

Data Mining :-

①

- ⇒ DM is the process of Extracting knowledge or insights from large amount of data using various statistical, & computational techniques.
- ⇒ The data can be structured, semi-structured or unstructured, and can be stored in various forms such as databases, data warehouses, and data lakes.
- ⇒ The primary goal of data mining is to discover hidden patterns and relationships in the data that can be used to make informed decisions or predictions.
- ⇒ This involves exploring the data using various techniques such as clustering, classification, regression analysis, association rule mining, and anomaly detection.
- ⇒ Data mining has a wide range of applications across various industries, including marketing, finance, health care and telecommunications.
 - eg:- In marketing data mining can be used to identify customer segmentation and target marketing campaigns.
 - In healthcare it can be used to identify risk factors for diseases and develop personalized treatment plans.
- ⇒ Database Data:- A db system also called database management system (DBMS), consists of a collection of interrelated data, known as database, and set of programs to manage and access the data.

- ⇒ A relational database is a collection of tables, each of which is assigned a unique name. Each table consists of attributes (columns or fields) and usually stores a large number of tuples (records or rows).
- ⇒ Each tuple in a relational table represents an object identified by a unique key and described by a set of attribute values.
- * ⇒ A semantic data model such as entity-relationship (ER) data model, is often constructed for relational databases. An ER model represents the database as a set of entities and their relationships.
- eg: Show me a list of all items that were sold in the last quarter
 Here Relational languages also use aggregate functions such as sum, avg, count, max, min,
- * Data Warehouse ⇒ A DW is a repository of information collected from multiple sources.
 - ⇒ DW are constructed via a process of data cleaning, data integration, data transformation, data loading periodic of data refreshing.
- ⇒ A data warehouse is usually modeled by a multi-dimensional data structure called a data cube, in which each dimension corresponds to an attribute or a set of attributes in the schema.

(4)

~~for ex. putting money in ATM~~

Transactional data: - If database captures a transaction, such as a customer's purchase, a flight booking, or a user's click on a web page.

A transactional database may have additional tables, which contain other information related to the transaction, such as item description, information about the salesperson or the branch, and so on.

e.g: "which items sold well together" This kind of j. market basket data analysis

Trans-ID	Items sold
T100	I1, I2, I3, I4, I5
T200	I2, I3
T300	-
-	-

- * Descriptive :- Descriptive mining tasks characterize properties of the data in a target dataset.
- * Predictive :- predictive mining tasks perform induction on the current data in order to make predictions

\rightarrow Basic statistical Description of data :-

The basic statistical description of data can be used to the properties of data.

* Measuring the Central Tendency :- (Mean, Median, Mode)

at Mean:- Mean is most common & effective numeric measure which is used to measure the center of data.

$$\text{Mean} = \bar{x} = \frac{\sum_{i=1}^N x_i}{N}$$

$$= \frac{x_1 + x_2 + \dots + x_N}{N}$$

Mean of salary

center of the hole data.

if weight are associated

$$\bar{x} = \frac{\sum_{i=1}^N w_i x_i}{\sum_{i=1}^N w_i} = \frac{w_1 x_1 + w_2 x_2 + w_3 x_3 + \dots + w_N x_N}{w_1 + w_2 + w_3 + \dots + w_N}$$

This is called weighted arithmetic mean or the weighted average.

b) Median : Median is the best measure for finding the center of the data.

e.g.- 40, 45, 47, 50, 54, 52, 57, 61, 64, 71, 110.

$$\text{center} = \frac{52+56}{2} = \frac{108}{2} = 54 \text{ center position / middle position}$$

c) Mode :- another measure for central tendency is mode.

d) Mean - Mode $\approx 3(X \text{ (mean)} - \text{median})$

e) Mid Range : To measure the central tendency of numeric data set arrange of largest & smallest

$$\text{in which } = \frac{30+110}{2} = \frac{140}{2} = 70.$$

attribute or a set of attributes.

Measuring the Dispersion of data

(3)

Range :- The difference b/w largest & smallest.

split the data distribution
Quantile o^o equal size consecutive sets

Quartile :- one fourth of data distribution

Percentile :- 100 quantities are more commonly referred to as per-

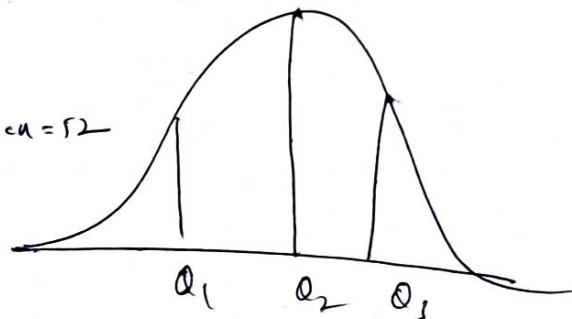
Interquartile range (IQR)

The distance between the first & third quartiles is a single measure of spread that gives the range covered by the middle half of the data. This distance is called the IQR

$$\boxed{IQR = Q_3 - Q_1}$$

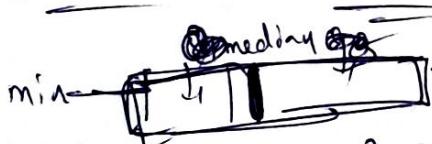
$$Q_1 = 42 \quad Q_3 = 65 \text{ & Median} = 52$$

$$IQR = 65 - 42 \\ = 13$$



* Five Number Summary

25th %, Median & 75th %



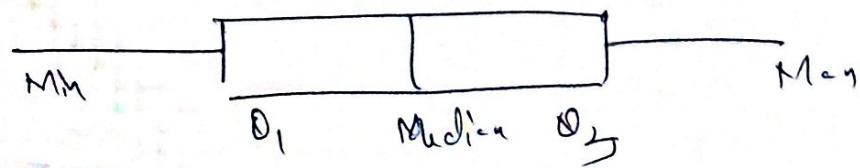
To identify the middle and what next value will come.

