

Step 2: Information  
↓  
processed data

Approaches:

- ✓ structured format
- ✓ Tables of rows and columns
- ✓ Data cleaning  
↓

Handling inconsistencies

Missing values

Noisy data

feature Extraction

e.g.: calculate average grades, participation levels, and study hours.

Visualization:

Graphs, charts, ..... to show trends such as the correlation between study habits and grades.

- ② 1. Handling right feature  
2. Incomplete data (filling).

Sampling \*\*\*  
↓

Step 3: Knowledge:

- ✓ To derive insights
- ✓ Truth

Descriptive analysis



use summary statistics

Trends → average grades & study times

Correlations

Study habit  $\propto$  performance

Predictive analysis

$y_d$

- ③ 1. Validating. (summarization ✓)  
2. ML ↗ Training  
Testing

Step 4: Intelligence:

Decision making & actions

1. personalized feedbacks
  2. Tutoring sessions
  3. Assignments
  4. Feedback analysis
- e.g.: Automated E-mails  
Dashboards

- ④ feedback is constructive &  
1. actionable  
2. Balancing the need for  
intervention with students autonomy

### ⑤ Knowledge:

Individuals interest / choice

④ Personalized recommendations.

③ cognitive science

Simulate user interactions, learn from user experience. | adapt recommendations to align with change of user interests and preferences.

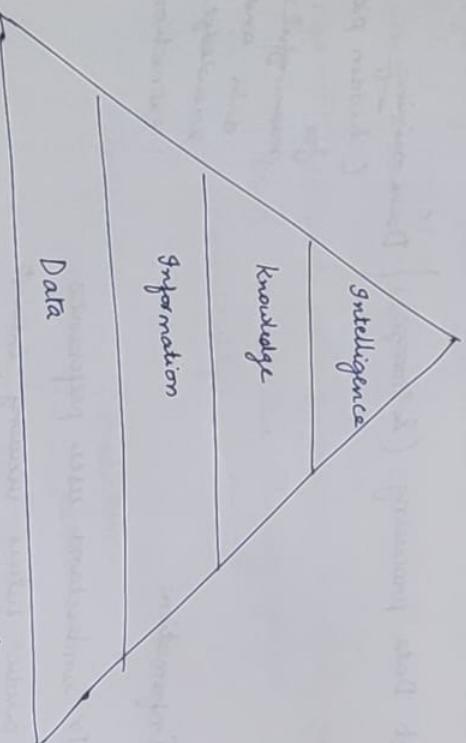


fig.11: Pyramid of data (hierarchy).

Evolution of database technology.  
1960's primitive file processing

DBMS

1970  
1980s

Advanced db systems. | web based.  
mid 1980's.

XML | web mining

Rdbms multimedia

Dm.

→ Extracting / mining knowledge from large amounts of data is Data mining / knowledge mining / knowledge extraction / data pattern analysis / data archaeology / data dredging / knowledge Discovery in databases. (KDD)

\* mining of gold from rocks / sand is referred as gold mining rather than sand / rock mining.

\* mining of gold from rocks / sand is referred as gold mining rather than sand / rock mining.

KDD process

1. Data cleaning (to remove noise and inconsistent data).
  2. Data integration (multiple data sources may be combined).
  3. Data Selection (Data relevant to analysis task area retrieved from db).
  4. Data Transformation (into forms, for mining by summary generation).
  5. Data mining (intelligent methods are applied in order to extract data patterns).
  6. pattern evaluation (interestingness measures)
  7. knowledge presentation (present knowledge to the user).

## \* Data Integration: characterization

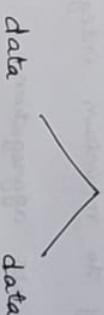
→ Data can be associated with classes / concepts.

Example:

All electronics store : < computer, printer > - classes  
Concepts of customers include < big spender >,  
< budget spender >

→ describes individual class / concept in precise and  
summarized manner.

→ Such ↑ descriptions are class / concept descriptions.



characterization



↓  
Summarizing the  
target class  
with another  
class).

Ex:

