



# INSTITUTE OF AERONAUTICAL ENGINEERING

(Autonomous)

Dundigal - 500 043, Hyderabad, Telangana

## Examinations Control Office

**Examination**

B TECH VI SEMESTER END EXAMINATIONS REGULAR JUNE 2025 REG UG20

**Month & Year**

1-Jun

**Date**

20/06/2025

**Course Name**

DATA MINING AND KNOWLEDGE DISCOVERY

**Course Code**

ACIC01

**E-Code**

7918

---

### Instructions to Evaluators

- ❖ Evaluators should spend at least 3-5 minutes on one answer booklet during the evaluation.
- ❖ Evaluators should cross check that marks are allotted for all the attempted questions.
- ❖ The marks should be assigned fairly according to the mark distribution specified in the scheme of evaluation.
- ❖ For questions that were attempted incorrectly, evaluators are required to award zero marks.
- ❖ The evaluator must give a proper justification in case of any mistakes identified in the marks provided.

## START WRITING FROM HERE

Q.No.

1.a

### Types of Attributes in a dataset:

A dataset is a data source which consists of information / data about the specific class. So let us see the type of attributes found in a dataset.

#### (i) Nominal Attributes -

These nominal attribute have fixed values, like an enum.  
examples:-

1. User Roles = {user, admin, manager etc,}
2. WeatherStatus = {cloudy, sunny, windy, rain, etc}
3. ReviewStatus = {pending, approved, denied}

Here as we see in the above example there are fixed values that can be in the Nominal Attributes. They are mainly useful in classifications.

#### (ii) Ordinal Attributes -

The ordinal attributes are ordered values. They specify the order in the dataset.

Q.No.

Examples :-

1. Marks Grades = { A+, A, B+, B, C+, C etc }  
 $A^+ > A > B^+ > B > C^+ > C$  etc,
2. TimeStamp = { 11:00PM, 12:00PM, 9:00AM } etc

Interval Attributes :-

Interval Attribute are order attribute with no "real zero", And each two consigutive values have a fixed interval between them.

Examples :-

1. Temperature = { 0°C, 5°C, 1°C, 100°C etc. }
2. Time { 0:00AM, 12:00PM, 6:50AM etc }

Ratio Attributes :-

Ratio attribute are order and have real zeros. They could be numeric value, or can be decimal value.

Examples :-

1. Marks = { 0, 10, 12, 5, 15, 21, 25 etc. }
2. Speed = { 0.0, 0.01, 0.51, 0.65, 1.51 etc }  
 (m/sec)





Q.No.

## Binary Attributes:-

As the name specifies the Binary attributes consist only Binary values (Yes or No) or (True or False).

Examples:-

1. isEligible = { Yes or No }

2. isPassed = { True / False }

Now, let us see a dataset above attribute types

index	UserId	Role	Grade	Marks	is Eligible
1	001	Student	A <sup>+</sup>	95	T
2	012	Teacher	A <sup>+</sup>	100	F
3	156	Student	C	61	F
4	119	Student	B <sup>+</sup>	73	T
5	017	Student	A	85	T
(Ordinal)	(Interval)	cardinal (Nominal)	cardinal	(Ratio)	(Binary)



Q.No.

1.b

Given,

Dataset of patient records.

AIM,

Build a predictive Model to predict patient readmission rates.

Process of Building a predictive Model :-

1. Data Cleaning :-

Dataset need to be cleaned by identifying the missing values and managing the outlier in the data. This will help

- Reduce Noise
- Improve data quality
- Minimizes Data Redundancy.

2. Data Integration :-

In this process, the data from various sources is integrated in to a single dataset.. This will help to decentralize the data, as the data is retrieved from various spot.



Q.No.

### 3. Data Transformation :-

In this step, the data is transformed into the useful format to make the process easy to build the model. Here techniques like Smoothing, Normalization and Text-vectorization are used.

### 4. Pattern Detection :-

By using the data mining algorithms the pattern that are interesting and relevant to the model are taken and irrelevant data is discarded.

