

UNIT-1

Data mining :-

extracting
of data.

Data mining refers to knowledge from large amount

Many other terms such as knowledge mining from data, knowledge extraction, pattern analysis, data archaeology, knowledge discovery in database (KDD) etc

Definition:-

"Data mining is the process of discovering interesting knowledge from large amount of data stored either in databases, data warehouse or any other information repository."

★ KDD process (Knowledge discovery from Data)

Knowledge discovery consist of an iterative sequence of the following steps:

1. Data cleaning -

To remove noise and inconsistent value (data).

2. Data integration -

where multiple data sources may be combined.

3. Data selection -

These data relevant to the analysis task are retrieved from the database.

4. Data transformation -

These data are transformed into form appropriate for mining. by ~~process~~

5. Data mining -

An essential process where intelligent methods are applied in order to extract data patterns.

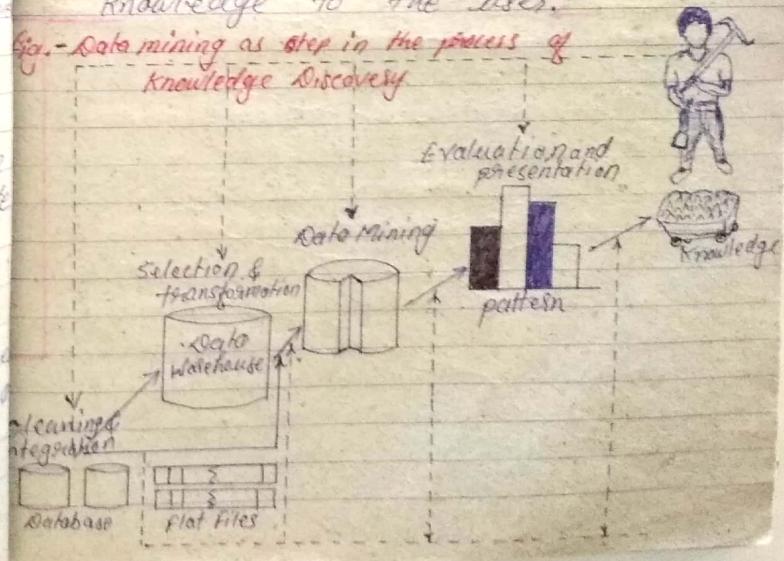
6. Pattern evaluation -

To identify the interesting patterns representing knowledge base on some interestingness measures.

7. Knowledge presentation -

Where visualization and knowledge representation technique are used to present the mined knowledge to the user.

Giz. - Data mining as step in the process of Knowledge Discovery



Data Collection and Database Creation
 (1960s and earlier)
 • Primitive file processing.

Database Management Systems
(1970s - early 1980s)
• Hierarchical and network database systems
• Relational database systems
• Data modeling tools: E-R model, etc.
• Indexing and accessing methods:
• B-tree, hashing etc.
• Query languages: SQL, etc.
• Transaction processing and query optimization.
• Online transaction processing (OLTP)

Evaluation of database System Technology

PAGE NO. 42
 DATE 31-07-2013

Web-based database
(1990s - present)
• XML-based database system
• Integration with information retrieval
• Data and information integration

Advanced Database Systems
(mid - 1980s - present)
• Advanced data models:
• Extended relational, object-oriented etc.
• Advance applications:
• Spatial, temporal, scientific, engineering, knowledge-based.
• Data warehousing and OLAP
• Data mining and knowledge discovery:
• Generalization, classification, association, clustering, frequent pattern, analysis, outliers, analysis, etc.
• Advance data mining applications:
• Stream data mining, web mining, text mining, web mining and society

New generation of integrated data and information System - (present - future)

PAGE NO. 43
 DATE 1-8-2013
 vintech.ans@gmail.com

Architecture of data mining System:

The architecture of a typical data mining system may have the following major components -

1. Database, datawarehouse, warehouse or other information repository.

2. Database or data warehouse server

3. Knowledge base -

This is the domain knowledge part it is used to guide the search or evaluate the results interestingness of resulting patterns.

4. Data mining engine.

It is essential to the data mining system and usually consists of set of function modules such as character, entity association, classification, clustering and outliers analysis etc.

