

11/01/2024

(1 hrs) Morning

## Data Mining

147

How to process the huge amount of data?

Beyond 15-20 attributes per 1000 records.

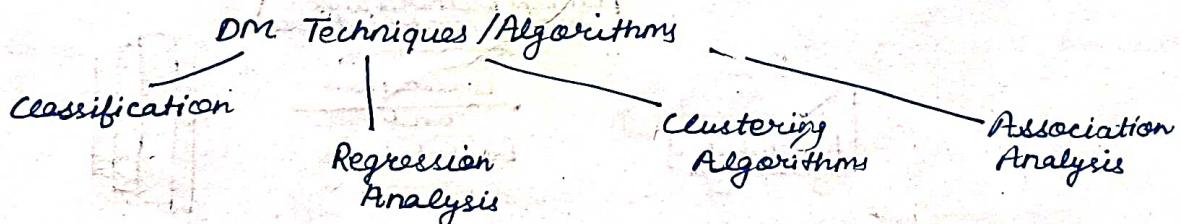
Although, electronics is developed to process huge data but it's meaningless if we don't have efficient algorithms that answer our query on such huge data.

ADAM

### Advanced Data Analysis Methods (ADAM)

It is the Data Mining, advanced course of DBMS.

Preprocessing - making data suitable for to be used in algorithms.



Data Presentation Tools - so Result more understandable to end-user.

Evaluation Parameters of Result

Either its appropriate or inappropriate.

Book -

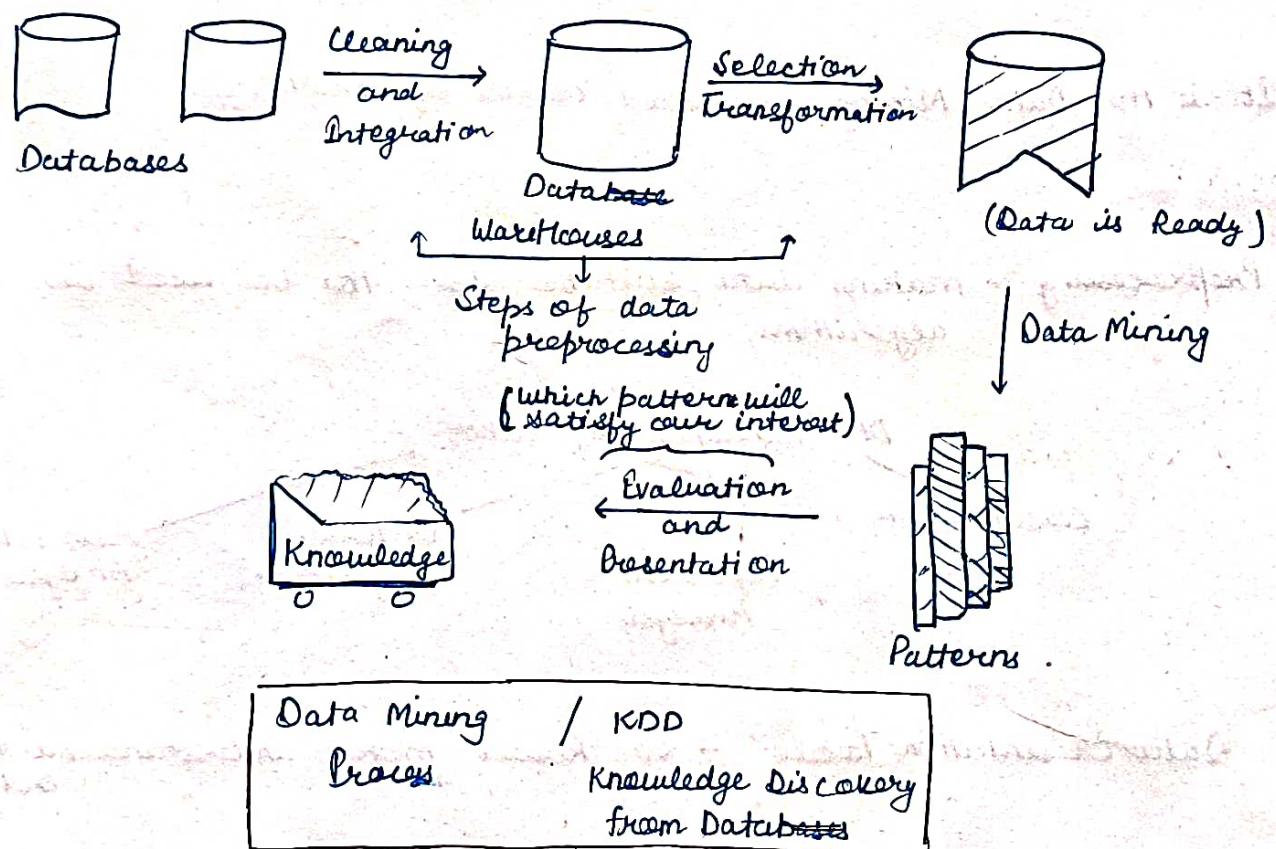
- Han & Kamber Data Mining Concepts and Techniques → Good concept & Technical language
- Vipin Kumar Introduction to Data Mining → Easy language

D B M S	Korth ↓ Easy	, Navath ↓ Technical
------------------	--------------------	-------------------------------

Data Mining is  
Process of finding interesting patterns from large amount  
of data.

We use data sets for storing huge data.

Representing Data Mining as a process, then it have states-



Data  
Very Initial Stage

Information  
At the intermediate stage

Knowledge  
At final stage,  
interesting patterns.

