

Data Mining Warehousing & Online Analytical processing.

Data warehousing :- decision support database - maintained separately - info processing providing a solid platform.

① Subject oriented :-

- * organized around major sub.

- * like → customer, product, sales.

- * focuses on **decision makers**

- * not on daily opt on transaction

- * subject issue → **structured** simple

- and **concrete view** (not)

- excluding not useful data.

② Integrated :-

- * combines multiple

- heterogeneous data sources**.

- [relational databases, flatfiles, etc]

- [online transaction records]

- * Data cleans, integrates.

- * converted data when

- moved to warehouse.

③ Non volatile :-

- * physically separate store - 2 transform

- * do not need transaction processing,

- recovery or concurrency control.

- * Needs - initial loading.

- access of data.

④ Time variant :-

- * Time horizon is longer

- * provides info from a

- historical perspective.

OLTP → (Database; operational) Online transaction processing

OLAP → (Warehouse) Online Analytical processing

Difference of

OLTP (Database)

vs OLAP (Warehouse)

① users

clerk, IT professionals

customer oriented

knowledge workers

market oriented

② function

day to day operations.

decision support

③ DB design

application oriented

subject oriented

④ data

current, upto date, detailed,
flat relational isolated.

historical, summarised,
multidimensional integrated,
consolidated

⑤ usage

repetitive

ad-hoc

⑥ access.

r/w.
index/hash on primary key

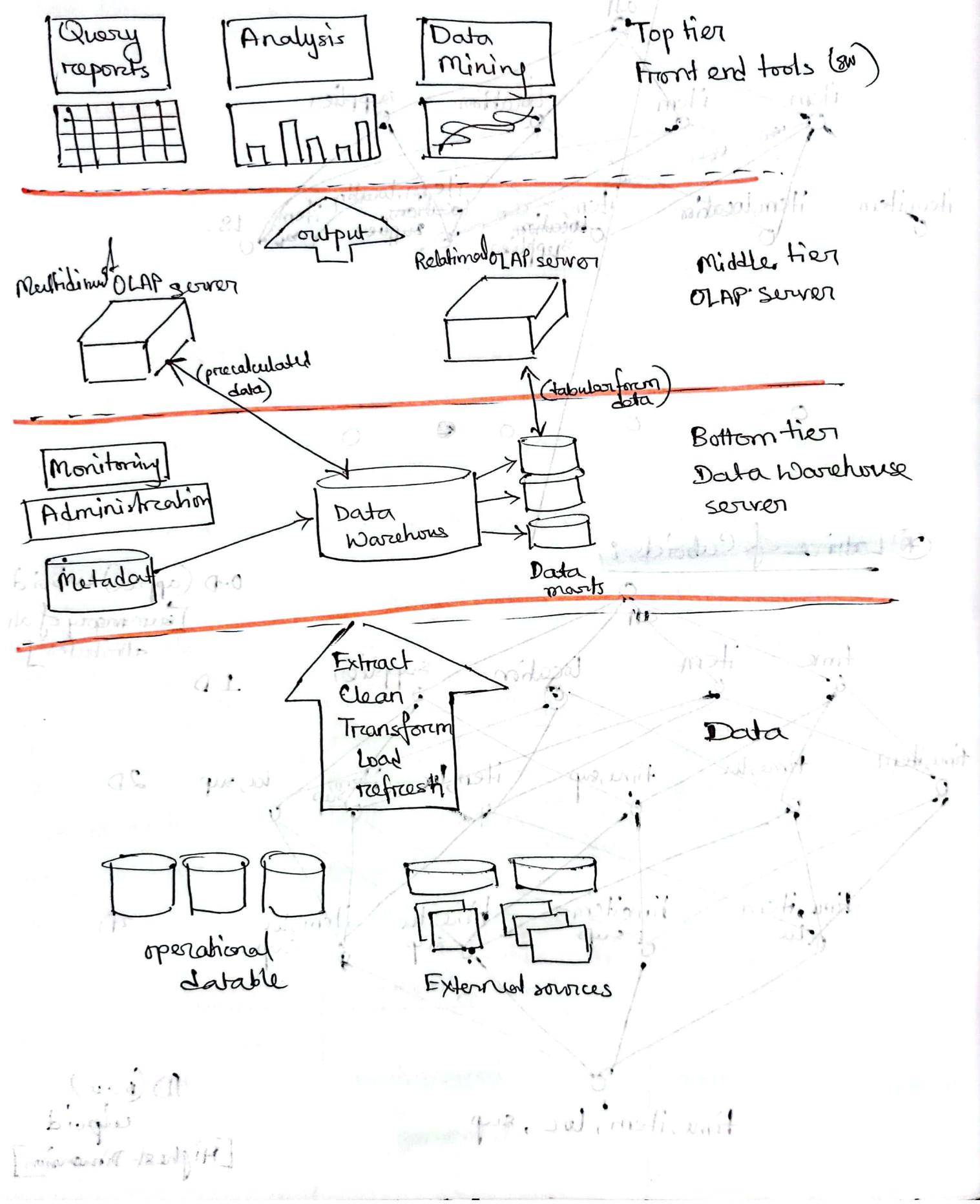
scans. read only

	OLTP	OLAP
⑦ unit of work	short and simple transaction	complex query
⑧ records accessed	tens 10s	millions M
⑨ users number	thousands 1000s	hundreds 100s
⑩ DB size	100 MB - GB	100 GB - TB
⑪ metric	transaction throughput	query throughput response
⑫ design	ERD (Entity Relationship model)	snowflake / star models
⑬ Access patterns	Concurrency control needed	No need, Read Only

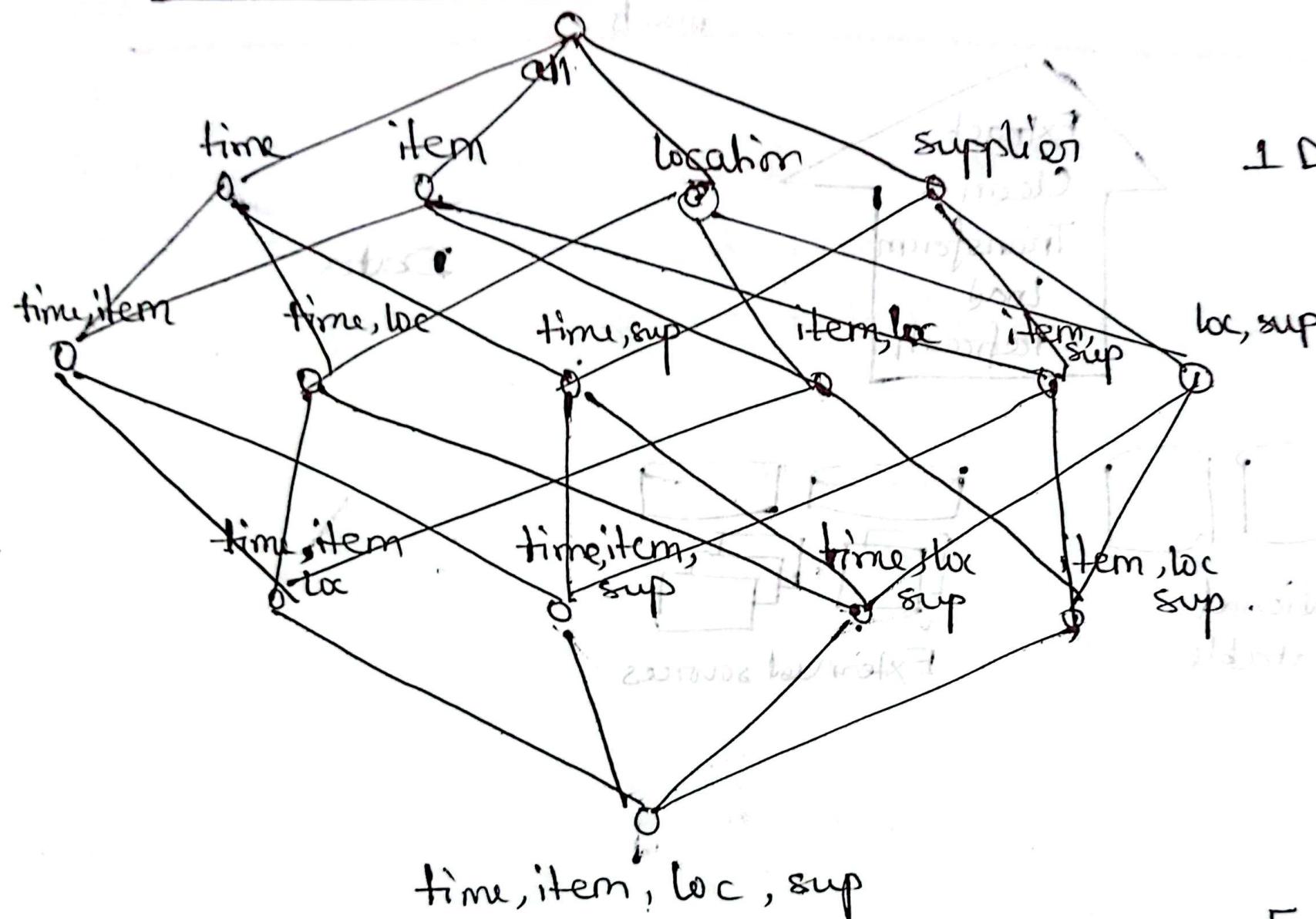
* Why separate Data Warehouse?

- ① For long term use
- ② Get historical data - no missing data.
- ③ DB - DB warehouse separate করলে data冗余 হবে যাবে
retrive করে computationally complex.
- ④ DB - ৩ (OLTP) detailed - ১ data রাখিব - not necessary.
- ⑤ Data access - ১ টা জন্য precalculation কোর্সে রাখিব calculate
করে access করা লাগবে but warehouse - ১ precalculated
data রয়ে গেলে possible so রাখিব calculate করা লাগবে না

Multi-Tiered Architecture (3tier)



Lattice of Cuboids:



0-D (apex) cuboid
[summary of all attributes]

1 D

2D

3D

4D (base)
cuboid
[Highest Dimension]

Computing number of cuboids.

each dimension $\rightarrow n$.

level $\rightarrow L$

$$\text{So, total number of cuboids} = \prod_{i=1}^n (L_i + 1)$$

Suppose a 3D cuboid id. \rightarrow location, item, time.

Location has attributes (district, division).

$$\text{So, } L_L = (2+1) = 3$$

Location (dist + div + all)

$$L_I = (1+1) = 2$$

item (item + all)

$$L_T = (1+1) = 2$$

Time (time + all)

$$\text{So total number of cuboids} = 3 \times 2 \times 2 = 12.$$

Group By:-

when attributes do not have sub attributes
so, 2D cuboid group by $= 2^n$ or $2^3 = 8$.