5. Pattern Evaluation Module

This component typically or employs interestingness measure and interact with the data mining module.

6. User Interface -

This module communication between users and data mining system allowing the users to interact with the system by specifying a data mining query or task.

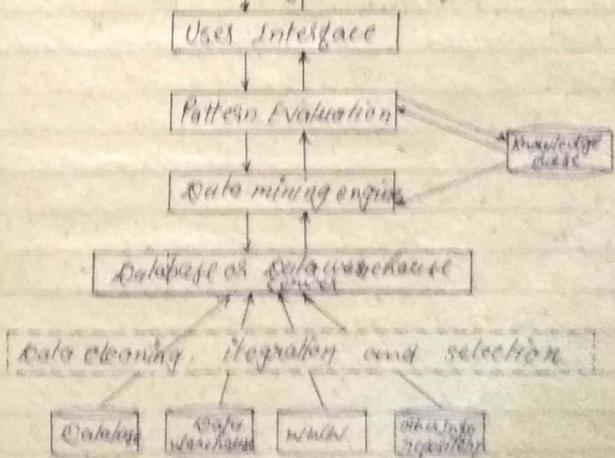


Fig. Architecture of a typical data mining system

Data type for Data mining:-

Data mining - on what kind of data?

These are number of different data stores on which mining can be performed. Data mining should be applicable to any kind of information repository.

- It include -
- ① Relational database
 - ② Data warehouses
 - ③ Transactional Database
 - ④ Advanced data analysis systems or
 - ⑤ Advanced Applications
 - Object recognition
 - spatial
 - Meta Temporal
 - Auto streams
 - WebDB
 - Heterogeneous DB

Relational Databases -

A database system also called RDBMS, consist of a collection of interrelated data known as database and a set of software program to manage and access the data.

A relational database is a collection of tables (relations), each of which assign a unique name, each table contains a set of attribute and large set of tuples. Each tuple in a relational table represent an object identified by a unique key and describable by a set of attribute value. A data model such as ER model is often constructed for relational databases. Example - All electronics company is described by the following relation:

customer, item, Branch, Employee

Data warehouses:-

To facilitate decision making, the data in the a data warehouse are organized around major subjects. It is usually modeled by a multidimensional database structure where each dimension correspond to one or a set of attributes in the schema.

Ex.

Suppose that all electronic a successful international company with branches around the world. Each branch has its own set of data. The president of all electronics is asked to provide an analysis of company's sales per item type, per branch for the 3rd quarter. This is a difficult task since relevant data spread out over several databases that are geographically located at the different site.

If all electronic had a data warehouse this task would be easy. A data warehouse is a repository of information collected from multiple sources stored under a specific schema and that is usually created at a single site.

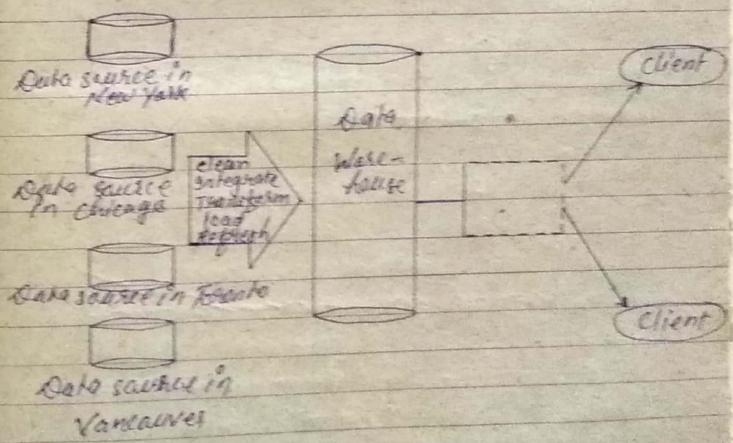


Fig.: Frame work of data warehouse for all electronics.

PAGE NO. 50
DATE 1-1-20

Transactional Databases :-

A transactional database consists of a file where each record represents a transaction. A transaction includes a unique transaction id and a list of item making up the transaction such as item purchased in a store.