### 5. Pattern Evaluation :-

Model evaluation metrics are used (F2-score etc.) to check the accuracy and efficiency of the model.

### 6. Model Representation :-

The frontend tools like, dashboards, Graphs are used to visualize the model and use it to predict the future events.



Q.No.

So, I would employ these following data mining techniques,

- (i) data cleaning
- (ii) data integration.
- (iii) data transformation.
- (iv) pattern recognition.
- (v) Model Evaluation
- (vi) and model representation,

to predict the patient readmission probabilities based on the dataset provided by the health-care Organization.

Q.No.

2.6

Given,

The value of the income attribute,

2000, 3000, 4000, 6000, 10,000.

Now, let us perform the Normalization methods specified.

(1) Do min-max - Normalization :-

In min-max normalization the data will be transformed between '0' - '1'

$$\text{minmaxNorm} = \frac{\text{item} - (\text{minimum\_value})}{(\text{max\_val}) - (\text{min\_value})}$$

(or) min\_value = 2000

max\_value = 10,000

$$\bullet \text{ for } 2000 \Rightarrow \frac{2000 - 2000}{10000 - 2000} = \frac{0}{8000}$$

$\Rightarrow 0$

$$\bullet \text{ for } 3000 \Rightarrow \frac{3000 - 2000}{10000 - 2000} = \frac{1000}{8000}$$

$\Rightarrow 0.125$



Q.No.

$$\bullet \text{ for } 4000 \rightarrow \frac{4000 - 2000}{10000 - 2000} = \frac{2000}{8000}$$

$$\Rightarrow 0.25$$

$$\bullet \text{ for } 6000 \rightarrow \frac{6000 - 2000}{10000 - 2000} = \frac{4000}{8000}$$

$$\Rightarrow 0.5$$

$$\bullet \text{ for } 10,000 \rightarrow \frac{10,000 - 2000}{10,000 - 2000} = \frac{8000}{8000}$$

$$\Rightarrow 1.0$$

$\therefore$  Min Max Normalized values  
 $\{0, 0.125, 0.25, 0.5, 1.0\}$

(ii) Z-Score Normalization

$$Z\text{-Score} = \frac{\text{item} - \text{mean}}{\text{S.D}}$$

S.D  $\rightarrow$  Standard deviation

$$\text{Mean} = \frac{\sum x_i}{n} =$$

$$= \frac{2000 + 3000 + 4000 + 6000 + 10,000}{5}$$



Q.No.

$$= \frac{5}{25000}$$

$$\therefore \bar{x} = 5000$$

$$\therefore \text{mean} = 5000$$

$$SD = \sqrt{(\sum (x_i - \bar{x})^2) / n}$$

$$= \sqrt{[(2000-5000)^2 + (3000-5000)^2 + (4000-5000)^2 + (6000-5000)^2 + (10000-5000)^2] / 5}$$

$$\Rightarrow \sqrt{\frac{3000^2 + 2000^2 + 1000^2 + 1000^2 + 5000^2}{5}}$$

$$\Rightarrow \sqrt{\frac{(3 + 2 + 1 + 1 + 5) \times (1000)^2}{5}}$$

$$= \frac{12000}{\sqrt{5}} = 536.6$$

$$SD = 536.6$$

$$a) z(2000) = \frac{(2000 - 5000)}{536.6} = \frac{3000}{536.6} = -5.6$$

$$b) z(3000) = \frac{(3000 - 5000)}{536.6} = \frac{2000}{536.6} = -3.72$$

$$c) z(4000) = \frac{(4000 - 5000)}{536.6} = \frac{1000}{536.6} = -1.86$$

$$d) z(6000) = \frac{(6000 - 5000)}{536.6} = \frac{1000}{536.6} = 1.86$$

$$e) z(10,000) = \frac{(10000 - 5000)}{536.6} = \frac{5000}{536.6} = 9.317$$

Z-Score Normalized value

$$\{-5.6, -3.72, -1.86, 1.86, 9.317\}$$

Q.No.