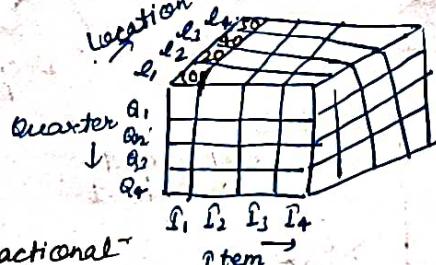
## Types of Data

- ① Relational Databases - Data stored in form of table Relational OLEP technology
- ② Data Ware House Data - " multi-dimensional format. querying operation  $\Rightarrow$  Roll up & Drill Down
- ③ Transactional Data - Record of what customers purchase in every visit.
- ④ Temporal Data / Sequence Data - Time is important constraint
- ⑤ Spatial Data - Position are important  
Temperature at different locations of world
- ⑥ Data Streams - Continuous data in various formats  
Temp. Sensors, Video Surveillance
- ⑦ Graph / Network Data - If feasible to draw graph that graph data we can get useful information.
- ⑧ Text Data
- ⑨ Multimedia Data
- ⑩ WWW

### Relational

Customer (name, cid, Age, DOB, itempurchase, income, itemp) purchase  
 Item (itemid, itemname, itemprice, itemq)  
 Sale (quarter, itemid, itemname, amountofSale)

### Data WareHouse



### Transactional

T <sub>1</sub>	Milk, bread, Coke
T <sub>2</sub>	Gurments, Milk, Bread
T <sub>3</sub>	Coke, Bread
T <sub>4</sub>	Coke

Data in compact form

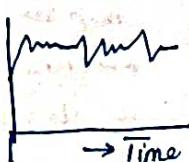
Pattern  
its some correlation or association among data

Reg. Milk with  
Bread is high

Market-Basket Analysis  $\rightarrow$  In a visit a customer purchase what things?

### Temporal

#### Stock Market



## Broad Categorisation of Data

Structured Data → That have a proper format  
eg. Relational data, Excel file

Unstructured Data → Not feasible to find the structure  
eg. Text data, PDF

Semi-structured Data → Contain both st. and unst. data.  
eg. Webpage with Text & Image

Unstructured → Structured

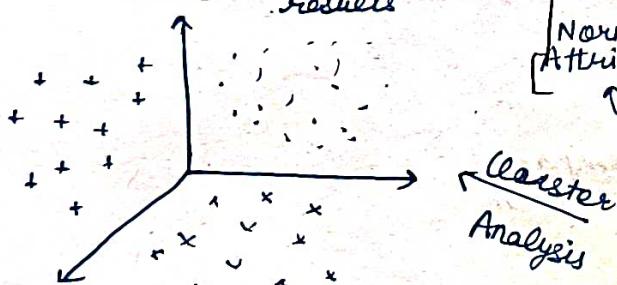
Get Result

Insert Value

They create model  
How do they  
do it?

## Data Mining Tasks

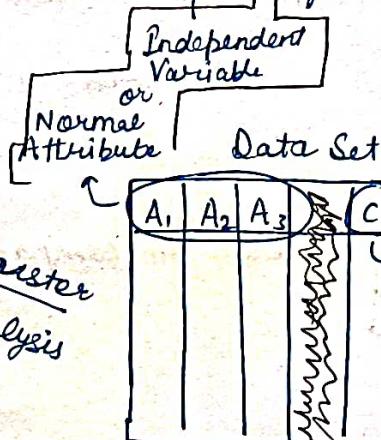
- ① Objective Predictive Tasks  
is to predict value of class attribute
- ② Descriptive Tasks  
Identify pattern in data  
then on basis of we summarize results



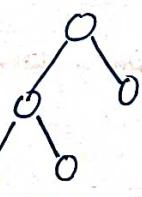
Data with similar structure/feature  
put in single cluster.

### Clustering Analysis

Represented by Random Variable  $X = \{A_1, A_2, A_3\}$



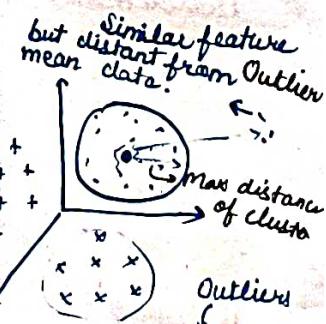
### Predictive Analysis



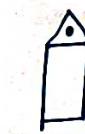
Decision-Tree Model

Class Attribute or Dependent Attribute

### Outlier Analysis



help to detect  
malicious activity  
as its different  
from regular  
data



→ Bread

Milk  
Helpful for  
inventory preparation,

### Association Analysis

# Data Mining Frequent Patterns Association and correlations

- Frequent Item sets
- Frequent Subsequences
- Frequent Substructures

Set of Items

Itemset = {Milk, Bread}

freq = 2

Frequent Itemset → Itemset has frequency greater than a threshold limit.

In Association Mining,  
we identify F.I.S. then apply algo.  
on them  
and get results.

e.g. - {Laptop, Dig. Camera}

high probability that he will also purchase 'memory card'.

If some substructure is frequent in graph, then it must be noted as its helpful.

e.g. - Diff. st. of Benzene

Classification, Regression Analysis

↓  
only applicable with discrete attribute value

Once you prepare data-set then algorithm predicts the class of new instance.

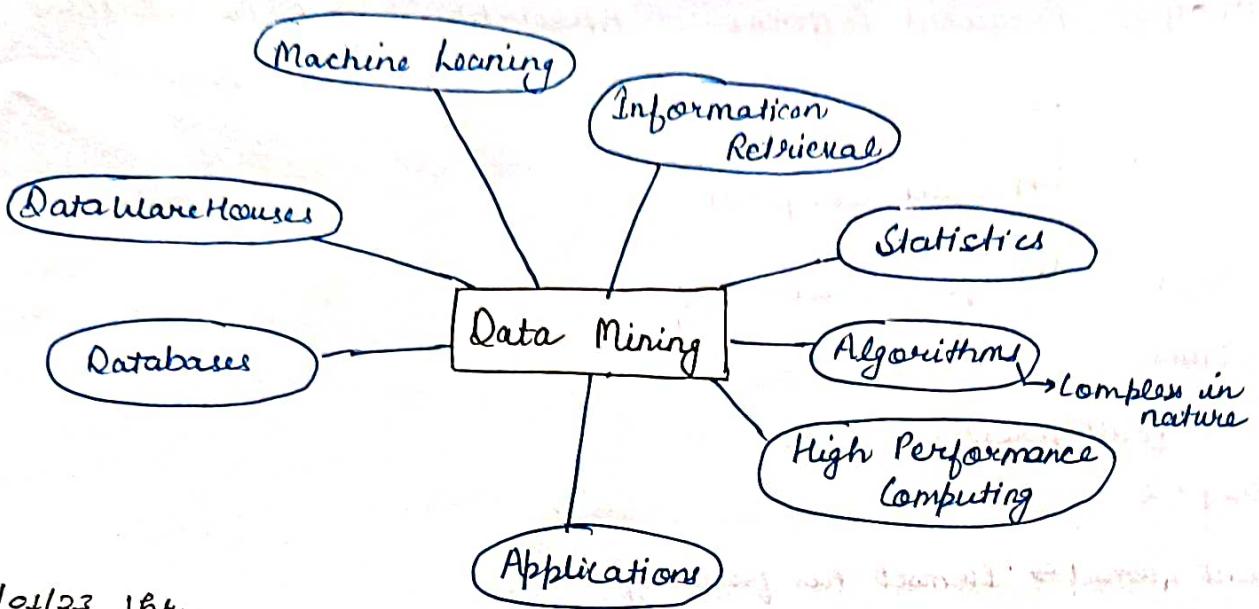
Prediction with continuous attribute value.