Trans.id	List of Item ids
T100	I ₁ , I ₂ , I ₃
T200	I ₄ , I ₅ , I ₆
:	:

Fig - Fragment of transaction database

The transactional database may have additional table associated with it, which contain other information regarding the sale, such as the date of transaction, the customer.id, the id number of the sales person of

PAGE NO. 51
DATE 1-1-20

amitabh.ouo1@gmail.com

the branch and so on.

Advanced data and Information System and Advanced Application:-

Relational database have been widely used in business applications. With the progress of database technology, various kinds of advanced applications have emerged to fulfil the requirement of users. The new database application include -

↳ Spatial Database (maps) -

It contains spatial related information. It include geographic databases (maps), very large scale integrating (VLIS) of computer aided design databases, satellite image databases etc.

Spatial data may be represented "Raster format" consisting of n-dimensional bit-map or pixels maps.

↳ Temporal Database:-

(Sequence databases
Time Series Databases)

The temporal database stores data that include time related attr. These attribute may involve certain stamps.

It's :-

- A sequence database has sequences of ordered events. for ex. include customer shopping sequences, web click streams and biological sequences.

- Time series DB stores seqs of values or event obtained over repeated measurement of time (say as hourly, daily, weekly) it includes data collected from the stock market, inventory control and the observation of natural phenomena (like temperature and wind):

Object Relational Database:-

object relational databases constructed based on object oriented data model.

Each object has associated which have the following:-

1. set of variable that describe object.
2. set of message that the object use to communicate.
3. set of method wh. where each method will the code.

Class 8 "Data mining Task Primitive"

Each user will have a data mining task mind that is some form of data analysis that he or she would like to have performed.

A data mining task can be specified in the form of data mining queries, which are input to the mining system. A data mining query is defined in terms of data mining task primitives. These primitives allow the user to interactively communicate with the data mining system.

The data mining primitives specify the following:-

1. The set of task relevant data to be mined:

This specifies the portion of the database or the set of data in which the user is interested. It includes database attribute or the dimensions of data warehouse.

2. The kind of knowledge to be mined:

This specifies the data mining functions to be performed such as - characterization, discrimination, association or correlation analysis, classification, prediction, clustering, outliers analysis.

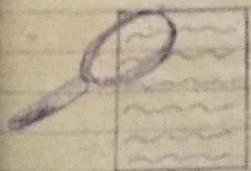
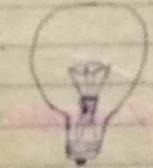
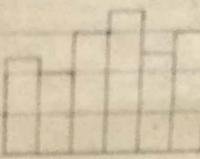
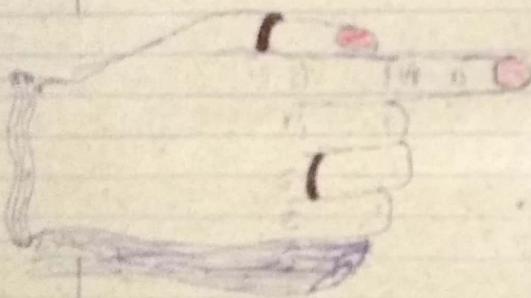
3. The background knowledge to be used in the discovery process:-

It includes concept hierarchy which is a popular form of background knowledge. It allows data to be mined at multiple level of abstraction.

4. The interestingness measure and threshold for pattern evaluation:-

It is used to guide the mining process to evaluate the discovered patterns. Different kind of knowledge may have different interestingness measure. For example - Interestingness measure for association rule may include confidence and support.

5. Visualization of discovered pattern
 The expected representation for visualizing the discovered patterns.
 Include various discovered patterns like table, chart, cubes etc.