It is also called boxplot.

* Five Number summary % of distribution counts of the Median (Q₂) The quartiles Q₁ and Q₃, and the smallest and largest

individual observations, written in the order of

Minimum, Q₁, Median, Q₃, Maximum



The variance of N observations, x_1, x_2, \dots, x_N

Variance is the variance of attribute x is

for numeric attribute x

$$\sigma^2 = \frac{1}{N} \sum_{i=1}^N (x_i - \bar{x})^2$$

\bar{x} is the Mean value

$$= \left(\frac{1}{N} \sum_{i=1}^N x_i^2 \right) - \bar{x}^2$$

standard deviation σ of the observations is the square root of the variance, σ^2

$$\sigma \approx \sqrt{\sigma^2}$$

⇒ Data Mining Functionalities:

(4)

* 5 functionalities

1 → Concept / class (definition) / descriptions:

Data is always associated with class / concepts

⇒ class / concepts can be done in '2' ways

↳ Data characterisation:

↳ Data description:

↳ Data characterisation:

Refers to the summary of the class / concept about which we are studying.

↳ → General overview.

2 → Data Description:

Compares the common features of the classes

Here we have '2' classes / concept's then it will compare the common features which are in between two classes or those two concepts. If we have any changes in the two classes (or concept) it will observe.

↳ : Bar chart, curves, graph etc.

2 → Mining frequent patterns, associations & correlations:

Mining : Data Mining is the process of sifting through large sets to identify patterns and relationships that can help solve the business problems through data analysis.

Anomaly : Something different, abnormal, peculiar, or not easily classified.

Frequent patterns: Things which are found most commonly in data items / objects

There are many ^{kind of} frequent patterns are there.

1. "frequent itemsets", 2) "frequent subsequences", & 3) "frequent subsequences":
Refers to a set of items that often appear together in a transactional data set.
Eg: Milk & bread. (which are frequently bought together in grocery stores by many customers.)
2. Frequent subsequences:
Frequently occurring subsequences, such as pattern that customers tend to purchase first laptop followed by a digital camera, and then memory card is a frequent sequential pattern.
3. Frequent substructure:
F-S can refer to different structure forms that may combined with itemsets or subsequences.
(Eg:- Graphs, trees, or lattices)

"Association Analysis": \rightarrow [relationship]

* It is a way of identifying the relationship between various items.

Eg:- Used to determine sales of items that are frequently purchased together.

How the data items are related to each other what is association between the data items.

$x \rightarrow$ clay pot $y \rightarrow$ kerosene lamp

→ Dry fruits are frequently purchased together items, next purchase
 together what are the items that are purchased together
Dry fruits & Chocolates
Dry fruits
 frequently
 ↓
 one data item
 One set
 together
 ↓
 Another data item
 another set

Correlation Analysis:

- It is a mathematical techniques.
- Shows how strongly pair of attributes are related together.
- Shows how strongly pair of attributes are related together.
- e.g. Tall people tend to have more weight

(1) This is one attribute

This is another attribute

Two attributes are strongly related to each other; & how they are related to each other.

Classification and regression for predictive analysis :-

Predicting the data → (Assuming the data)
 if data is missing we have to fill it or assume it.

Classification is process of finding a model that distinguishes data items

e.g. decision

Chocolate \in Dry Milk one group

Kit Kat second group

Two different types of chocolates will be classifying them. → with on their what you do name

$$S = \frac{I_{\text{bin}}}{\text{Total}}$$

→ Normalized char.

$$\frac{\text{Kit Kat}}{1} \text{ or } \frac{\text{Dry Milk}}{1}$$

- Here data items are classifying in to the names on based on the color of the wrapper
- Here what is the model? Wrapper ^{color} is the model your finding some model in order to distinguish the data items that you are having. This is the benefit.
- * "Decision tree" is a flowchart-like tree structure. based on decision tree we can easily classifying the data.

Regression

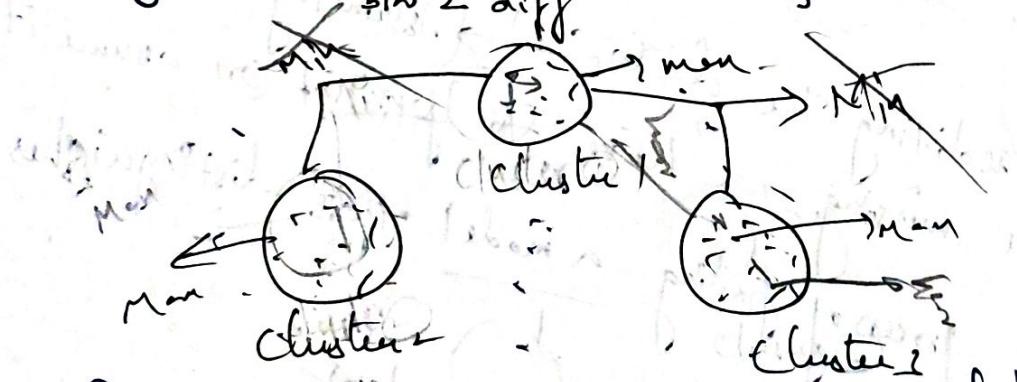
Statistical methodology that is used for numeric prediction of missing data.

$$\text{Ex: } 2, 3, 4, 6, 7, 8, 6, 11, 12 \rightarrow 15, 16, 17$$

→ Predicting the data

[done based on previous data]

4. Cluster Analysis: The data items are clustered based on the principle of maximising the ^{within the same} intraclass similarity and minimising the interclass similarity:



- Data represents the data so there data items are classified in to three clusters based on the similarity between those things.
- So whatever data sets are inside this cluster one. Whatever the data items are there all those data

items are similar to each.



- ~~1) Inside the cluster~~ the similarity between the items which are inside this cluster is maximum, and the items which are inside this cluster are also maximum.
- 2) But the similarity between the items which are inside this cluster should be minimum than what minimizing the inter class between two classes or between two clusters.

