

→ calculate average dissimilarity of the cluster obtained from step 3 if this value is less than current minimum & put all the medoids found in step 2 as the best set of medoids obtained found in step 2 as the best set of medoids obtained so far

→ Return to Step 1 to start the next iteration

② ENHANCED CLARANS

This method is different from PAM CLARA AND CLARANS. Thus method is produced to improve the accuracy of outlier CLARANS.

It is a two partitioning algorithm which is an nodes instead of selecting as random. Searching operations it is similar to CLARANS but there selected arbitrary nodes reduce the no. of iterations of CLARANS.

procedure

→ Input parameters num local & max neighbour. Initialize p to 1 & min cost to a large number.

→ calculating distance between each data points

→ If s has a lower cost set current + s & go to step

→ otherwise when i → max neighbour compare the lost of current with min cost. If the former is less than min cost set min cost of current & get best node to current increment i by 1. If i = num local but put best node & half others. else go to step 4.

Q)

A) PAM (Partition around Medoid)

PAM uses a k-means method instead method for clustering. It is very robust when compared to k-means in the presence of noise & outliers & mainly in the presence of noise & outliers & swap phase.

D) Build phase: This step is sequentially select k objects which is centrally located. These are objects to be used as medoids.

2) Swap phase: calculate the total cost for each pair of selected & non-selected object.

PAM procedure:

- 1) Input the dataset D
- Randomly select k objects from the dataset & calculate the Total cost T . For each pair of selected $s_i^0 \neq s_j^0$ non selected object s_h
 - For each pair if $T(s_i^0, s_j^0) > T(s_i^0, s_h)$ then s_i^0 is replaced by s_h .
 - Then find similar medoid for each non-selected object.
 - Repeat steps 2, 3 & 4 until find the medoids.

B) CLARA (Clustering large Applications)

- For $i=1$ to S repeat steps 2 to 5
- Draw a sample of $n+2k$ objects randomly from the entire data set & call PAM algorithm to find k medoids of the sample
- For each object o in the entire data set determine k -medoids which is most similar to

along with simplicity Naïve Bayes is known as out perform even highly sophisticated classification methods

Bayes theorem provides a way of calculating posterior probability $P(c|x)$ from $P(c)$, $P(x)$ & $P(c|x)$ look at the eqn below

$$\text{Posterior probability} \propto \frac{\text{Likelihood} \times \text{Class prior probability}}{P(c|x) \times P(c)}$$

↳ Bay Obj
↳ Bay Err
↳ Success
↳ Predictor prior probability of s

Posterior probability

$$P(c|x) = P(c_1|x)P(c_2|x) \dots P(c_n|x) \times P(c)$$

param

8) Discuss clustering based approaches to detect outliers to calculate

gains

=
↳ considered outliers The advantage of the clustering based approach is that small objects in cluster are considered outliers. The advantage of the clustering based approach is that they do not have to be supervised. Moreover clustering - based techniques are capable of being used in an incremental mode.

These are new kind of clustering. The based outlier detection approach have been proposed which are following.

Name	Age	Income	Location
Sandy Young	100	Picky	100% Native
"	100	Wise	100% Native

