

17/11/23

Data Mining & Warehousing:-

Data Mining :- extracting knowledge from data.

↳ supervised mining learning.

↳ unsupervised mining learning,

↳ Market Basket Mining Analysis - based on the products in the shopping cart of customers so that the brands can analyse.

→ Data Mining is the process of analysing data from different perspective and summarising it into useful information - info. that can be used to increase revenue, cut cost or both.

Data Mining Strategies:-

Semi Supervised learning

Supervised learning.

Unsupervised clustering.

Market Basket Analysis.

① Supervised learning! - When we are young we use induction to form basic concept definition, we see instances of concepts, representing animals, plants, buildings, etc.)
(Here, labels given to individual instances and choose what we believe to be the defining concept feature (attribute) from our own classification model. Later we use the model, we have developed to help us identify objects of similar structure. The name of this type of learning is supervised concept of learning or just supervised learning.

Supervised

Classification

Regression

i) Classification :- It is a supervised learning task where output is having defined labels.

Types - Binary, Multi Class, Multi Label.

Binary :- In binary classification the model predicts either 0 or 1, yes or no but in case of multi class classification the model predicts any one class from more than 2 classes.

In multi label classifi. any one instance can have more than one label or classes.

ii) Regression! - It is supervised learning task where output is having continuous value. It is generally divided into two parts - linear regression and logistic regression.

③ Semi-Supervised :- It is a supervised learning where the training data contains the very few labelled examples and a huge no. of unlabelled examples.

In this initially similar data is clustered along with an unsupervised learning algo. and further it helps to label the unlabelled data into labelled data.

Semi-supervised lies b/w sup. & unsup. learn'g.

③ Unsupervised clustering:- unlike supervised learning, unsupervised learning builds model from data without pre-defined classes. Data instances are grouped together based on the similarity scheme defined by the clustering system.

Clustering is the process of grouping the data into clusters, so that objects within a cluster have high similarity in comparison to one another but are very dissimilar to objects in other clusters.

18/1/23

Temporal database :- column of any table has attribute of time.

Time Series Databases :- there is always a fixed set

→ Data Mining can be applied on many kinds of data :-

(i) Transactional Databases :- These consist of a file where each record represent a transaction. A transaction typically includes a unique trans. Id number & a list of items making up the transaction.

(ii) Temporal Databases :- It stores and manages time varying data, eg: financial, medical, government etc; It maintains historical info.; incorporate notion of time in the system.

(iii) Time Series Databases: In a time series, there is a sequence of data points, measured typically at successive points in time, spaced at uniform time intervals.

→ A time series database stores sequence of values of events obtained over repeated measurement of time. Eg - hourly, daily, weekly.

(iv) Spatial Databases :- These databases contain spatial related info. Eg - geographical map databases, satellite image databases.

→ A spatial database that stores objects that change with time is called a spatio-temporal database.

Data Cleaning:-

↳ Missing Values.

↳ Noisy Data

(a) Missing Values:- (i) Ignored the tuple

when few values are missing then don't ignore, else ignore the table.

↳ This is usually done when several values attributes with missing values are present on the class label & missing.

(ii) Fill in the missing values manually.

↳ In general this approach is time consuming & may not be feasible given a large dataset with many missing values.

(iii) Use a global constant to fill in the missing values.

↳ Replace all missing values by the same constant such as label like unknown, ∞ , 0, etc. The

↳ The problem with this is that the mining program may mistakenly conclude that they form a interesting concept since they all have a value in common.

(iv) Use attribute mean to fill missing values.

↳ Missing values are replaced by mean value of column.

(v) Use the attribute mean for all samples belonging to the same class as the given tuple.

id	d_1	d_2	d_3	d_4	Income	Class
					100	C1
					200	C1
					?	C1
					350	C2
					450	C2
					mean value of C1 only.	Class

(vi) use the most probable value to fill in the missing values.
↳ other methods to fill values. (like regression)

(b) Noisy Data :-

smooth

(i) Binning :- ~~smooth~~ the sorted data values by consulting it's neighbourhood, that is the value around it.

The two common techniques that are used are as follow —

(i).(a) — Data Smoothing By Bin means.

(i).(b) — Smoothing By Bin Boundaries.

(i).(a). 4, 8, 15, 21, 21, 24, 25, 28, 34.

Bin 1: 4, 8, 15 \rightarrow [9, 9] \rightarrow equal size bins

Bin 2: 21, 21, 24 \rightarrow [22, 22, 22] partitioning.
 \rightarrow data should be mixed.

Bin 3: 25, 28, 34 \rightarrow [29, 29, 29]

↳ replace all values by bin ~~means~~ mean.

(ii).(b)

Bin 1:

\rightarrow 4, 4, 15

Bin 2:

\rightarrow 21, 21, 24

Bin 3:

\rightarrow 25, 25, 34.

replace all non boundary numbers with minimum difference boundary number

(ii) Regression:-

Data can be smooth by fitting the data to a function such as with regression.

(iii) clustering:-

Outliers may be detected by clustering, where similar values are organised into groups or clusters.

1/2/23

Mean, Median, Mode:-

Q. 13, 18, 13, 14, 13, 16, 14, 21, 13.

$$\text{Mean} = \frac{135}{9} = 15.$$

$$\text{Mode} = 13.$$

Median:-

arranging in ascending order - 13, 13, 13, 13, (14), 14, 16, 18, 21.

$$\text{Median} = 14.$$

Q. 10, 12, 13, (16), (17), 18, 19, 21.

$$\text{Mean} = \frac{15.75}{9} = 1.75.$$

$$\text{Median } \cancel{\text{mean}} = \left(\frac{n+1}{2} \right) \text{th term, } \frac{16+17}{2} = 16.5.$$

$$\text{Mode} = 3 \text{Median} - 2 \text{Mean}$$

$$= 3(16.5) - 2(15.75)$$

$$= 18.$$

age	frequency
1-5	200
5-15	450
15-20	300
20-50	1500
50-80	700
80-110	44

Find the median of age.

$$\text{Sol}^7 \quad \text{Median} = L_1 + \left(\frac{N/2 - (\sum \text{freq})_L}{\text{freq}_{\text{median}}} \right) \text{width}$$

where L_1 = lower boundary of median interval.

N = No. of values in the entire data set.

