 10.6 - 2.8 & -2.6 \\ -2.6 & 3.67 - 2.8 \end{bmatrix} \begin{bmatrix} x \\ y \end{bmatrix} = 2.8 \begin{bmatrix} x \\ y \end{bmatrix}$$

$$\begin{bmatrix} 7.8 & -2.6 \\ -2.6 & 0.87 \end{bmatrix} \begin{bmatrix} x \\ y \end{bmatrix} = 2.8 \begin{bmatrix} x \\ y \end{bmatrix}$$

$$\begin{bmatrix} 7.8x - 2.6y \\ -2.6x + 0.87y \end{bmatrix} = \begin{bmatrix} 2.8x \\ 2.8y \end{bmatrix}$$

$$7.8x = 2.6y$$

$$y = 3x$$

$$0.87y = 2.6x$$

$$y = 3x$$

$$5x = 2.6y$$

$$y = 1.92x$$

$$2nd \text{ eigenvector} = \begin{bmatrix} 1 \\ 1.92 \end{bmatrix} \quad \begin{pmatrix} 1 \\ 3 \end{pmatrix}$$

$V \rightarrow W$  (2D space)  
 (3D space) (4, 3)  
 $(2, 1, 4)$   
 To transform you need a  
 "transform" matrix

Make eigenvector unit length

$$\frac{1}{\sqrt{1^2 + (-4.73)^2}}$$

$$= \frac{1}{\sqrt{23.37}} = \frac{1}{4.83} = 0.2$$

$$1^{\text{st}} \text{ EV} = \begin{bmatrix} 0.2 \\ -0.97 \end{bmatrix}$$

$$2^{\text{nd}} \text{ EV} = \begin{bmatrix} 0.46 \\ 0.88 \end{bmatrix}$$

$$\begin{bmatrix} 0.46 & 0.2 \\ 0.88 & -0.97 \end{bmatrix}$$

$$M_{2 \times 2} \begin{bmatrix} 0.46 & 0.2 \\ 0.88 & -0.97 \end{bmatrix} = \begin{bmatrix} \cdot & \cdot \end{bmatrix}_{6 \times 2}$$

ask if orig matrix  $2 \times 2$

$$\begin{bmatrix} -2 & -1 \\ -4 & 1 \\ 2 & -3 \\ 6 & -1 \\ 0 & 3 \\ -2 & 1 \end{bmatrix} \cdot \begin{bmatrix} 0.46 & 0.2 \\ 0.88 & -0.97 \end{bmatrix} = \begin{bmatrix} \cdot & \cdot \end{bmatrix}_{6 \times 2}$$

Eigenvector

$$\begin{array}{cc} \text{PC1} & \text{PC2} \\ -1.8 & -1.37 \\ -0.96 & -1.77 \\ -1.72 & 3.31 \\ 1.88 & 2.17 \\ 2.64 & -2.91 \\ 0.04 & -1.37 \end{array}$$

$\text{Var}_1 = \text{Var}_2$

represent 98% var in data so reduce to 2% dimensions to only these 3

$$1, 2, 3, \dots, 100 \rightarrow P_{46}, P_{32}, P_{67}, \dots, P_1$$

$\text{PC}_1, \text{PC}_{100}$  Most variance shifted to some PC values

$$\text{Total var}^n = \frac{10.6}{(x,x)} + \frac{3.67}{(y,y)} \approx 14$$

Now var<sup>n</sup> shifted to some PCs

PCA → dimension reduc<sup>n</sup> technique

Var<sup>n</sup> moved in only a few directions (less dimension spaces) → represent overall data that reduced dimension maximally represent your data

~~Ch-3~~ Ch-3

### \* Sampling

(Technique to ↓ size of data)

selecting some samples from the given data

Types

① Random sampling:

→ select 10k randomly

Not good choice, maybe 10k belong to the same class

② Sampling w/o replacement (better training than ↓)

③ " w/ " (redundancy possible)

④ Stratified Sampling

cds per data dist<sup>n</sup>, select samples

250 ← {----- Partition dataset & select from each

250 ← {----- Skewed data (biased towards majority class)

↳ sampling → balance

80k → 20k → 60k + removed

→ Undersampling (Majority class instances are eliminated)

→ Oversampling (Regenerate minority instances w/o elim. maj) Generate 60k pseudo instance

→ Hybrid Sampling

Ch-8

## 1 Data Mining / Learning Classification algorithms

- ① Supervised (2 step-process → training & testing)  
on basis of similarity they get labelled
- ② Unsupervised (No training → categor. given data for training purpose)
- ③ Semi supervised (

②

	item1	item2	item3	
c <sub>1</sub>	✓			
c <sub>2</sub>		✓		

item1 item2 item3 compare similarity of each value

- ① Feedback mech

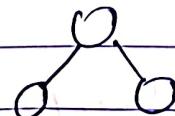
Learns from errors

Types:

If value of class attribute the discrete → classification  
continuous → regression

X <sub>1</sub>	X <sub>2</sub>	X <sub>3</sub>	class
✓	✓	-	+
✓	-	✓	-
✓	✓	✓	?
✓	✓	✓	?

→ supervised



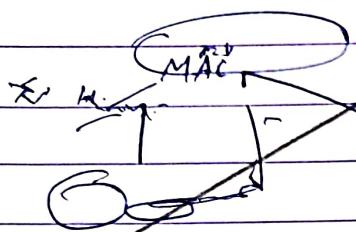
sweeper wise ex.

~~Decision~~

$D = \{ \text{designed} \}$  won't all get  $P$ .

$H_1$   $H_2$

$$D = \{ 1 \}$$



$$D_1 = \{ 1, 2, 3 \}$$

$$D_2 = \{ \}$$

$$S = 37 \text{ m}$$

$$D = \{ \}$$

$$32.40$$

$$D_3 = \{ 5, 3, 16, 10, 4 \}$$

$$D_1 = \{ \}$$

$$D = \{ \}$$
  
AGE

$$\leq 30$$

$$31-40$$

$$> 40$$

$$D_1 = \{ 3, 1, 2, 8, 9, 11 \}$$
  
student

NO

$$D_4 = \{ 1, 2, 8, 3 \}$$

YES

$$D_2 = \{ 3, 7, 12, 13 \}$$
  
all samples  $\rightarrow$  same class

$$D_3 = \{ 4, 5, 6, 10, 14 \}$$

Credit Rating

Excellent

Fair

$$D = \{ 1 \}$$

$$D = \{ 2 \}$$

$$D = \{ 3 \}$$

$$D_5 = \{ 9, 11 \}$$
  
YES

$$\{ 11 \}$$

$$\{ 9 \}$$

$$\{ 13 \}$$

$$\{ 14 \}$$

$$\{ 15 \}$$

$$\{ 16 \}$$

$$\{ 17 \}$$

$$\{ 18 \}$$

$$\{ 19 \}$$

$$\{ 20 \}$$

$$\{ 21 \}$$

$$\{ 22 \}$$

$$\{ 23 \}$$

$$\{ 24 \}$$

$$D_5 = \{ 16, 14, 3 \}$$

$$NO$$

$$D_7 = \{ 4, 5, 10 \}$$

$$YES$$

Performance good if all attributes are discrete attribute (Ranking) ~~continuous~~

Selected attribute  $\rightarrow$  selected randomly

If

(12, 8)

other classes → either split further

After the other nodes → or use voting

see Dec. Tree

Discussions

Madhusmita

## \* Attribute Selection Measures

### Information Gain

Gain Ratio

Gini ~~Ratio~~ Index

$$H(Y) = - \sum_{i=1}^m p_i \log(p_i) \quad p_i = P(Y=y_i)$$

Eg - Hello world (amt of info in these 2 words → given by above formula)

$$H=1$$

$$E=1$$

$$L=3$$

$$O=2$$

$$W=1$$

$$R=1$$

$$D=\frac{1}{10}$$

$$- \left[ \frac{1}{10} \log_2 \frac{1}{10} + \frac{1}{10} \log_2 \frac{1}{10} + \frac{2}{10} \log_2 \frac{3}{10} + \right.$$

$$\left. \frac{2}{10} \log_2 \frac{2}{10} + \frac{1}{10} \log_2 \frac{1}{10} + \frac{1}{10} \log_2 \frac{1}{10} + \frac{1}{10} \log_2 \frac{1}{10} \right]$$

$$= 0.30$$

$k$  = No. of diff. class labels

$m_i$  = No. of values in the  $i^{th}$  interval of a partition

$m_{ij}$  = No. " " of class  $j$  in interval  $i$

Entropy of  $p_i$  of  $i^{th}$  interval

$$E_i = - \sum_{j=1}^k p_{ij} \log_2 p_{ij}$$

$$p_{ij} = \frac{m_{ij}}{m_i}$$

Randomness/Impurity in data  $\rightarrow$  Entropy

\* Pure partition

$$D = \{10\} \xrightarrow{7+3-} \text{(10 instances)}$$

$$D_1 = \{3-4\}$$

$$D_2 = \{7+4\}$$

Objective of decision tree:

To classify sample in pure partition

Impurity involved

$$D = \{10\}$$

$$D_1 = \left\{ \begin{array}{l} 2- \\ 1+ \end{array} \right\}$$

$$D_2 = \left\{ \begin{array}{l} (-) \\ 6+ \end{array} \right\}$$

$\xleftarrow{\text{Impurity}}$  Amt of impurity/randomness

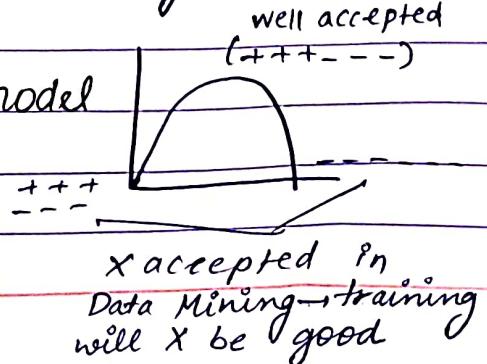
Got this partit<sup>n</sup> due to selec<sup>n</sup> of some attribute

$$\begin{array}{c}
 \text{12} \\
 \boxed{\phantom{000}} \\
 9 \\
 \boxed{++++} \quad \boxed{\phantom{00}} \quad \boxed{\phantom{00}}
 \end{array}
 \quad
 \begin{array}{c}
 \text{+ve class} \quad \text{-ve class} \\
 - \left[ \frac{4}{6} \log_2 \frac{4}{6} + \frac{2}{6} \log_2 \frac{2}{6} \right] \\
 = [0.66 \times -0.59 + 0.33 \times -1.59] \\
 = 0.66 \times 0.18 \\
 = -[-0.38 - 0.52] \\
 = 0.9
 \end{array}$$

Imp +++-- equal dist<sup>n</sup> 50% probability  $\rightarrow$  value of entropy is maximum (1)

++++++ 100% prob.  $\rightarrow$  0 entropy/impurity

$\rightarrow$  Pure part<sup>n</sup> X good for training model  
 $\rightarrow$  Biased towards that class



## Attribute Selection Measure : Information Gain (ID3/C4.5)

Expected Inform<sup>n</sup>(entropy)

$$\text{Info}(D) = - \sum_{i=1}^m p_i \log_2(p_i)$$

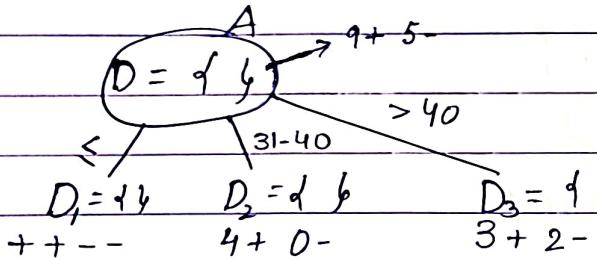
Total info in original data (Amt of impurity in data before partition)  
 $m \rightarrow$  no. of class labels

Inform<sup>n</sup> needed

$$\text{Info}_A(D) = \sum_{j=1}^{|D|} |D_j| \times \text{Info}(D_j) \quad \text{remaining in splitted part}$$

inform<sup>n</sup> gained

$$\text{Gain}(A) = \text{Info}(D) - \text{Info}_A(D)$$



$$\text{Info}(D) = I(9, 5) = -\frac{9}{14} \log_2 \frac{9}{14} - \frac{5}{14} \log_2 \frac{5}{14} = 0.94$$

