

UNIT-II

INTRODUCTION TO Data Mining

Introduction:-

Data mining is one of the most useful techniques that helps entrepreneurs, researchers, and individuals to extract valuable information from huge sets of data. Data mining is also called as Knowledge Discovery in Database (KDD). The knowledge discovery process includes Data cleaning, Data integration, Data selection, Data transformation, Data mining, Pattern evaluation and knowledge presentation.

What is datamining?

The process of extracting information to identify patterns, trends and useful data that would allow the business to take the data-driven decision from huge sets of data is called Datamining.

Datamining is the process of investigating hidden patterns of information to various perspectives for categorization into useful data which is collected and assembled in particular areas such as data warehouses, operational data

databases, helping decision making and often data requirement to eventually cost-cutting and generating revenue.

Data mining is a process used by organizations to extract specific data from huge databases to solve business problems.

It primarily turns raw data into useful information KDD Process :-

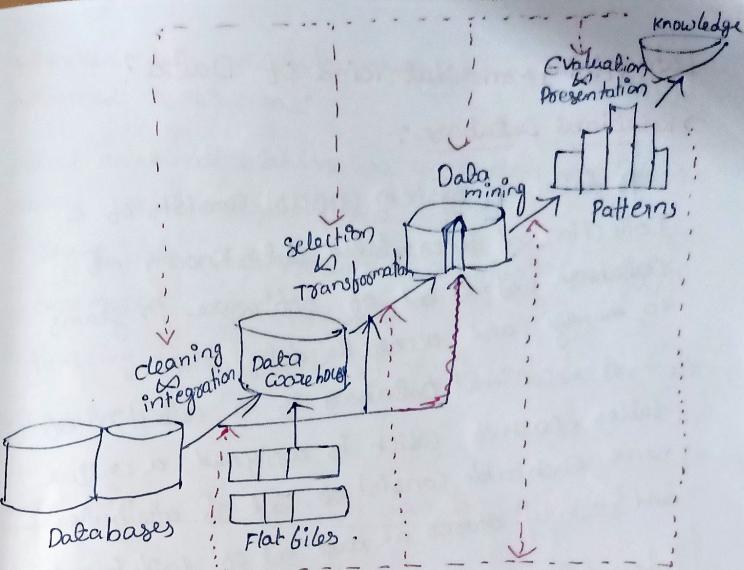
1. Data cleaning - to remove noise and inconsistent data
2. Data integration - where multiple data sources may be combined.
3. Data selection - where data relevant to the analysis task are retrieved from the database
4. Data transformation - where data are transformed and consolidated into forms appropriate for mining by performing summary or aggregation operations.
5. Data mining - an essential process where intelligent methods are applied to extract data patterns
6. Pattern evaluation - to identify the truly interesting patterns representing knowledge based on interestingness measures.
7. Knowledge presentation - where visualization and knowledge representation techniques are used to present mined knowledge to users.

Advantages of Datamining :-

- The Datamining technique enables organizations to obtain knowledge based Data.
- Datamining enables organizations to make lucrative modification in operation and production.
- Compared with other statistical data applications, datamining is a cost efficient.
- Datamining helps the decision making process of an organization.
- It is a quick process that makes it easy for new users to analyze enormous amount of data in short time.

Disadvantages :-

- There is a probability that the organization may sell useful data of customers to other organization for money.
- Many data mining analytics software is difficult to operate and needs advance training to work on.
- The Data mining techniques are not precise so that it may lead to severe consequences in certain conditions.



Datamining as a step in the process of Knowledge Discovery.

KDD :- The term KDD stands for knowledge discovery in Databases. It refers to the broad procedure of discovering knowledge in data and emphasize the high-level applications of specific Data mining techniques.

The Main objective of KDD Process is to extract information from data in the context of large databases. It does this by using data mining algorithms to identify what is deemed knowledge.

Data Mining - on what kind of Data?

→ Relational Database :-

A database system (DBMS) consists of a collection of interrelated data known as Database and a set of software programs to manage and access the data.

A relational Database is a collection of tables, each of which is assigned a unique name. Each table consists of set of attributes and usually stores a large set of tuples (records). When data mining is applied to relational database, we can go further by searching for trends or data patterns.

Datawarehouse :-

A datawarehouse usually modeled by a multidimensional data base structure where each dimension corresponds to an attribute or a set of attributes in the schema.

Transactional data :-

Transactional Database consists of a file where each record represents a transaction. A transaction typically includes a unique transaction identifier number and a list of the items making up the transaction.

Advanced data and Information systems and Advanced Applications:

object relational databases are constructed based on an object-relational data model.

Sequence data
 ↓
 stockmarket

↓
 data streams, spatial data (maps),
 continuous transmit

engineering design data (IC's), hyper text, multimedia,
 web data etc.

→ what kind of Patterns Can be mined?

Descriptive :- It describes the general properties of data

Predictive :- it makes predictions based on current data.

Data mining functionalities are used to specify the kind of patterns to be found in data mining tasks.

→ Concept / class Description : Characterization and discrimination eg: electronic store

→ Mining frequent patterns, associations and correlation

→ Classification and prediction.

→ Cluster analysis

→ Outlier Analysis - out of group.

→ Evolution Analysis - changes over time.

Example for characterization :-

A data mining system should be able to produce a description summarizing the characteristics of customers who spend more than \$1000 a year at All Electronics store.

The result should be General profile of the customers

- 40-50 years old
- Employed
- Excellent credit rating

Example for Data discrimination :-

A data mining system should be able to compare two groups of All Electronics store customers.

- Regular customers
 - Irregular customers
- Comparative profile of the customers.
 - 80% of the customers who frequently purchase computer products are between 20 and 40 only.
 - Have a university education.
 - Whereas 60% of the customers who infrequently buy such products are either seniors or youth.
 - They have no university degree.

