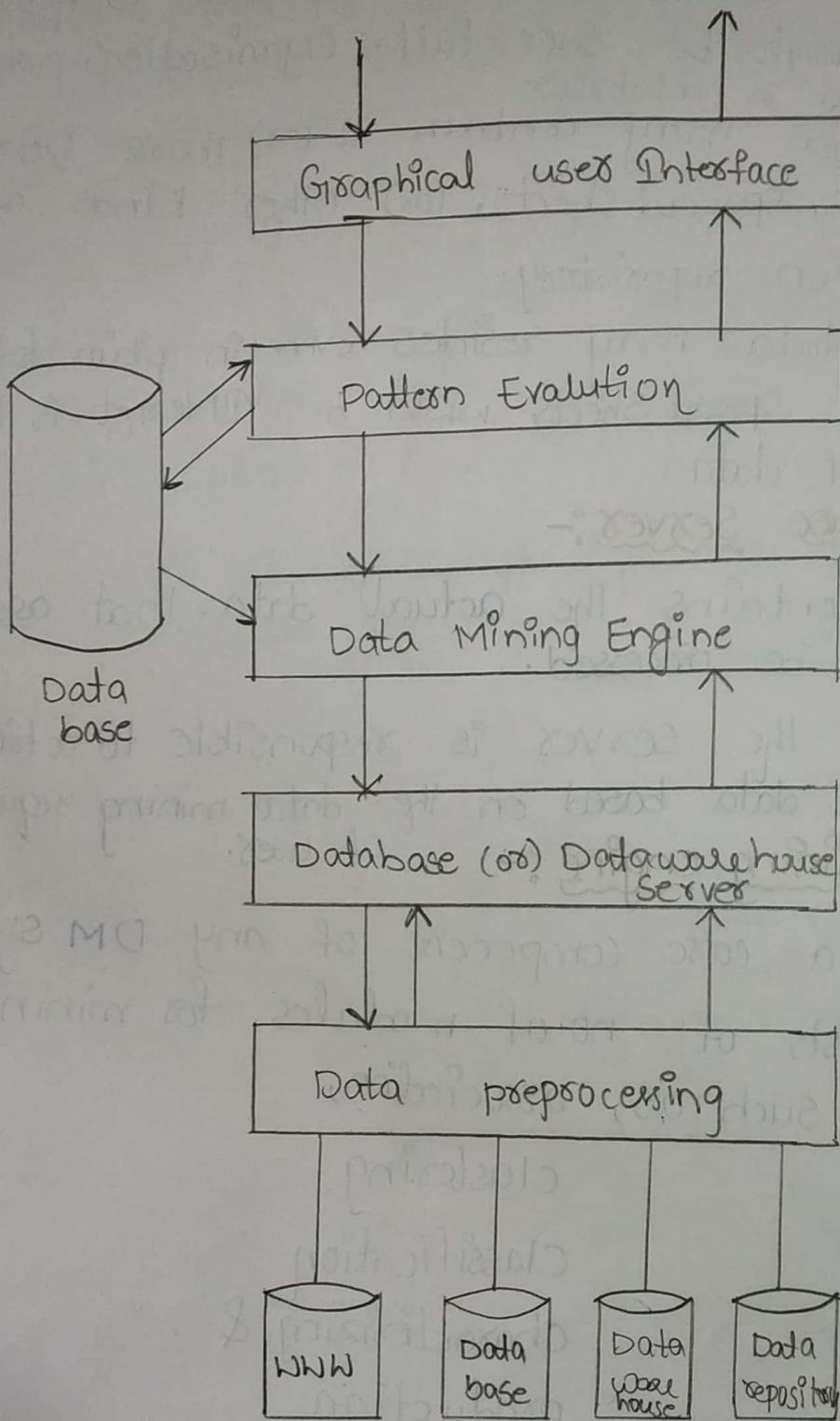


* Data Mining Architecture:-



Data Sources :-

Data base, DW house & WWW & text files & other documents are actual sources

- of data. You need historical data is needed
- * Large volumes of historical data is needed for DM to be successful. organisations manually store data in databases.
 - * DW house may contain 1 (or) more DBases text files, spread sheets, (or) other kinds of information repository.
 - * Some data may resides even in plain text files or spread sheets www or Internet is big source of data.
 - * DW house server :-
 - * It contains the actual data. that are ready to be processed.
 - * Hence the server is responsible to retrieving the relevant data based on the data mining request of users.
 - * Data mining engine :-
 - * It is a core component of any DM system. it consists of no. of modules for mining tasks. such as,
 - association,
 - clustering
 - classification
 - characterizing &
 - production
 - * Data mining Engine :- Pattern Evaluation Modules :-
 - * It is mainly responsible for the measure of interestingness of the patterns by using the a

threshold value.

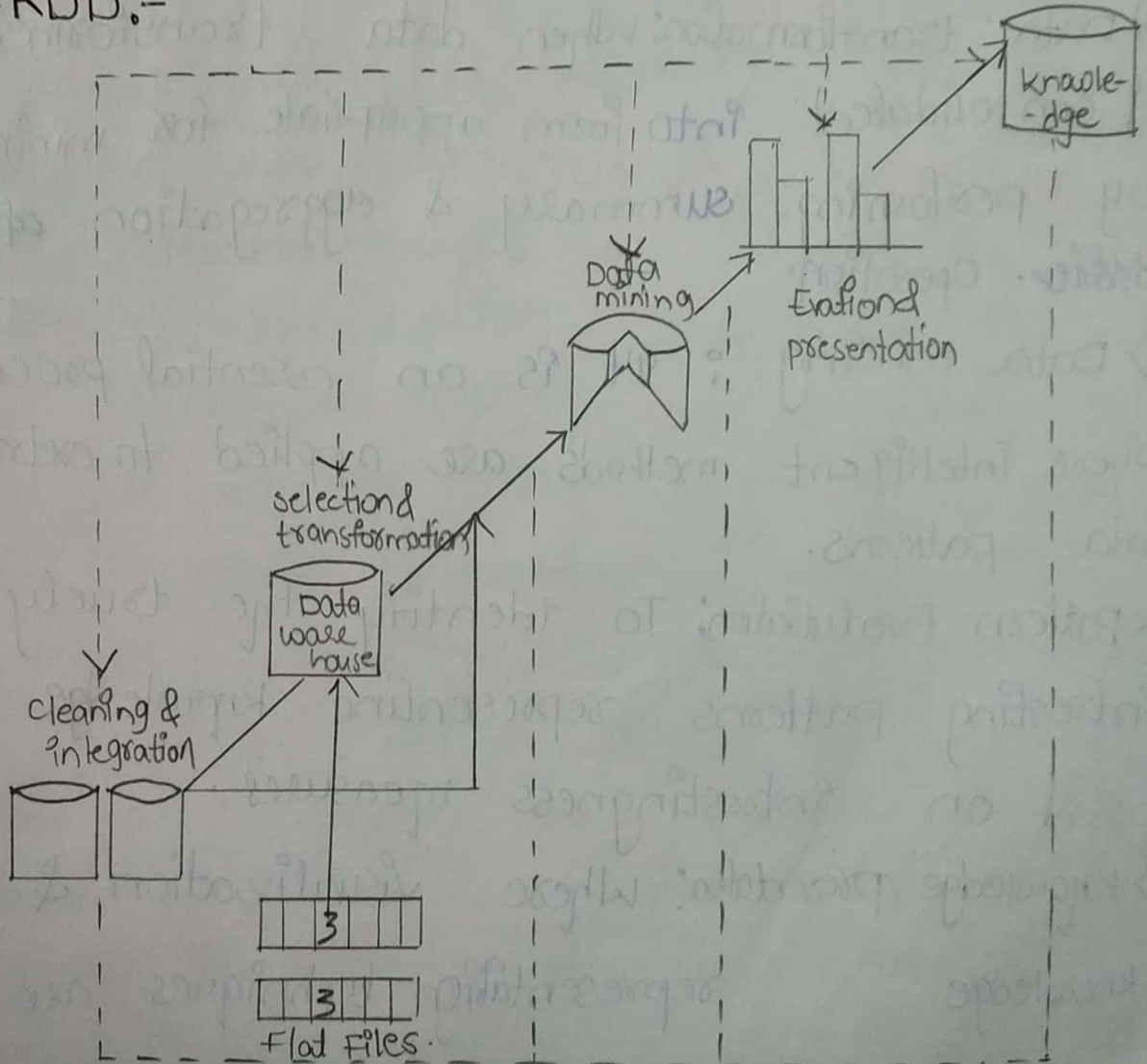
* It interacts with the DM Engine to focus the search towards data mining engine to focus the search towards interesting patterns.