Schema :-

Star

- ① Central fact table + denormalized dimensions
- ② Flat, no hierarchy table
- ③ Few joins
- ④ Fast
- ⑤ redundant data
- ⑥ easy maintenance
- ⑦ Readable
- ⑧ Dashboards, BI tools

Snowflakes

- ① Central fact table + normalized dim
- ② normalized + split into subtables
- ③ more joins
- ④ slow
- ⑤ high (less redundancy) storage efficiency
- ⑥ complex
- ⑦ moderately readable
- ⑧ storage optimization, integrity focus.

Fact Constellation

- multiple fact table + shared dim
- shared multiple fact table
- many
- moderate
- med to high storage efficiency
- complex
- hard for beginners.
- Complex systems with multiple fact types

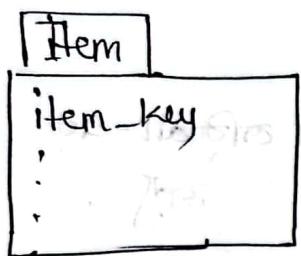
usage

simple + fast

reduce redundancy + normalize.

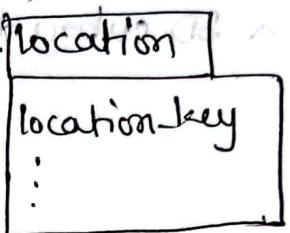
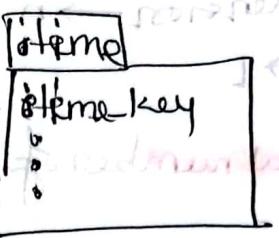
multiple types of facts.

* Star Schema

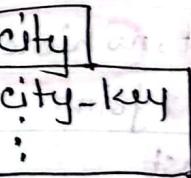
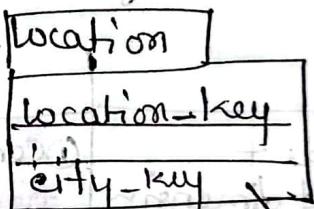
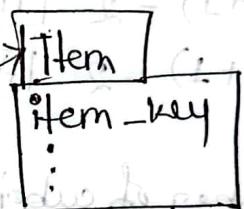
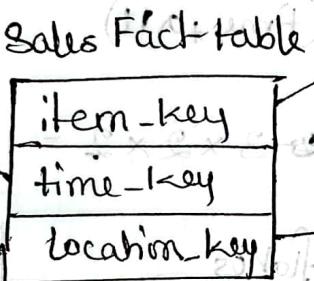
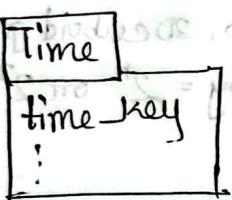


Sales Fact table

time-key
item-key
location-key
:



* Snowflake

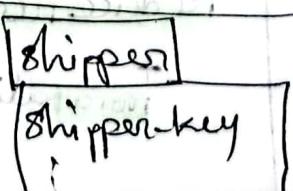
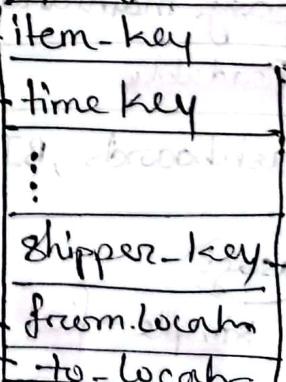
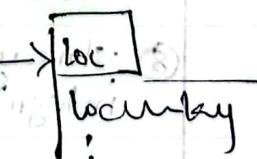
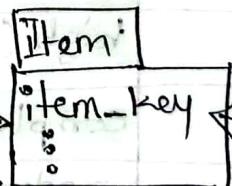


* Fact constellation :-



Sales Fact table

item-key
time-key
location-key
:



Data Cube: (Measures)

① Distributive: result derived by applying the fn to n aggregate values is same as that derived by applying the fn on all the data without partitioning.

Eg: count(), sum(), min(), max().

② Algebraic: if it can be computed by an algebraic fn with m arguments (m is a bounded int), each of which is obtained by applying a distributive aggregate fn

Eg: avg(), min-N(), standard-deviation()

③ Holistic:

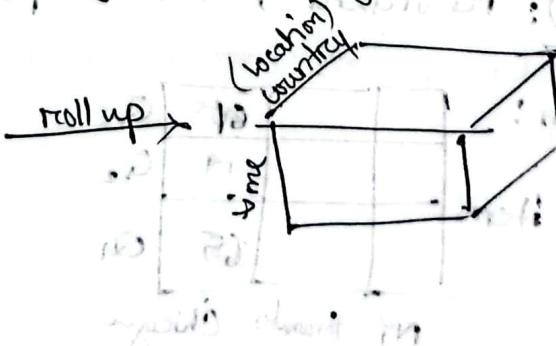
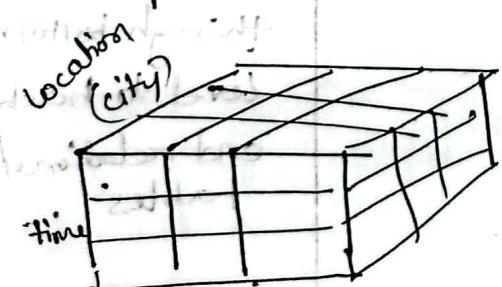
i) pending memory size : no constant bound on the storage size.

ii) Non distributive : partition wise median तर्क रखते join करने का correct result प्राप्त करना - full dataset - go to median नहीं आवश्यक

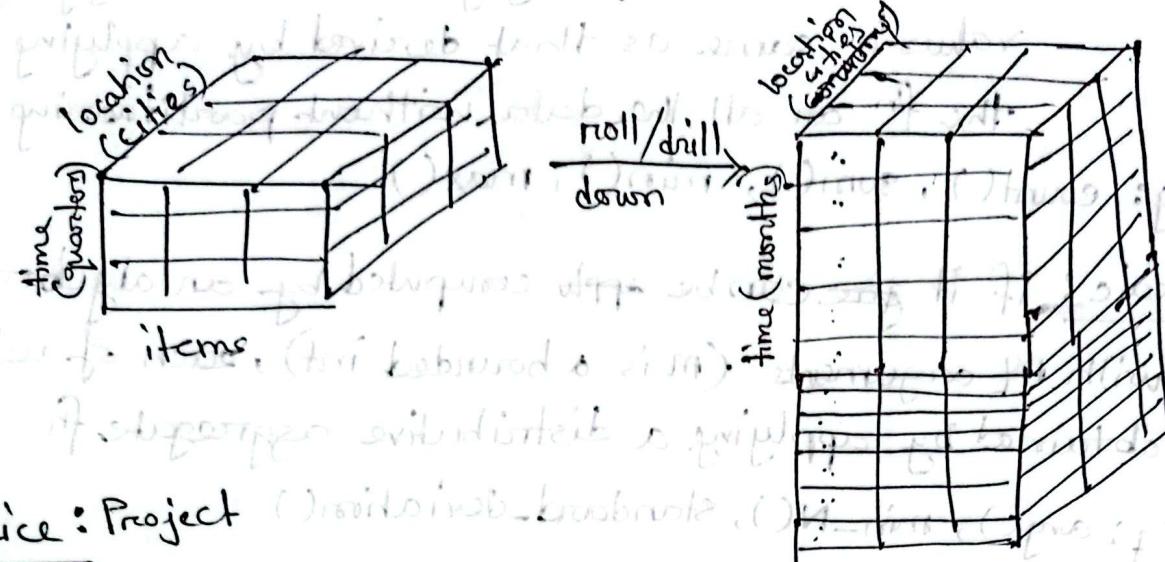
Eg: median(), mode(), rank().

OLAP operations :-

① Roll-up / drill-up : climb to higher hierarchy or dim reduction.



② Roll down / Drill down: introduce new dim, from high to lower level summarization



③ Slice: Project

	NY	65	
Toronto	45		
Chicago			
Q1	605	14	15
Q2			
Q3			

Items

Slice NY		3 cities	
Q1	Toronto	NY	Chicago
605	14	15	

④ Dice: Select

	Toronto	45	
Chicago			
Q1	605	14	
Q2			

⑥ Drill across:

more than 1 fact table

⑦ Drill through:

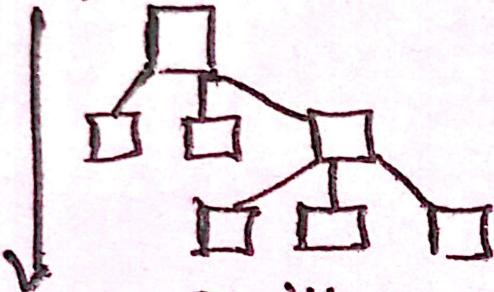
through bottom level to back end relational tables

⑤ Pivot (rotate): reorient, 3D to 2D planes

pivot to item:

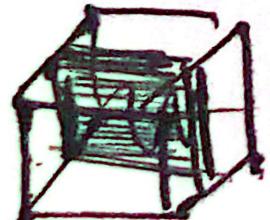
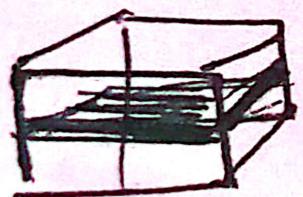
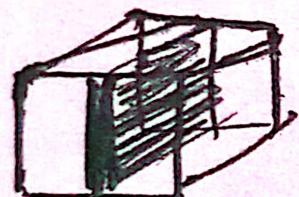
Items		605	Q3
		14	Q2
		65	Q1
NY	Toronto	Chicago	

drill down

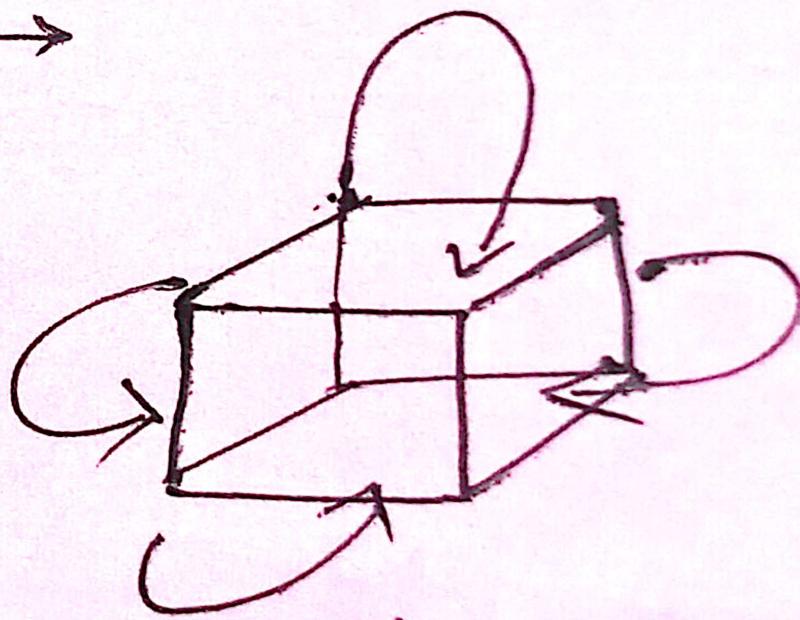


↑ rollup

Drill across



slice



Dicing.