Mining frequent Patterns, Associations and Correlation

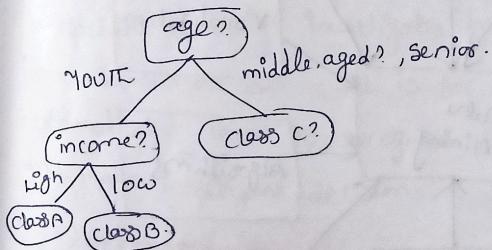
→ Frequent Patterns as the name suggest are patterns that occur frequently data

buys("X"; "computer") → buys("software")

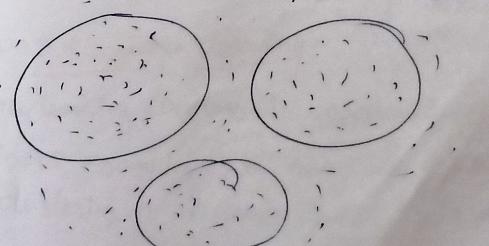
Classification & Prediction :-

Classification is the process of finding a model that describes and distinguishes data classes or concepts for the purpose of being able to use the model to predict

The class of objects whose class label is unknown

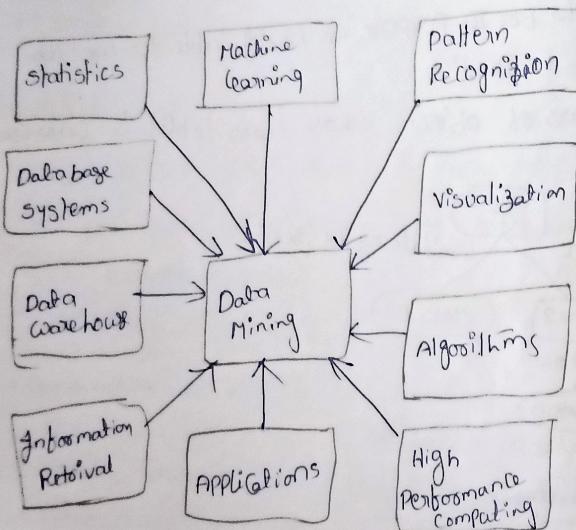


Cluster analysis :-



which technologies Are used ?

Data mining has incorporated many techniques from other domains such as statistics, machine learning, pattern recognition, databases & database warehouse systems, information retrieval, visualization, algorithms, high performance computing & many application domains.



statistics:-

statistics studies the collection, analysis, interpretation or explanation, and presentation of data. Data mining has an inherent connection with statistics.

A statistical model is a set of mathematical functions that describe the behaviour of the objects in a target class in terms of random variable and their associated probability distributions.

statistical models are widely used to model data and data classes.

Machine learning:-

Machine learning investigates how computers can learn based on data. ML is for computer programs to automatically learn to recognize complex patterns and make intelligent decisions based on data.

a) supervised learning:

The supervision in the learning comes from the labeled examples in the training dataset.

b) unsupervised learning:

The learning process is unsupervised since the input examples are not class labeled. we may have clustering to discover classes within the data.

Semi-supervised Learning :-

using both labeled & unlabeled (half, half)

Active learning :-

The algorithm is designed in such a way that the desired output should be decided by the algorithm itself.

Algorithm :-

An algorithm can be written for any application in data mining.

Application :-

Applications can be built & used as a technology.

Pattern Recognition :-

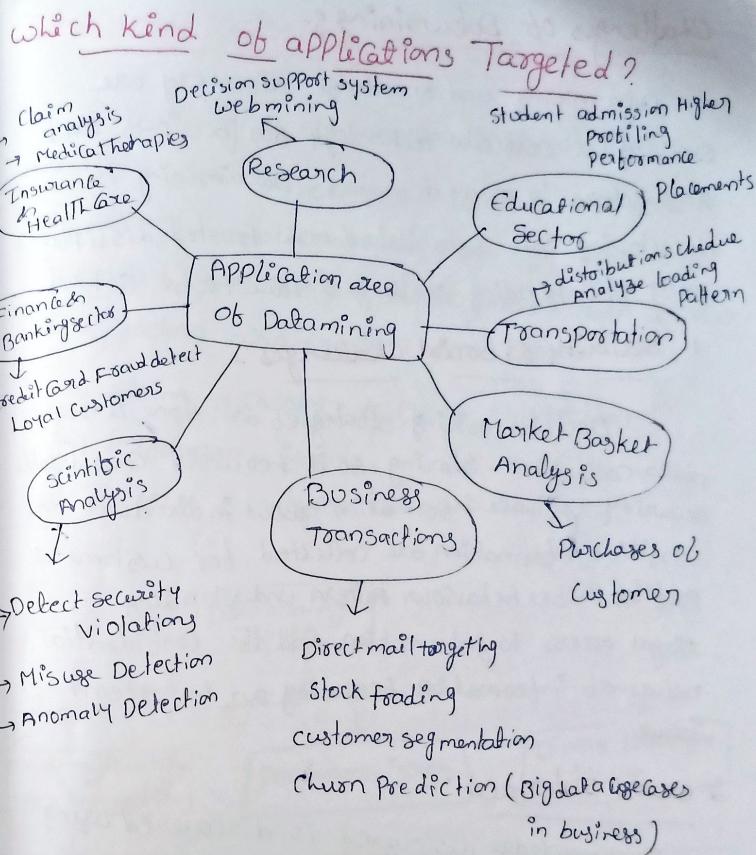
By observing various patterns we can perform data mining.

Visualization :-

To depending analysis of data by putting into graph we do visualization.

Database technology :-

using database SQL we extract data from database.



Challenges of Datamining :-

Data mining and knowledge discovery are evolving a crucial technology for business and researchers in many domains. Datamining is developing into established and trusted discipline. Many still pending challenges have to be solved.