* GUI :-

* GUI module communicates b/w the users & DM system. This module helps the users to use the system efficiently without knowing the real complexity behind the concepts.

* When the user specifies the ~~ass.~~ query or a task, and display the result in an easily understandable manner.

* KDD :-



* Knowledge Discovery in Data

* KD process is shown in the fig as interactive steps.

* The steps involved in KDD as

- 1) Data cleaning: To remove noise and inconsistent data.
- 2) Data Integration: Where multiple sources may be combined
- 3) Data selection: When data relevant to the analysis ~~tasks~~ are selected from the DB
- 4) Data transformation: When data transformed & consolidated into form appropriate for mining by performing summary & aggregation of ~~data~~ operation.
- 5) Data Mining: It is 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 & knowledge representing techniques are

used to present mined knowledge to users.
 * The preceding view shows data mining as one step in the knowledge discovery process,

* However, in Industry the term DM is used to refer KD process. \therefore We adopt a broad view of data mining functionality.

* What is data Mining?

Data mining is defined as extracting information from huge sets of data. In other words, we can say that data mining is the procedure of mining knowledge from data. The information or knowledge extracted, so can be used for any of the applications-

- 1) Market Analysis.
- 2) Fraud detection.
- 3) Customer Retention.
- 4) production control.
- 5) Science Exploration.

* What is the need for preprocessing of data?

In complete noisy & inconsistent data are common place properties of large scale real world databases & data warehouses. In complete data can occur for a no. of reasons. Attributes of interest may not

always be available, such as customer information for sales transition data. Other data may not be included simply because it was not considered imp at the time of entry relevant data may not be recorded due to misunderstanding. Data that were inconsistent with other recorded data may have been overlooked. Missing data, particularly for tuples with missing values for some attributes may need to be inferred.

* Data mining

DM is a process of extracting knowledge from a volume of data. It refers to a way of finding significant & useful information from an organisation. DM is a tool to convert data into information.

- * These are 3 reasons
- 1) knowledge discovery
 - 2) data visualization
 - 3) data correction.

* 1) It helps us to identify the invisible correlation points and trends all the data available in the DW house.

2) Data visualization:-

The main objective of this is to normalized large volumes of data to find the sensible way to display the data.

3) Data correction:-

The main objective of the DC is to identify & correct the incomplete & inconsistent data by removing errors.

* DATA MINING TASKS:-

These are divided into 2 categories.

1) Predictive tasks.

2) Descriptive tasks.

1) Descriptive tasks.

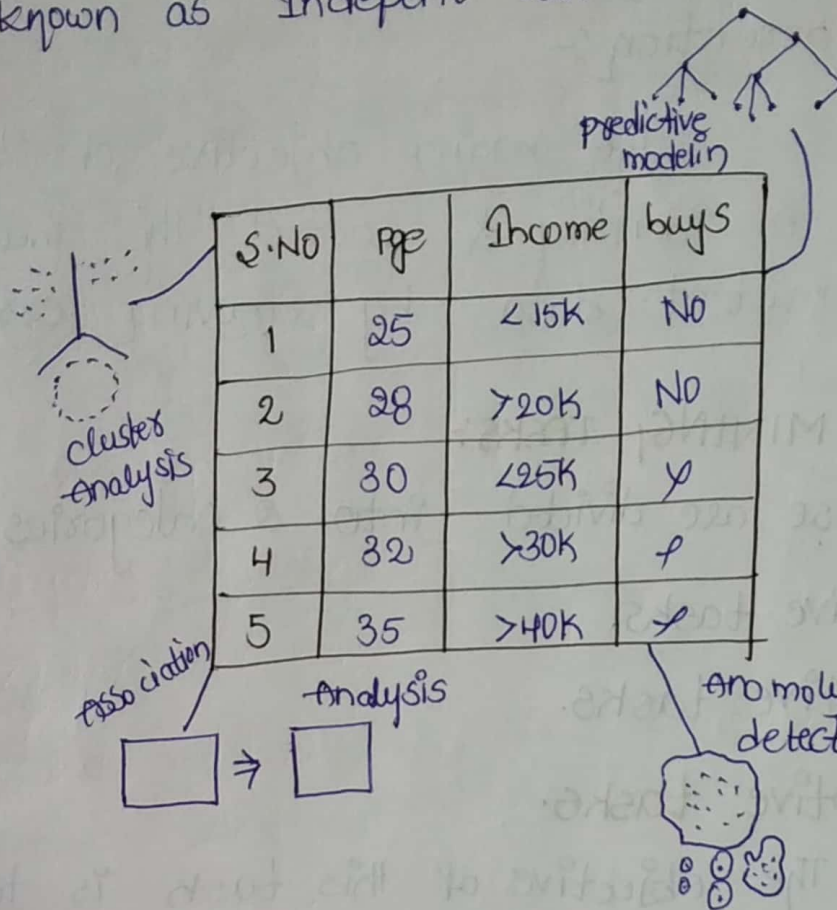
The objective of this task is to divide the pattern and [co-relation and clusters anomalies] thus summarizes the underlying relationship in data. It req Post processing techniques to validate the results.

2) predictive Task:

The objective of this task is to predict the value of the a particular attribute based on the value of the other attributes.

• Attribute i.e., to be predicted

is known as dependent (or) target variable and attributes used for making the prediction is known as Independent variable.

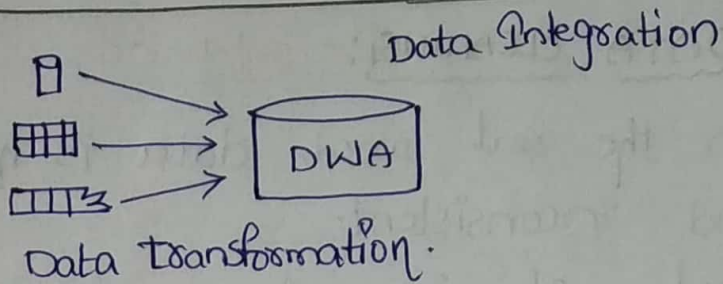


* Data Integration :-

When the data from the different multiple sources are combined from a coherent store.

Data Transforming :-

Where the data is converted from one form to the other data is normalized, aggregated and Generalized.