(iii) Decimal Scaling

$$\text{Scaled} = \frac{\text{item}}{10^i}$$

where "i" is the (n-1)

$$n = 5, i = 4$$

$$\text{Scale}(2000) = \frac{2000}{10000} = 0.2$$

$$\text{Scale}(3000) = \frac{3000}{10000} = 0.3$$

$$\text{Scale}(4000) = \frac{4000}{10000} = 0.4$$

$$\text{Scale}(6000) = \frac{6000}{10000} = 0.6$$

$$\text{Scale}(10,000) = \frac{10000}{10000} = 1.0$$

Decimal Scaled Values

{ 0.2, 0.3, 0.4, 0.6, 1.0 }

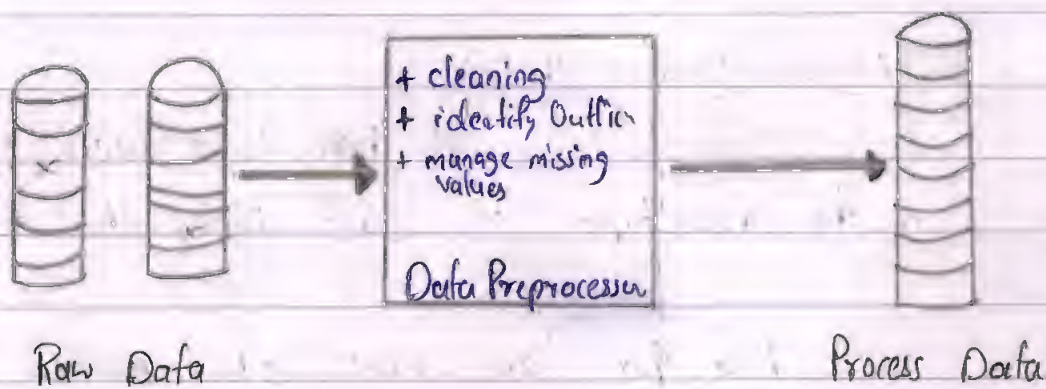


Q.No.

2.a

## Data Cleaning :-

Data cleaning is the important step in data preprocessing. Data cleaning involves reducing noisy data in the data set while identifying the Outliers. Many raw data-sets have inconsistent data with many missing values and redundant information. Here in data cleaning such missing and inconsistent data is handled. Data cleaning will improve the quality of the data and reduce the redundancy.



## Steps and Cleaning tools :-

- \* Perform description analysis using tools like pandas in python
- \* Find the missing values in the dataset `df.dropna().sum()`
- \* Fill the missing values, `df.dropna() → df.fillna()`

we can also specify the method

Q.No.

## Handling Missing Values:-

Missing values cause the noise in the data, so they need to be handled.

### (i) Replacing with Mean:-

We replace the missing values with mean of the total values in the column.

### (ii) Replace with a Constant:-

We will calculate a constant value and put in the place of the missing value.

### (iii) Remove (ignore) the missing:-

If there is a value missing in the row (tuple) we just ignore the tuple.

And there are many other methods to handle the missing values in the data set like medians and "Predicting Model to predict the missing values".





Q.No.

### 3.a. Online Analytic Processing (OLAP) :-

Online Analytic processing servers are used as a bridge between the data-warehouses and the frontend tools. The OLAP servers, will perform the analytical operations on the data stored in the Storage tier of the architecture where Data Marts, Metadata repositories and data-warehouse are located. And the analyzed information is ready to be used by the frontend tools to visualize the Information.

There are three Online Analytic processing server architectures, they are

#### (i) Relational Online Analytic Processing (ROLAP) :-

The ROLAP servers are built upon the Relational database System, which consist of the Normalized data, like Tables, attribute, and relations. And high level complex SQL (Structured query language) queries are used to perform the operations, on the ROLAP Server.



Q.No.

~~ROLAP~~ ROLAP Advantages -

- Easy to implement
- Suitable for small data warehouses

ROLAP DisAdvantages -

- Inefficient for large data
- Complex queries of SQL

(ii) MultiDimensional Online Analytic Processing (MOLAP):-

The MOLAP servers are flexible and used for operation on Multi-dimensional data querying. Like data stored in the form of Cubes (data cubes) and Data Mining Query Language (DMQL) are used to communicate with the gateway.

