

5/8/24

UNIT - I

Date: / / /

Objective : To analyse the student's Performance

Target audience : (age) < 5, 11 >

Step-1 : Data Collection (Raw fact)

Surveys, Questionnaires / forms → Primary data
 Personal & social info → personal

GPA, Attendance, sports, Course interest - Secondary data

Wearable devices - (Smart devices / watches)

1) Health - Stepcounting, Latent needs / patient data

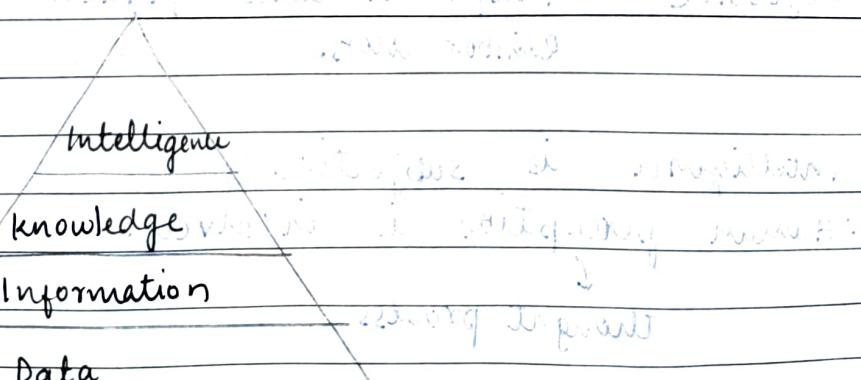
Physical / Mental - sleeping patterns

Step-2 : Information - Processed data.

(statistics)

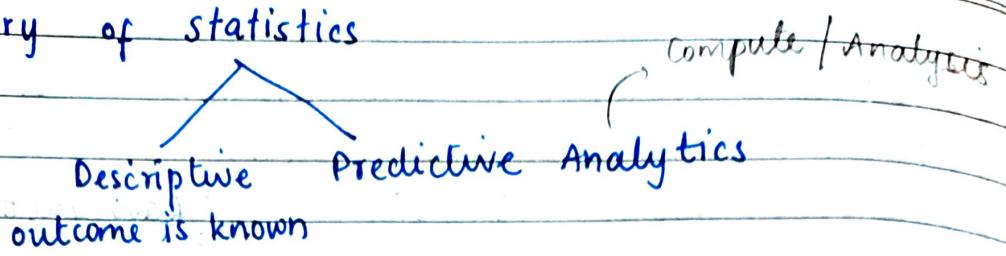
Graphs, Tables, charts

6/8/24 Step-3 : Knowledge



pyramid of Data

Summary of statistics



Correlation analysis - study hours & marks

among two/more entities with measuring parameters

Step - 4 : Intelligence - decision making
Agents / Assistant

Chatbots , QPAs , smart devices , robots .

Maths - Personalised / tutoring sessions .

Continuous evaluation and monitoring .

Main drawback - training AI . (issues of Subjectivity)
Explainability based on context .

subjective - Cannot be defined in a single way
on a scale . [nCr permutations .]

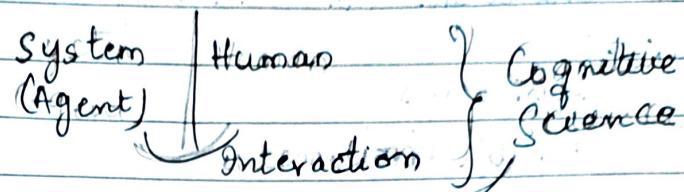
objective - True / False stmts , Fixed , Binary format
Rainbow colors , Rule based approach .

Intelligence is subjective if
Human perception is involved .

thought process .

IQ is objective measure .

biasd problem



- * Challenges:
 1. Data privacy & Ethical concerns are challenges while collecting data
 2. Achieving high Response rate is another challenge (Biasness)

3. Issues w information

- Null values / missing data
- duplicates (names without roll nos) /

data inconsistency

S.No	Name	%	Company
1	K. Vinay	?	Google
2	K. Vinay	80	Google

if no S.No
Name becomes duplicate

missing data

- data Quality / Unwanted data.

Address in Report card

4. Challenges of knowledge:

- Summarization (content usage)

- Outdated

Trend - seasonal changes - ~~pre~~, during & post

Pattern - Buying eraser after buying pencil

5. Challenges in Intelligence

- Biasness
- Subjective.

Classification :

1. Binary classification 2. Multi classification.

Review of a movie
Feedbacks.

o/p: High / Moderate / Low
Ratings

Prediction: Based on classified o/p we predict.

Ex: Movie duration, rating, cast, etc.

Predicting Placement score from CGPA

(mining)

7/8/24 Data mining : Extracting knowledge from large amount of data

approach to solve a problem

Synonyms : Data mining, knowledge mining, Data archaeology, knowledge representation, Data pattern analysis, Data dredging, knowledge extraction.
KDD - knowledge discovery in databases

Todo why is it called Data Mining? & not knowledge mining

Machine learning, data mining, pattern recognition

1) Data mining is a process of extracting

→ Movie Recommendation System

data → Genres, Watchtime, Language, Period

Actors preferences, Age group, favourite movie

Education / Agriculture

* Evolution of Database Technology.

2. Objective: Evolution of Health monitoring system.

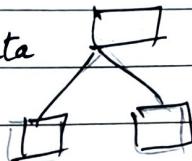
small clinics

1) Patient attributes (pid, pname)

Issue with Storing large volumes of data.

2) Tree format

Storage issue if data is too vast.



① csv files / Excel
flat files.

② hierarchical / Network database

3) Individual tables for [patients] [doctors] [appointments]

- RDBMS .

2. Future (4)

4) Online appointments

{ Web database
Webmining.

5) Images (unstructured)

⑤ NoSQL / MongoDB.
→ MySQL
(after structuring)

6) Multispeciality hospitals

(large volumes of data)

Visualization - reports / profile

⑥ Amazon, Google
for processing big data

7) Dashboards (Tableau) - Assistance

⑦ NLP

8) Advanced

- Data mining, ML, DL,

AI (based on datatype)