* Data Reduction :-

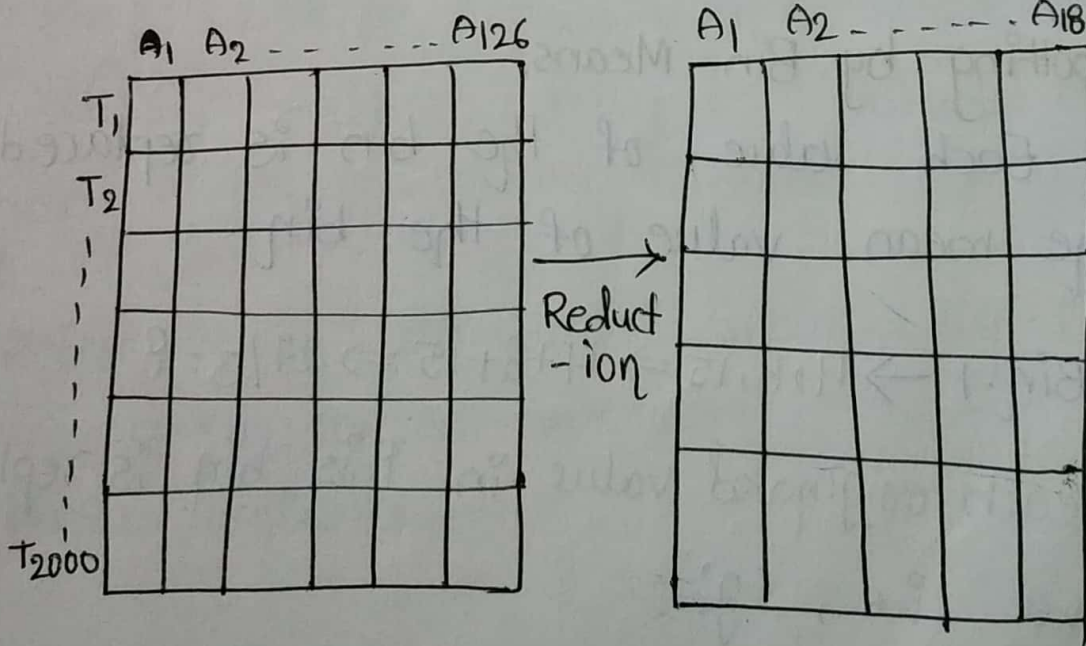
Which shows the reduced

representation of the data in the Dwtl i.e., smaller in volume but produces some analytical results.

Data Transformation
 $-2, 32, 100, 59, 48 \rightarrow 0.02, 0.32, 1.0, 0.59, 0.48$

→ Data preprocessing techniques can improve the quality of data this by improve the array's efficiency of mining process.

→ Data preprocessing is a important step in the KDD process, since quality divisions are based on quality on data.



DATA REDUCTION ATTRIBUTES.

* DATA CLEANING:

1) In the real world, data is noisy/incomplete and inconsistent.

2) Data cleaning is used to fill the missing values, smooth the noisy data and convert the inconsistencies.

1) MISSING VALUES:-

The following are the methods to fill the missing values of an attributes.

(a) Ignore the tuple:-

→ When the class label is missing

→ This method is not effective, if tuple contains several attributes with missing values.

(b) Smoothing by Bin Means:

Each value of the bin is replaced by the mean value of the bin.

Eg:- Bin.1 \rightarrow 4, 8, 15 \Rightarrow $4+8+15 \Rightarrow 27/3 = 9$

i.e., each originated value in this bin is replaced by value i.e., '9'.

Smoothing by Bin means.

* Data Reduction Techniques:-

These Techniques can be applied to obtain a reduced representation of the dataset that is much smaller in volume, yet closely maintains the integrity of the original data. Data reduction includes,

* Data cube aggregation

* Dimension reduction

* Data Compression

* Discretization

* Numerosity Reduction

→ Regression

→ Histograms

→ Clustering

→ Sampling

→ Data cube aggregation where aggregation operations are applied to the data in the construction of datacube.

→ Attribute subset selection, where irrelevant, weakly relevant or redundant attributes or dimensions may be detected and removed.

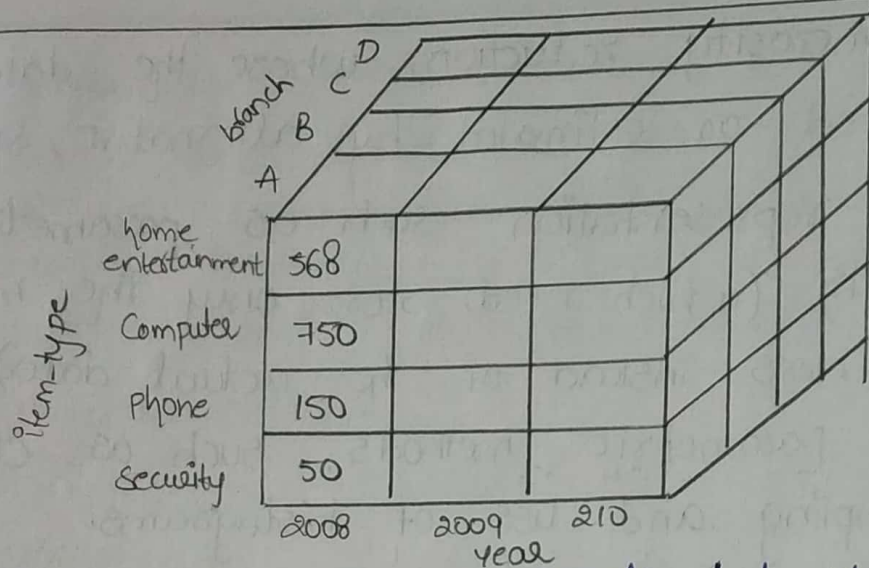
→ Dimensionality Reduction, where encoding mechanisms are used to reduce the data set size.

Examples :- Wavelet Transforms principal Components Analysis.

- Numerosity reduction, where the data are replaced or estimated by alternative, smaller data representation such as parametric models (which need, store only the model parameters instead of the actual data) or non parametric methods such as clustering, sampling and use of histograms.
- Discretization and concept hierarchy generation where raw data values for attributes are replaced by ranges or higher conceptual levels. Data discretization in the form of numerosity reduction that is very useful for the automatic generation of concept hierarchies.

* Data cube Aggregation:-

Reduce the data to the concept level needed in the analysis. Queries regarding aggregated information. The following figure shows a data cube for multidimensional analysis of status data with respect to annual sales per item type for each branch.



→ Each cells holds an aggregate data values corresponding to the data point in multidimensional space. Data cubes provide front access to pre computed, summarized data, there by benefiting m-line analytical processing as well as data mining.

→ The cube created at the lowest level of abstraction is referred to as the base cuboid. A cube for the highest level of abstraction is the apex cuboid. The lowest level of a data cube (base cuboid). Data cubes created for varying levels of abstraction are sometimes referred to as cuboids, so that a "data cube" may instead refer to a lattice of cuboids each higher levels of abstraction further reduces the resulting data size.

→ The following database consist of sales per quarter for the years 1997-1999

Year=1999	
Year=1998	
Year=1997	
Quarter	Sales
Q ₁	\$224,000
Q ₂	\$408,000
Q ₃	\$350,000
Q ₄	\$586,000

Year	Sales
1997	\$1,568,000
1998	\$2,356,000
1999	\$3,594,000

Suppose, the analyzer interested in the annual sales rather than sales per quarter, the above data can be aggregated so that the resulting data summarizing the total sales per year instead of per quarter. The resulting data is smaller in volume without loss of information necessary for analysis task.

* Dimensionality Reduction :-

It reduces the dataset size by removing irrelevant attributes. This is a method of attribute subset selection are applied. A heuristic method of attribute of subset selection is.

explained here:

Attribute sub selection/feature selection:

Feature selection is a must for any data mining product. That is because, when you build a data mining model, the dataset frequently contains a more information than is needed to build the model. For example, a data set may contain 500 columns that describe characteristics of customers, but perhaps 50 of those columns are used to build a particular model. If you keep the unneeded columns while building the model more CPU and memory are required during the training process, and more storage space is required for the completed model.