(Amt of impurity in given dataset)  
 Find impurity in respective part

Find

$$D_1 = -2 \log_2 \frac{2}{5} - 3 \log_2 \frac{3}{5}$$

normalise to avoid biasing effect on diff. dist<sup>n</sup>

$$\frac{5}{14} [D_1] + \frac{4}{14} \left[ -\frac{4}{4} \log_2 \frac{4}{4} - \frac{0}{4} \log_2 \frac{0}{4} \right] +$$

$$\frac{5}{14} \left[ -\frac{3}{5} \log_2 \frac{3}{5} - \frac{2}{5} \log_2 \frac{2}{5} \right] = 0.694$$

$$\text{Gain}(\text{age}) = \text{Info}(D) - \text{Info}_A(D) = 0.246 \rightarrow \text{max}$$

$$\text{Gain}(\text{income}) = 0.029$$

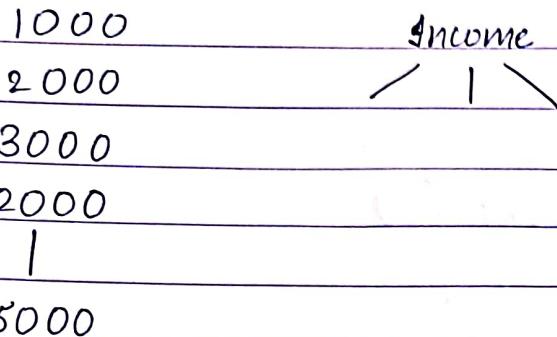
$$\text{Gain}(\text{student}) = 0.151$$

$$\text{Gain}(\text{credit\_rating}) = 0.048$$

Decision Tree → good for discrete attributes

If continuous attribute in data,

Income



75 70 85 60 120 100 95 90 125 220

No No NO Yes Yes Yes NO NO No No  
60 70 75 85 90 95 100 120 125 220

① Sort in ↑ order

	class
60	NO
70	NO

② Append the class label corr. to the row  
(above respective value)

③ Consider any value < 60 & > 220

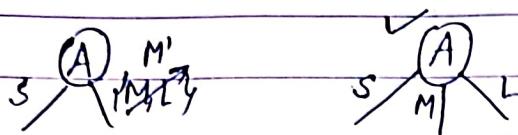
	NO	NO	NO	Y	Y	Y	N	N	N	N	
<	>	<	>	<	>	<	>	<	>	<	>
Yes	3	-0	3	0'3	0'3	112	211	310	310	1	1
No	7	1	2	5	3	4	3	4	4	3	1
imp	1	1	1	1	1	1	1	1	1	1	1

classes  
see label  
is NO so  
shift it

You get dist b/w diff intermed. pt.

Calculate gain → min value → select that as split point

"Inform" gain measure is biased towards attributes w/ a large no. of values (2 ~~discrete values~~ < 10 discrete values)



Normalise gain w/ cost

(pay cost for no. of splits)

$$\text{GainRatio}(A) = \text{Gain}(A) / \text{SplitInfo}(A)$$

$$\text{SplitInfo}(D) = -\sum_{j=1}^{(\text{no. of parts})} \frac{|D_j|}{|D|} \times \log_2 \frac{|D_j|}{|D|}$$

$$\begin{aligned} \text{Age} \rightarrow & \leftarrow 30 \rightarrow -\frac{5}{14} \log_2 \frac{5}{14} \\ & -\frac{4}{14} \log_2 \frac{4}{14} \quad \text{sum of all} \\ & -\frac{5}{14} \log_2 \frac{5}{14} \end{aligned}$$

$$Gini(D) = 1 - \sum_{j=1}^n p_j^2$$

$$gini_A(D) = \frac{|D_1|}{|D|} gini(D_1) + \frac{|D_2|}{|D|} gini(D_2)$$

+++-

$$1 - \left[ \left( \frac{3}{5} \right)^2 + \left( \frac{2}{5} \right)^2 \right]$$

$$\Delta gini(A) = gini(D) - gini_A(D)$$

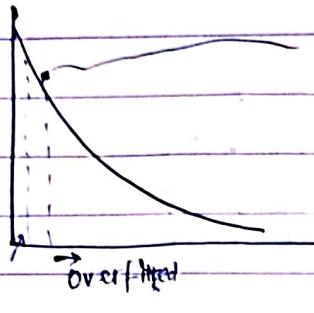
Consider smallest gini value for attribute selection

0 nodes  $\rightarrow$  error  $\uparrow$

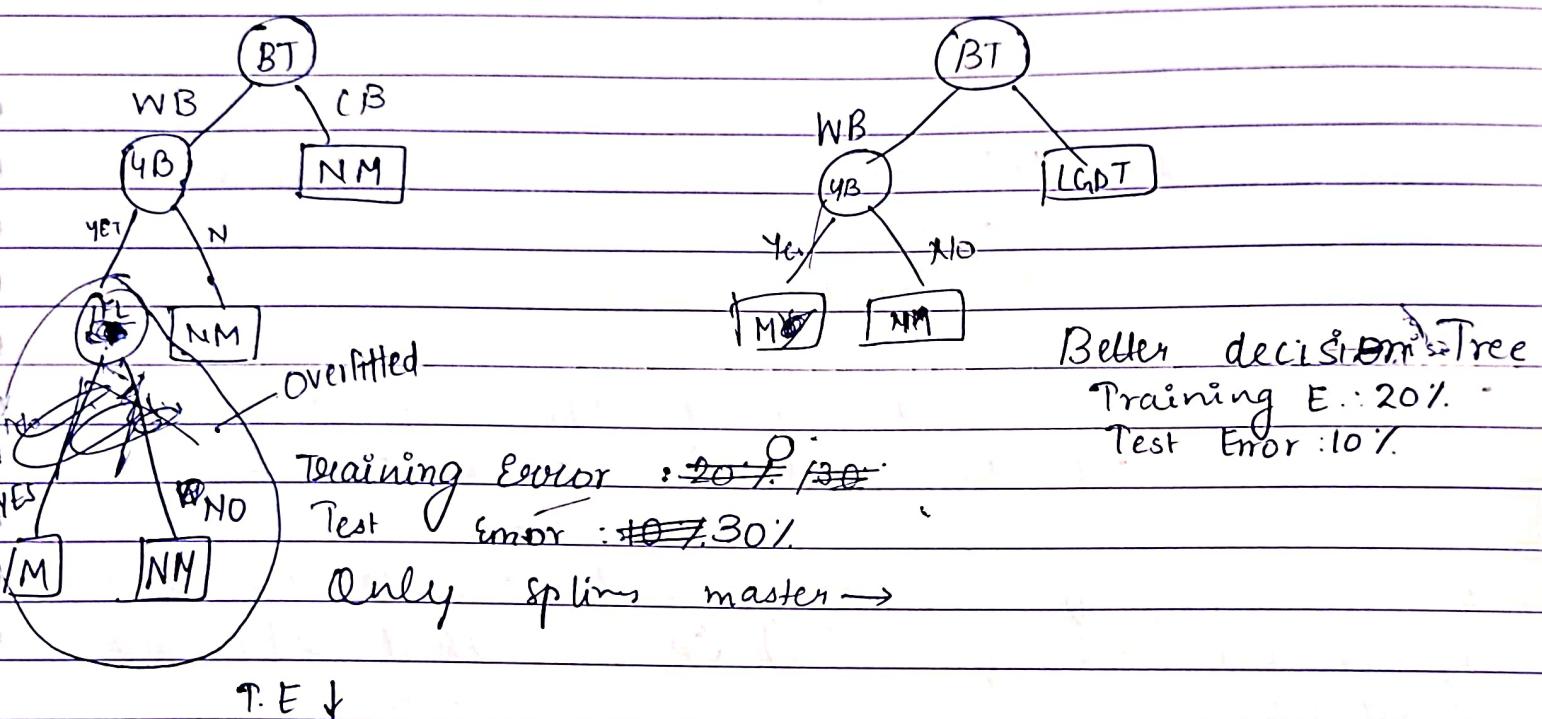
as you  $\uparrow$  no. of nodes  $\rightarrow$  Both errors  $\downarrow$

After a pt if we  $\uparrow$  no. of nodes, testing error

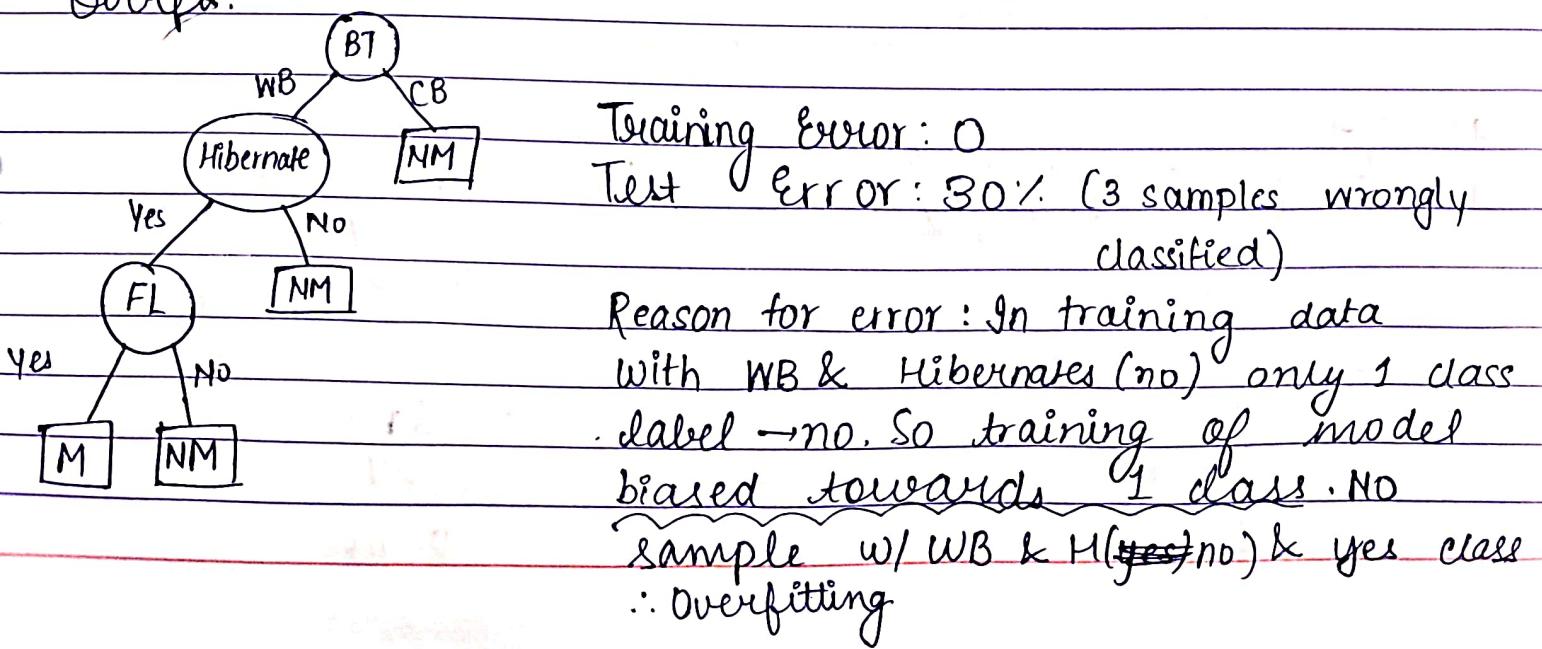
Both Train & Test error should be  
↓ within a certain node



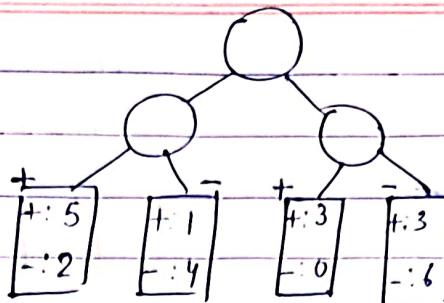
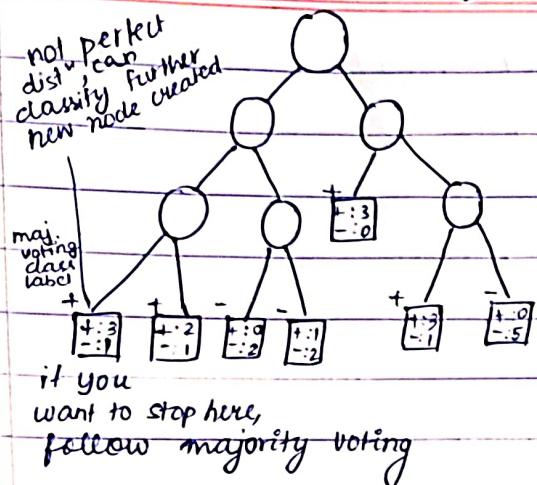
Underfitting → ↓ dataset



not always you can ↓ error → left w/  
10% error remains in Deci. Tree due to lesser.)  
Overfit.