1. Security & Social Challenges :-

Decision making strategies are done through data collection - sharing, so it requires considerable security. Private information about individual and sensitive information are collected for customers profiles, user behaviour pattern understanding. Illegal access to information and the confidential nature of information becoming an important issue.

2. User Interface :-

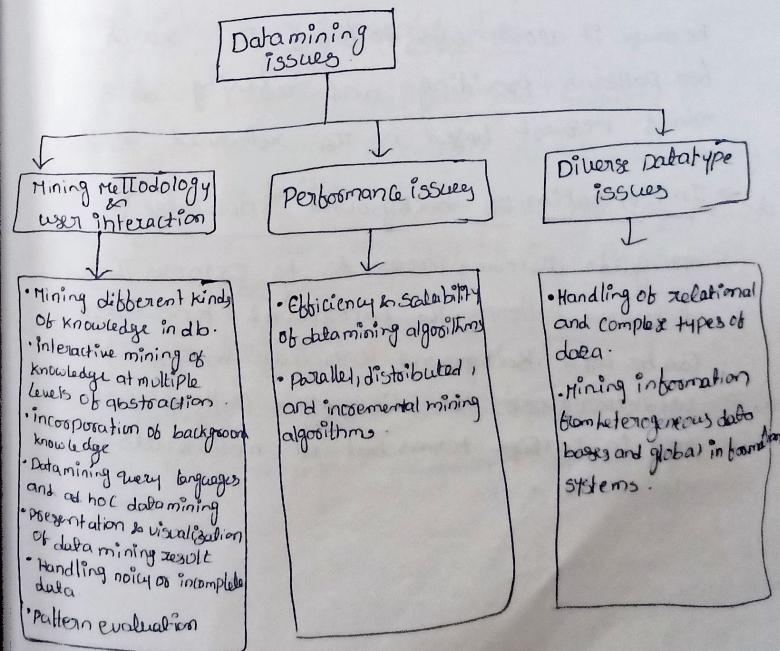
The knowledge discovered is discovered using datamining tools is useful only if it is interesting and above all understandable by the user.

From good visualization interpretation of data, mining results can be eaged and helps better understand their requirements.

Challenges & issues in Datamining :-

Data mining is not an easy task, as the algorithms used can get very complex and data is not always available at one place. It needs to be integrated from various heterogeneous data sources. These factors also create some issues.

- Mining Methodology & User interaction
- Performance issues
- Diverse Datatype issues



Mining Methodology & User interaction issues:

- It refers following kinds of issues:
 - Mining different kinds of knowledge in databases:
Different users may be interested in different kinds of knowledge. Therefore it is necessary for data mining to cover a broad range of knowledge discovery task.
 - Interactive mining of knowledge at multiple levels of abstraction:
The Data mining process needs to be interactive because it allows user to focus the search for patterns, providing and refining data mining request based on the returned result
 - Incorporation of background knowledge:
To guide discovery process & to express the discovered patterns, the background knowledge can be used. Background knowledge may be used to express the discovered patterns not only in concrete terms but at multiple levels of abstraction.

Datamining query languages and adhoc datamining:
Datamining query language that allows the user to describe adhoc mining tasks, should be integrated with a data warehouse query language and optimized for efficient & flexible Data mining.

Presentation & visualization of datamining results:
Once the patterns are discovered it needs to be expressed in high level languages, and visual representation. These representations should be easily understandable.

Handling noisy or incomplete data:

The data cleaning methods are required to handle the noise and incomplete objects while mining the data regularities. If the data cleaning methods are not there then accuracy of the discovered pattern will be poor.

Pattern evaluation:

The pattern discovered should be interesting because either they represent common knowledge or lack novelty.

Performance issues:-

There can be performance-related issues:-

Efficiency & scalability of data mining algorithms:-

In order to effectively extract the information from huge amount of data in databases, data mining algorithm must be efficient and scalable.

Parallel, distributed and incremental mining algorithms:-

The factors such as huge size of databases, wide distribution of data, and complexity of data mining methods motivate the development of parallel and distributed data mining algorithms. These algorithms divide the data into partitions which is further processed in a parallel fashion. Then the result from the partitions is merged. The incremental algorithms, update databases without mining the data again from scratch.

Diverse Datatype issues:-

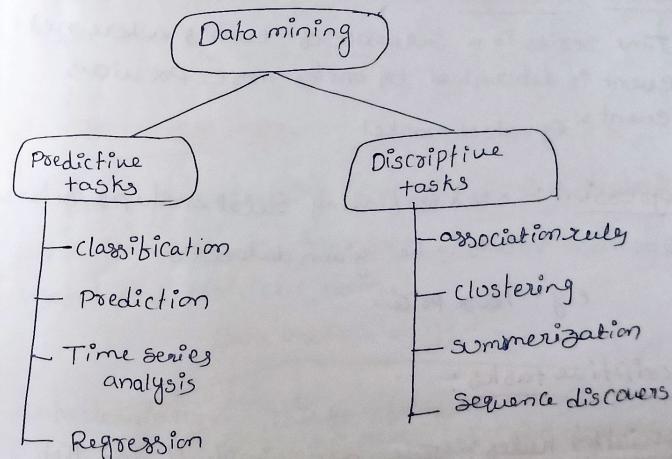
Handling of relational and complex types of data:-

The database may contain complex data objects, multimedia data objects, spatial data, temporal data, etc. It is not possible for one system to mine all these kinds of data.

Mining information from heterogeneous databases and global information systems:-

The data is available at different data sources on LAN or WAN. These data sources may be structured, semi-structured or unstructured. Therefore mining the knowledge from them adds challenges to data mining.

Data mining Tasks:-



Prediction:- Prediction uses same variables or field in the database to predict unknown or future values of other variables.