→ In which select a minimum set of features such that the probability of distribution of different classes given the values for those features is as close as possible to the original distribution given the values of all features.

→ Basic heuristic methods of attribute subset selection include the following techniques

some of which are illustrated below.

01) step-wise forward selection:-

The procedure starts with an empty set of attributes. The best of the original attributes is determined and added to the set. At each subsequent iteration or step, the best of the remaining original attributes is added to the set.

2) step-wise backward elimination:-

The procedure starts with full set of attributes. At each step, it removes the worst attribute remaining in the set.

3) Combination forward selection and backward elimination:-

The step wise forward selection and backward elimination methods can be combined, where at each step select the best attribute and removes the worst from among the remaining attributes.

4) Decision tree induction:-

It constructs a flow chart like structure where each internal (non leaf) node denotes a test on an attribute, each branch corresponds to overcome of the test, and each external (leaf) node denotes a class prediction.

At each node, the algorithm chooses the "best" attribute to partition the data into individual classes. When decision tree induction is used for attribute subset selection, a tree is constructed from the given data. All appearing in the tree form the reduced subset of attributes.

Forward Selection	Backward Elimination	Decision tree Induction
<p>Initial attribute set: $\{A_1, A_2, A_3, A_4, A_5, A_6\}$</p> <p>Initial reduced set: $\{\}$ $\Rightarrow \{A_1\}$ $\Rightarrow \{A_1, A_4\}$ \Rightarrow Reduced attribute set: $\{A_1, A_4, A_6\}$</p>	<p>Initial attribute set: $\{A_1, A_2, A_3, A_4, A_5, A_6\}$ $\Rightarrow \{A_1, A_3, A_4, A_5, A_6\}$ $\Rightarrow \{A_1, A_4, A_5, A_6\}$ \Rightarrow Reduced attribute set: $\{A_1, A_4, A_6\}$</p>	<p>Initial attribute set: $\{A_1, A_2, A_3, A_4, A_5, A_6\}$</p> <p>$\Rightarrow$ Reduced attribute set: $\{A_1, A_4, A_6\}$</p>

Fig:- Greedy methods for attribute subset selections.

* Wrapper approach / filter approach:-

The mining algorithm itself is used to determine subset then it is called wrapper approach or filter approach. wrapper approach leads to greater accuracy since it optimizes the evaluation measure of algorithm while removing attributes.

* Data Compression :-

In Data Compression, data encoding or transformations are applied so as to obtain a reduced or "compressed" representation of original data. If the original data can be reconstructed from the compressed data without any loss of information, the data compression technique used is called lossless. If instead, we can reconstruct only an approximation of the original data, then the data compression technique is called lossy. Effective methods of lossy data compression.

* wavelet transforms

* principal Component analysis.

* Wavelet transforms :- (Discrete wavelet transformation).

It is a form of data compression well suited for image compression. The discrete wavelet transform DWT is a linear signal processing technique that, when applied to a data vector D , transforms it to a numerically different vector, D_0 , of wavelet coefficients.

→ The general algorithm for a discrete wavelet transform is as follows.

- * The length L of the input data vectors must be an integer power of two. This condition can be met by padding the data vectors with zeros as necessary.
 - * Each transform involves applying two functions:
 - data smoothing
 - Calculating weighted difference.
 - * The two functions are applied to pairs of the input data, resulting in two sets of data of length $L/2$.
 - * The two functions are recursively applied to the sets of data, obtained in the previous loop, until the resulting data sets obtained are of desired length.
 - * A selection of values from the data set obtained in the above iterations are designated the wavelet coefficients of transformed data.
→ If wavelet coefficient are larger than some as a specified threshold then it can be retained. The remaining are set to zero.
- Haar and Daubechies are two popular wavelet transforms.

* Principal Component Analysis (PCA) :-

It is also called as Karhunen-Loeve (K-L) method.

Procedure :-

- (give N data vectors from k dimensions, find $CL = k$ orthogonal vectors that can be best used to represent data.
 - The original data set is reduced (projected) to one consisting of N data vectors on c principal components (reduced dimensions)
 - Each data vector is a linear combination of the c principal component vectors.
 - Works for ordered and unordered attribute
 - Used when the no. of dimensions is large.
- The principal components new set of axes give important about variance. Using the strongest components one can reconstruct a good approximation of the original signal.

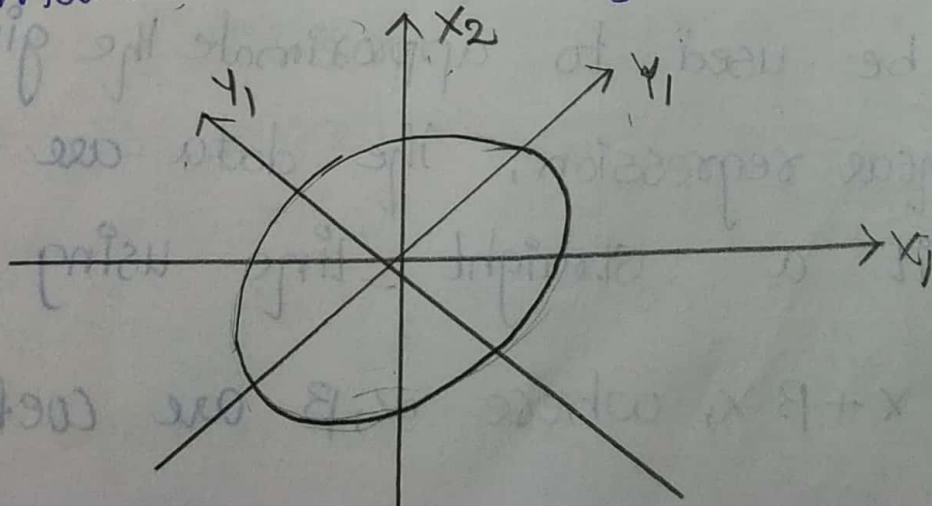


fig: principal Component analysis. Y_1 & Y_2 are the first two principal Components for the given data.

* Numerosity Reduction:

Data volume can be reduced by choosing alternative smaller forms of data. This tech. can be

* parametric method

* Non parametric method.

* parametric method:-

Assume the data fits some model, then estimate model parameters and store only the parameters, instead of actual data.

Non parametric Method: In which histogram, clustering and sampling is used to store reduced form data.

Numerosity reduction techniques:

01) Regression and log linear model:

→ Can be used to approximate the given data.

→ In linear regression, the data are modeled to fit a straight line using

$$Y = \alpha + \beta X, \text{ where } \alpha, \beta \text{ are coefficients}$$

Multiple regression:

$$Y = b_0 + b_1 x_1 + b_2 x_2$$

Many non linear functions can be transformed into the above.

log linear model:

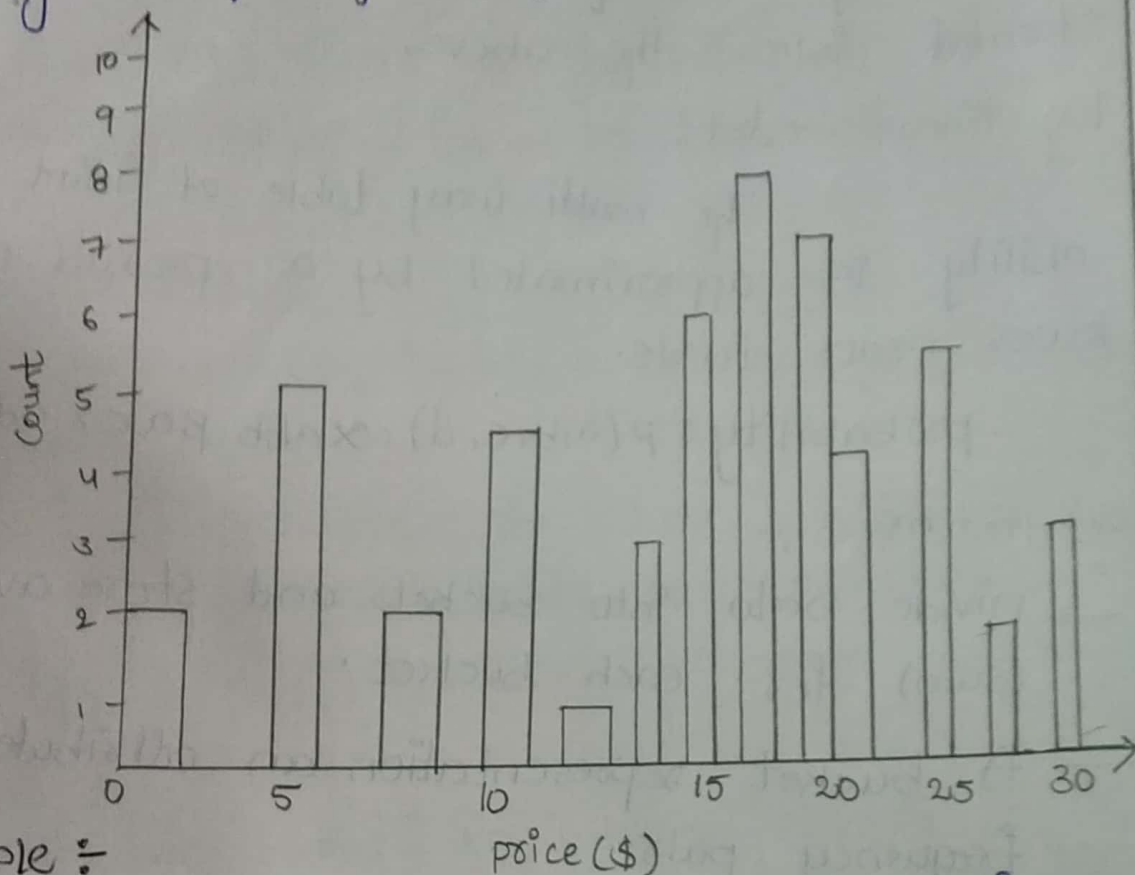
The multi way table of joint probability is approximated by a product of lower order table.

$$\text{probability: } P(a, b, c, d) = \alpha a b b \beta a c \times a d \gamma b c d$$

02) Histogram:-

- Divide data into buckets and store average (sum) for each bucket.
- A bucket representation an attribute value frequency pair.
- It can be constructed optimally in one dimension using dynamic programming.
- It divides up the range of possible values in a data set into class or groups. For each group, a rectangle bucket is constructed with a base length equal to the range of values in that specific group, and an area proportional to the no. of observation falling into that group.

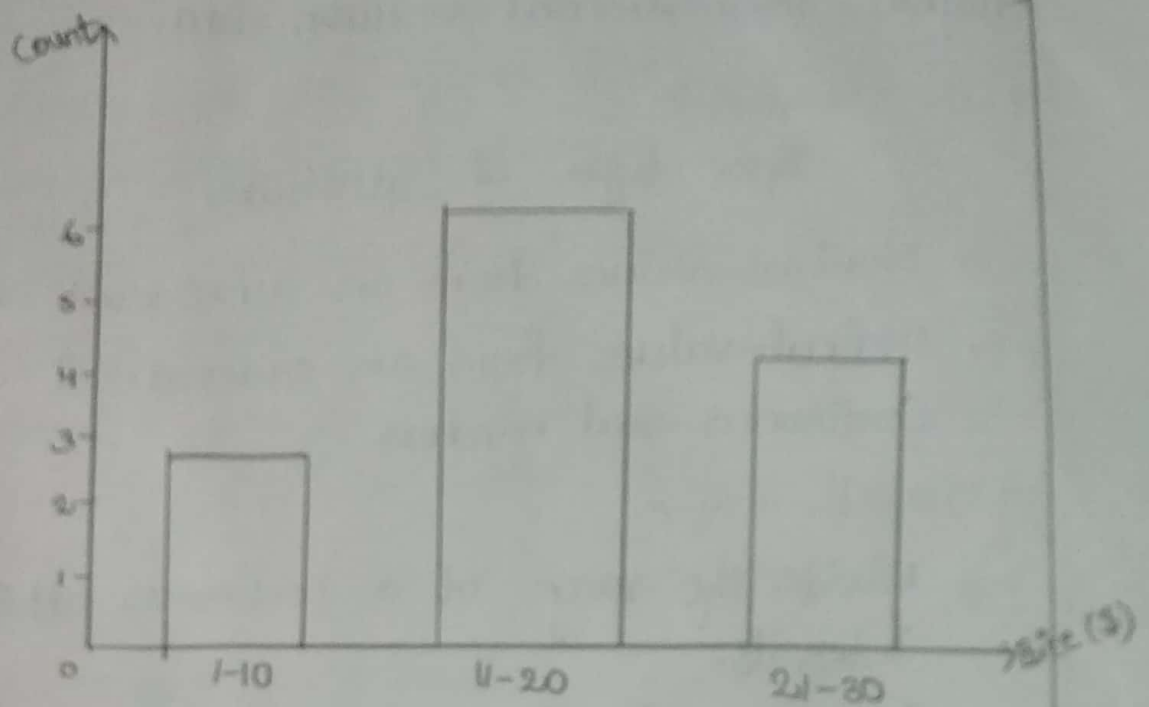
→ The buckets are displayed in a horizontal axis while height of a bucket represents the average frequency of the values.



Example :-

The following data are a list of prices of commodity of sold items. The numbers have been sorted. 1, 1, 5, 5, 5, 5, 5, 8, 8, 10, 10, 10, 10, 12, 14, 14, 14, 15, 15, 15, 15, 15, 15, 15, 18, 18, 18, 18, 18, 18, 18, 20, 20, 20, 20, 20, 21, 21, 25, 25, 25, 25, 25, 28, 28, 30, 30, 30.

The buckets can be determined based on the following partitioning rules, including the following.



1) Equi-width:-

Histogram with bars having the same width.

2) Equi-depth:-

Histogram with bars having the same height.

3) V-optimal:-

Histogram with least variance \square (count values).

4) Max diff: Bucket boundaries defined by user specified threshold.

V-optimal and Max Diff histogram tend to be the most accurate and practical. Histograms are highly effective at approximating both sparse and dense data, as well as

highly skewed, and uniform data.

* Discretization :-

Three types of attributes

- Nominal-values from an unordered set
- Ordinal-values from an ordered set
- Continuous-real numbers.

* Discretization :-

- Divide the range of a continuous attribute into intervals.
- Some classification algorithms only accept categorical attributes.
- Reduce data size by discretization
- prepare for further analysis.

Discretization and Concept hierarchy.

Discretization

- Reduce the number of values for a given continuous attribute by dividing the range of attribute into intervals. Interval labels can then be used to replace actual data values.
- Concept hierarchies.
- Reduce the data by collecting and replacing

low level concepts (such as numeric values for the attribute age), by higher level concepts (such as young, middle age or seniors).

Discretization and Concept hierarchy generation for numeric data.

