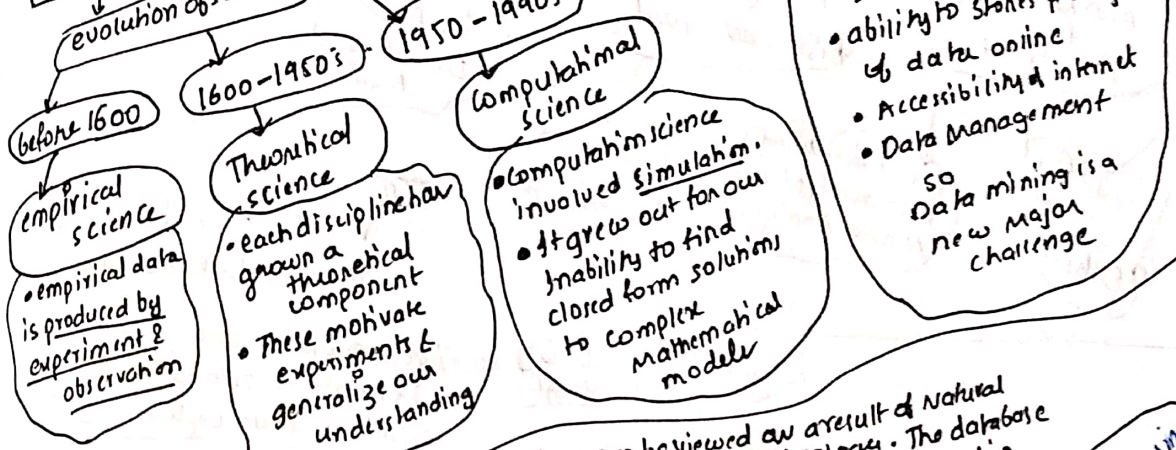


Data Mining - UNIT - I

Motivation



Data mining as an evolution of Information Technology

Data mining can be viewed as a result of natural evolution of Information technology. The database and data management industry evolved in development of several critical functionalities

1960's era

- Data collection, database creation, (IMS) Information management system and network DBMS
- primitive file processing

1970's

- DBMS, relational DBMS
- Hierarchical & network database systems
- Relational DBMS
- SQL
- Transaction control, concurrency & recovery
- OLTP (Online Transaction processing)

1980's

- Advanced data models
- Advanced DBMS: Spatial, temporal, etc.
- Advanced Analysis
- OLAP (Online analytical processing) & data warehousing & data mining

Data warehousing

- Subject oriented & integrated
- non volatile
- Time variant

1980's to 2000

- Data mining & data warehousing
- Web databases
- XML-based data systems

Present - Future:

- New generation of integrated data & information systems

Data mining

(Knowledge extraction in database)

"extraction or mining interesting (implicit, previously unknown & useful) information / knowledge / patterns from large amounts of databases"

Alternate Name

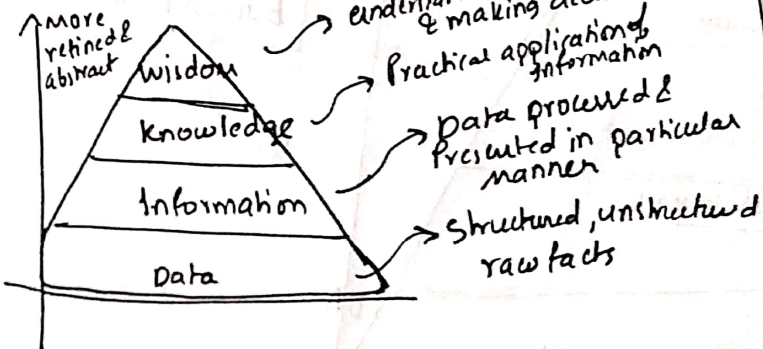
- Knowledge discovery
- Knowledge extraction
- data Analysis
- Information Harvesting

What is not data mining

- Deductive query processing
- Expert systems
- Statistical programs

Information Hierarchy / DIKW pyramid

refers to a class of models for representing the structural or functional relationships between data, information, knowledge & wisdom