$(\sum \text{freq})_L$ = sum of frequencies of all the intervals that are lower than the median interval.

$\text{freq}_{\text{median}}$ = frequency of median interval.

width = width of the median interval.

N = sum of all frequencies.
= 3194.

$N/2 = \underline{\underline{1597}}$ lies in $(20-50)$.
∴ Median interval = $(20-50)$.

$$L_1 = 20.$$

Median $\Rightarrow (\sum \text{freq}) / 2 = 950$.

$$\text{width} = 30.$$

$$(\text{freq})_{\text{median}} = \underline{\underline{1500}}.$$

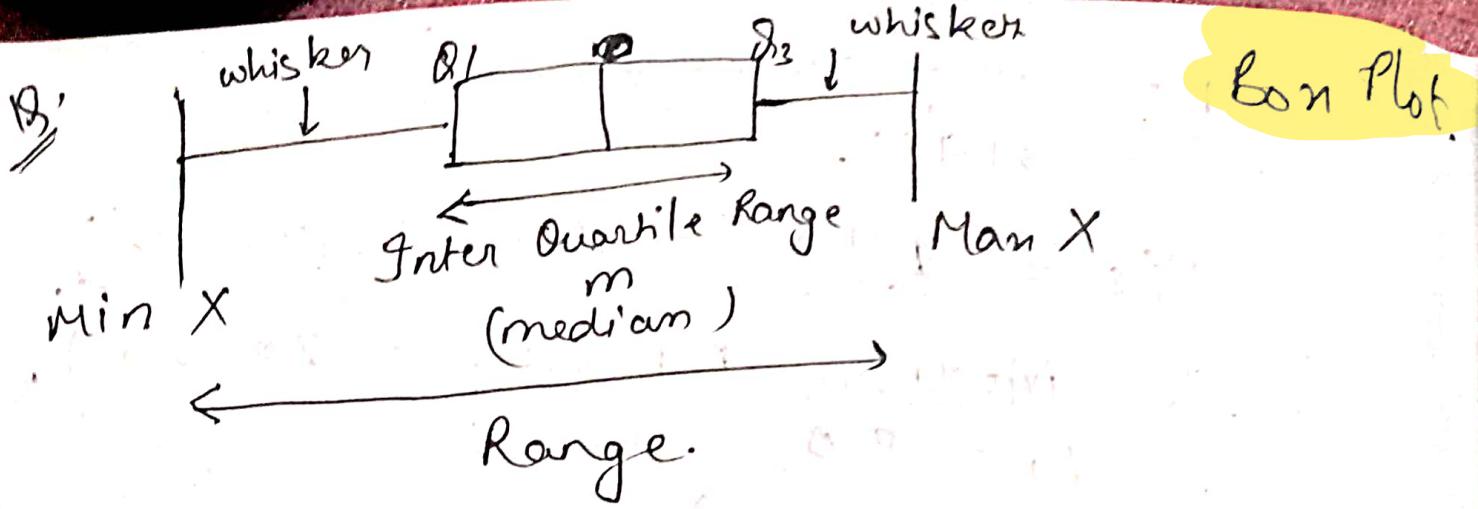
$$\text{Median} = 20 + \left(\frac{1597 - 950}{1500} \right) \times 30.$$

$$= 20 + \frac{647 \times 30}{1500}$$

$$= 20 + \frac{19410}{1500}$$

$$= 20 + 12.94$$

$$\boxed{\text{Median} = 32.94}$$



Data, 76, 79, 76, 74, 75, 71, 85, 82, 82, 79, 81.

→ arrange in increasing order.

→ median find.

21, 24, 25, 26, 26,

left half median

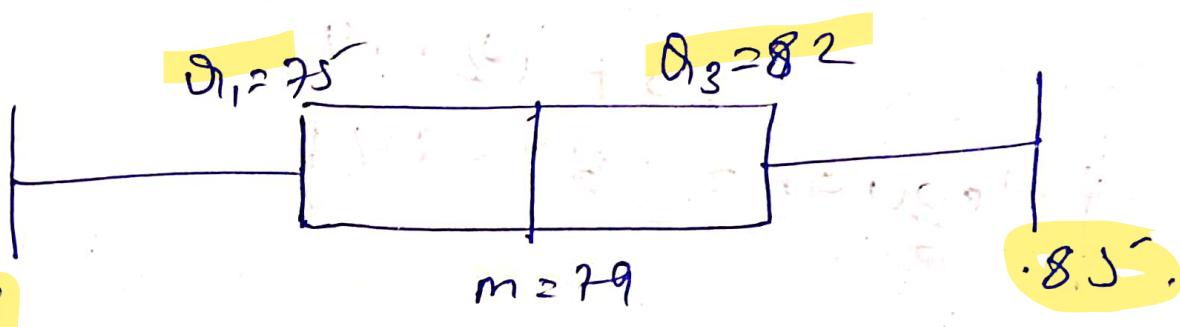
79

79, 81, 82,
82, 85.

right half
median

$$\text{Median} = 79 \Rightarrow m = 79$$

Measures



Q. calculate the standard deviation of
4, 2, 5, 8, 6.

mean $\bar{x} = \frac{25}{5} = 5$

$s = \sqrt{\frac{\sum(x - \bar{x})^2}{N}}$

$$= \sqrt{\frac{1+9+0+9+1}{5}}$$

$$= \sqrt{\frac{20}{5}} = \sqrt{4} = 2.$$

Q: Correlation Analysis

Pearson Correlation :

- Given two attributes, such analysis can measure how strongly one attribute implies the other, based on the available data.
- The value of pearson correlation lies b/w $[-1, 1]$ i.e. Pearson ($-1 \leq r_{A,B} \leq 1$)
- If $(r_{A,B} > 0)$, then A & B are (+ very) correlated. i.e. when the value of A ↑ then the value of B also ↑.
 - ↳ Higher the value, stronger the correlation.
- If $(r_{A,B} = 0)$, then A & B are independent and there is no correlation b/w them.
- If $(r_{A,B} < 0)$, then A & B are (-very) correlated.

$$H_{A,B} = \frac{\sum_{i=1}^N (a_i - \bar{A})(b_i - \bar{B})}{\sqrt{\sum_{i=1}^N (a_i - \bar{A})^2} \sum_{i=1}^N (b_i - \bar{B})^2}$$

Q. Tree Height (y) Trunk Diameter (x)

35	8
49	9
27	7
33	6
60	13
21	7
45	11
51	12

$$g_1 = 0.886$$

Solving Mean -

χ^2 (chi-square) Test

- For categorical data a correlation relationship b/w two attributes, A & B, can be discovered by χ^2 test. Suppose A has 'c' distinct values i.e. A_1, A_2, \dots, A_c & B has 'd' distinct values i.e. (b_1, b_2, \dots, b_d) . Let

→ let (A_i, B_j) denotes the event that attribute A takes the value A_i and B takes the value B_j . This is represented by a distinct slot in the table. This table is known as contingency Table. The χ^2 value is

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

where, O_{ij} = observed frequency (i.e. actual count of the joint event).