* So within the cluster the similarity between the objects should be maximum.

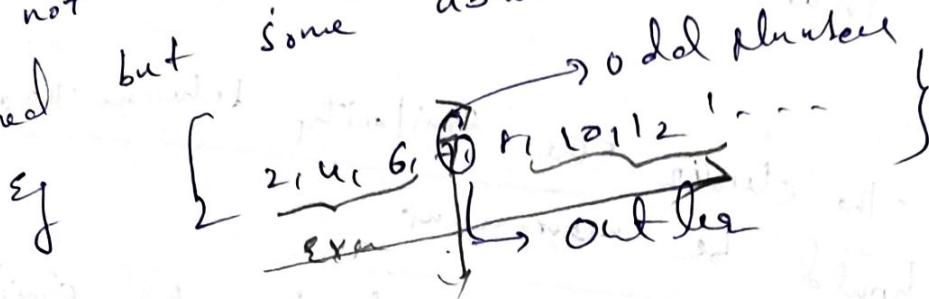
- 3) And between two different clusters the similarity between the objects should be minimum.
- 4) Analysis of these clusters is called as cluster analysis.
- 5) "Analysis of these clusters" is called as cluster analysis.

Market basket
Association
Outlier
Cluster

5- Outlier Analysis (anomaly mining) (Q)

⇒ Among the data items in a db, there may be some items which do not follow the general behaviour of data.

Those data items → Outlier [noise / exception]
or Anomaly:- which do not obey the rules & which are not according to the rule which we have defined but some abnormalities we can say like that



Proximity Measures for Binary Attributes

(7)

Let i, j are two objects and p total number of attributes. $(0, 1)$ either F or T

* Symmetric Binary Attribute:

e.g.: $M/F \rightarrow$ equal importance.

[All object's]

* Asymmetric Binary Attribute: "Unequal importance"

More important for one / other one if not much import
rule:

(M/F)

$$\sum p = q + r + s + t$$

similarity calculation

$q - 1$ for both objects.

$r - 1$ for $i, 0$ for j } same.

$s - 0$ for $i, 1$ for j

$t - 0$ for both i and j

$\leftarrow N/p$
+ve
else
different

first

		Obj j	Sum	
Obj i	1	1. 0.	same	$= 0, 0, 0, 0$
		2. 0. (S) t	$q+r$	$- 0, 0, 0, 0$
		sum $q+r+s+t$	$p \rightarrow$ Total No. of Objects	

(Contingency Matrix)

Dissimilarity for symmetric binary attribute

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

$$\rightarrow 1 - \frac{m}{p} = \text{similarity}$$

Dissimilarity for asymmetric binary attribute (all values are not binary here)

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

$\left. \begin{array}{l} \text{Asymmetric binary attr} \\ r > t \end{array} \right\} \begin{array}{l} p + \text{highlighted} \\ - N \text{抱着} \end{array}$

't' both are negative importance [N]

$$\rightarrow t[0, 0] \uparrow \text{ignored.}$$

\Rightarrow Then Asymmetric Binary Similarity.

$$\text{Sim}(i,j) = 1 - \frac{r+s}{\sqrt{r+s+t}}$$

$$= \frac{\sqrt{r+s+t} - r - s}{\sqrt{r+s+t}}$$

$$\text{Sim}_{\text{jaccard}}(i,j) = \frac{r}{\sqrt{r+s+t}}$$

let the values

Eg: $r \& p.$ to be set to 1, 'N' & 'negative' set to 0

refer patient table

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

dissimilarity for asymmetric A-4

$0 \rightarrow 1, 0$

$s \rightarrow 0, 1$
n y

$t \rightarrow 11$
yy

$$\text{dis(jack, jim)} = \frac{1+1+0}{1+1+0}$$

$$= \frac{2}{2} = 0.67$$

0.23

$$D_3 (\text{Jim, Mary}) = \frac{\alpha_{ij}}{\alpha_{\text{total}}} = \frac{1+(+) \textcircled{8}}{1+1+1+1} \cdot \frac{1}{4} = 0.75$$

$$D_4 (\text{Jack, Mary}) = \frac{\alpha_{ij}}{\alpha_{\text{total}}} = \frac{1+1}{1+1+1+1} = \frac{1}{3} = 0.33$$

\rightarrow Jack & Mary \rightarrow same disease
 out of three patients, Jack & Mary are the most likely to have a similar disease.
 because distance is less

\rightarrow Proximity analysis for ordinal variables:

Here we use ranking.

Name	Rank	α_{ij}	Z_{ijf}	$\alpha_{ij} \cdot \frac{(n-1)}{3}$
Jack	Excellent	4	$(4-1)/3 = 1$	
Mary	Better	3	$(2-1)/3 = 0.67$	
Jim	Cool	2	$(2-1)/3 = 0.33$	
Pat	Enough	1	$(1-1)/3 = 0$	

Total count M_f = 4 no. of entries (rankings)

$$Z_{ijf} = \alpha_{ij} \cdot \frac{(M_f - 1)}{3}$$

$$(\alpha_{ij} - 1)/3$$

$$d(i,j) =$$

$$|x_{i1} - x_{j1}| + |x_{i2} - x_{j2}|$$

$$+ \dots + |x_{ip} - x_{jp}|$$

Mannheim distance

$$= 4-1 = 3$$

	Jack	Mary	Jim	Pat	
Jack	0				$\rightarrow 4-1 = 3$
Mary	0.33	0			$= 4-0.33 = 3.67$
Jim	0.67	0.33	0		$\frac{3-1}{3} = \frac{2}{3} = 0.67$
Pat	1	0.67	0.33	0	$\frac{2-1}{3} = \frac{1}{3} = 0.33$
→ gantik Na	0.67	0.33	0.33	0.67	1

Dissimilarity for Attributes of Mixed types:-

	Test I Nominal	Test II Ordinal	Test III Numeric
1	Code A	Excellent	45
2	Code B	Fair	22
3	Code C	Good	64
4	Code D	Excellent	28

[Nominal
ratio]

or Nominal $\frac{P-M}{P}$ =

Dissimilarity measure = $d(i,j) = \frac{P-M}{P}$

$M = \text{No. of matches}$
 $P = \text{No. of variables}$.