- Task relevant data
 - o Database or data warehouse name
 - o Database table or data warehouse cubes
 - o conditions for data selection
 - o Data grouping criteria

- Knowledge type to be mined
 - o Characterization
 - o Discrimination
 - o Association / Correlation
 - o Classification / Prediction

- Background knowledge
 - o Concept hierarchy

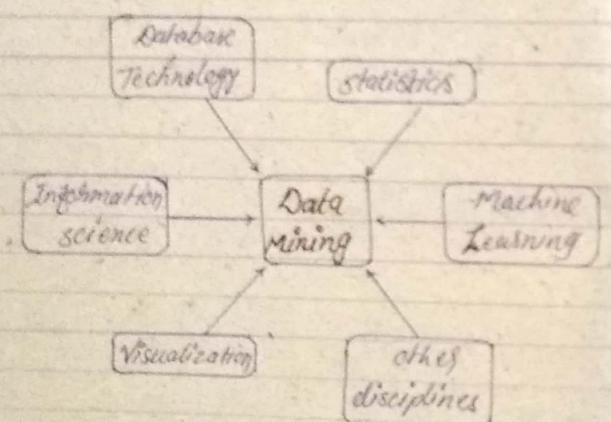
- Pattern interestingness measures
 - o Simplicity
 - o Certainty (Confidence)
 - o Utility (support)
 - o Novelty

- Visualization of discovered patterns
 - o Rules, tables, charts, graphs,
 - o decision trees; and cubes
 - o Drill down and roll up

Fig. Primitives for specifying a data mining task

classification of Data mining Systems

Data mining system can be categorized according to various criteria as follows -



ff. Data mining as confluence of multiple disciplines

1. Classification according to the kind of database used
2. Classification according to kind of knowledge
3. Classification according to the kind of techniques utilized.
4. Classification according to the application area

Data mining functionality :- (Kind of pattern can be mined)

DM functionality are used to specify the kind of patterns to be found in DM task. DM task can be categorized into two categories -

1. Descriptive -

Descriptive mining task characterize the general properties of data in the database.

2. Predictive -

Predictive mining task perform inference on the current data in order to make prediction.

DM functionality :-

* Concept / Class Description -

(Characterization & Discrimination)

Data can be associated with classes or concept. It can be useful to describe

PAGE NO. 61
DATE: 17/10/2022

individual classes/concept in summary concise and precise terms. Such description of a class or concept are called concept/class describing and can be derived via -

- 1. Data characterization (target class)
- 2. Data discrimination (Contrasting)
- 3. Both characterization & Discrimination

- i - Data characterization :- (target class)

Data characterization is a summarization of the general characteristics of a target class of data. The data corresponding to the user specified are collected by a database query.

There are several methods for effective data summarization and characterization like - data summarization based statistical measure, the data based OLAP roll-up operation can used to perform data summarization along a specified dimension.

PAGE NO. 62
DATE: 17/10/2022
vinith.ca07@gmail.com

An attribute oriented induction technique can be used to perform characterization.

The output of data characterization can be presented in various forms like pie chart, bar chart, curves, multidimensional data cube, multidimensional tables etc

Example :-

A data mining system should be able to produce a description summarizing the characteristic of customers who spend more than \$10000 a year at all electronics.

The result could be the general profile of the customers such as - their age, employeed and have excellent credit rating

- Data Discrimination - (contrasting class)

Data discrimination is a comparison of the general feature of the target class data object general feature of the object from one a set of contrasting classes. The target class and contrasting class can be specified by the user or the corresponding data object retrieved through database queries.

Example -

A data mining system should be able to compare two groups of "all electronics" customer such as those who shop for computer product regularly (more than twice in a month), vs those who rarely shop for computer products (less than 3 times in a year). The resulting description provides a general comparative profile of the customers such as 80% of the cu-

who frequently purchase computer product are between 20-40 years old and have a university education whereas 60% of the customers who infrequently buy such product are either senior or youth and have no university degree

★ Mining frequent patterns. Association, & correlations

Frequent patterns, are patterns that occur frequently in data. There are many kinds of frequent patterns like - itemsets,

1. Frequent itemsets -

It refers to the set of items that frequently appear together in a transactional data set such as milk, bread, butter

2. Subsequence -

A frequently occurring subsequence such as the pattern that

PAGE NO: 65
DATE: 1/1/20

customers tend to purchase i.e. A person first purchase a P.C. then Anti-virus and followed by a digital camera and then memory card.

3. Sub-structure:-

A substructure can refer to different structural forms such as graph tree or lattice which may be combined with items or subsequences.

Mining Association and Correlation

~~mining~~ Association rule mining finds interesting association or correlation relationship among a large set of data item.

for example - Suppose us marketing manager of 'all electronic' you would like to determine who items are frequently purchase together in the same transaction.