Data Warehouse Design Process

- 2 Approaches:
- ① Top-Down :- starts with overall design + planning (modular)
 - ② Bottom-Up :- " " experiments + prototypes (rapid)

Methods :- (swe - 20. pov).

- ① Waterfall :- linear process -> stage - ২ টাকা কর্তৃত - each stage 100% complete না হলে, অন্য স্টেজ কর পাব না।
- ② Spiral : - short + quick turn around, interval time between stages is short so 100% stage complete এবং রাঠে জম্প করা মাত্র

Steps :-

- ① Choose a business process to model (কোন type business i.e. sales, marketing, orders, invoices etc.).

- ② Choose grain of the bp (data warehouse / fact table - ফিট ফর শিল্প)

- ③ Choose dimensions (attribute) for each fact table

- ④ " " measures that will populate each fact table.

Applications :-

- ① Info processing - graph, table, chart

- ② Analytical processing - Cube, ROLAP, MOLAP etc.

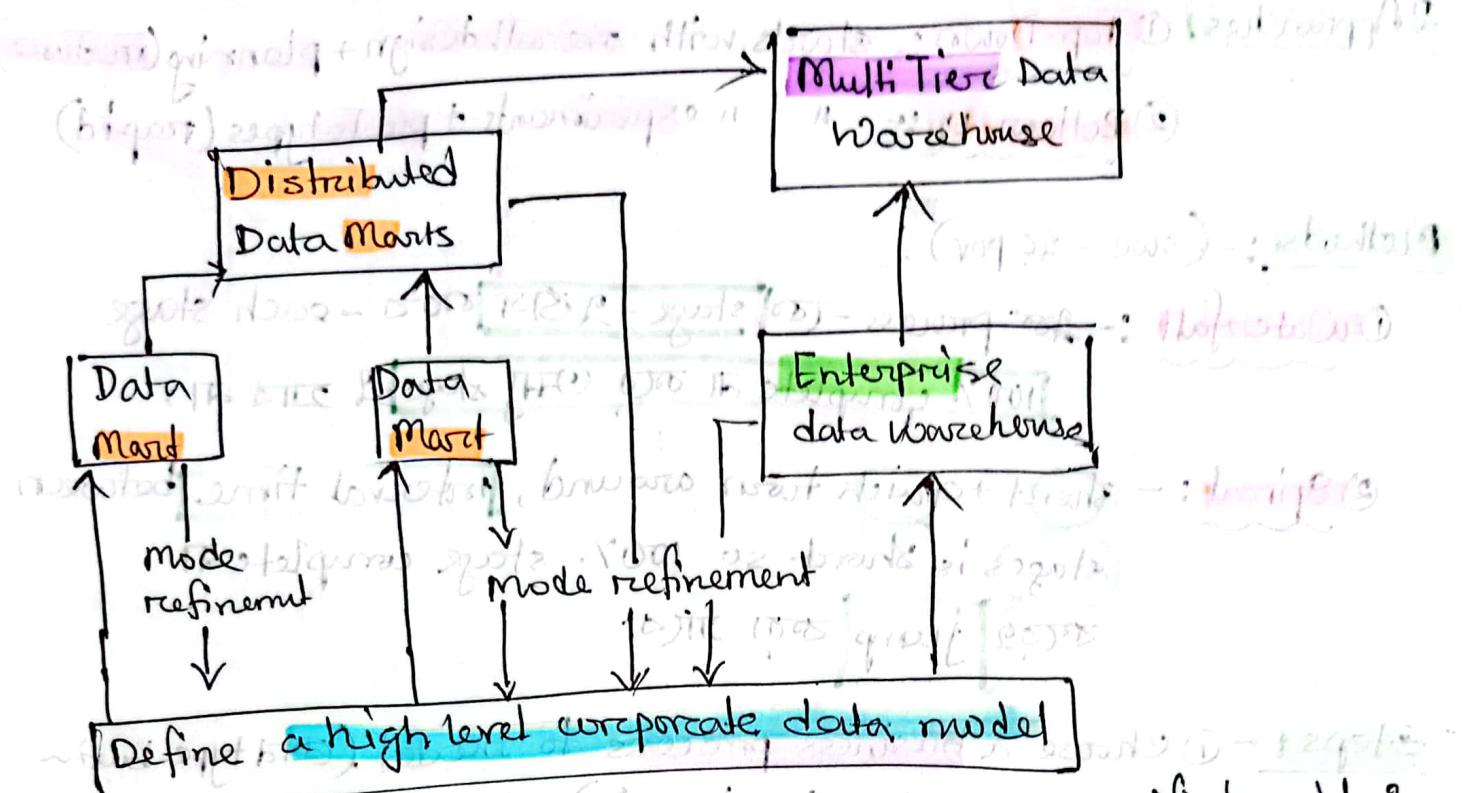
- ③ Data mining - new ফেলে info করা কর্তৃত, new pattern.

Information system এর উপর একটি উচ্চতা এবং গভীরতা আছে।

বিদ্যুতের ক্ষেত্রে বিভিন্ন পদক্ষেপ আছে।

১. Database

Datawarehouse Development : A Recommended Approach



(Different formats - unified model -
convert - does precalculated
aggregated value - for data cube)

Partial Materialization

- ① FULL Materialization :- সব cuboid তারে রাখা যাবে
- ② Partial :- কিছু "কিছু" রাখা যাবে না "not all"
- ③ No :- only main cuboid র রাখা

* যদি আর শান্ত রাখলে - operation র কাজ খুব সহজ
but storage waste & the computation time faster.

Compute Cube Operators : The operators used to op in the mathematical operation of data cuboids.
→ mean(), median() etc.

OLAP Indexing:

① **Bitmap indexing:** It is difficult to search a string amongst so much data. So we can represent a value using a bit creating separate columns for each value and assigning existence of the value row wise.

Example:

Base table

Cust	Region	Type
1	Asia	Retail
2	Europe	Dealer
3	Asia	Dealer
4	Russia	Retail

bitmap index according to Region

Rec ID	Asia	Europe	Russia
1	1	0	0
2	0	1	0
3	1	0	0
4	0	0	1

i = Asia
0 = Russia

bitmap

- * fast search
- * complex architecture.

according to type

	Retail	Dealer
1	1	0
2	0	1
3	0	1
4	1	0

② **Join indexing:** can be done in fact tables only.

of fact tables

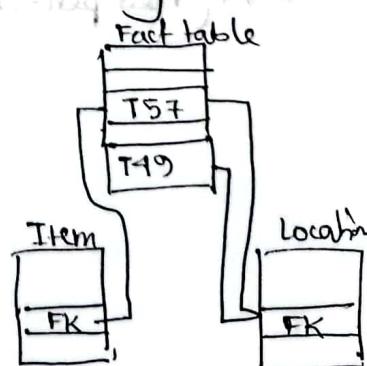
* Join the primary key with foreign keys of other tables.

* Fact table do not have details, so searching is easy.

Steps: * input a key into the system

** system searches the key in fact table

** that key is joint with the foreign keys of other table, thus system goes to the location.



④ Conversion from Data Cube to Table :

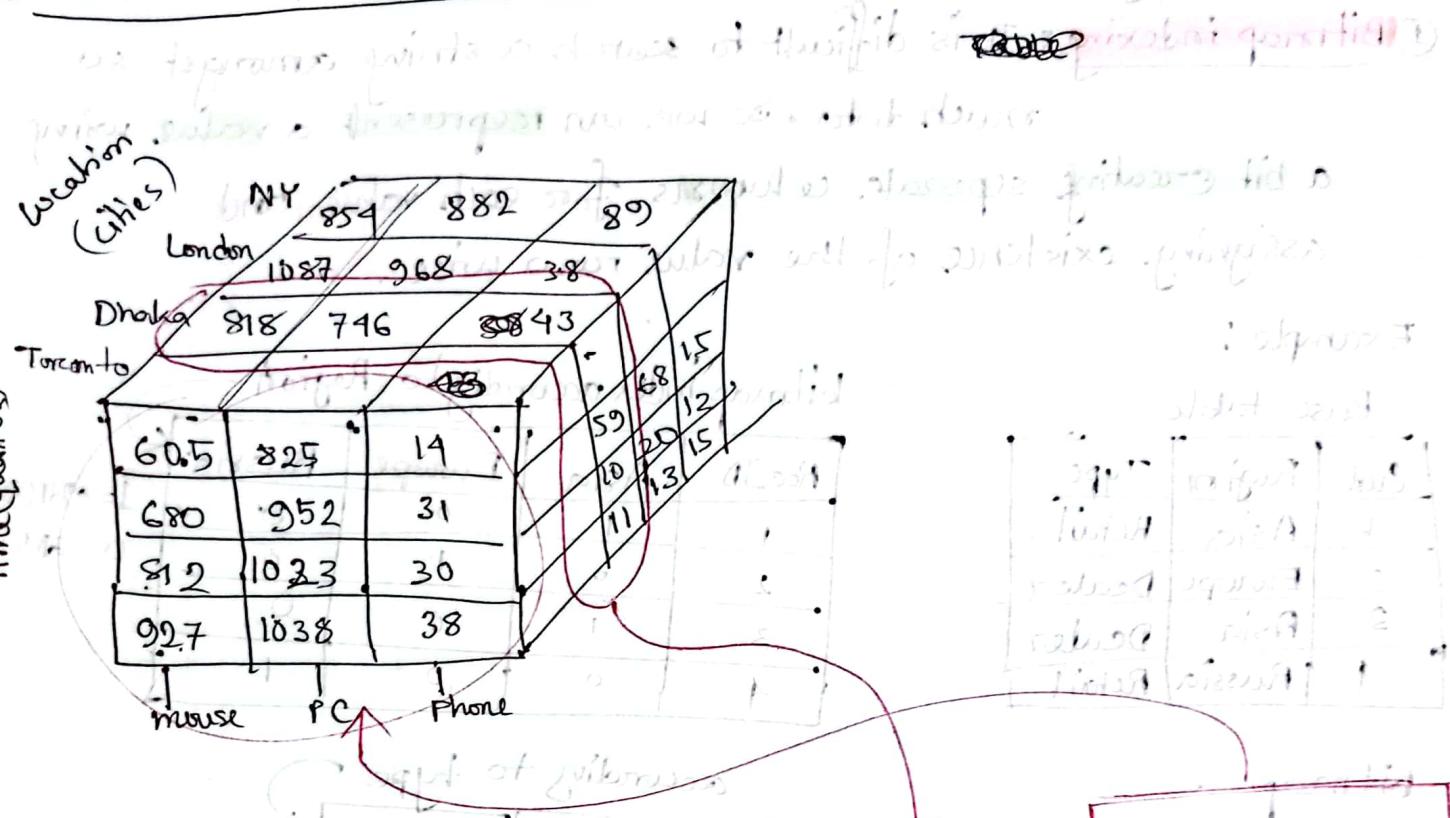


Table:

Time	location NY			London			Dhaka			Toronto		
	item	item	item	item	item	item	item	item	item	item	item	item
Q1	mouse	PC	Phone	mouse	PC	Phn	mouse	PC	Phn	mouse	PC	Phn
Q2	854	882	89	1087	968	38	818	746	93	605	825	19
Q3	634	712	15	1033	321	68	352	402	59	680	952	31
Q4	512	892	12	113	481	20	415	411	10	812	1023	30

* * * (Red pen - To see the cube not visible.)