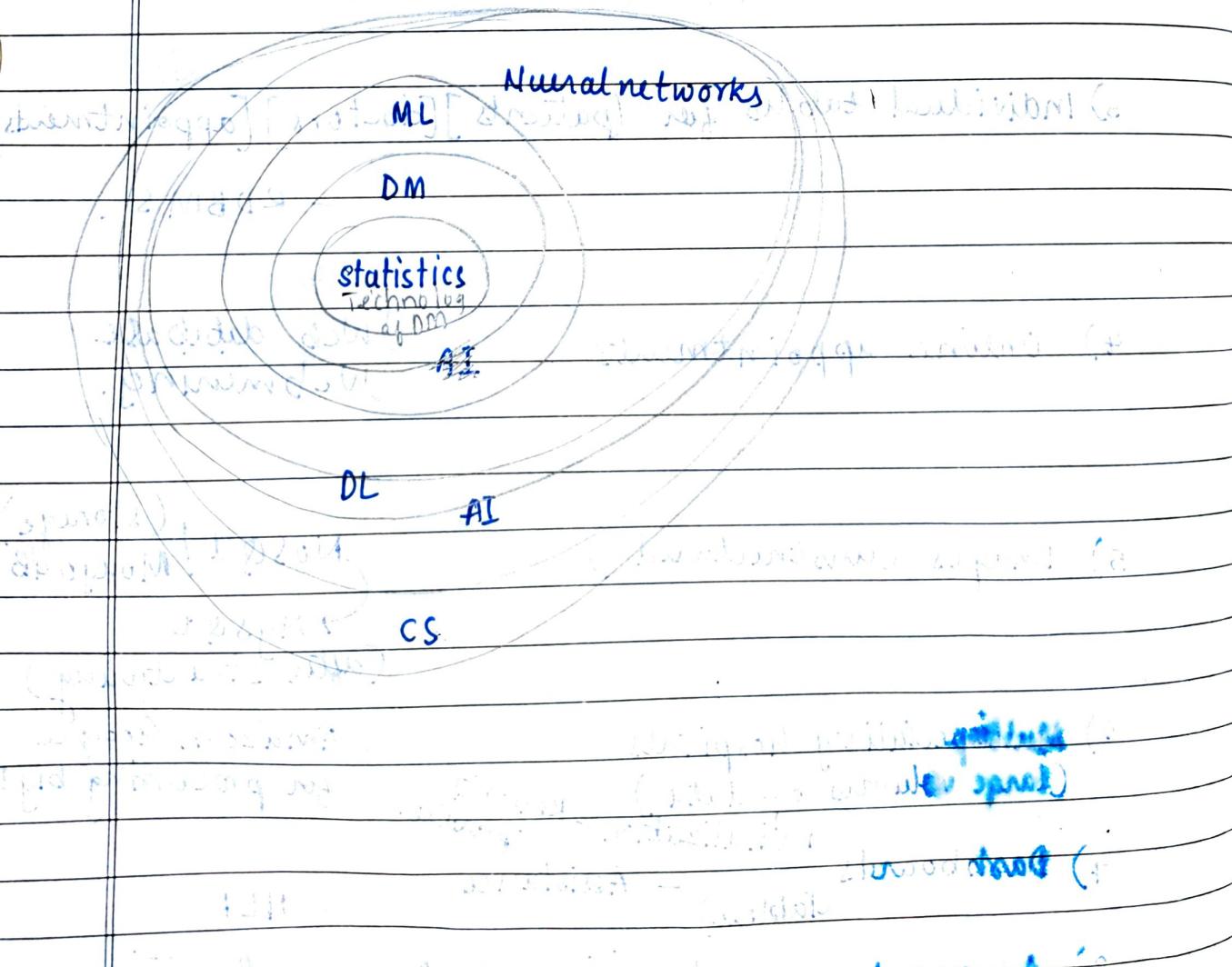
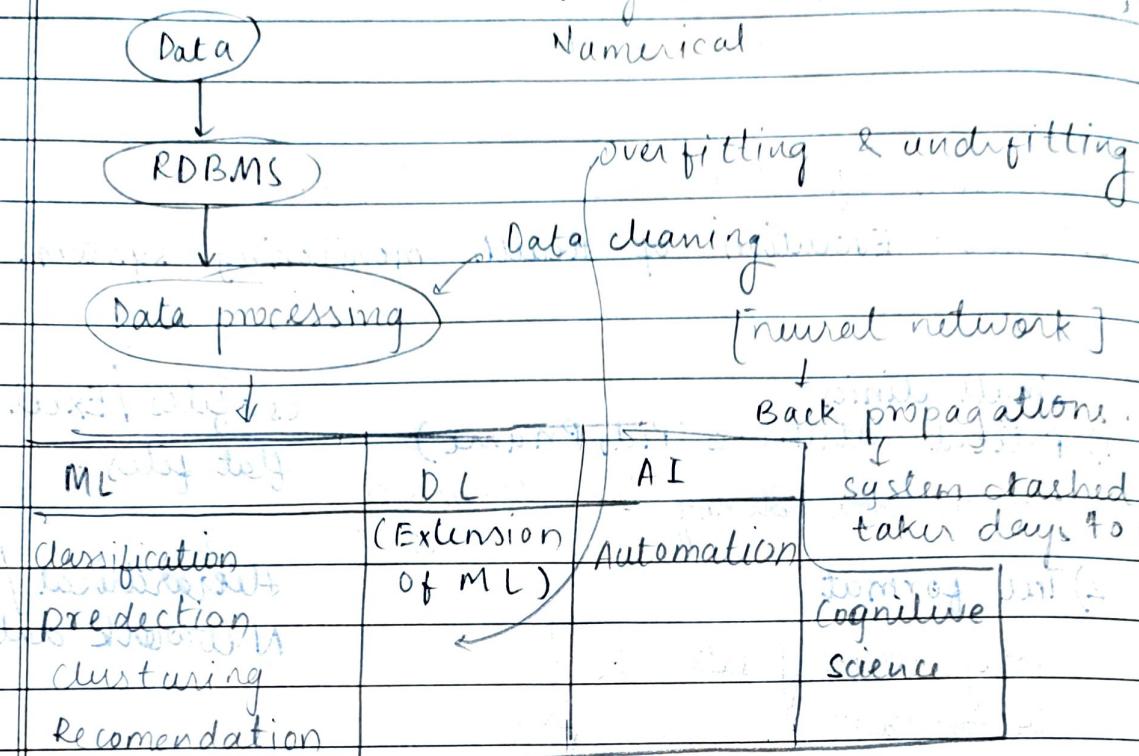
Robotic

(advanced).

types of Data

Image, text, audio, video, sign

Numerical



12/8/24

SM

* Knowledge Discovery Data processes.

1. Data cleaning : Its a procedure to remove the outliers replacing the null values , handling the issues of (to remove outliers , null values , missing data) noisy data) inconsistency , identifying the outlier basing on context of data

2. Data integration : In this step various formats & sources clinical data - investigations are done of data (try to perform action) are combined/integrated behavioural data -

Data preprocessing

These both are merged together.
[text + numeric , image + text , audio + signals , video + signals , . . .]

3. Data selection :

collecting data that is required / relevant to analysis is retrieved from database.

conversion

4. Data transformation : [image → numeric]

to convert one form of data into another.

→ [summary of statistical data is generated & we convert different forms of data to similar datatype.]