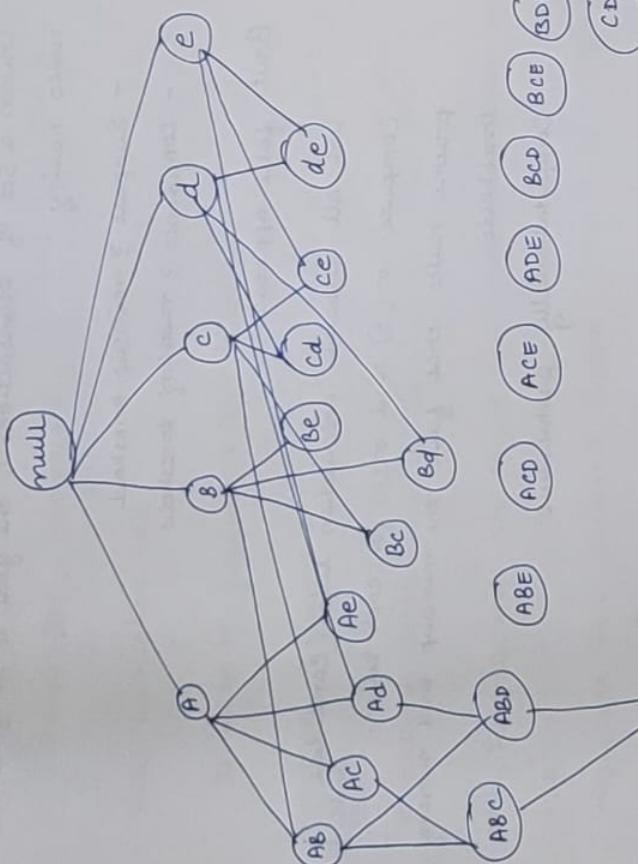
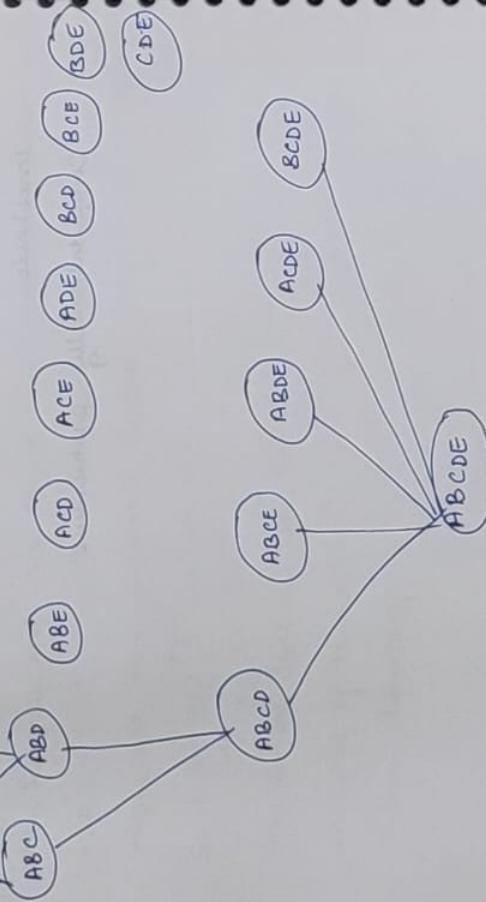
Given  $d$  items,  $2^d$  possible candidate itemsets are possible

$A \ B \ C \ D \ E$  — itemset

$\left. \begin{array}{l} AB \ AC \ AD \ AE \\ BC \ BD \ BE \\ CD \ CE \\ DE \end{array} \right\}$  — 2 itemsets

$\left. \begin{array}{l} AB \ AC \\ AB \ AD \\ AB \ AE \\ ABC \end{array} \right\}$

$\left. \begin{array}{l} AB \ BE \\ ABD \end{array} \right\}$



Point Status

$P_1$  noise B

$P_2$  core

$P_3$  noise B

$P_4$  noise B

$P_5$  core

$P_6$  noise B

$P_7$  noise B

$P_8$  noise B

$P_9$  noise X

$P_{10}$  noise B

$P_{11}$  core

$P_{12}$  noise B

② Apply abScan algorithm with Similarity threshold  $\theta$  of 0.8 (using the similarity matrix) to the given data points and minpts  $\Rightarrow 2$  (min. no. points in cluster)

what are core, border and noise points in set of points in the table.

$P_1$   $P_2$   $P_3$   $P_4$   $P_5$

$P_1$  1.00 0.10 0.41 0.55 0.35

$P_2$  0.10 1.00 0.64 0.47 0.98

$P_3$  0.41 0.64 1.00 0.44 0.58

$P_4$  0.55 0.47 0.44 1.00 0.26

$P_5$  0.35 0.98 0.85 0.26 1.00

$P_1$  :-

$P_2, P_5$

$P_2, P_5$

$P_1$ :  
 $P_5, P_2, P_3$

$P_5$

$P_1$  noise

$P_2$  core

$P_3$  core

$P_4$  noise

$P_5$

core

## Association Rule Mining Task:

$$\begin{aligned} \text{Support} &\geq \text{minSupp} \\ \text{confidence} &\geq \text{minConf} \end{aligned}$$

$\{ \text{minSupp} = 30\%, \text{minConf} = 50\% \} \rightarrow$  Valid association rule

Valid association rule.

It must be frequent + high confidence.

Purchase of cricket bat and ball might be frequent but purchasing togetherness is necessary for valid association rule.

May be milk and rice might have high  $\sigma(S)$  support count but might not have purchased together (not having  $\sigma(C)$  confidence) So it is not valid Association rule

How do we discover the rules?

millions of transactions?

Pattern discovery.

Given a set of transactions  $T$ , the goal is to find all rules having

- Support  $\geq$  minsup threshold
- confidence  $\geq$  minconf threshold

Brute force approach:

list all possible association rules (candidate rules).  
compute  $\sigma(S)$  and  $\sigma(C)$  for each rule.  
prune rules that fails the minsup and minconf thresholds  
computationally prohibitive.