MOLAP Advantages:-

- ~~Hard to use~~ easy to use
- Suitable for more data.

MOLAP DisAdvantages -

- Inefficient for small data
- And need good understanding of DMQL.



Q.No.

(iii) Hybrid Online Analytic Processing - (HOLAP) :-

HOLAP uses both the features of ROLAP and MOLAP, i.e., the data can be in the normalized format like tables and relations, or data could be multi-dimensional and stored in the data-cubes.

Advantages of HOLAP :-

→ Effective and efficient for both small and heavy data, and provide accurate results.

DisAdvantages of HOLAP :-

→ Complex to implement and use.

→ Need experience to maintain the server.

\* The OLAP server can perform the following operation

→ Roll up (Top to Bottom to up)

→ Drill down (Top to Bottom)

→ Slice and Dice (Single dimension and sub dimension)

→ Pivot (Rotate for a new perspective)



Q.No.

3.b

Given,

The data warehouse consists of three dimensions.

- Time
- Doctor
- Patient.

and, two Measures

+ Count

+ charge

(i) Classes of schema that are popularly used for modelling data warehouse:-

(a) Star Schema :-

In star schema there is one fact table and multiple dimension tables referring to the fact-Table.

Fact-Table :- The table that stores the measurable values

Dimension-Table :- The table that stores the descriptive data





Q.No.

(b) Snow Flake Schema :-

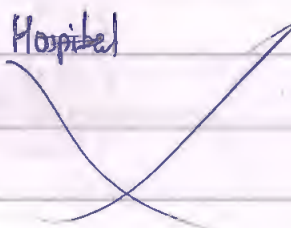
Snow flake schema is an extension of the star schema, that have multiple sub dimension connection to the existing dimension with only one single fact table.

(c) Galaxy Schema :-

The galaxy or fact constellation schema is the complex merge of multiple snowflake star schemas. In the Galaxy schema there are two or more fact tables that share the same dimension tables.

