

FORMULA LIST

$$* \text{ Mean } (\bar{x}) = \frac{1}{n} \sum_{i=1}^n x_i$$

* Empirical Formula for unimodal numeric data -
 Mean - mode = 3 (mean - median)

$$* \text{ Midrange} = (\text{Max value} + \text{Min value}) / 2$$

* Measuring dispersion of data :-

↳ Inter Quartile Range (IQR) = Q₃ - Q₁

↳ Five Number Summary = (min, Q_1 , median, Q_3 , max)

↳ To calculate outliers:

$$Q_1 - (1.5 \times IQR) \quad \text{and} \quad Q_3 + (1.5 \times IQR)$$

$$\hookrightarrow \text{Variance} (\sigma^2) = \frac{\sum_{i=1}^n (x_i - \mu)^2}{n}$$

$$\hookrightarrow \text{Standard Deviation (SD)} = \sqrt{\text{Variance}} = \sqrt{\frac{\sum (x - \bar{x})^2}{n}}$$

* Measuring data similarity and dissimilarity :

↳ Similarity - Value is higher when objects are more alike often falls in the range (0,1).

↳ Dissimilarity - (Need to calculate the distance between objects). Lower value when objects are more alike.

For similar objects value = 0

↳ Proximity is nothing but similarity or dissimilarity between objects.

↳ Dissimilarity data matrix :-

0				
$d(2,1)$	0			
$d(3,1)$	$d(3,2)$	0		
:	:	:		
$d(n,1)$	$d(n,2)$...	0	

↳ Proximity measure for normal attributes :-

$$d(i,j) = \frac{p-m}{p} \quad \left\{ \begin{array}{l} p = \text{Total no. of attributes} \\ m = \text{no. of matching attribute values} \end{array} \right.$$

↳ Proximity measure for binary attributes :-

→ For symmetric binary variables :-

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

		Object j		sum
		1	0	
Object i	1	r	s	q+r
	0	s	t	q+s
	sum	r+t	s+t	p

$$q=11, r=10, s=01, t=00$$

⇒ For asymmetric Binary variables -

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

↳ Proximity measure of numeric data -

↳ Euclidean Distance -

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

↳ Manhattan Distance -

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

↳ Cosine Similarity -

$$\cos(\pi, y) = \frac{\pi \cdot y}{\|\pi\| * \|y\|}$$

$\pi \cdot y$: Vector dot product
 $\|\pi\|$: Magnitude of vector

NOTE: If the cosine value is closer to 1, the objects are more similar.

* HANDLE NOISY DATA -

- ↳ Binning Method : (1) Sort the data & partition into bins
 (2) Smooth data by bin means, bin median, bin boundaries, etc.

2) Equal-Depth (Freq) : Partition the data in such a way that each bin has same no. of data

3) Equal-width : Divide the data such that each bin has some width.

To calculate width of a bin :

$$\text{width} = \frac{\text{Max value} - \text{Min value}}{\text{No. of bins}}$$

4) Clustering : Partition the data along the 2 biggest gaps in the data.

* NORMALIZATION -

1) Min-Max Normalization :

$$v' = \frac{v - \text{min}_A}{\text{max}_A - \text{min}_A} (\text{new_max}_A - \text{new_min}_B) + \text{new_min}_B$$

v = value given to normality

min_A & max_A : Min & Max value of given attribute range

new_max_A & new_min_B : New values given for the attribute range.

2) Z-score Normalization :

$$v' = \frac{v - \mu_A}{\sigma_A} \quad \left[\begin{array}{l} \mu_A : \text{Mean of attribute} \\ \sigma_A : \text{S.D. of attribute} \end{array} \right]$$

3) Decimal Scaling :

$$v' = \frac{v}{10^j} \quad [j = \text{calculated as per the max no. of digits of the range}]$$

CHAPTER 3 : FREQUENT MINING PATTERNS

* Applications -

- ↳ Market Basket Analysis (Association Analysis)
- ↳ Telecommunication
- ↳ Credit Cards / Banking Services
- ↳ Medical Treatments
- ↳ Basketball - Game Analysis.

* ASSOCIATION RULE PROBLEM :-

- ① Given a database of transaction.

Transaction	Items
t ₁	Bread, Jelly, Peanut Butter
t ₂	B, PB
t ₃	B, Milk, PB
t ₄	Beer, B
t ₅	Beer, M

Q# Find all the association rules :

$n = 4$	S	a
$B \Rightarrow PB$	60%	75%
$PB \Rightarrow B$	60%	100%
$Beer \Rightarrow B$	20%	50%
$PB \Rightarrow Jelly$	20%	33.3%
$J \Rightarrow PB$	20%	100%
$J \Rightarrow M$	0%	0%

* ASSOCIATION RULE DEFINITIONS

* Support : $(\text{Item A} + \text{Item B}) / \text{Entire Dataset}$

* Confidence : $(\text{Item A} + \text{Item B}) / \text{Item A}$

* APRIORI ALGORITHM

① Eg: A dataset D has 4 transactions. Let min support be 50%. & min confidence is 80%.

Find: All frequent item set & also generate association rule using Apriori Algorithm.

$$\begin{aligned}\text{Min Support count} &= \text{min supp threshold} * \text{Total no. of trans} \\ &= 50/100 * 4 \\ &= 2\end{aligned}$$

Database D

TID	Items	Itemset	Supp.
100	1 3 4	$\{1\}$	2
200	2 3 5	$\{2\}$	3
300	1 2 3 5	$\{3\}$	3
400	2 5	$\{4\}$	1 \rightarrow 2nd iteration because min supp was 2
		$\{5\}$	3 \downarrow min supp was 2 calculated above ($8/16 > 2$)

C ₂ Itemset	sup.	C ₂ Itemset	sup.
$\{1 2\}$	1	$\{1 2\}$	2
$\{1 3\}$	2	$\{1 3\}$	3
$\{1 5\}$	1	$\{2 3\}$	3
$\{2 3\}$	2	$\{2 5\}$	
$\{2 5\}$	3	$\{3 5\}$	
$\{3 5\}$	2		

$\rightarrow L_2$ Itemset	sup	C_3 Itemset	Itemset	sup
$\{1, 3\}$	2	$\{2, 3, 5\}$	$\{2, 3, 5\}$	2
$\{2, 3\}$	2			
$\{2, 5\}$	3			
$\{3, 5\}$	2			

∴ Frequent item sets are

$$L = \{2, 3, 5\}$$

Now for strong association rule:

Generate non-empty subset for L

$$S = \{2\}, \{3\}, \{5\}, \{2, 3\}, \{2, 5\}, \{3, 5\}$$

Find $S \rightarrow (L-S)$

$$\text{For eg: } \{2\} \rightarrow \{3, 5\}$$

→ calculated start
at the start

antecedent	$S \rightarrow (L-S)$	Support	Confidence	Confidence (%)
	$\{2\} \rightarrow \{3, 5\}$	2	$\frac{\{2, 3, 5\}}{\text{count}} \rightarrow \{2\}$	66.66
	$\{3\} \rightarrow \{2, 5\}$	2	$\frac{\{2, 3, 5\}}{\text{count}} \rightarrow \{3\}$	66.66
	$\{5\} \rightarrow \{2, 3\}$	2	$\frac{\{2, 3, 5\}}{\text{count}} \rightarrow \{5\}$	66.66
	$\{2, 3\} \rightarrow \{5\}$	2	$\frac{\{2, 3, 5\}}{\text{count}} \rightarrow \{2, 3\}$	100
	$\{2, 5\} \rightarrow \{3\}$	2	$\frac{\{2, 3, 5\}}{\text{count}} \rightarrow \{2, 5\}$	66.66
	$\{3, 5\} \rightarrow \{2\}$	2	$\frac{\{2, 3, 5\}}{\text{count}} \rightarrow \{3, 5\}$	100

Minimum confidence threshold = 80%.

∴ Strong association rules are -

$$\{2, 3\} \rightarrow \{5\}$$

$$\{3, 5\} \rightarrow \{2\}$$

Sup
2

- ② Consider the following database with minimum support count = 60%. Find all the frequent items set using apriori algorithm and also generates strong association rules if minimum confidences = 50%.

Transaction ID	Items Bought
T ₁	{M, O, N, K, E, Y}
T ₂	{D, O, N, K, E, Y}
T ₃	{M, A, K, E}
T ₄	{M, U, C, K, Y}
T ₅	{C, O, O, K, I, E}

Hint: O is bought 4 times in total, but it occurs in just 3 transaction.