↓  
for Predictive Analysis

Data Set

A <sub>1</sub>	A <sub>2</sub>	A <sub>3</sub>	Result
-	-	-	Pass
-	-	-	Pass
-	-	-	Fail
-	-	-	Pass
-	-	-	Pass

→ 0.7 Classified in true groups  
 → 0.3  
 → 0.15  
 → 0.42  
 → 0.87  
 0-1 continuous range of values  
 So, Regression analysis algorithm needed.  
 If new student comes.  
 on basis of grades algo predicts which class it belongs.



18/01/23 1st  
(Morning)

### Attributes / Features / Variables / Dimensions

Data objects have some characteristics which are known as attributes.

Categorical Attributes (qualitative) somehow representing some categorical quality.		Nominal (=, ≠)	Employee id	Operations mode
		Ordinal (<, >)	T-Shirt Size (Small, Medium, Large) Grades (B+, A, A+)	median, mode
	Numeric (quantitative)	Interval (+, -)	app. diff. is defined Calender date Marks 30 45	mean median mode
		Ratio (*, /)	weight length 30kg 60kg	geometric mean harmonic mean

If a data can be represented in matrix (2-d form) or multidimensional form so that some mathematical operation can be defined on these.

45 is 15 more than 30  
30 is 15 less than 45

Depend on the nature of problem in which category does the attribute.

## Binary Attribute

→ Symmetric Binary

Gender  
Male  
Female.

→ Asymmetric Binary

Covid-19

+

-

-

-

-

## Discrete Attribute

### Continuous Attribute

{C, C+, B, B+, H, A+} Grade

Height (0-200 cm)

151.26

140.023

Binary discrete Continuous

qualitative (nominal or ordinal)

quantitative (ratio or interval.)

- ① Time is Terms of Month (Binary, qualitative, Ordinal)
- ② Brightness measured by a light meter. (Continuous- quantitative,
- ③ Brightness measured by People's Judgement (Discrete, qualitative, ordinal)
- ④ Bronze, Silver & Gold Medals.
- ⑤ No. of patients in a hospital
- ⑥ Number of Books
- ⑦ Military rank
- ⑧ Temp in °C or °F → Interval
- ⑨ Temp in Kelvin. → Ratio →

Ref. Point  
well-defined

# Data Similarity & Dissimilarity Measures

19/01/27  
Morning (2 hrs)

Data Matrix

Dissimilarity Matrix

No of  
Attr.

No of Data  
Item types

After

Putting/Representing data in matrix  
we can perform no. of operation  
and derive the result.

	A <sub>1</sub>	A <sub>2</sub>	A <sub>3</sub>	A <sub>4</sub>	...	A <sub>100</sub>
1						
2						
3						
4						
:						
50,000						

No Similarity  $\leftrightarrow$   
Similarity (0, 1)  $\rightarrow$  Both are  
similar

Dissimilarity = 1 - similarity

Id	(nominal) Test 1	(ordinal) Test 2	(numeric) Test 3
1	Code A	Excellent	45
2	Code B	Fair	22
3	Code C	Good	64
4	Code A	Excellent	24

Dissimilarity Measurement for

nominal attributes

$$dis(i,j) = \frac{P-m}{P}$$

P  $\rightarrow$  no. of attr. of nominal types

m  $\rightarrow$  no. of matches in their values

$$\begin{array}{ccccc} & A_1 & A_2 & A_3 & \\ i & \text{Code A} & \text{Code B} & \text{Code C} & \end{array} \quad \begin{array}{c} \frac{3-2}{3} = \frac{1}{3} \\ \hline j & \text{Code A} & \text{Code B} & \text{Code D} & \end{array}$$

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

$$\begin{bmatrix} 0 & & & \\ 1 & 0 & & \\ 1 & 1 & 0 & \\ 0 & 1 & 1 & 0 \end{bmatrix}$$

$$d(2,1) = \frac{1-0}{1} = 1$$

## Dissimilarity measurement for Binary attribute

	1	0
i	+	q r t
0	s	t

$$P = q + r + s + t$$

↳ total no. of attributes

$$i = (1, 1, 0, 1)$$

$$j = (0, 1, 1, 0)$$

$$r=2 \quad q=1 \quad s=1 \quad t=0$$

$$d(i,j) = \frac{q+r}{q+r+s+t}$$

$$\text{sim}(i,j) = \frac{q+t}{q+r+s+t}$$

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

Simple Matching Coefficient (SMC)

Symmetric Binary will be better use

For Asymmetric Binary Attribute

Given 19 test

$$\text{sim}(i,j) = \frac{q+r}{q+r+s+t}$$

Both can be used  
Jaccard coefficient  $\frac{q+r}{q+r+s+t}$  depends on reg.

$$i = (1, 0, 0, 0, 0, 0, 0, 0, 0, 0)$$

$$j = (0, 0, 0, 0, 0, 0, 1, 0, 0, 1)$$

C	P	N
P	N	
N	N	
N	N	
N	P	
N	N	N
N	P	

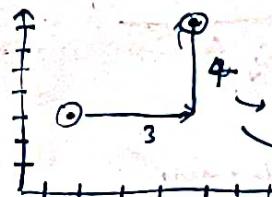
$$\left. \begin{aligned} SMC &= \frac{7}{10} \\ J &= \frac{0}{10} \end{aligned} \right\} \text{Numerical Similarity}$$

## Dissimilarity Measurement for Numeric Attribute

- Euclidean Distance → Shortest dist. b/w two points  $\Rightarrow$  Shortest Dissimilarity
- Manhattan Distance → Block dist./City distance
- Supremum Distance → Max. of both axis

$$i(2,3) \quad j(5,7)$$

$$d(i,j) = \sqrt{(5-2)^2 + (7-3)^2}$$



→ Total is Manhattan distance  
Max. of both is 4 which is Supremum distance.

Normalisation is performed to reduce the huge range & differences.