E_{ij} = expected frequency of (A_i, B_j) .

$$E_{ij} = \frac{\text{count}(A=a_i) \times \text{count}(B=b_j)}{N}$$

where, N = no. of data tuples

→ Degree of freedom = $(r-1) \times (c-1)$:

	male	female	Total.
1. fiction	250 (90)	200 (360)	450.
2. Non-fiction	50. (310)	1000 (840)	1050.
Total	300	1200	1500.

$$e_{11} = \frac{\text{count (male)} \times \text{count (fiction)}}{N}$$

$$= \frac{340 \times 450}{1500} = 90.$$

$$e_{12} = \frac{\text{count (female)} \times \text{count (fiction)}}{N}$$

$$= \frac{1200 \times 450}{1500} = 360.$$

$$\chi^2 = \frac{(250-90)^2}{90} + \frac{(50-310)^2}{310} + \frac{(200-360)^2}{360} + \frac{(1000-840)^2}{840} = 507.$$

significance level — 5%.] = 3.84
(table $D.F = 1$. will be given):

χ^2_{cal}

if $\chi^2 < 3.84$ then these values
are not correlated.

Data Transformation:-

- Smoothing
- Aggregation
- Generalization

Attribute Construction

(New attributes are constructed from given set of attributes)

- 1) First step is smoothing which works to remove noise from the data, such techniques include binning, regression, clustering.
- 2) Aggregation:- Where summary or aggregation operations are applied to the data. Eg. Daily sales may be aggregated to compute monthly or annual sales.
- 3) Generalisation:- Where low level primitive data are replaced by higher level concept through the use of concept hierarchies.
Eg. City can be generalized to a higher level concept like state or country

4) Normalization - Where the attribute data are scaled so as to fall within specified range
 Eg. -1 to 1 or 0 to 1.

5) Attribute Construction - Where new attributes are constructed from a given set of attributes.

Min-Max Normalization

$$V' = \frac{V - \min_A}{\max_A - \min_A} \cdot (\text{new_max}_A - \text{new_min}_A) + \text{new_min}_A$$

where \max_A and \min_A are maximum and minimum values of an attribute A

Min-Max Algo map a value V of attribute A to a value V' in the range new_min_A to new_max_A

Suppose the minimum and maximum value of attribute income are 12000 and 98000 map the income in range 0.0 to 1.0 using min_max Norm find the value of income 73600.

$$V' = \left(\frac{73600 - 12000}{98000 - 12000} \right) (1 - 0) + 0$$

$$= \frac{61600}{76000}$$

$$= \underline{0.716}$$

Z-score Normalization

The values of an attribute A are normalized based on the mean and standard deviation of A.

$$V' = \frac{V - \bar{A}}{S_A}$$

The mean and standard deviation of the values for the attribute income are respectively 51000 and 16000. Use with Z-score normalization transform the values of income to 73600.

$$V' = \frac{73600 - 51000}{16000}$$

$$= \frac{22600}{16000} = 1.4125$$

Data Warehouses - OLAP

Online Analytical

by Managerial level (Decision
Makers)

Databases - OLTP

Online Transactions

(Lower Level Staffs)

Data Warehouse :-

A Database

According to William H Inmon. A

Data warehouse is a subject

SINT

Oriented, Integrated, Time variant
and Non-Volatile Collection of
Data in Support of Management
decision making process.

- 1) **Subject Oriented** - A Data warehouse is organised around major subject areas like customer, supplier, etc.
- 2) **Integrated** - A Data Warehouse is usually constructed by integrating multiple heterogeneous sources.
- 3) **Time Variant** → Data are stored to provide information from a historical perspective. Every key structure in the datawarehouse contains explicitly or implicitly an element of time.

J.t. usually requires only two operations in data accessing. Initial loading of data and accessing of data.

location - Belrar

item (type)

time (Quata)	Home Entertainment	Computer	Home Phone
Q ₁	605	365	12
Q ₂	625	205	13
Q ₃	425	195	19
Q ₄	365	825	29

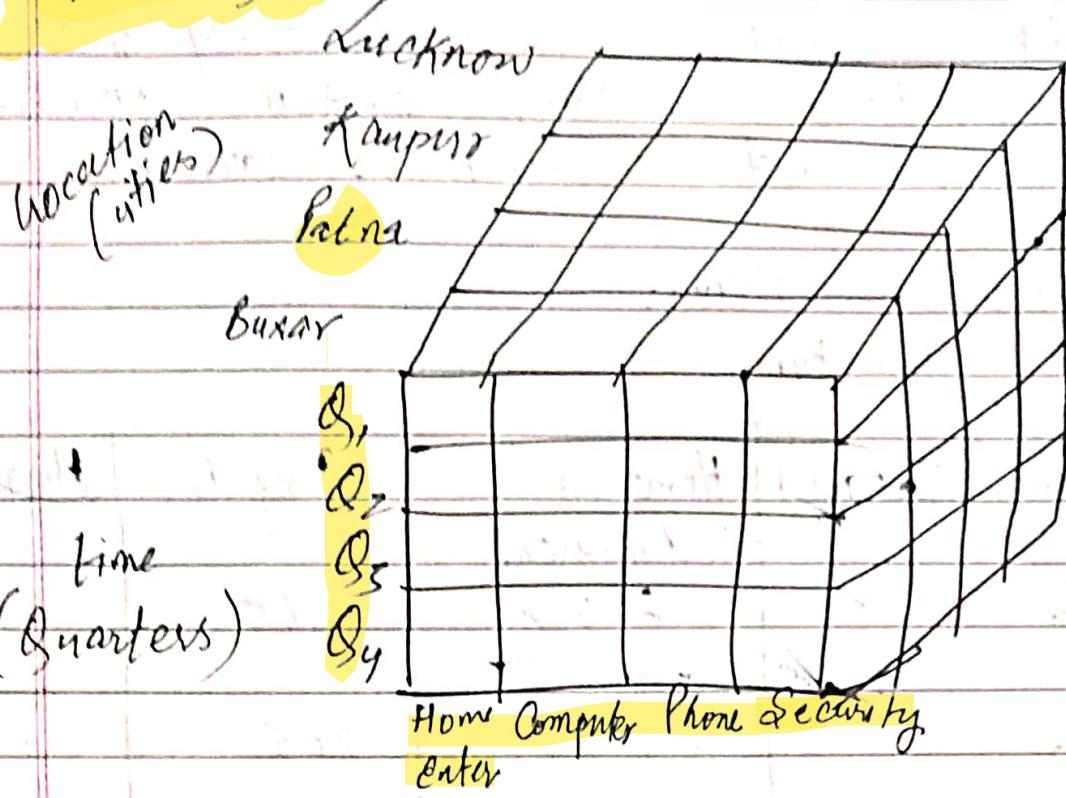
location - "Patna"