more refined & abstract

Wisdom

Knowledge

Information

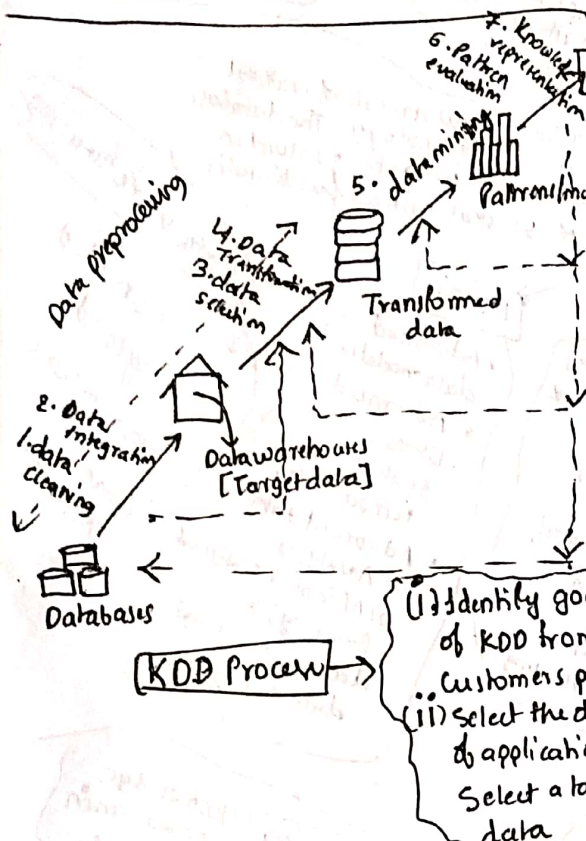
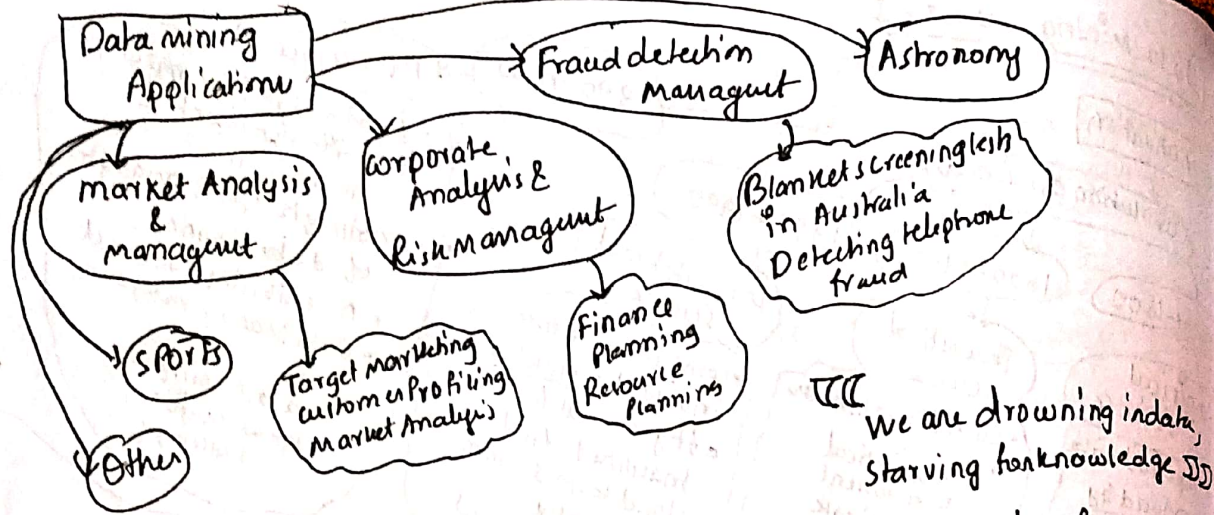
Data