Pear Table data.

Max-Min = 42

$$d(i,j) = \begin{bmatrix} 0 & & & \\ 23/42 & 0 & & \\ 19/42 & 1 & 0 & \\ 21/42 & 2/42 & 40/42 & 0 \end{bmatrix}$$

Dissimilarity measure for ordinal attribute

Convert to numeric attribute

{fair, good, excellent}  
1 2 3

$$d(i,j) = \begin{bmatrix} 0 & & & \\ 2/2 = 1 & 0 & & \\ 0.5 & 0.5 & 0 & \\ 0 & 1 & 0.5 & 0 \end{bmatrix}$$

Formula

$$\frac{r_{ij} - 1}{m_f - 1}$$

More standard normalisation

$$\text{Fair} = \frac{1-1}{3-1} = \frac{0}{2} = 0$$

$$\text{Good} = \frac{2-1}{3-1} = 0.5$$

$$\text{Excellent} = \frac{3-1}{3-1} = 1$$

How to integrate result -

$$d(i,j) = \frac{\sum_{f=1}^F \epsilon_{ij}^{(f)} d_{ij}^{(f)}}{\sum_{f=1}^F \epsilon_{ij}^{(f)}}$$

$\epsilon_{ij}$  = weights corresponding to dissimilarity

$$= \frac{1 \cdot d(2,1) + 1 \cdot d(2,1)^2 + d(2,1)^3}{3}$$

-0 → missing or zero value

≠ 1 → otherwise

## Cosine Similarity

Finds the similarity between two documents.

$$(x, y) = \frac{x \cdot y}{\|x\| \|y\|}$$

→ Very popular used in text mining.

Document Matrix -  $\begin{matrix} \theta = 0^\circ \\ \theta = 90^\circ \end{matrix}$

or

Term Frequency

Matrix

↓

(Basis of Text Mining)

	(not)	(tea)	(bread)	(gas)	...	Term 1000
Term 1	5	3	0	-	-	2
Term 2	20	6	0	-	-	12
Term 3	3	3	5	-	-	6
...	1	1	1	-	-	1
Term 1000	2	6	2	-	-	1

## Correlation Coefficient

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

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

$$\text{std}(x) =$$

24<sup>th</sup> Jan '24  
Morning (2hr)

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

+ve Corr.  $\rightarrow 1$

0 Corr.

-ve Corr.

$\rightarrow -1$

Neutral Correlation

## Data Preprocessing Techniques

- ① Data Cleaning → inconsistency, missing values, incomplete data. After cleaning, we get req. data.
- ② Data Integration → diff. sources Relational format, Non-Rel. format → Integrated in standard format.
- ③ Data Reduction → (Below)
- ④ Data Transformation → Normalise our values to get them in standard format.

Rel. Format

A <sub>1</sub>	A <sub>2</sub>	...	A <sub>100</sub>
1	2	...	100

Reduce to no. of attribute  
or and

Reduce the no. of entries

Numeric Binary

A <sub>1</sub>	A <sub>2</sub>	...	A <sub>100</sub>
1	2	...	100

⇒ No loss in data

## ① Data Cleaning

### ① Missing Values

A <sub>1</sub>	A <sub>2</sub>	A <sub>3</sub>	A <sub>4</sub>	A <sub>5</sub>	A <sub>6</sub>
10	0	25	-	-	-
20	1	26	-	-	-
30	0	-	-	-	-
40	1	25	-	-	-
50	1	27	-	-	-

- a. Ignorance Approach -

ex 2-3 missing value out of 5,000 records  
so ignore those tuples

sufficiently large

- b. Putting Random Value -

- May <sup>create</sup> Noisy, Outlier data

- c. Put the mean value

- Either global mean

- Either mean of adjacent neighbours

- d. Either selecting global constants like 0 etc.

- e. Prediction of value. {Much better, but also have some complexity}

- Decision tree method.

Using available values we create decision tree  
then the output gives us the value (unknown).

- f. Linear Regression Method

## ② Handling Noisy Data

Random errors in data is noise.

### a. Data Smoothing Method.

Smooth out the data

i) Binning method / Basket method (Data in sorted order)

25, 26, 29, 35, 40, 42, 45, 50, 55, 100

Baskets

1: 25, 26, 29

Mean  $x_1$  Replace all value with mean value

2: 35, 40, 42

$x_2$   $x_2$   $x_2$

3: 45, 50, 55

$x_3$   $x_3$   $x_3$

A <sub>3</sub>
25
26
100
25
27

Beyond our expectation

Noise entry are very less.

Mean value we get for every group is near to filter value.

### (ii) Method of Regression.

## ② Data Integration

### ① Entity Identification Problem-

Customer

Id			

Sales

Item	Cost	Id	

Rename those attr. in common names.

### ② Redundancy and Correlation-

is not duplicacy

↳ Repetition

Same tuple many time

Some attribute may be redundant.

ex

Age	DOB

is redundant

as using current date

& d.o.b. we can calc.

use.

During integration, no redundancy in data.

(b) Using the correlation of attribute

If corr. coeff. very high then one attr. can be substituted by other.

Chi-square test - ( $\chi^2$ )

For nominal attr.

A  $\rightarrow$  C distinct values  
( $a_1, a_2, \dots, a_C$ )

B  $\rightarrow$  r distinct values  
( $b_1, b_2, \dots, b_r$ )

$(A_i, B_j)$   $\rightarrow$  Joint event that attribute A take on value  $a_i$  and attribute B take on value  $b_j$   
(e.g.  $(A=a_1, B=b_1)$ )

$$\chi^2 = \sum_{i=1}^C \sum_{j=1}^r \frac{(O_{ij} - E_{ij})^2}{E_{ij}}$$

$O_{ij}$  = Observed freq. of the joint event  $(A_i, B_j)$

$E_{ij}$  = Expected freq. of  $(A_i, B_j)$

n = No. of data tuples

$$E_{ij} = \frac{\text{count}(A=A_i) \times \text{count}(B=B_j)}{n}$$

		Gender		
		Male	Female	
Reading Habits	Fiction	300(250)	200(200)	500
	Non-fiction	100(90)	400(380)	500
		400	600	1000

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

$$\text{cov}(A, B) = E(A - \bar{A}) E(B - \bar{B})$$

$$\text{cov}(A, B) = E(A - \bar{A})(B - \bar{B})$$

$$E(A) = \bar{A} = \frac{\sum_{i=1}^n a_i}{n}$$

$$E(B) = \bar{B} = \frac{\sum_{i=1}^n b_i}{n}$$

$$\text{corr. coeff.} = \frac{\text{covar.}(A, B)}{\text{std.}(A) \text{ std.}(B)}$$

If cov. is true, if value of A changes above mean then its more likely that B will also change. In case of the cov.