Soln: Min support count = Min support threshold * Total no. of transactions

$$= \frac{60}{100} \times 5 \\ = 3 //$$

Database D

Transaction ID	Itemset	Scans	Itemset	Supp.
T ₁	{M, O, N, K, E, Y}	→	{A}	1
T ₂	{D, O, N, K, E, Y}	→	{C}	2
T ₃	{M, A, K, E}	→	{B}	1
T ₄	{M, U, C, K, Y}	→	{DE}	4
T ₅	{C, O, O, K, I, E}	→	{IY}	1
		→	{K}	5
		→	{M}	3
		→	{N}	2
		→	{O}	3
		→	{U}	1
		→	{Y}	3

Itemset	Supp.	Itemset	Itemset	Supp.
$\{E\}$	4	$\{E, K\}$	$\{E, K\}$	4
$\{K\}$	5	$\{E, M\}$	$\{E, M\}$	2
$\{M\}$	3	$\{E, O\}$	$\{E, O\}$	ScanD
$\{O\}$	3	$\{E, Y\}$	$\{E, Y\}$	2
$\{Y\}$	3	$\{K, M\}$	$\{K, M\}$	3
		$\{K, O\}$	$\{K, O\}$	3
		$\{K, Y\}$	$\{K, Y\}$	3
		$\{M, O\}$	$\{M, O\}$	1
		$\{M, Y\}$	$\{M, Y\}$	2
		$\{O, Y\}$	$\{O, Y\}$	2

↓

Itemset	supp.	Itemset	Itemset	Itemset
$\{E, K, O\}$	3	$\{E, K, O\}$	$\{E, K\}$	3
		ScanD		
			$\{E, O\}$	2
			$\{K, M\}$	3
			$\{K, O\}$	3
			$\{K, Y\}$	3

∴ Frequent item sets are

$$L = \{E, K, O\}$$

Now for strong association rule :

Generate non-empty subset for L

$$S = \{E\}, \{K\}, \{O\}, \{E, K\}, \{O, K\}, \{O, E\}$$

Find $S \rightarrow (L-S)$

$S \rightarrow (L-S)$	Support	Confidence	Confidence (%)
$\{E\} \rightarrow \{O, K\}$	3	3/4	75%
$\{O\} \rightarrow \{E, K\}$	3	3/3	100%
$\{K\} \rightarrow \{E, O\}$	3	3/5	60%
$\{O, K\} \rightarrow \{E\}$	3	3/3	100%
$\{E, K\} \rightarrow \{O\}$	3	3/4	75%
$\{E, O\} \rightarrow \{K\}$	3	3/3	100%

Minimum confidence threshold = 50%.

Strong association rules are -

$$\{OKY \rightarrow E\}$$

$$\{E OY \rightarrow \{K\}\}$$

$$\{OY \rightarrow \{EKY\}\}$$

* APRIORI ADVANTAGES / DISADVANTAGES :

Advantages

- Uses large itemset property
- Easy to implement
-

Disadvantages

- Requires many database scan

* FP GROWTH

- Improves apriori to a big extent.
- Frequent Itemset mining is possible without candidate generation.
- Only 'two' scans are required.

STEPS :-

- ① Build a compact data structure called the FP-tree
 - Built using 2 passes over the data-set.
- ② Extracts frequent item sets directly from the FP tree.

Eg: ~~Database~~ Find the frequent item set. (support count = 2)

① Database D

TID	Items
T100	I ₁ , I ₂ , I ₅
T200	I ₂ , I ₄
T300	I ₂ , I ₃
T400	I ₁ , I ₂ , I ₄
T500	I ₁ , I ₃
T600	I ₂ , I ₃
T700	I ₁ , I ₃
T800	I ₁ , I ₂ , I ₃ , I ₅
T900	I ₁ , I ₂ , I ₃

Itemset	support
{I ₁ }	6
{I ₂ }	7
{I ₃ }	6
{I ₄ }	2
{I ₅ }	2

Descending order

② Itemset

TID	Items
T200	I ₂ , I ₄
T100	I ₁ , I ₂ , I ₃
T300	I ₂ , I ₃
T400	I ₁ , I ₂ , I ₄
T500	I ₁ , I ₃

Itemsets

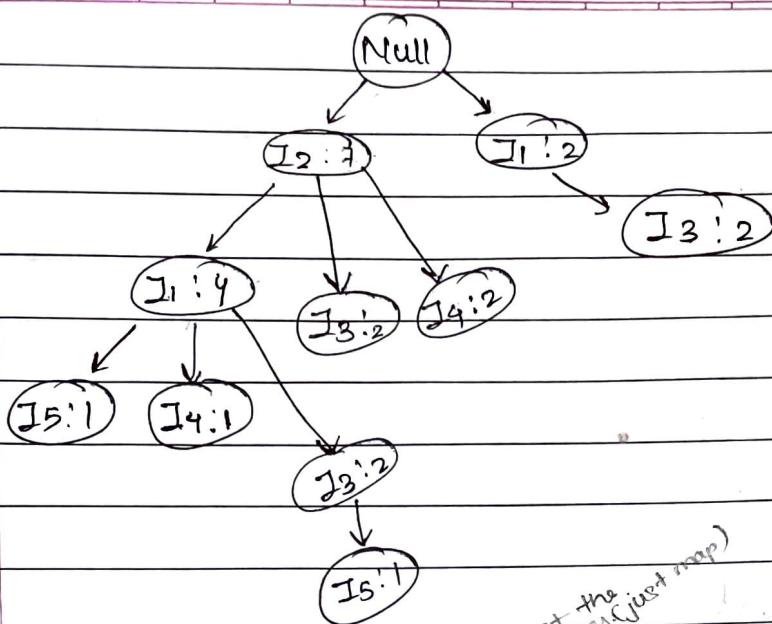
Itemsets	Support
{I ₂ }	7
{I ₁ }	6
{I ₃ }	6
{I ₄ }	2
{I ₅ }	2

(this order
should be
followed
while
dealing)

③ ID

TID	Item
T100	I ₂ , I ₁ , I ₅
T200	I ₂ , I ₄
T300	I ₂ , I ₃
T400	I ₂ , I ₁ , I ₄
T500	I ₁ , I ₃
T600	I ₂ , I ₃
T700	I ₁ , I ₃
T800	I ₂ , I ₁ , I ₃ , I ₅
T900	I ₂ , I ₁ , I ₃

[∴ I₂ > I₁ > I₃ in this arrangement]



Item	Conditioned Pattern Base	Conditional FP-Tree	Frequent Pattern Generated
I ₅	{I ₂ , I ₁ :1}, {I ₂ , I ₁ , I ₃ :1}	{I ₂ :2, I ₁ :2}	{I ₂ , I ₅ :2}, {I ₁ , I ₅ :2}, {I ₂ , I ₁ , I ₅ :2}
I ₄	{I ₂ , I ₁ :1}, {I ₂ :1}	{I ₂ :2}	{I ₂ , I ₄ :2}
I ₃	{I ₂ :2} (RST), {I ₂ , I ₁ :2}, {I ₂ :4, I ₁ :2}, {I ₁ :2} (KSD)	{I ₂ , I ₃ :4}, {I ₂ , I ₃ :2}, {I ₁ :2}	{I ₂ , I ₃ :2}, {I ₁ :2}
I ₂	{I ₂ :4}	{I ₂ :4}	{I ₂ , I ₁ :4}
I ₁	IGNORE AS NO BRANCH		

Q) Draw FP tree

T₁: b, e

T₂: a, b, c, e

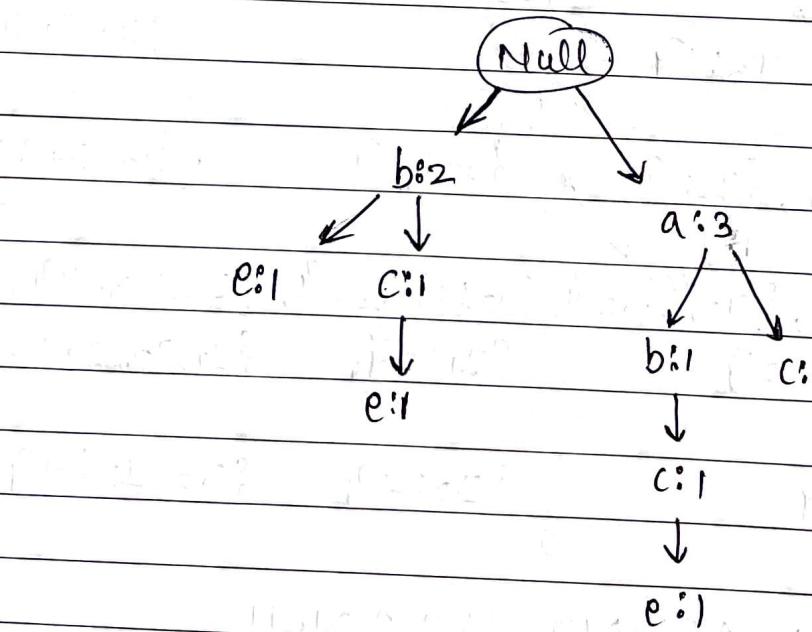
T₃: b, c, e

T₄: a, c

T₅: a

Soln: Items Support

a	3
b	3
c	3
e	3



③ Draw the FP tree. Support count = 3.