PAGE NO: 66
DATE: 1/1/20
rimlesh.ouer@gmail.com

$\text{buys}(x, \text{"computers"}) = \text{buys}(x, \text{"software"})$

[Support = 1%]
Confidence = 50%

where, x is a customer. A confidence of certainty of 50% means that if a customer buys a computer there is a 50% chance that he/she will buy software as well.

A 1% support means that one 1% of all the transaction under analysis showed that computers and software were purchased together.

This association rules involve a single attribute or predicate that is 'buys'. Association rule that contain a single predicate are referred to as single dimension association rule.

The another association rule which are related to purchase of 'all electronics'.

age > 20 & income > 100000
buys (a. "CD player")
[support = 2%
confidence = 60%]

That is an association between one or more attributes i.e. age, income, and buys. So the above rule can be referred to as a multicriteria association rule.

* classification and Prediction

Classification is the process of finding a model (or function) that describe and distinguish data classes or concept, for the purpose of being able to use the model to predict the class of object whose class label is unknown. This derived model is based on the analysis of a set of training data (i.e. data object whose class label is known). Rules.

* How is the derived model presented?

The derived model may be presented in various form such as - classification rules (if...then), Decision tree, neural networks etc.

Example:- (classification & prediction)

Suppose, as a sales manager of 'All Electronics', you would like to classify a large set of items in the store based on three kinds of response to a sales campaign, good response, mild response and no response.

We would like to derive a model for each of these three classes based on the descriptive features of the item such as price, brand, place mode, type and category.

PAGE NO. 69
DATE: 1/20

The resulting classification should maximize distinguish each class from the others, presenting an organized picture of the data set.

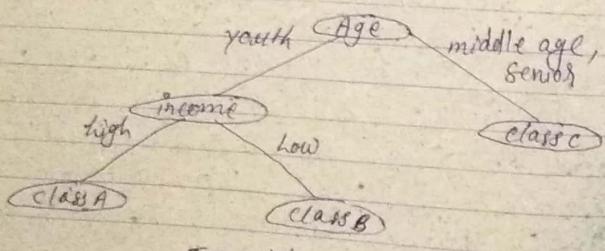
$\text{age}(x, \text{"youth"}) \text{ and } \text{income}(x, \text{"high"}) \rightarrow \text{class}(x, \text{"A"})$

$\text{age}(x, \text{"youth"}) \text{ and } \text{income}(x, \text{"low"}) \rightarrow \text{class}(x, \text{"B"})$

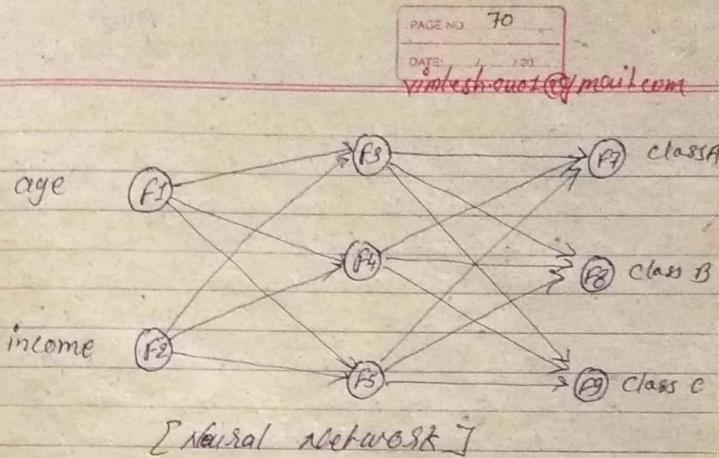
$\text{age}(x, \text{"middle age"}) \rightarrow \text{class}(x, \text{"C"})$