If A is above mean, B is below mean -ve covariance.

### ③ Data Reduction

① Feature Selection or Feature Subset Method

② Feature Reduction



Merging



Combine similar attr. & merge them

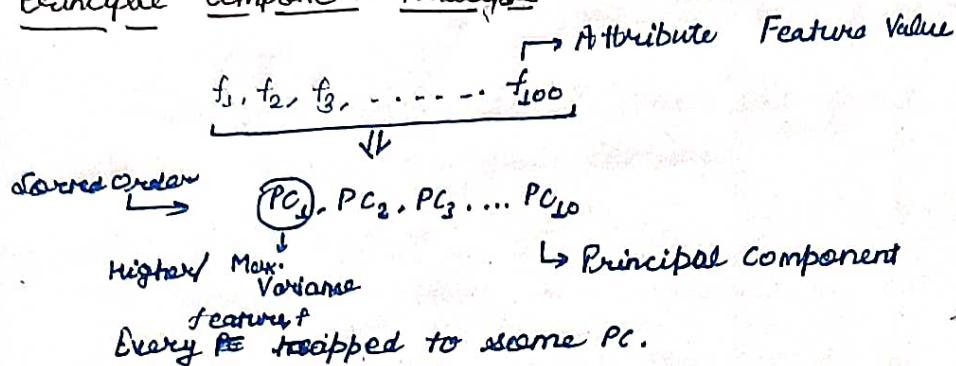
→ Selecting some important attr. we prefer the Ranking Method.

Principal Component Analysis (PCA).

$$f(A') = \frac{w_1 A_1 + w_2 A_2 + w_3 A_3}{w_1 + w_2 + w_3}$$

Data Reduction Methods 31<sup>st</sup> Jan'21  
(2 hr) Morning

### Principal Component Analysis



Data Matrix =  $\begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_{100} \end{bmatrix}$

$f_1, f_2, \dots, f_{100}$

Arrang PC acc. to their Data Variance level.

Variance is less  $\Rightarrow$  Not contributing in the classification

Data Matrix.

eg.  $x_i \begin{bmatrix} f_1 & f_2 & f_3 & f_4 & f_5 \\ 5 & 6 & 7 & 9 & 2 \\ 2 & 7 & 2 & 7 & 9 \\ 7 & 7 & 1 & 7 & 3 \\ 4 & 7 & 2 & 7 & 1 \\ 1 & 7 & 5 & 7 & 6 \end{bmatrix}$

\* How to map feature to P.C.?

Eigen vectors and Eigen values

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

Symn.

Mat.

$$\text{Variance} = \frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})^2$$

Covariance ( $f_1, f_2$ )

$$= \frac{1}{n} \sum_{i=1}^n (f_i - \bar{f}_1)(f_i - \bar{f}_2)$$

Prepare Covariance Matrix from Data Matrix

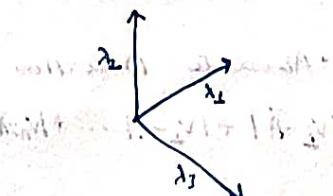
Inc one, other inc  $\Rightarrow$  pos. cov.

Inc one, other dec  $\Rightarrow$  -ve cov.

$$\text{Cov. Matrix} = \begin{bmatrix} \sum f_{11} & \sum f_{12} & \sum f_{13} \\ \sum f_{21} & \sum f_{22} & \sum f_{23} \\ \sum f_{31} & \sum f_{32} & \sum f_{33} \end{bmatrix}$$

Map our feature  
↓  
Eigenvector.

Highest eigen value  
 $\lambda$   
Highest PC.  
Max. Variance.



$\lambda \rightarrow$  represent amount of variance in that axis

$$D = \begin{bmatrix} f_1 & f_2 \\ 2 & 3 \\ 4 & 4 \\ 3 & 6 \\ 5 & 4 \\ 6 & 2 \\ 6 & 3 \\ 8 & 4 \\ 6 & 6 \end{bmatrix}$$

Cov. Matrix =

## Normalization -

- Min-Max Normalization - Preserve relationship in your prev. data.

$$V_i' = \frac{V_i - min_A}{max_A - min_A} (new\_max_A - new\_min_A) + new\_min_A$$

Referensi: [www.mechanics.tufts.edu/~mleach/EEG/EEG.html](http://www.mechanics.tufts.edu/~mleach/EEG/EEG.html)

eq. 10  
28  
42  
76  
91 → 25

- ## • Z-score Normalization

$$V_C' = \frac{V_i - \bar{A}}{\sigma_A} \rightarrow \text{standard deviation}$$

Is applicable when the real boundaries are not known.

There may be noise & outlier data.

- ## • Mean Absolute Deviation

$$V_C' = \frac{V_C - A}{S_A} \rightarrow \text{Absolute Deviation}$$

$$S_n = \frac{1}{n} (|V_1 - \bar{A}| + |V_2 - \bar{A}| + \dots + |V_n - \bar{A}|)$$

## Normalization

by decimal scaling.

- Normalization by decimal scaling

$$V_i' = \frac{V_i}{10^7}$$

## Classification 2nd Feb' 24 Morning (L14)

Identify / Classify some group of entity or the classify entity into groups is known as classification problem.

We have some data, from that we create model which is used to classify the new instance.

These models are algorithms.

Requirement in data is it should be labelled. In General class attribute tells about the class of each instance.

Two phases -

- Training Phase → with given data → Prepare Model
- Testing Phase → with new inputs → we do prediction.

Model

$$y = f(x)$$

Outcome Input

You can use different models. i.e. which one is best.  
But the accuracy should be very high.  
for good prediction.

\* Value of class attribute should be discrete.

Supervised - The class of all data are known.  
we know either we are correct or not.

Unsupervised - Class of data instances are unknown.

Training with unlabelled data.

Semi-supervised - Mix of both type data

## Decision Tree

Customer Name	Age	G	Braome	Defender
Ram	56	M	2000	Yes
Syam	53	M	3000	Yes
Sita	52	F	2500	No
Shiv	42	M	3200	No
Abhi	47	M	2700	Yes

Labelled Data

Discrete Class Attribute

Generate Decision Tree

Input: Data Partition D

attribute-list, attribute-selection method

Output: Decision Tree.