$P=1$

	CA	CB	CC	CD
CA	0			
CB		0		
CC			0	
CD	0		1	0

	Ordinal attribute	Nominal	ordinal	Numeric
1	CA	S		45
2	CB	F		22
3	CC	G		64
4	CD	E		28

$$M_f = 3 \{ s_i, F_i, hood \}$$

$$2if = (\gamma_{if} - 1) / M_f - 1$$

TII	rank	rif	qif
1	SA	3	$(2-1)/2 = 1$
2	Fair	1	$(P-1)/2 = 0$
3	hood	2	$(2-1)/2 = 0.5$

	Code A(1)	Code B(1)	Code C(1)	Code D(1)
Code A(1)	0			
Code B(1)		1	0	
Code C(1)	0.5	0.5	0	
Code D(1)	0	1	0.5	0

	T Nominal	T4 order	Total Nm
1	LAI	S	97
2	CD	F	22
3	CC	Ho	64
4	CD	S	25

Markov.

Code1	0	1		
Code2	23	0		
Code3	19	42	0	
Code4	17	6	26	0

(9)

$$Min = 0$$

$$Max = 42$$

Divide by 42 for
normalization.

Data Mining

(10)

(1)

Similarity and dissimilarity :-

- * More important - used in clustering, some classification, anomaly detection.

- * Similarity :- Numerical measures of how alike two data objects are.
- i) Higher when objects are more alike.
- ii) Is highest when objects are more alike.
- iii) often falls in the range [0,1]

- * Dissimilarity :- Numerical measure of how different are two data objects.
- i) Lower when objects are more alike.
- ii) Minimum dissimilarity of often 0.
- iii) Upper limit varies.

→ Proximity refers to a similarity or dissimilarity.

Sim | dis for objects with single attribute.

e.g. are the attribute values for two data objects

Dissimilarity.

similarity

Attribute type

Nominal

$$d = \begin{cases} 0 & \text{if } p=q \\ 1 & \text{if } p \neq q \end{cases}$$

$$s = \begin{cases} 1 & \text{if } p=q \\ 0 & \text{if } p \neq q \end{cases}$$

ordinal

$$d = \frac{|p-q|}{(n-1)}$$

$$s = 1 - \frac{|p-q|}{n-1}$$

Interval or ratio

$$d = (p-q)$$

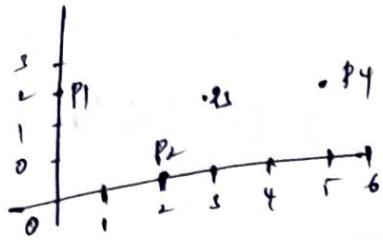
$$s = -d, s = \frac{1}{1+d} (0)$$

$$s = 1 - \frac{d - \min(d)}{\max(d) - \min(d)}$$

and area

Euclidean Distance

$$\text{dist} = \sqrt{\sum_{k=1}^n (P_k - Q_k)^2}$$



Point	x	y
P ₁	0	2
P ₂	2	0
P ₃	3	1
P ₄	5	1

(P₁, P₂)
(P₁, P₃)

$$(P_1, P_2) =$$

$$\sqrt{(x_1 - x_2)^2 + (y_1 - y_2)^2}$$

$$(0-5)^2 + (2-1)^2$$

$$= \sqrt{25+1}$$

$$= \sqrt{26}$$

$$\sqrt{(0-2)^2 + (2-0)^2} = \sqrt{(2)^2 + (1)^2} = \sqrt{4+1} = \sqrt{5}$$

$$\sqrt{26}$$

$$= 2.828$$

$$(P_1, P_3) = \sqrt{(2-3)^2 + (0-1)^2} = \sqrt{(1)+(1)} = \sqrt{2} =$$

$$= \sqrt{(0-3)^2 + (2-1)^2} = \sqrt{9+1} = \sqrt{10} = 3.162$$

	P ₁	P ₂	P ₃	P ₄
P ₁	0	2.828	3.162	5.099
P ₂	2.828	0	1.414	3.162
P ₃	3.162	1.414	0	2
P ₄	5.099	3.162	2	0

Distance Matrix

{distance point to itself} = 0

$$(P_1, P_4) :- \sqrt{(0-5)^2 + (2-1)^2} = \sqrt{(5)^2 + (1)^2} = \sqrt{25+1} = \sqrt{26} = 5.099$$

$$(P_2, P_3) = \sqrt{(2-3)^2 + (0-1)^2} = \sqrt{(-1)^2 + (-1)^2} = \sqrt{2} = 1.414$$

$$(P_2, P_4) = \sqrt{(2-5)^2 + (0-1)^2} = \sqrt{(3)^2 + (1)^2} = \sqrt{9+1} = \sqrt{10} = 3.162$$

$$(P_3, P_4) = \sqrt{(2-5)^2 + (1-1)^2} = \sqrt{(2)^2} = \sqrt{4} = \sqrt{2^2} = 2$$

Minkowski Distance

* Minkowski distance is a generalization of Euclidean Distance
Given two objects p & q .

(H)

$$\text{dist} = \left(\sum_{k=1}^n |p_k - q_k|^r \right)^{\frac{1}{r}}$$

Where r is a parameter, n is the number of dimensions and p_k & q_k are, respectively, the k th attributes of data objects p and q .

$\rightarrow r=1$ city block (Manhattan distance).

$\rightarrow r=2$:- Euclidean distance.