$\text{age}(x, \text{"senior"}) \rightarrow \text{class}(x, \text{"C"})$

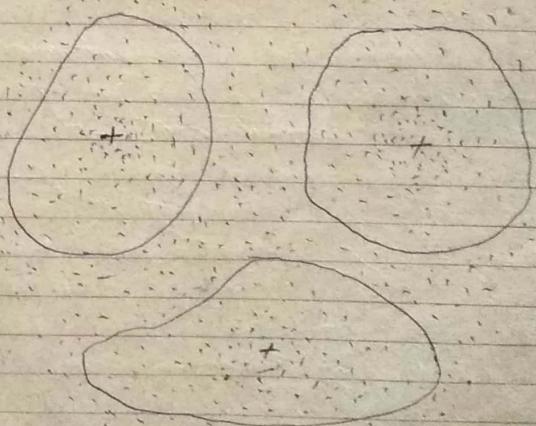
[if... then rule]



[Decision Tree]



★ Cluster Analysis & Outlier Analysis 08/08/2019



cl
classification and prediction which analyse class labeled data object. Clustering analyse data objects without consulting a known class label

In general, the classes are not present in the training data simply because they are not known to begin with. Clusters can be used to generate such labels.

The objects are clustered or grouped based on the principle of maximizing the intra-class similarity and minimizing the inter-class similarity.

Example [Cluster Analysis]

cluster analysis can be performed on electronic customer data in order to identify homogeneous sub populations of customers. The cluster may represent individual

target group for marketing. Following figure shows a 2D plot of customers with respect to customers location in a city.

Outlier Analysis -

A database may contain data objects they do not comply (contain) with the general behavior or model of the data. These objects are called noise. DM method classify outliers as noise or exceptions.

Data mining issues :-

1. Mining methodologies and user interaction issue

2. Performance issue

3. issue relating to the diversity of database types.

★ Mining methodology and user interaction issues -

These reflects the kind of knowledge mined, the ability to mine knowledge at multiple granularities, the use of domain knowledge, ad hoc mining and knowledge visualization.

• Mining different kind of knowledge in databases :-

Data mining should cover a wide spectrum of data analysis and knowledge discovery tasks, including data characterization, discrimination, association and correlation analysis, classification,

PAGE NO. 73
DATE: 12/12/20

PAGE NO. 74
DATE: 12/12/20

prediction, clustering, outliers analysis and evolution analysis.

• Interactive mining of knowledge at multiple levels of abstraction :-

Because it is difficult to know exactly what can be discovered within a database, the data mining process should be interactive.

• Incorporation of background knowledge :-

Background knowledge or information regarding the domain under study, may be used to guide the discovery process.

• Data mining query language and ad hoc data mining :-

Data mining query language should be integrated with a database or data warehouse query language and optimized for efficient and flexible data mining.

• Presentation and visualization of data mining results :-

Discovered knowledge should be expressed

PAGE NO. 75
DATE: 12/12/20

in high level languages, visual representations, other expressive forms so that the knowledge can be easily understood and directly used by humans.

• Handling noisy or incomplete data:-
Data cleaning method and analysis methods that can handle noise as required, as well as outliers mining method for the discovery and analysis of exception cases.

• Pattern evaluation - (the interestingness problem)
The use of interestingness measures or user-specified constraints to guide the discovery process and reduce the search space is another active area of research.

• Performance issues:-
These includes: efficiency, scalability, and parallelization of data mining algorithms.

PAGE NO. 76
DATE: 12/12/20

- Efficiency and scalability of data mining algorithms :-

To effectively extract information from a huge amount of data in databases, data mining algorithm must be efficient and scalable. or,

The running time of data mining algorithm must be predictable and acceptable in large databases.

- Parallel, distributed, and incremental mining algorithms :-

The huge size of many databases, the wide distribution of data, and the computational complexity of some data mining methods are factors motivating the development of parallel and distributed data mining algorithm.

High cost of data mining process promotes the need for incremental data mining algorithms that incorporate database updates without having to mine the entire data again "from scratch".

* Issue relating to the diversity of database types:

- Handling of relational and complex type of data:

It is not realistic to expect system to mine all kinds of data given the diversity of data types and diverse goals of data mining. Therefore, one may expect to have different data mining techniques for different kinds of data.

- Mining information from heterogeneous databases and global information systems

Discovery of knowledge from diverse sources of structured, semistructured, and unstructured data with diverse data semantics poses great challenges to data mining. Mining may help disclose high-level data regularities in multiple heterogeneous databases that are unlikely to be discovered by simple query systems and may improve information exchange and interoperability of heterogeneous databases.