5. Data mining : Extracting meaningful insights from large

6. Pattern evaluation : Interesting measures . support & confidence

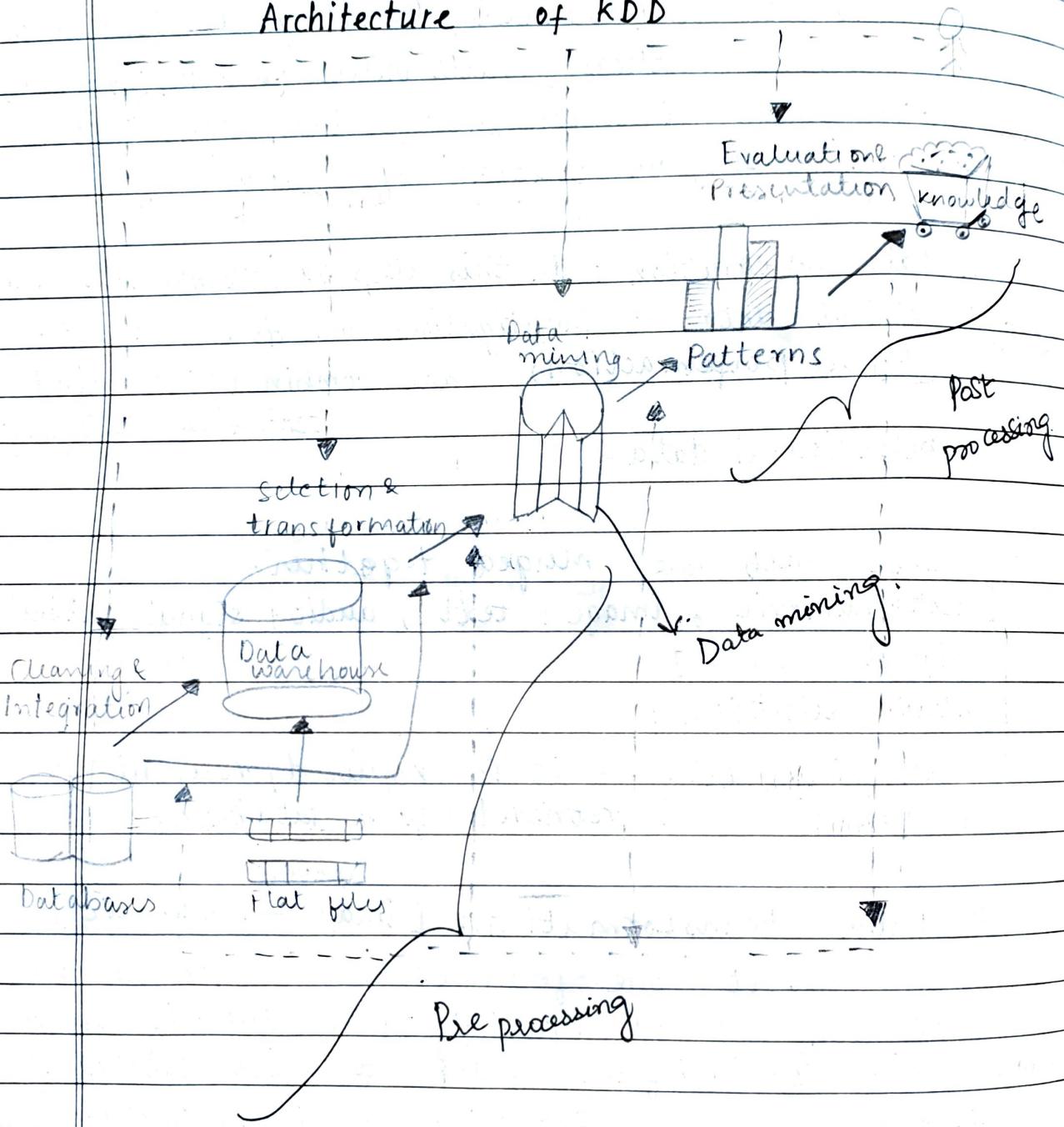
post processing

Post processing

7. knowledge presentation : visualization.

→ tableau

Architecture of KDD



(not always structured), Query processing
 Database always relies on transactional data
 (Present data)

Datawarehouse (Historic data) } Both storage,
 retrieval / mgmt

knowledge base always in a structured formats.

Applied by Intelligent agents [Robots]

Used for decision making & Business Intelligence

Diff types of databases can be

- 1) Numerical data, categorical, text, image, audio data, time series
- 2) RDBMS, NoSQL, Data warehouses, Distributed DB.

what kinds of data can be trained?

Illustrate various databases wrt datamining.

13-8-24

Classification of Attributes

Attribute / Entity / Object

Emp ID	Encl	Ename	...

1. Categorical data : [Qualitative data]

- a) Nominal : the values of the nominal attribute are different names i.e., they provide the enough info to distinguish one object with

idea

other object: $< = \neq >$

zipcodes, empidnos, eyecolor, gender are examples of nominal.

- b) Ordinal value: The values of ordinal attributes provide info to order the objects ($<, >$), ranking.

{Good, better, best}, {High, Medium, low}
Basing on context we try to segregate it from higher to lower.

interval

Numeric data : [Quantitative]

- a) Interval data: For interval attributes the diff b/w values are meaningful where a unit of measurement exists. (+, -)

Ex: Thermometer (Fixed value) + low fever | high fever
Calendar dates

- b) Ratio data: For ratio variables both differences ratio are meaningful (\times, \div)
{Electric current, Temp in Kelvin, Count values, Age, Mass, Length} Physical values.

Discrete and Continuous Variables / data / attributes:
↳ no floating | decimal exists

{1, 2, 3, 4, ..., ∞ }

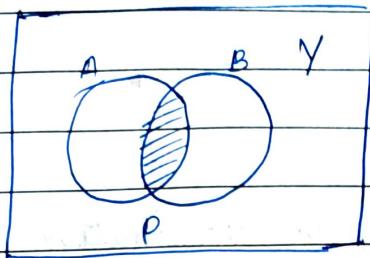
continuous - floating | decimal exists
{1, 1.1, 1.2, 1.3, ..., 2, ..., ∞ }

14/8/24

Sampling

Represents entire population

→ $\langle a-z, A-Z \rangle, \langle 1, 2, \infty \rangle,$
 $\langle \text{colleges} \rangle$



$A \cap B$ - sampling / ep

Covid - 19:

Problems → spreadful, No diagnostics, Symptoms are changing. Variations

Based on samples
Age, Gender, localities, Genetic
diabetes, B.P.

they made a sample space

Central Limit theorem

s ∈ P. if the soln works for subset's
then the soln works for whole population 'P'

Precautions

Descriptive statistics

Inferential statistics.
→ Experimentation.

i. Outcome is known

2. Past data

3. Summary / characteristics.

(Precaution - covid-19)

(After vaccination the reactions)

- It is a branch of stats that involves organization summary of data → It deals with generalisation about a population based on sample data & it draws the conclusion from sample space's
 - Usually, mean, median & mode are used to describe the avg values of data & are referred to be central tendency in nature → Applications are regression analysis, hypothesis tests. (Assumptions)
(Reln b/n independent & dep var's)
 - We can measure shape & size of probability distribution.
- Normal dist.

mean
= median
= mode

Data is uniform

Example of Covid -19 :

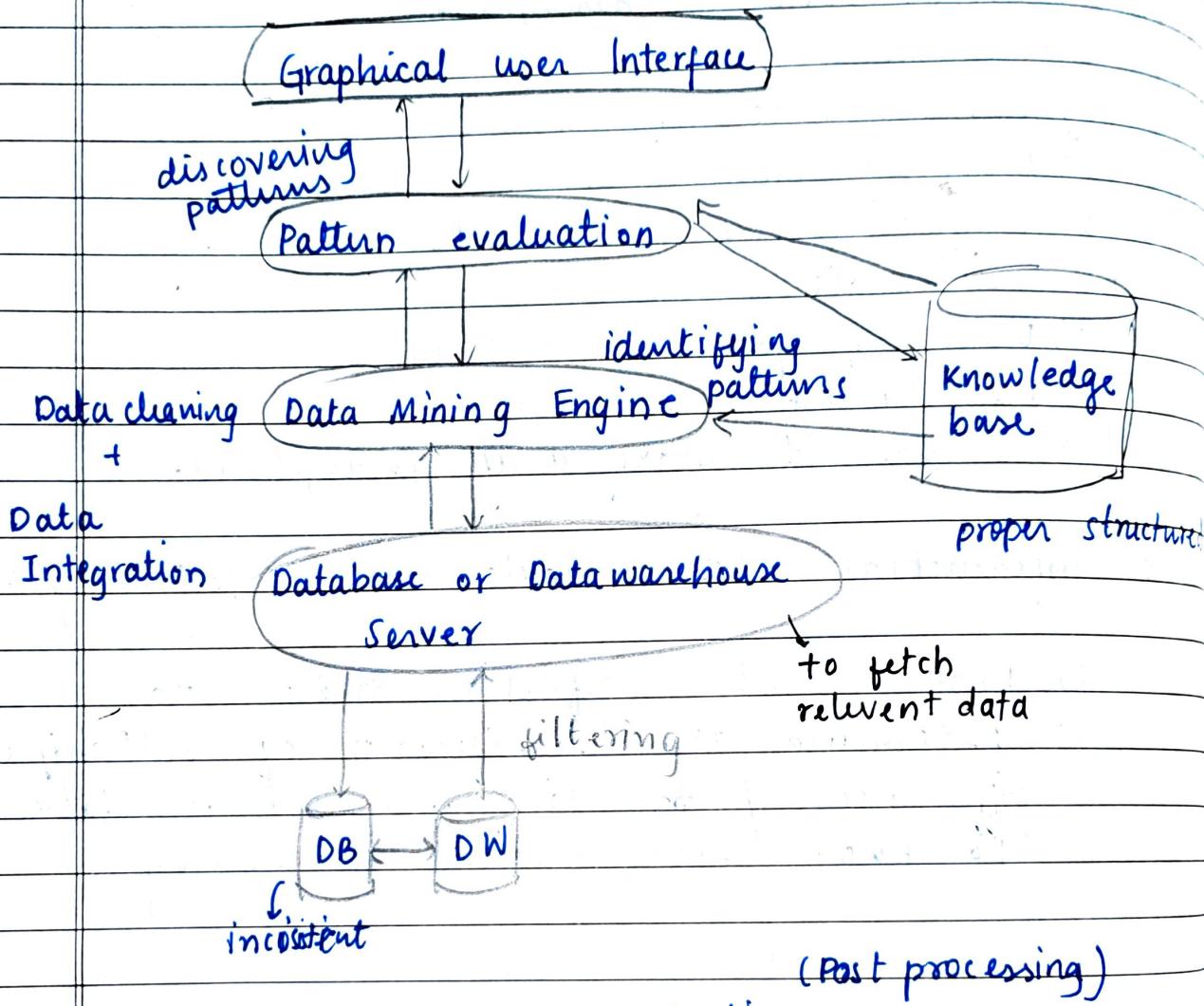
- * In case of covid-19 descriptive stats deals with the following
 - 1) provides info of symptoms for covid
 - 2) dist of cases across various regions
 - 3) Frequency of complaints given by individuals
 - 4) Usage of preventive measures like wearing gloves, masks, vaccinated, and so on
 - 5) it doesn't make any assumptions.