item (type)

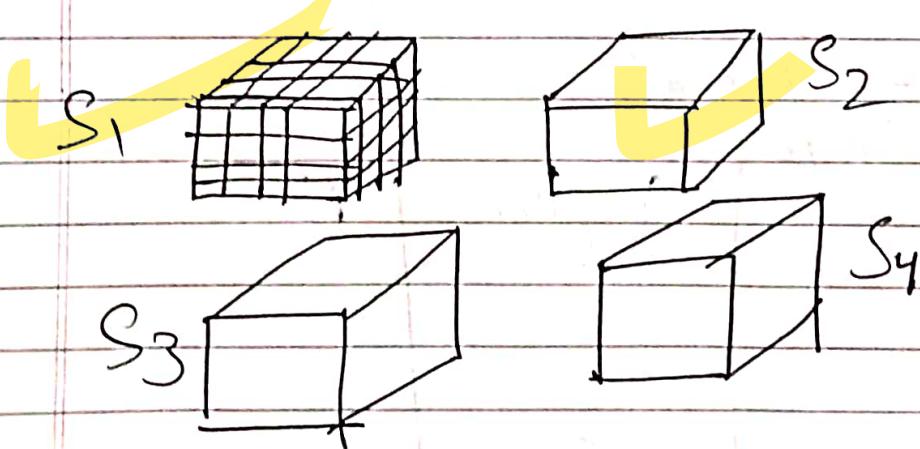
time	Home	Computer	Phone
Q ₁	395	1600	
Q ₂	700	300	
Q ₃	900	900	
Q ₄	400	400	

location - "Kangra"

Data Cube



If we want to increase dimensions for eg. by adding the Supplier attribute make one cube for each supplier.



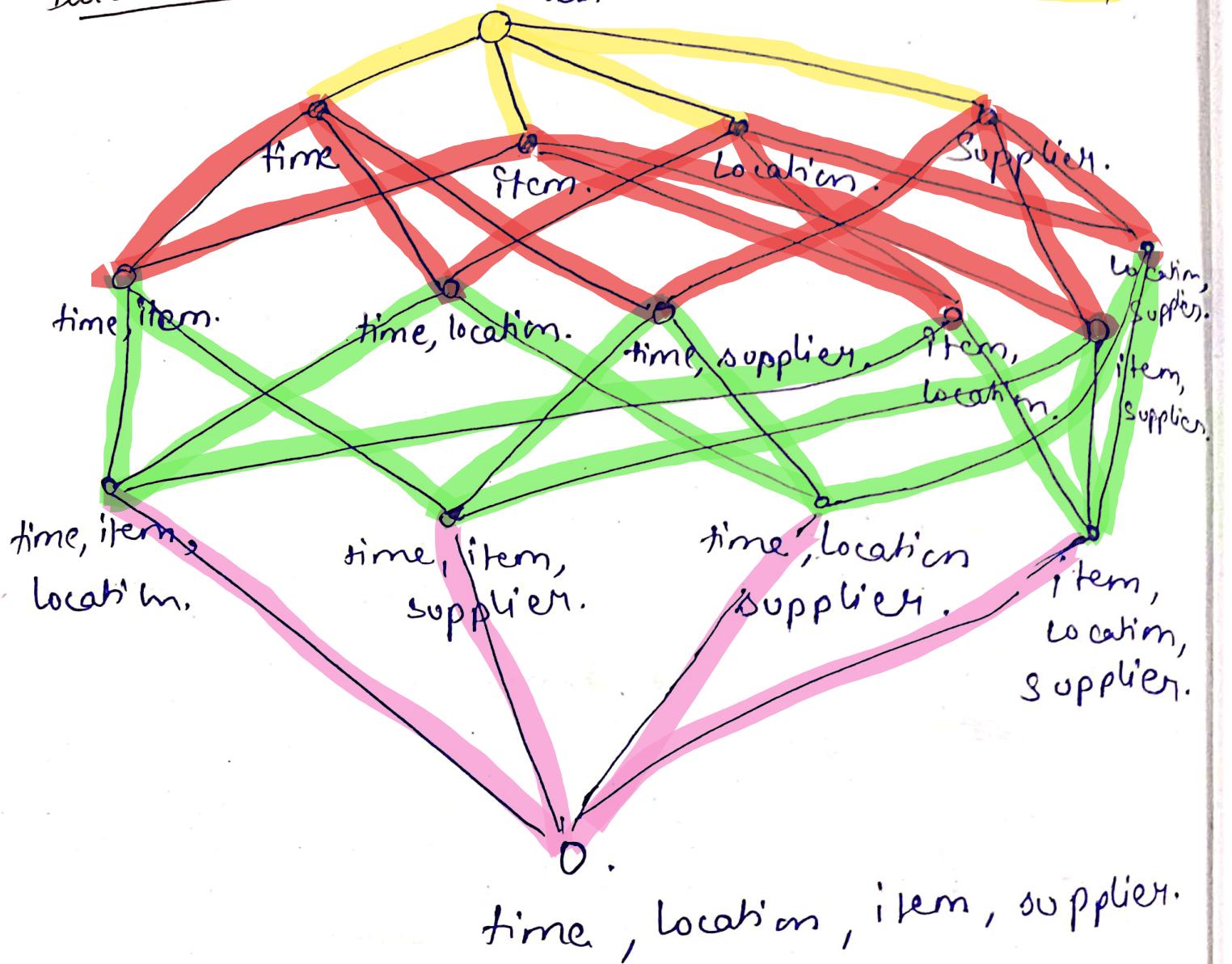
If we continue in this way the n dimensional data can be viewed as n-1 dimensional cubes

8/2/23 :-

Data Warehouse :-

(Lattice of cuboid). .