Q. How to find amt of training & test error in the 2 decision trees?



$$\text{Training Error} = \frac{2+1+0+3}{24} \times 100 = 25\%$$

$$\begin{aligned}\text{Training Error} &= 1+1+0+1+0+1+0 \\ &= \frac{4}{24} \times 100 \\ &= 16\%\end{aligned}$$

+ve class label, 1  
-1 misclassified

This training error can be assumed to be equal to test error as we don't have test set

Misclassification of sample is not considering complexity of model (no. of nodes). Include that for better evaluation

$$E_g(T) = \frac{\sum_{i=1}^k (e_i(T) + \Omega_i(T))}{N(T)}$$

↑ no. of leaf node  
and of error on leaf node  
cost assoc'd w/ that node  
Cover fitted DT may  
have many nodes  
so consider cost of  
creating new node  
in DT)

a)  $\Omega = 0.5$

$$\begin{aligned}DT_1 &= \frac{4 + 7 \times 0.5}{24} \quad \text{no. of nodes} \\ &= 0.3125 \quad \text{penalty for each node}\end{aligned}$$

$$\begin{aligned}DT_2 &= \frac{6 + 4 \times 0.5}{24} \\ &= 0.333\end{aligned}$$

b)  $\Omega = 1$

$$\begin{aligned}DT_1 &= \frac{4 + 7 \times 1}{24} \\ &= 0.458\end{aligned}$$

$$\begin{aligned}&= \frac{6+4}{24} \\ &= 0.416\end{aligned}$$

$\downarrow$  node  $\downarrow$  error good

$\Omega$  can be used to control overfitting in your decision tree. (tuning parameter)

$\Omega \rightarrow$  penalty assoc'd on each node

If you consider creation of new nodes & you are not getting ~~the~~ benefit, then no need of more overfitting

Test.  $X = (\text{Home Owner} = \text{No}, \text{Marital Status} = \text{Married}, \text{income} = 120) = ?$

A -  $P(Y = \text{yes}) = \frac{3}{10} \quad P(Y = \text{no}) = \frac{7}{10}$  continuous ↑

$$P(X|Y = \text{yes}) = P(\text{HO} = \text{No} | Y = \text{yes}) \times P(\text{MS} = \text{Married} | Y = \text{yes}) \times P(\text{income} = 120) = ?$$

$$= \frac{3}{3} \times \frac{0}{3} \times P(\text{income} = 120) = 0$$

$$\begin{aligned} P(X|Y = \text{no}) &= P(\text{HO} = \text{No} | Y = \text{no}) \times P(\text{MS} = \text{married} | Y = \text{no}) \times P(\text{income} = 120 | Y = \text{no}) \\ &= \frac{4}{7} \times \frac{4}{7} \times P(\text{income} = 120 | Y = \text{no}) \end{aligned}$$

$$P(\text{income} = 120 | Y = \text{yes}) = \frac{1}{\sqrt{2\pi} q_j} \exp^{-\frac{(120 - \mu_j)^2}{2q_j^2}} = 0.007254 \text{ (assume)}$$

$$P(A|Y = \text{no}) = 0.0072$$

$$= \frac{4}{7} \times \frac{4}{7} \times 0.0072 = 0.0023$$

$$\begin{aligned} P(Y = \text{No} | X) &= P(X|Y = \text{No}) \times P(Y = \text{No}) \\ &= 0.0023 \times \frac{7}{10} \\ &= 0.0016 \xrightarrow{\text{maxm}} \text{No} \checkmark \end{aligned}$$

### Problem w/ Naive Bayes Classifier

If var are cond<sup>nd</sup> indep then its a v. good classifier & frequently used

Since it is a product of cond<sup>nd</sup> probab.

If any cond<sup>nd</sup> probab becomes zero, then whole prob. becomes 0, so decision gets biased towards another

Test.  $X = (\text{Home Owner} = \text{No}, \text{Marital Status} = \text{Married}, \text{income} = 120) = ?$

A -  $P(Y = \text{yes}) = \frac{3}{10} \quad P(Y = \text{no}) = \frac{7}{10}$  continuous ↑

$$P(X|Y = \text{yes}) = P(\text{HO} = \text{No} | Y = \text{yes}) \times P(\text{MS} = \text{Married} | Y = \text{yes}) \times P(\text{income} = 120) = ?$$

$$= \frac{3}{3} \times \frac{0}{3} \times P(\text{income} = 120) = 0$$

$$\begin{aligned} P(X|Y = \text{no}) &= P(\text{HO} = \text{No} | Y = \text{no}) \times P(\text{MS} = \text{married} | Y = \text{no}) \times P(\text{income} = 120 | Y = \text{no}) \\ &= \frac{4}{7} \times \frac{4}{7} \times P(\text{income} = 120 | Y = \text{no}) \end{aligned}$$

$$P(\text{income} = 120 | Y = \text{yes}) = \frac{1}{\sqrt{2\pi} q_j} \exp^{-\frac{(120 - \mu_j)^2}{2q_j^2}} = 0.007254 \text{ (assume)}$$

$$P(A|Y = \text{no}) = 0.0072$$

$$= \frac{4}{7} \times \frac{4}{7} \times 0.0072 = 0.0023$$

$$\begin{aligned} P(Y = \text{No} | X) &= P(X|Y = \text{No}) \times P(Y = \text{No}) \\ &= 0.0023 \times \frac{7}{10} \\ &= 0.0016 \xrightarrow{\text{maxm}} \text{No} \checkmark \end{aligned}$$

### Problem w/ Naive Bayes Classifier

If var are cond<sup>nd</sup> indep then its a v. good classifier & frequently used

Since it is a product of cond<sup>nd</sup> probab.

If any cond<sup>nd</sup> probab becomes zero, then whole prob. becomes 0, so decision gets biased towards another

Joint Probability  
Conditional "

$$P(x, y)$$

$$P(x|y)$$

$$P(x, y) = P(y|x) P(x) = P(x|y) P(y)$$

$$P(y|x) = \frac{P(x|y) P(y)}{P(x)}$$

Bayes Theorem

	x				C
	$x_1$	$x_2$	$x_3$	$x_4$	
S	-		+		y
M	-		+		
S		-			y
M	-		+		
L			-		

Training Phase:  $P(y|x)$

Test Phase  $x'$   $P(y'|x') = P(y=\text{Yes}|x') = 0.25$ ,  
 $P(y'=\text{No}|x') = 0.33$

Max value  $\rightarrow$  That will be  
the prediction

$$P(y=\text{Yes}|x) = \frac{P(x|y=\text{Yes}) \times P(y=\text{Yes})}{P(x)}$$

$$P(y=\text{No}|x) = \frac{P(x|y=\text{No}) \times P(y=\text{No})}{P(x)}$$

Evaluate & compare

(denominator same so no need to compare)

$$P(y=\text{Yes}) = \frac{3}{5} \quad P(y=\text{No}) = \frac{2}{5}$$

reading skill arm length Age

↑ in arm length  $\rightarrow$  reading skill when Age is varied  
if Age fixed, reading & arm cond<sup>nal</sup> independent

$$P(x|y=y) = \prod_{i=1}^d P(x_i|y=y)$$

If var are cond. indep their  
prob = product of individual  
cond<sup>nal</sup> prob.

Joint Probability  
Conditional "

$$P(x, y)$$

$$P(x|y)$$

$$P(x, y) = P(y|x) P(x) = P(x|y) P(y)$$

$$P(y|x) = \frac{P(x|y) P(y)}{P(x)}$$

Bayes Theorem

	x				C
	$x_1$	$x_2$	$x_3$	$x_4$	
S	-		+		y
M	-		+		
S		-			y
M	-		+		
L			-		

Training Phase:  $P(y|x)$

Test Phase  $x'$   $P(y'|x') = P(y=\text{Yes}|x') = 0.25$ ,  
 $P(y'=\text{No}|x') = 0.33$

Max value  $\rightarrow$  That will be  
the prediction

$$P(y=\text{Yes}|x) = \frac{P(x|y=\text{Yes}) \times P(y=\text{Yes})}{P(x)}$$

$$P(y=\text{No}|x) = \frac{P(x|y=\text{No}) \times P(y=\text{No})}{P(x)}$$

Evaluate & compare

(denominator same so no need to compare)

$$P(y=\text{Yes}) = \frac{3}{5} \quad P(y=\text{No}) = \frac{2}{5}$$

reading skill arm length Age

$\uparrow$  in arm length  $\rightarrow$  reading skill when Age is varied  
if Age fixed, reading & arm cond<sup>nal</sup> independent

$$P(x|y=y) = \prod_{i=1}^d P(x_i|y=y)$$

If var are cond. indep their  
prob = product of individual  
cond<sup>nal</sup> prob.

$$\begin{aligned}
 A &= 2 \times f_+ = 50 & f_- = 5 \\
 R(R_1) &= 2 \times \left[ f_+ \log_2 \left( \frac{f_+}{e_+} \right) + f_- \log_2 \left( \frac{f_-}{e_-} \right) \right] \\
 &= 2 \times \left[ 50 \log_2 \left( \frac{50}{34.375} \right) + 5 \log_2 \left( \frac{5}{20.625} \right) \right] \\
 &= 2 \times \left[ 50 \log_2 (1.45) + 5 \log_2 (0.24) \right] \\
 &= 2 \times [50 \times 0.53 + 5 \times (-2.05)] \\
 &= 2 \times 16.25 \\
 &= 32.5
 \end{aligned}$$

$$\begin{aligned}
 R(R_2) &= 2 \times \left[ 2 \times \log_2 \left( \frac{2}{125} \right) + 0 \right] \\
 &= 2 \times 2 \times \log_2 1.6 \\
 &= 4 \times 0.678 \\
 &= 2.712
 \end{aligned}$$

$$\begin{aligned}
 R(R_1) &= 32.5 & \checkmark & (\text{higher one is better}) & (R_1 > R_2) \\
 R(R_2) &= 2.712 & & & (R_1 \text{ better rule})
 \end{aligned}$$

\* Foil Gain new term  
in exam will  
be covered by R1

$R_1 : A \rightarrow + P_0, n_0$  } which one is better? what gain you get in  
 $R_2 : A \wedge B \rightarrow + P_1, n_1$  } the extended rule

$$\text{Foil-gain} = P_1 \times \left( \log_2 \frac{P_1}{P_1+n_1} - \log_2 \frac{P_0}{P_0+n_0} \right) \quad (\text{coverage diff. in acc.})$$

$$D = +: 100$$

$$-: 60$$

$$\begin{aligned}
 R_1 : A \rightarrow + & & 50 & 5- & 4 & \text{coverage } \frac{5}{60}, \text{ if rule is} \\
 & & & & & \text{accepted}
 \end{aligned}$$

$$\begin{aligned}
 R_2 : A \wedge B \rightarrow + & & 1 & & & \\
 & & & & &
 \end{aligned}$$

Find  
gain if > th  
then accept

$$\begin{aligned}
 A &= 2 \times f_+ = 50 & f_- = 5 \\
 R(R_1) &= 2 \times \left[ f_+ \log_2 \left( \frac{f_+}{e_+} \right) + f_- \log_2 \left( \frac{f_-}{e_-} \right) \right] \\
 &= 2 \times \left[ 50 \log_2 \left( \frac{50}{34.375} \right) + 5 \log_2 \left( \frac{5}{20.625} \right) \right] \\
 &= 2 \times \left[ 50 \log_2 (1.45) + 5 \log_2 (0.24) \right] \\
 &= 2 \times [50 \times 0.53 + 5 \times (-2.05)] \\
 &= 2 \times 16.25 \\
 &= 32.5
 \end{aligned}$$

$$\begin{aligned}
 R(R_2) &= 2 \times \left[ 2 \times \log_2 \left( \frac{2}{125} \right) + 0 \right] \\
 &= 2 \times 2 \times \log_2 1.6 \\
 &= 4 \times 0.678 \\
 &= 2.712
 \end{aligned}$$

$$\begin{aligned}
 R(R_1) &= 32.5 & \checkmark & (\text{higher one is better}) & (R_1 > R_2) \\
 R(R_2) &= 2.712 & & & (R_1 \text{ better rule})
 \end{aligned}$$

\* Foil Gain new term  
in exam rule  
we covered by R<sub>1</sub>

R<sub>1</sub>: A → + P<sub>0,n<sub>0</sub></sub> { which one is better? what gain you get in  
 R<sub>2</sub>: A <sup>AB</sup> → + P<sub>1,n<sub>1</sub></sub> } the extended rule

$$\text{Foil-gain} = P_1 \times \left( \log_2 \frac{P_1}{P_1+n_1} - \log_2 \frac{P_0}{P_0+n_0} \right) \quad (\text{coverage diff. in acc.})$$

$$D = +: 100$$

$$- : 60$$

$$\begin{aligned}
 R_1: \quad A \rightarrow + & \quad 50 \quad 5-4 \quad \text{coverage } \frac{5}{60}, \text{ if rule is} \\
 & \quad \text{coverage } \frac{4}{60}, \text{ if rule is}
 \end{aligned}$$

$$GR_2: A \uparrow B \rightarrow + \quad 1$$

Find  
gain if > th  
then accept

## Issues:

How to know 1 rule is better off than the other

### Training dataset

Eg - +: 100  
- : 60

$R_1$ : 50+      5-

$R_2$ : 2+      0-

Accuracy = out of covered, how many have the correct class. (same class as the rule)

