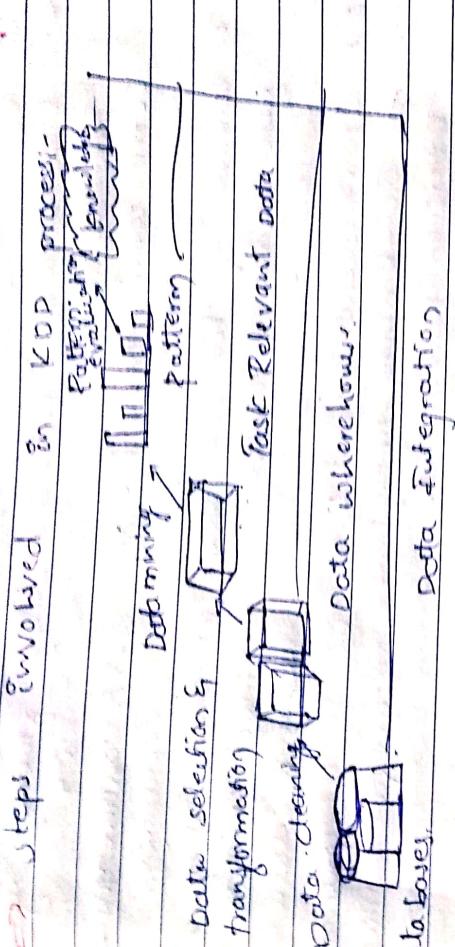


Describe the steps involved in mining discovered knowledge.



- (1) Data Cleaning: The noise & inconsistent data is removed.
- (2) Data Integration: multiple Data Source are combined.
- (3) Data Selection: Data relevant to the analytical task are retrieved from database.
- (4) Data transformation:- Data is transformed or consolidated into forms appropriate for mining by performing Summary or aggregating operation.
- (5) Data mining:- Intelligent methods are applied in order to extract data patterns.

08:53 - 08:55 2020-5-5

- (6) Pattern Evaluation: Data patterns are evaluated
- (7) Knowledge presentation: Visualization of Knowledge

Pr. representation techniques are used to present mind knowledge to user.

Q. Different type of business processing

A) DTP (Data Processing) -
DTP contains analytical for DLT P routine transaction processing.

It consists of historical data. ① consists of operational data from various businesses.

Current data

It is subject oriented. Used for data mining, analytics, used for business tasks, precision marketing etc.

2) The data is used in planning problem solving & decision making

③ The data is used to perform day to day fundamental operations.

a) It reveals a snapshot ④ It provides a multi dimensional view of present business tasks

different business task.

5) Large amount of data is stored typically in TB.PB. relatively small as the historical data is added

Ex: GB, MB.

6) Relatively slow as the amount of data involved is large. ⑤ very fast as the queries may take hours. S.I. of data

⑥ It only need backup from ⑦ backup & recovery time to time as compared process is maintained to DTP
→ need regularly.

2020-5-5 08:52

- ③ This data is generally managed by CEO, mgm. , my clerks, manager
- ④ Only read & rarely write
- ⑤ Both read & write operation

- ⑥ Explain star schema with an example of sales

- In the star schema, the center of the star can have fact table & a number of associated dimension tables . It is known as star schema as its structure resembles a star. The star schema in the simplest type of data ware house schema it is also known as star join schema & is optimized for querying large data sets.
- Each dimension in a star schema is represented with only one dimension table .
- * This dimension table contains the set of attributes
- * The following diagram shows the sale data of a company with respect to the four dimensions , namely time, item, branch, & location .
 - * There is a fact table ~~hosted~~ at the center . It contains the keys to each fact dimension .
 - * The fact table also contains the attribute namely dollars sold & units sold .

location
dimension - from
dimensions of the location - very
branch : tree
branch name
branch-type
province or state
country.

2020-5-5 08:52

Q) write a note on meta-rule guided mining

→ "How are metarules useful?" meta rules allow user to specify the syntactic form of rules that they are interested in mining. The rules forms can be as constraints. To help improve the efficiency of the mining process metarule less may be based on the analysis experience, expectations, or intention regarding the data or may be automatically generated based on the data base schema.