$\rightarrow r=\infty$:- "Supremum" (L_∞ norm, L_∞ norm) distance.

$$d(x,y) = \lim_{r \rightarrow \infty} \left(\sum_{k=1}^n |x_k - y_k|^r \right)^{\frac{1}{r}}$$

Manhattan
distance
formula

L_1	p_1	p_2	p_3	p_4
p_1	0	4	4	6
p_2	4	0	2	4
p_3	4	2	0	2
p_4	6	4	2	0

$\sqrt{1}$ power

Point	a	y
p_1	0	2
p_2	2	0
p_3	3	1
p_4	5	1

L_2	p_1	p_2	p_3	p_4
p_1	0	2.828	3.162	5.099
p_2	2.828	0	1.414	3.162
p_3	3.162	1.414	0	2
p_4	5.099	3.162	2	0

$(p_i - p_j)$ distance
 $(p_i - p_j)$

Euclidean Distance.
formula $\sqrt{2}$ /
power of 2

and among the given

Manhattan Distance

$$P_1 = (0, 2) = (x_1, y_1)$$

$$P_2 = (2, 0) = (x_2, y_2)$$

$$\text{Distance} = |x_1 - x_2| + |y_1 - y_2| \\ = |0 - 2| + |2 - 0| \\ = 2 + 2 = 4$$

Distance from P_1 to P_2 is 4

And distance from P_2 to P_1 is 4.

→ Similarly for the other points.

Common Properties of Distance

Common properties of distance are $d(p, q) \geq 0$ for all p, q and $d(p, p) = 0$ only if $p = q$.

1) $d(p, q) \geq 0$ for all p, q (positive definiteness)

2) $d(p, q) = d(q, p)$ for all p, q (symmetry)

3) $d(p, r) \leq d(p, q) + d(q, r)$ for all points p, q, r (Triangle Inequality).

(Triangle Inequality). between distance (dissimilarity) between

where $d(p, q)$ is the distance

points (data objects) p & q . Matching and coefficients

[similarity SMC versus Jaccard]

$$P = 100\ 000\ 000\ 000$$

$$0\ 0\ 0\ 0\ 0\ 0\ 1\ 0\ 0\ 1$$

$M_{01} = 2$ (the no's of attributes when p was 0 and q was 1)

$M_{10} = 1$ (the no's of attributes when p was 1 and q was 0)

$M_{00} = 2$ (the no's of attributes when p was 0 and q was 0)

$M_{11} = 0$ (the number of attributes where p was 1 and q was 1)

$$SMC = \frac{(M_{11} + M_{00})}{(M_{01} + M_{10} + M_{11} + M_{00})} \quad (P2)$$

$$\rightarrow \frac{(0++)}{(2+1+0++)} \quad |$$

$$= \frac{2}{10}$$

$$= 0.2 = 0$$

$$J = \frac{(M_{11})}{(M_{01} + M_{10} + M_{11})}$$

$$= 0 [(2+1+0) = 0]$$

$SMC = \frac{\text{No of matches}}{\text{no of attributes}}$

$$= \frac{(M_{11} + M_{00})}{(M_{01} + M_{10} + M_{11} + M_{00})}$$

$$J = \frac{\text{number of " matches}}{\text{no of not-both-zero attributes values.}} \quad |$$

$$= \frac{(M_{11})}{(M_{01} + M_{10} + M_{11})}$$

Cosine similarity

If d_1 & d_2 are two document vectors, then

$$\text{cos}(d_1, d_2) = \frac{(d_1 \cdot d_2)}{\|d_1\| \|d_2\|}$$

where \cdot indicates vector dot product and $\|d\|$ is

the length of vector d .

$$\text{Ex: } d_1 = 3205000200$$

$$d_2 = 111$$

$$d_2 = 1000000102$$

$$d_1 \cdot d_2 = 3*1 + 2*0 + 0*0 + 5*1 + 0*0 + 0*0 + 2*1 + 0*0 + 0*2 = 7$$

$$\|d_1\| = \sqrt{3^2 + 2^2 + 0^2 + 5^2 + 0^2 + 0^2 + 2^2 + 0^2 + 0^2} = \sqrt{35} = 5.91$$

$$\|d_2\| = \sqrt{1^2 + 1^2 + 1^2 + 0^2 + 0^2 + 0^2 + 0^2 + 1^2 + 0^2 + 2^2} = \sqrt{10} = 3.16$$

$$\cos(d_1, d_2) = 0.3150$$

$$\cos(d_1, d_2) = \frac{(d_1, d_2)}{\|d_1\| \|d_2\|}$$

$$\cos(x^f) = \frac{5}{(6.487)(2.245)}$$

$$\approx 0.3150$$

Extended Jaccard Coefficient (Tanimoto)

$$T(p, q) = \frac{p \cdot q}{\|p\| + \|q\| - p \cdot q}$$

$\cos(x^f) = 0$ indicate both are dissimilar

$\cos(x^f) = 1$ indicates both are similar

→ Problem : Euclidean distance is

(13)

	A ₁	A ₂	
x ₁	1.5	1.7	0.1414
x ₂	2	1.9	0.6708
x ₃	1.6	1.8	0.2226
x ₄	1.2	1.5	0.6082

let $p_1(x_1, y_1)$ and $p_2(x_2, y_2)$

euclidean Distance =

$$\sqrt{(x_2 - x_1)^2 + (y_2 - y_1)^2}$$

$$x = (1.4, 1.6) = \sqrt{(1.4 - 1.5)^2 + (1.6 - 1.7)^2} = 0.1414 \Rightarrow 0.1414$$