(ii) Star schema :-

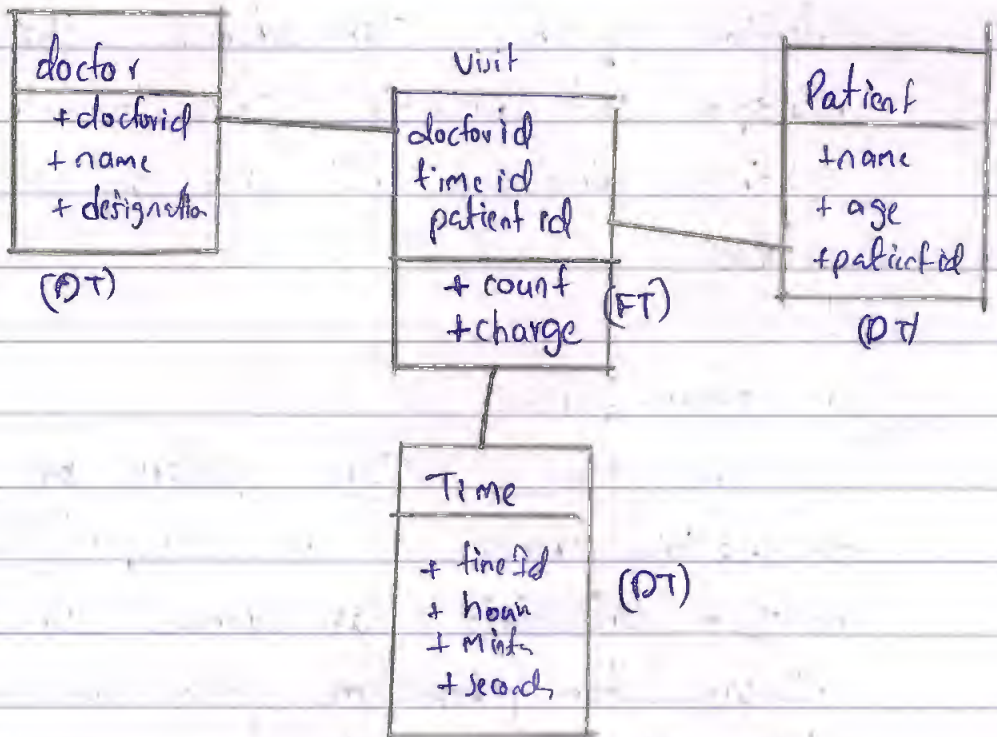
Patient



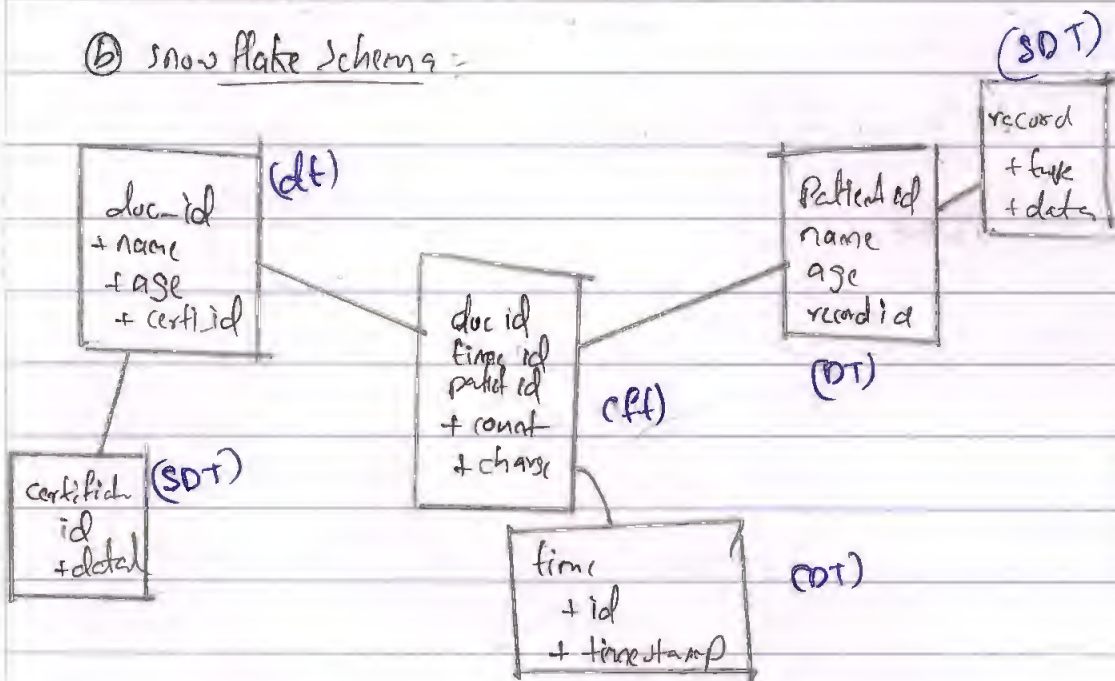


Q.No.