Association Rule for Data Mining :-

Support $\rightarrow x\%$, there are $x\%$ case in the training data of and condition where $(A \wedge B) \rightarrow (A \vee B)$

Confidence $\rightarrow y\%$, there is a probability of $y\%$ that a new case meeting the condition A will cause the condition B which means $(A \rightarrow B)$.

Association rule :-

$$\text{Support } (A \rightarrow B) = P(A \cup B) = \frac{\text{count}(A \cup B)}{\text{Total no. of transactions}}$$

$$\text{confidence } (A \rightarrow B) = P(B|A) = \frac{\text{count}(A \cap B)}{\text{count}(A)}$$

Suppose, support = 2%
confidence = 60% for {milk, bread} \rightarrow table-1

Table

ID	Item (frequent)
1	{milk, bread}
2	{bread, milk, butter}
3	...
4	{bread, milk}
5	...
10	{milk, bread}

Total		Count		Count	
001	001	001	001	001	001
002	002	002	002	002	002
003	003	003	003	003	003
004	004	004	004	004	004
005	005	005	005	005	005
006	006	006	006	006	006
007	007	007	007	007	007
008	008	008	008	008	008
009	009	009	009	009	009
010	010	010	010	010	010
Total		Count		Count	
001	001	001	001	001	001
002	002	002	002	002	002
003	003	003	003	003	003
004	004	004	004	004	004
005	005	005	005	005	005
006	006	006	006	006	006
007	007	007	007	007	007
008	008	008	008	008	008
009	009	009	009	009	009
010	010	010	010	010	010

$$\text{confidence} = \frac{2}{10} \rightarrow (milk \rightarrow bread)$$

Binning: Partitioning large number of data points into separate bins

3 techniques:

① Equal Frequency : $\frac{[10, 11, 14]}{3}, \frac{[16, 17, 18]}{3}$

" width : $[10-14], [15-19]$

② By median : median ব্যাস দিয়ে সরবরাহ করো $[10, 12, 14] \rightarrow [12, 12, 12]$

③ By boundary : সূচী boundary-এ পরিষেবা

যদি যদি যমান $[10, 12, 15] \rightarrow [10, 10, 15]$

Correlation for Nominal Data : (Chi square).

Step 1: Get expected frequency

Step 2: $\chi^2 = \sum \frac{(Actual - Expected)^2}{Expected}$

Degree of freedom = $(row - 1) \times (column - 1)$

Significance level $\rightarrow 0.001$ or 0.01 (or given)

If χ^2 (chi-square) $>$ table of chi square (given), [Rejected]

	male	female	total
fiction	250	200	450
nonfiction	50	1000	1050
total	300	1200	1500

Expected value:

$$E_{male, fiction} = \frac{\text{Total}_{male} \times T_{fiction}}{T_{m,f}} = \frac{300 \times 450}{1500} = 90$$

$$E_{male, nonfic} = \frac{300 \times 1050}{1500} = 210.$$

$$E_{female, fic} = 360 ; E_{female, nonfic} = 840.$$

$$\chi^2 = \frac{(250 - 90)^2}{90} + \frac{(50 - 210)^2}{210} + \frac{(200 - 360)^2}{360} + \frac{(1000 - 840)^2}{840} \\ = 507.93$$

will compare χ^2 with given chi square (কোম্প হল rejected)
no correlation

Correlation for Numeric Data:

① covariance (কোভিয়েন্স)

② correlation " " = $\frac{\text{cov}(A, B)}{\sigma_A \cdot \sigma_B}$

$$\text{cov}(A, B) = \frac{\sum (A_i - \bar{A})(B_i - \bar{B})}{\text{Total}}$$

Examp:-

ID	A(\$)	B(\$)
1	\$ 6	20
2	5	10
3	4	14
4	3	5
5	2	5

existing position of liquid

$$[(p_1)_{\text{nw}} \quad (p_2)_{\text{e}}] \\ (p_1)_{\text{w}} \quad (p_2)_{\text{ne}}$$

$$(p_{11}) - (p_{12})_3 = (p_{11})_{\text{nw}}$$

$$\frac{p_{11}-p_{12}}{1-p_{11}} = (p_{11})^{n-1}$$

$$\frac{p_{11}-p_{12}}{1-p_{11}} = (p_{11})^{n-1}$$

covariance

$$E_A = \frac{6+5+1+3+2}{5} = \$14$$

$$E_B = \frac{20+10+14+5+5}{5} = \$10.80$$

$$\text{cov}(A, B) = \frac{(6 \times 20) + (5 \times 10) + (1 \times 14) + (3 \times 5) + (2 \times 5) - (14 \times 10.80)}{5}$$

$$\sigma_A = \sqrt{\frac{\sum (A_i - \bar{A})^2}{n}} = 2$$

and $\sigma_B = 32.56$.

$$\text{Corr}(A, B) = \frac{\text{cov}(A, B)}{\sigma_A \cdot \sigma_B} = \frac{7}{2 \times 32.56} = 0.107.$$

$\pi_{AB} > 0$; positively correlated (higher the stronger); $\pi_{AB} < 0$; negatively correlated.

PCA (Principal Component Analysis)

Given, $[x_1 \dots x_n]$ and $[y_1 \dots y_n]$

Step 1: Normalize (Z-score)

$$\text{① Real value} - \bar{x} \quad \text{mean} = \frac{\sum x_i}{n}$$

$$\text{② standard deviation} \quad SD = \sqrt{\frac{\sum (x_i - \bar{x})^2}{n-1}}$$

$$\text{③ Normalize, } z = \frac{x - \bar{x}}{SD} \quad \text{or, } z = \frac{x - \bar{x}}{\sigma}$$

$$\text{④ Normalized value} - \bar{z} \quad \text{mean, } \bar{z} = \frac{\sum (z_i)}{n}$$

Step 2: Covariance Matrix

$$\begin{bmatrix} \sigma^2(x) & \text{cov}(x,y) \\ \text{cov}(x,y) & \sigma^2(y) \end{bmatrix} \leftarrow$$

$$\text{cov}(x,y) = \frac{\sum (xy)}{n-1} - (\bar{x} \cdot \bar{y})$$

$$\sigma^2(x) = \frac{\sum (x^2)}{n-1} - (\bar{x})^2$$

$$\sigma^2(y) = \frac{\sum (y^2)}{n-1} - (\bar{y})^2$$

(x,y)	(x,x)	(y,y)
1,1	1,1	1,1
1,2	1,2	2,2
2,1	2,1	1,2
2,2	2,2	2,2

Step 3: PC1 অন্বে করা

$$C - \lambda I = 0.$$

$$\text{or, } \begin{bmatrix} \sigma^2(x) & \text{cov}(x,y) \\ \text{cov}(x,y) & \sigma^2(y) \end{bmatrix} - \begin{bmatrix} \lambda & 0 \\ 0 & \lambda \end{bmatrix} = 0.$$

$\lambda = ?$ (Highest value of λ is PC1).

Step 4: Projections.

$$\Leftrightarrow (C - \lambda \vec{I}) \vec{v} = 0$$

ultimately -

$$\begin{bmatrix} - & - \\ - & - \end{bmatrix} \begin{bmatrix} \otimes V_1 & V_1 \\ V_2 & V_2 \end{bmatrix} = \begin{bmatrix} 0 \\ 0 \end{bmatrix}$$

$$V_1 = ?$$

$$V_2 = ?$$

if $v_1 v_2 = -v_1 v_2$

if $v_1 v_2 = v_1 v_2$

$$\begin{bmatrix} t \\ -t_2 \end{bmatrix}$$

$$\rightarrow v_1 = \cancel{a} t$$

$$v_2 = -t_2$$

Let $r_1 = r_2 = t$

$$S_0, \vec{v} = ?$$

$$\hat{V} = \frac{\vec{V}}{|\vec{V}|} = ? \rightarrow \begin{cases} V_1 \\ V_2 \end{cases}$$

$$P_i = V_1 x_i + V_2 y_i$$

■ WVT (Wavelet Transformation)

JWT (Wavelet Transformation)
Dataset \rightarrow subset $\xrightarrow{\text{Core feature}}$ High frequency set (Core feature)
 \downarrow Low " "
Useless features removed

Given Dataset, $D = \{n, y, z, m\}$.