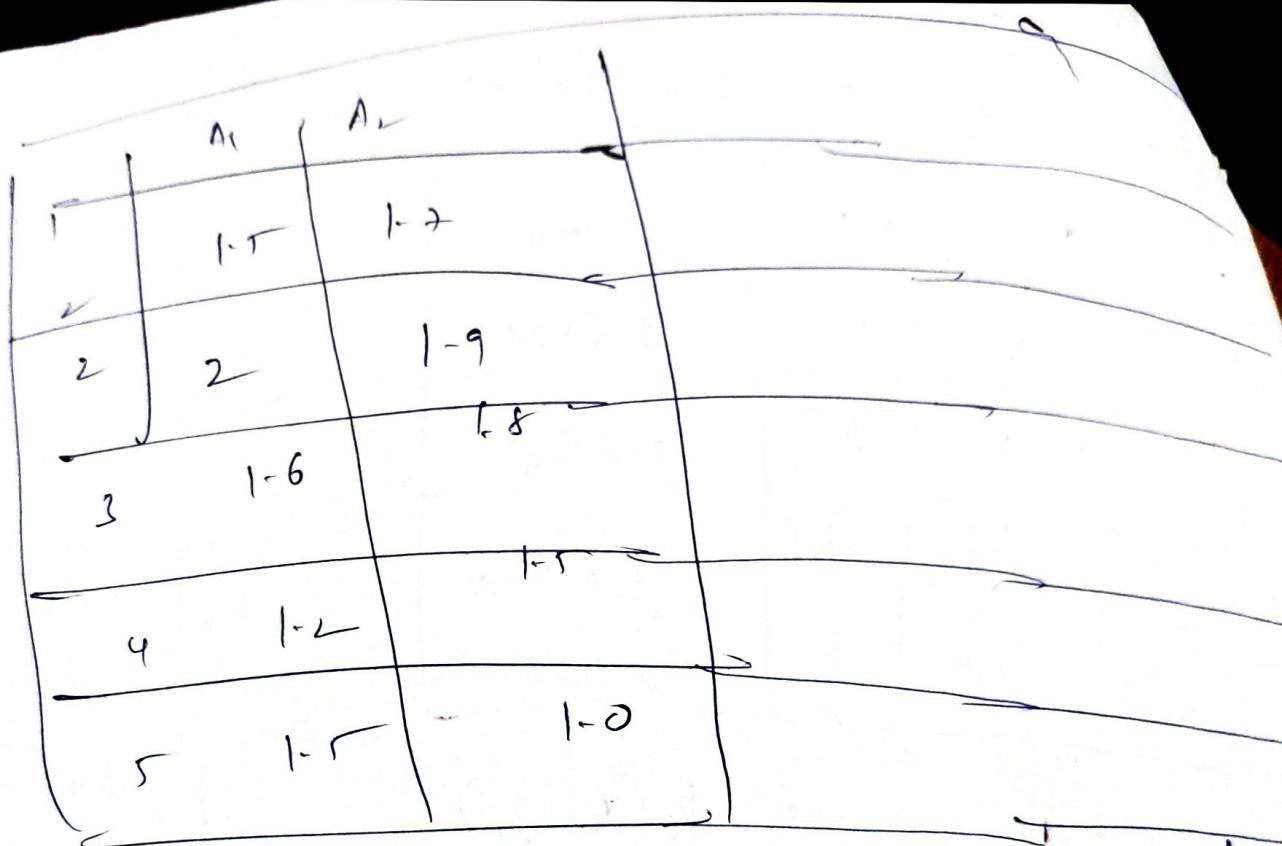
→ Manhattan Distance :

	A ₁	A ₂	Manhattan Distance
1	1.5	1.7	$ 1.4 - 1.5 + 1.6 - 1.7 = 0.2$
2	2	1.9	$= 0.9$
3	1.6	1.8	$= 0.4$
4	1.2	1.5	$= 0.3$
5	1.5	1.0	$= 0.5$