T₁: f, a, c, d, g, i, m, p

T₂: a, b, c, f, l, m, o

T₃: b, f, h, j, o

T₄: b, c, k, s, p

T₅: a, f, c, e, l, p, m, n

Soln: Items Support

a 3

b 3

c 4

d 1

e 1

f 3

g 1

h 1

i 1

j 1

k 1

l 2

m 3

n 1

o 2

p 3

s 1

Items support

a 3

b 3

c 4

f 4

m 3

p 3

↓ Descending order

Items support count

f 4

c 4

a 3

b 3

m 3

p 3

↓

TID Items

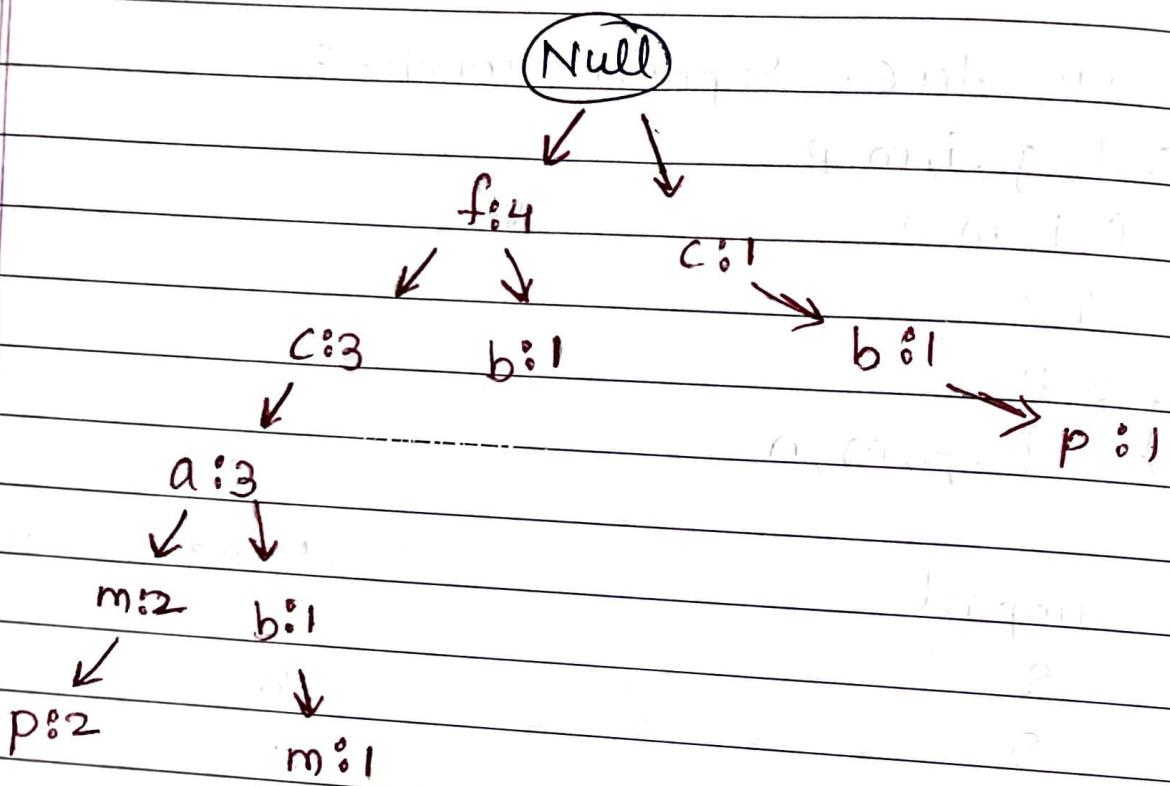
T₁ f, c, a, m, p

T₂ f, c, a, b, m

T₃ f, b, ~~p~~

T₄ c, b, p

T₅ f, c, a, m, p



CHAPTER 4: CLASSIFICATION TECHNIQUES

* Classification : (Definition).

- ↳ Process of finding a model that describes data classes.
- ↳ For the purpose of being able to use the model.
- ↳ To predict the class of objects whose class label is unknown.

* SUPERVISED CLASSIFICATION : (Type 1)

- ↳ Set of possible classes is known in advance.

* UNSUPERVISED CLASSIFICATION : (Type 2)

- ↳ Set of possible classes is not known. After classification we can try to assign a name to that class. Unsupervised classification is called clustering.

NOTE : Classification consists of assigning a class label to a set of unclassified class.

* WHY CLASSIFICATION? A MOTIVATING APPLICATION

Eg:

- Credit Approval
 - A bank wants to classify its customers based on whether they are expected to pay back their approved loans.
 - The history of past customers is used to train the classifier.
 - The classifier provides rules, which identify potentially reliable future customers.
 - Classification rule :
 - If age = "31...40" & income = high & credit-rating = excellent.

- Future customers
- Paul : age = 35, income = high \Rightarrow Excellent credit rating
- John : age = 20, income = medium \Rightarrow Fair credit rating

* TYPICAL APPLICATIONS OF CLASSIFICATION :

- Credit approval
- Target Marketing
- Medical Diagnosis
- Fraud Detection

* CLASSIFICATION - A TWO STEP PROCESS

① MODEL CONSTRUCTION :

- Describing a set of predetermined classes.
- The model is represented as classification rules, decision trees, or mathematical formulas.

② MODEL USAGE :

- For classifying future or unknown objects.
- Test sample is compared with the classified result from the model.

Classification Model can be in the following forms :

- ↳ IF - THEN Rules
- ↳ A decision Tree
- ↳ Neural Network

.... (Not in syllabus)

CLASSIFIER TECHNIQUES :

* BAYESIAN

- A Bayes classifier is a simple probabilistic classifier.
- A naive Bayes classifier will be studied.

NAIVE BAYES CLASSIFICATION -

- ① Eg : Predict the class label for an unknown sample "x" using Naive Bayesian classification.

'x' = (outlook = Sunny, Temperature = Cool, Humidity = High, Wind = Strong).

Play Tennis : Training Examples

(Target
↓ attribute)

Day	Outlook	Temperature	Humidity	Wind	Play Tennis
D1	Sunny	Hot	High	Weak	No
D2	Sunny	Hot	High	Strong	No
D3	Overcast	Hot	High	Weak	Yes
D4	Rain	Mild	High	Weak	Yes
D5	Rain	Cool	Normal	Weak	Yes
D6	Rain	Cool	Normal	Strong	No
D7	Overcast	Cool	Normal	Strong	Yes
D8	Sunny	Mild	High	Weak	No
D9	Sunny	Cool	Weak	Weak	Yes
D10	Rain	Mild	Normal	Weak	Yes
D11	Sunny	Mild	Weak	Normal	Strong
D12	Overcast	Mild	Strong	High	Strong
D13	Overcast	Hot	Weak	Normal	Weak
D14	Rain	Mild	Strong	High	Strong

Soln:Learning Phase :

$$P(\text{Play} = \text{Yes}) = 9/14$$

$$P(\text{Play} = \text{No}) = 5/14$$

Outlook	Play = Yes	Play = No	Temp	Play = Yes	Play = No
Sunny	2/9	3/5	Hot	2/9	2/5
Overcast	4/9	0/5	Mild	4/9	2/5
Rain	3/9	2/5	Cool	3/9	1/5