① Length = 2^n (padding ~~zero~~ if need)

② High freq set = difference of adjacent pairs $\rightarrow \{n-y, z-m, \dots\}$

$$\text{Low } f_q \text{ set} = \arg \min_{\{x_i\}} \sum_{i=1}^n \left| \frac{x_i+4}{2}, \frac{z+m}{2}, \dots \right|$$

③ will consider how far set again to separate:

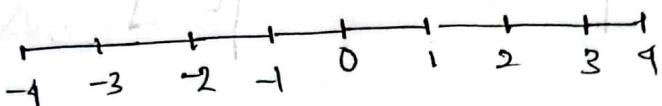
$$H. \text{ f.g set} = \left\{ \frac{x+y}{2}, \frac{z+m}{2}, \dots \right\}.$$

$$L \cdot f_{q \text{ set}} = \left\{ \frac{\left(\frac{n+4}{2} \right) + \left(\frac{z+m}{2} \right)}{2}, \dots \right\}$$

④ Result = { L₁, L₂, High frequency set }

$$= [\text{low fq set - } 20, \text{ first 20 value} + \text{high fq set - } 20, \text{ 20 value}]$$

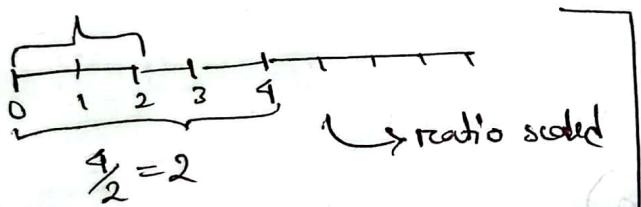
■ Interval Scale :- no true zero point in the difference of an attribute



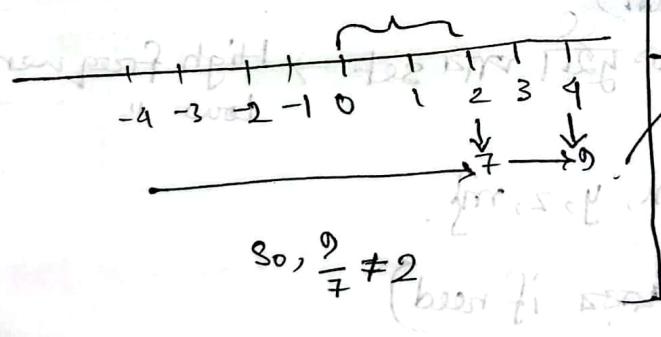
ex, 0°C ≠ no heat

OK = no heat

■ Ratio scale :- ratio calculation must be meaningful.



$$10^2 + 10^2 = 200$$



■ Mean = $\frac{\sum x_i}{n}$ ← using direct method

■ Mode : value frequent

$$\text{weighted mean} = \frac{\sum (w_i x_i)}{\sum w_i}$$

1 2 3 3 3, 4 5
unimodal

(sensitive to outliers)

1 2 3 3 4 4 5
Bimodal

Mode

Index 0 1 2 3 4
 Dataset length odd = $\frac{N-1}{2} \rightarrow 1, 2, 3, 4, 5, 6, 7, 8, 9, 10$. $\left(\frac{5-1}{2}=2\right)$

" " even = $\frac{N-2}{2}, \frac{N-2}{2} + 1 \rightarrow 1, 2, 3, 4, 5, 6, 7, 8, 9, 10$. $\left(\frac{6-2}{2}=2\right)$
 $\frac{3+1}{2} = 3.5$ median.

For table, median = $L_i + \frac{\left(\frac{N}{2} - (\sum f_{\text{eq}})_{\text{low}}\right)}{f_{\text{eq, median}}} \times W$.

$\frac{N}{2}$ = cumulative f_{eq} \rightarrow go to row is median row.

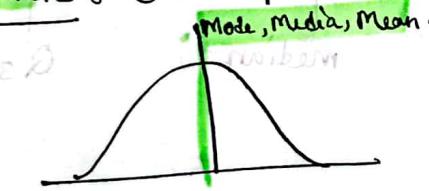
L_i = median row - 10 min range.

$(\sum f_{\text{eq}})_{\text{low}}$ = median - 10 cum. f_{eq} .

$f_{\text{eq, med}}$ = median - 10 f_{eq}

width = interval.

Symmetric: Central point with both side equal distribution



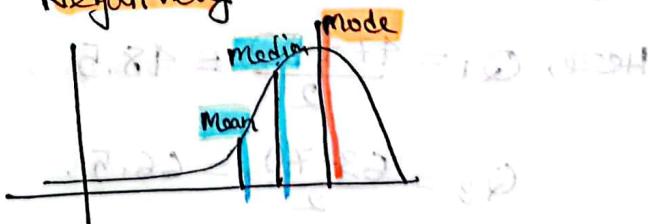
Qualitative two new strengthen rules

Asymmetric: Due to outliers

Positively skewed



Negatively skewed



Moderately skewed: outlier (out) not high or low (unimodal).

mean - mode ≈ 3 (mean - median)

$$3I = 3.81 - 3.00 = 12 - 3.63 = 9.38$$

$$3I = 3.81 \times 2.1 = 4.027 \times 2.1 = 8.46$$

Range: Subtraction of min-max of a dataset.

Quantiles: points that partitions dataset equally = the 4th of a dataset

2 Quantile \rightarrow $\frac{1}{\textcircled{1}} \frac{2}{\textcircled{2}} \frac{3}{\textcircled{3}} | \frac{4}{\textcircled{1}} \frac{5}{\textcircled{2}} \frac{6}{\textcircled{3}}$

3 Quantile \rightarrow $\frac{1}{\textcircled{1}} \frac{2}{\textcircled{2}} | \frac{3}{\textcircled{2}} \frac{4}{\textcircled{3}} | \frac{5}{\textcircled{3}} \frac{6}{\textcircled{4}}$

Quartile \rightarrow $\frac{1}{\textcircled{1}} \frac{2}{\textcircled{2}} \frac{3}{\textcircled{3}} | \frac{4}{\textcircled{2}} \frac{6}{\textcircled{3}} \frac{5}{\textcircled{4}} | \frac{7}{\textcircled{3}} \frac{8}{\textcircled{4}} | \frac{10}{\textcircled{4}} \frac{11}{\textcircled{5}} \frac{12}{\textcircled{6}} \frac{13}{\textcircled{7}}$

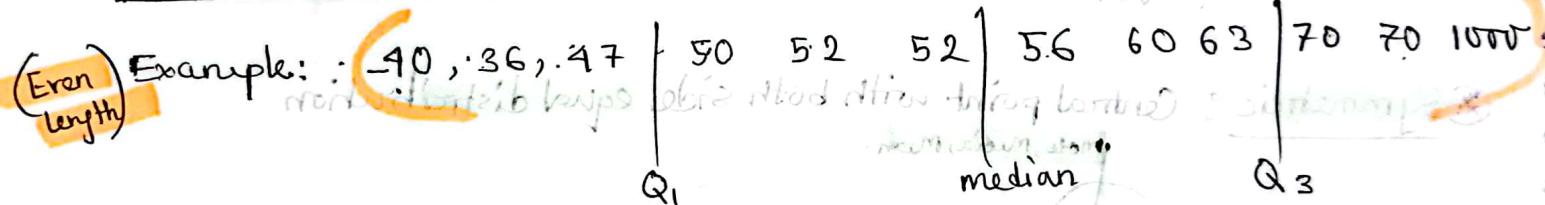
percentile \rightarrow 99 points that separates datasets into 100 equal portions

1 percentile = 1 quartile = 25 percentile

4 " = 100 percentile (= 1 percentile).

Inter Quartile Range:

$$IQR = Q_3 - Q_1 ; \text{ outliers } < 1.5 \times IQR \text{ (outliers)}$$



which datapoints are outliers?

(* sorted at first sort and after that)

$$\text{Here, } Q_1 = \frac{47+50}{2} = 48.5$$

$$Q_3 = \frac{63+70}{2} = 66.5$$

$$1.5 \times IQR = 1.5 \times 18.5 =$$

$$+ 5 \times IQR = 1.5 \times 66.5 = 99.75$$

$$IQR = Q_3 - Q_1 = 66.5 - 48.5 = 18$$

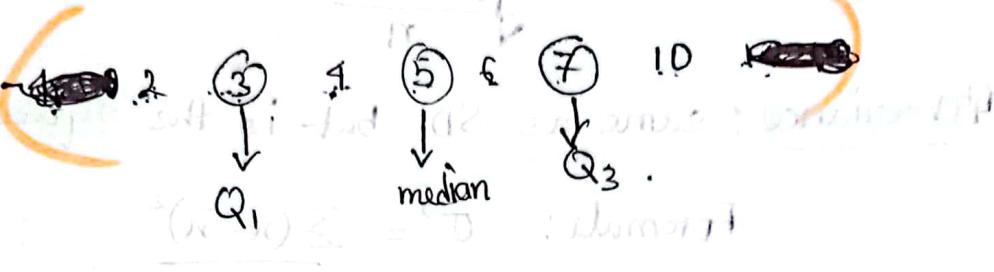
$$\text{So, } 1.5 \times IQR = 1.5 \times 18 = 27$$

$$Q_3 = 27 + 66.5 = 93.5 \quad (\text{more than } 93.5)$$

$$Q_1 = 48.5 - 27 = 21.5 \quad (\text{less than } 21.5)$$

Ans: outliers :- -40, 1000.

(odd length) Example:



$$Q_1 = 3$$

$$Q_3 = 7$$

$$\text{So, IQR} = 7 - 3 = 4$$

$$1.5 \text{ IQR} = 6$$

$$\text{Now, } Q_1 = 3 - 6 = -3 \quad \text{and} \quad Q_3 = 7 + 6 = 13$$

$$Q_3 = 7 + 6 = 13$$

so, value less than -3 and more than 13 are outliers.

~~Actual outliers~~ 2. (tr).

Q1, median, Q3 on sorted data.

Step-1 mark Q₁, Q₃, median on sorted data.

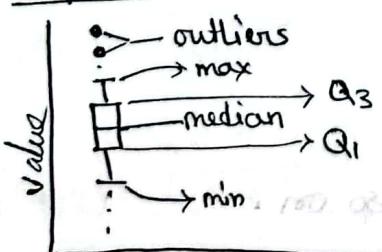
Step-2 IQR = Q₃ - Q₁ → box starts from here, $\frac{10-4}{2} = 3$ → ends at 7.

$$Q'_1 = Q_1 - 1.5 \text{ IQR}$$

$$Q'_3 = Q_3 + 1.5 \text{ IQR}$$

Step 3 find out less than Q'₁ and more than Q'₃ → outliers.

Boxplot : To visualize dataset - 20 value



minimum value = 0.06 and maximum value = 0.08

Standard deviation: Measures how spread out the number of data points are from the mean value

$$\sigma = \text{Formula: } \sqrt{\frac{\sum (x - \bar{x})^2}{n}}$$

Variance: same as SD but is the square of SD

$$\text{Formula: } \sigma^2 = \frac{\sum (x - \bar{x})^2}{n}$$

Proximity Measures: similarity/dissimilarity in the data.

$$\text{Data Matrix} = \begin{bmatrix} L_1 & L_2 & L_3 & \dots \\ S_1 & S_2 & S_3 & \dots \end{bmatrix}$$

Dissimilarity Matrix $x =$

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

* For Nominal attributes (colors, genders, bg, city, marital status etc.).

$$d(i,j) = \frac{P-m}{P} \quad \left| \begin{array}{l} P = \text{total attribute considered} \\ m = \text{no. of matched attribute} \end{array} \right.$$

Example:

ID	test 1	test 2
1	codeA	A
2	codeB	A
3	codeC	B
4	codeA	C

$$d(2,1) = \frac{2-1}{2} = 50\%$$

$$d(4,3) = \frac{2-0}{2} = 100\% = 100\%$$

thus $d(2,1) \Rightarrow 50\% \text{ dissimilar}$

$d(4,3) \Rightarrow 100\% \text{ dissimilar. and so on.}$

* For Ordinal (size, ratings, levels, etc) status etc) numbering $\frac{1}{2124}$

Step-1: Numeric representation.