① Apply DBScan algorithm to the given data points

and

Create the clusters with

minpoints = 4 and epsilon ( $\epsilon$ ) = 1.9

Data points

$P_1: (3, 4)$   $P_2: (4, 6)$

$P_3: (5, 5)$   $P_4: (6, 4)$

$P_5: (7, 3)$   $P_6: (6, 2)$

$P_7: (9, 2)$   $P_8: (8, 4)$

$P_9: (3, 3)$   $P_{10}: (2, 6)$

$P_{11}: (3, 5)$   $P_{12}: (2, 4)$

$$\epsilon: 1.9$$

$P_1: P_2, P_{10}$

$P_2: P_1, P_3, P_{11}$

$P_3: \cancel{P_1} P_2, P_4$

$P_4: P_3, P_5$

$P_5: P_4, P_6, P_7, P_8$

$P_6: P_5, P_7$

$P_7: P_5, P_6$

$P_8: P_5$

$P_9: P_{12}$

$P_{10}: P_1, P_{11}$

$P_{11}: P_2, P_{10}, P_{12}$

$P_{12}: P_9, P_{11}$

Euclidean distance.

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

	$P_1$	$P_2$	$P_3$	$P_4$	$P_5$	$P_6$	$P_7$	$P_8$	$P_9$	$P_{10}$	$P_{11}$	$P_{12}$
$P_1(3, 4)$	0											
$P_2(4, 6)$	1.41 ✓	0	0									
$P_3(5, 5)$	2.83	1.41 ✓	0									
$P_4(6, 4)$	4.24	2.83	1.41 ✓	0								
$P_5(7, 3)$	5.66	4.24	2.83	1.41 ✓	0							
$P_6(6, 2)$	5.83	4.47	3.16	2.00	1.41 ✓	0						
$P_7(9, 2)$	6.40	5.00	3.61	2.24	1.00 ✓	1.00 ✓	0					
$P_8(8, 4)$	5.83	4.47	3.16	2.00	1.41 ✓	2.83	2.24	0				
$P_9(3, 3)$	4.00	3.16	2.83	3.16	4.00	3.16	4.12	5.10	0			
$P_{10}(2, 6)$	1.41 ✓	2.00	3.16	4.47	5.83	5.66	6.40	6.32	3.16	0		
$P_{11}(3, 5)$	2.00	1.41 ✓	2.00	3.16	4.47	4.24	5.40	5.10	2.00	1.41 ✓	0	
$P_{12}(2, 4)$	3.16	2.83	3.16	4.00	5.10	4.47	5.39	6.00	1.41 ✓	2.00	1.41 ✓	0

Given a set of transactions find rules that will predict the occurrence of an item based on the occurrences of other items in the transaction.

Item	Market basket transactions
1	Bread, milk
2	Bread, Diaper, Beer, Eggs
3	Milk, Diaper, Beer, Coke
4	Bread, Milk, Diaper, Beer
5	Bread, Milk, Diaper, Coke

$$\begin{aligned} \text{Eg: } & \{\text{Diaper}\} \rightarrow \{\text{Beer}\} \\ & \{\text{milk, bread}\} \rightarrow \{\text{Eggs, coke}\} \\ & \{\text{Beer, Bread}\} \rightarrow \{\text{milk}\}. \end{aligned}$$

Observe patterns in form of rules

LHS:                    RHS:  
Itemset                Itemset

? Association?

Cohesive ✓ exists  
↓  
bought together purchased

discount

Bread, milk  $\rightarrow$  eggs(free) may be:

$$\frac{40}{5} > \text{min\_sup}$$

$$\frac{0.4}{5/20} = 0.16$$

Eg: train ticket  $\rightarrow$  taxi ticket

Valid rules?

Data driven. (no hypothesis).

Frequent itemset.

itemset:  $\{\text{milk, bread, Diaper}\}$

k-itemset:  $\{k_1, k_2, \dots, k_{n-1}\}$  if  $n > k$

Support count ( $\sigma$ ):

frequency of occurrence of an itemset

$$\sigma(\{\text{milk, bread, Diaper}\}) = 2$$

Support:

fraction of transactions that contain an itemset.

$$S(\{\text{milk, bread, Diaper}\}) = 2/5$$

Frequent itemset.

An itemset whose support is  $\geq \text{minSup}$ .

Association rule:

An implication expression of the form  $X \rightarrow Y$ , where  $X$  and  $Y$  are itemsets.

Example

$$\{\text{milk, Diaper}\} \rightarrow \{\text{Beer}\}$$

Rule Evaluation Metrics:

Support ( $S$ ):

$$S = \frac{\sigma(\{\text{milk, Diaper, Beer}\})}{|T|} = \frac{2}{5} = 0.4$$

Confidence ( $C$ )

Measures how often items in  $Y$  appear in transactions that contains  $X$

$$C = \frac{\sigma(\{\text{milk, Diaper, Beer}\})}{\sigma(\{\text{milk, Diaper}\})} = \frac{2}{3} = 0.67$$

## Visualizations of discovered patterns

→ Here we use / apply task relevant data.

Rules, tables, reports, charts, graphs, dts, etc.

clustering  
functions

decision tree:

Drill down and roll ups.

Ex:

- A user may specify a selection on items at

→ All Electronics using concept rule " home entertainment", even though individual items in db might not be stored acc. to type,

→ Rather at lower concepts such as "CD player", "TV", or "VCR", "Radio", "LED", ...

→ A concept hierarchy on item that specifies lower level concepts  
of "TV", "CD player", ..., can be used in collection of task relevant data.