Humidity	Play = Yes	Play = No	Wind	Play = Yes	Play = No
High	3/9	4/5	Strong	3/9	3/5
Normal	6/9	1/5	Weak	6/9	2/5

Test Phase :

- Given a new instance,

$x' = (\text{Outlook} = \text{Sunny}, \text{Temperature} = \text{Cool}, \text{Humidity} = \text{High}, \text{Wind} = \text{Strong})$

- MAP RULE

$$\begin{aligned}
 P(x'| \text{Yes}) &= P(\text{Outlook} = \text{Sunny} / \text{Yes}) * P(\text{Temp} = \text{cool} / \text{Yes}) * \\
 &\quad P(\text{Humidity} = \text{High} / \text{Yes}) * P(\text{Wind} = \text{Strong} / \text{Yes}) * \\
 &\quad P(\text{Yes}) \\
 &= 2/9 * 3/9 * 3/9 * 3/9 * 9/14 \\
 &= 0.0053
 \end{aligned}$$

$$\begin{aligned}
 \hookrightarrow P(x'|No) &= P(\text{Outlook} = \text{Sunny}/No) * P(\text{Temp} = \text{cool}/No) * \\
 &\quad P(\text{Humidity} = \text{High}/No) * P(\text{Wind} = \text{Strong}/No) * \\
 &\quad P(\text{No}) \\
 &= 3/5 * 1/5 * 4/5 * 3/5 * 5/14 \\
 &= 0.0206 //
 \end{aligned}$$

~~V. INN:~~ \Rightarrow Given the fact $P(x|Yes) < \cancel{P(x|No)}$, we label ~~No~~ x to be "Play Tennis = No." (... Last statement has 2 Marks).

(2) Naive Bayesian classification eg :

- Predict a class label of an unknown sample using Naive Bayesian classification on the following training dataset from all electronics customer database.
- The unknown sample is -
 $x' = \{\text{age} = "<= 30", \text{Income} = "\text{Median}", \text{Student} = "\text{Yes}",$
 $\text{Credit rating} = "\text{fair"}\}$

Age	Income	Student	Credit Rating	Buys computer
≤ 30	High	No	Fair	No
≤ 30	High	No	Excellent	No
$31 \dots 40$	High	No	Fair	Yes
> 40	Medium	No	Fair	Yes
> 40	Low	Yes	Fair	Yes
> 40	Low	Yes	Excellent	No
$31 \dots 40$	Low	Yes	Excellent	Yes
≤ 30	Medium	No	Fair	No
≤ 30	Low	Yes	Fair	Yes
> 40	Medium	Yes	Fair	Yes

≤ 30	Medium	Yes	Excellent	Yes
31...40	Medium	No	Excellent	Yes
31...40	High	Yes	Fair	Yes
> 40	Medium	No	Excellent	No

Soln: Learning Phase:

$$P(\text{Buys-computer} = \text{Yes}) = 9/14$$

$$P(\text{Buys-computer}) = 5/14$$

Age	Play = Yes	Play = No	Income	Play = Yes	Play = No
≤ 30	2/9	3/5	High	2/9	2/5
31...40	4/9	0/5	Medium	4/9	2/5
> 40	3/9	2/5	Low	3/9	1/5

Student	Play = Yes	Play = No	Credit_Rating	Play = Yes	Play = No
Yes	6/9	1/5	Excellent	8/9	3/5
No	4/9	4/5	Fair	6/9	2/5

Test Phase:

- Given a new instance,

$x' = \{\text{age} \geq \text{"}< 30\text{"}, \text{Income} = \text{"Medium"}, \text{Student} = \text{"Yes"}, \text{Credit-Rating} = \text{"Fair"}\}$

- MAP RULE

$$\hookrightarrow P(x' | \text{Yes}) = P(\text{Age} = \text{"}< 30/\text{Yes}\text{"}) * P(\text{Income} = \text{"Medium}/\text{Yes}\text{"}) \\ * P(\text{Student} = \text{Yes}) * P(\text{Credit-Rating} = \text{Yes}) \\ * P(\text{Yes})$$

$$= \frac{2}{9} * \frac{4}{9} * \frac{6}{9} * \frac{3}{9} * \frac{9}{14}$$

$$= 0.028 //$$

→ $P(x' | N_0) = P(\text{Age} = < 30 | N_0) * P(\text{Income} = \text{Medium} | N_0)$

$* P(\text{Student} | N_0) * P(\text{Credit_Rating} | N_0)$

$* P(N_0)$

$$= \frac{3}{5} * \frac{2}{5} * \frac{1}{5} * \frac{2}{5} * \frac{5}{14}$$

$$= 0.007 //$$

Since $0.028 > 0.007$,

∴ The Naïve Bayesian classifier predicts Buyers computer = "Yes" for sample x .

(Wally
lengthy
attempt
in exam) *

DECISION TREE

- (Classification)

- Builds classification models in the form of a tree structure.
- It breaks down data into smaller subsets.

Eg:

- ① Using given training data set, create classification model using decision tree.

(Table same as Eg 1 of Naïve Bayes classification on pg 19 in the N.B.)

Soln: $E(S) = \sum_{i=1}^c -P_i \log_2 P_i$ → No. of values. ($E \rightarrow \text{Entropy}$)

Play Golf

Yes	No
9/14	5/14

Entropy (Play Golf) = Entropy (5, 9)

$$= \text{Entropy}(0.36, 0.64) \quad \begin{matrix} 5/14, 9/14 \\ \downarrow \end{matrix}$$

$$= -(0.36 \log_2 0.36) -$$

$$-(0.64 \log_2 0.64)$$

$$= 0.94$$

target

$$\Rightarrow E(T, x) = \sum_{c \in X} P(c) E(c)$$

		Play Golf		
		Yes	No	
OUTLOOK	Sunny	3	2	5
	Overcast	4	0	4
	Rainy	2	3	5



probability

From the table.

$$E(\text{PlayGolf}, \text{Outlook}) = P(\text{sunny}) * E(3, 2) + P(\text{Overcast}) * E(4, 0) + P(\text{Rainy}) * E(2, 3)$$

$$= (5/14) * 0.971 + (4/14) * 0.0 + (5/14) * 0.971$$

$$= 0.693$$

\Rightarrow

		Play Golf		
		Yes	No	
TEMPERATURE	Hot	2	2	4
	Mild	4	2	6
	Cold	3	1	4

$$E(\text{PlayGolf}, \text{Temp}) = P(\text{Hot}) * E(2/4, 2/4) + P(\text{Mild}) * E(4/6, 2/6) \\ + P(\text{Cold}) * E(3/4, 1/4)$$

$$= 4/14 * 1 + 6/14 * 0.92 + 4/14 * 0.81 \\ = 0.91$$

Similarly for
 $E(4/6, 2/6)$ &
 $E(3/4, 1/4)$

Similarly,

$$E(\text{Humidity}) = 0.788$$

$$E(\text{Wind}) = 0.892$$

$$[E = \text{Entropy}(0.5, 0.5) \\ = - (0.5 \log_2 0.5) - (0.5 \log_2 0.5) \\ = 1]$$

$$\Rightarrow \boxed{Gain(T, X) = \text{Entropy}(T) - \text{Entropy}(T, X)} \rightarrow \text{FORMULA OF INFORMATION GAIN.}$$

$$\circ G(\text{Play Golf}, \text{Outlook}) = E(\text{PlayGolf}) - E(\text{PlayGolf}, \text{Outlook}) \\ = 0.940 - 0.693 \\ = 0.247$$

$$\circ G(\text{PlayGolf}, \text{Temp}) = E(\text{PlayGolf}) - E(\text{PlayGolf}, \text{Temp}) \\ = 0.940 - 0.91 \\ = 0.030$$

$$\circ G(\text{PlayGolf}, \text{Humidity}) = E(\text{PlayGolf}) - E(\text{PlayGolf}, \text{Temp}) \\ = 0.940 - 0.788 \\ = 0.152$$

$$\begin{aligned}
 P(G(\text{Play Golf}, \text{Wind})) &= E(\text{Play Golf}) - E(\text{Play Golf, Wind}) \\
 &= 0.940 - 0.0892 \\
 &= 0.048
 \end{aligned}$$

\therefore The highest value is the outlook,

\therefore Root Node will be the outlook.