$R_1$

Let measurement lie towards +ve class

$$\text{Accuracy}(R_1) = \frac{50}{55} = 90\%$$

$$(R_2) = \frac{2}{2} = 100\% \quad \left. \begin{array}{l} \text{Problem w/ Accuracy} \rightarrow \text{X good} \\ \text{choice if data imbalanced} \end{array} \right.$$

### Likelihood Ratio statistic

$$R = 2 \sum_{i=1}^k \frac{f_i \log_e \left( \frac{f_i}{e_i} \right)}{\text{accuracy of rule } \left( \frac{f_i}{e_i} \right)}$$

$k$  = no. of classes

$f_i$  = observed frequency of class  $i$  examples covered by the rule

$e_i$  = expected freq<sup>u</sup> of a rule that makes random predictions

$$R_1 \quad e^+ = \frac{55 \times 100}{160} \quad e^- = \frac{55 \times 60}{160} \rightarrow 60 \text{ samples -ve} \quad \left. \begin{array}{l} \text{total coverage/freq} \\ \text{Total rule} \end{array} \right. \quad \left. \begin{array}{l} \text{55 covered} \\ \text{samples -ve} \end{array} \right. \quad \left. \begin{array}{l} \text{expected frequencies} \end{array} \right.$$

$$R_2 \quad e^+ = 2 \times \frac{100}{160} \quad e^- = 2 \times \frac{60}{160} \quad \text{probab. belonging to +ve/-ve}$$

$$\text{Accuracy}(R) = \frac{1}{2} = 50\%$$

1 Coverage means acceptability of rule is ↑

$$R_2(A_1=a_1) \Rightarrow - (1, 2, 4)$$

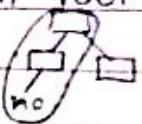
conflict w/R<sub>1</sub> if A<sub>1</sub>=a<sub>1</sub> in prediction

- ① ↑ priority w/o rule w/ more attribute tests
- ② one w/ ↓ misclassification rule

To get rules:

From dataset prepare a DT. For every path from root to leaf, write a rule for it

$$(Age=young) \wedge (\text{student}=no) \Rightarrow \text{buys-computer}=no$$



Issue: Rule set may have redundant / conflicting rules will be observed in rule set

Rule set must be
 

- (one sample covered by one rule)
- mutually exclusive & exhaustive
- (some samples cannot be covered by ≥ 1 rule)
- (all samples in dataset must be covered by at least 1 rule)

## 2 Sequential covering method

Rules are learnt 1 at a time

Write some rule for a class

$$X \Rightarrow Y$$

1, 2 satisfies → remove

threshold → Eg- accuracy should be  $\geq 60\%$

$$x \wedge z \Rightarrow Y \quad \text{Acc} = 40\% X$$

$$x \Rightarrow Y \quad \text{Acc} = 80\% \checkmark$$

## Issues:

How to know 1 rule is better off than the other

### Training dataset

Eg - +: 100  
- : 60

$R_1$ : 50+      5-

$R_2$ : 2+      0-

Accuracy = out of covered, how many have the correct class. (same class as the rule)

$R_1$

Let measurement lie towards +ve class

$$\text{Accuracy}(R_1) = \frac{50}{55} = 90\%$$

$$(R_2) = \frac{2}{2} = 100\% \quad \left. \begin{array}{l} \text{Problem w/ Accuracy} \rightarrow \text{X good} \\ \text{choice if data imbalanced} \end{array} \right.$$

### Likelihood Ratio statistic

$$R = 2 \sum_{i=1}^k \frac{f_i \log_e \left( \frac{f_i}{e_i} \right)}{\text{accuracy of rule } \left( \frac{f_i}{e_i} \right)}$$

$k$  = no. of classes

$f_i$  = observed frequency of class  $i$  examples covered by the rule

$e_i$  = expected freq<sup>u</sup> of a rule that makes random predictions

$$R_1 \quad e^+ = \frac{55 \times 100}{160} \quad e^- = \frac{55 \times 60}{160} \rightarrow 60 \text{ samples -ve} \quad \left. \begin{array}{l} \text{total coverage/freq} \\ \text{Total rule} \end{array} \right. \quad \left. \begin{array}{l} \text{55 covered} \\ \text{samples -ve} \end{array} \right. \quad \left. \begin{array}{l} \text{expected frequencies} \end{array} \right.$$

$$R_2 \quad e^+ = 2 \times \frac{100}{160} \quad e^- = 2 \times \frac{60}{160} \quad \text{probab. belonging to +ve/-ve}$$

- ↪ node is error good
- ↪ can be used to control overfitting in your decision tree (tuning parameter)
- ↪ penalty also on each node

If you consider creation of new nodes & you are not b getting benefit, then no need of more overfitting

## \* Overfitting & Tree Pruning

Lack of samples

Noise in data

2 ways to handle overfitting :

- 1) Prepruning (check cond' of of while adding nodes to DT)
- 2) Postpruning (construct normally, then check cond', if accuracy ↑ after deleting nodes then do it)

## \* Rule based classifier

Prepare rule set from data set

	$A_1$	$A_2$	$A_3$	C
1	$a_1$	$a_1$	$x_1$	+
2	$a_1$	$a_1$	$x_2$	-
3	$a_1$	$a_3$	$x_3$	+
4	$a_1$	$a_2$	$x_2$	-

Rule set | used to decide class label of test set

Test       $A_1 = a_1, A_2 = a_2, A_3 = a_3 \& C = ?$

Rule :  $X \Rightarrow Y$

Antecedent Consequent

R<sub>1</sub>:  $(A_1 = a_1) \wedge (A_2 = a_2) \Rightarrow +$

$\text{Coverage}(R) = \frac{|X|}{|D|} \text{ no of antecedents}$

$\text{Accuracy}(R) = \frac{\text{Correct}}{\text{Covers}}$

coverage → how many instances your rule covers in dataset {1, 2, 4}

$\text{Coverage}(R) = \frac{|X|}{|D|} = \frac{2}{3} = 66.7\% = 67\%$  (4 was added later)

accuracy → whatever the coverage, what part how many have same class label as the rule

$$\frac{n_c + mp}{n+m}$$

To handle it, M-estimation

$$P(x_i | y_j) = \frac{n_c + mp}{n+m} \quad \text{cannot be 0}$$

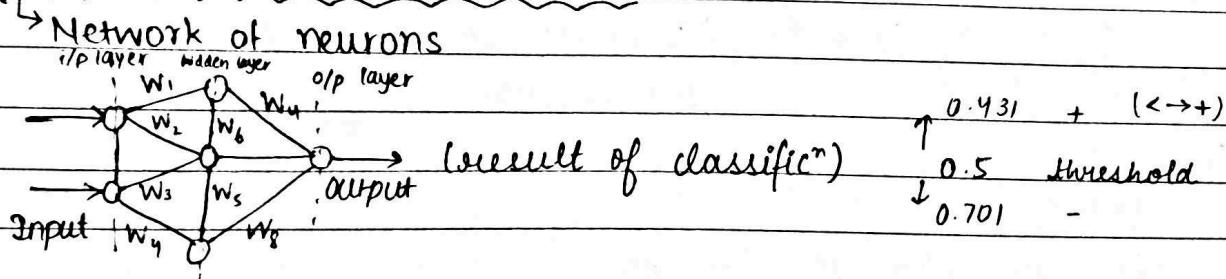
n: Total no. of instances from class  $y_j$  (yes class)

$n_c$ : No. of training examples from class  $y_j$  that take on the value  $x_i$

m: equivalent sample size

p: user defined parameter <sup>tuning param.</sup> (will be given in Qn)

## \* Neural Network Classifier



Objective: Get optimum weights corr. to diff. links in your network.

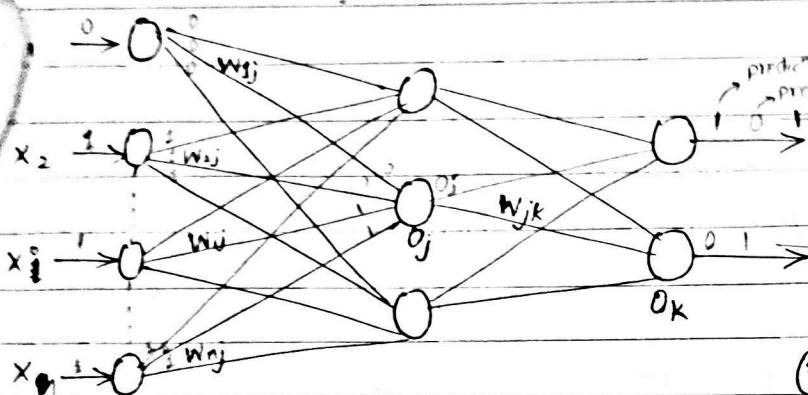
Do updates in wt based on instances in your sample

1 → . . .  
2 → . . .  
3 → . . .  
n → . . .  
no. of nodes ↑ → complexity ↑

set of nodes  
i/p layer: through which we receive i/p from dataset/env  
o/p layer: " " " " o/p  
hidden layer: b/w i/p & o/p layer

## Multilayer Feedforward Neural Network

Fully connected n/w



① Initialize all the weights and the bias in the network

- ② While terminating conditions is not satisfied
- ③ for each training tuple  $X$  in  $D$
- ④ // Propagate the i/p forward
- ⑤ for each i/p layer unit  $j$

⑥  $O_j = I_j$  simply fwd the received i/p

⑦ for each hidden or o/p layer unit  $j$

⑧  $I_j = \sum_i w_{ij} O_i + O_j$  // Compute the net i/p of unit  $j$  wrt the previous layer  $i$

$$⑨ O_j = \frac{1}{1 + e^{-I_j}}$$

⑩ // Backpropagating the error

⑪ for each unit  $j$  in the o/p layer

$$⑫ Err_j = O_j(1-O_j)(T_j - O_j)$$

⑬ for each unit  $j$  in the hidden layer from the last to the first hidden layer

$$⑭ Err_j = O_j(1-O_j) \sum_k Err_k w_{jk} \text{ // Compute the error wrt the next higher layer } k$$

from which error received  
from last to that node next higher layer  $k$

⑮ for each weight  $w_{ij}$  in the n/w

$$⑯ \Delta w_{ij} = (\lambda) Err_j O_i$$

$$⑰ w_{ij} = w_{ij} + \Delta w_{ij}$$

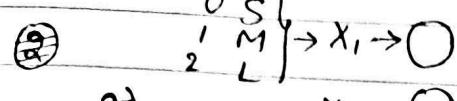
⑱ for each bias  $O_j$  in n/w

$$⑲ \Delta O_j = (\lambda) Err_j$$

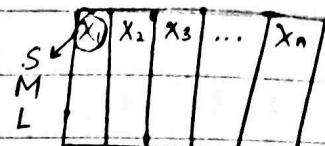
$$⑳ O_j = O_j + \Delta O_j$$

NN varies due to Entropy & information gain  
see → how to handle continuous attributes?