Method: ① Create a Node N

② If tuples in N are all of the same class C then

③ return N as a leaf node labelled with the class C;

④ If attribute-list is empty then

⑤ return N as a leaf node labelled with the majority class in D

⑥ Apply attribute-selection-method (D, attribute-list) to find the best splitting criterion;

⑦ Label Node N with splitting criterion

⑧ if splitting attribute is discrete-value and multiway split allowed then

⑨ attribute-list  $\leftarrow$  attribute-list  $\leftarrow$  splitting-attribute

⑩ for each outcome j of splitting criterion // Partition the tuples and grow subtrees for each partition

⑪ let  $D_j$  be the set of data tuples in D satisfying outcome j

⑫ if  $D_j$  is empty then

⑬ attach a leaf node labelled with the majority class in  $D_j$  to node N

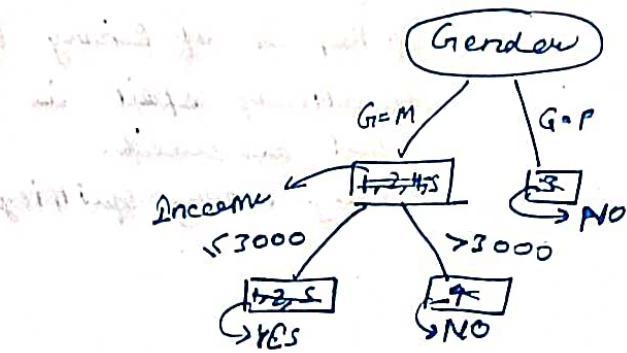
⑭ else attach the node returned by generate-decision-tree ( $D_j$ , attribute-list) to node N

⑮ end for

⑯ Return N

## Characteristics

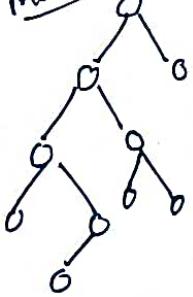
- Based on Greedy Algorithm
- Recursive in nature
- Divide and Conquer Approach



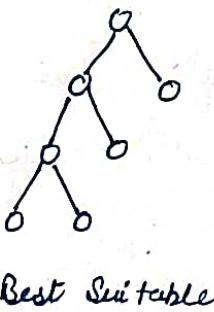
## 2 problems

- Overfitting
- Underfitting

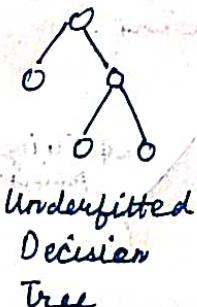
07 Feb '24  
Morning (2nd)



Overfitted Decision Tree



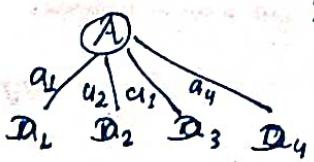
Best Suitable



Underfitted Decision Tree

So, we have to make a tradeoff.

## Discrete Attribute



If multiway split is allowed.

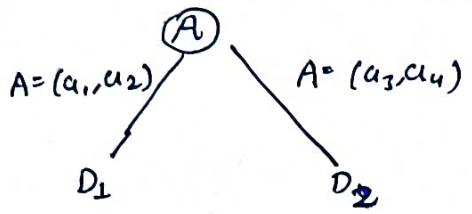
## Ordinal



If multiway split is allowed.

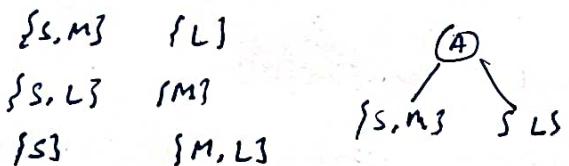
## Discrete Attribute and Binary Tree

$(a_1, a_2)$        $(a_3, a_4)$



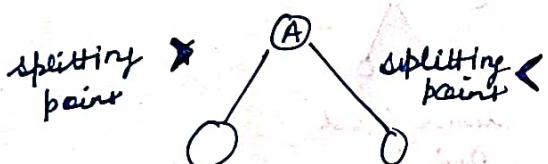
If Req. is of Binary Tree,  
multiway split is  
not allowed.  
Only binary splitting

In case of ordined attr.



## Continuous Attribute

Splitting point



## Attribute Selection Measure

- ① Information Gain → Based on concept of Information Theory
- ② Gain Ratio
- ③ Gain Index

Measure  
the amount  
of information

Min value  
of gain  
i.e. error.

↓  
Entropy

Measurement of information  
in some words or sentences

$$= - \sum_{i=1}^n p_i \log(p_i)$$

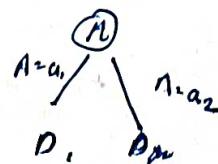
Direct Sigh

$$\begin{array}{ll} D = 1 & 1/1 \\ i = 2 & 2/1 \\ n = 2 & 2/1 \\ C = 1 & 1/1 \\ S = 2 & 2/1 \\ g = 1 & 1/1 \\ h = 2 & 2/1 \\ \hline \end{array}$$

We put in  
formulae  
get some  
result.

Info  
Gain

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



$$\text{Info}_A(D_i) = \frac{|D_i|}{|D|} \text{ Info}(D)$$

Our objective is to select attribute  $A$  that gives max. info. gain.

A <sub>1</sub>	A <sub>2</sub>	A <sub>3</sub>	C
			YES
			NO
		YES	YES
		YES	NO
		NO	NO

