

2019

1) Explain the theoretical foundation of data mining with examples

Theoretical foundation of data mining

According to probability theory data mining mining find the patterns that are interesting only to the extent that they can be used in the decision making process of some enterprise therefore data mining is the study of performance induction on data base.

2) Data Reduction : The basic idea of this theory is reduction data representation which trade accuracy for speed in response to the need to obtain quick approximates answers to queries on very large databases. Some of the data reduction technique are as follows

- Significantly value decomposition, wavelets,
- Regression, log-linear models, histograms,
- clustering, sampling, construction of index trees.

3) Data Compression : The basic idea of this theory is to compress the given data by encoding in term of the following.

- Bits, Association rules, Decision Trees, outliers

4) Pattern Discovery . The basic idea of this theory is to discover patterns occurring on a database following are the areas that contributed to this theory.

- machine learning Neural Network
- classification mining, e.g. sequential pattern matching, clustering.

(E) microeconomic view As per this theory a database schema consists of data patterns that are stored in a database therefore data mining is the task of performing induction on database

⑥ Inductive Database : A part from the DS oriented techniques, there are statistical techniques available for data analysis these techniques can be applied to scientific data from economics & social science as well.

2. Explain the application of data mining

- Data mining in Finance.
we have to increase customer loyalty collecting and analyzing customer behavior
- collect data also one needs to help bank that predict customer behaviour and to launch relevant service of products

② Data mining in Health care:-
Basically it provides government regular for competitor information that can find competitive advantage Although it supports the R&D process And then goto market

2020-5-5 09:11

start to get with applied sciences & its information at every place.

Data mining for Intelligence is related to money laundering hidden data trafficking etc. Also, help in identifying intrusion detection with high level of anomaly detection & with high identify suspicious activity from a day one basically convert text - based crime reported into word processing files.

④ Data mining in Telecommunication:

In this, data mining gains a competitive advantage and reduce customer churn by understanding demographic characteristics of predicted customer behavior.

⑤ Data mining for energy:

As data mining capture weak signals of potentially threatening events. Also identify previously undefined patterns connections

⑥ Data marketing & sales
Basically it enables business to understand the hidden pattern inside historical purchasing transaction data thus helping in planning launching new marketing compliances. Generally we use for market basket analysis. Retail

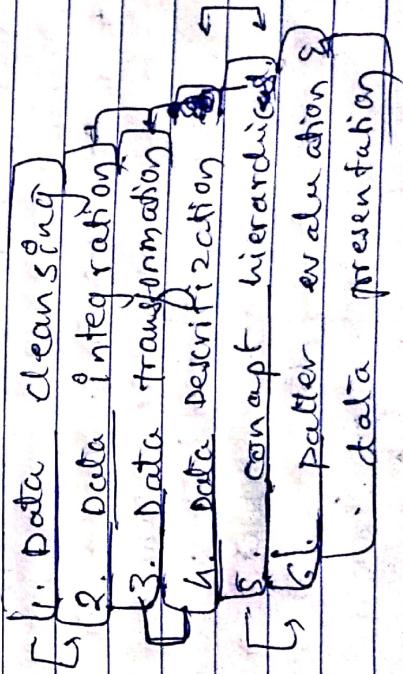
Companies behaviour buying pattern

① Data mining in E-commerce:-

Many E-commerce companies are using data mining business Intelligence to offer cross-sell through their websites. One of the most famous of these is, of course Amazon.

② Data mining in Biological data analysis.
Semantic integration of Histochemistry.
distributed genomic proteinic databases of
Alignment, indexing, Simpanity search
of comparative analysis multiple
redundant sequences. Association of
Path analysts. visualization tools on
genetic data analysis biological data
mining - is very important part of
bioinformatics

③ Describe the different stages of DM process.



2020-5-5 09:

① Data Cleaning

This is very initial stage in the case of DM where the classification of the data becomes an essential component to obtain final data analysis. It involves identifying and removal of inaccurate & possibly data from a set of tables, data base of a record set.

② Data Integration:-

It is a technique which involves the merging of the new set of information with the existing set. The source may, however, involve many data bases, or flat files sets,

③ Data Transformation :-

This requires the transformation of data within formats generally from the source system to the target system some strategies

clustering

(5) Concept hierarchies:- They minimize the data by replacing collecting low level concept from high level concept. The multidimensional data with multiple levels of abstraction are designed by concept hierarchies. The methods are binning, Histogram, analysis, cluster, analysts, etc.

