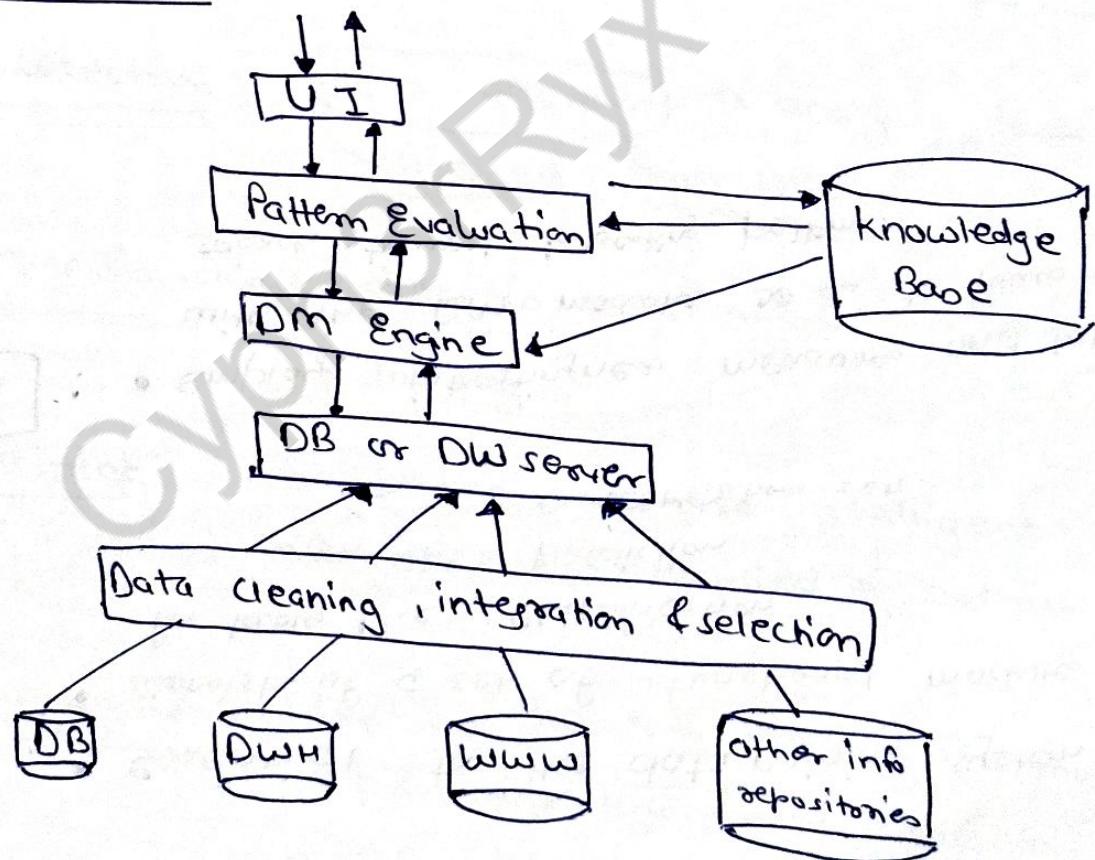


Annexure No :

DATA MININGData mining :

- ↳ Extraction of knowledge from large amount of data.
- ↳ referred as Knowledge Discovery from data.

Architecture of DM system

Knowledge Base

- Domain that is used to search or evaluate the interestingness of patterns.
- This contains user's previously used data and beliefs to assess a pattern's interestingness.

DM engine

- essential to the data mining system
- consists of a set of functional modules for tasks like Classification
Prediction
Characterization etc.

Pattern evolution

- employs interestingness measures and interacts with the DM modules so as to focus the search toward interesting patterns.

* KDD (Knowledge Discovery from Data)

→ Preprocessing operations.

↳ Required to make pure data in DWH before using it in DM processes.

- 1) Data cleaning: Remove noisy & inconsistent data.
- 2) Data integration: multiple sources are combined.
- 3) Data selection: data relevant to analysis task are retrieved from DB.
- 4) Data transformation: data is transformed into forms for mining.
- 5) Data mining: A process where intelligent methods are applied to extract data.
- 6) Pattern evaluation: identify the truly interesting patterns depicting knowledge based on some interestingness measures.
- 7) Knowledge Patterns: visualization and knowledge representation techniques are used to present the mined knowledge to the user.

Issues in DM:

① Performance:

- As data multiply these factors become critical
- the algos of DM to extract huge amount of data
- the run time & execution time of DM algo must be short and predictable.

② Parallel, distributed & incremental mining algo:

- Algo's partitions data into pieces

Each piece is processed in parallel by searching for patterns.
The patterns will eventually merge together.

③ Diverse DB:

(i) Handling of complex type data:

- one DM system can't mine all kind of data

(ii) Mining dynamics of global data repos:

- multiple sources of data are connected by internet forming a huge network of data.

(4) Mining methodology:

- (i) mining in multidimensional space
- (ii) mining various and new kind of knowledge
- (iii) handling noisy data
- (iv) Pattern evaluation

APPLICATION of DM:

- ↳ market analysis
- ↳ fraud detection & management
- ↳ customer retention
- ↳ quality analysis
- ↳ customer feedback analysis
- ↳ Risk analysis & management
- ↳ Target marketing
- ↳ Web Analysis
- ↳ Text mining from documents, emails, news letters.

Disadvantages of DM:

- PRIVACY**
- ↳ Collection of user data
 - ↳ data leak/breach
 - ↳ misuse of data collected
 - ↳ Business politics.

Classification of Dm:

- (i) Database Technology
- (ii) Statistics
- (iii) ML model
- (iv) Information Science
- (v) Visualizations
- (vi) Other disciplines

Sources of data:

- ① www
- ② organization
- ③ survey & services
- ④ Hard papers
- ⑤ IOT device

DM architecture:

(1) No coupling:

- ↳ don't use any service from DB or DWH
- ↳ data extracted from source
- ↳ Poor architecture for DM
- ↳ Processed → DM algos

(2) Loose coupling:

- ↳ uses services from DB or DWH.
- ↳ memory based DM,
- ↳ data extraction → Processing → DM algo → stores in system.

(3) Semi tight coupling:

- ↳ links its DB or DWH and uses its features like sorting and indexing for DM
- ↳ intermediate result can be stored in DB for better performance.

(4) Tight coupling:

- ↳ DB or DWH is treated as information retrieval piece.
- ↳ provides scalability, high performance & integrated information.

Business Intelligence:

set of methodologies

processes

architecture

technologies.

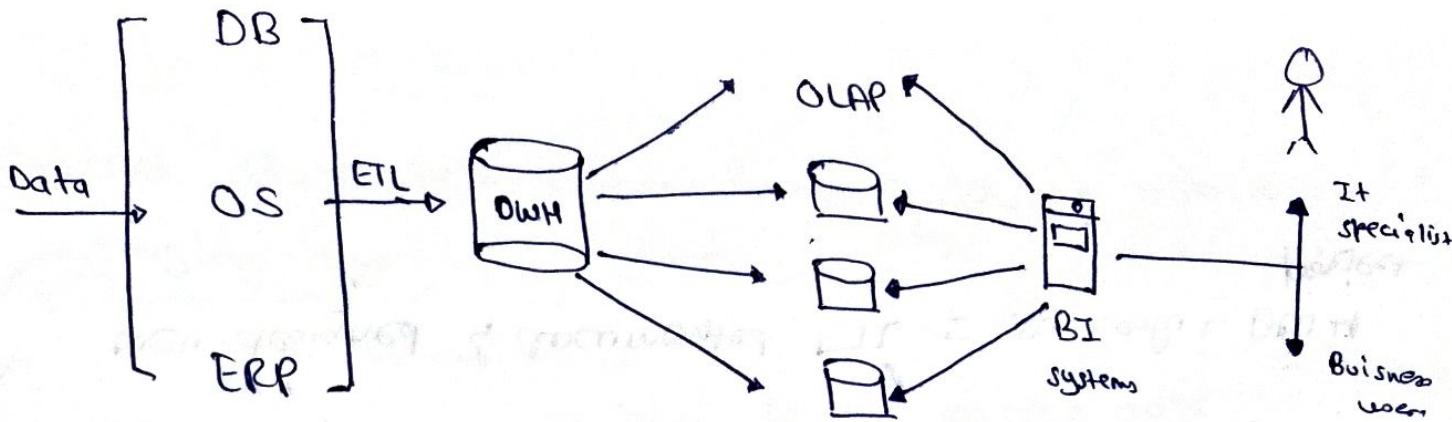
transform raw data
into meaningful
& useful
information

that allows
business users

to make

to achieve a
top position
in competition

efficient
business
decisions.



- Data from OS, DB & ERP were extracted & transformed into a more consumable form.
- Data from a DWH is then loaded to OLAP cubes.
- Data marts are also filled with data from DWH.
- OLAP cubes facilitates the analysis of data.
- With the BI system, IT specialist can setup the system for BI user so that they can analyze & access the data easily from BI system.

ETL Process:

Extract - Transform - Load

Process of data integration.

Takes up

large
volume

of raw data

Convert it
for

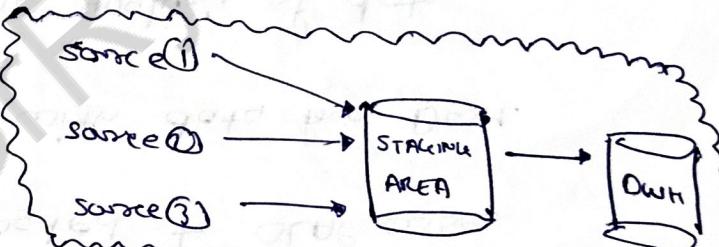
analysis

Data
to
DWH

LOADS

Why do we need ETL?