Object	satisfaction level (test 2)	1
1	Excellent (3)	1
2	Fair (1)	0
3	Good (2)	0.5
4	Excellent (3)	1

Let, B fair = 1
good = 2

excellent = 3

Min Max Normalization
 $\frac{Y - \text{min}}{\text{max} - \text{min}}$

Step-2: Normalize $\Rightarrow Z_{i,f} = \frac{\text{Value}_i - 1}{\text{Max}_{i,f} - 1}$

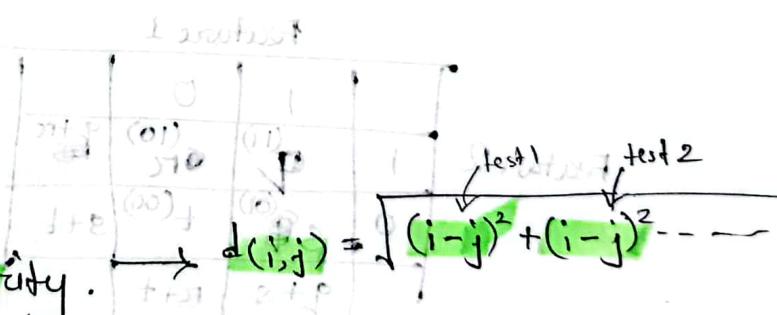
Here, Maximum feature, $\text{Max}_f = 3$.

$$\text{So, } Z_{1,2} = \frac{3-1}{3-1} = \frac{2}{2} = 1$$

$$Z_{2,2} = \frac{1-1}{3-1} = \frac{0}{2} = 0$$

$$Z_{3,2} = \frac{2-1}{3-1} = \frac{1}{2} = 0.5$$

$$Z_{4,2} = \frac{3-1}{3-1} = 1$$



Step-3: Calculate dissimilarity.

$$d(2,1) = \sqrt{(1-0)^2} = 1 \rightarrow 100\% \text{ dissimilar}$$

$$d(3,1) = \sqrt{(0.5-0)^2} = 0.5 \rightarrow 50\% \text{ dissimilar}$$

$$d(4,1) = \sqrt{(1-1)^2} = 0 \rightarrow 0\% \text{ dissimilar}$$

$$d(3,2) = \sqrt{(0.5-0)^2} = 0.5 \rightarrow 50\%$$

$$d(4,2) = \sqrt{(1-0)^2} = 1 \rightarrow 100\%$$

$$d(1,3) = \sqrt{(1-0.5)^2} = 0.5 \rightarrow 50\%$$

* For numeric : Each datapoint - rowwise distance কেবল কর্তৃত,

3 approaches:

$$① \text{Euclidean} : d(i,j) = \sqrt{(x_{i1} - x_{j1})^2 + (x_{i2} - x_{j2})^2 + \dots + (x_{ip} - x_{jp})^2}$$

② Manhattan : absolute value usage

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

③ Minkowski : generalized one.

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

if $h=1$, then Manhattan

$$\text{if } h=2 \text{ " Euclidean" } d(i,j) = \frac{s}{\epsilon} = \frac{1-\epsilon}{1+\epsilon} = \epsilon$$

* For binary :

Step-1 : Contingency table রীত।

		Feature 1		
		1	0	
Feature 2	1	q	r	$q+r$
	0	s	t	$s+t$
		$q+s$	$r+t$	

For symmetric attribute (ভাল state- $\frac{1}{2}$ important):

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

For asymmetric attribute (যদি একটা state নয় তবে না):

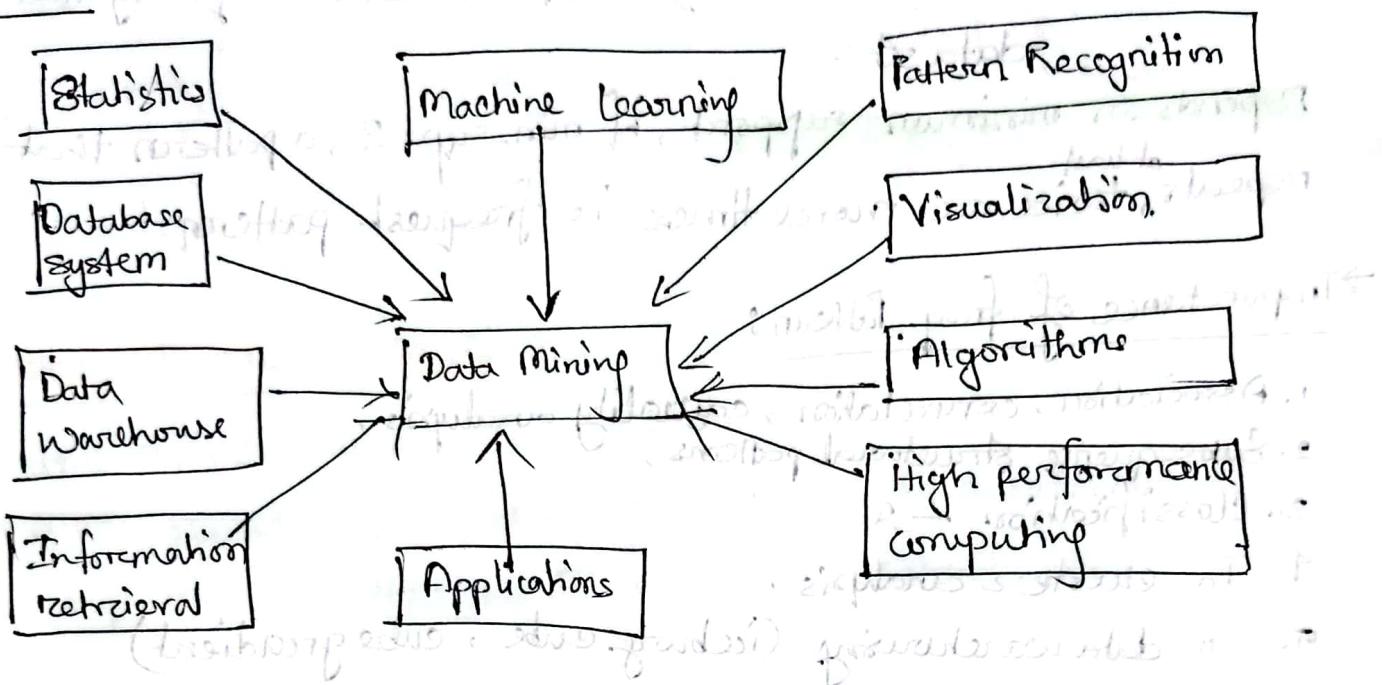
$$d(i,j) = \frac{r+s}{q+r+s}$$

$[t = (0,0) \text{ state নাইবে না}]$

Theory:

① "Data Mining adopts techniques from various domains"

Ans: Different areas best suited to data mining



Describe:

Statistics: Hypothesis testing, regression.

ML: Decision tree, clustering, neural nets.

Database Systems: Efficient data retrieval, storage.

Info retrieval: handle unstructured data

Pattern: complex data trends

visualize: meaningful representation

② Support $\rightarrow 5\%$. $(A \wedge B) \rightarrow C$. if $A \wedge B$, then C

Support $\rightarrow 5\%$ conf $\rightarrow 85\%$ why?

Support $\frac{5}{100} = 5\%$ condition true

* support $\frac{5}{100} = 50$ कर्तव्य $(A \wedge B)$ condition true

* conf $\frac{85}{100} \times 50 \approx 43$ कर्तव्य \rightarrow actually does C while $(A \wedge B)$ true.

(किमी 7 कर्तव्य $(A \wedge B)$ true but C true कर्तव्य)

Chap-6

Frequent Pattern: A pattern (a set of items, subsequences, substructures etc) that occurs frequently in a data set.

Depends on minimum support, if min-sup=2, a pattern that repeats at least twice or more times is frequent pattern

→ Importance of freq. Pattern:

1. Association, correlation, causality analysis.

2. Subsequence structural patterns,

3. Classification →

4. In cluster analysis.

5. " data warehousing (iceberg cube, cube gradient)

6. Semantic data compression.

7. Broad applications.

Frequent itemset Formulae:

$$\text{Support } (A \rightarrow B) = P(A \cup B) = \frac{\text{count } (A \cup B)}{\sum \text{transactions}}$$

$$\text{Confidence } (A \rightarrow B) = P(B|A) = \frac{\text{count } (A \cup B)}{\text{count } (A)}$$

Association Rule:

Support: Probability that a transaction contains AUB.

Confidence: Conditional probability that a transaction having A also contains B.

(B and C had count (A ∩ B ∩ C) = 0.7 and C has 0.9)

Example:

TID	Itemset
1	{bread, butter, banana}
2	{Mango, berry, banana}
3	{bread, butter, banana}
4	{mango, banana}
5	{bread, butter, oil}

$U = \{bread, butter\}$
 $A = \{bread, butter\}$
 $B = \{banana\}$.

$$\text{support} = \frac{c(A \cup B)}{\sum \text{trans}} = \frac{2}{5}$$

$$\text{conf} = \frac{c(A \cup B)}{c(A)} = \frac{2}{3}$$

(*) Closed Frequent Itemset :

* will consider support count

* यदि एक उपसेट का sup-count अपने उपसेट से अधिक है तो उपसेट नहीं।

sup-count - २०, similar इसीलिए, तो उपसेट नहीं।

Ex Suppose,

① {milk, bread, butter, oil} — 2 times (sup-count)

② {milk, bread, butter} — 3 times

③ {milk, bread} — 3 times.

उपर्युक्त ② उपसेट उपसेट ① का है।

again ③ उपसेट उपसेट ① और ② का है।

but ③ ② का sup-count समान है, इसीलिए हम ③ को रखते हैं।

Closed \rightarrow sub = sup(support count)

Maximal \rightarrow sup supercount \geq min-supcount

subset नहीं है

Maximum Frequent Itemset

- * will only consider minimum sup-count
 - * If superset not frequent than subset, then keep subset
 - * If " is frequent than subset, remove subset.

Super subset - \geq sup-count \geq min sup count; subset \leq

* Challenges of Apricori Method:

- ① Multiple scans.
 - ② Huge # of candidates.
 - ③ Tedious workload of support counting for candidates!