Outlook = Overcast

Refer the main table
(Mini table)

Temp	Humidity	Windy	Play Golf
Hot	High	False	Yes
Cool	Normal	True	Yes
Mild	High	True	Yes
Hot	Normal	False	Yes
Hot	High	False	Yes

NOTE : Keep calculating the entropies of the sub-division also until you get the final leaf node.

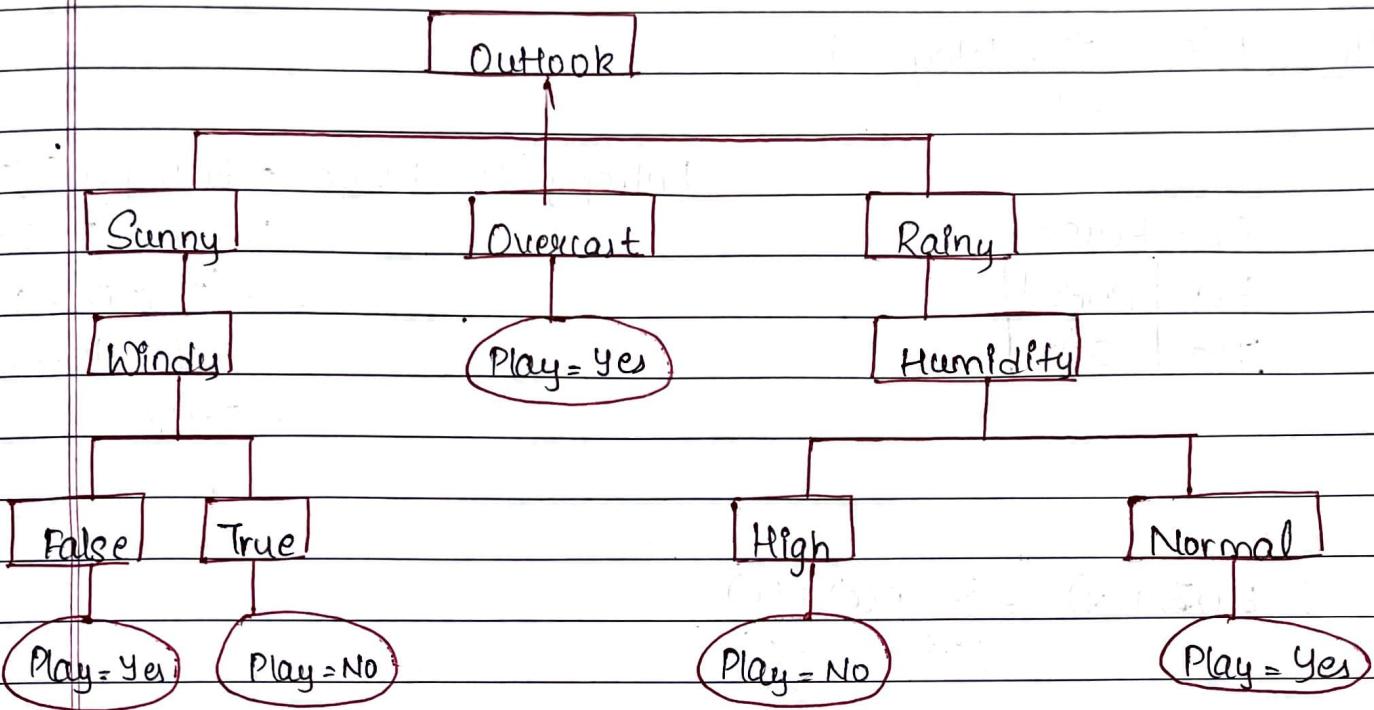
R₁ : IF(Outlook = Sunny) AND (Windy = False) THEN Play = Yes

R₂ : IF(Outlook = Sunny) AND (Windy = TRUE) THEN Play = No

R₃ : IF(Outlook = Overcast) THEN Play = Yes

R₄ : IF(Outlook = Rainy) AND (Humidity = High) THEN Play = No

R₅ : IF(Outlook = Rain) AND (Humidity = Normal) THEN Play = Yes



DECISION TREE

- (2) Create classification model using decision tree for the following training data set.

SR.NO.	Income	Age	Own House
1	Very High	Young	Yes
2	High	Medium	Yes
3	Low	Young	Rented
4	High	Medium	Yes
5	Very High	Medium	Yes
6	Medium	Young	Yes
7	High	Old	Yes
8	Medium	Medium	Rented
9	Low	Medium	Rented
10	Low	Old	Rented
11	High	Young	Yes
12	Medium	Old	Rented

Soln: $E(S) = \sum_{i=1}^C -P_i \log_2 P_i$

Own House

Yes	Rented
7/12	5/12

$$\text{Entropy (Own House)} = \text{Entropy}(7/12, 5/12)$$

$$= \text{Entropy}(0.58, 0.41)$$

$$= -(0.58 \log_2 0.58)$$

$$- (0.41 \log_2 0.41)$$

$$= 0.98$$

$$\Rightarrow E(O, T, X) = \sum_{C \in X} P(C) E(C)$$

		Own House		
		Yes	Rented	
INCOME	Low	0	3	3
	Medium	1	2	3
	High	4	0	4
	Very high	2	0	2

$$E(\text{Own House, Income}) = P(\text{Low}) * E(0, 3/3) + P(\text{Medium}) * E(1/3, 2/3) \\ + P(\text{High}) * E(4/4, 0) + P(\text{Very high}) * E(2/2, 0)$$

$$= \frac{3}{12} \times 0 + \frac{3}{12} \times 0.18 + \frac{4}{12} \times 0 + \frac{2}{12} \times 0$$

$$= 0.23$$

		Own House		
		Yes	Rented	
AGE	Young	3	1	4
	Medium	3	2	5
	Old	1	2	3

$$\begin{aligned}
 E(\text{Own house, Age}) &= P(\text{Young}) * E(3/4, 1/4) + P(\text{Medium}) * \\
 &\quad E(3/5, 2/5) + P(\text{Old}) * E(1/5, 2/5) \\
 &= \frac{4}{12} \times 0.81 + \frac{5}{12} \times 0.97 + \frac{3}{12} \times 0.92 \\
 &= 0.90
 \end{aligned}$$

~~∴ the higher value is age~~

~~∴ The root note will be the age.~~

$$G(O, A) = E(O) - E(O, A)$$

$$= 0.98 - 0.90$$

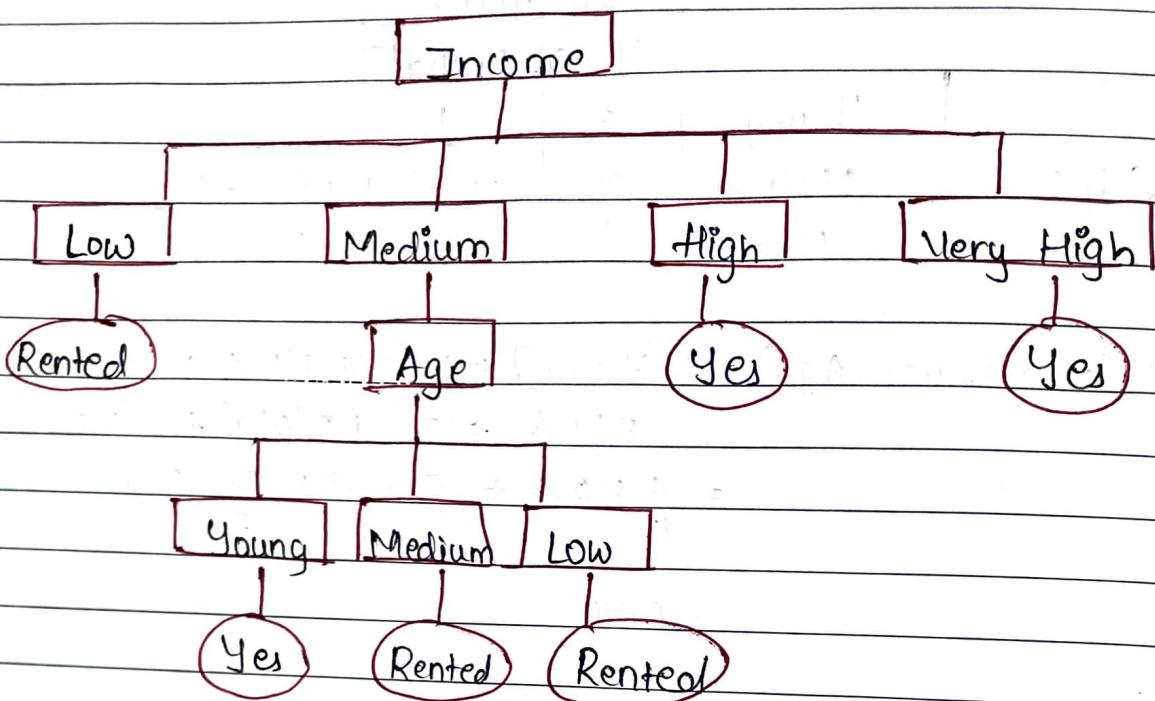
$$= 0.08$$

$$G(O, I) = E(O) - E(O, I)$$

$$= 0.98 - 0.823$$

$$= 0.75$$

~~∴ Income has the higher gain, ∴ The root note will be the income.~~



* PREDICTION

(Data Mining Technique)

- We can predict the continues value of response variable with the help of predictor variable.
- It can be done with the help of statistical technique of regression.
- It assumes the data to fit in some kind of function & involve study of those function.
- Most widely used approach for numeric prediction is regression.