- * Inferential stats deals with following:

- 1) Deal w the infection rate, to
- 2) Determine the significance of relationships, differences / trends observed in data
(to know how vaccine is working)

* *Important Question*

Typical Data Mining Architecture



- Focuses on pattern evaluation whereas KDD focuses on knowledge discovery.

(pre processing)

Although we have knowledge base, we use DB as it is dynamic.

Miner's job is:

interesting patterns from where we can derive something.

from all the patterns generated from data mining Engine interesting patterns are identified.

Data Mining Engine:

- Provides the functionality of classification, clustering, prediction, outlier analysis, association analysis, evolution analysis.

Knowledge base engine:

- Discovering the hidden patterns which are at the measure of interestingness.
- It includes concept hierarchy
- It may also include the users belief which relies on multiple hypothesis.

Pattern evaluation module:

Thresholds are set to filter out the discovered patterns.

GUI:

Its a communication channel b/w user and client, A query is given, searching takes place through exploratory datamining & then reports are generated based on visualizations.

representing theoretical in numerical form.

19/8/24

*Classification of Mining Techniques :

Frequent patterns finding \Rightarrow Data mining

Case 1 : customer id, product purchase, sales - investment
profit / loss %

discount, trends,

primary factor - customer visiting regularly.

2) Text mining :

case - 2 : Online shopping.

customer reviews, ratings, opinions, feedbacks

+ve/+ve (text)

feedback - +ve/-ve/neutral

sentiment - functionalities.

text - opinions online - Navigating patterns

text / opinion mining (Amazon recommendation sys).

1. Data Mining :

Ex : Customer id, purchase date, product id, quantity purchased.

2. When you try to apply any sentiment on text / any review. It's termed as opinion mining

opinion mining is subset of text mining

3. Web mining:

case - 3: Social media / Group communities.

- Marketing, Advertising.

update latest info - (linked in)

large investment \Rightarrow investment patterns.

4

network / Graphs mining

(searching becomes easy)

bfs, dfs, nodes, links, dynamic.

bfs: nodes, links,

edges, weights.

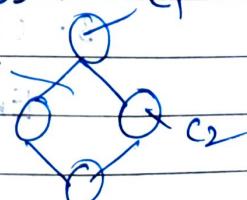
Ex - Amazon reviews Clique stream analysis,
topic modelling

We focus mostly on navigating pattern
hence termed it as web mining.

4. Graph mining:

Social network analysis online media,
Group communities, influential users.
for targeted marketing.

frequent buyers!



5. Trajectory mining:

Any movement patterns analysis of traffic flow, GPS

Main objective of trajectory mining is to optimize the traffic flow.

case - 4: Moving objects / traffic analysis

Trajectory mining

6. Association mining:

1. Data mining on which kind of data can be mined?

A- 1) We can mine the data on the general

1. Relational Database:

Rows, Cells, Entities.

2. Transactional Database:

wrt. Individual customers

wrt. Product similarity

Association Analysis.

Transaction Id	Products
10001	Tea, Milk, Honey, Coffee

3. Object - oriented database:
Encapsulation & Reuse

Each entity is considered as object

4. Object-Relational Database :

- handles complex queries.
- Uses concept of class hierarchy.

5. Spatial database: GPS, Map database,

Medical images - (Raster format)

↓ satellite images

(To find out park near to house)

6. Temporal database:

Focuses on timestamp ordering.

7. Timeseries database:

Focuses on periodic trends.

(unstructured)

8. Text database : www , large documents

9. Multimedia database :

audio + video + text + speech combining.

Any content deliver.

↗ historic data

10. Heterogeneous & legacy database

It combines systems such as relational, object oriented, N/w based, hierarchical approaches, file systems, spread sheets & multimedia.

Data Characterization & Discrimination.

1. Data can be associated with class / concept.
 Ex: In a retail shop. set of items considering
 whole computer, printer, ... → Class.
 Specifying **bigspenders, budget spenders** → Concept.

Data characterization: Summarizing data.

2. It also describes individual class / concept in precise & summarised manner.
 such type of summarised manner / description are called concept hierarchies.
 It also explains features of data.
- Summarizing the data under study is called target class

Examples of Data characterization

It customer purchasing a product, sales, avg sales

Data Discrimination:

It compares target class with another class from the set of items.

Computer ↗ s/w - Functionalities

Computer ↗ H/W - CPU, mouse, keyboard, ...

HP computers	Dell
90% Buyers frequent	10% Buyers infrequent

1A

* Data cleaning:

1. Missing values 2. Noisy data

1. Missing values :

1. Ignore the tuples.

if there are 10 diff columns out of which 10th col represent % & 9 cols are missing

Rollno	name	mark1	m2	m3	m4	...	per
-	-	-	-	-	-	-	90% ((ignore it))

2. Fill the missing values manually:

If its a small dataset, the missing values can be filled manually ex: Avg of marks.

3. Use a global constant such as 'unknown' or '-oo'

If all the values are missing from that column, the it takes "unknown" as a pattern

4. Use the attribute mean to fill the missing values:

Ex: Avg income of all electronics customers is \$28,000 we use this value to replace the missing - , height of missing individual.

5. Use the attribute mean for all the samples belonging to same class as given tuple.
- Ex: If classifying cust acc to credit risk, replace missing value w any income value for customers in the same credit risk category as that of given tuple.

6. Use the most probable value to fill the missing values.

Regression model, decision tree induction methods

2. Noisy Data:

Noise is a random error / variance in a measured variable.

Given a numeric attribute such as price, how can we smooth the data to remove the noise, this could be achieved through clustering mechanism.

Random Space



distance based approach

If we have + & * we try grouping
+'s on 1 side & *'s on another

But if we also have Σ & α
then + & * form 1 group & Σ, α as
one

3 Regression:

Regression model best fits the relation
b/w independant & dep. var's & is
formed from the base point,

$$y = mx + c.$$

4. Combine Human & Computer introspection:

$\Sigma = \{ \text{alphabet} \}$ if 'o' & '7' are input
they are

Data Transformation & Integration

→ Smoothing of Data - removing noisy data

We have 3 diff approaches

1. binning
2. clustering
3. Regression.