understanding of why & making decisions

Practical applications of information

Data processed & presented in particular manner

Structured, unstructured raw facts



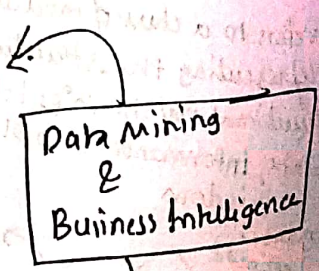
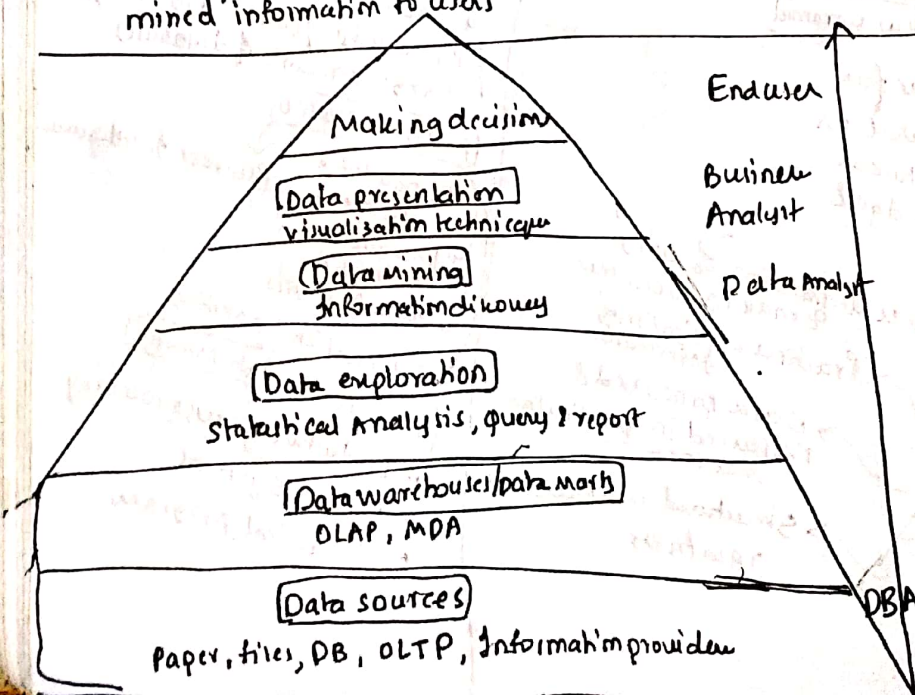
Explosive growth of data from Terabytes to zettabytes is typically a data explosion

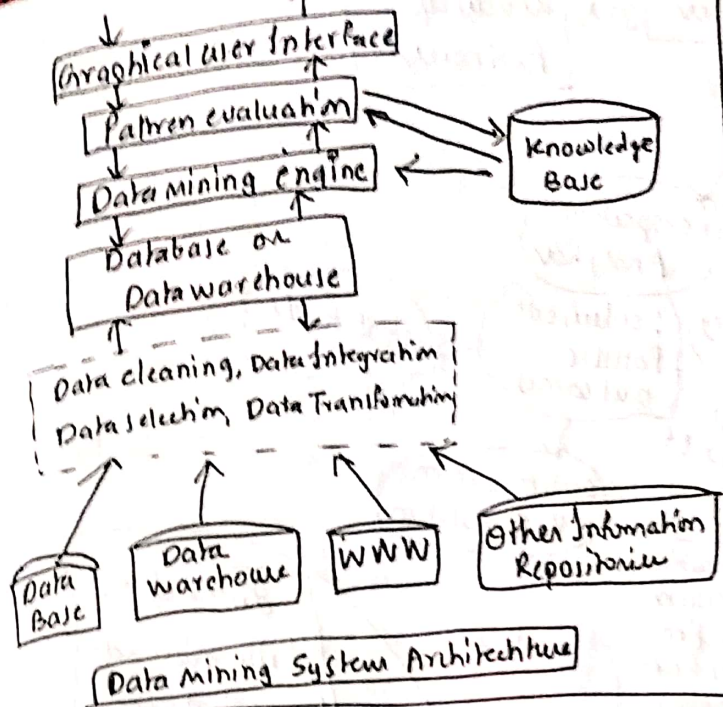
1. Data cleaning: - To remove noise & inconsistent data.
 2. Data Integration: - where multiple data sources may be combined.
 3. Data selection: - where data relevant to analysis tasks are retrieved from databases.
 4. Data Transformation: - where the data are transformed and consolidated into new and common formats which are appropriate for mining.
- *Data Pre processing: steps 1 to 4
5. Data mining: - an essential process where intelligent methods are applied to extract data patterns or desired results.

7. Knowledge representation: - where visualization & knowledge representation techniques are used to present mined information to users

"data mining step may interact with a user on a database".

6. Pattern evaluation: - To identify truly interesting patterns and representing knowledge based on interesting measures.





- Knowledge base is domain knowledge that is used to guide the search or evaluate the interestingness of a resulting pattern
- Data mining engine consists of functional modules
 - characterization
 - Association & Co-relation Analysis
 - Prediction etc.



Data mining Functionalities

(or) Kind of Patterns

What data is telling!

are three types

Descriptive Analytics

: Insight into Past

Business Intelligence

- Cluster Analysis
- Outlier Analysis
- Evolution Analysis
- Mining frequent Patterns
- Associations
- Co-relationships
- Characterization & Discrimination

Predictive Analytics

: understanding the future

Forecasting

- Classification
- Regression
- Prediction
- Time Series Analysis

Perspective Analytics

: advise on possible outcomes

Optimization & Simulation

What to do?