(6) Pattern evaluation of data presentation:-
The data is presented in an efficient manner, the client, as well as the customer can make use of it in the best possible way after going through the above set of stages the data then is presented in forms of graph & diagram & there by understanding it with minimum statistical knowledge.

(7) Discuss about typical OLAP operations on multidimensional data with an ex:

OLAP Online : Analytical processing server
It is a software technology that allows user to analyze information from multiple database systems at the same time. It is based on multidimensional data model & allows the user to query on multi-dimensional data. OLAP database are divided into one or more cubes. These cubes are known as hyper

- Cubes
OLAP operations : They are five basic analytical operations that can be performed on an OLAP cube.

- ① Drill - Down
In drill down operation the less detailed data is converted into highly detailed data if can be done say:-
→ moving down in the concept hierarchy
+ adding a new dimension.



② Roll - Up
It is just opposite of the drill - down operation. It performs aggregation on the OLAP cube. It can be done by :-

- Climbing up in the concept hierarchy
→ Reducing the dimensions.

2020-5-5 09:10

- Creating OLAP cube
- Creating an OLAP plan with Vertica
 - Dia is It selects a sub-cube from the OLAP cube by selecting two or more dimensions To the cube given in the overview section a sub-cube is selected by selecting following dimensions with criteria.
 - location = "A" or "B"
 - item = "q" or "q₁"
 - Item = car or "Bus"
location B
 - Using It selects single elements from the OLAP cube which results in a new sub-cube creation in the cube gives In the overview section Silia is performed on the dimensions Time = "q₁".
- 2020-5-5 09:10
- for Business Plan
Plan Vertical

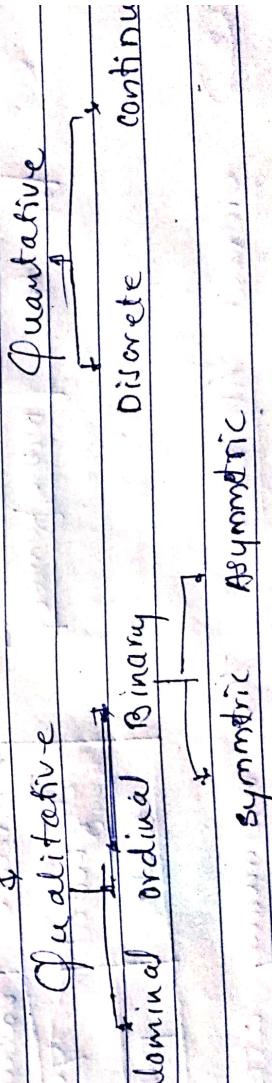
Pivot: It is also known as rotation operation as it rotates the current view of to get a new view of the representation in the sub-cube obtained after the slice operation performing pivot operation given a new view of it.

Customer ID	Name	Address	Phone No.	City	State	Country
C1	John Doe	123 Main St	123-4567890	New York	NY	USA
C2	Jane Smith	456 Elm St	987-6543210	Los Angeles	CA	USA
C3	Mike Johnson	789 Oak St	543-2109870	Chicago	IL	USA
C4	Sarah Davis	234 Pine St	876-5432100	Boston	MA	USA

Ques: Define attribute Explain different types of attribute with example.

An attribute can be seen as a data field that represents characteristics referred to a data object. For a customer object attributes can be customer ID, address etc. we can say that a set of attributes used to describe a given object are known as attribute. A vector or feature vector.

Type of attribute



- ① Nominal Attribute related to names
- The value of a nominal attribute are name of things, some kind of symbols representing some values of nominal attributes that's why nominal category or state of attribute also referred as categorical attribute.
- Eg there is no order (rank position) among values of nominal attributes

sex	Attribute	values
	colours	black, white, brown,
	Categorical Data	Lecturer, professor Assistant professor

Binary Attribute binary data has only 2 values / states for ex: yes or no, effected or unaffected true or false.

i) Symmetric: both value are equally important (Result)

Attribute	values	Attribute	values
Gender	male, female	Queen	defeated

2020-5-5 09:10

ordinal Attribute: - The ordinary attribute contains values that have a meaningful sequence or ranking reduces (order)

between them, but the magnitude b/w values is not actually known, the order of values that shows what is important but don't predict how important it is.

Attribute	Values
Grade.	A, B, C, D, E
Basic pay scale	16, 17, 18

Quantitative Attribute

① Numeric:

if Numeric Attribute is quantitative because it is a measurable quantity, represented in integer or real values. Numerical attributes are of 2 types Interval and ration.

② Discrete: Discrete data have finite values. It can be numerical and can also be categorical form. These attribute has finite or countably infinite set of values.