Description:- Characterize the general properties of data and used for finding patterns that describes the data.

Classification:-

Used to classify data into one of several predefined classes. e.g.: Loan $\begin{cases} \text{Accept} \\ \text{Reject} \end{cases}$ Based on data

Prediction:- Predictive task come up with a model from the available data set that is helpful in predicting unknown or future values of another dataset
e.g.: Based on medical reports Doctor predict the disease

Time series analysis:-

Time series is a sequence of events where next event is determined by one or more previous events
e.g.: stock market

Regression:- used for finding relationships between two more values in the given datasets
e.g.: house price.

Disciptine tasks:-

Association Rules:- These are simple if and then statements that help to discover relationships between datasets

e.g.: If person buy milk Then he buy bread

clustering:- It is a division of information into groups of connected objects. Cluster is used to identify similar data.
e.g.: Library books

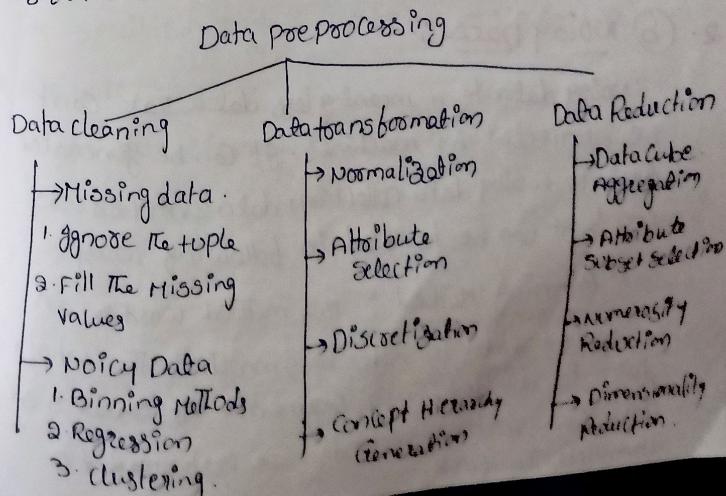
summarization:- Used for finding compact description for subset of data.

Sequence discovery:-

Helps to discovery similar patterns in transactional data over time.

Data Preprocessing:-

Data Preprocessing is a datamining technique which is used to transform the raw data in a useful and efficient format.



1. Data cleaning:-

The data can have many irrelevant and missing parts. To handle this part, data cleaning is done. It involves handling of missing data, noisy data etc.

(a) Missing Data:-

This situation arises when some data is missing in the data. It can be handled in various ways:

1. Ignore the tuples:-

This approach is suitable only when the dataset we have is quite large and multiple values are missing within a tuple.

2. Fill the Missing values:-

There are various ways to do this task. You can choose to fill the missing values manually, by attribute mean or the most probable value.

(b) Noisy Data:-

Noisy data is a meaningless data that can't be interpolated by machines. It can be generated due to faulty data collection, data entry errors etc. It can be handled in following ways.

1. Binning Method:- This method works on sorted data in order to smooth it. The whole data is divided into segments of equal size and then various methods are performed to complete the task.

2. Regression:-

Here data can be made smooth by fitting it to a regression function. The regression used may be linear (having one independent variable) or multiple (having multiple independent variables.)

3. Clustering:-

This approach groups the similar data in a cluster. The outliers may be undetected or it will fall outside the clusters.

② Data transformation:-

This step is taken in order to transform the data in appropriate forms suitable for mining process.

1. Normalization:-

It is done in order to scale the data values in a specified range (-1.0 to 1.0 or 0.0 to 1.0).

2. Attribute selection:-

In this strategy, new attributes are constructed from the given set of attributes to help the mining process.

3. Discretization:-

This is done to replace the raw values of numeric attribute by interval levels or conceptual levels.

4 Concept Hierarchy Generation:

Here attributes are converted from lower level to higher level in hierarchy for example "City" can be converted to "Country".

③ Data Reduction:

Data mining is a technique that is used to handle huge amount of data while working with huge volume of data, analysis become harder in such cases. In order to get rid of this, we use data reduction technique. It aims to increase the storage efficiency and reduce data storage and analysis costs.

The various steps of data reduction are:

1. Data Cube Aggregation:

Aggregation operation is applied to data for the construction of the data cube.

2. Attribute subset selection:

The highly relevant attributes should be used rest all can be discarded. For performing attribute selection, one can use level of significance and P-value of the attribute. The attribute having P-value greater than significance level can be discarded.

3. Numerosity Reduction:

This enable to store the model of data instead of whole data. For eg: Regression models.

4. Dimensionality Reduction:

This reduce the size of data by encoding mechanisms. It can be lossy or lossless. If after reconstruction from compressed data, original data can be retrieved, such reduction are called lossless reduction else it is called lossy reduction.

More effective methods of dimensionality reduction are:

→ Wavelet transforms &

→ PCA (Principal Component Analysis)

Data Cleaning:

Data cleaning routines attempt to fill in missing values, smooth out noise while identifying outliers, and correct inconsistencies in the data.

→ Missing values:

i) Ignore tuple:

which row having more missing values

can ignore that

ii) Fill in the missing values manually

iii) Use a global constant to fill in the missing values (e.g. NA)

iv) Use a measure of central tendency for the attribute to fill the missing value

v) Use the most probable value to fill in the missing value (e.g. decision tree)

S. No	Name	Occupation	Branch	Date	Price	Address	Pin code
1	Ramesh	Govt	TPG	10-May		TPG	53420
2	George	Self	TNK	11-May	2500	TNK	
3	Ponam	Private	TPG	11-May	500	TPG	53410
4	Govind	Private	VJWD	12-May	600	VJWD	52000
5	Ravi						
6	Ramesh	Business	VJWD	14-May	1400	VJWD	52000
7	Bhuwan	Business	TPG	14-May	2300	TPG	53420
8	Kiran	Govt	TNK			TNK	53310
9	Self	RSY	15-May	1100	RSY		53320
10	Yogi	Business	TNK	16-May	1800	TNK	53420

→ Noisy Data:

Noise is a random error or variation in a measured variable

Approaches in noisy data:

- (i) Binning
- (ii) Regression
- (iii) Outlier analysis

→ Binning:

a) Partition into equal frequency bins

b) Smoothing by bin means

c) Smoothing by bin boundaries.