① As many features in dataset, consider that many nodes as I/P layer



b)  $x_2 \rightarrow \text{O}$



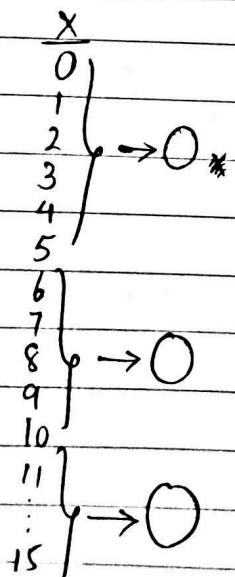
② For diff. values consider diff. nodes

$$\begin{matrix} S \rightarrow & x'_1 & x'_2 & x''_1 \\ M \rightarrow & 0 & 1 & 0 \\ L \rightarrow & 0 & 0 & 1 \\ \text{(or } 0 & 0 & 0 \text{)} \end{matrix}$$

Diff types of attributes

$x_2 \rightarrow \text{O}$

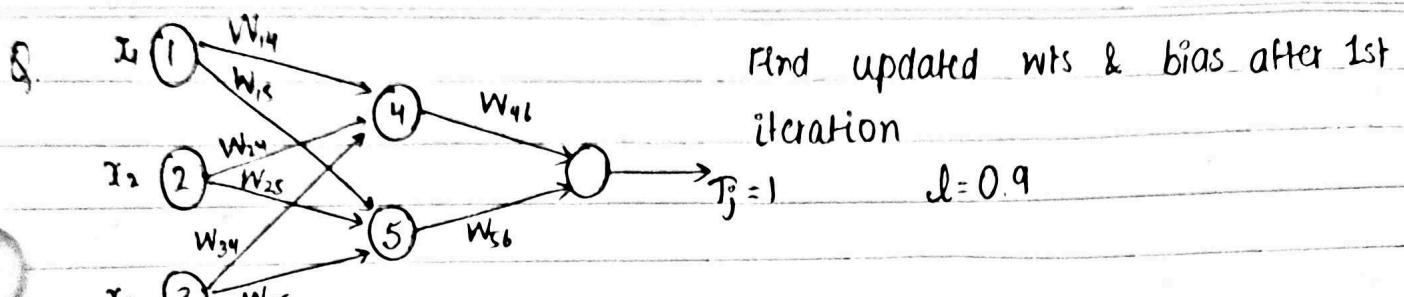
③ Continuous attribute



Node for each range

error

	$x_1$	$x_2$	$x_3$	$\dots$	$x_i$	$C$
1	0	1	0	-	1	0.531
2						0.421



	$x_1$	$x_2$	$x_3$	$w_{14}$	$w_{15}$	$w_{24}$	$w_{25}$	$w_{34}$	$w_{45}$	$w_{46}$	$w_{56}$	$o_4$	$o_5$	$o_6$	
①	1	0	1	0.2	-0.3	0.4	0.1	0.1	-0.5	0.2	-0.3	-0.2	-0.4	0.2	0.1

Ans -  $I_4 = \sum_i w_{ij} o_i + \theta_j$        $O_4 = \frac{1}{1+e^{0.7}} = 0.332$

$$\begin{aligned}
 &= w_{14}O_1 + w_{24}O_2 + w_{34}O_3 + \theta_4 \\
 &= 0.2 \times 1 + 0.4 \times 0 + (-0.5) \times 1 + (-0.4) \\
 &= -0.7
 \end{aligned}$$

$$\begin{aligned}
 I_5 &= w_{15}O_1 + w_{25}O_2 + w_{35}O_3 + \theta_5 & O_5 &= \frac{1}{1+e^{-0.1}} = 0.525 \\
 &= -0.3 + 0.1 \times 0 + 0.2 \times 1 + 0.2 \\
 &= 0.1
 \end{aligned}$$

$$\begin{aligned}
 I_6 &= w_{46}O_4 + w_{56}O_5 + \theta_6 & O_6 &= \frac{1}{1+e^{0.105}} = 0.474 \\
 &= -0.3 \times 0.332 - 0.2 \times 0.525 + 0.1 \\
 &= -0.105 \text{ (rounded off)}
 \end{aligned}$$

$$\begin{aligned}
 Err_6 &= O_6(1-O_6)(T_6 - O_6) \\
 &= 0.474(1-0.474)(1-0.474) = 0.1311
 \end{aligned}$$

$$\begin{aligned}
 Err_5 &= O_5(1-O_5)(Err_6 W_{56}) \\
 &= 0.525(1-0.525)(0.1311)(-0.2) = -0.0065
 \end{aligned}$$

$$\begin{aligned}
 Err_4 &= O_4(1-O_4)(Err_6 W_{46}) \\
 &= 0.332(1-0.332)(\cancel{0.474} 0.1311 \times (-0.3)) = -0.0087
 \end{aligned}$$

$$\Delta W_{14} = (d) Err_{14} \cdot \sigma_j \\ = 0.9 \times (-0.008)$$

$d$ , weight 1

$$W_{ij} = W_{ij} + \Delta W_{ij} \\ = 0.2 + 0.0078$$

role in prediction

$$\Delta O_4 = (d) Err_4$$

$$O_4 = O_4 + \Delta O_4 \\ = -0.4 + \dots$$

## \* KNN Classifier

K Nearest Neighbour

Prepare model  $\rightarrow$  test it on data  $\rightarrow$  too much time in training

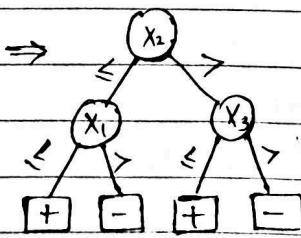
Classifiers  $\swarrow$  Lazy learner/classifier (Bayesian, Decision Tree)

Eager learner/classifier

Reduced training time

Construction of training model is not done until you have the test sample (postponed)

$x_1$	$x_2$	$x_3$	$C$
-	-	-	+
-	+	-	+



Objective: training sample must map on the model

Test:

$x_1$	$x_2$	$x_3$	$\Rightarrow$	+
+	-	-		

(They try to find a sample matching the test sample in training sample then return that class (comparison algorithm fast))

Drawback: If test sample doesn't match any training sample

Sol: Instead of finding actual match, find the

$$\Delta W_{ij} = (l) E_{ir_4} O_i \\ = 0.9 \times (-0.0087) \times (1) = 0.0078$$

$$W_{ij} = W_{ij} + \Delta W_{ij} \\ = 0.2 + 0.0078$$

$$\Delta O_4 = (l) E_{ir_4}$$

$$O_4 = O_4 + \Delta O_4 \\ = -0.4 + \dots$$

## \* KNN Classifier

(K Nearest Neighbour)

Prepare model  $\rightarrow$  test it on data  $\rightarrow$  too much time  
in training

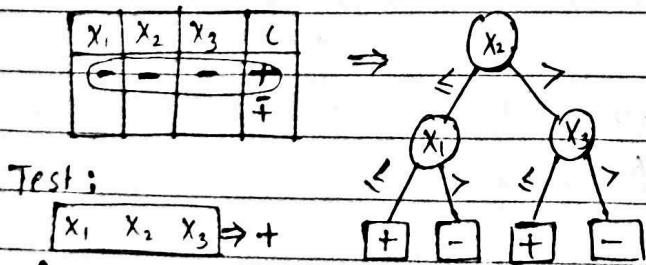
classifiers

Lazy learner/classifier (Bayesian, Decision Tree)

Eager learner/classifier

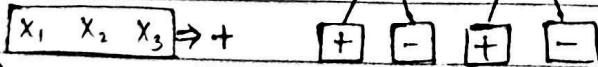
reduced training time

Construction of training model is not done until you have the test sample (postponed)



Objective: training sample must map on the model

Test:



(They try to find a sample matching the test sample in training sample then return that class (comparison algo fast))

Drawback: If test sample doesn't match any training sample

Soln: Instead of finding actual match, find the

closest training sample w/ test sample

Test similarity (sim measure as per type of data)

$T_1 \rightarrow 0.1 \rightarrow +$        $K=3$  (3 nearest neighbours)

$T_2 \rightarrow 0.2 \rightarrow -$       3 closest neighbours considered

$T_3 \rightarrow 0.3 \rightarrow +$

$T_4 \rightarrow 0.4$

$T_5 \rightarrow 0.5$

Predict based on no. of nbrs  
you select

$T_1 \rightarrow +$   
 $T_2 \rightarrow -$   
 $T_3 \rightarrow +$

majority voting  $\Rightarrow +$

↳ class label of the training set

$K=2$

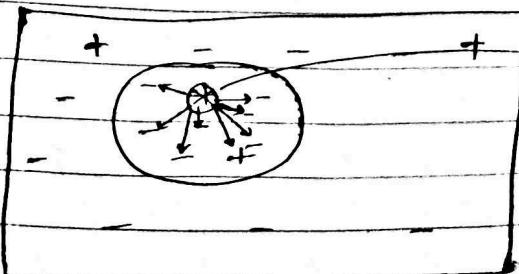
$\rightarrow +$  or  $\rightarrow -$  anyone (or any tie breaking cond'  
if you define)

Problem: If you don't select  $K$  value wisely, its performance  
won't be good

$K$  v. small : overfitting

$K$  v. large : misclassification or underfitting

So many variants to select  $K$  optimally



Test sample

We follow majority, but here see many samples in circle are v. far & only some at close distance

Sometimes weight is defined,  
 $w = \frac{1}{d(x, x')^2}$  If distance is 1, weight 1

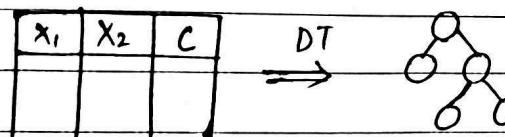
so closer sample is playing bigger role in prediction  
 Majority voting → all samples playing equal role

$$= w \cdot \underbrace{(x - y')}_\text{output} \rightarrow 1 \text{ or } 0 \text{ predict}$$

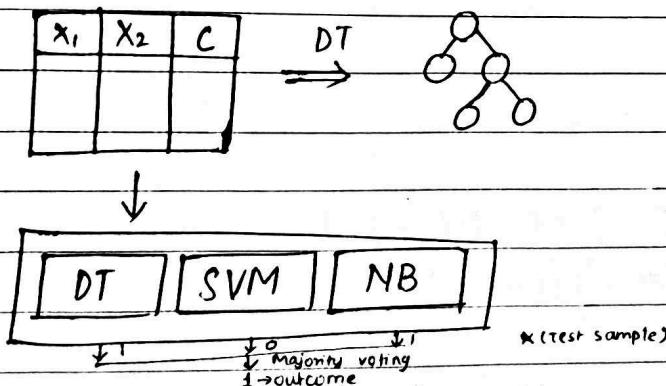
Multiply w/ weight then predict

### \* Ensemble classifier

Instead of using a single classifier, we can use diff classifiers on the same layer.



Ensembling : mixing diff classifiers



Ensemble method gives good predictions as compared to indiv. classifiers

Approaches to do Ensembling

- ① Feature level ensembling
- ② Data " "
- ③ Algorithm " "

①

X <sub>1</sub>	X <sub>2</sub>	X <sub>3</sub>	X <sub>4</sub>	X <sub>5</sub>	X <sub>6</sub>

Finding a subset of features used in 1 classifier  
 $\{x_1, \dots, x_6\} \rightarrow \{x_1, x_2, x_3\}, \{x_2, x_3, x_6\}, \{x_4, x_5, x_6\}$

↓  
DT

↓  
SVM

↓  
NB

Take care what kind of data is good for a certain type of classifier.

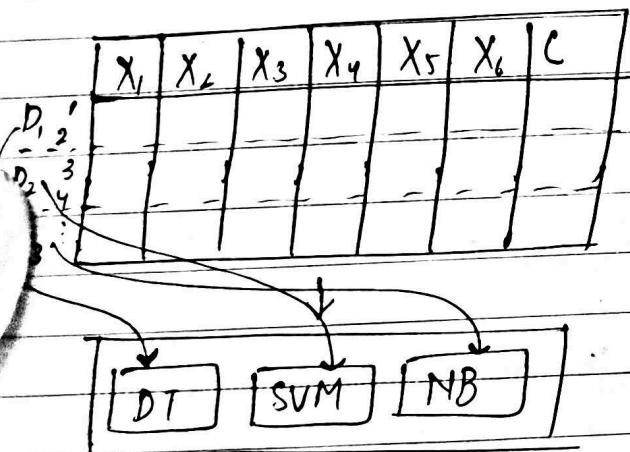
DT good if discrete attributes

NB " " conditionally independent attributes

To get result:

- 1) Majority Voting
- 2) Weighted aggregation (wt. to features/classifier)

## ② Data level Ensembling



## ③ Algorithm level

DT SVM NB, DT SVM KNN, DT SVM Rule based

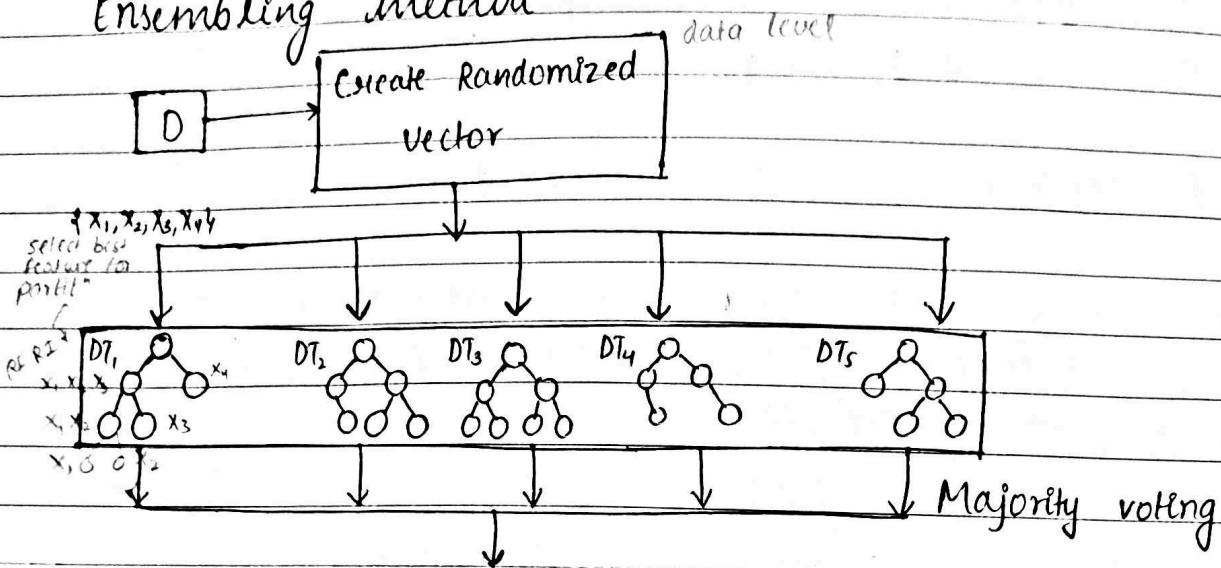
Diff. combinations of classifiers

\*Disadv of Ensembling: Complexity (lot of classifiers, redundancy of features in feature level) For some app" it may not be affordable