- ① Analysing data
- ② Transactional DB can't answer complex business questions
- ③ Provides method of moving the data from various sources into a DWH.
- ④ Data source changes → DWH automatically updated.
- ⑤ well designed & documented ETL = successful DWH project



Process ① : Extraction :

Extracted data sets comes from a source into a

Staging area

↳ acts as buffer

↳ b/w DWH & source data.

↳ also used for data cleaning & organization

↳ because data might corrupt
as it comes from various
sources.

Challenge:

↳ How ETL handles structured & unstructured data?

↳ Create a custom solution to assist in filtering
the unstructured data.

Process ② : Transformation :

- Data from multiple sources is normalized & converted into single system format.

Methods: Cleaning, filtering, sorting, splitting, joining

Process ③ : Loading

- Finally transformed data is loaded into DWH from staging area.
- Depending upon business need a data can be loaded in batches or all at once.

D W H

① Any type of data.

② Can expand capacity.

③ DM tools

④ Any format

⑤ All kind of data

⑥ Designed for analyzing data.

⑦ uses OLAP

⑧ Tables are denormalized

⑨ subject oriented

⑩ Techniques: Data Modelling

D B

① Structured data only

② can't expand capacity

③ No tools

④ Only tabular format

⑤ number of text

⑥ Designed for recording the data

⑦ uses OLTP

⑧ Tables are normalized

⑨ application oriented

⑩ Techniques: ER model

B I

DW/H

① Generating Business insights.

② O/p is data visualization, reporting & dashboards.

③ Audience : Manager & data analysts

④ Tools : dataline



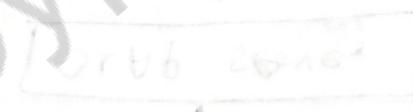
↓



↓



↓



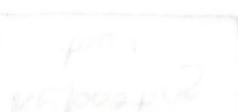
↓



↓

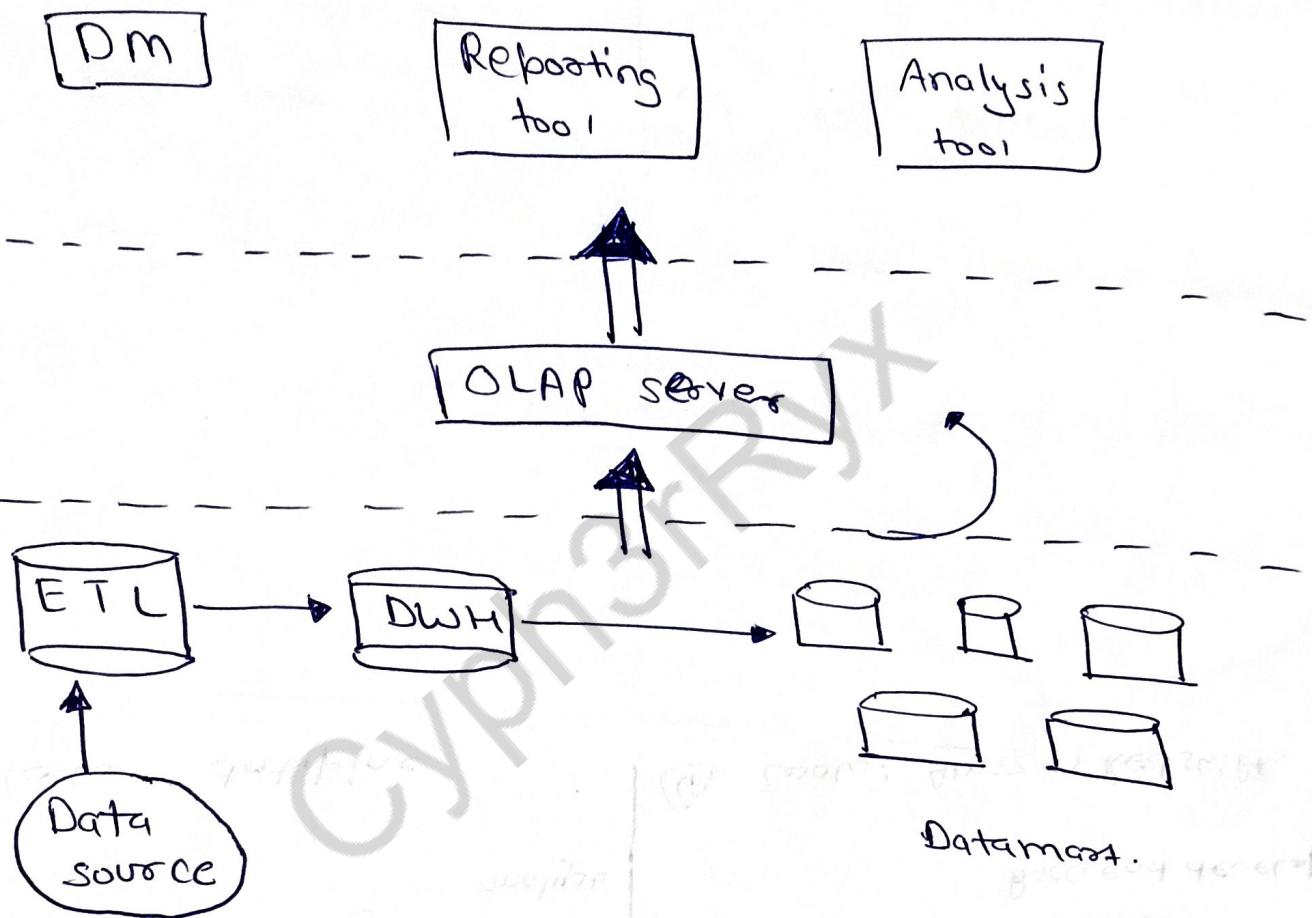


↓



DWH Architecture

3 tier architecture



Bottom tier consists of DWH, Data marts, ETL, source data.

Middle tier consists of OLAP server.

Top tier consists of Data mining, data analyzer & reporting tool.

OLAP server in middle is implemented using

① Relational OLAP model (ROLAP):

DBMS that maps fact on multidimensional data to standard relational operations.

② Multidimensional OLAP model (MOLAP):

directly implements multidimensional information & operations.



OLAP

Online Analytical Processing

a category of software tools which provides analysis of data for business decisions

Primary goal: Data Analysis

not processing.

e.g. Amazon analyse customer's buying list to give them a personalised homepage to retain them more.

* OLTP:

Online Transaction Processing

administer day to day transaction of an organization.

Primary goal:

Data Processing

not analysis

e.g. Sending a text message

or

Adding a book to shopping cart,

OLAP

- ① Analysis oriented
- ② Informational processing.
- ③ Accuracy is maintained of old data.
- ④ Focus on information out.
- ⑤ Uses DWH
- ⑥ Database Query management system
- ⑦ Subject oriented
- ⑧ Volume is stored in TB, PB
- ⑨ Processing time = long
- ⑩ Only read operation
- ⑪ Improves efficiency of business analysis.

OLTP

- ① Transaction oriented
- ② Operational processing
- ③ Only current data is maintained
- ④ focus on data in
- ⑤ uses DBMS
- ⑥ Database modifying system
- ⑦ Application oriented
- ⑧ Volume stored in MB, GB
- ⑨ Processing time = fast
- ⑩ Both read & write operation
- ⑪ Improves efficiency of user's productivity.

OLAP server operations:

- ① Roll up:
 - summarize data
 - by climbing up hierarchy or
 - by dimension reduction.
- ② Drill down:
 - reverse of roll up
 - from higher level summary to lower level summary
- ③ Slice & dice:
 - Project & select.

→ Slice is performed for the dimension from a given cube and provides a new subcubes.

→ Dice selects 2 or more dimensions from a given cube and provides a new subcube.
- ④ Pivot:
 - aka "rotation".
 - rotates the axes in view in order to provide an alternative presentation of data.

* Type of OLAP servers:

(1) ROLAP : Relational OLAP

↳ placed betⁿ relational back-end server & client front-end tools

↳ store and manage DWH, uses extended relational DBMS.
↳ it includes optimization of each DBMS backend & tools.

(2) MOLAP: Multi-dimensional OLAP

↳ array based multidimensional storage engines for multi.dim. view of data.
↳ storage utilization may be low because of MD stores.