John Meek good loco
- locomotives

song	young	1000 NISPA'
Bill	"	
Pick	M-a	Wgh
Tar	"	Hesky
Sung	Shor	Soft

Q) Explain the process of predicting a class label using Naïve Bayes classifier

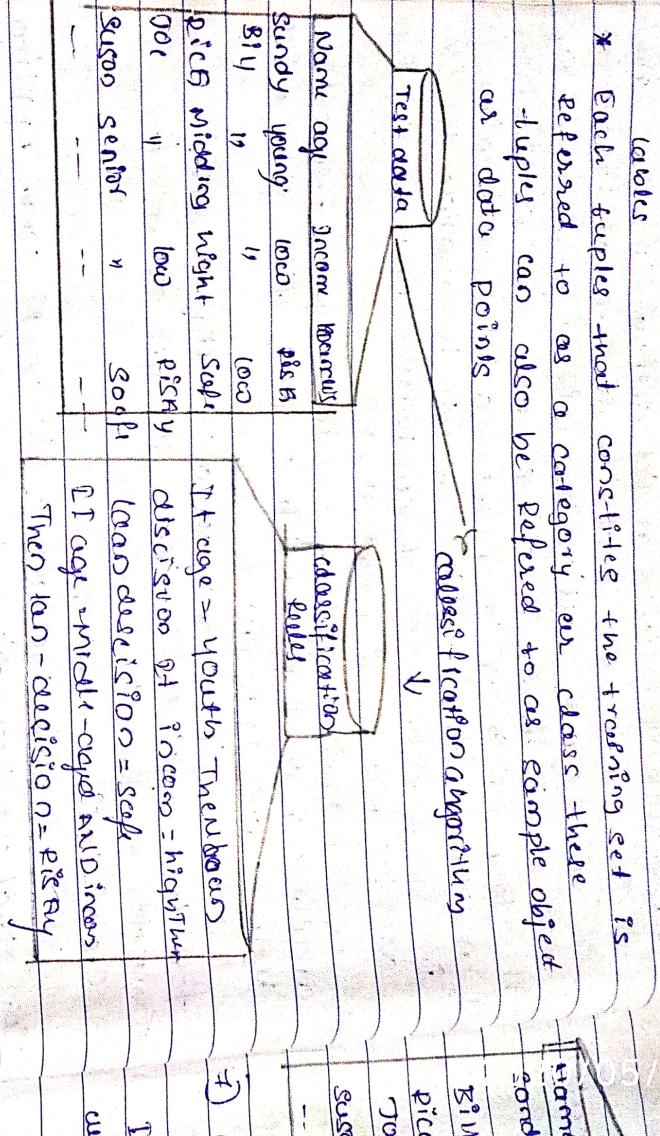
It is a classification based on Bayes Theorem with an assumption of independence among predictors. To simple terms a naïve Bayes classifier assumes that the presence of a particular feature is not related to the presence of any other feature.

For example a fruit may be considered to be an apple if it is red, round, about 3 inches in diameter. Etc. If these features depend on each other or upon the existence of the other feature all of these properties independently contribute to the probability that this fruit is an apple & that is why it is known as naïve.

Naïve Bayes model is easy to build & particularly useful for very large data sets.

or data base tuples either associated class labels

- * Each tuple that constitutes the training set is referred to as a category or class. These tuples can also be referred to as sample objects or data points.



⑤) Using classifier for classification:-

- In this step the classifier is used for classification. Here the test data is used to estimate the accuracy of classification rules can be applied to the test data tuples if the accuracy is considered acceptable.

- 6) By applying the data classification process
classification models predict categorical class labels
for ex:- we can build a classification model to categorize
the bank loan applications either safe or risky
following are the examples of cases where
the data is used to analyse a particular classification
- D A bank loan officer wants to analyze the data
in order to know which customers (loan application)
are risky or which are safe.
- 2) a marketing manager at a company needs to
analyze a customer with a given profile who will
buy a new computer
- These examples are risky or safe
for loan applications data and yes or no for
marketing data
- Classification work :-
- With the help of the bank loan application
that we have discussed above let us understand
the working of classification - the data classification
process includes the following steps
- > D Building the classifier or model
D Using classifier for classification
- => D Building classifier or model :-
- * This step is the learning step or the learning set
or the learning phase
- * In this step the classification algorithm builds
the classifier
- * The classifier is built from the training set made

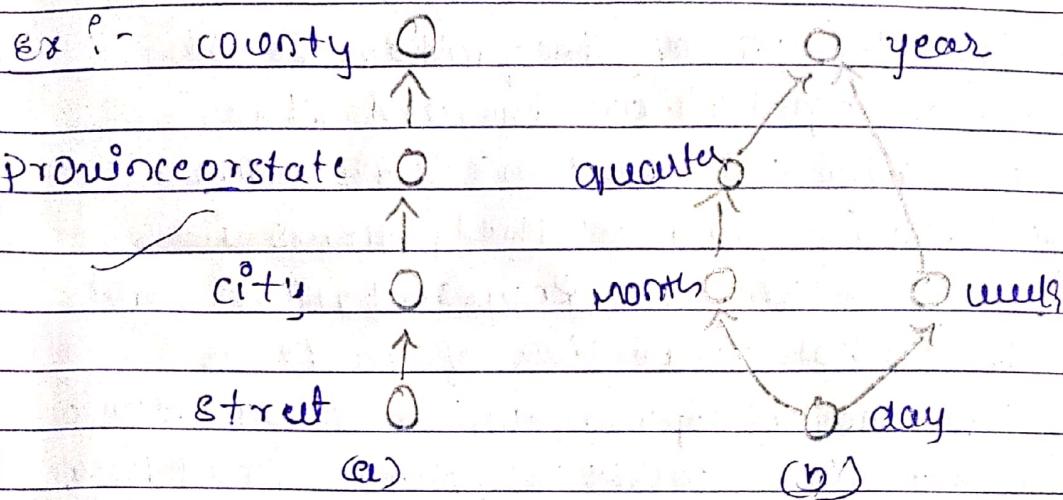
Given a transaction data base let = 0
 $\{o_1, o_2, \dots\}$ be the set of items that appear in the set of items in transaction t_i for any itemset P . Let $O(P)$ be the corresponding set of items $O(P) = \{t \in T | O(P) \geq \text{min sup}\}$ where min sup is a user specified threshold. The task of frequent itemset mining is to find all the frequent itemsets.