* TYPES OF REGRESSION -

- ① Linear
- ② Multiple Linear
- ③ Non-Linear

⇒ LINEAR REGRESSION (Single predictor variable)

- Data is modeled using straight line.

- FORMULA = $y = \alpha + \beta x$ → (Independent variable)

(Dependent variable)

where $x \rightarrow$ predictor variable

$y \rightarrow$ Response variable

⇒ MULTIPLE LINEAR REGRESSION (Multiple predictor variable)

- FORMULA = $y = \alpha + \beta_1 x_1 + \beta_2 x_2 + \beta_3 x_3$

⇒ NON-LINEAR REGRESSION

- If the given response variable & predictor variable have got polynomial relationship then it is called non-linear regression.

- FORMULA = $y = \alpha + \beta_1 x + \beta_2 x^2 + \beta_3 x^3$

* How to calculate α & β .

① Using least square method

$$\beta = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{\sum_{i=1}^n (x_i - \bar{x})^2}$$

$$\alpha = \bar{y} - \beta \bar{x}$$

$\bar{x} \rightarrow$ Mean value of x

$\bar{y} \rightarrow$ Mean value of y

Eg

①	Midterm (x)	Final year (y)
	45	60
	70	70
"	60	54
	84	82
	75	68
	84	76

Find:

- 1) Equation of prediction (linear regression formula)
- 2) What will be the final year marks if the midterm marks is 10?

Sln:	x	y	$(x - \bar{x})$	$(y - \bar{y})$	$(x - \bar{x})(y - \bar{y})$	$(x - \bar{x})^2$
	45	60	-24.67	-8.33	205.5	608.608
	70	70	0.33	1.67	0.55	0.1089
	60	54	-9.67	-14.33	138.57	93.508
	84	82	14.33	13.67	195.9	205.348
	75	68	5.33	-0.33	-1.76	28.408
	84	76	14.33	7.67	109.9	205.348
					$\Sigma = 648.66$	$\Sigma = 1141.32$

$$\bar{x} = 69.67$$

$$\bar{y} = 68.33$$

$$1) \beta = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{\sum_{i=1}^n (x_i - \bar{x})^2}$$

$$\beta = \frac{648.66}{1141.32}$$

$$\beta = 0.568 //$$

$$\alpha = \bar{y} - \beta \bar{x}$$

$$= 68.33 - 0.568 \times 69.67$$

$$= 28.76 //$$

Equation of prediction = $y = \alpha + \beta x$

$$y = 28.76 + 0.568 x$$

2) For $x = 40$, we get $y = ?$

$$y = 28.76 + 0.568 x \quad (\because x = 40)$$

$$\therefore y = 28.76 + 0.568 \times 40$$

$$y = 51 //$$

$$\therefore y = 51 \text{ when } x = 40.$$

* MODEL EVALUATION AND SELECTION

- Methods for estimating a classifier's accuracy
 - Holdout method, random subsampling
 - Cross-validation
 - Bootstrap
- Comparing classifier
 - ROC curves

(Q → Explain different methods that can be used to evaluate & compare the accuracy of different classification algorithm) (Ans is the above points).

Q → Explain Bagging & Boosting of classification

* CLASSIFIER EVALUATION METRICS : CONFUSION MATRIX

Confusion Matrix :

Actual class \ Predicted class		C_1	$-C_1$
C_1	C_1	True positives (TP)	False negative (FN)
	$-C_1$	False positives (FP)	True negative (TN)

Eg :

Actual class \ Predicted class	Buy-computer = Yes	Buy-computer = No	Total
Buy-computer = yes	6954 (FP) 6954412 (FP)	46 (FN) 2588 (TN)	7000
Buy-computer = no	7366	2634	3000
Total			10000

$$\text{Accuracy} = \frac{\text{TP} + \text{TN}}{\text{Total}}$$

$$= \frac{6951 + 2588}{10000}$$

$$= 95\%$$

* FORMULAS:

^{IMP} ① Classifier accuracy = $\frac{\text{TP} + \text{TN}}{\text{Total}}$

^{IMP} ⑥ Specification = $\frac{\text{TN}}{\text{TN} + \text{FP}}$

^{IMP} ② Error Rate = $\frac{\text{FP} + \text{FN}}{\text{Total}}$ OR 1 - Accuracy

^{IMP} ③ Precision = $\frac{\text{TP}}{\text{TP} + \text{FP}}$

^{IMP} ④ Recall / Sensitivity = $\frac{\text{TP}}{\text{TP} + \text{FN}}$

^{Not imp:} ⑤ F (Measure) (F₁ or F-score) = $\frac{2 \times \text{precision} \times \text{recall}}{\text{precision} + \text{recall}}$

Eg:

Actual class \ Predicted class	cancer = yes	cancer = no	Total
Cancer = yes	90	210	300
Cancer = no	140	9560	9700
Total	230	9770	10000

1) Accuracy = $\frac{TP + TN}{Total}$

$$= \frac{90 + 9560}{10000}$$

$$= 96.5\% //$$

2) Error Rate = $\frac{FP + FN}{Total}$ OR Accuracy

$$= 1 - 96.5\%$$

$$= 1 - 0.965$$

$$= 0.035$$

$$= 3.5\% //$$

3) Precision = $\frac{TP}{TP + FP}$

$$= \frac{90}{90 + 40}$$

$$= 0.391 = 39.1\% //$$

4) Recall = $\frac{TP}{TP + FN}$

$$= \frac{90}{90 + 210}$$

$$= 0.3$$

$$= 30\% //$$

5) Specificity = $\frac{TN}{TN + FP}$

$$= 9650 / 9650 + 140$$

$$= 0.9856$$

$$= 98.56\% //$$

MODULE 5: CLUSTERING

* CLUSTER ANALYSIS :

- ↳ Cluster : a collection of data objects.
- ↳ Similar to one another within the same cluster.
- ↳ Dissimilar to the objects in other clusters.
- ↳ Cluster Analysis - Finding similarities between data according to the characteristic found in the data and grouping similar data objects into clusters.

* EXAMPLES OF CLUSTERING APPLICATIONS :

- ↳ Marketing
- ↳ Land use
- ↳ Insurance
- ↳ City-planning
- ↳ Earth-quake studies

* QUALITY : ~~WHAT IS GOOD CLUSTERING ?~~

- A good clustering method will produce high quality clusters with
 - ↳ High intra-class similarity
 - ↳ Low inter-class similarity

* CLUSTERING

K-Means Clustering - (Steps) -

- ① Partition objects into k nonempty subsets.
- ② Compute seed points as the centroids of the clusters of the current partition.
- ③ Assign each object to the cluster with the nearest seed point.

④ Go back to step 2, stop when no more new assignment.