Some of the data mining functionalities

1. Characterization & description
2. Mining of frequent patterns
3. Associations & Co-relationships
4. Classification & regression
5. Clustering analysis
6. Outlier Analysis and Evolution Analysis

Characterization & discrimination

↳ description of a class or a concept is called class / concept description

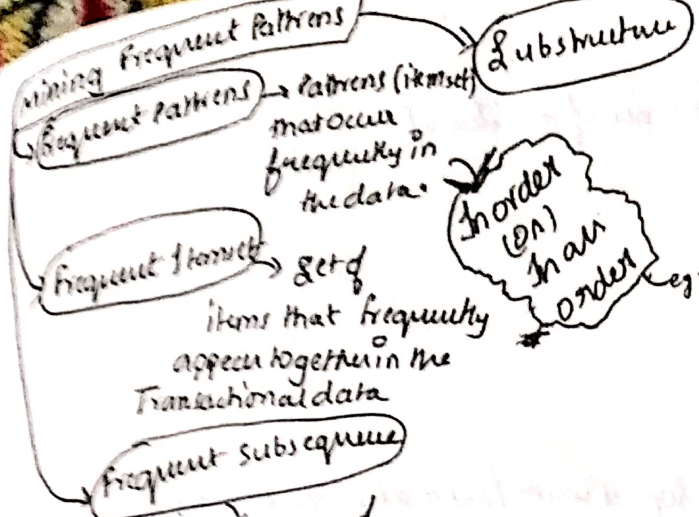
These descriptions can be derived from

(i) Data characterization: is a summarization of general characteristics / features of a target class of data
output forms: pie chart, bar chart

(ii) Data discrimination: is a comparison of general features of the target class data objects with the general features of objects from one or a set of contrasting classes.
output forms: Same as above
 Discrimination descriptions expressed in rule form are referred to as discriminant rules

What technologies are used





• A frequent pattern mining search is recovering relationships in a given dataset
 • A frequent mining of patterns leads to a discovery of Association or - correlation items between the items in the dataset

Association rule is a data mining Algorithm discovered by IBM company

Association rule identifies the association between two items based on their occurrence

Itemset

A collection of one or more items
 eg:- milk, bread & diaper.

eg:- In a given set of transactions, predict the transaction or find the rules that will predict the occurrence of item based on the occurrence of another item.

Mostly used in market-basket transactions.

Market - Basket Analysis

eg:-

TID	itemset
1	Bread, milk
2	Bread, diaper, beer, coke
3	Milk, diaper, beer, coke
4	Rc car, battery

k-item set

An itemset which contains k items.

{milk} → {bread}
 {milk, bread} → {eggs, coke}

Frequent itemset

An itemset whose support \geq threshold \geq min sup

Any itemset whose support value \geq (min support) min sup which is called as frequent

Association rules

Support Count

A support count is represented using σ , frequency of an occurrence of an itemset.

Support & Support count are Synonymous

σ (milk, bread, diaper) = 2

A fraction of transactions that contain an itemset is called as support
 eg: {milk, bread, diaper} = $\frac{2}{5}$

22-12-18

③ Transactions
eg:- Booking, buying

ID	Item
1	TV, mobile
1	1

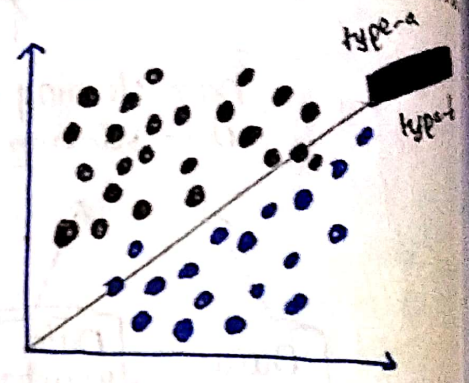
- Types of databases
- Determining functionality
 - descriptive
 - Predictive

- ④ Other kinds of data
- Time-related/sequence data
 - data stream (video surveillance)
 - Spatial data (maps)
 - hypertext data etc

Classification and prediction

- A classification is used to find a model or a function which describes/distinguish between class & concept.
- Different approaches of classification

eg:- classification of types of butterfly



- → data points of type-a based
- → data points of type-b based

① If-then-else :-
 → Based on the features
 If (wing color = red)
 classify as (type-a);

② Decision tree :-
 → using a decision tree and an algorithm to classify.



- Based on the training data, we extract the features
- These training data contains labeled examples of which type of butterfly
- This is given to the model to train the model.

③ Mathematical formula :-
 → using mathematical formula & relations we can classify.