All.



→ Given a set of dimensions, we can generate a cuboid, for each of the possible subsets of given dimensions. The result would form a lattice of cuboids, each showing a data at different level of summarisation.

The figure shows the lattice of cuboids forming a data cube for the dimension - time, item, location & supplier.

- The cuboid that holds the lowest level of summarisation is called the base cuboid.
- The zero dimⁿ cuboid which holds the highest level of summarisation is called the open cuboid.

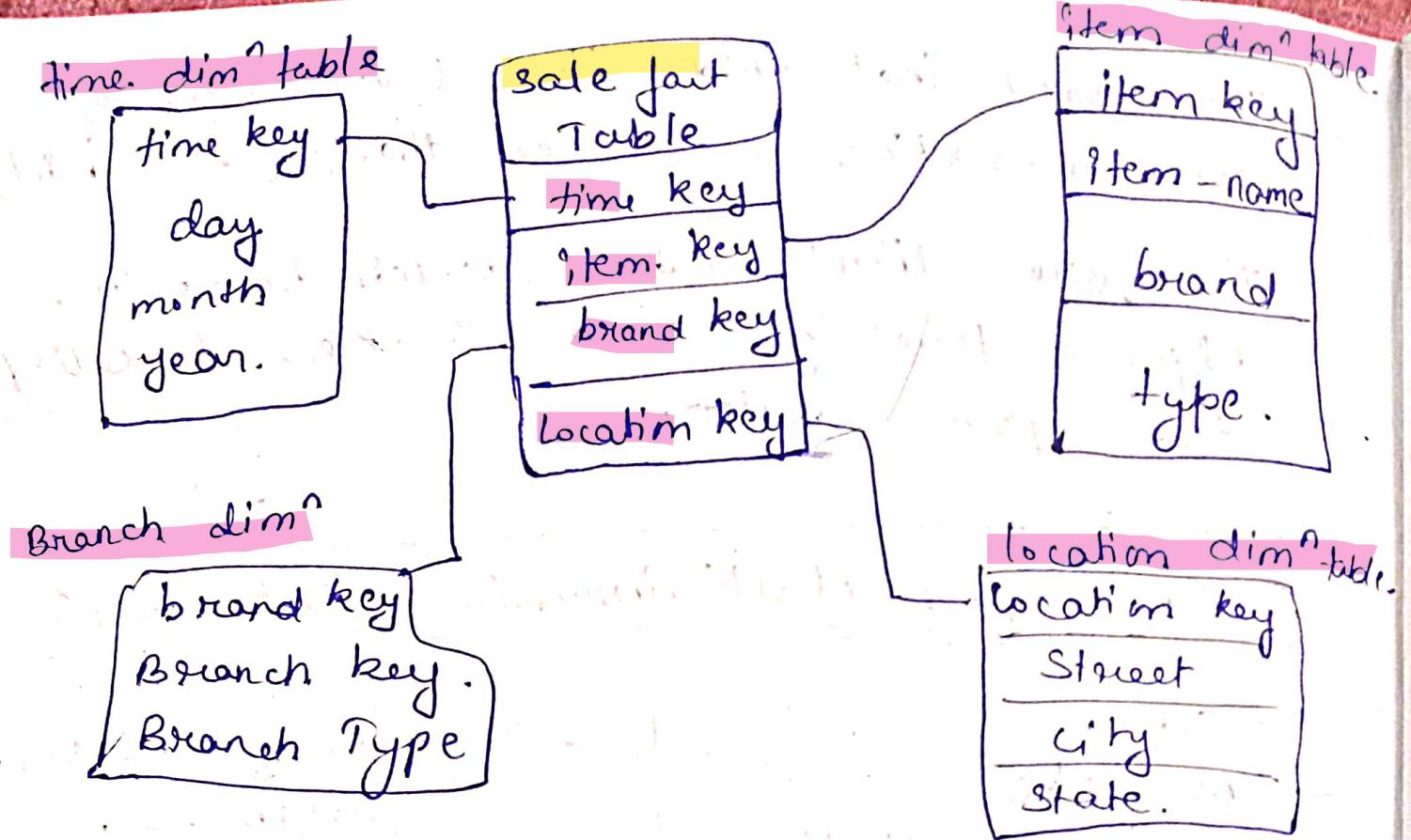
Schemas for Multidimensional Databases:-

- Star
 - Snowflake
 - Fact - Constellation.
-] 3 schemas for Multi-D databases.
- The E-R model is commonly used in the design of relational databases, where a database schema consists of set of entities and the relationship b/w them.
 - The most common data models for data warehouse are
 - Star
 - Snowflake
 - Fact - Constellation

Star :-

Sale Part Table
time key
item key
brand key
location key

P.T.O.



→ The most common model is star schema in which data warehouse contain —

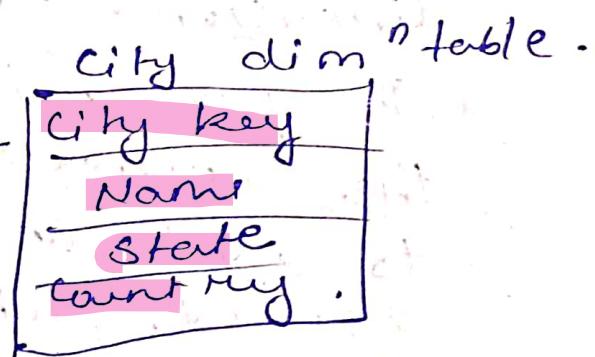
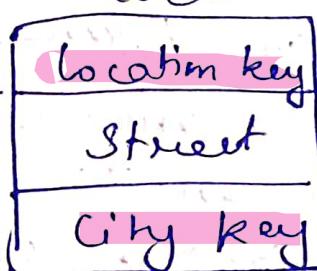
(i) a large central fact table containing a bulk of data with usually no redundancy.