(3) HOLAP: Hybrid OLAP

↳ combo of ROLAP & MOLAP

↳ offer high scalability & of ROLAP faster computation of MOLAP

(4) Specialized SQL servers:

↳ provides advance query language & query processing support for SQL queries over star and snowflake schemas

$$1. \underline{\text{Mean}} = \frac{(\text{total summation})}{n}$$

2. Mode = number with the highest frequency.

3. Median: Middle score that divides the total scores into two equal halves.

$$4. \underline{\text{Variance}}: S^2 = \frac{1}{n-1} \sum_{i=1}^n (x - \bar{x})^2$$

$$\text{S.D}: S = \sqrt{\frac{1}{n-1} \sum_{i=1}^n (x - \bar{x})^2}$$

Ex. 64, 70, 38, 96, 42, 47, 40, 39, 46, 50

x	$(x - \bar{x})$	$(x - \bar{x})^2$
64	-5.2	27.04
50	0.8	0.64
38	11.2	125.44
96	46.4	2100.24
42	-7.2	51.84
47	-2.2	4.84
40	-9.2	84.64
39	-10.2	104.04
46	-3.2	10.24
50	0.8	0.64

$$\boxed{\bar{x} = 49.2}$$

$$\boxed{\sum (x - \bar{x})^2 = 2600.4}$$

$$S^2 = \frac{2600.4}{9} = 288.7$$

$$S = \sqrt{288.7} = \sqrt{289}$$

$$= 17$$

$$\underline{\text{Mean}} = 49.2$$

$$SD = 17$$

$$49.2 - 17 = 32.2$$

$$49.2 + 17 = 66.2$$

Perks will spend in this range only. b/w 32.2 - 66.2 \$

* Data Preprocessing:

- ① Data Transformation: Transforms the received data in a form that it can be mined.
- ② Data Cleaning: Clean out noisy data & fill in missing values
- ③ Data Reduction: Obtain a reduced representation of the dataset that is much smaller in volume, yet maintains the integrity of original data.
- ④ Data Integration: Combines data from multiple sources into DWH.

DATA CLEANING:

- Real world data is noisy, incomplete & inconsistent.
→ fills data, make data consistent & remove noisy data.
- ① Missing values:

- ↳ ignore the tuple
- ↳ fill in missing values manually
- ↳ fill in the most probable value
- ↳ use another attribute which data is with you

Noisy data:

↳ Binning : Smoothes a sorted data value by consulting its neighbourhood values around it.

↳ Regression: A technique that conforms a data values to a fn.

↳ linear : fits 2 attribute in best line so that one attribute can predict other.

↳ multiple: fits more than 2 attributes & data is fit in multidimensional surface.

Outlier

Outlier analysis:

Outlier detected by clustering.

Similar values are organized in groups.

* Data Transformation:

data is transformed in a form such that it can be mined

(i) Smoothing:

↳ removes noise from the data

(binning, regression)

(ii) Aggregation:

↳ summary or aggregate values are applied
(constructing a data cube)

(iii) Generalization:

↳ low level concepts are replaced with high level concepts.

e.g. street → city / country

(iv) Normalization:

↳ Attribute values are normalized by scaling their values so that they fall in specified range

e.g. 2
40
500
1
3
900

All the values are different and nowhere near to each other so if we normalize the data in such a way that its range becomes (0-1).

Method: 1 Min. - max. Normalization

$$v' = \frac{v - \min(x)}{\max(x) - \min(x)}$$

v = original attribute value

$\min(x)$ = min. value of all the elements = 1

$\max(x)$ = max. value of all the element. = 900

for e.g. we need v' of 2:

$$v' = \frac{2 - 1}{900 - 1} = \frac{1}{899} = 0.001112$$

In this way we can normalize the data for all elements.

Method - (II) Z-score normalization:

(Z-score mean normalization)

$$Z = \frac{X - \bar{x}}{SD_x}$$

$$SD = \sqrt{\text{Variance}}$$

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

x	$x - \bar{x}$	$(x - \bar{x})^2$
2	-239	57,121
40	-201	60,401
500	259	67,081
1	-240	57,600
3	-238	58,164
900	654	434,281

$$SD = \sqrt{\frac{713124}{n-1}}$$

$$= \sqrt{118,854.67}$$

$$SD = 344.75$$

$$\bar{x} = 241$$

$$\sum (x - \bar{x}) = 713124$$

$$Z_1 = \frac{2 - 241}{344.75} = \frac{-239}{344.75} = -0.693$$

$$Z_2 = \frac{40 - 241}{344.75} = \frac{-201}{344.75} = -0.583$$

$$Z_3 = \frac{500 - 241}{344.75} = \frac{259}{344.75} = 0.7512$$

$$Z_4 = \frac{1 - 241}{344.75} = \frac{-240}{344.75} = -0.696$$

* Data reduction:

↳ to get a reduced representation of data set
that is much smaller in volume.

(1) Data cube aggregation:

Aggregation opⁿ are applied to the data in the construction of a data cube.

(2) Attribute subset selection:

irrelevant attributes are removed & relevant are selected.

(3) Discretization:

raw data values are replaced by ranges or higher conceptual levels.

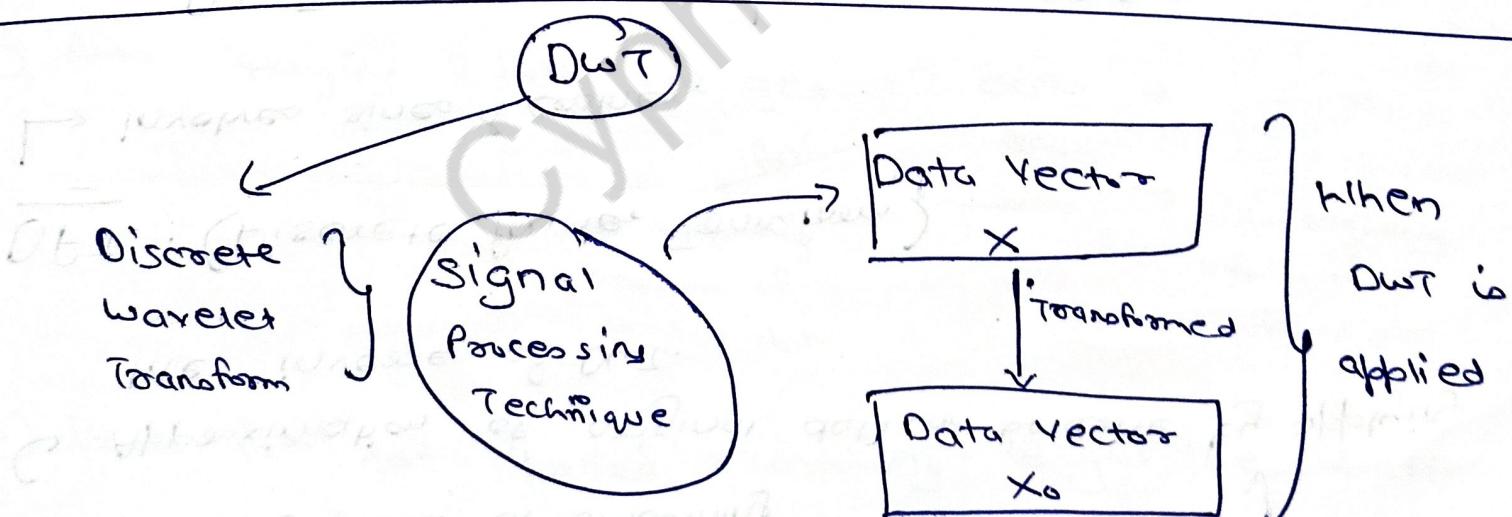
Annexure No.:

DMBI

- Data cleaning
- Data Transformation
- Data Reduction
- Data Integration

Data Cleaning

- fill missing data
- smooth & noisy data
- resolve data inconsistencies



Both X & X_0 must have same length

Advantages:

- ↳ resultant data is very sparse.
- ↳ computation of data is very fast.
- ↳ Able to remove noise from the data.
- ↳ reduces task of smoothing.
- ↳ Approximation of original data can be done by applying the inverse of DWT.

DFT : (Discrete Fourier Transform)

- ↳ involves sines & cosines.

DWT

- DWT is independent
- DWT gains lossy compression.
- Provides accurate wavelet coefficients.
- Occupies less space

DFT

- DFT is based on the results of DWT
- DFT don't get lossy compression
- don't provide accurate wavelet coefficients.
- Occupies more space

Annexure No :

Hierarchical Pyramid Algorithm:

- ↳ halves data at each iteration.
- ↳ the data is reduced to half so speed of computation increases.

Method:

- ① input vector : L ($L = \text{integer & power of 2}$)
- ② apply 2fx^n for each transform
 - ↳ first performs data smoothing
 - ↳ second finds weighted difference
- ③ After applying 2fx^n we get 2 data set
 - ↳ first → low frequency version of original data
 - ↳ second → high frequency version of original data
- ④ Both fx^n s are applied recursively to get the resultant vector of length 2
- ⑤ Then wavelet coefficients are assigned to the transformation data vectors finally.

DATA MINING Applications:

- ① Spatial mining
- ② Multimedia mining
- ③ Temporal mining
- ④ Web mining
- ⑤ Text mining.

Spatial mining

- for spatial model DM
- uses geographical analysis
- analysts use this analysis for creating B.I
- requires specific techniques and resources to get data into relevant & useful format.
task: search for spatial patterns.

Multimedia mining:

- data is in audio, video, images, graphics
- mining of image, text & audio.
- methods are:
 - (i) similarity search
 - (ii) multidimensional analysis
 - (iii) classification & prediction
 - (iv) association mining

Temporal mining:

- ↳ Time ordered data
 - ↳ eg: Sales, web logs, calls, etc.
- ↳ Reprise repeated measurements.
- ↳ mining methods require modifications to handle temporal relationships
- ↳ Time ordered data contribute to prediction, given the history of events.
- ↳ Time ordered data often link certain events to specific pattern

Web mining:

- ↳ automatic discovery & extract information from web documents.
- ↳ info. can be extracted from server logs & web activity
- ↳ mining: web activity
- ↳ Categories are:
 - ↳ web content mining
 - ↳ web structure mining
 - ↳ web usage mining.

Text mining:

- ↳ retrieve high quality info. from text.
- ↳ available data is in unstructured format.
- ↳ useful when it is necessary to extract information from large datasets.