- Sometimes it is required for the user to specify the link e.g.: sales of certain items may be closely linked to particular events such as Halloween or to particular groups of customers. [may not be a part of general data analysis request].

## II. Kind of knowledge to be mined.

data set specified by the user.

Strong semantic ties to enhance initial data set.  
Rank attributes to evaluate.

## I: Task relevant Data

- \* 1<sup>st</sup> primitive is specification of data on which mining is to be performed

- \* Typically user interacts with subset of database.

Reason:

- No. of patterns generated could be exponential w.r.t. db size.

- \* If dm-task is to study associations between items frequently purchased at All Electronics by customers in India Canada, task relevant data can specify

- The name of db or dw (~~All~~ AllElectronics-db)

- names of tables / cubes containing relevant data (customer, purchases, items\_sold).

- conditions for selecting relevant data.

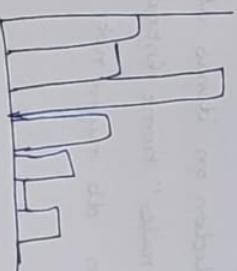
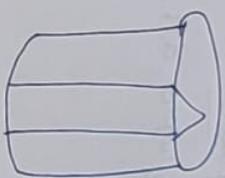
- (e.g. retrieve data pretraining to purchases made in China for current year).

- Relevant attributes/~~dimensions~~ dimensions.

- Income and age from customer table.

## Task relevant data

Db or DW name  
conditions for [deadlock]  
data selection.



knowledge type to be more.

characterisation

clustering

classification

prediction

Discrimination & Association

## Background knowledge

Concept hierarchies

user beliefs about relationships in data

pattern interestingness

measures

certainty

simplicity

## ④ Combined and inspection

→ **Noisy data**: binning  
**Variable**:  
 → Given a numeric attribute such as, say, price, how can we "smooth" out the data to remove the noise.

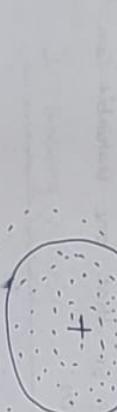
→ **Binning**  
 Data mining: customer-id, purchase, product id, quantity, Date

→ **Outlier analysis**  
 Given a numeric attribute such as, say, price, how can we "smooth" out the data to remove the noise.  
 ⑤ clustering

⑤ clustering



Outlier analysis



⑥

Outliers may be detected by clustering, where similar values are organized into groups ("clusters").

"clusters".

④ Regression

$$y = mx + c$$

Info theoretic measure identifies outlier patterns in handwritten character db: too ideal for classification.

Mislabelled characters

Pattern based on  $\theta$

expressed by human

## ⑤ Binning

Data mining:

customer-id, purchase, product id, quantity, Date

Patterns: frequent buyers

Text mining:

customer feedbacks, reviews, sentiment analysis, ratings, comments

opinion mining patterns

online shopping, customer experience with

Amazon reviews

web mining

sentiment analysis, topic modeling, clickstream, navigation patterns

clickstream, navigation patterns

Social netw relationships, fb, online media traffic Group communities, influential users mining