→ avg. marks, fail %

→ Data aggregation : summarize data from different forms & levels (Granularity) & generalizing (only on sample data)

→ Data Generalization : (on Entire data / population)

→ Data Normalization :

Generally we represent data from (0-1) or z-score [scaling]

→ Performing analysis is easy

Expected problems.

5M * Concept Hierarchies

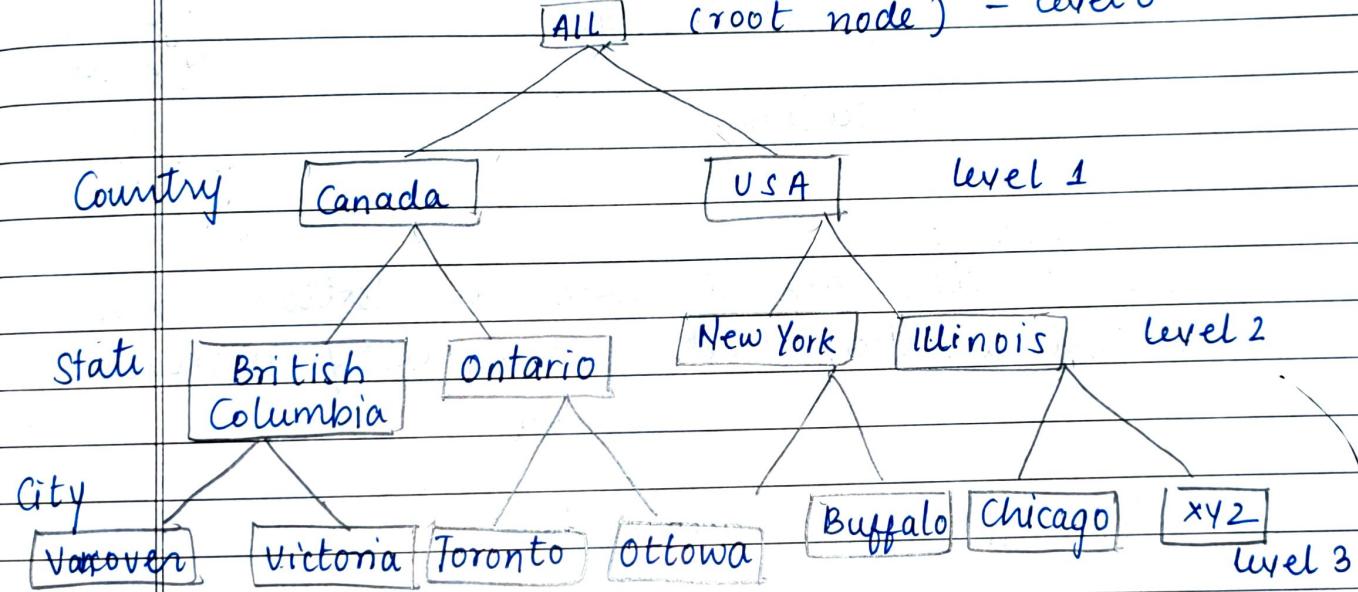
Concept hierarchy defines sequence of mappings from a set of low level concepts to higher level in turn resulting in general concepts.

- subsets are low level whereas supersets are high level

deciding factor
continent / location

Dimension attribute

(root node) - level 0



Concept Hierarchy for attribute dimension

Dimension attribute is described by attribute no., street, city, ... country

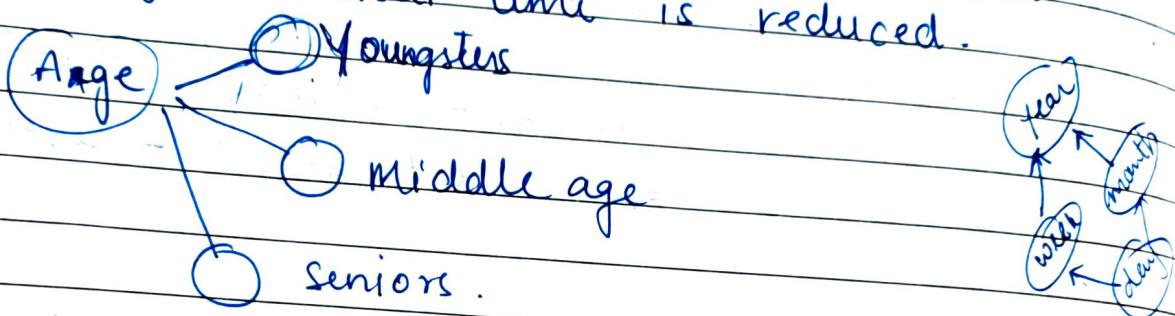
These attributes are related to total order

forming a concept hierarchy such as
street < city < state < country (A.O.)

↓
higher granularity
multiple "

Concept Hierarchy is the total/partial order for time dimension
day < week < month < year < decade
& this can also be classified as seconds.

Advantage is we are generalising the data.
searching traversal time is reduced.



drawbacks :

1. Construction time is more at bottom levels.
2. Implementation level - complexity increases

Graph
(Directional)

4 diff types of Concept h

1. Schema hierarchy
2. Set grouping hierarchy
3. Operation based hierarchy
4. Rule based hierarchy

Preprocessing & Discovering patterns can be automated
but evaluating patterns can't be automated
(human intervention is required) Knowledge Engineering

Data Engineer - Analysis Reports/projects visualization
Raw → actual

10M

tasks

Data mining primitives : (5M)

steps for Data mining tasks / query:

1. Task relevant data (5m)
2. what kind of knowledge to mine? (5m)
3. what type of background knowledge is required?
4. what kind of interesting measures we require for patterns discovery in Data mining.
5. How would you represent the data (mined) ?

Myth: Data mining systems can autonomously ~~pick~~ out all the valuable knowledge that is embedded in a given large database with human intervention.

1. Task relevant data.

Ex. 1) A group of boys segregated rank wise

- Physical appearance /
- Marks
- Skills

The selection depends on End user biased cond'n → colour, religion

Specificity is → Placement requires skills
→ NCC / Army requires physical

Ex-3) If a data mining task is to study the associations b/w the items that are frequently purchased at all electronics by the customers in Canada, what could be specified in task relevant data.

* dirty data: even after preprocessing, if no meaningful patterns are discovered.

1. Task relevant data can be specified thru name of the database / data warehouse all_electronics_db. which includes tuples, tables, containing info such as purchase rate, customer details, items sold.

cno	item	phno

Item no	name	pn

2. We probably define the conditions for task relevant data.

Ex - Retrieve data i.e. pretraining to purchase made in china for the current year.

3. Relevant attributes / dimensions.

Ex: name & price from item table
income & age from customer table

2. What kind of knowledge to be mined?

→ Electronic store

Customer purchasing computer is likely to buy a printer
associated (Reln).

Printer is dependent on computer
User of probability - 60% likely to happen
defined Confidence - 27% already happened
(not fixed)

→ based on Association Analysis.

3. What kind of Background knowledge is required

By using Concept hierarchy, they represent low level data & high level data

→ Generalization is used. used in Backend.

1) Fast Relevant data

→ Database / Datawarehouse, conditions for data selection.

2) Knowledge to be mined

→ Characterization & discrimination.

→ Classification & prediction.

→ Clustering & association.

3) Background knowledge to be mined

- 1. Concept hierarchies
- 2. User beliefs about relationships in data

4) Pattern interesting measures

- 1. Utility⁽²⁾, Certainty⁽¹⁾, simplicity⁽⁴⁾, novelty⁽³⁾, support & confidence

★ ★ Presentation & Visualisation of discovered patterns

I. Rule based.

$\text{age}(x, \text{'young'}) \text{ and } \text{income}(x, \text{High}) \Rightarrow x \in Y_2$.

Define rule w certain conditions.

space
complexity
generalized
version.