Text mining:

- ↳ info. retrieval
- ↳ natural language processing

Clustering:

- Unsupervised learning: • considered as self learning process
es. student with all study material but no faculty to guide.
- discovering patterns from data w/o any labels.

Cluster Analysis: form groups with some similarity.

es. grouping of students studying similar subjects.

- Partitioning methods: division of n items into set of k cluster such that sum of squared distance is minimized.

$$E = \sum_{i=1}^k \sum_{p \in C_i} (p - c_i)^2; \quad c_i = \text{centroid of cluster.}$$

- 2 methods: (i) k mean
(ii) k medoids.



Annexure No.:

$(k=2 \text{ i.e. 2 clusters})$

K-mean:

- S-1 Take mean value
- S-2 find nearest no. of mean and put it in cluster.
- S-3 Repeat one and two until we get same mean.

e.g. $k = \{2, 3, 4, 10, 11, 12, 20, 25, 30\}$

$k = 2$

mean₁ = 4

mean₂ = 12

[Random mean no.]

Now take element ① - m₁ & ② - m₂ whichever is closer to m₁ or m₂ put it in.

e.g. $2 - 4 = -2$ $2 - 12 = -10$

-2 is closer than -10 so put it in

(i) $k_1 = \{2, 3, 4\}$; $k_2 = \{10, 11, 12, 20, 25, 30\}$

New m₁ = $(2+3+4)/3 = 3$ New m₂ = $108/6 = 18$

(ii) $k_1 = \{2, 3, 4, 10\}$; $k_2 = \{11, 12, 20, 25, 30\}$

$$m_1 = 4.75$$

or ⑤

$$m_2 = 19.6$$

or ⑥

$$(iii) U_1 = \{2, 3, 4, 10, 11, 12\}$$

$$m_1 = 7$$

$$U_2 = \{20, 25, 30\}$$

$$m_2 = 25$$

$$(iv) U_1 = \{2, 3, 4, 10, 11, 12\}$$

$$m_1 = 7$$

$$U_2 = \{20, 25, 30\}$$

$$m_2 = 25$$

We got same mean both time so we need to
stop

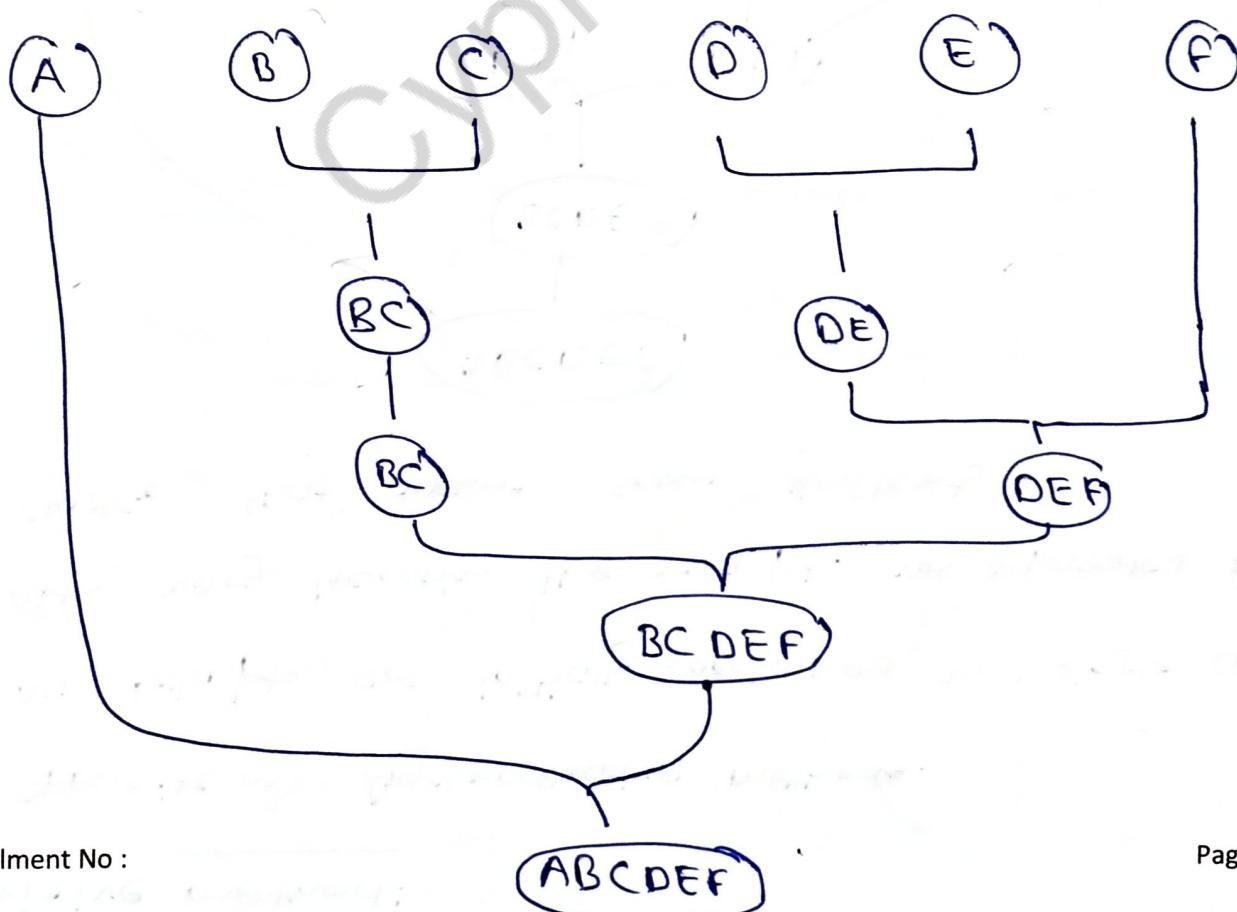
∴ Final answer : $U_1 = \{2, 3, 4, 10, 11, 12\}$

$$U_2 = \{20, 25, 30\}$$

Annexure No.:

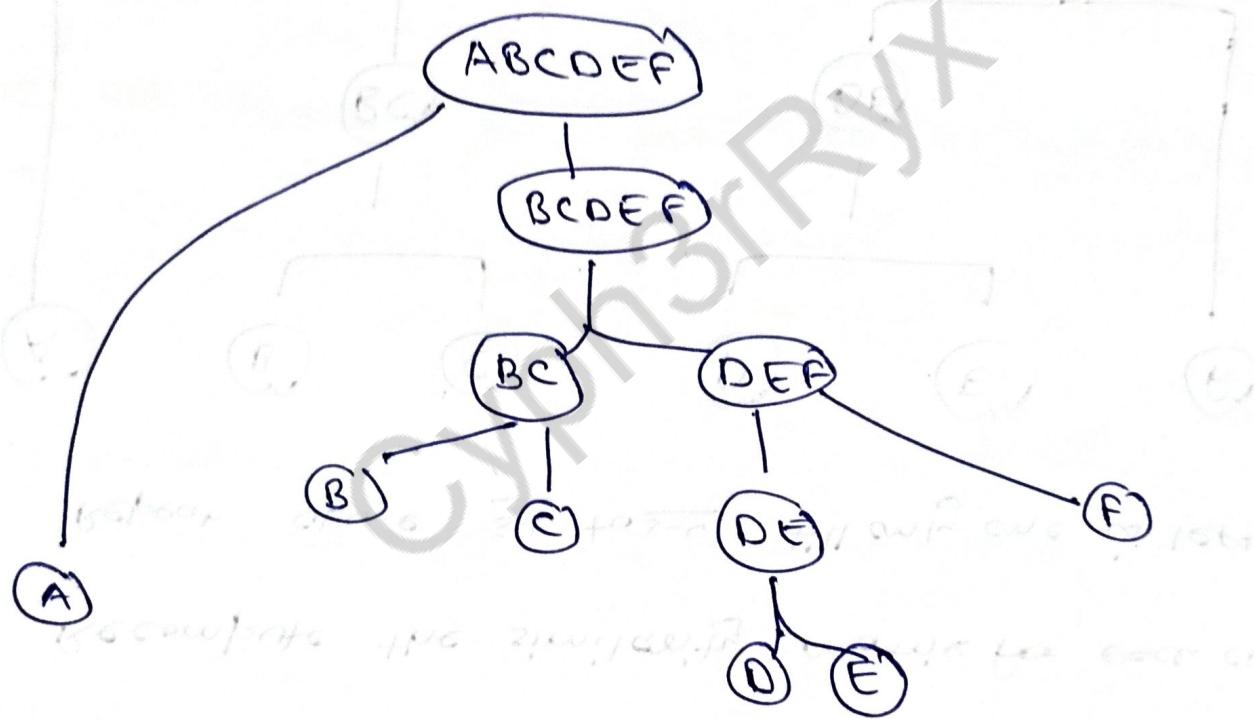
~~Objectives:~~Agglomerative methods:

- S-1 Calculating similarity among clusters.
- S-2 Every datapoint is taken as an individual cluster.
- S-3 Merging clusters with higher proximity among each others.
- S-4 Recompute the similarity matrix for each cluster.
- S-5 Repeat above S-1 to S-4 till only one is left.



Divisive method:

- Opposite of Agglomerative methods.
- All data pts are initially considered as a single cluster.
- After every iteration the data pt. are separated from the cluster that doesn't show similarity



Density Based Scan

(DB scan) :

Spatial
Clustering of Application
with Noise.