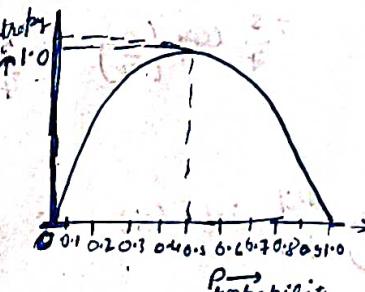
G:3  
C<sub>1</sub>:3

G:0  
C<sub>1</sub>:6

G:1  
C<sub>1</sub>:5

$$C_0 = -\frac{1}{2} \log\left(\frac{1}{2}\right) - \frac{1}{2} \log\left(\frac{1}{2}\right)$$

$$= \frac{1}{2} \leftarrow \text{Maximum Value of entropy}$$



If we write in decision tree

Info = 0

$$\text{Info Gain} = \frac{1}{2} - 0$$

$\approx 0.5 \leftarrow \text{Max. Gain}$

G:3  
G<sub>1</sub>:3

G:0      1:6

G:3  
G<sub>1</sub>:3

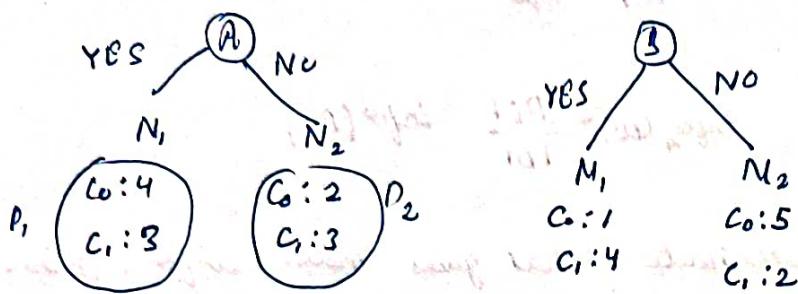
G:1      C<sub>1</sub>:5

$$\text{Info} = -\frac{1}{6} \log\left(\frac{1}{6}\right) - \frac{5}{6} \log\left(\frac{5}{6}\right)$$

$$= \log 6 - \frac{5}{6} \log 5$$

$C_0 : 6$
$C_1 : 6$

$$\text{info}(D) = L$$



$$P_1 = -\frac{4}{7} \log\left(\frac{4}{7}\right) - \frac{3}{7} \log\left(\frac{3}{7}\right) = 0.985$$

$$P_2 = -\frac{2}{5} \log\left(\frac{2}{5}\right) - \frac{3}{5} \log\left(\frac{3}{5}\right) = 0.970$$

$$P_3 = -\frac{1}{5} \log\left(\frac{1}{5}\right) - \frac{4}{5} \log\left(\frac{4}{5}\right) = 0.721$$

$$P_4 = -\frac{5}{7} \log\left(\frac{5}{7}\right) - \frac{2}{7} \log\left(\frac{2}{7}\right) = 0.863$$

$$(P) \frac{|D_i|}{|D|} \text{info}(D_i) = \frac{7}{12} \text{info}(P_1) + \frac{5}{12} \text{info}(P_2)$$

$$= 0.97875$$

$$\hookrightarrow \text{Infor Gain} = 1 - 0.97875 = 0.02125$$

$$(3) \frac{|D_i|}{|D|} \text{infor}(D_i) = \frac{5}{12} \cdot 0.721 + \frac{7}{12} \cdot 0.863$$

$$= 0.8038$$

$$\hookrightarrow \text{Infor Gain} = 1 - 0.8038 = 0.1962$$

Greater.

So, better attribute to be chosen is B.

## Income

75, 70, 85, 60, 120, 100, 95, 100, 90, 125, 220

## Score

NO	NO	NO	n's	YES	YES	NO	NO	NO	NO
60	70	75	85	90	95	100	120	125	220
✓	✓	✓	✓	✓	✓	✓	✓	✓	✓
55	65	72	80	87	92	97	110	122	172
=									

$$\begin{array}{c} 155 \\ \hline 6:10 & 6:3 \\ 6:0 & 6:7 \end{array}$$

$$\begin{array}{c} 72 \\ \hline 0 & 3 \\ 2 & 5 \\ \hline 3 & 4 \end{array}$$

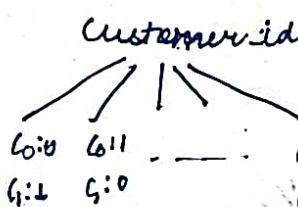
$$\begin{array}{c} 565 \\ \hline 0 & 3 \\ 1 & 0 \end{array}$$

$$\begin{array}{c} 80 \\ \hline 0 & 3 \\ 3 & 4 \end{array}$$

9<sup>th</sup> Feb '24  
Morning (1-4 hrs)

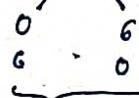
## Gain Ratio -

$$\text{Information Gain} = \text{info}(D) - \sum_{i=1}^n \text{info}_A(D_i)$$



If multiway split allowed.

A 3,3



Pure  
Classification

If pure class.,  
then impurity is  
least.

Least Entropy

Least Randomness.

We can divide on basis of distinct  
values the attribute have.  
instead of classifying in terms of att..

$$\text{Gain Ratio} = \frac{\text{Gain}}{\text{Split Info}}$$

sum of weights of each partition

$$= - \sum_{j=1}^v \frac{|D_j|}{|D|} \log_2 \left( \frac{|D_j|}{|D|} \right)$$

computed for all attribute.

whose is more, the more appropriate it is.

### Gini Index -

(Misclassification  
Error)

$$\text{Gini}(D) = 1 - \sum_{i=1}^m p_i^2$$

$$p_i = \frac{|C_i, D|}{|D|}$$

→ No. of tuples in your dataset

Q.

	C <sub>0</sub> : 3	C <sub>1</sub> : 3
① /	C <sub>0</sub> : 0	C <sub>1</sub> : 3
	C <sub>0</sub> : 6	C <sub>1</sub> : 0
	C <sub>0</sub> : 2	C <sub>1</sub> : 1

Pure classification

Compute Misclassification error.

$$P_1 = 1$$

$$P_1 = \frac{2}{6}, P_2 = \frac{3}{6}, P_3 = \frac{1}{6}$$

$$\text{Gini}(D) = 0$$

$$\text{Gini} = \frac{24}{27}$$

## Errors in Decision Tree

16<sup>th</sup> Feb  
Morning (- 1 hrs)

Training Error. → Error in design of model.

Test Error. → How much error the model gives on unseen records?