* Impròve Apriori

- Improving Apriori

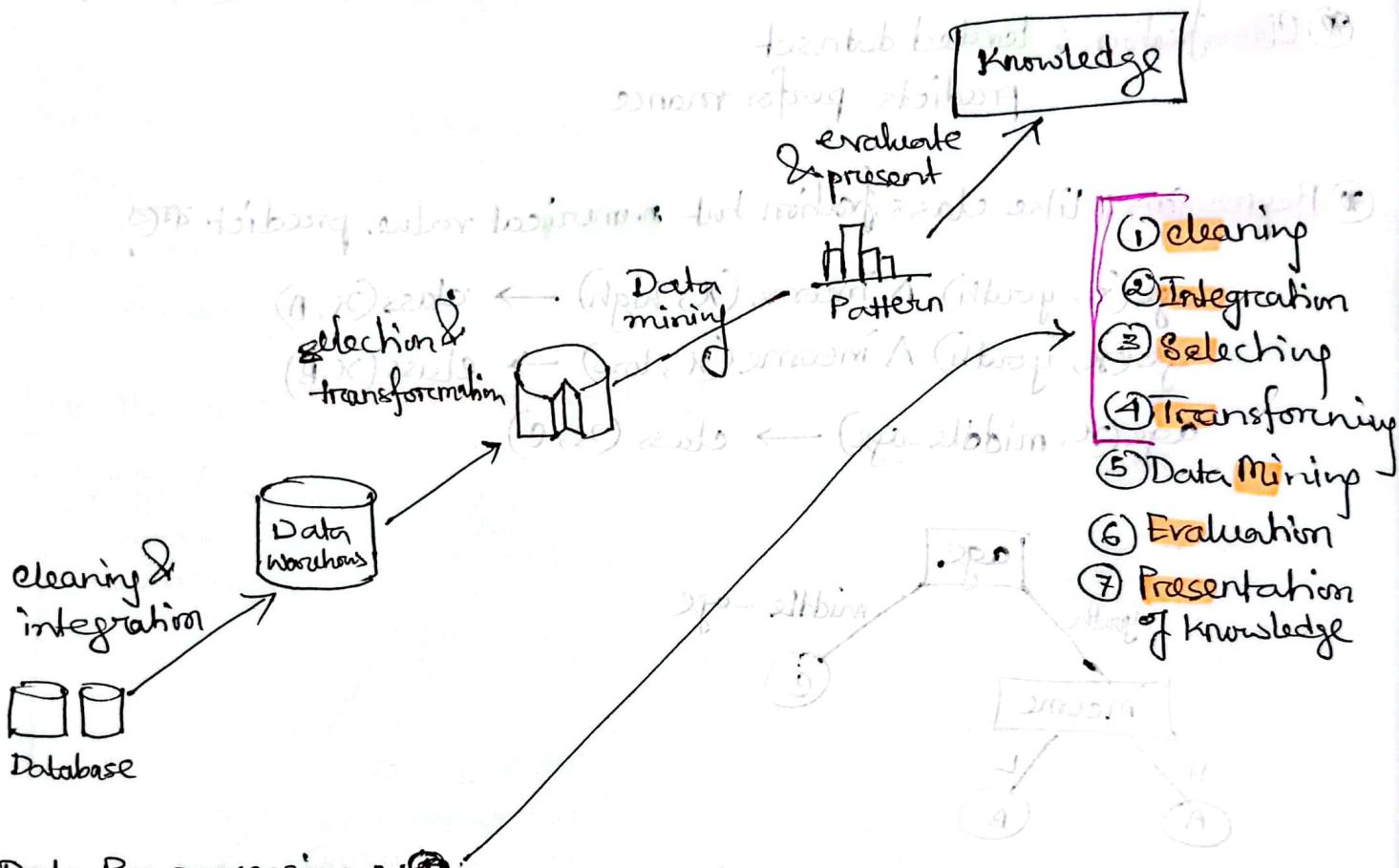
 - ① Reduce passes of transaction database scans (twice-local, global)
 - ② Shrink no. of candidates.
 - ③ Facilitate support counting of candidates.

(tristes tristes) que - des < bœufs
tristes que - des < bœufs que < bœufs

Chapter - 1

Data Mining: Extracting interesting patterns or knowledge from huge amount of data.

Knowledge Discovery from Data (KDD)



④ Data Pre processing :

④ What kind of data can be mined?

1. Data Streams / sensor data
2. Time series
3. Structure, graphs, networks
4. Obj-relational databases
5. Heterogeneous data.
6. Spatial - spatiotemporal data
7. Text
8. Audio, video, web.

④ Patterns that can be mined:

1. Generalization
2. Characterization, Discrimination
3. Frequent Pattern
4. Association & Correlation
5. Classification & regression,
6. clustering analysis
7. Outlier analysis.

- ④ Frequent patterns : ① Itemset $\rightarrow \{$ bread, butter $\}$ आवश्यक वाले विषयों का सेट
 ② Subsequence $\rightarrow \{$ PC की तरफ, पत्रों की OS विषयों की सौम्यता
 ③ Substructures $\rightarrow \{$ chemistry, geometrics etc.
 $2\text{H}_2 + \text{O}_2 \xrightarrow{\text{mostly}} \text{H}_2\text{O}$

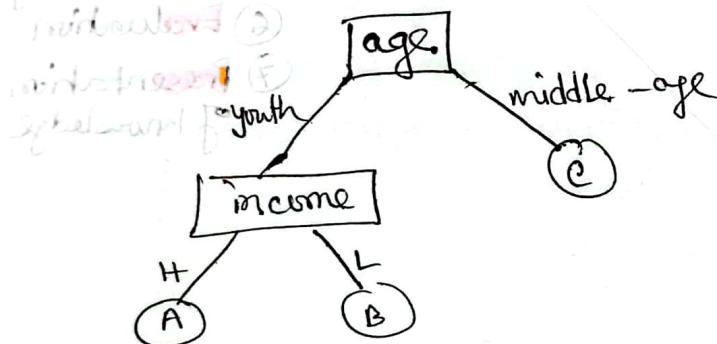
④ Classification : labeled dataset
 predicts performance

④ Regression : like classification but numerical value predict रखें

age(x, youth) \wedge income(x, high) \rightarrow class(x, A)

age(x, youth) \wedge income(x, low) \rightarrow class(x, B)

age(x, middle-age) \rightarrow class(x, C)



* too much branch = bad performance.

④ Cluster Analysis :

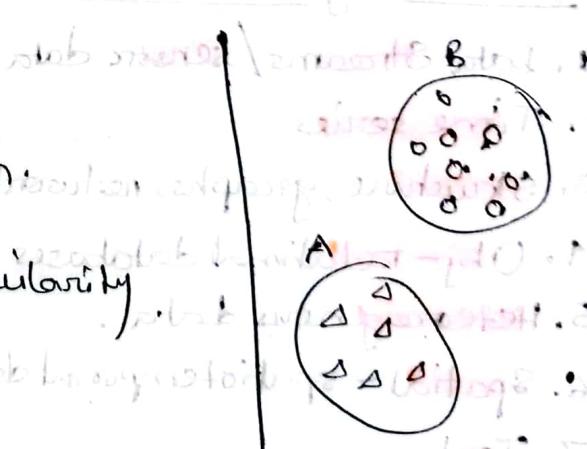
1. Unsupervised learning.

2. Groups plots in graph to find pattern.

3. Principle : maximize intra class similarity

minimize inter class similarity

4 Common character \rightarrow cluster datapoints.

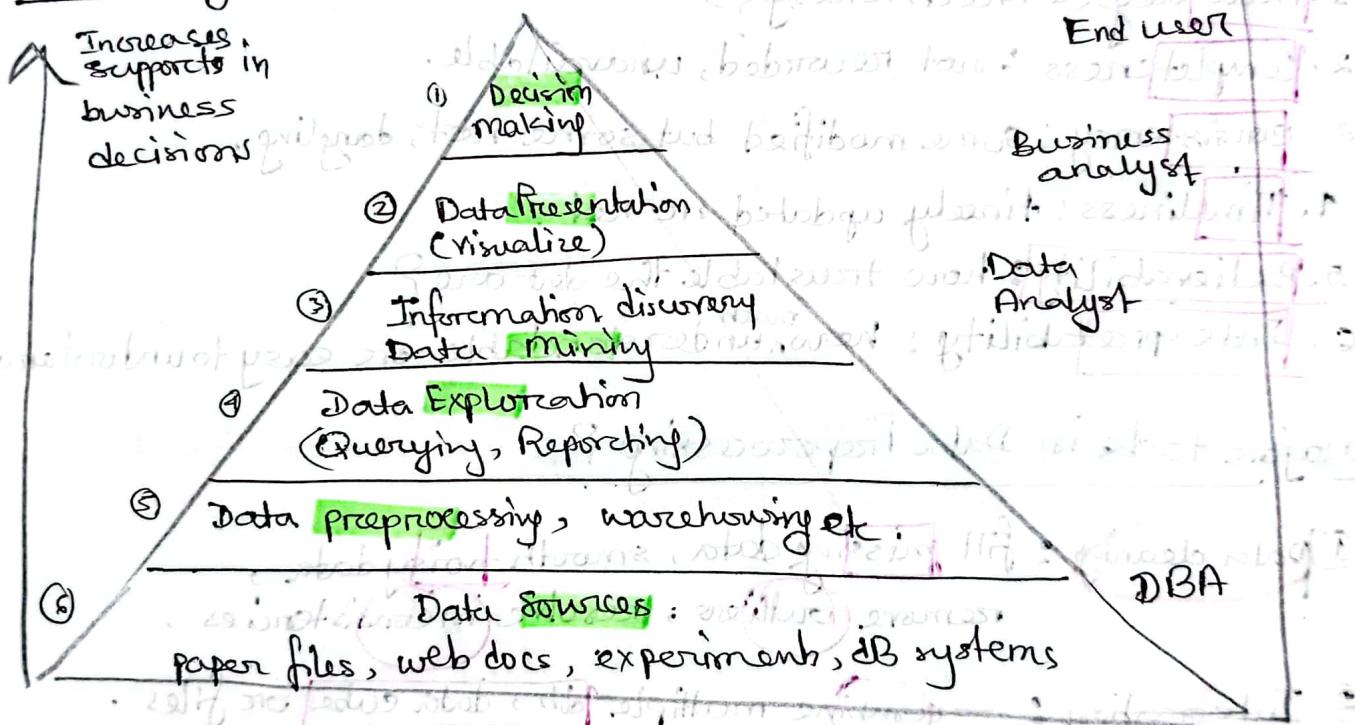


- * Outlier Analysis:
 1. data not complying with general behaviors
 2. Noise
 3. Exception.
 4. Methods: By product of clustering / regression analysis.
 5. Fraud detection, rare events analysis.

④ Are all patterns interesting?

- No
- interesting - when human understand it
- Valid on new data.
- potential.
- novel, unique.

⑤ Data mining In Business Intelligence



⑥ Medical Data Mining

- Statistics, ML
- preprocess data (feature extraction, dimension reduction)
- Classification, clustering
- post processing for presentation.