- G Data objects are clustered based on density (μ/v)
- C, 2 inputs $\Rightarrow \epsilon$ & min. pts()

(ϵ) -radius of circle formed with data objects as center.

min. pts() is the minimum no. of data points inside the circle.

3 Types of data points:

① Core point: it should satisfy the condn of min. pts()

② Boundary point: neighbour of core

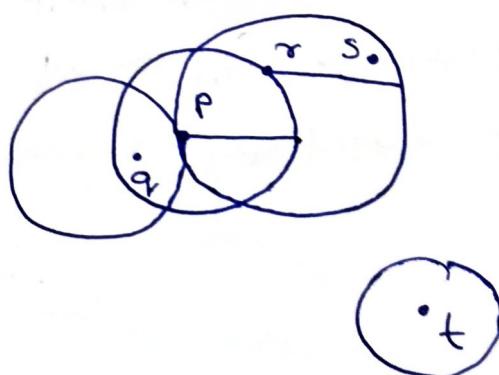
③ Noise point: not core nor boundary

Eg. min. pts = 3

p, r \rightarrow core

q, s \rightarrow boundary

t \rightarrow noise.



Evaluation of Clustering:

① Clustering Tendency:

- non uniformity among data points is vital
- measuring the probability of data points generated by uniform data distribution.

Hypothesis (H)

Null hypothesis

: Non uniform data

Alternate hypothesis

: Random data

$H > 0.5 \rightarrow$ data contains cluster

∴ rejects null hypothesis.

H closer to 0 → no clustering tendency.

② Number of Cluster:

Depends on

Distribution shape
scale in data set
clustering resolution

2 approach

Domain knowledge

initial knowledge
on forming no. of
clusters.

Data driven approach

Empirical

elbow

③ Quality of Clustering:

Characteristic of Cluster :

intra		inter.
btw 2		btw 2
data pts.		data cluster.

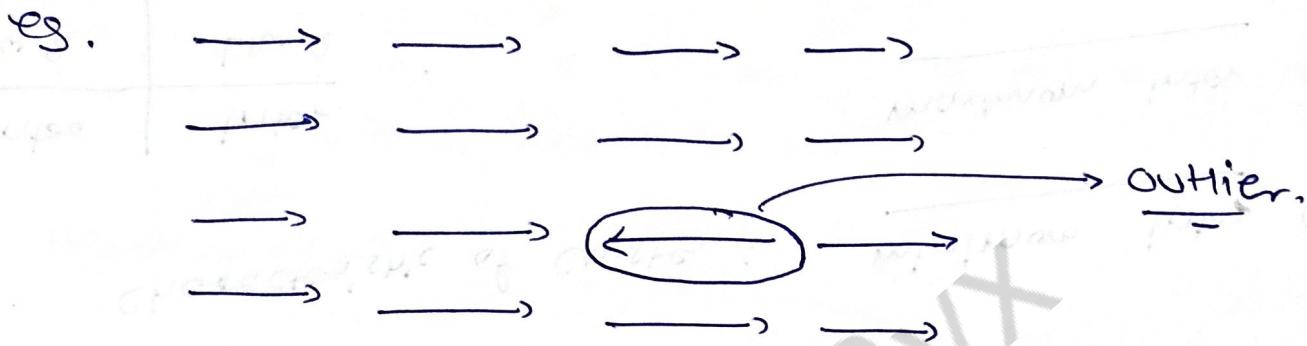
minimum intra cluster
distance
maximum inter
cluster distance

Extrinsic measures : True labels
required

Intrinsic measures : not required

Outlier Analysis:

Outlier is the data object which don't obey the general behavior.



Outlier detection: Process of identifying the outliers and then removing it.

2 Approaches

Statistical

↳ Probability based

- Parametric

- Non Parametric

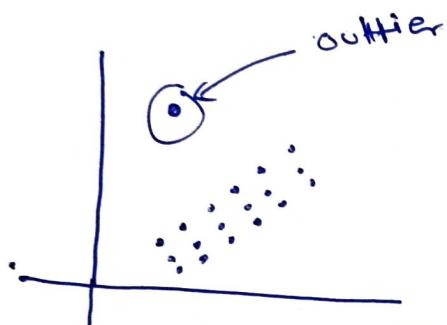
Proximity

- ↳ Location based
- Density based
 - Distance based
 - Grid based
 - Deviation based

Types of outliers:

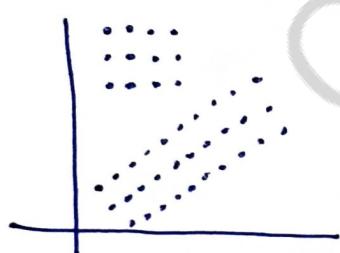
① Global / Point outliers:

when a single data object deviates from rest of the data set



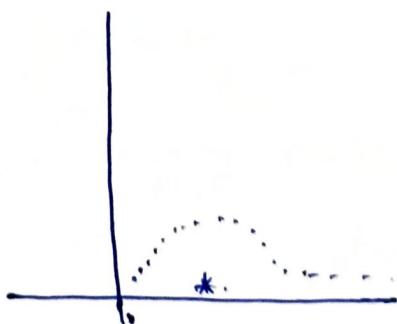
② Collective outlier:

when a group of data points deviates from the rest of the data set.



③ Conditional outlier:

deviation due to specific conditions



Interestingness of pattern:

from millions of DM patterns which of them are really interesting?

3 Questions.:

- (1) what makes the pattern interesting?
- (2) Can DM system generates all of the interesting patterns?
- (3) Can DM system generates only interesting patterns

Ans. ① A pattern is interesting if it is

- easily understood by humans
- valid on new / test data
- potentially useful.

Ans. ②

Refers to completeness of DM system

As it is not possible to generate all interesting patterns

Ans. ③

Refers to optimization of DM system.

Difficult to generate only interesting patterns.

Classification of DM system:

(1) Based on mined DB:

- (i) Relational
- (ii) Transactional
- (iii) Object relational
- (iv) Data warehouse

(2) Based on knowledge mined:

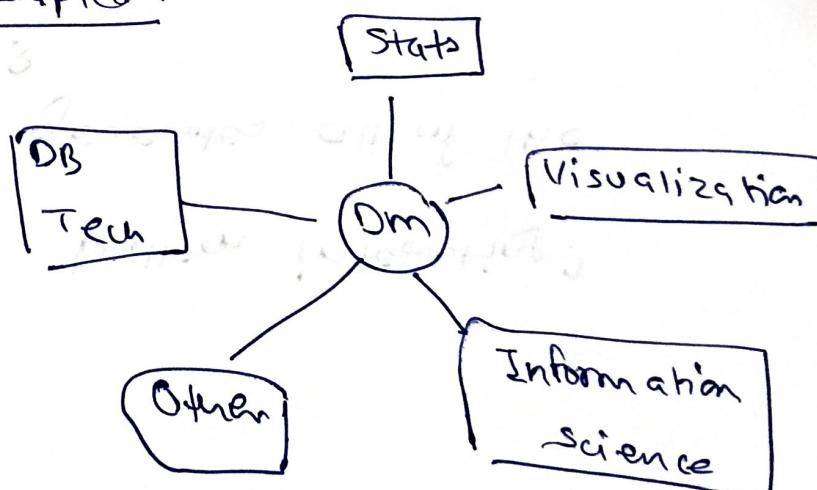
- (i) Characterisation
- (ii) Discrimination
- (iii) Association & correlation analysis
- (iv) Classification
- (v) Prediction
- (vi) Outlier Analysis
- (vii) Evolution Analysis

(3) Based on techniques used:

ML, Statistics, Neural Networks, etc.

(4) Based on application adapted:

- (i) finance
- (ii) DNA
- (iii) Email
- (iv) Stock market
- (v) Telecommunication



Frequent Pattern:

The pattern that appears frequently in dataset
eg. Milk & Bread

Market Basket Analysis:

Process of analysing customer buying habits, by finding the associations between the different items that a customer will place in their basket.

Strategies: ① Placing them together
 ② Placing them at 2 different ends

- This analysis will help sellers to plan their shelf space for increased sales.

C, frequent patterns are represented by association rules.
ex. Computer & Antivirus.

Computer \Rightarrow Antivirus software

Support \Rightarrow 2%

Confidence \Rightarrow 60%

Support: identifies how frequently a rule is applied
to given dataset.

$$S(P \rightarrow Q) = \frac{f(P \cup Q)}{N}$$

*N = Total
Transactions.*

Confidence: defines frequent occurrence of items of Q in transactions.

$$C(P \rightarrow Q) = P(B|A) = \frac{\text{Support}(A \cup B)}{\text{Support}(A)}$$

APRIORI Algorithm,

By R. Agrawal &

(1997)

R-soileant

Shows how objects are associated with each other.

Objective : To generate an association.

Eg.

Min. support = 50%.

Threshold confidence = 70%.

ID	Items
100	1 3 4
200	2 3 5
300	1 2 3 5
400	2 5

Now check how much time each item is occurring

i.e. 1 occurs $\frac{2}{4}$ time
 2 occurs. $\frac{3}{4}$ time
 3 $\frac{3}{4}$ " "
 4 $\frac{1}{4}$ "
 5 $\frac{1}{4}$ "

$\Rightarrow \text{Support} / 20 (\text{total})$

Item	Support	min. Support
1	2	$2/4 = 50\%$
2	3	$3/4 = 75\%$
3	3	$3/4 = 75\%$
4	1	$1/4 = 25\% (X)$
5	3	$3/4 = 75\%$

Enrollment No :

Cancel out the item which has support less than suggested
 i.e. item 4 $\Rightarrow 25\%$.