Eg: 6, 10, 17, 22, 22, 25, 27, 30, 36

Partition into equal frequency bins

Bin 1: 6, 10, 17

Bin 2: 22, 22, 25

Bin 3: 27, 30, 36

$$m = \frac{\text{Sum of all terms}}{\text{no. of terms}}$$

Smoothing by bin means finding mean value

Bin 1: 11, 11, 11

Bin 2: 23, 23, 23

Bin 3: 31, 31, 31

smoothing by bin boundaries

smoothing by b

Bin1 : 6, 6, 17

Bin2 : 22, 22, 25

Bin3 : 27, 27, 36

→ Regression :-

Linear regression involves finding the "best" line to fit two attributes so that one attribute can be used to predict the other.

Multiple linear regression is an extension of linear regression where more than two attributes are involved and the data are fit to multidimensional surface.

→ outlier analysis:-

outlier may be detected by clustering for example, where similar values are organized into groups, or "clusters". Intuitively, values that fall outside of the set of clusters may be considered outliers.

Data Transformation:-

It is a data pre processing technique that transforms (or) ^{consol} consolidate the data into alternate forms appropriate for mining.

- (i) smoothing: Remove the noise from data [binning, regression, clustering]
- (ii) Aggregation: summary or aggregate function. It is used to constructing a data cube
- (iii) Generalization: low level concepts are replaced with high level. e.g. city to country
- (iv) Normalization:

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

e.g. {2, 40, 500, 1, 3, 900}

we should change these values into 0-1 range
that means we can normalize the attribute values
In this we have two types. Cse 103

- (i) Min-Max normalization:

$$v' = \frac{v - v_{\min}}{v_{\max} - v_{\min}} \quad v' - \text{new value}$$

v_{\min} - minimum value of attribute
 v_{\max} - maximum value of attribute
 v - original attribute value.

(ii) Z-score normalization.

$$v' = \frac{v - \bar{x}}{\sigma}$$

\bar{x} - mean of all attributes
 σ - standard deviation
 v - original value
 v' - new value

e.g.: min-max normalization.

8
10
15 } Marks. It will scale the data into
20 0 & 1

$$v' = \frac{v - \min_n}{\max_n - \min_n}$$

In the given attribute $\min_n = 8$
 $\max_n = 20$

(i) original value $v = 8$

$$v' = \frac{8 - 8}{20 - 8} = \frac{0}{12} = 0$$

new value = 0

(ii) 10

$$v' = \frac{10 - 8}{20 - 8} = \frac{2}{12} = \frac{1}{6} = 0.16$$

(iii) 15

$$v' = \frac{15 - 8}{20 - 8} = \frac{7}{12} = 0.58$$

(iii) 20

$$v' = \frac{20 - 8}{20 - 8} = \frac{12}{12} = 1$$

new value = 1

new values are

8	0	new min = 0
10	0.16	new max = 1
15	0.58	
20	1	

e.g.: Z score normalization:

$$v' = \frac{v - \bar{x}}{\sigma}$$

mean of marks = total of all attributes / no of attributes

$$\begin{aligned} \text{standard deviation} &= \sqrt{\frac{\sum (x - \bar{x})^2}{n}} \\ &= \sqrt{\frac{\sum (\text{every individual marks} - \text{mean of marks})^2}{n}} \\ &= \end{aligned}$$

8
10
15
20 } marks

$$\text{mean of marks} = \frac{8+10+15+20}{4} = 13.25$$

$$\begin{aligned}\text{standard deviation} &= \sqrt{\frac{(8-13.25)^2 + (10-13.25)^2 + (15-13.25)^2 + (20-13.25)^2}{4}} \\ &= \sqrt{\frac{(-5.25)^2 + (-3.25)^2 + (1.75)^2 + (6.75)^2}{4}} \\ &= \sqrt{\frac{27.56 + 10.56 + 3.06 + 45.56}{4}} \\ &= \sqrt{\frac{86.74}{4}} = \sqrt{21.6} = 4.6\end{aligned}$$

$$\text{standard deviation} = 4.6$$

$$(i) S = \frac{8-13.25}{4.6} = -1.14$$

$$(ii) 10 = \frac{10-13.25}{4.6} = -0.70$$

$$(iii) 15 = \frac{15-13.25}{4.6} = 0.36$$

$$(iv) 20 = \frac{20-13.25}{4.6} = 1.46$$

Data Reduction:

It is a Postprocessing technique that helps in obtaining reduced representation of data set from the available data set.

- Integrity of the original data should be maintained after reduction in data volume.
- It should produce same analytic result as on original data.

Data cube aggregation:- It is a process for which information is gathered and expressed in a summary form.

e.g.: Year 2017

Half Year	Sales
H1	500
H2	300

Year 2018

Half Year	Sales
H1	600
H2	100

after aggregation

Year	Sales
2017	800
2018	700

Dimensionality reduction:-

It eliminates the redundant attributes which are weakly important across the data
 $DOB \rightarrow age$

- (i) stepwise forward selection
- (ii) stepwise backward selection
- (iii) Decision tree induction.

i) Stepwise forward selection

If represents the original data in the compressed or reduced form by applying data encoding (or) transformation. In this we have Lossless & lossy.