mining

Trajectory mining

movement patterns

→

targeted marketing

↓

traffic flow

GPS, vehicles, objects moving

forensics purchased traffic flow optimization

↑

1st

forensics purchased traffic flow optimization

↑

2nd

forensics purchased traffic flow optimization

↑

3rd

forensics purchased traffic flow optimization

↑

4th

forensics purchased traffic flow optimization

↑

5th

forensics purchased traffic flow optimization

↑

6th

forensics purchased traffic flow optimization

↑

7th

forensics purchased traffic flow optimization

↑

8th

forensics purchased traffic flow optimization

↑

9th

forensics purchased traffic flow optimization

↑

10th

forensics purchased traffic flow optimization

↑

11th

forensics purchased traffic flow optimization

↑

12th

forensics purchased traffic flow optimization

↑

13th

forensics purchased traffic flow optimization

↑

14th

forensics purchased traffic flow optimization

↑

15th

forensics purchased traffic flow optimization

↑

16th

forensics purchased traffic flow optimization

↑

17th

forensics purchased traffic flow optimization

↑

18th

forensics purchased traffic flow optimization

↑

19th

forensics purchased traffic flow optimization

↑

20th

forensics purchased traffic flow optimization

↑

21st

forensics purchased traffic flow optimization

↑

22nd

forensics purchased traffic flow optimization

↑

23rd

forensics purchased traffic flow optimization

↑

24th

forensics purchased traffic flow optimization

↑

25th

forensics purchased traffic flow optimization

↑

26th

forensics purchased traffic flow optimization

↑

27th

forensics purchased traffic flow optimization

↑

28th

forensics purchased traffic flow optimization

↑

29th

forensics purchased traffic flow optimization

↑

30th

forensics purchased traffic flow optimization

↑

31st

forensics purchased traffic flow optimization

↑

1st

forensics purchased traffic flow optimization

↑

2nd

forensics purchased traffic flow optimization

↑

3rd

forensics purchased traffic flow optimization

↑

4th

forensics purchased traffic flow optimization

↑

5th

forensics purchased traffic flow optimization

↑

6th

forensics purchased traffic flow optimization

↑

7th

forensics purchased traffic flow optimization

↑

8th

forensics purchased traffic flow optimization

↑

9th

forensics purchased traffic flow optimization

↑

10th

forensics purchased traffic flow optimization

↑

11th

forensics purchased traffic flow optimization

↑

12th

forensics purchased traffic flow optimization

↑

13th

forensics purchased traffic flow optimization

↑

14th

forensics purchased traffic flow optimization

↑

15th

forensics purchased traffic flow optimization

↑

16th

forensics purchased traffic flow optimization

↑

17th

forensics purchased traffic flow optimization

↑

18th

forensics purchased traffic flow optimization

↑

19th

forensics purchased traffic flow optimization

↑

20th

forensics purchased traffic flow optimization

↑

21st

forensics purchased traffic flow optimization

↑

22nd

forensics purchased traffic flow optimization

↑

23rd

forensics purchased traffic flow optimization

↑

24th

forensics purchased traffic flow optimization

↑

25th

forensics purchased traffic flow optimization

↑

26th

forensics purchased traffic flow optimization

↑

27th

forensics purchased traffic flow optimization

↑

28th

forensics purchased traffic flow optimization

↑

29th

forensics purchased traffic flow optimization

↑

30th

forensics purchased traffic flow optimization

↑

31st

forensics purchased traffic flow optimization

↑

1st

forensics purchased traffic flow optimization

↑

2nd

forensics purchased traffic flow optimization

↑

3rd

forensics purchased traffic flow optimization

↑

4th

forensics purchased traffic flow optimization

↑

5th

forensics purchased traffic flow optimization

↑

6th

forensics purchased traffic flow optimization

↑

7th

forensics purchased traffic flow optimization

↑

8th

forensics purchased traffic flow optimization

↑

9th

forensics purchased traffic flow optimization

↑

10th

forensics purchased traffic flow optimization

↑

11th

forensics purchased traffic flow optimization

↑

12th

forensics purchased traffic flow optimization

↑

13th

forensics purchased traffic flow optimization

↑

14th

forensics purchased traffic flow optimization

↑

15th

forensics purchased traffic flow optimization

↑

16th

forensics purchased traffic flow optimization

↑

17th

forensics purchased traffic flow optimization

↑

18th

forensics purchased traffic flow optimization

↑

19th

forensics purchased traffic flow optimization

↑

20th

forensics purchased traffic flow optimization

↑

21st

forensics purchased traffic flow optimization

↑

22nd

forensics purchased traffic flow optimization

↑

23rd

forensics purchased traffic flow optimization

↑

24th

forensics purchased traffic flow optimization

↑

25th

forensics purchased traffic flow optimization

↑

26th

forensics purchased traffic flow optimization

↑

27th

forensics purchased traffic flow optimization

↑

28th

forensics purchased traffic flow optimization

Solution  
→ inconsistent data:

functional dependencies.

Noisy data example.

(Pg no 141) (ascending order sorting is important).

1. Age values for data tuples are  
13, 15, 16, 16, 19, 20, 20, 21, 22, 25, 25, 25, 25,  
30, 33, 33, 35, 35, 35, 36, 40, 45, 46, 52, 70.

Step 1: Sort the data.

Step 2: Partition the data into equidepth bins of  
depth 3 ( $\lceil \frac{12}{3} \rceil$ ) since

Bin 1: 13, 15, 16 Step 3: calculate the arithmetic  
Bin 2: 16, 19, 20 mean of each bin

Bin 3: 20, 21, 22

Step 4:  
Replace each of the values in  
each bin by the arithmetic mean  
calculated for the bin.