(ii) A set of smaller attendant tables (dimⁿ tables), one for each dimension.

→ The schema graph resembles star.

Snowflake :-

→ Same diagram as star but location dimⁿ table.



→ The snowflake schema is a variant of star schema where some dimⁿ tables are normalized, thereby further splitting of data into further tables.

Fact - Constellation :-

→ It is collection of star schema by help of dimⁿ tables.

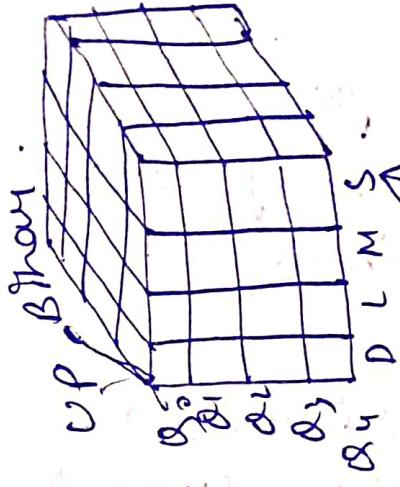
→ Sophisticated applications may require multiple fact tables to share dimension table. This kind of schema can be viewed as a collection of stars and hence is called as galaxy schema or a fact constellation schema.

14/2/23.

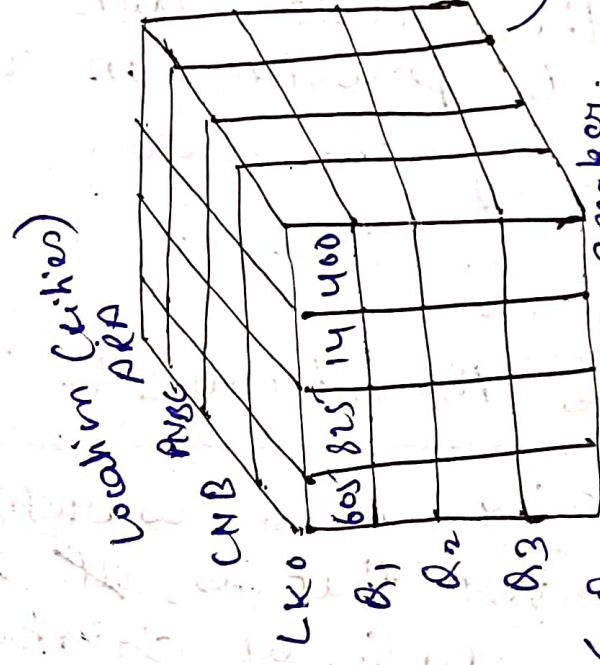
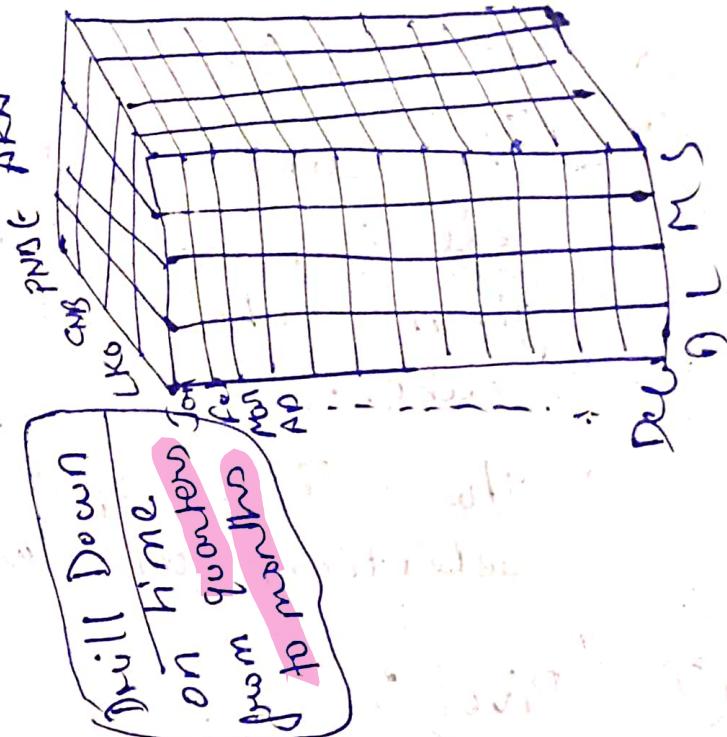
- Data Cube :- A data cube allows data to be modelled & viewed in multiple dimension. In general terms → In general terms, dimensions are the perspective or entities with respect to which an organization wants to keep related. Each dimn may have a table associated with it called a dimension table. → Fact Table :- contains the name of the fact or keys of each related dimension table.

OLAP operations :-

- ① Roll up - The roll up performs the aggregation on a data cube either by climbing up or dimensionality reduction.
- ② Drill Down - The drill down is reverse of roll up. It navigates from less detailed data to more detailed data. It can be realized either by stepping down a concept hierarchy for a dimension or introducing additional dimensions.



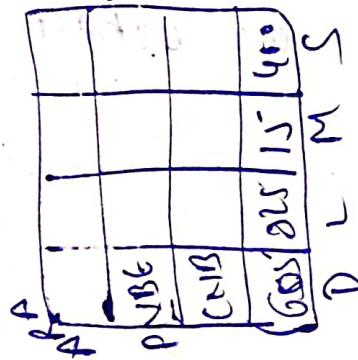
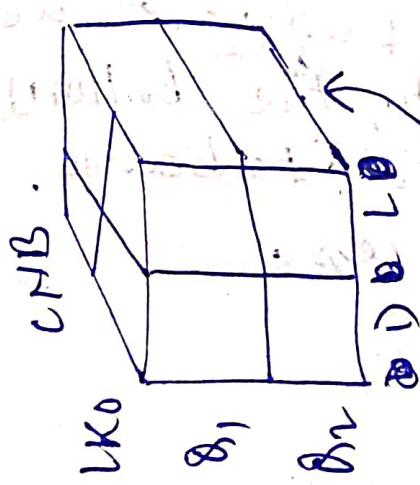
Roll-up
on allocation
criteria



Q4 Desk top mobile speaker
top top file
stem (type).