Lossless :- If original data can be reconstructed from compressed data without losing any information

Lossy :- If reconstructed data is the approximation of compressed data.

Two effective methods of dimensionality Reduction are:

→ Wavelet transforms :- (DWT) - discrete wavelet transform

It is a linear signal processing technique that when applied to data vector transforms it numerically different vector x^t of wavelet coefficient

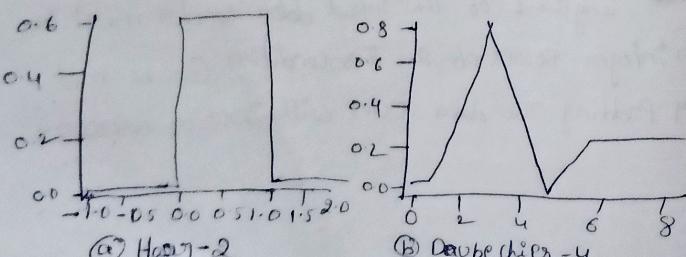
→ When applying this technique we consider each tuple as n-dimensional data vector depicting n measurements made on the tuple from n different database attributes.

→ It removes noise without smoothing the data.

There are several families of DWTs, some popular wavelet transforms include the Haar-2, Daubechies and Daubechies-6. The general procedure for applying a discrete wavelet transform uses a hierarchical pyramid algorithm that halves the data at each iteration, result is fast computational speed. The method follows

1. The length L of the input data vector must be an integer power of 2. This condition can be met by padding the data vector with zeros as necessary.

- Each transform involves applying two functions. The first applies some data smoothing, such as a sum of weighted average. The second performs a weighted difference, which acts to bring out the detailed features of the data.
- 3. The two functions are applied to pairs of data points in X , that is to all pairs of measurements (x_{2i}, x_{2i+1}). This result in two datasets of length $L/2$. In general these represent a smoothed or low-frequency version of the input data and the high frequency content of it.
- 4. The two functions are recursively applied to the data sets obtained in the previous loop until the resulting data sets obtained are of length 2.
- 5. Selected values from datasets obtained in the previous iteration are designated the wavelet coefficients of the transformed data.



A matrix multiplication can be applied to the input data in order to obtain the wavelet coefficients where the matrix used depends on the given DWT.

Principal Components Analysis:

Principal components analysis as a method of dimensionality reduction.

The data to be reduced consists of tuples or data vectors described by n attributes or dimensions.

PCA is also called Karhunen-Loeve (or K-L method) searches for K n -dimensional orthogonal vectors

that can best be used to represent the data, where $K \leq n$. The original data are thus projected on to a much smaller space, resulting in dimensionality reduction.

PCA combines the essence of attributes by creating an alternative, smaller set of variables.

The initial data can then be projected onto this smaller set.

The basic procedure is:

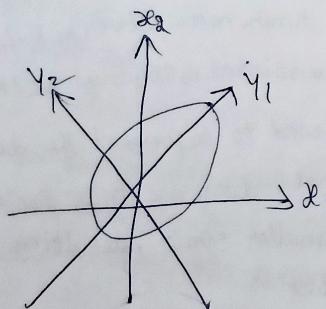
1. The input data are normalized, so that each attribute falls within the same range. This step helps ensure that attributes with large domains will not dominate attributes with smaller domains.

2. PCA computes K orthogonal vectors that provide basis for the normalized input data. These are unit vectors that each point in a direction perpendicular to the others.

These vectors are referred to as the Principal Components.

3. The Principal Components are sorted in order of decreasing "significance" (or) strength.

The Principal Components essentially serve as a new set of axes for the data. They provide important information about variance.



The sorted axes are such that the axis shows the most variance among the data, the second axis shows the next highest variance.

The first two principal components y_1 and y_2 for the given set of data originally mapped to the axes x_1 and x_2 . This information helps identify groups or patterns within data.

4. Because the components are sorted in decreasing order of "significance," the data size can be reduced by eliminating weaker components that is those with low variance.

Feature subset selection:

→ Another way to reduce the dimensionality is to use only a subset selection of the features.

→ While it might seem that such an approach would lose information, this is not the case if redundant features are present.

→ Redundant features duplicate much or all of the information contained in one or more other attributes. For example, the purchase price of a product and the amount of sales tax paid contain much of the same information.

→ Irrelevant features contain almost no useful information for the data mining task at hand. For instance, students' ID numbers are irrelevant to the task of predicting student's grade point averages.

- Redundant and irrelevant features can reduce classification accuracy and the quality of the clusters that are found.
 - while some irrelevant and redundant attributes can be eliminated immediately by using common sense or domain knowledge, selecting the best subset of features subsequently requires a systematic approach.
 - the ideal approach to feature selection is to try all possible subsets of features as input to the datamining algorithm of interest, and then take the subset that produces the best results. this method has the advantage of reflecting the objective and bias of the datamining algorithms that will eventually be used. unfortunately since the number of subsets involving n attributes is 2^n such an approach is impractical in most situations. There are three standard approaches to feature selection: embedded, filter & wrapped.
1. Embedded approaches:-
- Feature selection occur naturally as part of the datamining algorithms. The algorithm specifically, during the operation of the data mining algorithm, the algorithm itself decides which attribute to use and which to ignore.

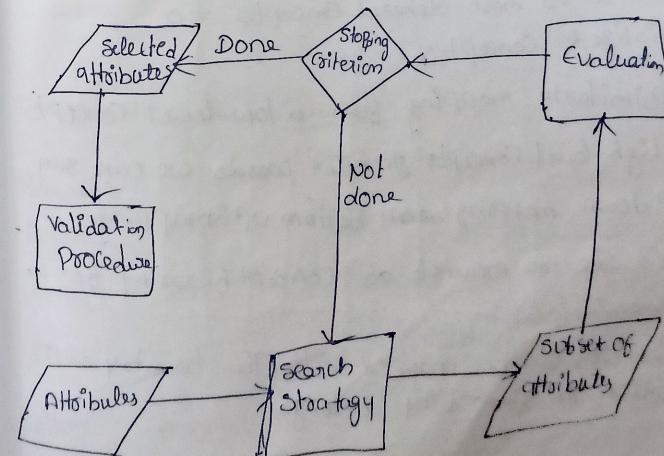
2. Filter approaches:

Features are selected before the datamining algorithm is run, using some approach that is independent of the data mining task.