There have been many scalable methods developed for frequent pattern mining [3.13.9] however the real bottleneck of the problem is not at the efficiency but at the usability.

To solve this problem is not natural to explore how to "compress" the patterns i.e. find a concise & succinct representation that describes the whole collection of patterns.

Two major approaches have been developed in this direction lossless compression & lossy approximation the former represented by the closed frequent itemset (16.18.19) emphasis too much on the supports of patterns so that its compression power is quite limited the latter represented in the maximal frequent itemsets (S, T, U) $O(P)$

- * Many concept hierarchies are implicit within the database schema
 - * A concept hierarchy that is a total or partial among attributes in a database schema is called schema hierarchy
 - * Concept hierarchies that are common to many applications may be predefined in the data mining system
 - * Data mining should provide users with the flexibility to tailor predefined hierarchies according to their particular needs



Q) Explain how pattern compression can be achieved by pattern clustering.

=> Frequent - patterns (or item sets) mining has been a forced research theme in a data mining due to its broad applications at mining association (2-3) correlation (6) generalizations (7) sequence pattern (4), episodes (5) spatial periodicity (12) emerging pattern (8) & many other important data mining tasks

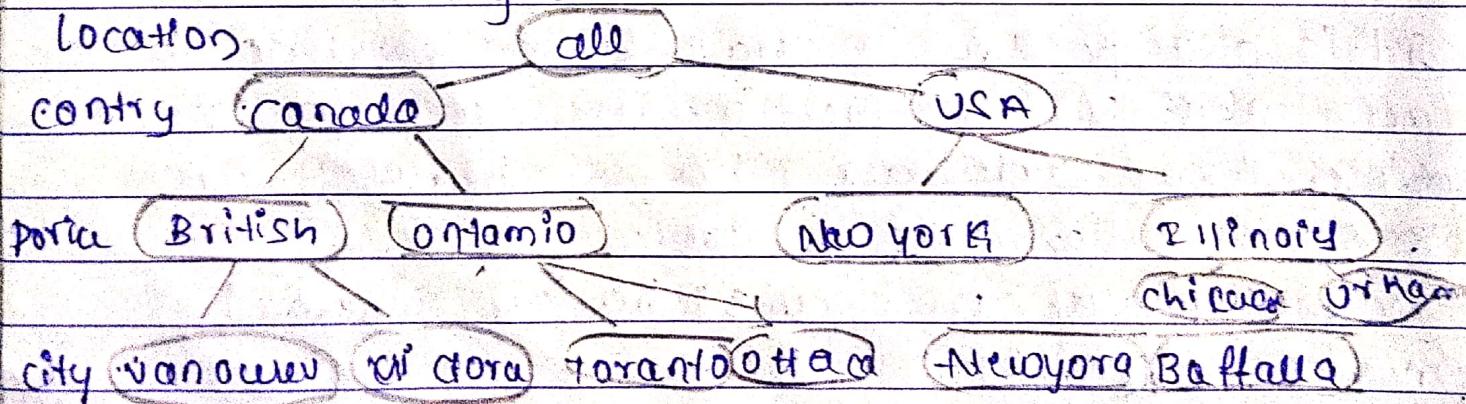
to the dimensions

=> Disadvantages :-

- * Fact constellation solution is difficult to maintain
- * complexity of the schema involved due to the No. of aggregations

Q) With Example illustrate how concept hierarchies are useful in OLAP

- => A concept hierarchy defines a sequence of mapping from a set of low-level concepts to higher level more general concepts
- * consider a concept for the dimension location
- * city values for location include van couver, toronto, new york and chicago
- * Each city, however can be mapped at the province or state to which it belongs
- * The provinces and states can also be mapped to the country to which they belong
- * These mappings a concept hierarchy for the dimension location mapping a set of low level concepts to high-level more general concepts
- * This concept hierarchy is illustrated in below figure.