Page No:

Remaining itemset = $(1, 2, 3, 5)$

from pairs. $\Rightarrow (1, 2), (1, 3), (1, 5), (2, 3), (2, 5), (3, 5)$

Items	Support	Min. Support
$(1, 2)$	1	$1/4 = 25\% \text{ (X)}$
$(1, 3)$	2	$2/4 = 50\%$
$(1, 5)$	1	$1/4 = 25\% \text{ (X)}$
$(2, 3)$	2	$2/4 = 50\%$
$(2, 5)$	3	$3/4 = 75\%$
$(3, 5)$	2	$2/4 = 50\%$

Check how many times $(1, 2)$ occurs together and same goes for all items.

Remaining itemset = $(1, 3), (2, 5), (2, 3), (3, 5)$

Now from triplets

i.e $(1, 3, 5)$ $(1, 2, 5)$ $(1, 2, 3)$, $(2, 3, 5)$

Check the similar items betn dataset and only form triplet with those.

item set	support	min. support
(1, 2, 3)	1	$\frac{1}{6} = 16.6\% \text{ (X)}$
(1, 2, 5)	1	$\frac{1}{6} = 16.6\% \text{ (X)}$
(1, 3, 5)	1	$\frac{1}{6} = 16.6\% \text{ (X)}$
(2, 3, 5)	2	$\frac{2}{6} = 33.3\%$

Only one triplet remains (2, 3, 5)

Now let calculate support & confidence.

$$\text{Confidence} = \frac{\text{support}(A \cup B)}{\text{support}(A)}$$

Using (2, 3, 5) we need to find/generate association rules.

Rules	Support	Confidence
$(2 \wedge 3) \rightarrow 5$	2	$\frac{2}{2} \Rightarrow 100\%$
$(3 \wedge 5) \rightarrow 2$	2	$\frac{2}{2} \Rightarrow 100\%$
$(2 \wedge 5) \rightarrow 3$	2	$\frac{2}{2} \Rightarrow 100\% \text{ (X)}$
$2 \rightarrow (3 \wedge 5)$	2	$\frac{2}{2} \Rightarrow 100\% \text{ (X)}$
$5 \rightarrow (2 \wedge 3)$	2	$\frac{2}{2} \Rightarrow 100\% \text{ (X)}$
$3 \rightarrow (2 \wedge 5)$	2	$\frac{2}{2} \Rightarrow 100\% \text{ (X)}$
Enrollment No:		Page No:

$$(2 \wedge 3) \rightarrow S \Rightarrow \text{confidence} = \frac{S((2 \wedge 3) \cup S)}{S(2 \wedge 3)} = \frac{2}{2} = 100\%$$

(by previous
table)

$$S(A \cup B)$$

$$(2 \vee 3) \rightarrow S$$

Do it for all values.

$$(S \vee S) \rightarrow S$$

Remove all the rows which has confidence less than threshold confidence.

i.e 70%

confidence = support(SAB)

now after applying support of confidence

and we get support(SAB) = 60%

$$(S \wedge S) \rightarrow S$$

support = 60%

$$(A \wedge S)$$

support = 60%

$$(A \wedge S)$$

support = 60%

FP Growth Algorithm:

Frequency Pattern Growth is an efficient & scalable method for mining the complete set of FP using a tree structure for storing information about FP called FP tree.

Eg.

Min. support = 30 %,

Item sets are given as below.

ID	Items.
1	E, A, D, B
2	D, A, E, C, B
3	C, A, B, E
4	B, A, D
5	D
6	D, B
7	A, D, E
8	B, C

ID	Frequency	Priority
A	5	3
B	6	1
C	3	8
D	6	2
E	4	4

How to write priority?

- more frequency = more priority
- same frequency = FCF

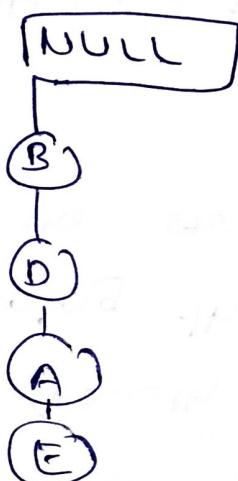
Now rewrite ID based on priority

(B, D, A, E, C)

order items according to their priority (B D A E C)

ID	Item	ordered items.
1	E A D B	B D A E
2	D A E C B	B D A E C
3	C A B E	B A E C
4	B A D	B D A
5	D	D
6	D B	B D
7	A D E	D A E
8	B C	B C

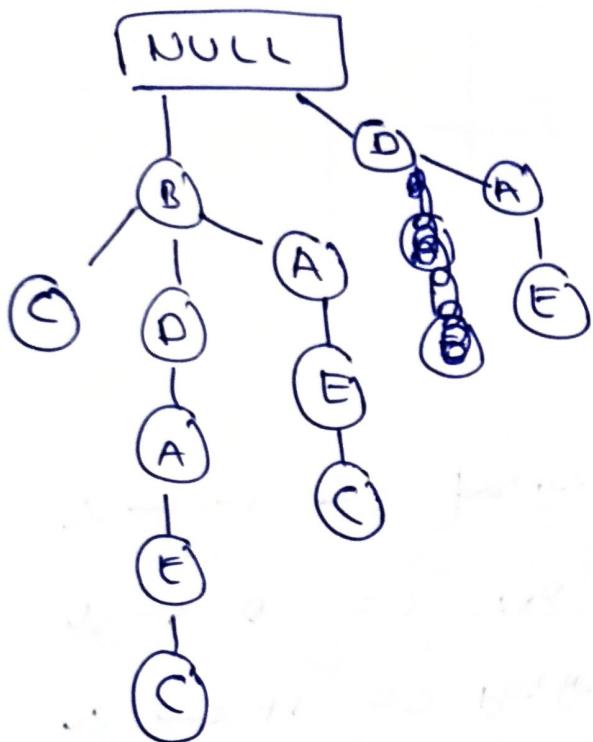
Now let's make a tree with root node NULL



for ordered item - (1)

B D A E

B - 1
D - 1
A - 1
E - 1



$B \rightarrow 1, 2, 3, 4, 5$
 $D \rightarrow 1, 2, 3, 4$
 $A \rightarrow 1, 2, 3$
 $E \rightarrow 1, 2$
 $C \rightarrow 1$
 $A \rightarrow 1$
 $E \rightarrow 1$
 $C \rightarrow 1$
 $D \rightarrow 1, 2$
 $A \rightarrow 1$
 $E \rightarrow 1$
 $C \rightarrow 1$

Final Answer with FP Tree.

- ↳ only 2 passes over dataset
- ↳ compresses dataset
- ↳ no candidate generation
- ↳ FP tree may not fit in memory
- ↳ FP tree is expensive

Characteristics
of FP growth.

* Correlation Analysis:

Used to measure the relationship b/w 2 variables

$$r_{A,B} = \frac{\sum (A - \bar{A})(B - \bar{B})}{(n-1) \sigma_A \sigma_B}$$

$r_{A,B}$ = Karl Pearson Correlation coefficient

\bar{A}, \bar{B} = mean of A & B

σ_A, σ_B = standard deviation of A & B

n = no. of tuple in DB.

$r \rightarrow 3$ values $(0, -1, +1)$

$r \rightarrow +1 \Rightarrow$ Perfect positive correlation

$r \rightarrow 0 \Rightarrow$ No correlation

$r \rightarrow -1 \Rightarrow$ perfect negative correlation

Eg.

A	B
20	8
12	34
9	4

$$\bar{A} = \frac{20 + 12 + 9}{3}$$

$$= 13.66$$

$$\bar{B} = \frac{8 + 34 + 4}{3}$$

$$= 15.33$$

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

$$\therefore \sigma A = \sqrt{\frac{(20 - 13.66)^2 + (12 - 13.66)^2 + (9 - 13.66)^2}{3-1}}$$

$$\therefore \boxed{\sigma A = 5.68}$$

$$\sigma B = \sqrt{\frac{\sum (B - \bar{B})^2}{n-1}} = \sqrt{\frac{(8 - 15.33)^2 + (34 - 15.33)^2 + (4 - 15.33)^2}{3-1}}$$

$$= \boxed{16.28}$$

Now in main
formula

$$r_{AB} = \frac{\sum (A - \bar{A})(B - \bar{B})}{(n-1) \sigma A \sigma B}$$

$$= \frac{[(20 - 13.66)(8 - 15.33) + (12 - 13.66)(34 - 15.33) + (9 - 13.66)(4 - 15.33)]}{(3-1)(5.68)(16.28)}$$



negative
correlation

Classification & Prediction:



finding a good model to predict the class of object whose class labels are unknown

e.g.: grouping of patients based on the medical records.

Predicting a missing/unknown value based on past/corrent data

e.g., predicting the correct treatment for a person based on their medical condition.

categorization of new data with the help of past data.

Classification has 2 stages:

① model construction

② model usage,

Training data	Mark	Result
	4	Pass
	3	fail
	2	fail
	6	Pass
	7	Pass
	8	Pass

Classification Algorithm

if mark $\leq 3 \Rightarrow$ result = fail

if mark = 7; result = ?