Meta-rule-guided mining

Suppose that as a market analyst for All Electronics, you have access to the data describing customer (such as customer age, address, & credit rating) as well as the list of customer transactions. You are interested in finding correlations between customer traits & the items that customers buy. However, rather than finding all of the association rules reflecting these relationships

You can handle the data in following ways

- pure & random data
- labeled data
- specialized data (automobile, banking etc.)
- form of data you want
- at last such a system will be able to take care of itself

where it is to be predicted whether
there are substantial risk elements in
the given situation. This kind of
problem is suitable especially if decision
and its consequences are both of two
categories only i.e. it is going to happen
Typically it can be solved by first
of all to be considered for problem
situation with respect to information &
a sufficient set of many test cases.

Classification is a process of building a model
process of automatically selecting the
production using information from database
classification

=> classification

It is a process of building a model
that describes the data & makes the computer
the process is to be able to make this
model to predict the class of object which
class label is unknown. This learned model
is based on the quantity of data &
classifying data.

2020-5-5 08:52

Information Gain: This measure is based on pioneering work by Claude Shannon on information theory, which studied the value of information content of message. Let node N represent the tuple with the highest information gain is chosen as the splitting attribute for the node N . The expected information needed to classify a tuple in D is given by

$$\text{Given by } \text{Info}(D) = - \sum_{i=1}^n p_i^e \log_2 p_i^e$$

Where p_i^e is the probability that an arbitrary tuple in D belongs to class C_i & is estimated by $\frac{|C_i|}{|D|}$. Info(D) is the average amount of information needed to identify the class label of a tuple in D . Info(D) is also known as the entropy of D . The expected information required to classify a tuple from D , based on the partitioning by attribute A is calculated by.

Therefore $\text{Info}(D) = \sum_j |D_j| \cdot \text{Info}(D_j)$

Information gain is defined as the difference b/w the original information requirement & information obtained after partitioning on A .

$$\text{Gain}(A) = \text{Info}(D) - \text{Info}(D_A)$$

2020-5-5 08:51

change their location step by step
to one may notice that the counters
been generated, as a result of this loop
as the record now could. A loop has
done before. If we some data set points
like new, initial, a boundary. Not to be
this previous step will have direct
boundary of the cluster resulting from
the need to be calculate k new uninitially
early group stage or alone if the point
is passing this step is completed and an
it to the next could when no point
point belonging to a given data set of which
will also the next step is to take each
as well as possible. For a given point each
will be the latter cluster is the place from
be cause of different condition cannot different
which is placed in a boundary. So
case one for each cluster there are others
when the main idea is to define a
cluster a given data set through a certain
the possible follows a sample of such way
also the well known clustering problem
similarly generated learning algorithm that
new step is due to the
new clustering algorithm.
What effect the cluster is equal to
which is the process of forming a group

of points in a region of space.

2020-5-5 08

until no more changes are done or in other words until do not move any more finally this algorithm aims at minimizing an objective function known as squared error function given by

$$J(v) = \sum_{i=1}^n \sum_{j=1}^m (x_i^j - v_j)^2$$

where $\|x_i^j - v_j\|^2$ is the Euclidean distance between x_i^j & v_j
 x_i^j is the j^{th} data point in i^{th} cluster
 v_j is the j^{th} cluster center.

Algorithm steps for k-means clustering:-

let $x = (x_1, x_2, \dots, x_n)$ be the set of data points and $v = (v_1, v_2, \dots, v_k)$ be the set of centers.

- 1) Randomly select k cluster centers.
- 2) Then calculate the distance b/w each data point and cluster centers.
- 3) Assign the data point to the cluster center whose distance from the cluster is minimum of all the cluster centers.

Recalculate the new cluster center using

$$v_i = \left(\frac{c_i}{\sum_{j=1}^k c_i} \right) \sum_{j=1}^k x_i^j$$

where c_i represents the no. of data points in cluster

2020-5-5 08:51

- ④ No data point was reassigned from one cluster to another repeat from step 3.

Scalability: the general mechanics for obtaining reliable classification accuracy evaluates

Random partition into train and test part:

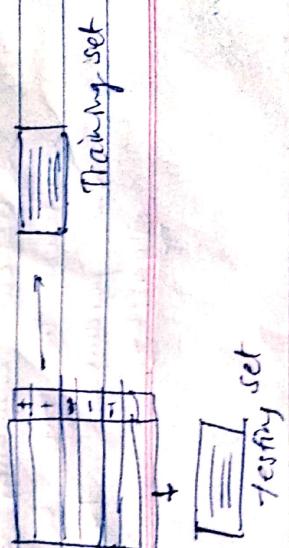
- ① Hold out
 - + use this independent data set, e.g. training set (2/3) test set (1/3) random Sampling
- ② Repeated hold out.
 - a k-fold cross-validation.
 - b randomly divide data set into k subsamples
 - c use $k-1$ subsamples as training data of one solo sample as test data - repeat k times
 - d leave-one-out for small size data.