II Table based.

Age	income	class	count
		A	
		B	
		C	

data characterization.

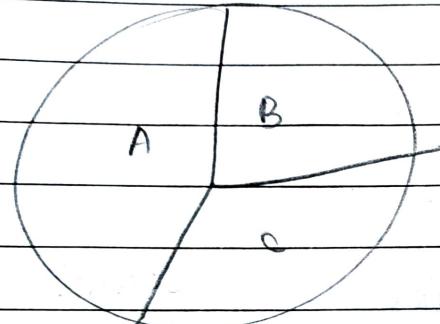
III Cross tab specific concept

age	income		class		
	high	low	A	B	C
young					
old					
adult					

used for data discrimination

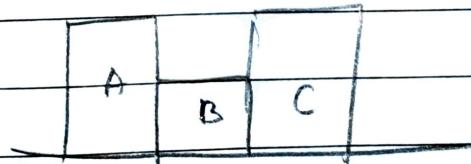
IV

Pie-chart



To highlight some features.

(V) Bar chart

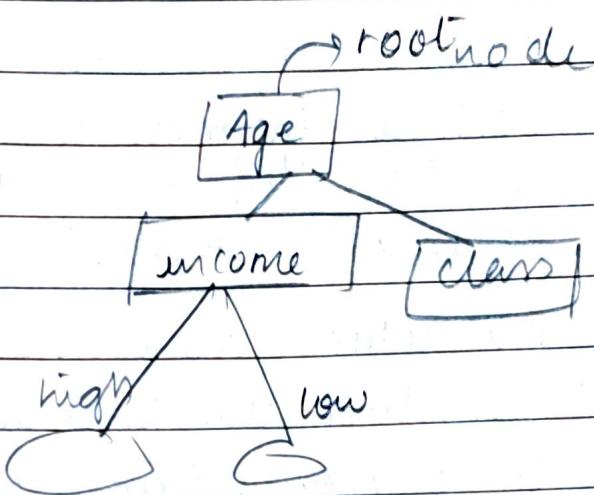


VI Decision tree

tree data structure

ML model

Works very well
for large data

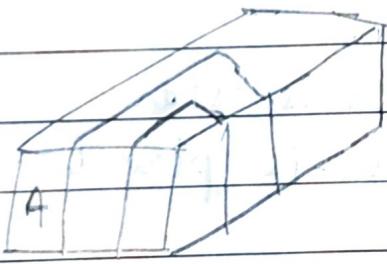


VII

Data cubes - Data warehouses in xyz

axes

Dimension
space



each class
represent
1 dimension

27/8/24

Date: 1 / 1

Knowledge Representation methods.

→ o/p is decision making.

What is knowledge?

(Derived information (from unstructured data).
↳ is also called knowledge.)

Advantage: Analysis

Knowledge representation (KR) is process of extracting the knowledge from derived info, to perform certain tasks for analysing, basing on performance speed evaluating its measures & generating the summary inform of reports, visualisation, graphs, ...

Truth / Logic / Myths:

↓ whatever has mathematical approach

Sun is star - truth → Earth is spherical

sun is planet ↗ false ↘ myth → Earth is flat

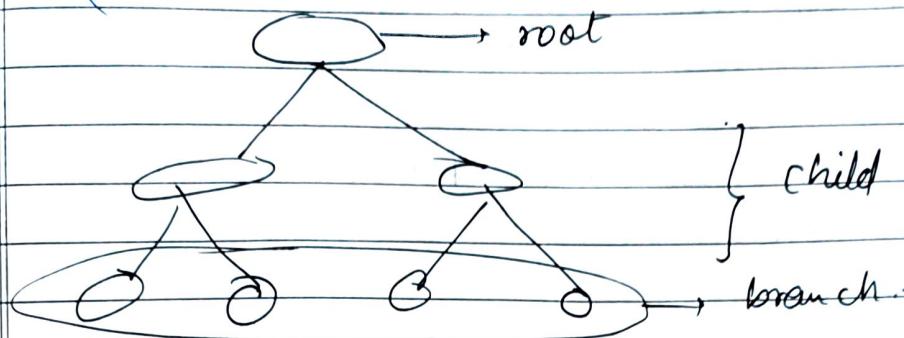
Acc. to Mythology, Sun is God (plant) → truth
(SunGod Ra)

Acc to Flat theory, Earth is flat - - truth.

Various forms:

1. Decision Tree: i) Used for classification & Regression problem.
derived from BST.
searching becomes easy as well as cleaning
- 2) It represents knowledge as a tree like structure where each node to the attribute & each branch represents a decision

outcome of the decision.



•

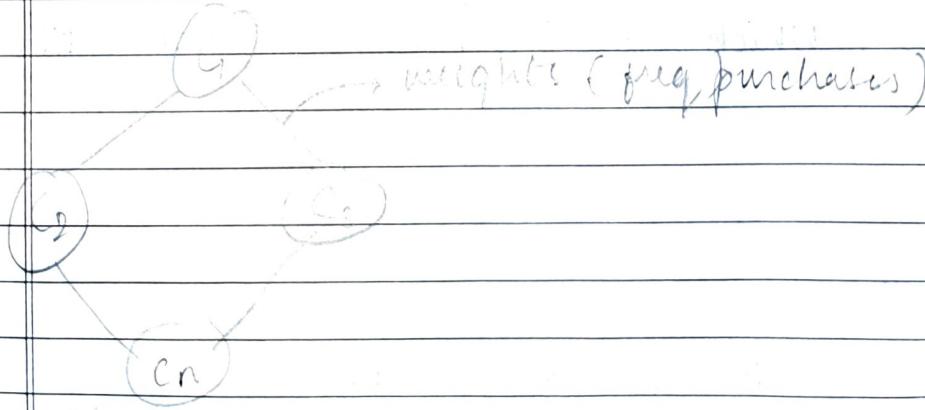
Example: A decision tree for classifying whether a customer will buy a product basing on the attributes like age, income / previous purchase history.

2. Association Rules:

3. Graphs & Networks:

↓
easy to represent large amt of data.

(Social N/w)



Graph represents data as nodes & edges
The main applications of graphs are

Recommendation systems.

Social Network analysis

4. Rule based Systems : Chat GPT

If - then.

Traffic.

- Maintaining Consistency.

Disadv : Hard Threshold. (no feasibility).
(Objective exams)

Example: In a medical diagnosis sys, a rule might be 'If the patient has high fever & a sore throat then patient might have severe cough'

5. Ontologies and Semantics :

Object = which describes behaviour of entity.

Protege editor

Manual efforts are more

A is a good boy
Entity Refn object

3/9/24

UNIT - 2

Date: / /

Market Basket Analysis.

In late 1980's, the Although, the products are present in the shop, the consumer didn't buy them because they are unorganised.

Market Basket.

{milk, bread} → {butter}

: Rule

Association /
coherence Analysis.

- Given a set of transactions, find the rules that will predict the occurrence of an item based on occurrence of other item in the transaction. This concept is referred as Market basket analysis / transactions.

Tid	Tname
1	bread, milk
2	bread, diaper, chocolate, eggs
3	milk, diaper, chocolate, coke
4	bread, milk, diaper, chocolate
5	bread, milk, coke, diaper.

bread - item {bread} - itemset

K-items set : $k_d = \{k_1, k_2, k_3, k_4, \dots\}$.

- Support count (σ): Frequency of occurrence of an item set.

Ex: $\tau_{milk, bread, diaper} = 2$.

3. Support (s): It is fraction of transaction that contains as item set. $(s) = 2/5$

4. Frequent item set: count
An item set whose support is \geq minimum support is called a frequent item set
min support = threshold value = 0.