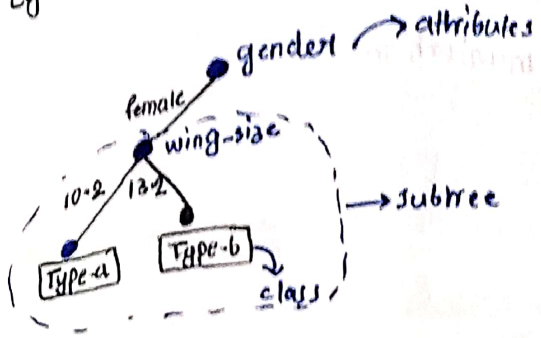
eg:-
 Age(x, youth) & income(x, 10000) → (x, class A)
 Age(y, mid-age) & income(x, 10000) → (x, class B)

"The above problem is having a linearly separable solution, so it is a linearly separable problem"

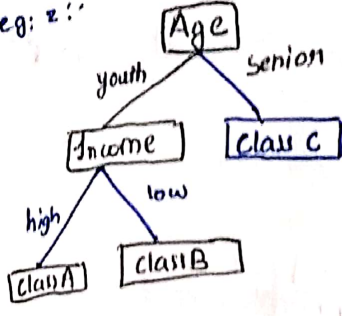
④ Neural Networks.

⇒ Decision Tree :- A decision tree may look like a flow chart or a tree diagram with each node containing an attribute value which is a test attribute value. And every sub-tree represents outcome of test value. And leaf nodes represent the classes (labels).

eg:- decision tree for type of butterflies



eg: 2:-



⇒ Prediction:-

Whereas the prediction takes the continuous data and predict the future values.

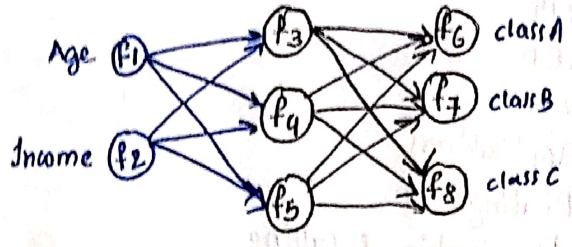
- prediction is a continuous valued function.
- This can be used to predict unavailable, missing numerical values.

⇒ clustering analysis:-

- Data objects which are not known class labels or unknown class labels or without knowing class labels.
- class labels are not present in the Training data, for that classes generate the clusters.
- Objects are clustered or grouped based on the similarities.
- Maximizing intra class similarity and minimizing the inter-class similarity.

⇒ Neural Networks :-

- It contains a collection of Neurons
- And they contains weights and connection between units

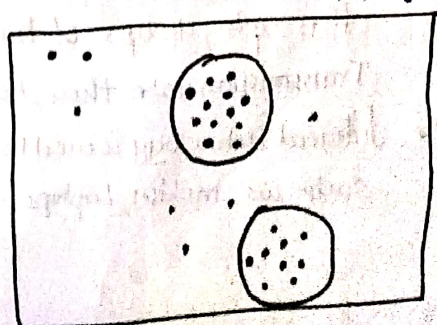
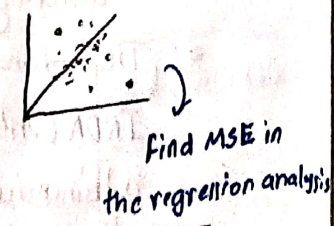


Classification is performed only on

- Numerical dataset
- categorical data (eg: male or female)
- ordered events (eg: historical events)
- Discrete data (sequences)

⇒ Regression analysis :-

- A Regression Analysis is a Analysis in which we find the error rate using statistical methods.
- Regression analysis contains statistical methods which is used for Numerical Prediction
- Identification of trends on the available data
- Relevance analysis is done before classification & prediction to identify the attributes.