(ii) @ Star - schema :-

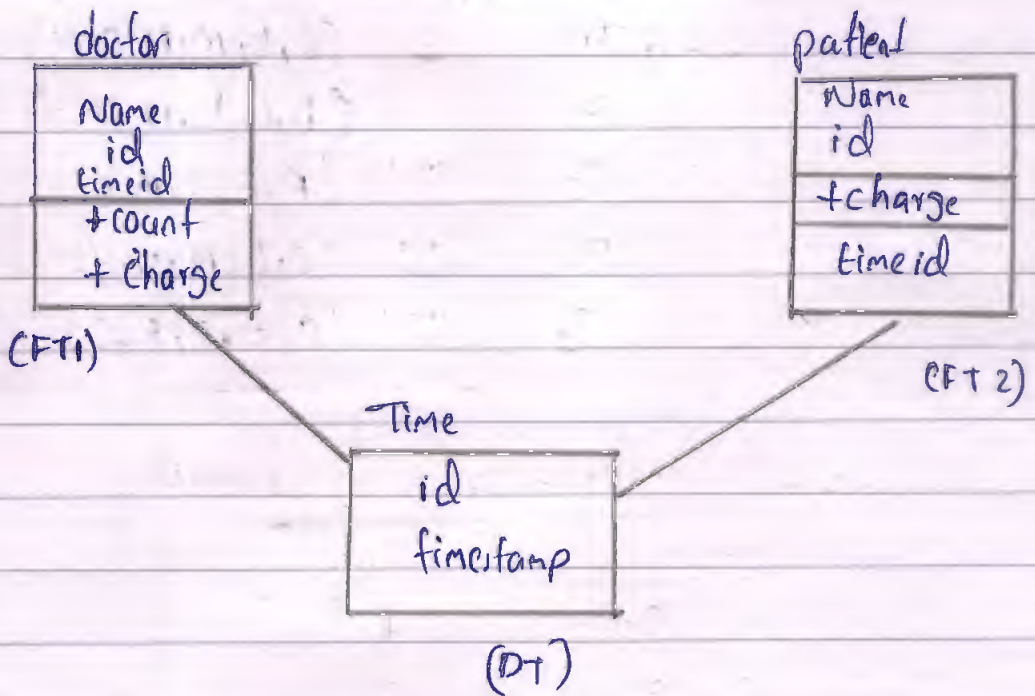


(b) snowflake schema :-



Q.No.

© Galaxy or fact constellation schema:-







Q.No.

5.b Construction of  $\lambda$ -frec for given transaction.

Given,

Minimum Support Threshold (MST) = 3.

Transaction id	item
T <sub>1</sub>	{E, K, M, N, O, Y}
T <sub>2</sub>	{D, E, K, N, O, Y}
T <sub>3</sub>	{A, E, K, M}
T <sub>4</sub>	{C, K, M, U, Y}
T <sub>5</sub>	{C, E, I, K, O}

Step 1  $\Rightarrow$  write frequency of each item

item	count
A (x)	1
C (x)	2
D (x)	1
E	4
I (x)	0
K	5
M	3
N (x)	2
O	3
U (x)	1
Y	3

remove item with  
count less than MST  
= 3

Q.No.

Step 2 write priority

item	count	priority
E	4	2
K	5	1
M	3	3
O	3	4
Y	3	5

$\Rightarrow [K > E > M > O > Y]$

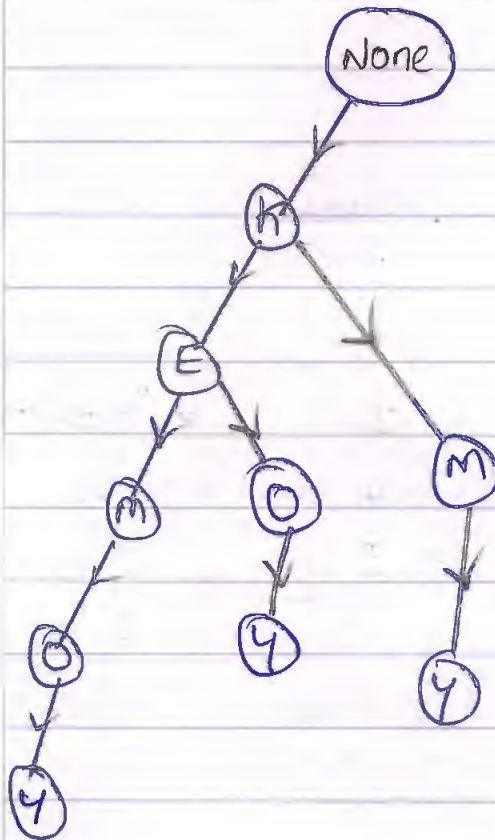
Step 3 Write Transaction in order of priority of item

id	items in Order
T <sub>1</sub>	K, E, M, O, Y
T <sub>2</sub>	K, E, O, Y
T <sub>3</sub>	K, E, M
T <sub>4</sub>	K, M, Y
T <sub>5</sub>	K, E, O

Step 4 :

Now, our next step is to draw the  
Fp-free.