Bin 4: 22, 25, 25

Bin 5: 25, 25, 30

Bin 6: 33, 33, 35

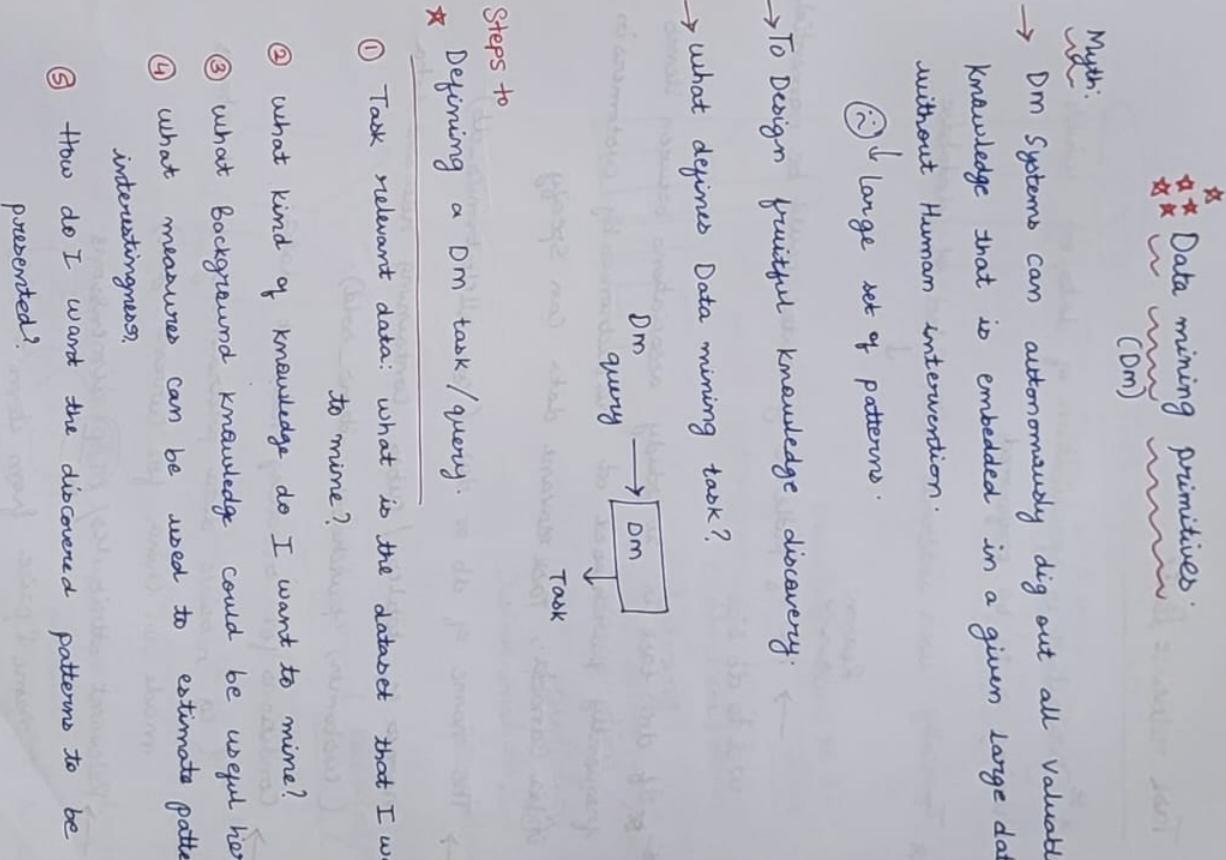
Bin 7: 35, 35, 35

Bin 8: 36, 40, 45

Bin 9: 46, 52, 70

## Data Integration & Transformation:

- ✓ Smoothing : Removes noise from data. ↗ clustering
- ✓ Aggregation : Summary (daily sales data). at multiple granularities.
- ✓ Generalization : Age ↘ Sheet (country location)
- ✓ Normalization →  $-1.0 \text{ to } 1.0$   $0.0 \text{ to } 1.0$  ↗ Regression
- ✓ Attribute construction. ↗ income / age customer database
- ✓ Age, income → new Income-to-Age-Ratio ↗ new predictions ... (29, 29, 34) - group to 20-30 30-40 ↗ without Human intervention.
- Data Reduction : ↗ large set of patterns.  
To obtain a reduced representation of dataset, to ↓ Volume  
yet closely to maintain integrity of original data.



## Exploratory data mining

✓ A process to discover patterns, relationships, and insights in data

✓ No Hypothesis

✓ Statistical measures:

① Identifying patterns

2. Detect anomalies / outliers

3. understands data distribution

4. Generate hypotheses for further predictions

for

## Techniques

① Data visualization

2. Summary statistics, data aggregation

3. correlation analysis & feature selection

## Heterogeneous databases

\* combines systems such as relational, object oriented, hierarchical, XML, spreadsheets, multimedia, file systems.

## Data mining - on what kind of data?

### Relational databases

E-R model

Transactional databases: consists of a file where each record rep. a transaction: trans-id, items purchased in store.

t-id	Item-id
------	---------

Object oriented databases:

each entity is object.

All electronics eg: individuals, customers:

photo (employee).

Object relational databases

handles large data.

Inheritance:

spatial databases:

Geographical, rmp databases, medical images

Raster format:

Pixels maps with other methods

House located near specific point, mountain temp (altitude).

Temporal databases:

vector

Event occurring shows relational data, having time attributes

Time series databases:

Sequence of values over time

periodically i.e. stock exchange.

Text databases: word descriptions, documents, error, www, unstructured.

Multimedia databases: audio, video, image.

Content based retrieval, voice-mail systems.

## \* Data mining functionalities

1. concept class description: characterization & discrimination
  2. Association Analysis
  3. classification and prediction
  4. cluster
  5. outlier
  6. evolution
- \* DATA PREPROCESSING
- [Chapter 2 in detail]

→ use the attribute mean to fill in the missing value  
let us say, average income of All Electronics customers is \$ 28,000. use this value to replace the missing value for income.

→ use the attribute mean for all samples belonging to the same class as given tuple: for Eg: if classifying customers acc. to credit-risk, replace the missing value with average income value for customers in the same credit risk category as that of the given tuple.

→ use the most probable value to fill in the missing value.

Regression, inference based tools like decision tree.

induction

finds "best" line to fit 2 variables so that one variable can be used to detect the other.

class | concept —  $\frac{1}{2}$  printer...