Evaluation on "large" data, hold-out.

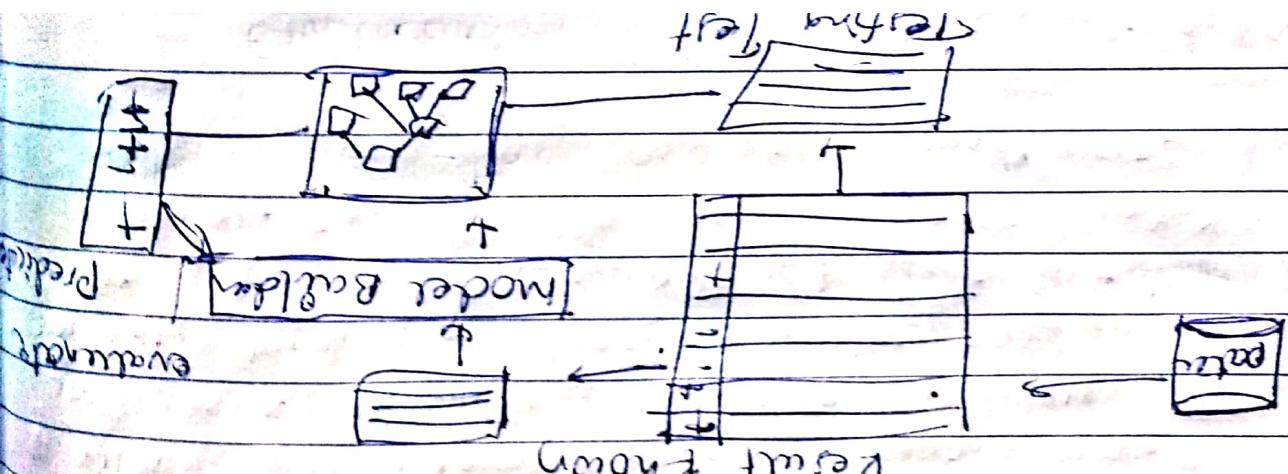
- Simple evaluation is sufficient.
- Randomly split data into training & test sets (usually 2/3 for train, 1/3 for test)

- ④ Build a classifier using the train set of evaluate it using the test set.

- Steps: Split data into train and test sets

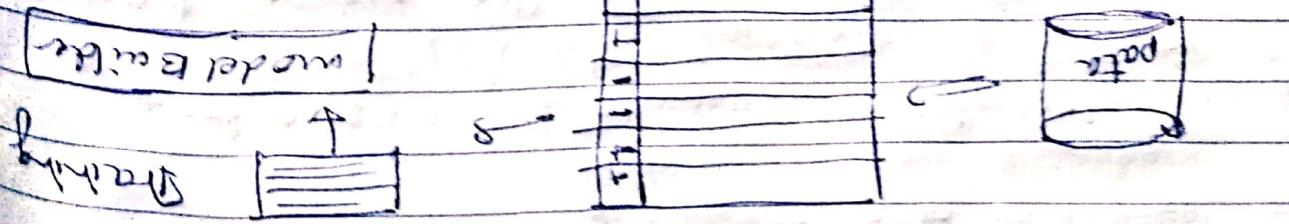


2020-5-5 08:51



Step 1: Build a model on a training set.

Training Set



Step 2: Build a model on a training set

$$Dy = \frac{1}{3} (x^3 - y^3)$$

Binary Distance

- cat of categorical variables file data set.
 - 0 & 1 when there is a mixture of num
 - numerical data when all the numerical variables are mixed.
 - cat variables file having up the issue of
 - cat numerical variables for the future of category
 - different measure we can only valid for
 - share also be noted that all share

$$\frac{b(h-x)}{\sqrt{1-h^2}} \quad \text{Binary Distance}$$

$$\frac{b(h-x)}{\sqrt{1-h^2}} \quad \text{Binary Distance}$$

- distance between two points in the case of categorical variables
 - the categories are mapped by a function from the categorical variables to numerical values
 - the class need mapping if the categories are not binary
 - no doubt, the one having assigned
 - different value to each category

2020-5-5 08:51

- Always need to determine the value of ψ which may be complex same time
- The computation cost is high because calculating the distance of all the data points for all the training samples.