- Binning (see sections before)
- Histogram analysis (see sections before)
- clustering analysis (see sections before)
- Entropy based discretization
- Segmentation by natural partitioning

Entropy based Discretization :-

- Given a set of samples S , if S is partitioned into two intervals S_1 & S_2 using boundary T , the entropy after partitioning is

$$E(S, T) = \frac{|S_1|}{|S|} \text{Ent}(S_1) + \frac{|S_2|}{|S|} \text{Ent}(S_2)$$

- The boundary that minimizes the entropy function over all possible boundaries is selected as binary discretization.
- The process is recursively applied to partitions obtained until some stopping criterion

is met, e.g.

$$E(n(s)) - E(T, s) > \delta$$

→ Experiments show that it may reduce data size and improve classification accuracy.

* Linear Regression:

Data are modeled to fit a straight line

→ Often uses the least square method to fit the line

$$Y = \alpha + \beta X$$

→ Two parameters, α and β specify the line and all to be estimated by using the data at hand.

→ Using the least squares criterion to the known values of $Y_1, Y_2, \dots, X_1, X_2, \dots$

* Multiple Regression:-

It allows a response variable Y to be modeled on a linear function of multidimensional feature vectors.

$$Y = b_0 + b_1 X_1 + b_2 X_2$$

→ Many non linear functions can be transformed into the above.

* log-linear model :-

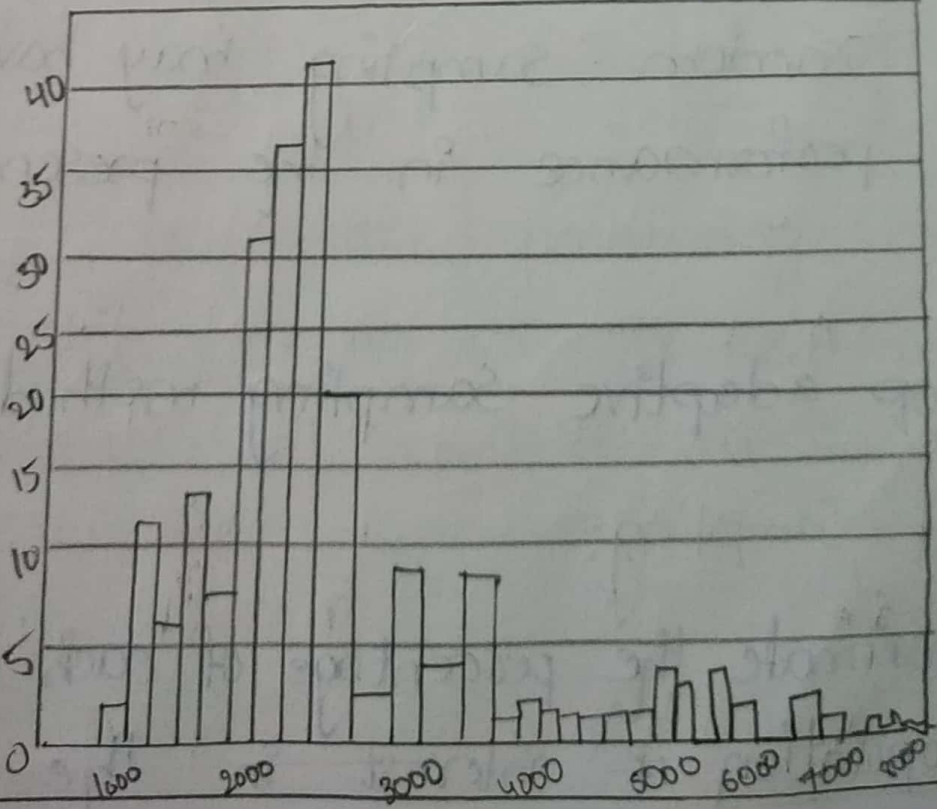
Approximates discrete multi dimensional probability distribution.

→ The multi way tables of joint probabilities is approximated by a product of lower order tables

Probability : $p(a,b,c,d) = \alpha_{abc} \alpha_{abd} \alpha_{bcd}$

* Histograms :

- A popular data reduction technique
- Divide data into buckets and store average sum for each bucket.
- Can be constructed optimally in one dimension using dynamic programming.
- Related to quantization problem.



* clustering :-

- partition data set into clusters, and one can store cluster representation only.
- can be very effective if data is clustered but not if data is "smeared".
- Can have hierarchical clustering and be stored in multi dimensional index tree structure.
- There are many choices of clustering definitions and clustering algorithms, further detailed in chapter 8.

* Sampling :-

- Allow a mining algorithm to run in complexity that is fractionally sub linear to the size of data.
- Choose a representative subset of the data
- Simple random sampling may have very poor performance in the presence of skew.

Develop adaptive sampling methods.

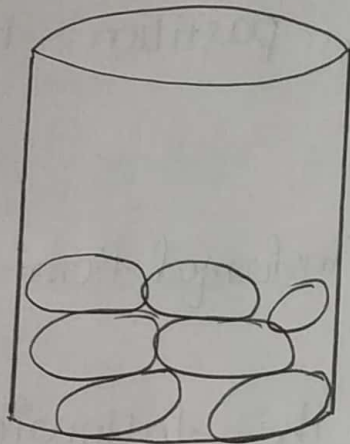
Stratified Sampling:

- Approximate the percentage of each class / sub population of interest in the overall

database.

- used in conjunction with skewed data.
- Sampling may not reduce database I/O's (Page at a time)

Sampling:

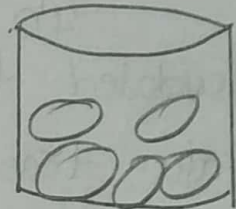


Raw data

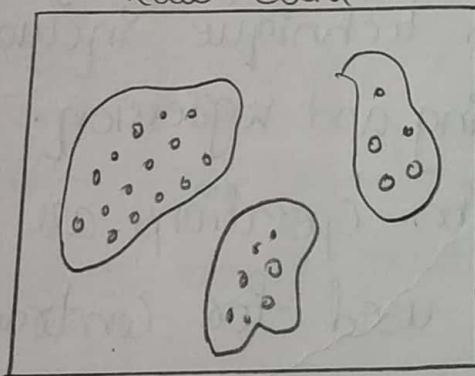
SR SWR
(Simple random sample without replacement)



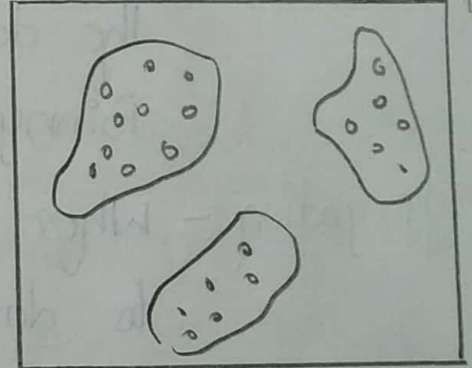
SRSWR



Raw data



cluster / stratified sample



segmentation by natural partitioning:

3-4-5 rule can be used to segment numeric data into relatively uniform "natural" intervals.

→ If an interval covers 3, 6, 7, 9 distinct values

at the most significant digit, partition the range into 3 equiwidth intervals.

→ If covers 2, 4, 8 distinct values at the most significant digit, partition range into 4 intervals

→ If it covers 1, 5 or 10 distinct values at the most significant digit, partition the range into 5 intervals.

DATA TRANSFORMATION :-

In this data is transformed from 1 consolidated form to other.

→ Data transformation involves the following methods.

1) Smoothing :- It is used to remove the noisy from the data such technique include Binny, clustering and regression.

2) Aggregation :- Where Aggregation operations are applied to data i.e., used for constructing the data cubes at multiple levels.

3) Generalization :- Where the low level data are replaced by high level concept using concept hierarchies.

4) Normalization :- Where the attribute data falls within a small specific range

such as 0.0 to 1.0.