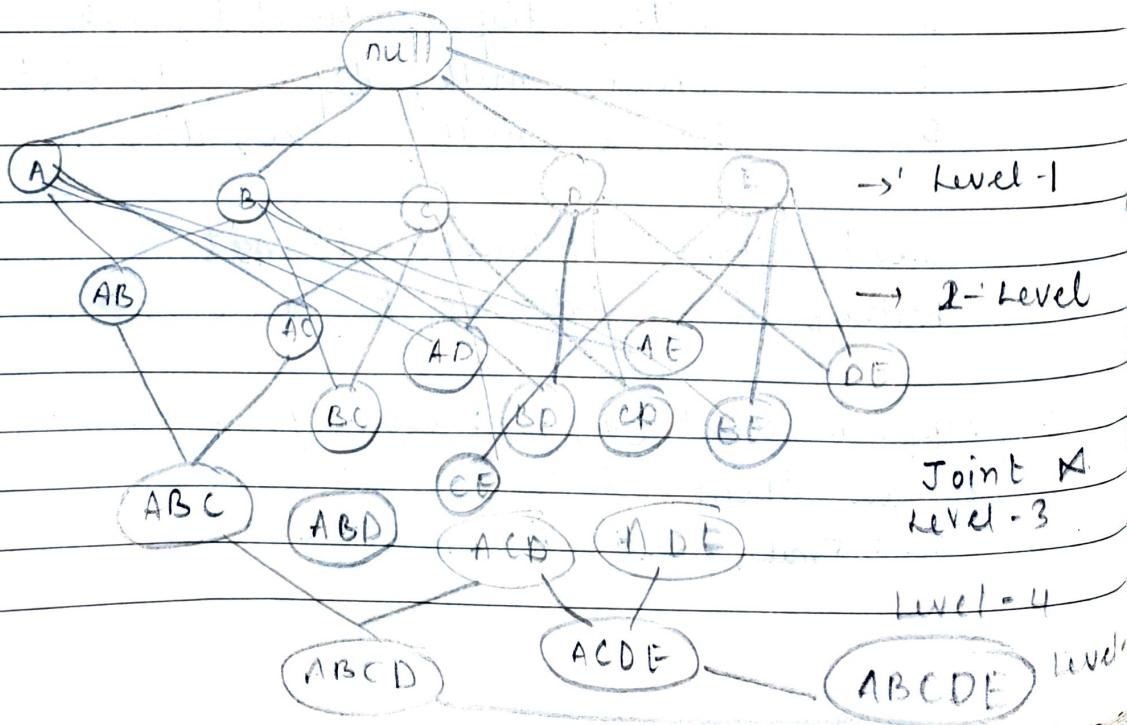
5. Confidence: (c)

It measures how often items in Y appear in transaction that contain X

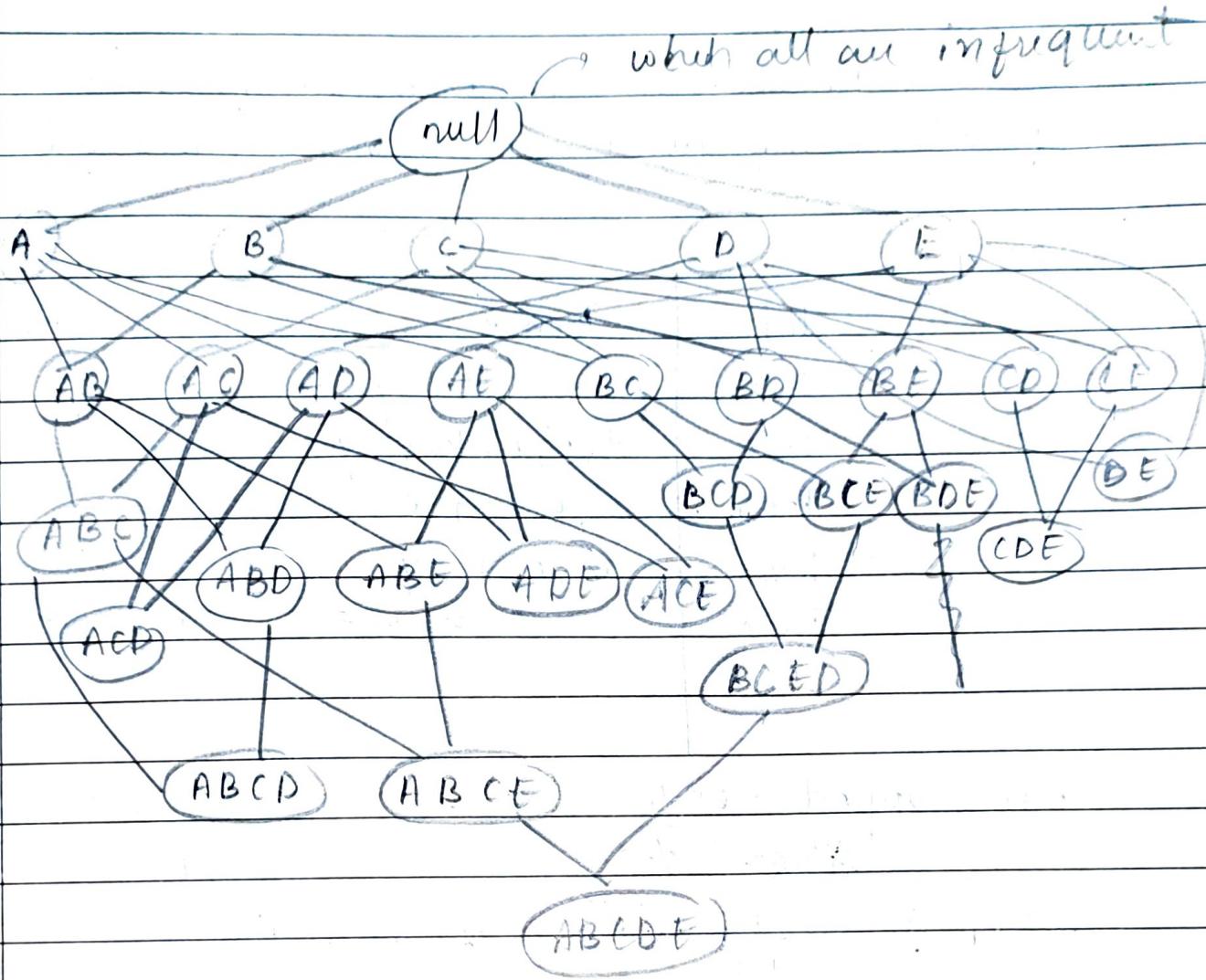
→ If milk is frequently brought & Rice is frequently brought separately, it need not be that {milk, Rice} is a frequent item set.

* Brute Force method / Approach:

If we have d items, possible no. of itemsets are 2^d . (no repetition).



→ Drawback:
Computationally prohibited / expensive



4/9/24 drawback:

No track of count

If A is infrequent then all the sets consisting of A are inconsistent.

→ We need not draw a tree for lower freq.

Continuous semantic

Pruning is difficult.

* Apriori Algorithm

Apriori principle:

A is infrequent AB will never be frequent
If AB is infrequent A may be infrequent.

Example:

Tid	Iname	(m) 3 ✓
T ₁₀₀	m, o, n, k, e, y	< o > 4 ✓
T ₂₀₀	d, o, n, k, e, y	< n > 2
T ₃₀₀	m, a, k, e	< k > 5 ✓
T ₄₀₀	m, c, k, y	< e > 4 ✓
T ₅₀₀	c, o, o, k, i, e	< y > 3 ✓ < d > 1 < a > 1 < c > 1 < i > 1

$$\text{min support} = 60\%.$$

$$\text{min confidence} = 80\%.$$

$$\text{support} = \frac{60}{100} \times 5 \xleftarrow{\text{(no of transactions)}} 3.$$

take sets having \geq support value.