Learn, Analyze, Generate classification

* Decision Tree : (ID3 Algorithm)

- (1) In the given dataset choose a target attribute.
- (2) Calculate I.G (info. gain) of target attribute

$$I.G = \frac{-P}{P+N} \log_2 \left(\frac{P}{P+N} \right) - \frac{N}{P+N} \log_2 \frac{N}{P+N}$$

- (3) for remaining attributes, find entropy

$$\text{Entropy} = I.G \times \text{Probability.}$$

$$E(A) = \sum \frac{P_i + N_i}{P+N} I(P_i, N_i)$$

- (4) Calculate Gain = $I.G - E(A)$

e.g.

Age	Completion	Type	Profit
old	Yes	s/w	Down
old	No	s/w	Down
old	No	h/w	Down
mid	Yes	s/w	Down
mid	Yes	h/w	Down
mid	No	h/w	up
mid	No	s/w	up
new	Yes	s/w	up
new	No	h/w	up
new	No	s/w	up

S1

Target Attribute (\Rightarrow Profit)

$z = -5 + 10x_1 + 6x_2$

$\text{Profit} = \frac{1}{2} z^2 - 5z - 10x_1^2 - 6x_2^2$

$\text{Profit} = \frac{1}{2} (-5 + 10x_1 + 6x_2)^2 - 5(-5 + 10x_1 + 6x_2) - 10x_1^2 - 6x_2^2$

$\text{Profit} = 25 - 50x_1 - 60x_2 + 100x_1^2 + 120x_1x_2 + 36x_2^2$

S-2

=

$$I.G = \frac{-P}{P+N} \log_2 \left(\frac{P}{P+N} \right) - \frac{N}{P+N} \log_2 \left(\frac{N}{P+N} \right)$$

$$P \Rightarrow \text{Count (down)} \Rightarrow 5$$

$$N \Rightarrow \text{Count (up)} \Rightarrow 5$$

$$= \frac{-5}{10} \log_2 \left(\frac{5}{10} \right) - \frac{5}{5+5} \log_2 \left(\frac{5}{5+5} \right)$$

$$= \left(\frac{1}{2} \log_2 (2^{-1}) + \frac{1}{2} \log_2 (2^1) \right)$$

$$= \left[\frac{1}{2} \times (-1) \log_2 (2) + \frac{1}{2} \times (-1) \log_2 (2) \right]$$

$$= \left(\frac{1}{2} \times -1 + \frac{1}{2} \times -1 \right) = -\frac{1}{2} + (-\frac{1}{2}) = -1$$

$\boxed{-1}$

$$\boxed{I.G = -1}$$

S-3 Entropy Calculation for remaining attributes:

$$E(A) = \frac{\sum P_i N_i}{P+N} I(P_i N_i) \quad [I_G \times \text{Probability}]$$

(Age) : Prepare table for each attribute

Rows: values of attribute

Columns: values of target attribute.

	down	up
old	3	0
mid	2	2
new	0	3

$$\text{Entropy} = I_G \times \text{Probability}$$

$$I_G(\text{old}) = - \left[\frac{3}{6} \log \left(\frac{3}{6} \right) + \frac{3}{6} \log \left(\frac{1}{3} \right) \right] = 0$$

$$\text{Prob.} = 3/6$$

$$\sum (\text{old}) = 0 \times 3/6 = 0$$

$$I.C_E(\text{mid}) = - \left[2/_{10} \log(2/_{10}) + 3/_{10} \log(3/_{10}) \right] = 1$$

$$E(\text{mid}) = I \times 4/_{10} = 0.4$$

$$I.C_E(\text{new}) = - \left[0/_{3} \log(0/_{3}) + 3/_{3} \log(3/_{3}) \right] = 0$$

$$E(\text{new}) = 0 \times 3/_{10} = 0$$

$$\text{Entropy (Age)} = 0 + 0.4 + 0 = 0.4$$

Gain = $I.C_E - E(A) = 1 - 0.4 = 0.6$

Same goes for Completion and Type

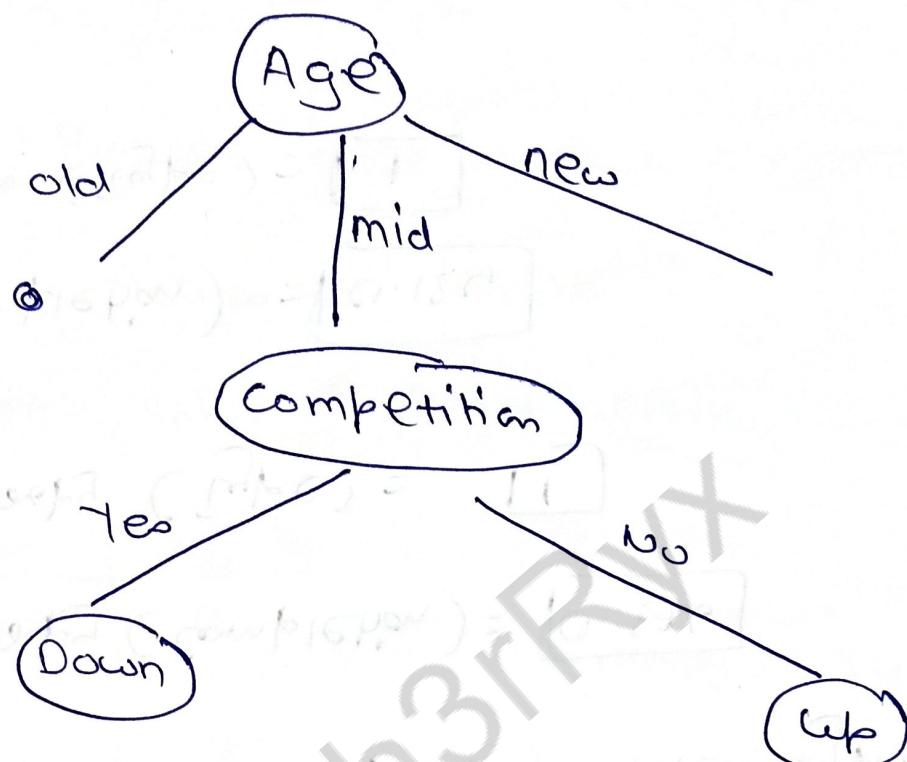
$$\therefore \text{Entropy (Completion)} = 0.876$$

$$\therefore \text{Entropy (Type)} = 1$$

$$\text{Gain(Completion)} = 0.124$$

$$\text{Gain } \text{Type} = 1$$

Highest gain \rightarrow root node



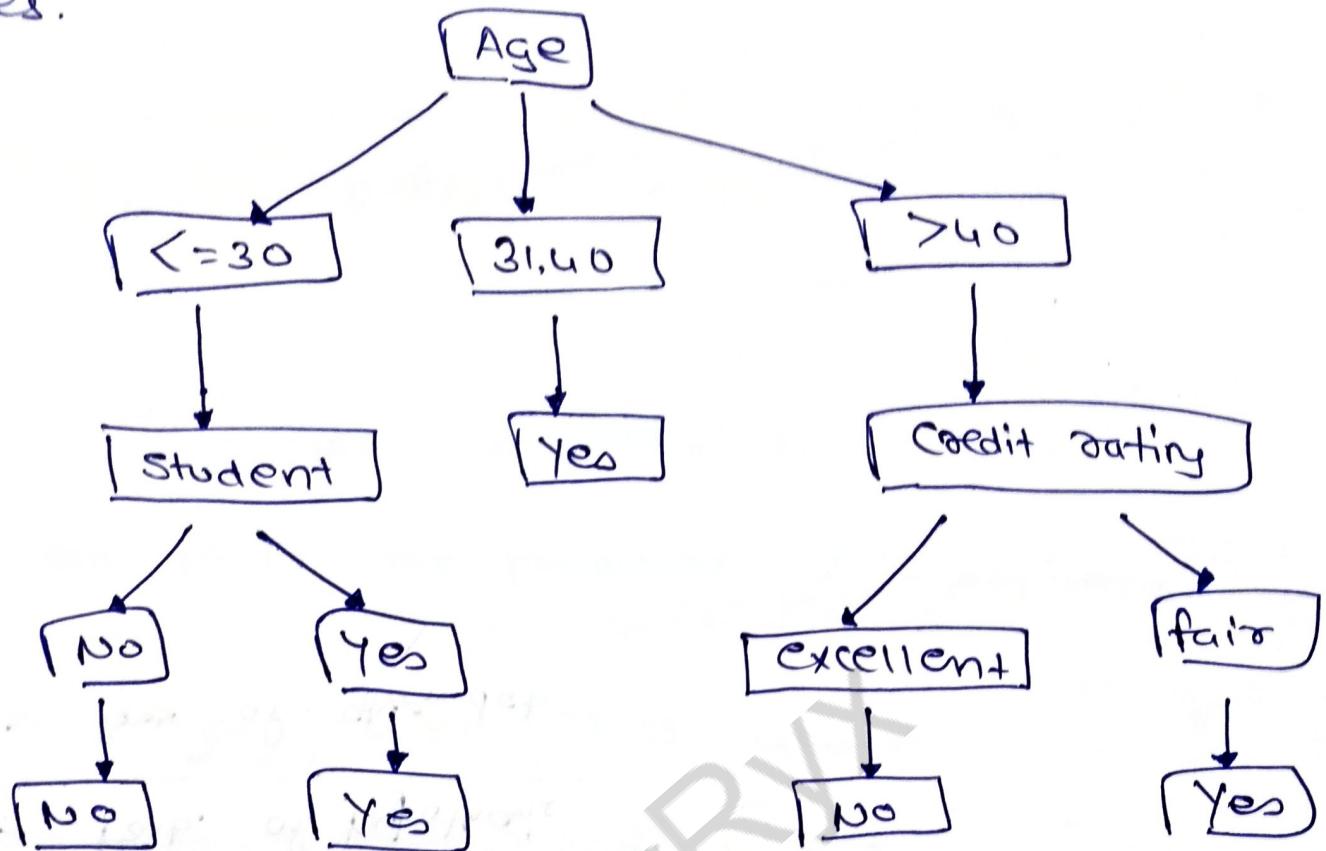
$$E(\text{Gain}) = \alpha \times 3^0 = 10$$

$$\text{Gain}(\text{new}) = -[\beta^3 \text{Gain}(3^3) + 3\beta^2 \text{Gain}(3^2)] = 0$$

$$E(\text{Gain}) = 2 \times 3^0 = 10$$

$$\text{Gain}(\text{mid}) = -[\delta^3 \text{Gain}(3^3) + 8 \text{Gain}(3^2)] = 10$$

es.