3. Wrapped approaches:

These methods use the target datamining algorithm as a black box to find the best subset of attributes in a way similar to that of the ideal algorithm described above but typically without enumerating all possible subsets.

An architecture for Feature subset selection:



Discretization & Binarization:-

Discretization is used to transform the attributes that are in continuous format.

Data Discretization converts a large number of data values into smaller one, so that data evaluation and data management becomes very easy.

Eg: we have an attribute of age with the following values.

Age : 10, 11, 13, 14, 17, 19, 30, 31, 32, 38, 40, 42, 70, 72, 73, 75

Attribute	age 1	age 2	age 3
After Discretization	10, 11, 13, 14, 17, 19	30, 31, 32, 38, 40, 42	70, 72, 73, 75
Young			
Mature			
old			

Hierarchy Generation data discretization and concepts:-

A Concept Hierarchy represents a sequence of mapping with a set of more General Concepts to specialized concepts.

Similarly mapping from a low-level concepts to high-level concepts. In other words we can say top down mapping and bottom up mapping.

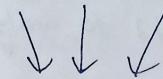
Let's see an example of concept hierarchy for the dimension location.

Each city can be mapped with the country with which the given city belongs.

For example Delhi can be mapped to INDIA and INDIA can be mapped to Asia.

TOP-DOWN MAPPING:-

TOP down mapping starts from top with General Concepts and move to the bottom to the specialized concepts.

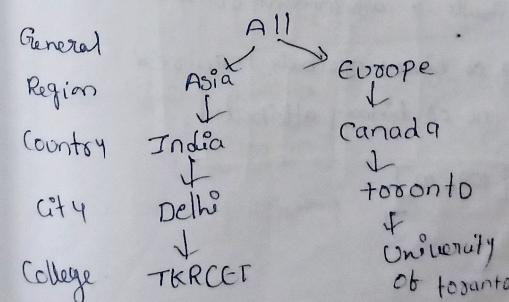


BOTTOM UP MAPPING:-

Bottom up mapping starts from bottom with specialized concept and move to the top to the Generalized concepts.



Concept Hierarchy Generation :-



Binaryization:-

Binaryization is used to transform both the discrete attributes and the continuous attributes into binary attributes in data mining with respect to feature selection.

Best Binaryization approach is the one that produces the best result for the data mining algorithm that will be used to analyze the data.

Simple techniques for binarization:

- Assigning numerical value
- Finding number of binary attribute required
- Conversion into binary

eg: say There is an categorical attribute with 'm' number of values.

→ Assigning numerical value:
Number assigned will be between [0, m-1]

For ordinal attribute → assignment follows order

→ finding number of binary attribute required
say n be the no. of binary attributes

$$n = \lceil \log_2 m \rceil$$

Converting the number assigned into binary value.

e.g.: if number of binary attribute is $n=3$ in numbers

Then we can write three bit binary number,

$$2 = 010$$

if the number of binary attribute is $n=4$ in numbers Then

$$2 = 0010$$

eg: let us consider an example to learn.

{awful, poor, ok, good, great} - ordinal form

Attribute Values	Integervalue
awful	0
poor	1
ok	2
good	3
great	4

Identifying number of binary attributes;

$$n = \lceil \log_2 m \rceil$$

$$n = \lceil \log_2 5 \rceil = 3$$

Binary conversion

Attribute value	Integer values	x_1	x_2	x_3
awful	0	0	0	0
poor	1	0	0	1
ok	2	0	1	0
good	3	0	1	1
great	4	1	0	0

Binarization

If we have mutual relation overcoming the issues

$$\text{No. of binary attribute} = \text{No. of values}$$

attribute values	Integer value	x_1	x_2	x_3	x_4	x_5
awful	0	1	0	0	0	0
Poor	1	0	1	0	0	0
ok	2	0	0	1	0	0
good	3	0	0	0	1	0
great	4	0	0	0	0	1

Data transformation:

Decimal scaling

- It normalizes the values of an attribute by changing the position of its decimal point.
- The no. of points by which the decimal point is moved can be determined by the absolute maximum value of attribute A.
- A value v of attribute A is normalized to v' by computing

$$v' = \frac{v}{10^j}$$

where j is the smallest integer such that ~~max~~ $\max(|v'|) < 1$

Eg: suppose values of an attribute P varies from -99 to 99

The maximum absolute value $P = 99$

for normalization the values we divide the numbers by 100 (i.e $j=2$) (no of integers in largest number so that values come out to be 0.98, 0.97 and so on)

Binary conversion

Attribute value	Integer values	x_1	x_2	x_3
awful	0	0	0	0
Poor	1	0	0	1
OK	2	0	1	0
Good	3	0	1	1
Great	4	1	0	0

Binarization

If we have mutual relation overcoming the issues

No. of binary attribute = No. of values

Attribute values	Integer value	x_1	x_2	x_3	x_4	x_5
awful	0	1	0	0	0	0
Poor	1	0	1	0	0	0
OK	2	0	0	1	0	0
Good	3	0	0	0	1	0
Great	4	0	0	0	0	1

Data transformation:

Decimal scaling:

- It normalizes the values of an attribute by changing the position of their decimal points.
- The no. of points by which the decimal point is moved can be determined by the absolute maximum value of attribute A.
- A value v of attribute A is normalized to v' by computing

$$v' = \frac{v}{10^j}$$

where j is the smallest integer such that ~~max~~ $\max(|v'|) < 1$

Eg.: suppose values of an attribute P varies from -99 to 99

The maximum absolute value $P = 99$

for normalization the values we divide the numbers by 100 (i.e $j=2$) (as there are 2 integers in largest number so that values come out to be as 0.98, 0.97 and so on)

Measures of Similarity & Dissimilarity -

Basics :-

Similarity :-

Numerical measures of how alike two data objects are is higher when objects are more alike often falls in the Range [0, 1]

Dissimilarity :-

Numerical measures of how different are two data objects lower when objects are more alike. Maximum dissimilarity is often 0. Upper limit varies.