let $p_1(x_1, y_1)$ and $p_2(x_2, y_2)$

$$MD = |x_2 - x_1| + |y_2 - y_1|$$

and access in



$$n = (1.4, 1.6)$$

let $P_1(x_1, y_1)$ and
 $P_2(x_2, y_2)$

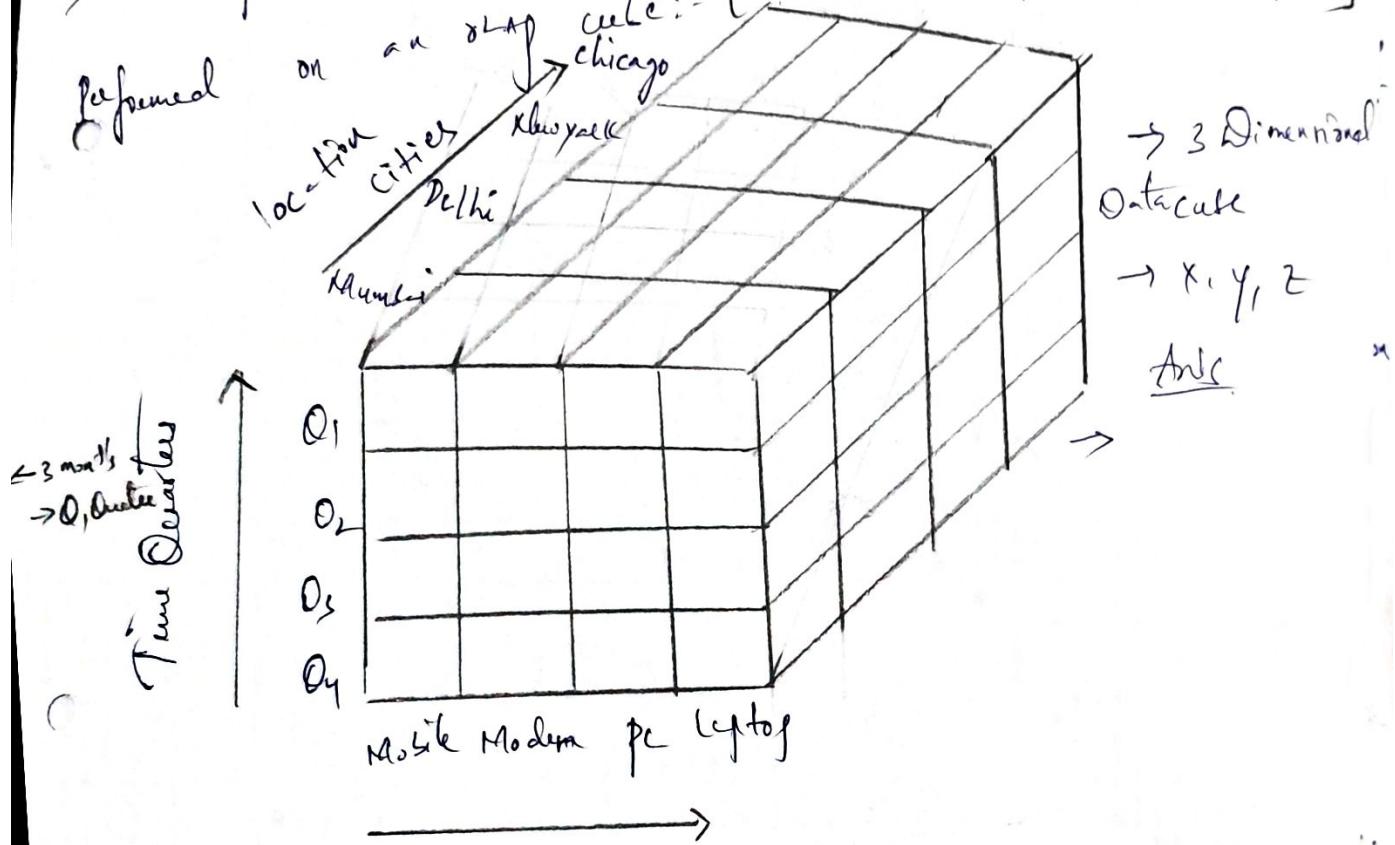
$$h=2$$

Minkowski distance

$$\sqrt{h^2 + (x_2 - x_1)^2 + (y_2 - y_1)^2}$$

"OLAP operations"

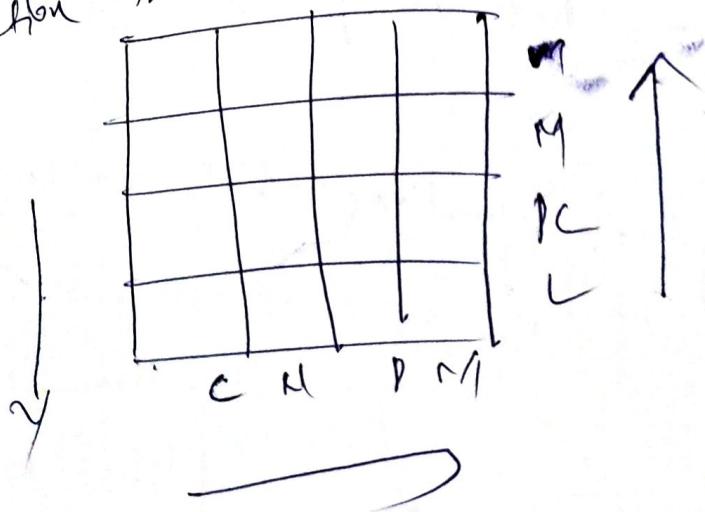
- They are based on multidimensional data model and allows the user to query on multidimensional data.
- OLAP cube is a structure that is optimized for proper data analysis.
- 5 operations basic analytical operations that can be performed on an OLAP cube:- [Drill Down, Rollup, Dice, slice, pivot]