- Rules:
- IF age = " ≤ 30 " AND student = "no" THEN buys computer = "no"
 - If age = " ≤ 30 " AND student = "yes" THEN buys computer = "yes"
 - If age = "31,40" THEN buys computer = "yes"
 - If age = " > 40 " AND credit rating = "excellent" THEN buys computer = "no"
 - If age = " > 40 " AND credit rating = "fair" THEN buys computer = "yes"

Bayes' Theorem:

$$P(A/B) = \frac{P(B/A) \cdot P(A)}{P(B)}$$

A = hypothesis

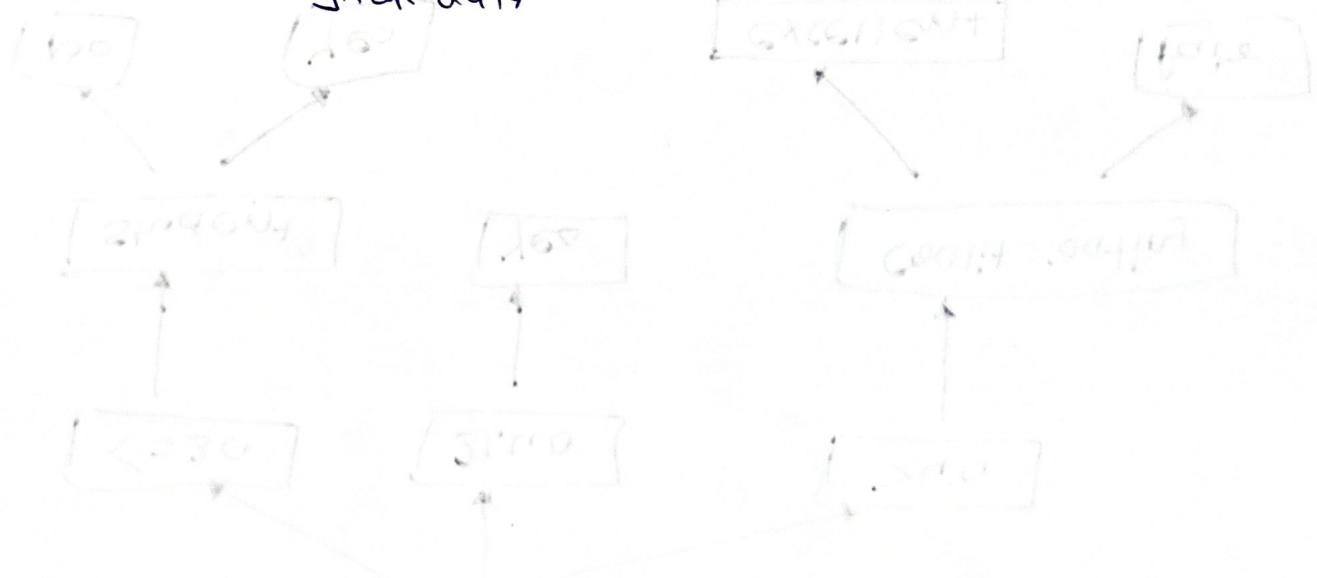
B = given data

$P(A/B)$ = finding prob. of hypothesis when prob. of training examples are given.

$P(B/A)$ = finding prob. of given data with prob. of hypothesis that is true

$P(A)$ = Prob. of hypothesis

$P(B)$ = Prob. of given data



Bayesian Classification:

$$P(A/B) = \frac{P(B/A) \cdot P(A)}{P(B)}$$

- Bayesian classifiers are statistical classifiers
- They can predict the probabilities of class items

Naive Bayes Theorem / Classification:

same formula $\Rightarrow P(A/B) = \frac{P(B/A) \cdot P(A)}{P(B)}$

e.g

fruit = {Yellow, Sweet, Long}

Fruit	Yellow	Sweet	Long	Total
Orange	350	400	0	650
Banana	400	300	350	1050
Others	50	600	50	180
Total	800	850	400	2050

We need to select a fruit that is yellow, sweet & long from given data.

$$\therefore P(\text{Yellow/orange}) = \frac{P(\text{orange/yellow}) \cdot P(\text{yellow})}{P(\text{orange})}$$

$$= \frac{\left(\frac{350}{800}\right) \cdot \left(\frac{800}{1200}\right)}{\left(\frac{650}{1200}\right)} = 0.53$$

$$\therefore P(\text{Sweet/orange}) = \frac{P(U/S) \cdot P(S)}{P(O)} = 0.69$$

$$\therefore P(\text{Long/orange}) = 0$$

Now $P(\text{Orange/fruit}) = (0.53) \times (0.69) \times (0) = 0$

$$P(\text{banana/fruit}) = (1) \times (0.75) \times (0.85) = 0.65$$

$$P(\text{others/fruit}) = (0.33) \times (0.66) \times (0.33) = 0.072$$

We will choose the fruit which has highest probability

i.e Banana

KNN Algorithm:

lazy learning

k - nearest neighbour algorithm

e.g. Given data query $\Rightarrow \mathbf{x} \Rightarrow (\text{maths} = 6; \text{CS} = 8)$
& $k = 3$

Classification = Pass / fail

Maths	CS	Result
4	3	F
6	7	P
7	8	P
5	5	F
8	8	P

We need to find
that, given query
 \mathbf{x} has 2 condn
for maths & CS
i.e G&F resp. So
if the student is
Pass or fail based
on that so let see.

Euclidean Distance $(d) = \sqrt{(x_{01} - x_{A1})^2 + (x_{02} - x_{A2})^2}$

x_0 = observed value

x_A = actual value.

$$d = \sqrt{(6-4)^2 + (8-8)^2} = \sqrt{4} = 2$$

$$d = \sqrt{(6-6)^2 + (8-8)^2} = \sqrt{0} = 0$$

$$d_1 = \sqrt{(6-4)^2 + (8-3)^2} = \sqrt{29} = 5.34$$

$$d_2 = \sqrt{(6-6)^2 + (8-7)^2} = \sqrt{1} = 1$$

$$d_3 = \sqrt{(6-7)^2 + (8-8)^2} = \sqrt{10} = 1$$

$$d_4 = \sqrt{(6-5)^2 + (8-5)^2} = \sqrt{10} = 3.16$$

$$d_5 = \sqrt{(6-8)^2 + (8-8)^2} = \sqrt{4} = 2$$

Now as $k=3$; we need to pick 3 neighbours.

Choose 3 values and all of them should be as near as they can and as minimum as they can.

i.e. ~~maximum~~ (1), (1), (2)

d_2	d_3	d_5
↓	↓	↓
Pass	Pass	Pass

3P / OF

↓

Given query is

Pass

Q: Given question $\Rightarrow x \rightarrow$ (max of 3+3)

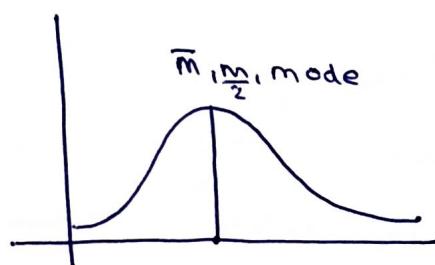
so + largest + minimum + 3+3

Data Skewness:

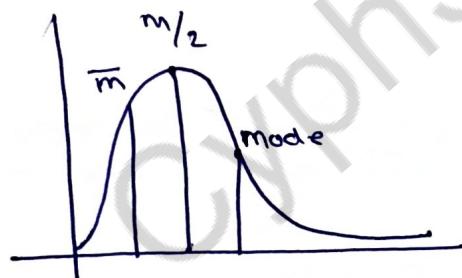
when most of the values are skewed to the left or right side from the median then it is called skewed

3 types:

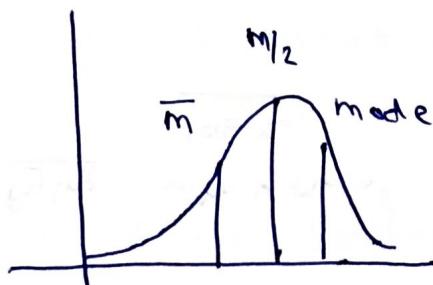
- (1) Symmetric : mean , median & mode at same pt.



- (2) Positively skewed : values are to the left from median.



- (3) Negatively skewed : values are to the right from median



Classification

- (1) **Supervised**
- (2) Classifying with help of labels.
- (3) more complex.
- (4) Naive Bayes, etc.

Clustering

- (1) **Unsupervised**
- (2) ~~grouping~~ with help of similarity.
- (3) less complex.
- (4) k means, etc.

Credits:

Cyph3rRyx

For More Handwritten Study Notes, Click [Here!](#)