Example:

① Explain k-means clustering and solve the following with $k=2$. $\{2, 3, 4, 10, 11, 12, 20, 25, 30\}$ (One dimensional)

Soln:

Step 1: $M_1 = -4$, $M_2 = 12$ (We can take any random pts)

$$K_1 = \{2, 3, 4\}$$

$$K_2 = \{10, 11, 12, 20, 25, 30\}$$

Step 2: New centroid points has to be calculated.

$$M_1 = (2+3+4)/3 = 3$$

$$\therefore M_1 = 3$$

$$M_2 = (10+11+12+20+25+30)/6 = 18$$

$$\therefore M_2 = 18$$

$$K_1 = \{2, 3, 4, 10\}$$

$$K_2 = \{11, 12, 20, 25, 30\}$$

Step 3: Reassign the centroid

$$M_1 = (2+3+4+10)/4 = 4.75$$

$$\therefore M_1 = 4.75$$

$$M_2 = (11+12+20+25+30)/5 = 19.6$$

$$\therefore M_2 = 19.6$$

$$K_1 = \{2, 3, 4, 10, 11, 12\}$$

$$K_2 = \{20, 25, 30\}$$

Step 4: Reassign the centroid

$$M_1 = 7$$

$$M_2 = 25$$

$$K_1 = \{2, 3, 4, 10, 11, 12\}$$

$$K_2 = \{20, 25, 30\}$$

Step 5: Reassign the centroid

$M_1 = 7$, $M_2 = 25$ (Same as previous in step 4)

We have similar means

Final clusters are $K_1 = \{2, 3, 4, 10, 11, 12\}$ & $K_2 = \{20, 25, 30\}$

(2) Cluster the given set of values $\{2, 3, 6, 8, 9, 12, 15, 18, 22\}$ into three clusters.

Soln:

Step 1: $M_1 = 2$, $M_2 = 3$, $M_3 = 6$ (Given)

$K_1 = \{2\}$, $K_2 = \{3\}$, $K_3 = \{6, 8, 9, 12, 15, 18, 22\}$

Step 2: $M_1 = 2$, $M_2 = 3$, $M_3 = 12.86$

$K_1 = \{2\}$, $K_3 = \{8, 9, 12, 15, 18, 22\}$

$K_2 = \{3, 6\}$

Step 3: $M_1 = 2$, $M_2 = 3$, $M_3 = 12.86$

$K_1 = \{2, 3\}$, $K_2 = \{6, 8, 9\}$, $K_3 = \{12, 15, 18, 22\}$

Step 4: $M_1 = 2.5$, $M_2 = 7.6$, $M_3 = 16.75$

$K_1 = \{2, 3\}$, $K_2 = \{6, 8, 9, 12\}$, $K_3 = \{12, 15, 18, 22\}$

Step 5: $M_1 = 2.5$, $M_2 = 8.75$, $M_3 = 18.33$

$K_1 = \{2, 3\}$, $K_2 = \{6, 8, 9, 12\}$, $K_3 = \{15, 18, 22\}$

Final mean values are $M_1 = 2.5$, $M_2 = 8.75$, $M_3 = 18.33$

& final clusters are $K_1 = \{2, 3\}$ $K_3 = \{15, 18, 22\}$
 $K_2 = \{6, 8, 9, 12\}$

* ALGORITHM FOR K-MEANS CLUSTERING METHOD -

Given K , the K -means algorithm is implemented in 4 steps :

- 1) Partition objects into K nonempty subsets.
- 2) Compute seed points as the centroids of the clusters.

Example

① Object

A
B
C
D
E

Attribute 1(x):

2
3
1
3
1.5

Attribute 2(y):

2
2
1
1
0.5

Euclidean Distance: $D(i,j) = \sqrt{(x_2 - x_1)^2 + (y_2 - y_1)^2}$

$$C_1 = A(2,2)$$

$$C_2 = C(1,1)$$

Distance Matrix -

D^0	A	B	C	D	E	Cluster	Centroid
0	1		1.41	1.41	1.51	C_1	GP_1
1.41		2.24	0	2	0.71	C_2	GP_2

$$G^0 = (0 < 1.4) (1 < 2.24) (0 < 1.41) (1.41 < 2) (0.71 < 1.58)$$

A	B	C	D	E	Cluster centroid
1	1	0	1	0	$C_1 - GP_1$
0	0	1	0	1	$C_2 - GP_2$

$$K_1 = \{A, B, D\}$$

$$K_2 = \{C, E\}$$

$$[A(x) + B(x) + D(x)], [A(y) + B(y) + D(y)]$$

Step 2: $C_1 = [(2+3+3)/3, (2+2+1)/3] = (2.67, 1.67)$

$$C_2 = [(1+1.5)/2, (1+0.5)/2] = (1.25, 0.75)$$

$$D^* = \begin{matrix} A & B & C & D & E \end{matrix} \text{ Cluster Centroid}$$

$$\begin{matrix} 0.75 & 0.47 & 1.79 & 0.75 & 1.65 \end{matrix} \text{ } C_1 - GP_1$$

$$\begin{matrix} 1.45 & 2.25 & 0.32 & 1.76 & 0.35 \end{matrix} \text{ } C_2 - GP_2$$

$$G' = \begin{matrix} A & B & C & D & E \end{matrix} \text{ Cluster centroid}$$

$$\begin{matrix} 1 & 1 & 0 & 1 & 0 \end{matrix} \text{ } C_1 - GP_1$$

$$\begin{matrix} 0 & 0 & 1 & 0 & 1 \end{matrix} \text{ } C_2 - GP_2$$

$$K_1 = \{A, B, D\}$$

$$K_2 = \{C, E\}$$

$$\therefore G^* = G'$$

Final clusters are: $GP_1 = \{A, B, D\}$

$$GP_2 = \{C, E\}$$

- ② Group the following objects into $K=2$ group of medicine based on the two factors (ph & weight index).

<u>Object</u>	Attribute 1(x)	Attribute 2(y)
Medicine A	1	1
B	2	1
C	4	3
D	5	4

$$C_1 = A(1,1), \quad C_2 = B(2,1)$$

$$D^o = \begin{array}{cccc} A & B & C & D \\ 0 & 1 & 3.6 & 5 \\ 1 & 0 & 2.83 & 4.24 \end{array}$$

$$G^o = \begin{array}{cccc} A & B & C & D \\ 1 & 0 & 0 & 0 \\ 0 & 1 & 1 & 1 \end{array}$$

$$K_1 = \{A\}, \quad K_2 = \{B, C, D\}$$

Step 2:

$$C_1 = A(1,1)$$

$$C_2 = [(2+4+5)/3, (1+3+4)/3] = (3.67, 2.67)$$

~~B2~~

$$D' = \begin{array}{cccc} A & B & C & D \\ 0 & 1 & 3.6 & 5 \\ 3.15 & 2.36 & 0.467 & 1.88 \end{array}$$

$$G' = \begin{array}{cccc} A & B & C & D \\ 1 & 1 & 0 & 0 \\ 0 & 0 & 1 & 1 \end{array}$$

$$K_1 = \{A, B\}$$

$$K_2 = \{C, D\}$$

Step 3:

$$C_1 = (1.5, 1)$$

$$C_2 = (4.5, 3.5)$$

$$D^2 = \begin{array}{cccc} A & B & C & D \\ 0.5 & 0.5 & 3.20 & 4.61 \\ 4.30 & 3.53 & 0.71 & 0.71 \end{array}$$

$$G^2 = \begin{array}{cccc} A & B & C & D \\ 1 & 1 & 0 & 0 \\ 0 & 0 & 1 & 1 \end{array}$$

$$K_1 = \{A, B\}$$

$$K_2 = \{C, D\}$$

$$\therefore G^1 = G^2$$

Final clusters are - $GP_1 = \{A, B\}$

$$GP_2 = \{C, D\}$$

* PROBLEM OF K-MEANS METHOD

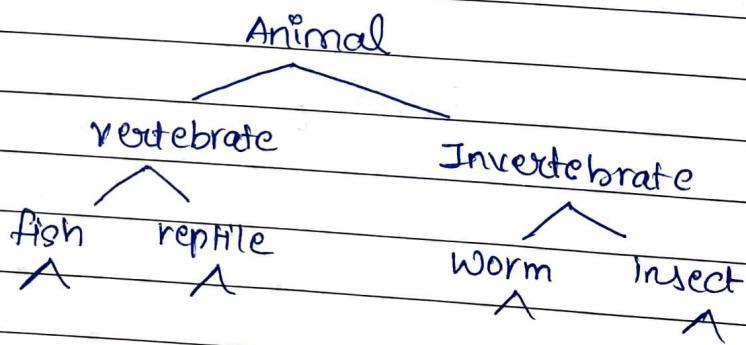
- K-means algorithm is very sensitive to outliers.
- K-Medoids - instead of taking the mean value of the object in a cluster as a reference point, medoids can be used, which is the most centrally located object in a cluster.