Q.No.



fp-tree

Frequent pattern :-

(i) None  $\rightarrow$  K  $\rightarrow$  M  $\rightarrow$  Y

(ii) None  $\rightarrow$  K  $\rightarrow$  E  $\rightarrow$  O  $\rightarrow$  Y

(iii) None  $\rightarrow$  K  $\rightarrow$  E  $\rightarrow$  M  $\rightarrow$  O  $\rightarrow$  Y



Q.No.

7.b. Given data point

$A_1(2,10)$ ,  $A_2(2,5)$ ,  $A_3(8,4)$ ,

$B_1(5,8)$ ,  $B_2(7,5)$ ,  $B_3(6,4)$

$C_1(1,2)$ ,  $C_2(4,9)$ ,  $C_3(4,9)$ .

and initial centroid  $A_1, B_1, C_1$

applying K-mean algorithm

Data points	$C_1(2,10)$	$(5,8)$	$(1,2)$	cluster
$A_1(2,10)$	0	3.60	5.1	1
$A_2(2,5)$	5	4.25	2.1	3
$A_3(8,4)$	2.6	5	6.5	2
$B_1(5,8)$	4.7	0	4.81	2
$B_2(7,5)$	7.5	3.6	6.2	2
$B_3(6,4)$	4.13	4.12	5.38	2
$C_1(1,2)$	1.8	7.21	0	3
$C_2(4,9)$	3.8	1.41	7.61	2
<del><math>C_3(4,9)</math></del>				

New mid point (centroid)

$C_1(2,10)$ ,  $(6,6.5)$ ,  $(3,7)$

Q.No.

data point	(2,10)	(6,6.5)	(3,7)	new	
(2,10)	0	4.1	2.6	1	1
(2,5)	5	4.2	1.8	3	3
(8,5)	8.6	2.8	5.8	2	2
(5,8)	5.7	1.6	2.2	2	2
(7,5)	7.5	2.1	3.2	2	2
(6,5)	5.13	1.6	3.2	2	2
(1,2)	1.8	5.8	(2,1)	1	3
(9,9)	3.8	1.56	(1,1)	3	2

repeat the step until new clste = old clste

we get,

$$\text{cluster 1} = \{A_1, C_1\}$$

$$\text{cluster 2} = \{A_2, B_1, B_2, B_3\}$$

$$\text{cluster 3} = \{A_2, C_2\}$$

Q.No.

7.4 Different types of data used in the cluster Analysis:-

Cluster Analysis is the process of dividing the data into different clusters based on the methods like, partitioning, hierarchy based, Grid and Model based.

The data is unsupervised and the clusters are detected based on the above methods. And Cluster Analysis can be done on any type of data type like

- Nominal data, used for simple clustering based on classified values.
- Ordinal data, Order of the data without real numerical value.
- Binary data like yes or no.
- Text data, we can perform cluster analysis on the text data also, too.

Many applications use to cluster data like Meta, Google.

- Image, also can be clustered into groups based on their feature.



Q.No.

## Typical Requirement of Clustering:-

- \* Data need to be preprocess and cleaned with out noise and missing.
- \* Data should be formatted like text to vector, and Images to matrix for easy clustering.
- \* Data should be unlabeled to perform clustering. if data is labeled then classification is best option.



Q.No.



Q.No.	







Q.No.	



Q.No.	





Q.No.	



Q.No.

$$Z_{score} = \frac{x_i - \mu}{SD}$$

$$SD = \sqrt{\frac{\sum (x_i - \mu)^2}{n}}$$

$$SD = \sqrt{\frac{(x_1 - \mu)^2 + (x_2 - \mu)^2 + \dots}{n}}$$

Part - Kne, Kna, Cl  
 Hi - Brich, cort  
 Dwell - 1030, 1040  
 GPR - 1050, 1060  
 GPR - 1070

GPR

- water
- order
- input
- bottle
- Binay

var  
 GPR  
 Celine  
 Binay  
 1070