Data characterization: Summarizing the data of the class under study (target class).

feature w/ ↑ in sales of slow producer (car)

discrimination: comparing target class | another set of class.

more unknown

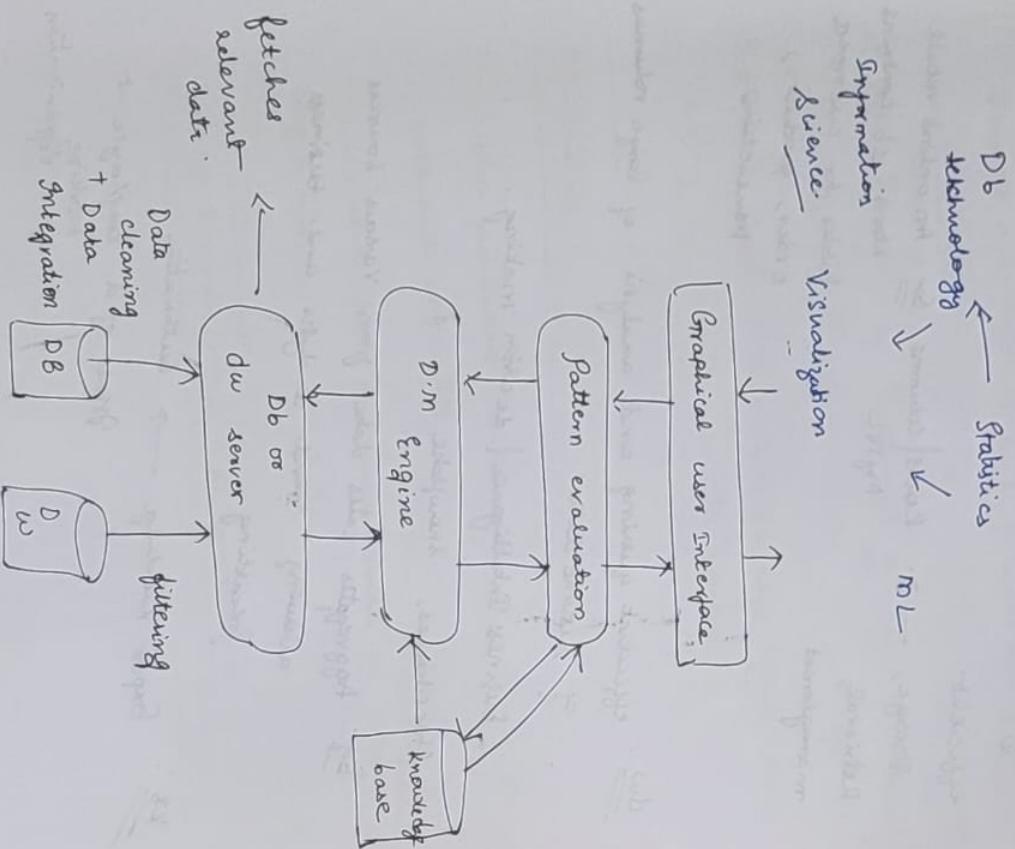
think "unknown as interesting pattern".