Attribute	Values
profession	Teacher, Businessman, Peon
ZIP code	300001, 110040

2020

- ③ Continuous: Continuous data have infinite no of states, continuous data is of float type there can be many values b/w 2 & 3

Attribute	Value
Height	5.4, 5.2 - etc
Weight	50, 33 - etc

→ Outlier detection → Supervised & unsupervised

- ① Supervised outlier detection method
 - Supervised outlier detection method
 - we model an outlier detection algorithm
 - classification problem samples examined domain experts used for training & testing
- Challenges
 - classes are unbalanced That is, the pre population of outlier is typically much smaller than that of normal object
 - for handling unbalanced classes can be used such as oversampling.
 - catch as many outlier as possible. it's recall is more important than accuracy (i.e.) net mislabeling novel objects as outliers.

- ② Unsupervised outlier detection methods → In some application scenarios, objects labeled as normal or outlier are not available. Thus: An unsupervised learning method has to be used.
- Unsupervised outlier detection method make implicit assumption:
 - The normal objects are somewhat clustered.

2020-5-5 09:10

challenges

- * Normal objects may not share any strong patterns, but the collective outlier may share high similarity, in or in a small area.
- * In case of normal activities are diverse, so do not fall into high quality clusters unsupervised methods may have a high false positive rate if may not many outliers be detected.

The latest unsupervised methods developed smart ideas to take outliers directly without explicitly detecting clusters.

(3)

- * Semi-supervised outlier detection method.
For many applications, although obtaining some labeled examples is feasible, there are often small the. If some examples are available, labeled normal objects are available:
 - * Use the labeled example and the proximate unlabeled objects to train a model for normal objects.
 - * Those not fitting the model of normal objects are detected as outlier.

If only some labeled outliers are available; a small not to labeled outliers may not cover the possible outliers well.

2020-5-5 09:10

- ④ Find k-th nearest neighbor
- Given n objects in a dataset, find k objects
in dataset to an object such that
distance from the rest of the objects
can be sorted by increasing or
decreasing order.
- In distance based classifier, distance metric
- rely on the computation of dist
- value based on closer distance
distance - based classifiers are highly
defined for k -dimensional datasets for our
value of k .
- ① Index based approach
- By maintaining the neighborhood
index scheme.
- Index: Index is used to record for each
nearest object within radius of current
object.
- * use $\mathcal{O}(k \cdot n \log n)$ neighbors of object
 - * more time is not an advantage
 - * complexity is $\mathcal{O}(n^2)$ due to computation
of $n \cdot n$ the number of objects in the
dataset
- ② Nested loop method
- The nested loop calculate distance
pertaining to all other objects to find the
distance k nearest neighbor related by
low complexity of $\mathcal{O}(kn^2)$

2020-5-5 09:10

Steps:

- (1) Divides the buffer space into two halves (first & second arrays)
- (2) Brute data into blocks & then feed two blocks into the arrays.
- (3) Directly computes the distance b/w each pair of objects inside the array or b/w arrays
- (4) Decide the outlier
- (5) Some computational complexity as the index-based algorithm.
- (6) Local Distance Based Outlier detection error:-
The previous outlier detection schemes are average when it comes to detecting outliers in real world scattered datasets

Steps:

- (1) for each object o in D retrieve its k nearest neighbors.
 - (2) calculate $LDOF$ for each object o .
The objects with $LDOF \geq LDOFs$ are directly discarded.
 - (3) According to their $LDOF$ values sort the objects.
- Output: highest $LDOF$ values of first n objects
- It uses concept of kNN distance.
 kNN anner decides distance of a object
→ It finds the nearest neighbour of the points & decides if it is an outlier or not
- 2020-5-5 09:10

→ It reduces the false detection of normal data as outlier.

(*) what is data cube measure ? discuss the three categories of data cube measure

→ Data cube is a structure that enable OLAP to achieves the multidimensional functionality. The data cube is used to represent data cubes - have categories of data called dimensional and measures.

→ Measure - represents some fact (or number) such as cost or units of service

→ Dimension : represents descriptive categories of data such as time or location
the three categories of data cube measures.

① Distributive :- if the result derived by applying the function on aggregate values is the same as that derived by applying the functions on all the data without partitioning.

Eg:

count(), sum(), min(), max().

② Algebraic : if it can be computed by an algebraic function with m arguments (where m is a bounded integer.) each of which is obtained by applying a distributive aggregate function.

(5) holistic if there is no constant bound
on the storage size needed to describe
a subaggregate

Eg: median(), mode(), rank().