* K-MEDOID

- Name of the algorithm - PAM (Partitioning Around Medoids)

* HIERARCHICAL CLUSTERING

- Build a tree-based hierarchical taxonomy (dendrogram) from a set of documents.
- Eg :



- Clustering obtained by cutting the dendrogram at a desired level.
- Each connected component forms a cluster.

* HIERARCHICAL CLUSTERING ALGORITHMS

- i) Agglomerative (bottom-up)
- ii) Divisive (Top-bottom)

(Not in syllabus)

AGGLOMERATIVE HIERARCHICAL CLUSTERING METHOD:

- ↳ Single Link algorithm
- ↳ Complete Link algorithm
- ↳ Average Link algorithm

* TYPICAL ALTERNATIVES TO CALCULATE THE DISTANCE BETWEEN CLUSTERS

- ↳ Single Link: Smallest distance between an element in one ~~other~~ cluster and an element in the other.
i.e, $\text{dis}(k_i, k_j) = \min(t_{ip}, t_{jq})$

↳ Complete Link: $\text{dis}(k_i, k_j) = \max(t_{ip}, t_{jq})$

↳ Average: $\text{dis}(k_i, k_j) = \text{avg}(t_{ip}, t_{jq})$

~~Not imp~~ ↳ Centroid: $\text{dis}(k_i, k_j) = \text{dis}(c_i, c_j)$

~~Not imp~~ ↳ Medoid: $\text{dis}(k_i, k_j) = \text{dis}(M_i, M_j)$

* Example.

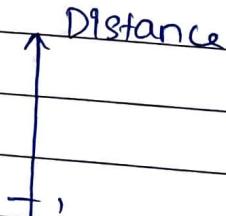
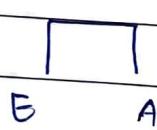
①	E	A	C	B	D
E	0				
A	1	0			
C	2	2	0	3	
B	2	5	1	0	
D	3	3	6	3	0

Solve the following eg using Single Link Hierarchical clustering.

Soln: Step 1: Identify the two clusters with the shortest distance in the matrix, and merge them together. Recompute, the distance matrix, as those two clusters are now in a single cluster, (no longer exist by themselves).

By looking at the distance matrix above, we see that E & A have the smallest distance from all. So, we merge those two in a single cluster, & recompute the distance matrix.

Dendrogram



Distance Matrix

$$\rightarrow \text{dist}((E,A), C) = \text{MIN}(\text{dist}(E,C), \text{dist}(A,C)) \\ = \text{MIN}(2, 2) \\ = 2$$

$$\rightarrow D((E,A), B) = \text{MIN}(D(E,B), D(A,B)) \\ = \text{MIN}(2, 5) \\ = 2$$

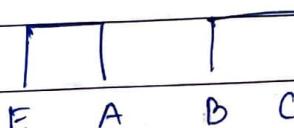
$$\rightarrow D((E,A), D) = \text{MIN}(D(E,D), D(A,D)) \\ = \text{MIN}(3, 3) \\ = 3$$

	E,A	C	B	D
E,A	0			
C	2	0		
B	2	1	0	
D	3	6	3	0

Step 2: Consider the distance matrix obtained in step 1
(given above)

∴ B,C distance is minimum, we combine B & C.

Dendrogram



Distance



Distance Matrix

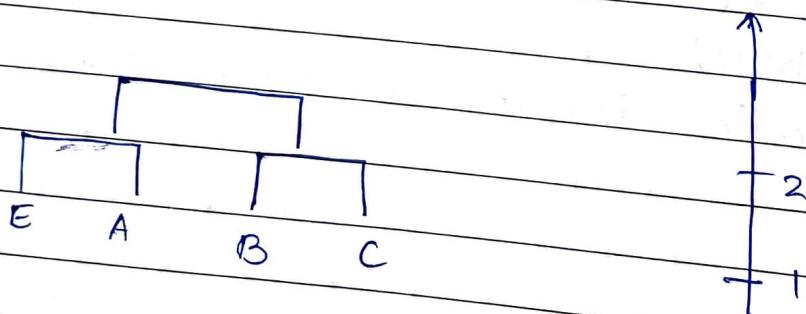
$$\begin{aligned}
 D(B,C), (EA)) &= \text{MIN}(D(BE), D(BA), D(C,E), D(CA)) \\
 &= \text{MIN}(2, 5, 2, 2) \\
 &= 2
 \end{aligned}$$

$$\begin{aligned}
 D(Bc), D) &= \text{MIN}(D(B,D), D(cD)) \\
 &= \text{MIN}(3, 6) \\
 &= 3
 \end{aligned}$$

	E,A	B,C	D
E,A	0		
B,C	2	0	
D	3	3	0

Step 3: Consider the distance matrix obtained in step 2
 (given above)
 ∵ (E,A) & (B,C) distance is minimum, we combine them.

Dendrogram

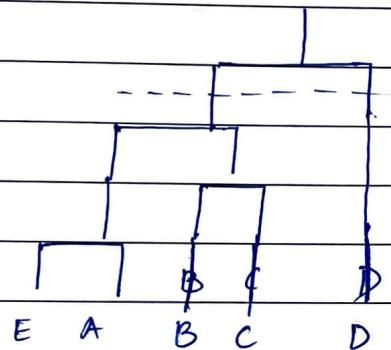


$$\begin{aligned}
 D(EA, BC, D) &= \text{MIN}(D(E,D), D(A,D), D(B,D), D(C,D)) \\
 &= \text{MIN}(3, 3, 3, 3) \\
 &= 3
 \end{aligned}$$

E, A, B, C D
 E, A, B, C 0
 D 3 0

Step 1 : Finally combine D with (EABC)

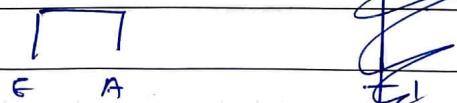
Final Dendrogram



(2)		E	A	C	B	D
E	0					
A	1	0				
C	2	2	0			
B	2	5	1	0		
D	3	3	6	3	0	

Solve the following eg using complete link Hierarchical clustering.

Soln: Step 1:



$$\begin{aligned}
 D((E, A), C) &= \text{MAX}(D(E, C), D(A, C)) \\
 &\rightarrow \text{MAX}(2, 2) \\
 &= 2
 \end{aligned}$$

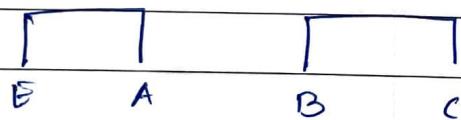
$$\begin{aligned}
 D((EA), B) &= \text{MAX}(D(E, B), D(A, B)) \\
 &= \text{MAX}(2, 5) \\
 &= 5
 \end{aligned}$$

$$\begin{aligned}
 D((EA), D) &= \text{MAX}(D(E, D), D(A, D)) \\
 &= \text{MAX}(3, 3) \\
 &= 3
 \end{aligned}$$

	E, A	C	B	D
E, A	0			
C	2	0		
B	5	1	0	
D	3	6	3	0

Step 2:

Dendogram



~~Distance~~

Distance Matrix

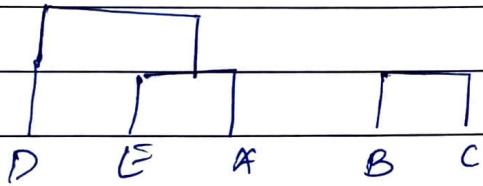
$$\begin{aligned}
 D(B, C), (EA) &= \text{MAX}(D(BE), D(BA), D(CE), D(CA)) \\
 &= \text{MAX}(2, 5, 2, 2) \\
 &= 5
 \end{aligned}$$

$$\begin{aligned}
 D(BC, D) &= \text{MAX}(D(B, D), D(C, D)) \\
 &\geq \text{MAX}(3, 6) \\
 &= 6
 \end{aligned}$$

	E, A	B, C	D	
E, A	0			
B, C	5	0		
D	3	6	0	

Step 3:

Dendrogram



$$\begin{aligned}
 D(EAD, BC) &= \text{MAX}(D(E, D), D(A, D), D(B, D), D(C, D)) \\
 &= \text{MAX}(2, 3, 6, 3) (2, 5, 6, 2, 2, 3) \\
 &= \text{MAX} = 6 //
 \end{aligned}$$

E, A, B, C	D
D	6

Step 4: Finally combine D with (E, A, B, C)

Final Dendrogram