100	50	5	0.00
25	12.5	1.25	0.05
6.25	3.125	0.3125	0.0015
1.5625	0.78125	0.078125	0.000391
0.390625	0.1953125	0.01953125	0.00009765625

S S K R P P R E U N S L K O .



A child's drawing of a teardrop-shaped balloon. The word "slice" is written vertically along the left side, and "time" is written vertically along the right side. The bottom part of the teardrop is shaded pink.

A simple line drawing of a person from the waist up. The person has a large, round head with a single vertical line for a mouth. They are wearing a pink, short-sleeved t-shirt. The drawing is done with black outlines on a white background.

③ Slice & Dice :-

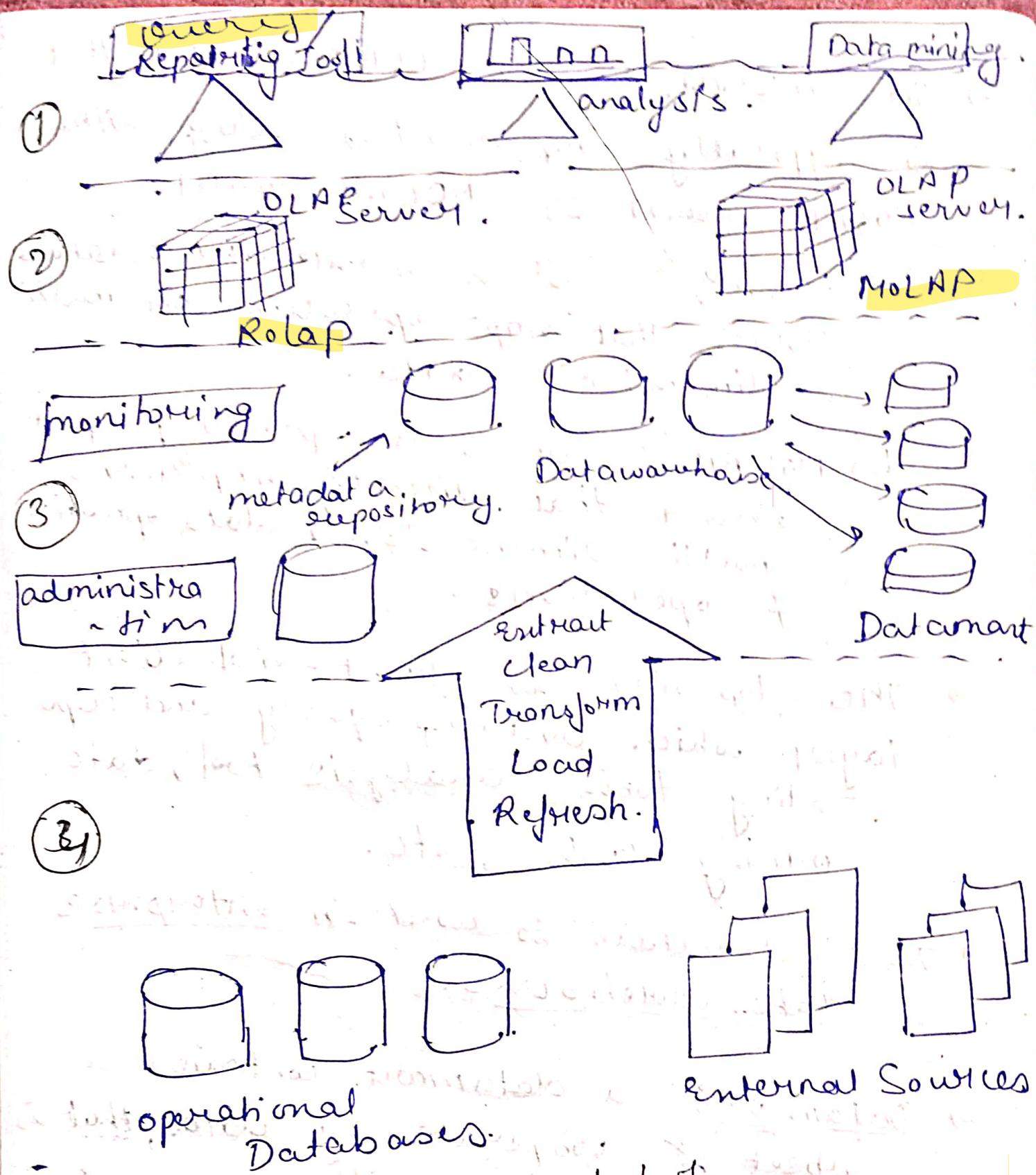
- The slice operation performs a selection on 1 dimension of the given cube, resulting in a sub cube.
- The dice operation performs a selection on two or more dimensions.

④ Pivot :-

It is a visualisation operation that rotates the data axis in view in order to provide an alternate presentation of the data.

Dataware House Architecture :- (Three Tier)

- The Bottom Tier is a warehouse database server that is almost always a relational database system. Backend tools & utilities are used to feed data into the bottom tier from operational databases or other external sources.



- ① Top Tier front end tools.
- ② Middle Tier OLAP servers.
- ③ Bottom Tier Datawarehouse Server.
- ④ DataMART

- The middle tier is OLAP server that is typically implemented using either ROLAP model or MOLAP model.
- ↳ ROLAP :- It is an extended relational DBMS that maps operations on multi-dimensional data.
- ↳ MOLAP :- It is a special purpose server that directly implements multi-dimensional data operations.
- The top tier is front-end client layer which contains query and reporting tools, analysis tool, data mining tools, etc.
- This structure is used in enterprise data warehouse.

- Datamart :- A datamart contains a subset of corporate data, that is of value to a specific group of users. Its scope is confined to specific selected subject.
- Depending on the source of data, datamarts can be categorised as - independent or dependent.

- Independent - datamarts are sourced from data captured from one or more operational system or external info providers.
- Dependent - datamarts are sourced differently from enterprise dataware house.

- Independent Data Marts are source from the data captured from one or more operational System or external information providers.
- Dependent Data Mart are source directly from enterprise dataware house.

Date - 15/02/2023 :-

Curse of dimensionality :-

Dataware house,

*.) ~~OLTP~~ ^{source} primary objective ?

OLAP — Analysis.

→ It is used to produce the information required by top level in less time.

OLTP system

single group By operation - producing single 1-D cuboid.

multiple group By oper. — multi-cuboid.