L₂

< m, o > 1

4

< m, k > 3

m	3
o	4
k	5
e	4

< m, e > 2

< m, y >

< o, k > 2

< o, e > 3

< k, e >

< k, y > 4.

L3

$$\begin{array}{ccc} \langle m, k, o \rangle & 1 & \\ \langle m, k, e \rangle & 2 & \Rightarrow \langle o, k, e \rangle \\ \langle o, k, e \rangle & 3 & \end{array}$$

No limit for association rules.

Interesting patterns : $\langle o, k, e \rangle$

Association rule Support Confidence confidence.

$$\begin{array}{lll} \text{O} \wedge \text{K} \Rightarrow \text{E} & 3 & \frac{3}{6} = 3/3 = 1 = 100\% \\ \text{O} \wedge \text{E} \Rightarrow \text{K} & 3 & \frac{3}{6} = 3/3 = 1 = 100\% \\ \text{E} \wedge \text{K} \Rightarrow \text{o} & 3 & \frac{3}{6} = 3/4 = 0.75 = 75\% \\ \text{e} \Rightarrow \text{o} \wedge \text{k} & 3 & \frac{3}{6} = 3/4 = 0.75 = 75\% \\ \text{k} \Rightarrow \text{o} \wedge \text{e} & 3 & \frac{3}{6} = 3/5 = 0.6 = 60\% \\ \text{o} \Rightarrow \text{e} \wedge \text{k} & 3 & \frac{3}{6} = 3/4 = 0.75 = 75\% \end{array} \quad \left. \begin{array}{l} \text{interesting} \\ \text{patterns.} \\ \text{having} \\ \text{than } 80\%. \\ \text{confidence} \end{array} \right\}$$

n = support count

$$\begin{array}{l} n! = 3! \\ = 6 \end{array}$$

* FP Growth Algorithm:

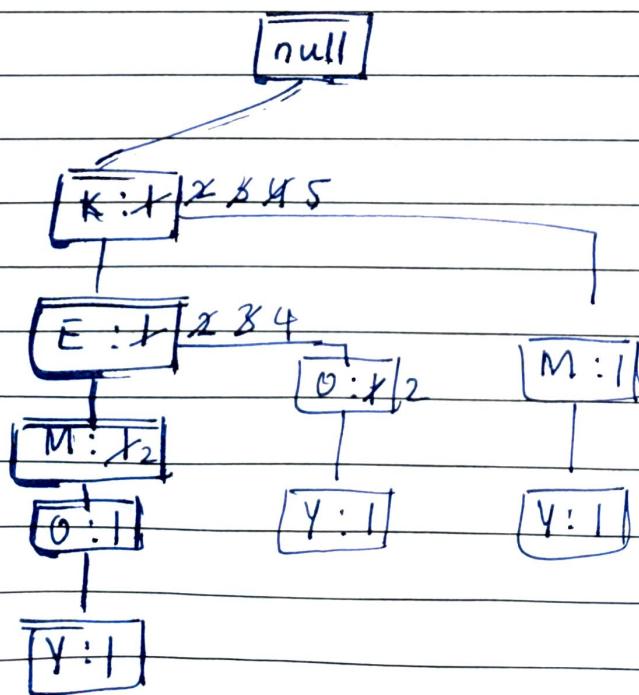
Mining the complete set of frequent patterns using a tree structure for storing info abt FP.

min support ≥ 3 .

Tid	Tname	<u>ordered list-item</u>
T ₁	e, n, m, e, K, y	{K, e, m, o, y}
T ₂	K, y, e, n, d, o	{K, E, O, Y}
T ₃	K, m, a, e	{K, E, M}
T ₄	m, y, c, k, u	{K, M, Y}
T ₅	c, i, e, k, o	{K, E, O}

K = 5, E = 4, M = 3, O = 3; Y = 3.

FP - Growth tree construction



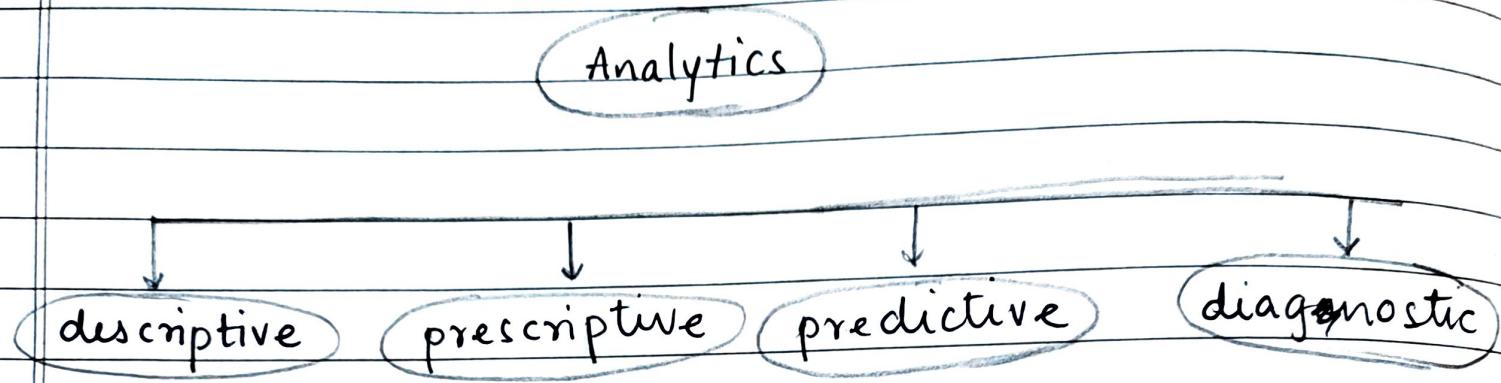
items	Conditional pattern base	Conditional fp-tree
Y	{ KEM:1, KE0:1, KM:1 }	{ K:3 }
O	{ KEM:1, KE:2 }	{ KE:3 }
M	{ KE:2, K:1 }	{ K:3 }
E	{ K:4 }	{ K:4 }
K	{ Ø }	{ Ø }

Items	FP-Growth
Y	{ K, Y:3 }
O	{ K, O:3, E, D:3 }
M	{ K, M:3 }
E	{ K, E:4 }
K	{ Ø }

Drawbacks of FP Growth :

- If its a large database tree construction becomes complex
- After every step we have to see association rules.

Unit - 3 Data modelling



Forecasting techniques

Heuristic
approach

Moving average

10/9/24

Date: / /

1. For each of the following provide an example of an association rule from the market basket domain that satisfies the following condition also describe whether such rules are subjectively interesting.

- 1) A rule that has high support & high confidence
- 2) A rule that has high support & low confidence.
- 3) low support & low confidence
- 4) A rule that has low support & high confidence

2.

2. Consider the following dataset (8m)

Customer_id	T-id	items bought
1	0001	{a, d, e}
1	0024	{a, b, c, e}
2	0012	{a, b, d, e}
2	0031	{a, c, d, e}
3	0015	{b, c, e}
3	0022	{b, d, e}
4	0029	{c, d}
4	0040	{a, i, o}
5	0030	{a, d, e}
5	0038	{a, b, e}

- a. Compute the support for itemsets $\{c\}$, $\{b, d\}$ and $\{b, d, e\}$ by treating each t-id as a market basket.
- b. Use the results of Part(a) to compute the for an association rules $\{b, d\} \rightarrow \{e\}$ & $\{e\} \rightarrow \{b, d\}$. Is confidence a symmetric measure?
- c. Repeat part (a) by treating each customer-id as Market basket. Each item should be treated as binary variable. if an item appears in atleast one transaction bought by customer & 0 otherwise

d. Use results in part c to compute the confidence for the association rule

$$\{b, d\} \rightarrow \{e\} \quad \& \quad \{e\} \rightarrow \{b, d\}$$