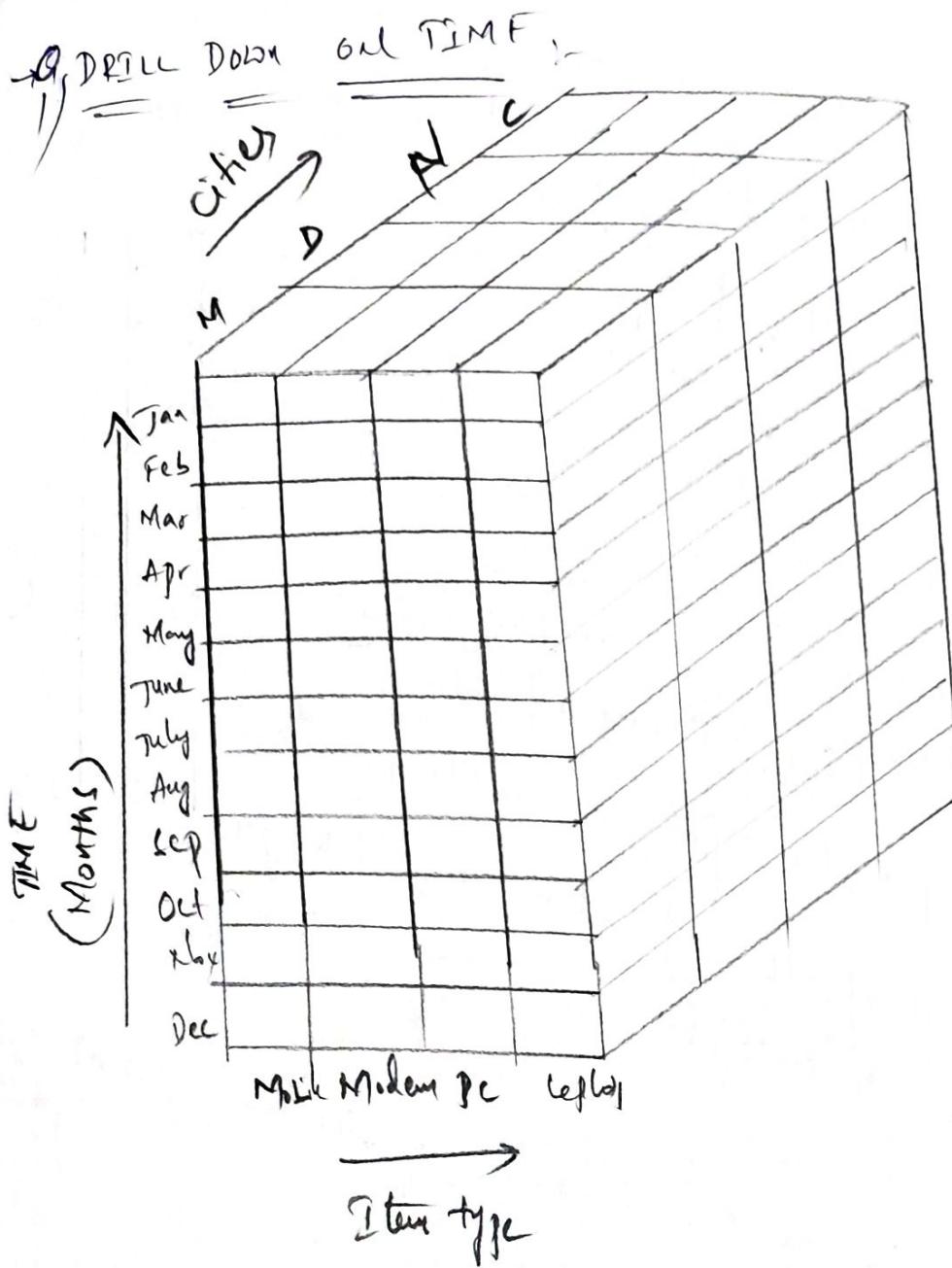


This type

Drill Out:- We want details information on more information from the down side cube

→ Pivot :- [Rotating any given date]
The positions of X, Y & Z are changed. & rotating of the
information then it will automatically reflected as
the new coordinates

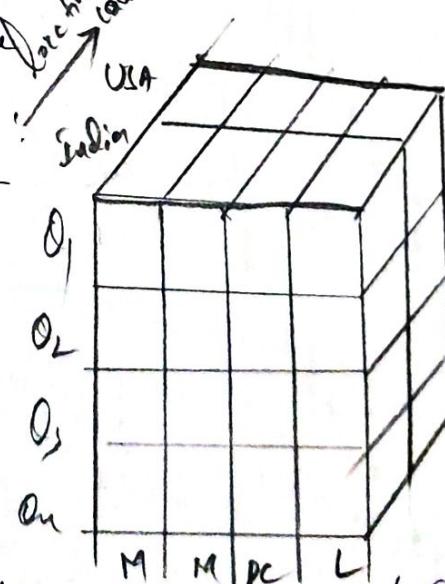




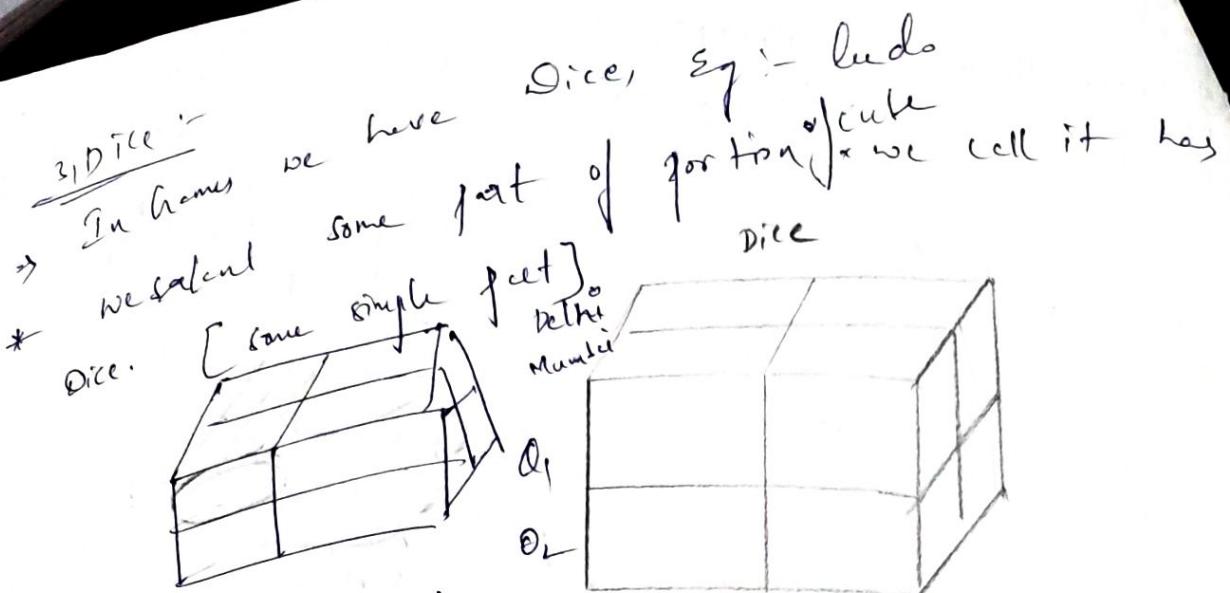
ii) Rollup :- Exact opposite side (or) top of the cube

operations we will perform for countries
Rollup on locations

cities are four → Chicago } USA
 → New York }
 → Delhi } India
 → Mumbai }



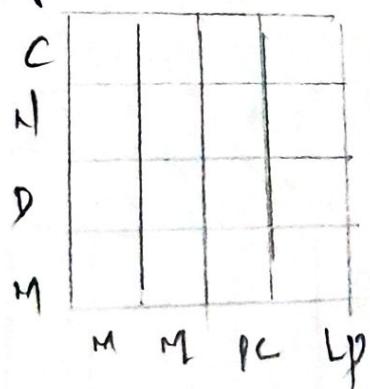
* We want generalization information. (here in terms of country)
 and access the data.



\rightarrow Here only we want the information for the (Q_1, Q_2) Quarters information as well as Number & Delivery of particular product mobile, modem.

\rightarrow Some part of cube information we taking here that information by the overall cube side.

\rightarrow 4) Slice :— Slice is cutting of the cube what ever information you need.



\rightarrow 'Q₁' Quarter information of the factors & product we want in that just one piece of cube compact of information needed.