5) Attribute construction:-

- Where the new attributes are constructed and added given from the given set of attributes.
- Smoothing is a form of data cleaning.
- Aggregation and Generalization are the forms of Data Reduction.
- This method is useful

* Normalization :-

The normalization strategy is also known as standardization helps us to avoid dependence on the choice of measurement units. As the measurement unit can effect the data analysis.

- When we try to express an attribute in smaller units (metres to inches) kg-pound etc..). It will lead to larger range for that tending to give attribute greater effect or weight.
- Normalization involves transforming the data to fall within a smaller (or) common range
Ex:- $[-1, 1]$ $[0.0 - 1.0]$.

Methods for Normalization :-

- 1) Minimum Maximum normalization
- 2) z-score (or) mean normalization
- 3) Decimal scaling (or) decimal normalization.

* Min - Max normalization :-

Let us assume an attribute 'A' with minimum A, maximum A and all maximum values of attribute respectively.

→ Let the attribute A have the values $V_1, V_2, V_3, \dots, V_i$

$$V_i' = \frac{V_i - \text{Min}_A}{\text{Max}_A - \text{Min}_A} (\text{new Max} - \text{new Min}_A) + \text{new Min}_A$$

* z-score normalization :-

When min & Max values are not known.

$$V_i' = \frac{V_i - \text{mean}_A}{\text{std. dev}^A}$$

* decimal scaling (or) decimal normalization :-

We convert the values in decimals by

$$V_i' = \frac{V_i}{10^j}, \text{ where } j \geq 1$$

* Measures of similarity & dissimilarity.

Object in cluster are similar to one another & dissimilar to two objects in another cluster.

Similarity & Dissimilarity measures are known measures of proximity.

Similarity measure for two objects i, j is equal to '0' then the objects are unlike.

The higher is the value of similarity measure the greater is the similarity b/w 2 objects.

Similarity measure is equal to '0'. means the objects are like (similar).

We commonly use two data by structure for representing similarity & dissimilarity.

1) Data matrix

2) dissimilarity matrix.

Data Matrix:-

Also known as object by attribute structure. It stores 'n' data objects with p attributes in $n \times p$ matrix.

Example:-

$$\begin{bmatrix} x_{11} & \dots & x_{1p} & \dots & x_{1P} \\ \vdots & & & & \\ x_{i1} & \dots & x_{ip} & \dots & x_{iP} \\ \vdots & & & & \\ x_{n1} & \dots & x_{np} & \dots & x_{nP} \end{bmatrix}$$

Dis-similarity Matrix :- It is object by object structure represented by $d(i,j)$ and similarity is called form dissimilarity as $s(i,j) = 1 - d(i,j)$ representation.

$$\begin{bmatrix} 0 & & & & \\ d(2,1) & 0 & & & \\ d(3,1) & d(3,2) & 0 & & 1 \\ \vdots & \vdots & \vdots & \dots & \vdots \\ d(n,1) & d(n,2) & \vdots & \dots & 0 \end{bmatrix}$$

$d(i,j) \Rightarrow$ diff b/w obj i, j

$d(i,j)$ is non-negative number

$$d(i,i) = 0$$

$$d(i,j) = d(j,i)$$

Proximity measure for nominal attribute :-

Example :- young senior middle

dissimilarity $d(i,j)$ of nominal attributes

$$= \frac{P-m}{P}$$

P = total no. of matches

m = no. of similar matches.

$$S(i, j) = \frac{1-P-m}{P} \Rightarrow \frac{m}{P}$$

Proximity measure for Binary attributes:-

Binary data has values 0 or 1.
 $d(i, j)$ for binary data is represented in

Contingency Matrix.

Contingency Matrix:-

		0_j	1_j	sum
1_i	q	x	$q+x$	
0_i	s	t	$s+t$	
sum	$q+s$	$x+t$	P	

q : $i=1, j=1$

x : $i=1, j=0$

s : $i=0, j=1$

t : $i=0, j=0$

total no. of attributes $p = q + r + s + t$.
 dissimilarity that is based on symmetric
 binary attributes is called Symmetric
binary dissimilarity.

If the two states are equally impor-
 -tant then they are termed as symm-
 -etric binary attributes.

$$\text{Dissimilarity } d(i, j) = \frac{r+s}{q+r+s+t}$$

$$\text{Similarity } s(i, j) = 1 - d(i, j) = \frac{q}{q+r+s}$$

↓
 which is equal to
 Saccard Co-efficient.

for asymmetric binary attributes, the
 dissimilarities termed as asymmetric binary
similarity where 't' is ignored.

$$d(i, j) = \frac{r+s}{q+r+s}$$

Dissimilarity for Numeric Data:-

It is calculated by euclidean,
 manhattan & minkowski distance. The most

popular is euclidian distance where.

$$i = x_{i1}, x_{i2}, x_{i3}, \dots$$

$$j = x_{j1}, x_{j2}, x_{j3}, \dots$$

$$d(i, j) = \sqrt{(x_{i1} - x_{j1})^2 + (x_{i2} - x_{j2})^2 + \dots}$$

Manhattan Distance :- $d(i, j) = |x_{i1} - x_{j1}| + |x_{i2} - x_{j2}| + \dots + |x_{ip} - x_{jp}|$

Minkowski Distance :- It is generalized form of Euclidean & Manhattan where $d(i, j) = h$.

$$h = \sqrt{|x_{i1} - x_{j1}|^h + |x_{i2} - x_{j2}|^h + \dots + |x_{ip} - x_{jp}|^h}$$

where $h \geq 1$

Eg real number.

* Short Answer Questions:-

1) What is data mining?

Data mining is defined as extracting information from huge sets of data. In other words, we can say that data mining is the procedure of mining knowledge from data. The information or knowledge extracted so can be used for any of the following applications.

- 1) Market Analysis
- 2) Fraud detection
- 3) Customer Retention
- 4) Production Control
- 5) Science Exploration.

2) Difference b/w data mining & OLAP?

Data mining	OLAP
1) Bottom up, discovery-driven	1) Top-down, query-driven
2) Requires no assumption	2) Repetitive testing of user originated theories.
3) No intensive human interaction with the db is required.	3) Requires a great deal of human interaction with the data base.

4) Runs mostly automatically

4) users must have a good idea about the information which she/he is looking for.

5) User interaction is limited to the selection of data mining alg & the selection of appropriate parameters.

5) User is in permanent interaction with the system.

6) Why is this happening? \leftrightarrow Is this true?
And.

What might happen if--?

3) What are the types & tasks that arrived during data mining?

1) \rightarrow The set of task-relevant data to be mined.

\rightarrow The kind of knowledge to be mined.

\rightarrow The background knowledge to be used in discovery.

\rightarrow The interestingness measures & thresholds for pattern evaluation.

\rightarrow The expected representation of visualizing the discovered pattern.

4) What is the need for preprocessing of data?

A) In Complete noisy & in-consistent data are common place properties of large real world databases & data ware houses.

Incomplete data can occur for a no. of reasons. Attributes of interest may not always be available, such as customer information for sales transition data.

Other data may not be included simply because it was not considered important at the time of entry relevant data may not be recorded due to misunderstanding.

Data may have been overlooked. Missing data, particularly for tuples with missing values for some attributes, may need to be inferred.

5) What say knowledge discovery?

Some people treat data mining some as knowledge discovery while others view data mining as an essential step in the process of knowledge discovery. Here is the list of steps involved in the knowledge discovery process.

- * Data cleaning
- * Data integration
- * Data transformation
- * Data mining
- * pattern Evaluation
- * Knowledge presentation.

Q) What is dimensionality reduction?