Challenge → Lightweight ensembling method

## 1. Random Forest Classifier

Ensembling method



RF - RI (Random Forest with Random Input)

RF - R+C (" " " combination)

### \* Characteristics of Random Forest

feature class

X	Y
0.2	1
0.3	1
0.1	-1
0.4	1
0.6	1
0.5	-1
0.7	1
1	-1

Probability Distribution func is used for all classifiers to create the Randomized vector

They use multiple decision trees as classifiers

Data & algo level ensembling

sampling w/ replacement (method to create randomized vector)

X	Y	used to train
0.2	1	DT <sub>1</sub> (have
0.2	1	their own
0.3	1	depth & no. of
0.4	1	nodes)
0.6	1	
0.6	1	
0.7	1	
0.2	1	

DT<sub>1</sub> → DT<sub>2</sub>    DT<sub>3</sub> → DT<sub>3</sub>

## ⇒ RF - RI

If dataset is large & has sufficient features → cond<sup>n</sup>  
Instead of looking for best partition in DT, they  
use random partition  $\{x_1, x_2, x_3, x_4\} \leftarrow \{x_1, x_2, x_3\}$   $\{x_4\}$

Objective: Randomization w/ feature & data

If dataset has ↓ features → many features are common  
over diff. decision trees → performance ↓ →  
use RF-RC to avoid this

## ⇒ RF - RC

$$\begin{matrix} x_1 & x_2 \\ \cancel{x_1} & \cancel{x_2} \end{matrix} \quad \begin{matrix} x_1 & x_2 \\ \cancel{x_1} & \cancel{x_2} \end{matrix}$$

create new features using linear comb of features

$$x_3 = Gx_1 + Gx_2$$

use them to create diff depth DTs

avoids redundancy of features → helps in randomization

## \* Association Analysis / Pattern Mining

2, 3, 5, 9, 10, 12, 3, 5, (10, 13, 15), 18, 19, (10, 13), 15, 21, 20, (10, 13, 15)

Recommender system uses pattern mining to give  
sugges"n. Eg - Buy shirt → can buy pant

Transac"n: visit of customer in market & the item he purchased  
Market-Basket: Record of items purchased & analysing it  
Analysis

## Terminologies

- 1) item : any product customer purchases in mkt visit (Beer, Nuts etc.)
- 2) item set : <sup>1-itemset</sup> {Beer} <sup>2-itemset</sup> {Beer, Diaper} <sup>3-itemset</sup> {Beer, Diaper, Egg} <sup>4-itemsset</sup> {Beer, Nuts, Diaper}
- 3) support count: No. of transac<sup>n</sup> where an itemset is present
- 4) support
- 5) confidence
- 6) frequent itemset

S.C.: no of transac<sup>n</sup> supporting an itemset & what is its count

of Beer, Diaper  $\Rightarrow S.C. = 3$

It is a pattern  $\rightarrow$  it appears in different transac<sup>n</sup>s

patterns in data are represented by an association rule



Beer  $\rightarrow$  Diaper (or Diaper  $\rightarrow$  Beer)

Nuts  $\rightarrow$  Milk

Beer, Nuts, Diaper

Beer, Coffee, Diaper

Beer, Diaper, Egg

Nuts, Eggs, Milk

Nuts, Coffee, Diaper, Eggs, Milk

Some patterns may be useful for the app, some not

$\Rightarrow$  Measures to judge usefulness of an association rule

$$\text{Support} = \frac{\sigma(X \cap Y)}{N} - \frac{\text{support count of } X \cap Y}{\text{total transactions}}$$

Out of total transac<sup>n</sup> in dataset how many have  $X \cap Y$  (in a group, how many in favour)

$$\text{support}(Beer \rightarrow Diaper) = \frac{3}{5} = 60\%$$

$$\text{Confidence} = \frac{\sigma(X \cap Y)}{\sigma(X)} - \frac{\text{support ct of } X \cap Y}{\text{support ct of } X}$$

Where Beer is present, diaper is also present

$$= \frac{3}{5} = 100\%$$

out of those  $\frac{3}{5}$  supporting what is the level of confidence, how many follow

Itemset is said to be frequent if its support > min defined support threshold

{Beer, Diaper} → 60% > Min support (of app)  
Eg - 40%

Objective: Identify frequent itemsets → make business decisions accordingly

{Beer, Diaper, Milk} → so many other itemsets can be formulated  
sub patterns → sub itemsets

{Beer, Diaper} & {diaper, milk} ...

Transform to binary dataset →

6 features →

Max<sup>m</sup> length pattern is 6

(freq itemset w/ longest length)

{Beer, Nuts ... Milk}

$2^6 - 1$  subsets / subpatterns

	$x_1$	$x_2$	$x_3$	$x_4$	$x_5$	$x_6$
1	1	1	1	0	0	0
2	1	0	1	1	0	0
3	1	0	1	0	1	0
4	0	1	0	0	1	1
5	0	1	0	1	1	1

100 features → so many subpatterns possible → checking sup & conf. for all is difficult. Solutions:

Closed itemset / closed pattern

min support = 2

{Beer, diaper} → support count = 3 (freq. ✓)

{Beer, diaper, Milk} → s.c. = 0 ≠ 3

i) frequent

ii) no superset has same support

A-

## Max-pattern

1) itemset frequent

2) no ~~superpattern~~ frequent superpattern / superset

Max<sup>m</sup> frequent itemset (supcount max<sup>m</sup>)

## Downward closure property

Eg- Tid	items	(a) Find frequent itemset with min. support = 2
1	{a,b,d,e}	(b) find all association rules belonging to the frequent itemsets you obtained in a
2	{b,c,d}	(c) Use confidence based pruning to remove all association rules having threshold < 80%
3	{a,b,d,e}	(d) Also draw a lattice based structure to prune the association rules as per confidence threshold of 30%.
4	{a,c,d,e}	
5	{b,c,d,e}	
6	{b,d,e}	
7	{c,d}	
8	{a,b,c} <sup>see</sup>	(e) Draw a hash-Tree structure for candidate 3-itemsets you obtained in Apriori generation step of (a) using the defined hash func H(P) = h(P) mod 3
9	{a,d,e}	<sup>do</sup>
10	{b,d}	

A - (a)	{a} 5 {b} 7 {c} 5 {d} 9 {e} 6	{a,b} ✓ 3 {a,c} ✓ 2 {a,d} ✓ 4 {a,e} ✓ 4 {b,c} ✓ 3 {b,d} ✓ 6 {b,e} ✓ 4 {c,d} ✓ 4 {c,e} ✓ 2 {d,e} ✓ 6	{a,b,c} ✓ 1 {a,b,d} ✓ 2 {a,b,e} ✓ 2 {a,b,c,d} ✓ 1 {a,b,c,e} ✓ 1 {a,b,d,e} ✓ 1 {a,c,d,e} ✓ 1 {a,c,e} ✓ 1 {a,c,d,e} ✓ 1 {a,d,e} ✓ 4 {b,c,d} ✓ 2 {b,c,e} ✓ 1 {b,c,d,e} ✓ 1 {b,d,e} ✓ 4	c <sub>3</sub> {c,d,e} → {a,b,d,e} 2 {a,b,c,d} 0 {a,b,c,e} 1 → {a,b,d,e}
---------	---	--	--	---

(a)  $\{a, b, d, e\}$

(c)

$$\text{length } K = 4 \quad 2^k - 2 = 2^4 - 2 = 14 \text{ rules}$$

$$\frac{\text{#abde}}{\text{#fabde}} = \frac{2}{5} = 0.4 \quad \checkmark$$

$$a \rightarrow bdc$$

$$\frac{2}{7} = 0.28$$

$$b \rightarrow ade$$

$$\frac{2}{9} = 0.22$$

$$d \rightarrow abe$$

$$\frac{2}{6} = 0.33 \quad \checkmark$$

$$e \rightarrow abd$$

$$\frac{2}{6} = 0.33 \quad \checkmark$$

$$ab \rightarrow de$$

$$\frac{2}{3} = 0.67 \quad \checkmark$$

$$ad \rightarrow be$$

$$\frac{2}{4} = 0.5 \quad \checkmark$$

$$ae \rightarrow bd$$

$$\frac{2}{4} = 0.5 \quad \checkmark$$

$$bd \rightarrow ae$$

$$\frac{2}{6} = 0.33 \quad \checkmark$$

$$be \rightarrow ad$$

$$\frac{2}{4} = 0.5 \quad \checkmark$$

$$de \rightarrow ab$$

$$\frac{2}{6} = 0.33 \quad \checkmark$$

$$abd \rightarrow e$$

$$\frac{2}{2} = 1 \quad \checkmark$$

$$abc \rightarrow d$$

$$\frac{2}{2} = 1 \quad \checkmark$$

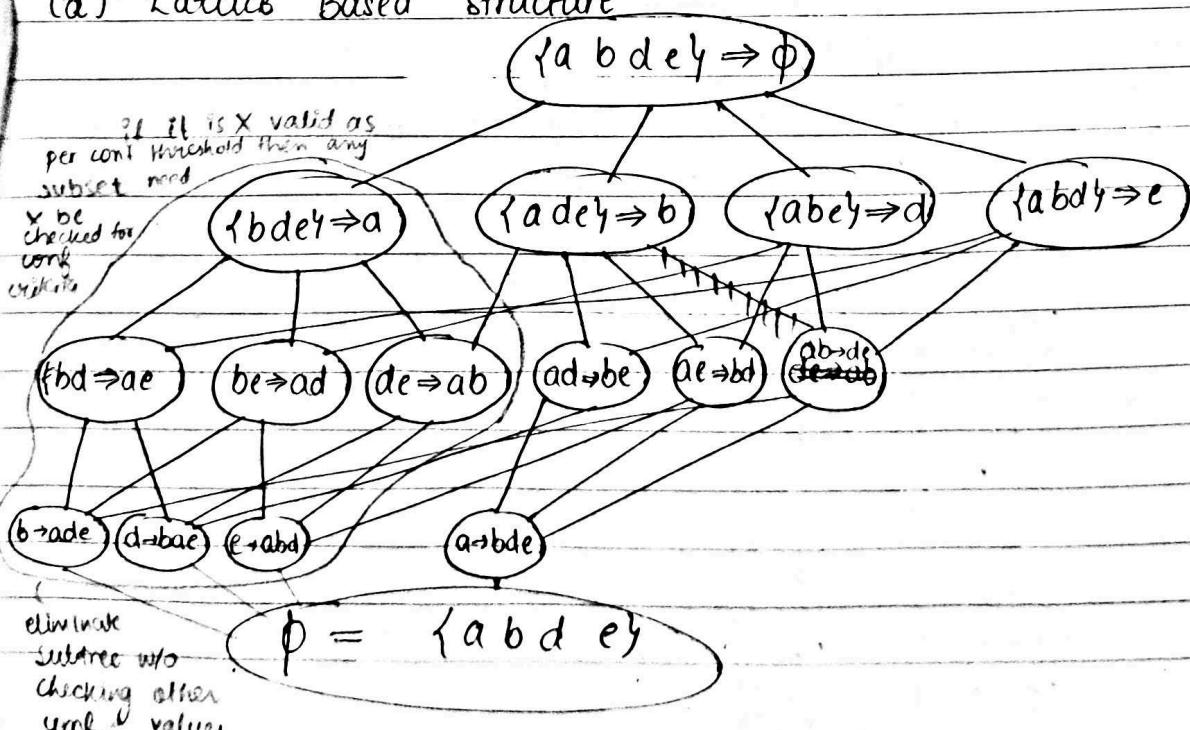
$$bde \rightarrow a$$

$$\frac{2}{4} = 0.5 \quad \checkmark$$

$$ade \rightarrow b$$

$$\frac{2}{4} = 0.5 \quad \checkmark$$

(d) Lattice Based structure



$$bde \rightarrow a \quad \text{conf}(R) = 0.5 = 60\% > 30\%$$

If conf = 75% X (suppose)

$$X \rightarrow Y - X$$

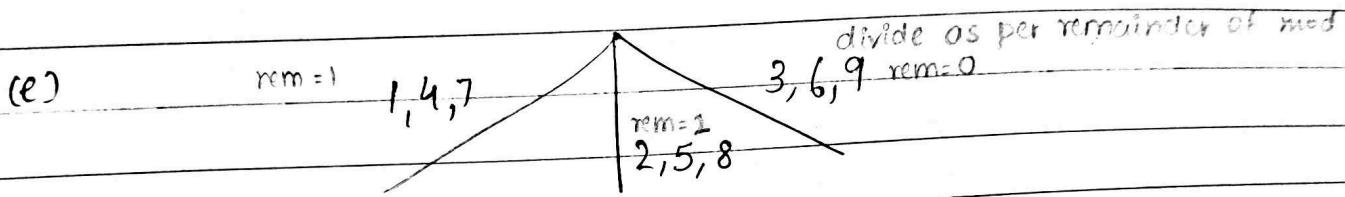
$$X' \subseteq X$$

$$X' \Rightarrow Y - X'$$

Prune subsets/subtree of  $bde \rightarrow a$

$adc \rightarrow b$	0.5 ✓	$abe \rightarrow d$ ✓	$abd \rightarrow e$ ✓
$bd \rightarrow ae$	0.33 ✓		

Lattice based procedure for pruning:



Raise candidate 3-itemsets

(1, 2, 3)

(4, 8, 2)

(1, 5, 3)

(2, 5, 6)

(9, 3, 6)

(2, 4, 5)

sorted  
order  
list: (5, 2, 6)  
(3, 8, 9)

(1, 5, 6)

(2, 4, 6)

(2, 3, 5)

(1, 3, 8)

(2, 9, 5)

Transaction buckets to check which transaction?