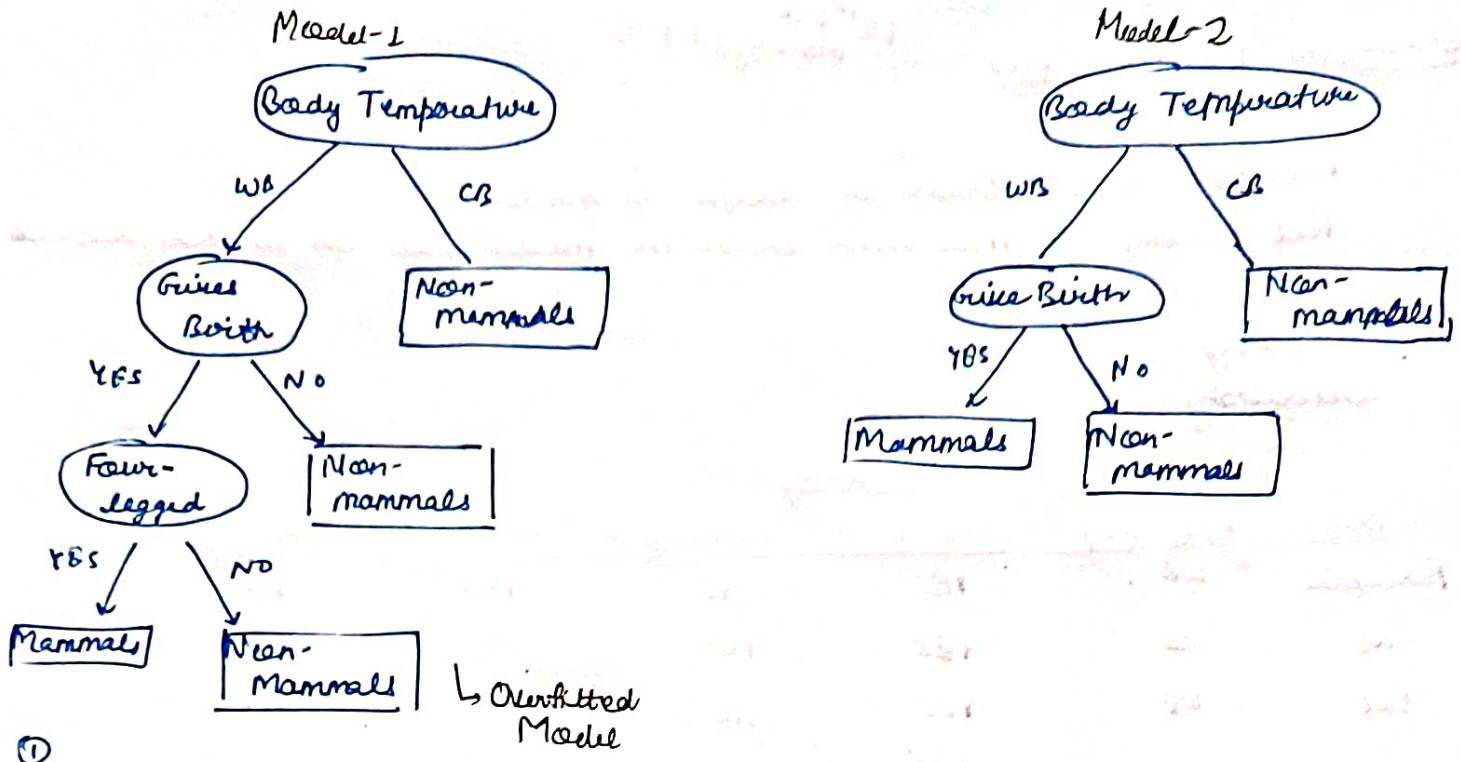
Overfitting  
Underfitting

### Training

	Name	Body Temp.	Gives Birth	Four-legged	Hibernates	Class Label
1	Porcupine	WB	YES	YES	YES	YES
2	Cat	WB	YES	YES	NO	YES
3	Bat	WB	YES	NO	YES	NO
4	Whale	WB	YES	NO	NO	NO
5	Salamander	CB	NO	YES	YES	NO
6	Komodo Dragon	CB	NO	YES	NO	NO
7	Python	CB	NO	NO	YES	NO
8	Salmon	CB	NO	NO	NO	NO
9	Eagle	WB	NO	NO	NO	NO
10	Guppy	CB	YES	NO	NO	NO

### Testing

	Name	Body Temp	Gives Birth	Four-legged	Hibernates	Class Label
1	Human	WB	YES	NO	NO	YES
2	Pigeon	WB	NO	NO	NO	NO
3	Elephant	WB	YES	YES	NO	YES
4	Leopard Shark	CB	YES	NO	NO	NO
5	Turtle	CB	NO	YES	NO	NO
6	Penguin	GB	NO	NO	NO	NO
7	Bee	GB	NO	NO	NO	NO
8	Dolphin	UCB	YES	NO	NO	YES
9	Spring Anteater	WB	NO	YES	YES	YES
10	Gila Monster	CB	NO	YES	YES	NO



①

1, 2, 3, 4, 5, 6, 7, 8, 9, 10

$$\text{Train error} = \underline{\underline{0}}$$

②

1, 2, 5, 6, 7, 8, 9, 10

$$\text{Train error} = \frac{3, 4}{10} = \underline{\underline{30\%}}$$

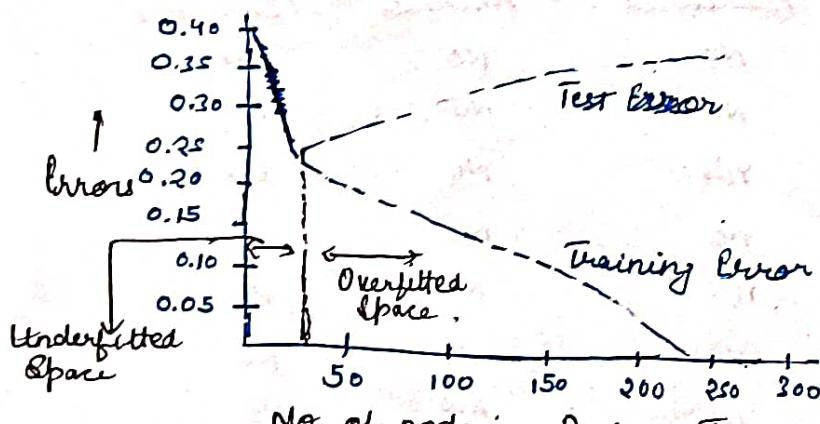
① 2, 3, 4, 5, 6, 7, 10

$$\text{Test error} = \frac{8, 9}{10} = \underline{\underline{30\%}}$$

②

1, 2, 5, 6, 7, 8, 9, 10

$$\text{Test error} = \frac{9}{10} = \underline{\underline{10\%}}$$



No. of node in Decision Tree →

If input size is very small, and our model is not properly trained, so high training error leads to testing error.

Underfitting

No. of nodes in Decision Tree ↗  
↓

Complexity of model inc.  
& cost also inc.

So, balance b/w