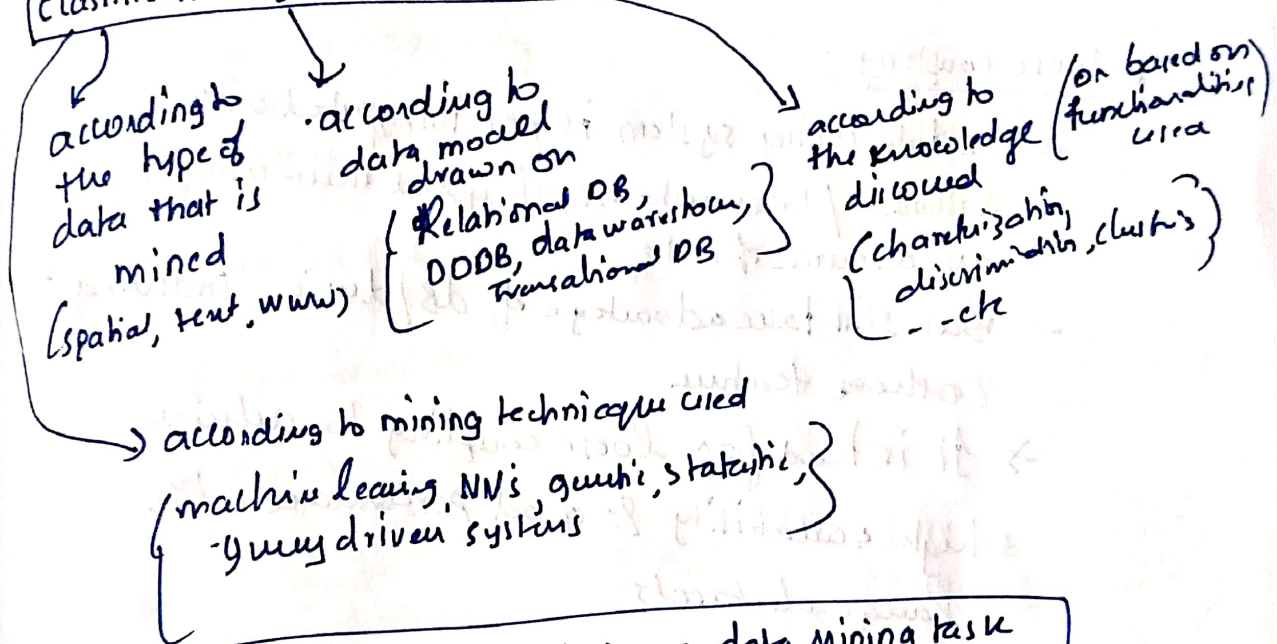


Are all Patterns Interesting?

1. It is easily understood by humans
2. valid on new or test data with some degree of certainty
3. Potentially useful
4. novel.
5. If it validates a hypothesis

There are subjective vs objective interesting measures:

Classification of data mining systems



Five task primitives of specifying a data mining task

1. Task relevant data { relevant (attributes, dimensions) }
2. Knowledge to be mined { specifying which data mining functionality to be used }
3. Background knowledge { Domain Knowledge }
4. Pattern interestingness measure { interestingness measures for certainty, simplicity, utility & novelty }
5. Visualization of discovered patterns { rules, tables, pie, bar chart, decision trees }

Integration of Data Mining system with a database or a Data warehouse system

— difference between OLAP & OLTP

• NO coupling :

→ data mining system sources such as flat files to obtain data for mining.

→ But no data ~~mining~~ functions are implemented in base

Process

→ This is a poor design choice.

• Loose coupling :

→ The data mining system is not integrated with database/data warehouse beyond their usage as a source of data

→ But still take advantage of DB/DW's indexing & other features

→ It is hard for loose coupling to achieve high scalability & good performance with large datasets

• Semi-tight coupling :

• Some of the primitive operations such as aggregation, sorting & pre-computation of statistical functions can be done within the database & during query. And these can be stored inside DB/DW's itself to promote high performance of data mining system

• Tight coupling :

→ Complete integration of DB/DW into Data mining System

→ high scalability & performance

— neglecting Technical & Implementation details. It is the best Architecture

Major Issues in data mining

Mining Methodology & user-interaction issues

- Interactive mining of knowledge at multiple levels of abstraction
- Incorporation of background knowledge
- Data mining query language & ad-hoc data mining
- Presentation & visualization of mining results
- Handling outlier or incomplete data
- Pattern evaluation

Performance issues

- Efficiency & Scalability of data mining algorithms
- Parallel, distributed & incremental algorithms

Issues related to diversity of database types

- Handling relational & complex data types
- Mining information from heterogeneous database

Types of Data sets

Record

- relational
- matrix
- Transactional data

graph & Network

- WWW
- social network
- molecular structure

Ordered

- Video data
- Temporal data
- sequential data

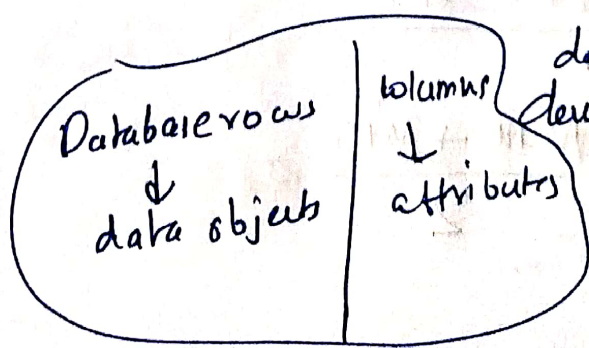
Spatial image & multimedia

Important Characteristics of structured data

- dimensionality
- sparsity
- Resolution
- distribution
 - centrality & dispersion