1, 2, 3, 5, 6

1+ 2, 3, 5, 6

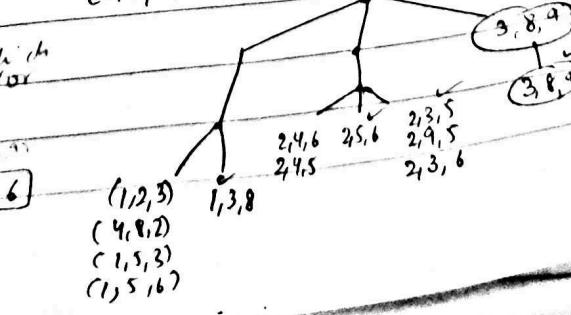
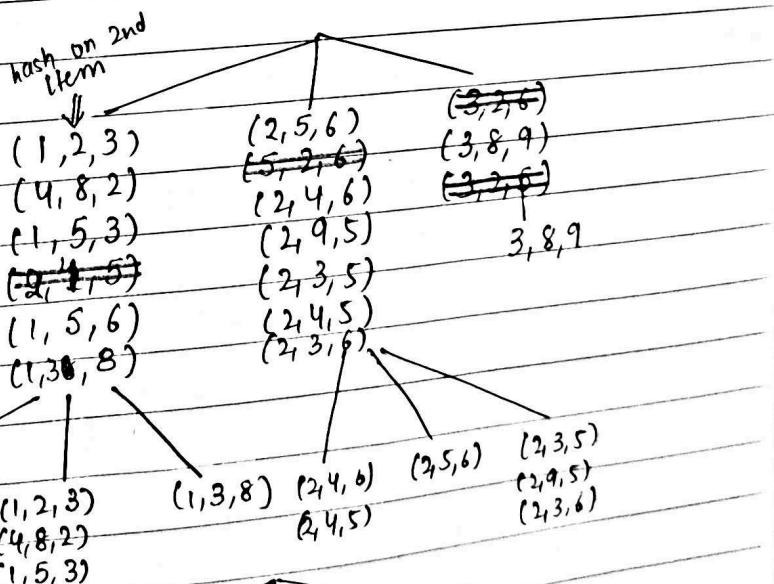
1 2 [3 5 6]

1 3 [5 6]

1 5 [6]

(1, 2, 3) hash on 1st item  $1 \mod 3 = 1 \rightarrow$  left of HT

(4, 8, 2)  $4 \mod 3 = 1 \rightarrow$  left



Issues in Apriori

improve w/ hash tree  $\rightarrow$  no. of comparisons

Trans id	Items	
T <sub>1</sub>	{M, O, N, K, E, Y}	
T <sub>2</sub>	{D, O, N, K, E, Y}	
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, J, E}	

Find all freq. Itemsets belonging to this dataset  
support = 3  
threshold

Item	C <sub>1</sub>	C <sub>2</sub>
{M}	3	{M, K} {M, O, K} 1
{O}	3	{O, K} {M, K, E} 2
{N}	2	{M, X} {M, X, Y} 2
{K}	5 $\Rightarrow$ {E} $\Rightarrow$ {Y}	{O, E} $\Rightarrow$ {O, K, E} 3 $\Rightarrow$ {O, K, E, Y}
{E}	4	{M, Y} {K, Y} {O, K, Y} 2
{Y}	3	{K, E, Y} 2
{D}	1	{O, Y} 2
{A}	1	{K, E} 4
{U}	1	{K, Y} 3
{C}	2	{E, Y} 2
{I}	1	

## \* FP Growth (Frequent Pattern Growth) algorithm

Trans id	Items	minsup = 2
T <sub>1</sub>	I <sub>1</sub> , I <sub>2</sub> , I <sub>5</sub>	Partition into multiple subsets
T <sub>2</sub>	I <sub>2</sub> , I <sub>4</sub>	for every subset we find the patterns
T <sub>3</sub>	I <sub>2</sub> , I <sub>3</sub>	computing sup. values again & merge
T <sub>4</sub>	I <sub>1</sub> , I <sub>2</sub> , I <sub>4</sub>	
T <sub>5</sub>	I <sub>1</sub> , I <sub>3</sub>	
T <sub>6</sub>	I <sub>2</sub> , I <sub>3</sub>	
T <sub>7</sub>	I <sub>1</sub> , I <sub>3</sub>	
T <sub>8</sub>	I <sub>1</sub> , I <sub>2</sub> , I <sub>3</sub> , I <sub>5</sub>	
T <sub>9</sub>	I <sub>1</sub> , I <sub>2</sub> , I <sub>3</sub>	

1. items	① Find freq. 1-itemsets	sup ct
{I <sub>1</sub> }	6	
{I <sub>2</sub> }	7	
{I <sub>3</sub> }	6	
{I <sub>4</sub> }	2	
{I <sub>5</sub> }	2	

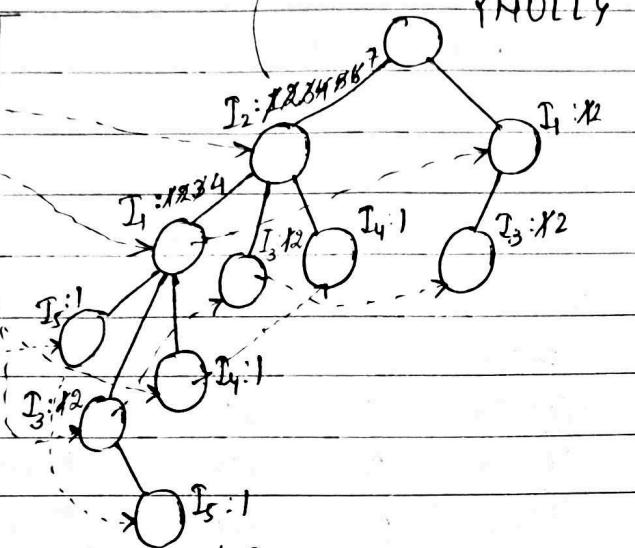
② sort in desc

L-order list
{I <sub>2</sub> }
{I <sub>3</sub> }
{I <sub>1</sub> }
{I <sub>4</sub> }
{I <sub>5</sub> }

If Qn asks freq. 2-itemsets then stop Apriori algo there

Every transac<sup>n</sup> in the dataset is considered as a branch in the FP tree.  $T_1 \rightarrow I_1, I_2, I_5 \rightarrow$  as per L-order  $\rightarrow I_2, I_1, I_5$   
 comes 1st time in this tree so its support is 1

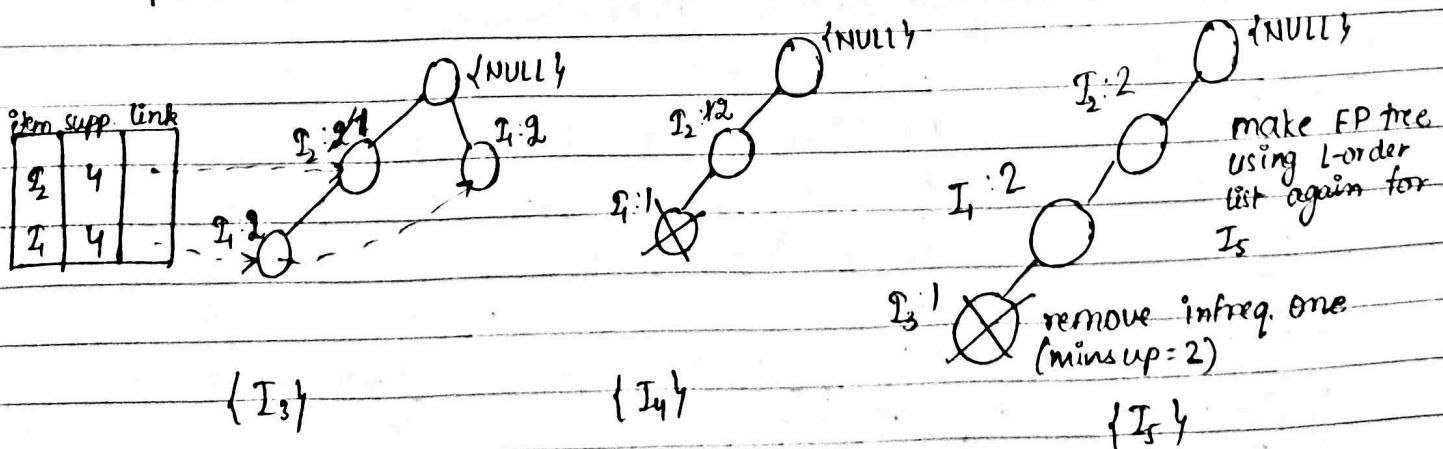
items	support	link
{I <sub>2</sub> }	7	-
{I <sub>1</sub> }	6	-
{I <sub>3</sub> }	6	-
{I <sub>4</sub> }	2	-
{I <sub>5</sub> }	2	-



how many paths to reach  $I_5$

item<sub>5</sub> support link  $\Rightarrow$  start w/ bottom of the list (min sup value)

items	conditional pattern base <sup>sup. value of <math>I_5</math></sup>	conditional FP-tree	Frequent Patterns
{I <sub>5</sub> }	{I <sub>2</sub> , I <sub>1</sub> : 1} {I <sub>2</sub> , I <sub>4</sub> , I <sub>3</sub> : 1}	I <sub>2</sub> : 2, I <sub>4</sub> : 2	{I <sub>2</sub> , I <sub>5</sub> : 2}
{I <sub>4</sub> }	{I <sub>2</sub> , I <sub>1</sub> : 1}	I <sub>2</sub> : 2	{I <sub>2</sub> , I <sub>4</sub> : 2}
* {I <sub>3</sub> }	{I <sub>2</sub> , I <sub>1</sub> , I <sub>3</sub> : 2} {I <sub>2</sub> , I <sub>1</sub> : 2} {I <sub>1</sub> : 2}	{I <sub>2</sub> : 4, I <sub>1</sub> : 2} {I <sub>1</sub> : 2}	{I <sub>2</sub> , I <sub>3</sub> : 4} {I <sub>2</sub> , I <sub>3</sub> : 4} (I <sub>3</sub> support 4 in the table) *{I <sub>2</sub> , I <sub>1</sub> , I <sub>3</sub> : 2} (min(I <sub>1</sub> , I <sub>2</sub> ) in the path w/ 4, 2 = 2)



$\{I_1\}$	$\{I_2\}$ corr. to $I_1$ support in the tree path	$\{I_2 : 4\}$	$\{I_2, I_1 : 4\}$
$\{I_2\}$			

Building FP tree ↓ complex than scanning whole dataset

### \* Strong Association Rule

$$N = 10,000 \text{ (Transac<sup>m</sup>')}$$

$$\text{Computer Games} = 6000$$

$$\text{Video} = 7500$$

$$\text{Computer Games and video} = 4000$$

Also rules are representation of freq. itemsets

$$\{X_1, X_2, X_3\}$$

$$X_1 \rightarrow X_2 X_3, X_2 \rightarrow X_1 X_3$$

Computer Games → Video support = 40% confidence = 66%

$$\text{support} = \frac{\text{Computer U Video}}{10000} = \frac{4000}{10000} = 40\% \quad \begin{cases} \text{satisfy both} \\ \therefore \text{this rule is useful} \end{cases}$$

$$\text{evaluate conf} = \frac{\text{Comp U Video}}{6000} = \frac{4000}{6000} = 66\% \quad \begin{cases} \text{satisfy both} \\ \therefore \text{this rule is useful} \end{cases}$$

Association ~~val~~ value

$$\text{Prob. of video} = \frac{7500}{10000} = 75\%$$

if cust. only purchase video prob  
is 75%. if comp + video, conf value  
is 66%. → some ↓ in % age. Negative  
asso<sup>n</sup> in this rule of video w/  
computer game

$$\frac{(A \cap B) / P(A \cup B)}{P(A)}$$

+ve asso → Strong assoc rule.

lift value  $lift(A, B)$

$$lift(A, B) = \frac{P(A \cup B)}{P(A)P(B)}$$

= 1 (then A, B are independent)  
= < 1 (negative association)  
=> 1 (positive)

If A & B is indep. then  $P(A) \times P(B) = P(A \cup B)$   
or  $A \cap B$ ?