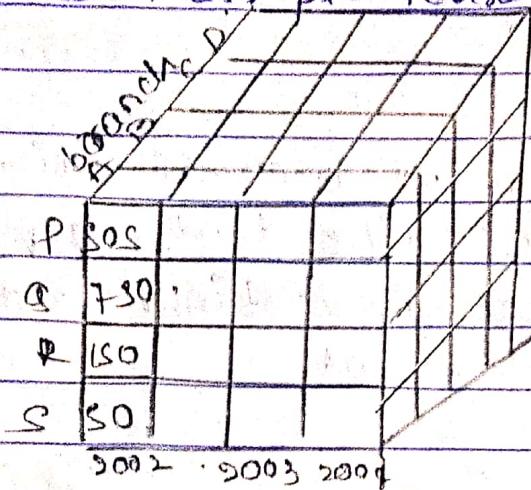
dimension table	fact table	dimension table	fact table
time - key	item - key	item key	time key
day	item - key	item name	time key
day of week	branch key	brand	supplier key
month	location - key	type	from location
greater	dollars sold	supplies key	to location
year	units sold	units shipped	

branch	location	shipper
dimension table	dimension table	dimension table
branch - key	location key	shipper key
branch - name	street	shipper name
branch - type	city	location key
	provincioristem	shipper type
	country	

- * The sales fact table is same as that in star schema
- * The shipping fact table table has the five dimensions. namely item key, time, key, shipper key from location to location
- * The shipping fact table also contains two measures, namely dollars sold & units sold
- * It is also possible to share dimension tables between fact tables for example, time, item and location dimension to both are shared between the sales & shipping fact table
- => Advantages of Fact consolidation Schema:
 - * provides a flexible schema
 - * different fact tables are explicitly assigned to

2020/05/05 10:09

Value corresponding to the data point is
Multipleimensional space concept hierachies may
exists for each attribute, allowing the analysis
of data at multiple levels of abstraction



The cube created at the lowest level of consideration is referred to as the base cuboid. The base cuboid should resemble an individual entity of interest, such as sales or customer.

The above diagram is that you have collected the data for your analysis. These data consist of sales per quarter for the years 2002 to 2005 rather than total per quarter.

⑥ Sequential patterns

This data mining technique helps to discover or identify similar patterns or trends in transaction data for upto period

⑦ Predictions:-

Predictions are used a combination of the other data mining techniques like trend sequences past events are instances to a Right sequence for predicting of future event

⑧ Discuss the need for data reduction with an illustrative example explain data cube aggregation technique

=> Need for data reduction :-

- * A database or data warehouse may store lot of data so it may take very long to perform data analysis and mining on such huge amounts of data
- * Data reduction techniques can be applied to obtain a reduced representation of the data set that is much smaller in volume but still contains critical information
- * complex data analysis or mining may take a very long time to run on the complete data if data cube aggregations i.e. aggregation operations are applied to the data is multidimensional aggregation information each cell holds an aggregated data

2) clustering :-

clustering analysis is a data mining technique to identify data that are like each other. This process helps to understand the differences and similarities between the data.

3) Regression :-

This analysis is the data mining method of identifying & analyzing the relationship between variables. It is used to identify the likelihood of a specific variable given the presence of other variables.

4) Association Rules :-

This data mining technique helps to find the association between two or more items. It discovers a hidden pattern in the data set.

5) outlier detection :-

This type of data mining technique refers to observation of data items in the data set which do not match an expected pattern or expected behaviour. This technique can be used in a variety of domains such as intrusion detection and fault detection etc. Outlier detection is also called outlier detection.

F

Appendix 3 Endpaper

CS181005

data mining

Q14

- D What is data mining? Discuss the techniques that influence the development of data mining methods.
- Data Mining is the process of discovering patterns in large data sets, involving methods at the interface of machine learning, statistics and data base systems. Data mining also known as knowledge discovery in databases refers to the non-trivial extraction of implicit, previously unknown and potentially useful information from data stored in databases.
- Data mining Techniques :-
- 1) classification
 - 2) clustering
 - 3) Regression
 - 4) outlier detection
 - 5) sequential patterns
 - 6) prediction
 - 7) association rule
- D Classification :-
- This analysis used to retrieve important and relevant information about data and method. This data mining method helps to classify data in different classes.
- D Clustering :-
- Clustering analysis is a data mining technique to identify data that are alike each