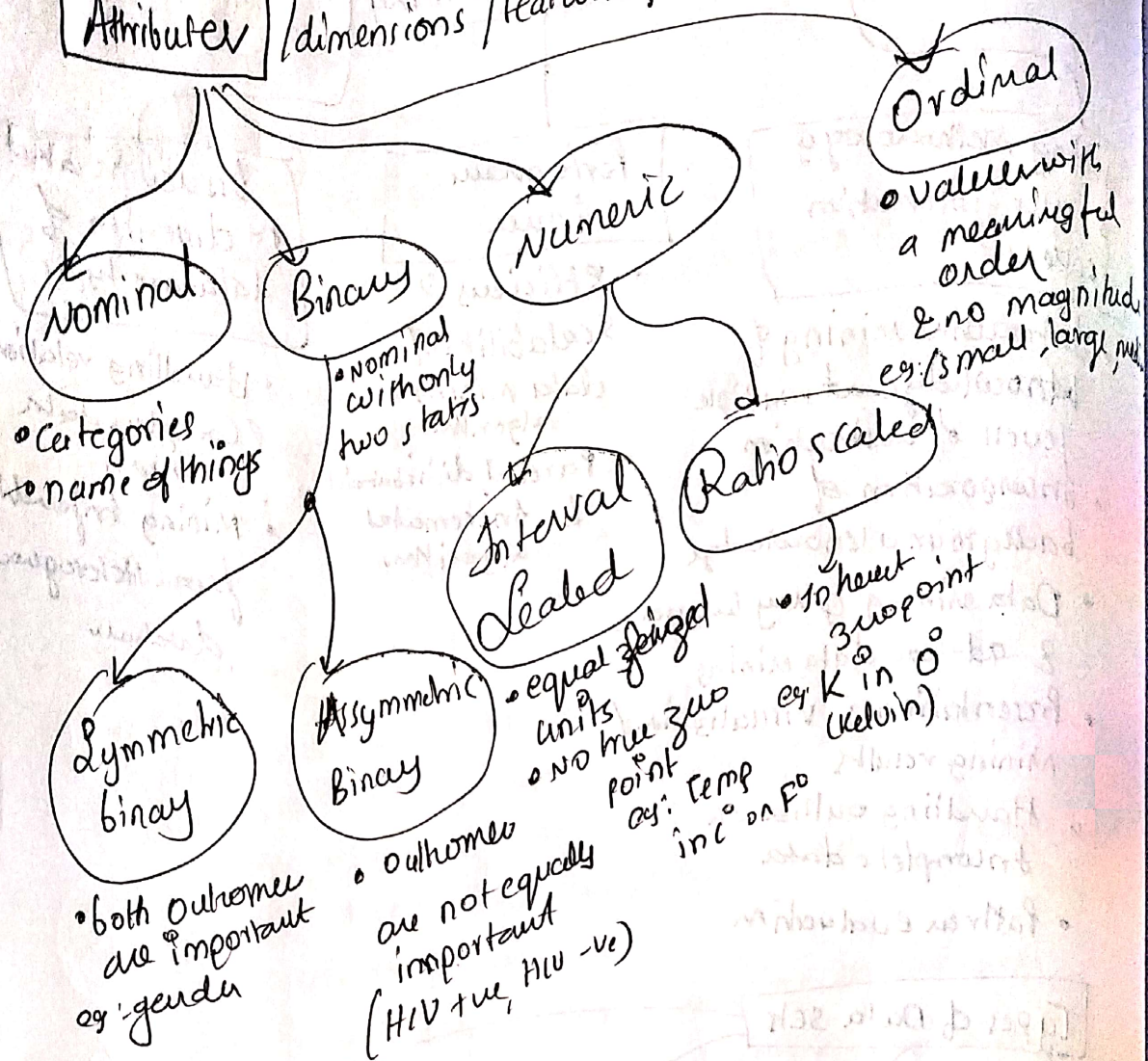
Data object

represents an entity



data objects are described by attributes

Attributes / dimensions / features / variables



Ordinal
 • values with a meaningful order
 & no magnitude
 eg: small, large, tall