Data matrix vs Dissimilarity matrix:-

SUPPOSE that we have n objects (eg: persons, items, courses) described by P attributes (eg: age, height, weight, gender)

The objects are $\mathbf{x}_1 = x_{11}, x_{12}, x_{13}, \dots, x_{1P}$
 $\mathbf{x}_2 = x_{21}, x_{22}, x_{23}, \dots, x_{2P}$ and so on

When x_{ij} is the value for object \mathbf{x}_i of the j th attribute.

Then Data matrix $\begin{bmatrix} x_{11} & \dots & x_{1P} \\ x_{21} & \dots & x_{2P} \\ \vdots & \ddots & \vdots \\ x_{n1} & \dots & x_{nP} \end{bmatrix}$

Similarity / Dissimilarity for objects with single attribute

p and q are the attribute values for two data objects

Attribute type	Dissimilarity	Similarity
----------------	---------------	------------

Nominal	$d = \begin{cases} 0 & \text{if } P = q \\ 1 & \text{if } P \neq q \end{cases}$	$s = \begin{cases} 1 & \text{if } P = q \\ 0 & \text{if } P \neq q \end{cases}$
---------	---	---

ordinal	$d = \frac{ P - q }{n-1}$	$s = 1 - \frac{ P - q }{n-1}$
---------	---------------------------	-------------------------------

Interval (00)	$d = P - q $	$s = -d, s = \frac{1}{1+d} (0)$
---------------	---------------	---------------------------------

$$s = 1 - \frac{d - \min_d}{\max_d - \min_d}$$

Dissimilarities between Data objects with multiple numeric attributes.

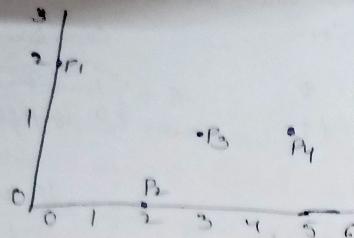
Euclidean Distance .

$$\text{dist} = \sqrt{\sum_{k=1}^n (p_k - q_k)^2}$$

where n is the no. of dimensions & p_k and q_k are respectively k^{th} attributes of data objects p and q .

→ Standardization is necessary, if scales differs.

Euclidean Distance:



Point	x	y
P ₁	0	2
P ₂	2	0
P ₃	3	1
P ₄	5	1

	P ₁	P ₂	P ₃	P ₄
P ₁	0	2.828	3.162	5.099
P ₂	2.828	0	1.414	3.162
P ₃	3.162	1.414	0	2
P ₄	5.099	3.162	2	0

Distance matrix

Minkowski distance:

Minkowski distance is a Generalization of Euclidean Distance. Given two objects p and q

$$dist = \left(\sum_{k=1}^n |P_k - Q_k|^r \right)^{1/r}$$

where r is a parameter, n is the number of dimensions, and P_k and Q_k are respectively, the kth attributes of data objects p and q.

e.g.: r=1 city block (Manhattan, taxicab, L₁ norm) distance

A common example of this is the Hamming distance which is just the number of bits that are different between two binary vectors.

r=2 Euclidean distance

r=∞, "supremum" (L_{max norm}, L_{infinity norm}) distance

This is the maximum difference between any attribute of the vectors.

$$d(x, y) = \lim_{r \rightarrow \infty} \left(\sum_{k=1}^n |x_k - y_k|^r \right)^{1/r}$$

Denote d by r will n i.e. all these distances are defined for all number of dimensions

Minkowski Distance

Point	x	y	L_1	P_1	P_2	P_3	P_4
P_1	0	2	P_1	0	4	4	6
P_2	2	0	P_2	4	0	2	4
P_3	3	1	P_3	4	2	0	2
P_4	5	1	P_4	6	4	2	0

	L_2	P_1	P_2	P_3	P_4
P_1	0	2.828	3.162	5.099	
P_2	2.828	0	1.414	3.162	
P_3	3.162	1.414	0	2	
P_4	5.099	3.162	2	0	

Manhattan Distance formula :-

$$D_m(x, y) = \sum_{i=1}^n |x_i - y_i|$$

P_1 at (x_1, y_1) and P_2 at (x_2, y_2)

it is $|x_1 - x_2| + |y_1 - y_2|$

$$= |0 - 2| + |2 - 0|$$

$$= 2 + 2$$

$$= 4$$

UNIT - III Association Rules

Association rule learning is a type of unsupervised learning technique that checks for the dependency of one data item on another data item and maps accordingly so that it can be more Profitable. It tries to find some interesting relations or association among the variable of dataset.

It is based on different rules to discover the interesting relations between variables in the database.

- Market basket analysis is a technique used by the various big retailer to discover the associations between items. We can understand it by taking an example of a supermarket as in a supermarket all products that are purchased together are put together.

For example if a customer buys bread, he most likely can also buy butter, eggs, or milk so these products are stored with in a shelf as mostly near by.

Customer 1	Customer 2	Customer 3	Customer 4
milk	milk	milk	Sugar
bread	bread	bread	eggs
butter	butter	butter	