value → fill with constant such as "unknown" or  
-∞. Not recommended

mining Task may

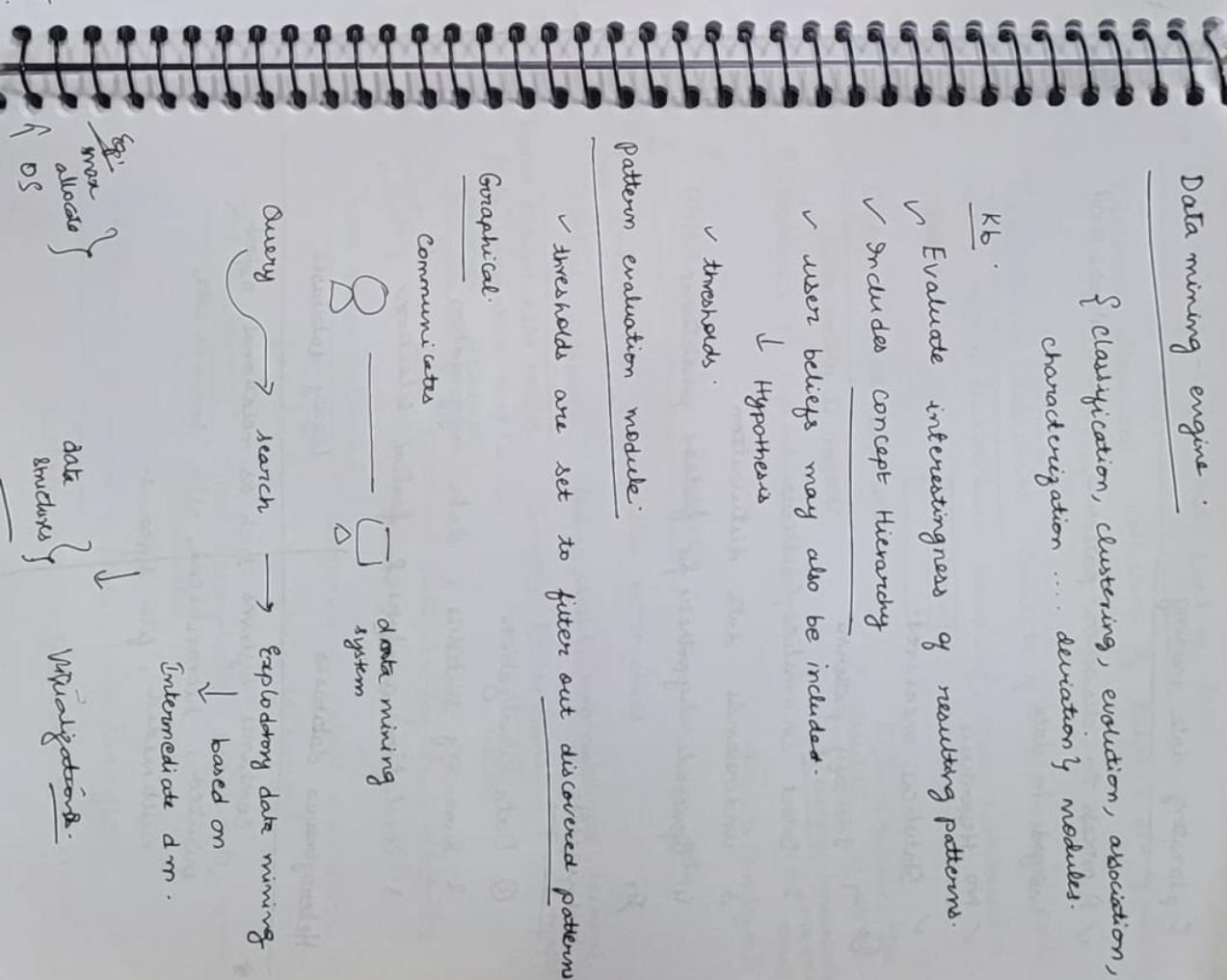
Multiple disciplines of data mining

Data mining engine



{ classification, clustering, evolution, association, characterization ... . . . . . deviation modules.

- ✓ Evaluate interestingness of resulting patterns.
  - ✓ Includes concept hierarchy
  - ✓ user beliefs may also be included.



# Introduction to data mining

Objective: To analyze the performance of students.

Step 1:

Data collection

Various students

How?

- Methods:
- ✓ Surveys & questionnaires: primary
  - ✗ Marks / Grades : secondary
  - ✓ wearable devices: To track sleep / physical activities.



Smart watches.

Challenges:

- Ensuring data privacy & ethical considerations
- Achieving high response rate and accurate self-reporting from students.

fact

Raw in nature is data

Real world data



Raw facts



Data



Dbms → Information

Data mining → knowledge [Patterns]

ML



DL

Intelligence: (solution)

↓  
AT

data pre-processing



data cleaning

Recommendations:-

Automation:-

## Step 5: cognitive science

Human intervention to Robot / AI / model:

Trends & patterns:

Trend:

→ changing over time

e.g. ↑ in average time students spend on online learning platforms over several semesters

Pattern:

Regularity / Repeated occurrence of specific events within a dataset / data.

example: students who participate in study groups tend to achieve higher grades consistently.

Classification problem

prediction problem.

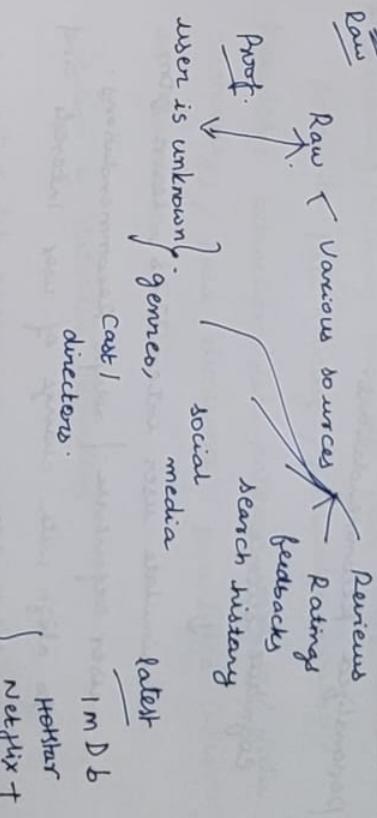
To understand user preferences.  
✓ predict future viewing, interests

✓ recommend similar movies based on past behaviours and ratings (Amazon recommendations).

Movie recommendation systems with AI and cognitive science.

② Movie recommendation system.

Raw Data collection.



③ Data processing: ( & storage) | Data mining: (C. hidden patterns)

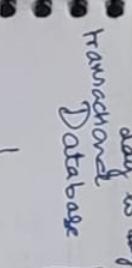
for meaningful data and knowledge extraction).

④ Information:

## Patterns

### Health care

1. Patients with high bp and diabetes are prone to heart disease.
2. Patients receiving treatment  $\hat{=}$  have a  $\uparrow$  recovery rate compared to  $\underline{\text{B}}$ .
3. Patients readmitted within 30 days often have a history of multiple chronic conditions.



historical processing  
data warehouse

data warehouse  
analytical processing

knowledge base.

efficient

storage, → rows | columns

for: An online retail store db contains

Retrieval,  
management

tables for customers,  
order, product,  
transactions.

dw: efficient querying and analysis of large volumes of historic data.

Business Intelligence | decision making

Star schemes, Snowflake

Aggregates sales data from various sources allowing trend analysis and business forecasting.

Amazon Example → recommendation analysis

K3: Expert knowledge → automation

Issues of a Employee at working organization.