Numerical measure of how much -ve or +ve association is there b/w 2 itemsets A & B

If  $lift(A, B) = 0.4 \rightarrow$  this amt of negative assoco req. to make it 1

Chi square test (hypothesis is considered that the 2 itemsets are indep.)  
it rejected → some asso is there

① Contingency table of observed values

		Computer Games	Comp. games	
		$e_{11}$	$e_{12}$	$e_{21}$
Video	Video	4000	3500	7500
	Computer Games	2000	500	2500
		6000	4000	10000

② Expected value

$$e_{11} = \frac{6000 \times 7500}{10000} = 4500$$

$$e_{21} = \frac{6000 \times 2500}{10000} = 1500$$

$$e_{12} = \frac{4000 \times 7500}{10000} = 3000$$

$$e_{22} = \frac{4000 \times 2500}{10000} = 1000$$

③  $\chi^2 = \frac{(observed - expected)^2}{Expected}$

= 507 (suppose)

signif =  $(2-1)(2-1)$  ( $2 \times 2$  contingency table)  
look in table



+ve, -ve asso find see how?  
can be asso? can't?

club diff items so that customer can buy them

Eg- DVD is purchased along w/ Computer games or video

$$N = 10k$$

$$\text{Comp. games} = 6000$$

$$\text{Video} = 7500$$

$$\text{Comp. games} \& \text{video} = 4000$$

$$\text{DVD} = 2500$$

$$\text{Comp. games} \& \text{DVD} = 2000$$

$$\text{Video} \& \text{DVD} = 3500$$

$$\text{Comp. games} \rightarrow \text{Video}$$

$$\text{Comp. games} \rightarrow \text{DVD}$$

$$\text{Video} \rightarrow \text{DVD}$$

$$\text{sup} = 40\% \quad \text{conf} = 66\%$$

$$\text{Ans- } ① \text{ sup} = \frac{\{ \text{C.G U DVD} \}}{N} = \frac{2000}{10000} = 20\%$$

$$\text{sup} = \frac{\{ \text{V U DVD} \}}{N} = \frac{3500}{10000} = 35\%$$

$$② \text{ conf} = \frac{\{ \text{C.G U DVD} \}}{6000} = \frac{2000}{6000} = 33\%$$

$$\text{conf} = \frac{\{ \text{C.G U V U DVD} \}}{7500} = \frac{3500}{7500} = 46.66\%$$

Both X satisfy

1 -ve 1 +ve  $\rightarrow$  go w/ +ve

Both +ve  $\rightarrow$  go w/ higher one

↳ calc. lift

If they satisfy sup & conf criteria  
lift value to check +ve / -ve assoc & choose b/w them

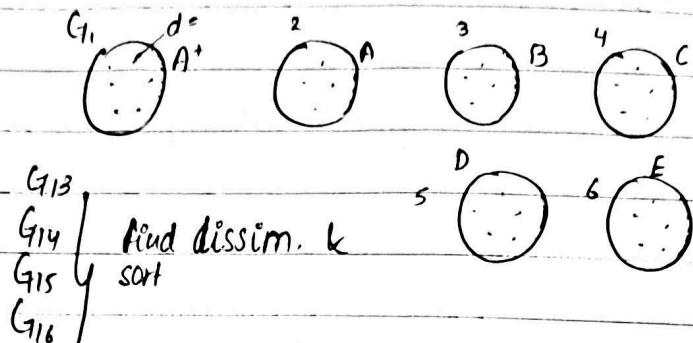
Dimensionality reduc<sup>n</sup> → Feature subset selec<sup>n</sup>  
Feature extraction (can use clustering)

$x_1 \ x_2 \ x_3$   $x_4 \ x_5 \ x_6$  Best group of  
 $x_1 \quad x_2$  features that  
can be clubbed

## \* Clustering

Manual labelling

X feasible w/ large dataset



Objective: Finding group of data object such that data objects in same group are very similar & those in other groups are very dissimilar

Unsupervised method

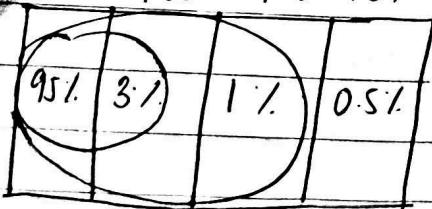
Clustering → automatic classification

Topic 1

D<sub>1</sub>  
D<sub>2</sub> D<sub>3</sub>

pal components can be grouped as clusters

PC1 PC2 PC3 PC4



Association analysis

	I <sub>1</sub>	I <sub>2</sub>	I <sub>3</sub>	I <sub>4</sub>	I <sub>5</sub>
1					
2					

I<sub>2</sub> I<sub>1</sub> I<sub>5</sub>

freq. itemsets

all instances having  
these items can be  
grouped

Compression → PCA lossy data compression

redundant attributes removed → lossless compression

KNN → Identify dist of 1 sample w/ other samples  
better handled by clustering

Outlier detect → sample very dissimilar w/ other groups

Intracluster var sum over all clusters  $\rightarrow$  should be minimal

Intracluster var min, Inter cluster max<sup>m</sup>

Eg- Network data analysis, Credit card fraud detect

M1:

(1)  
3  
(2)  
5  
(4)  
6

Intra cluster similarity = sum of similarity b/w  
elements of all clusters

M2:

(1)  
2  
(3)  
4  
(5)  
6

	I <sub>1</sub>	I <sub>2</sub>	I <sub>3</sub>	I <sub>4</sub>	I <sub>5</sub>
1					
2					
1000					

$1 \leq 3$  (3 clusters)

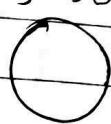
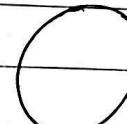
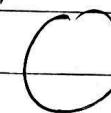
$$C_1 = 250$$

$$C_2 = 400$$

$$C_3 = 350$$

(no. of clust.) (no. of instances)

$C \ll N$



Partitioning approach:

How you partition data into various clusters

Hierarchical approach:  $\leftarrow$  Bottom up

Top down

1, 2, ..., 1000

Top Down

web page  $\rightarrow$  hyperlink on each page

1, 2, ..., 300

301, ..., 500 | 501, ..., 1000

11, ..., 75

76, ..., 300

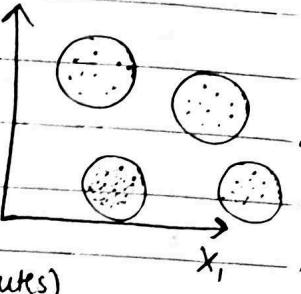
Bottom up

1 | 2 | 3 | ... | 1000

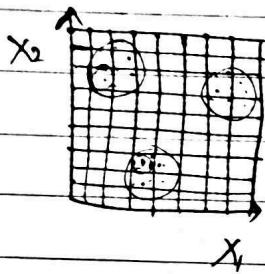
Similar  $\rightarrow$  make hierar

## Density-based approach

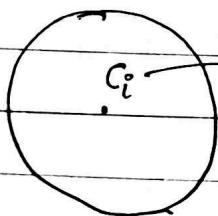
No. of data points within a region  $x_2$



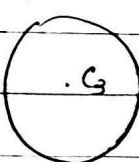
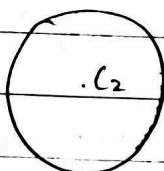
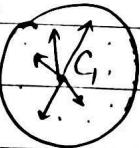
## Grid based approach



respective samples in cells  $\rightarrow$  aggregate func<sup>n</sup> defined



Representative point for the cluster



sum all squared distances (within cluster var)

Good clustering algo should have min value

① Randomly select 3 repres. samples

assume all are numeric samples

$X_1, X_2, X_3$  ② Distance of this pt wrt  $s_i \rightarrow$  select as per min distance the cluster it belongs to (Euclidean distance, Similarity)

1	-	-	-
2	-	-	-
3	-	-	-
...			
20	-	-	-

$k=3$

③ Update the representative sample  $\rightarrow$  Mean value  
Considered as rep. sample in next iteration



X Y  
Eg-  $A_1(2, 10)$

$A_2(2, 5)$

$A_3(8, 4)$

$B_1(5, 8)$

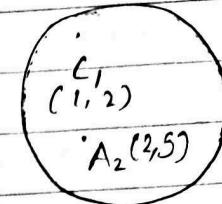
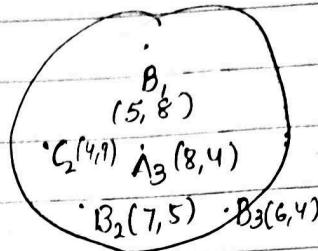
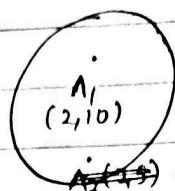
$B_2(7, 5)$

$B_3(6, 4)$

$G(1, 2)$

$C_2(4, 9)$

$K=3$

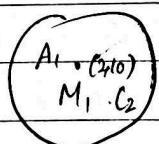


$$d(A_2, A_1) = \sqrt{5^2} = 5 \quad \text{Centroids}$$

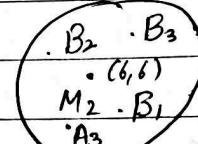
$$d(A_2, B_1) = \sqrt{9+9} = 3\sqrt{2} \quad M_1 = (6, 6) \quad M_3 = (1.5, 3.5)$$

$$d(A_2, C_1) = \sqrt{1+9} = \sqrt{10} \quad M_1 = (2, 10)$$

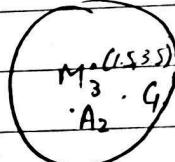
$$d(A_3, A_1) = \sqrt{36+36} = 6\sqrt{2}$$



$$d(A_3, B_1) = \sqrt{9+16} = 5$$



$$d(A_3, G) = \sqrt{49+4} = \sqrt{53}$$



$$d(B_2, A_1) = \sqrt{25+25} = 5\sqrt{2}$$

$$d(A_2, M_2) = \sqrt{16+1} = \sqrt{17}$$

$$d(B_2, B_1) = \sqrt{4+9} = \sqrt{13}$$

$$d(A_2, M_3) = \sqrt{0.5^2+1.5^2} = \sqrt{\frac{1}{4}+\frac{9}{4}} = \sqrt{5}$$

$$\therefore d(B_2, B_1) = \sqrt{4+9} = \sqrt{13}$$

$$d(A_3, M_1) = \sqrt{4+4} = \sqrt{8}$$

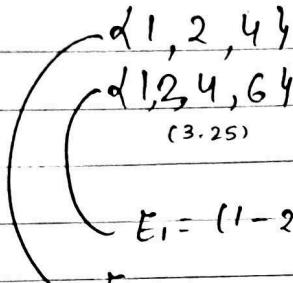
$$\text{Objective: } E = \sum_{i=1}^k \sum_{p \in C_i} (p - c_i)^2$$

↳ point

Ex-X: 1, 2, 4, 6, 7, 8, 12, 28

1 feature  $\rightarrow$  8 samples

$k=3$



{ 6, 7, 8 }  
 { 7, 8 } ↳  
 ( 7.5 )

{ 12, 28 }  
 { 12, 28 }

( 20 )

2 clusters  
sets

$$E_1 = (1-3.25)^2 + (2-3.25)^2 + (4-3.25)^2 + (6-7)^2 + (7-7)^2 + (8-7)^2 + (8-7)^2 =$$

$$E_2 =$$

$$C1 \rightarrow \text{Mean } \frac{7}{3} = 2.33$$

$$\frac{6+7+8}{3} = 7$$

$$\frac{12+28}{2} = 20$$

$$E_1 = 1.76 + 0.1089 + 1.67 + 1 + 1 + 128 = 133.53 \rightarrow \text{total error}$$

$$\begin{aligned} E_2 &= (3.25-1)^2 + (3.25-2)^2 + (4-3.25)^2 + (6-3.25)^2 + (7-7.5)^2 + (8-7.5)^2 + 128 \\ &= 5.0625 + 1.56 + 0.5625 + 7.5625 + 0.25 + 0.25 + 128 \\ &= 143.25 \end{aligned}$$

$E_1 < E_2$  so better

Drawbacks:

① K Means Clustering works fine for numerical data, not nominal attributes

both work abt same not

Mean X defined for nominal attribute  $\rightarrow$  variation is K-Mode clustering

② Very sensitive towards outliers (due to mean)