Discrete vs continuous Attributes

Discrete

finite or countable set of values
 eg: ZIP codes, Binary attributes

Continuous attributes

have real finite numbers represented using floating point variables

Formulas

$$\text{Mean} = \frac{x_1 + x_2 + \dots + x_n}{n}$$

$$\text{midrange} = \frac{\text{Max} + \text{Min}}{2}$$

Median = $\frac{1(n+1)}{2}$ when n = number of data values
 (If n is odd)

Other wise

If n = even

Median = $\frac{1}{2} \left[\frac{n}{2} + \frac{n}{2} + 1 \right]$

Mode: most repeated / most frequent value

- If two values occur frequently Bimodal
- If three " " " Trimodal
- If more " " " Multimodal

☺☺ The mean, median mode are collectively called as central tendency

Measures of dispersion

← Note

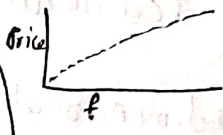
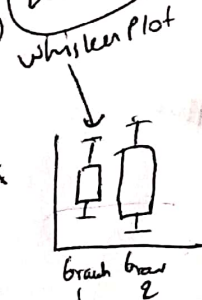
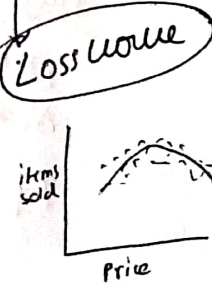
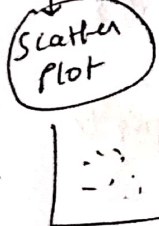
: A Box plot used to represent range, median, quartiles & Inner quartile range (IQR) $[Q_3 - Q_1]$

: Five number summary is

- (Min), Q_1 , (Median, Q_2), Q_3 , (Max)

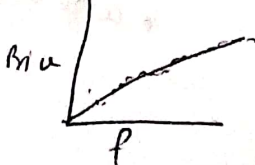
: Outliers either fall below $Q_1 - 1.5 \times IQR$ or above $Q_3 + 1.5 \times IQR$

Graphic display of basic descriptive statistics



$f = \frac{i - 0.5}{n}$

Quantile - Quantile plot



two combi -ve combi

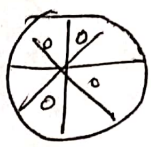
Data visualization Techniques

Pixel oriented

- data of n dimensions create n windows on screen, one for each dimension
 - values in dimensions are mapped to pixels
 - colors represent the corresponding value
- laying out pixels in circles

Geometric Projection visualization Techniques

- direct visualizations
- scatter plot
- Landscapes
- hyperslice
- parallel coordinates
- Projection Pursuit Techniques



Icon Based Visualization Techniques

- Chernoff faces
- stick figures
- tile bars
- shape coding
- color icons

Hierarchical Visualization Techniques

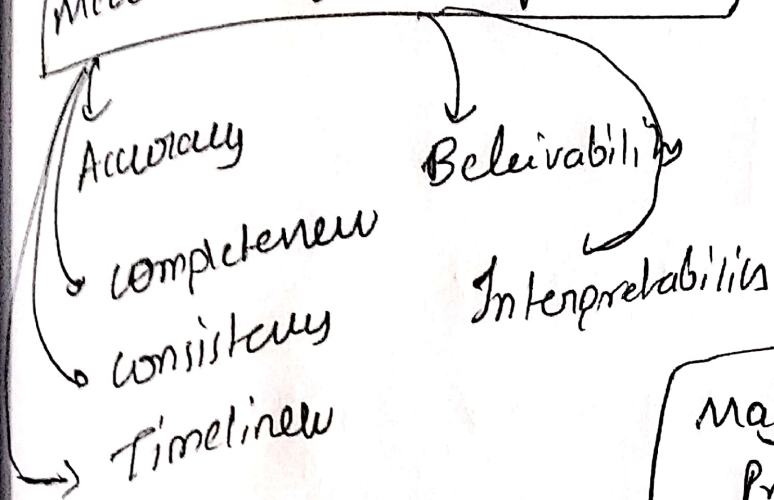
- dimensional scaling
- worlds within worlds
- tree maps
- cone trees

visualizing complex data & relations

- Tags, social networks, News, Typography

Measures of Similarity & Dissimilarity

Measures of Data Quality



Major Tasks in Data Preprocessing

- Data cleaning
 - fill in missing values,
 - smooth noisy data
 - identify or remove outliers
 - resolve inconsistencies
- Data Integration
 - integration of multiple databases, data cubes or files
- Data Reduction
 - dimensionality reduction
 - numerosity reduction
 - data compression
- Data Transformation & data discretization
 - normalization
 - concept hierarchy generation