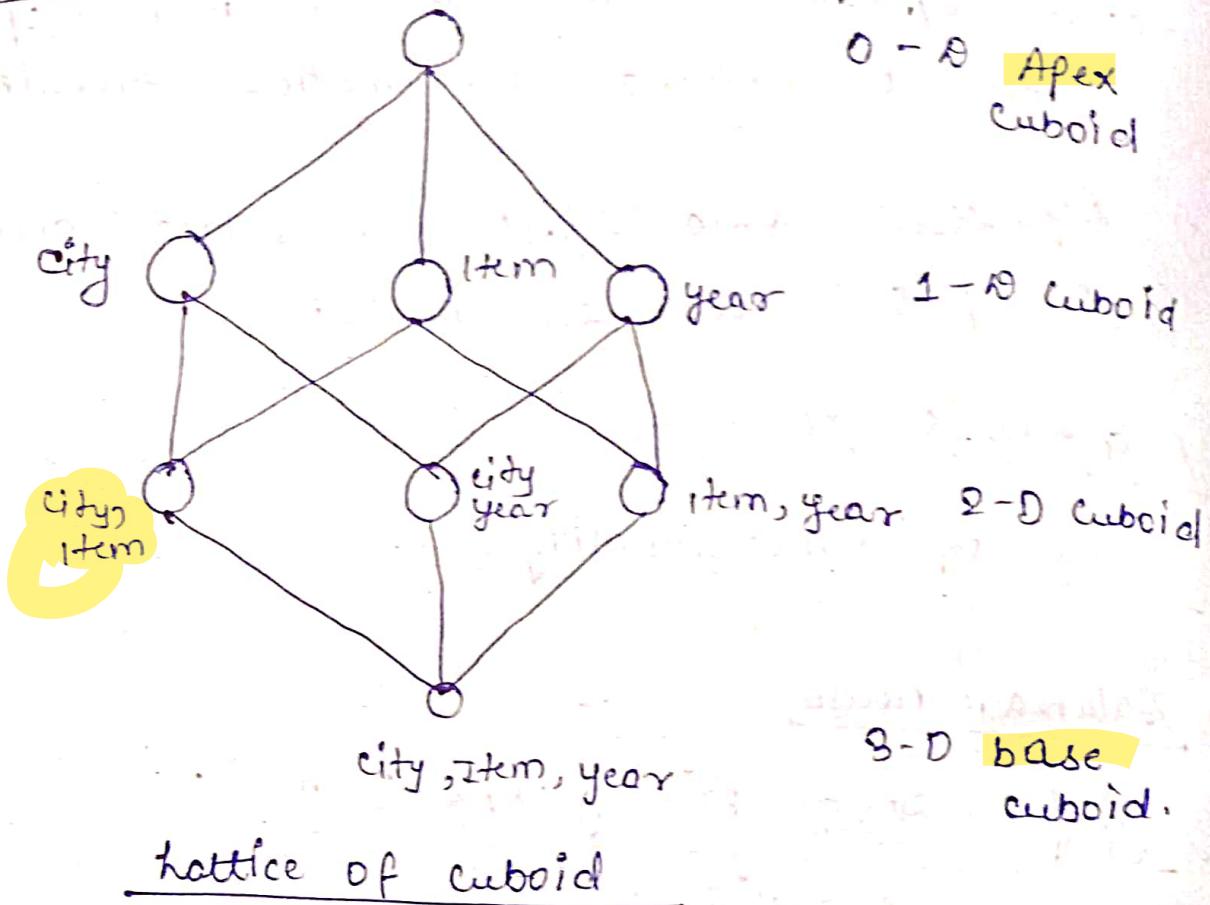
multi-cuboid — join — latent of cuboid.

curse of dimensionality →

no. of Aim \uparrow , no. of cuboid \uparrow

It is very difficult to store those.

Sales Data :-



- * The total no. of cuboid or group-by that can be computed for this data cube is 2^3 . (2^n)
- * These group-by forms a lattice of cuboid for the data cube.
- * The base cuboid contains all 3-dimension (city, item & year). It can return the total set for any combination of these 3-dimensions.
- * The Apex Cuboid / 0-D cuboid refers to the case where group-by is empty. It contains the total sum of all sales.

This pre-computation requires large amount of storage and it may explode when many of the dimensions have associated concept hierarchy.

For 'n' dimensional data cube, the total no. of cuboids that can be generated if the dimensions have hierarchies associated with it.

$$\text{Total no. of cuboid} = \prod_{i=1}^n (L_i + 1)$$

where L_i is the no. of level associated with dimension 'i'.

Ques.) If the data cube has 10 dimension and each dimensions have 5 levels associated with it, then find the total no. of cuboids generated.

Solution.

$$\begin{aligned} & \prod_{i=1}^n L_i + \prod_{i=1}^n 1 \\ & 10 \times 4 + 1 + 1 + 2 + 3 + 4 + 5 + 6 + 7 + 8 + 9 + 10 \times 10 \\ & = 40 + 1 + 10 \times 10 = 40 + 55 = 95. \end{aligned}$$

$$\text{no. of Cuboid} = \prod_{i=1}^{10} (4+1)$$

$$\begin{aligned} & = \prod_{i=1}^{10} 5 \\ & = 5 \times 5 \times 5 \times 5 \times 5 \\ & = \underline{\underline{5^{10}}}. \end{aligned}$$

Materialization of Cuboids :-

- No materialization
- Full , ,
- Partial , ,

* There are 3 choices for data cube materialization given a base cuboid.

1) No Materialization

→ Don't precompute any of the non-base cuboid but this may lead to extremely slow operations.

2) Full Materialization

→ Precompute all of the cuboids.

→ It typically requires huge amount of memory space in order to store all of the precomputed cuboids.

3) Partial Materialization

→ Selectively compute a proper subset of the whole set of possible cuboids.

→ It considers 3 factors :-

① Identify the subset of cuboids or sub-cubes to materialise.

② Use the materialised cuboids during query processing.

③ Efficiently update the materialized cuboid during load & refresh.

#. Association Rule :- Frequent pattern.

↳ Analyse market basket of customers.

↳ Association rule is usually written in form of
 $A \rightarrow B$

↳ support \rightarrow eg-2% (Mobile & tempered glass purchased)
↳ confidence \rightarrow 60% (if 10 mobile = 6 tempered).