- Mining Methodology:**
1. Mining various, new kinds of knowledge
 2. Mining in multidimensional space (sub-attribute is not)
 3. An interdisciplinary effort
 4. Boosting the power of discovery in a networked environment
 5. Handling noise, uncertainty, incompleteness of data
 6. Pattern evaluation
 7. Pattern/constraint guided mining

Chapter-3

Measures for data Quality

- 1. Accuracy (correct/wrong?)
- 2. Completeness : not recorded, unavailable.
- 3. Consistency : some modified but some not, dangling.
- 4. Timeliness : timely updated or not.
- 5. Believability : how trustable the data are?
- 6. Interpretability : how understandable are easy to understand?

Major tasks in Data Preprocessing

- ① Data cleaning : fill missing data, smooth noisy data, remove outliers, resolve inconsistencies.
- ② Integration : combine multiple db, data cubes or files.
- ③ Data reduction : Dimensionality reduction
Numerosity
Data compression
- ④ Data Transformation : - Normalization
- Concept hierarchy generation.

Data Cleaning:

MISSING DATA

1. Ignore the tuple (not good)
2. Fill in manually (not possible for huge data)
3. Use global constant
4. " central tendency (mean / median / mode) attribute wise

5. Central tendency for same attribute and same class.
6. Most probable value (using regression model).

NOISY DATA: (Random errors or variance in a measured variable is called Noise.)

due to - faulty data collection instruments

- data entry problems
- data transmission
- technology limitation
- inconsistency in naming convention

- duplicate records

- incomplete "

- inconsistent "

Handle → 1. **Binning** ← partition into equal frequency.
smooth by bin means / median / boundaries.

2. **Regression**

3. **Clustering** - detect + remove outliers

4. Computer ft. human inspection
(detect roots check roots)

(detects bad data & handles it). need other methods like RIAA
and its variants, statistical methods

④ Data Integration:

Handle redundancy: correlation analysis = find patterns
catalog of strong, high covariance

Ⓐ Nominal $\rightarrow \chi^2$ -chi square value = $\sum \frac{(actual - expected)^2}{expected}$
(large = related)

Ⓑ Numeric \rightarrow Pearson's product moment coefficient.

$$r_{A,B} = \frac{\sum (a - \bar{a})(b - \bar{b})}{\sqrt{\sum (a - \bar{a})^2} \cdot \sqrt{\sum (b - \bar{b})^2}}$$
$$r_{A,B} = \frac{\text{cov}(A, B)}{\sigma_A \cdot \sigma_B} = \frac{\sum (A \cdot B) - \bar{A} \cdot \bar{B}}{\sqrt{\sum (A - \bar{A})^2} \cdot \sqrt{\sum (B - \bar{B})^2}}$$

$r_{A,B} > 0$, positive correlation

$r_{A,B} < 0$, negative correlation

$r_{A,B} = 0$, independent

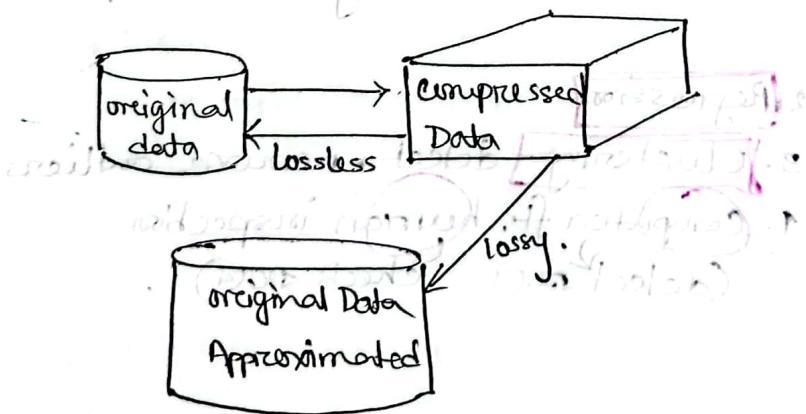
⑤ Data reduction:

① Dimensionality reduction:

1. Data compression:

a) Principle Component Analysis.

b) Wavelet Transformation.



2. Attribute subset selection. (Forward + Backward selection)

3. Attribute construction.

② Numerosity reduction :

a) Regression :- linear regression $[Y = wX + b]$

- multiple " $[Y = b_0 + b_1 X_1 + b_2 X_2]$

- log-linear models.

b) Histogram Analysis :

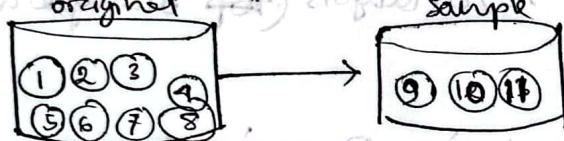
c) clustering

d) Sampling :- taking a set random sample from dataset Ω, N

(~~to increase speed, reduce overhead~~) ~~choose a random point~~ ~~choose a random point~~

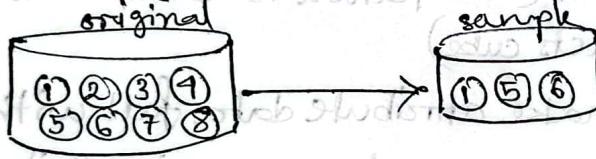
* SRWOR (Simple Random Sample without replacement)

→ selected object ~~is removed from the population~~

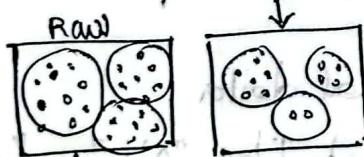


* SRWR (Simple Random Sample with Replacement).

→ selected object ~~not removed from the population~~



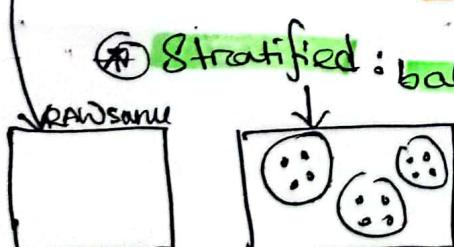
* Cluster :- main memory → process → part saved



- Full dataset is divided into ~~fold~~ and randomly sampled out.

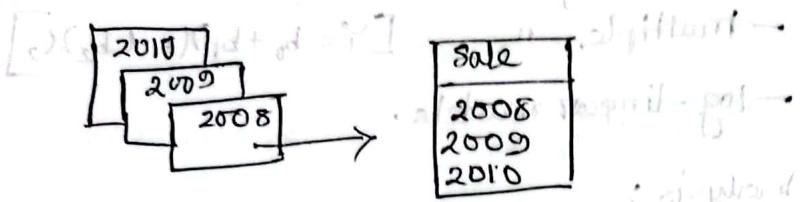
- SRWOR or SRWR → anyone could be used.

- imbalanced → 25% rate



* Stratified : balanced sample → each class same instance.

c) Data Cube Aggregate: Combine data into cube but may loss data



④ Data Discretization Methods:

1. Binning
2. Histogram Analysis (top down split, Unsupervised)
3. Clustering Analysis (bottom up merge, Unsupervised)
4. Decision tree analysis (split, supervised)
5. Correlation (e.g., χ^2) analysis (top unsupervised, bottom-up merge)

⑤ Data Transformation Strategies Overview:

1. Smoothing (remove noise: Binning, regression, clustering)
2. Attribute construction (new feature is constructed + added)
3. Aggregation (constructs cube)
4. Normalization (to make attribute data fall within small range)
5. Discretization (raw values of a numeric attribute are replaced by interval labels or conceptual labels.)
6. Concept hierarchy generation for nominal data ("streets" generalized to higher level concepts like "city" or "country")

■ Normalization:

$$\text{actual value} - \bar{x}$$
$$\frac{\text{actual value} - \min}{\max - \min}$$

* min-max

* Z-score

* decimal scaling.

* Ass. Rule:

$$\text{conf}(A \rightarrow B) = P(B|A) = \frac{\text{sup-ent}(A \cup B)}{\text{supent}(A)}$$

[y% chance of a new data of $A \rightarrow B$]

$$\text{support}(A \rightarrow B) = P(A \cup B) = \frac{\text{count}(A \cup B)}{\text{Total}} \quad \begin{array}{l} \text{In training data, nB case} \\ \text{among 100 sample of } A \rightarrow B \end{array}$$

* Apriori:

Base table

I ₁ , I ₂ , I ₅
I ₁ , I ₄
I ₂ , I ₃
I ₁ , I ₂ , I ₄
I ₁ , I ₃
I ₂ , I ₃
I ₁ , I ₃
I ₁ , I ₂ , I ₃ , I ₅
I ₁ , I ₂ , I ₃

$$S = 20\% \quad C = 80\% \quad n = 9$$

so, min-support $\left\lceil \frac{20}{100} \times 9 \right\rceil = 2$,

I₁ I₂ I₃ I₄ I₅

Step-1	
C ₁ (1 item)	(each count):
I ₁	6
I ₂	7
I ₃	6
I ₄	2
I ₅	2

$\rightarrow L_1$ pass count

Step-2 (2 items)

min-support

so, min-support

so, min-support

(2 <= count मान रखा)

Step-2 (2 items) (Pass eligible : I₁, I₂, I₃, I₄, I₅)

I ₁ , I ₂	9 ✓	I ₃ , I ₄	0
I ₁ , I ₃	9 ✓	I ₃ , I ₅	1
I ₁ , I ₄	1 ✗	I ₄ , I ₅	0
I ₁ , I ₅	2 ✓		
I ₂ , I ₃	4 ✓		
I ₂ , I ₄	2 ✓		
I ₂ , I ₅	2 ✓		

L ₂	
I ₁ , I ₂	9
I ₁ , I ₃	9
I ₁ , I ₅	2
I ₂ , I ₄	2
I ₂ , I ₅	2
I ₂ , I ₃	9

Step - 3 (3 items) (Eligible : I_1, I_2, I_3, I_4, I_5)

C₃

I_1, I_2, I_3	2
I_1, I_2, I_4	1
I_1, I_2, I_5	2
I_2, I_3, I_4	0
I_2, I_3, I_5	0
I_3, I_4, I_5	0

L₃.

I_1, I_2, I_3	2
I_1, I_2, I_5	2

I_1, I_2, I_3

I_n

I_1, I_2, I_5 eligible

I_2, I_3, I_5

ΣS

ΣS eligible

Step 4 - (4 items) (Eligible: I_1, I_2, I_3, I_5)

I_1, I_2, I_3, I_5	1 (not 2 or more)
----------------------	-------------------

thus for set = I_1, I_2, I_3
 I_1, I_2, I_5 .

1	2
F	F
2	3
C	C
S	F

Math ~~प्राची~~ - 9

Bin.

Bin1: 9 8 15

Bin2: 21 21 29

Bin3: 25 28 39

Bin mean:

Bin1: 9 9 9

22 22 22

29 29 29

Borders

Bin1: 9 9 15

Bin2: 21 21 29

Bin3: 25 25 39

1	2
3	4
5	6
7	8
9	10