In dimensionality reduction, data encoding or transformation are applied as to obtain a reduced or "compressed" represent of the original data. If the original data can be reconstructed from the compressed data without any loss of information, the data reduction is called lossless.

Q) What are the data preprocessing tasks?

- * Data cleaning
- * Data Integration
- * Data Transformation
- * Data reduction.

What is data cleaning?

Data cleaning is a technique that is applied to remove the noisy data & correct the inconsistency in data. Data

Cleaning involves transformations to correct the wrong data. Data cleaning is performed on a data preprocessing step while preparing the data from a data warehouse.

9) Issues during data Integration?

A) The issues to be considered during data integration are

(i) Matching

(ii) Redundancy

(iii) Detection & resolution of data value-conflicts.

10) strategies of data reduction?

A) (i) Data cube aggregation

(ii) Attribute subset selection

(iii) Dimensionality reduction

(iv) Numerosity reduction

(v) Discretization & concept hierarchy generation.

11) How do you clean the data? for missing values.

A) (i) Ignore the tuple

(ii) Fill in the missing value manually

(iii) Use a global constant to fill in the missing value

(iv) Use the attribute means to fill in the missing value.

(v) Use the attribute mean for all samples belonging to same class as given tuple.

(vi) Use most probably value to fill in the missing value. for Noisy data:-

1) Binning

2) Regression

3) Clustering

12) What is concept hierarchy & how it is useful in data mining?

A) A concept hierarchy for a given numerical attribute defines a discretization of the attribute. Concept hierarchies can be used to reduce the data by collecting & replacing low-level concepts with higher level concepts although detail is lost by such data. generalization, the generalized data may be more meaningful & easier to interpret.

What is data generalization?

It is process that abstracts a large set of task relevant data in a database from a relatively low conceptual levels

to higher conceptual levels & approaches for Generalization.

1) Data cube approach

2) Attribute oriented induction approach.

14) List the primitives for specifying a data mining table.

A) → The set of task-relevant data to be mined

→ The kind of knowledge to be mined

→ The background knowledge to be used in the discovery process.

→ The interestingness measures & threshold for pattern evaluation.

→ The respected representation for visualizing the discovered pattern.

15) What are the steps involved in KDD process.

A) 1) Developing & understanding of

→ The application domain

→ The relevant prior knowledge

→ The goals of the end users.

2) Creating a target data set

3) Data cleaning & pre processing

4) Data reduction & projection

5) Choosing the data mining tasks
 6) choosing the data mining alg.

7) Data mining

8) Interpreting mined patterns

a) Consolidating discovered knowledge.

16) Mention some of data mining techniques.

There are several ~~many~~ as major data mining techniques have been developing & using in data mining projects recently. Including association, classification, clustering, prediction, sequential patterns.

17) What is meta learning?

Meta learning is a sub field of Machine learning where automatic learning algorithms are applied on meta data about machine learning experiments.

18) Define summarization?

Summarization is a key data mining concept which involves techniques for finding a compact description of a data sets. ... i.e., data visualization of automated report generation. Clustering [13, 23] is another data mining techniques that is often used to summarize

Large datasets.

19) Classification of datamining system?

Datamining is an interdisciplinary field, the confluence of a set of disciplines, including database systems, statistics, machine learning, visualization & information science. Data mining system can be categorized a/c to various criteria.

20) Challenges in Data mining?

A) *Introduction.

* Noisy & Incomplete data

* Distributed data

* Complex data

* performance

* Incorporation of Background knowledge

* Data visualization.

* Data privacy & security.

21) Identify any 3 functionality of Data Mining?

The functionality of Data Mining are:

➤ Concept/class Description: Data characterization & Data Discrimination

2) Mining freq patterns, Associations & Correlations.

3) Association analysis.

22) Interpret major issues of Data mining?

The major issues of Data mining are:-

- (i) Limited and irrelevant information
- (ii) Noisy and irrelevant information
- (iii) Human interaction & knowledge
- (iv) Large data sets and high dimensionality
- (v) Uncertainty
- (vi) Dynamic updates.

23) Discuss relational databases?

A relational databases (RDB) is a collective set of multiple data sets organized by tables, records and columns. RDBs establish a well-defined relationship b/w database tables. Tables communicate and share information, which facilitates data searchability, organization and reporting.

RDBs use structured Query Language (SQL), which is a standard user application that provides an easy programming interface for database interaction.

24) state object-oriented databases?

An object-oriented database is a system

offering database management facilities in an object-oriented programming environment. Data is stored as objects and can be interpreted only using the methods specified by its class. An object class is a set of objects that share a common structure and a common behavior.

An object-oriented database is fundamentally created around an object-oriented data model.

25) Explain spatial databases?

A spatial database is a database that is optimised for storing and querying data that represents objects defined in a geometric space. Most spatial databases allow the representation of simple geometric objects such as points, lines and polygons. Some spatial databases handle more complex structures such as 3D objects, topological coverages & linear networks.

26) Explain the outlier analysis?

A database may contain data objects that do not comply with the general behavior or model of the data objects, which are grossly different from or inconsistent

with the remaining set of data, are called outliers.

→ Most data mining methods discard outliers as noise or exceptions.

27) Express what is a decision tree?

→ A decision tree is a flow-chart-like tree structure, where each internal node (non leaf node) denotes a test on an attribute, each branch represents an outcome of the test, & each leaf node holds a class label. The top most node in a tree is the root node.

28) List the types of data that can be mined?

→ database data.

→ data warehouse data

→ transactional data.

→ data streams.

→ Spatial data

→ engineering design data

→ hypertext and multi-media data.

→ graph and networked data.

29) Describe transactional data bases?

In general, each record in a transactional database captures a transaction, such as a customer's purchase, a flight booking, or a user's click on a

web page. A transaction typically includes a unique transaction identify number (trans-ID) and a list of the items making up the transaction; such as the items purchased in the transaction. A transactional database may have additional tables, which contain other information related to the transactions, such as item description, information about the salesperson or the branch, and so on.

30) Define multidimensional data mining?

Multidimensional data mining is an approach to data mining that integrates OLAP-based data analysis with knowledge discovery techniques.

→ It is also known as exploratory multidimensional data mining and online analytical mining (OLAM).

→ It searches for interesting patterns by exploring the data in multidimensional space.

31) Define Data cube?

A data cube is defined as a collection of dimensions & facts. Data cube allows the data to be modeled & viewed in multi-dimensions.

Ex:- 2D view.

32) Define data characterization?

Data characterization is a summarization of the general characteristics or features of a target class of data.

→ The data corresponding to the user-specified class are typically collected by a query.

33) Contrast heterogeneous databases and legacy databases.

A) Heterogeneous databases:-

→ A heterogeneous database consists of a set of interconnected, autonomous component databases.

→ The components communicate in order to exchange information and answer queries.

→ Component DB may differ greatly, hence difficult to assimilate the information.

Legacy databases:-

→ A legacy database is a group of heterogeneous databases that combines different kinds of data systems, such as relational (or) object-oriented databases, hierarchical DBs, N/W DBs, spreadsheets, multimedia DBs or file systems.

34) Differentiate classification & prediction?

Classification

→ Classification is the process of identifying to which category, a new observation belongs to on the basis of a training data set containing observations whose category membership is known.

→ In classification, the accuracy depends on finding the class label correctly.

→ In classification, the model can be known as the classifier.

→ A model or the classifier is constructed to find the categorical labels.

Prediction

→ Prediction is the process of identifying the missing or unavailable numerical data for a new observation.

→ In prediction the accuracy depends on how well a given predictor can guess the value of a predicted attribute for a new data.

→ In prediction, the model can be known as the predictor.

→ A model or a predictor will be constructed that predicts a continuous-valued function or ordered value.