

CS701 DATA WAREHOUSE AND DATA MINING

UNIT I

1. Define data warehouse?

A data warehouse is a repository of multiple heterogeneous data sources organized under a unified schema at a single site to facilitate management decision making .(or)A data warehouse is a subject-oriented, time-variant and nonvolatile collection of data in support of management's decision-making process.

2. What are operational databases?

Organizations maintain large database that are updated by daily transactions are called operational databases.

3. Define OLTP?

If an on-line operational database systems is used for efficient retrieval, efficient storage and management of large amounts of data, then the system is said to be on-line transaction processing.

4. Define OLAP?

Data warehouse systems serves users (or) knowledge workers in the role of data analysis and decision-making. Such systems can organize and present data in various formats. These systems are known as on-line analytical processing systems.

5. How a database design is represented in OLTP systems?

Entity-relation model

6. How a database design is represented in OLAP systems?

- Star schema
- Snowflake schema
- Fact constellation schema

7. Write short notes on multidimensional data model?

Data warehouses and OLTP tools are based on a multidimensional data model. This model is used for the design of corporate data warehouses and department data marts. This model contains a Star schema, Snowflake schema and Fact constellation schemas. The core of the multidimensional model is the data cube.

8. Define data cube?

It consists of a large set of facts (or) measures and a number of dimensions.

9. What are facts?

Facts are numerical measures. Facts can also be considered as quantities by which can analyze the relationship between dimensions.

10. What are dimensions?

Dimensions are the entities (or) perspectives with respect to an organization for keeping records and are hierarchical in nature.

11. Define dimension table?

A dimension table is used for describing the dimension. (e.g.) A dimension table for item may contain the attributes item_name, brand and type.

12. Define fact table?

Fact table contains the name of facts (or) measures as well as keys to each of the related dimensional tables.

13. What are lattice of cuboids?

In data warehousing research literature, a cube can also be called as cuboids. For different (or) set of dimensions, we can construct a lattice of cuboids, each showing the data at different level. The lattice of cuboids is also referred to as data cube.

14.What is apex cuboid?

The 0-D cuboid which holds the highest level of summarization is called the apex cuboid. The apex cuboid is typically denoted by all.

15.List out the components of star schema?

- A large central table (fact table) containing the bulk of data with no redundancy.
- A set of smaller attendant tables (dimension tables), one for each dimension.

16.What is snowflake schema?

The snowflake schema is a variant of the star schema model, where some dimension tables are normalized thereby further splitting the tables in to additional tables.

17.List out the components of fact constellation schema?

This requires multiple fact tables to share dimension tables. This kind of schema can be viewed as a collection of stars and hence it is known as galaxy schema (or) fact constellation schema.

18.Point out the major difference between the star schema and the snowflake schema?

The dimension table of the snowflake schema model may be kept in normalized form to reduce redundancies. Such a table is easy to maintain and saves storage space.

19.Which is popular in the data warehouse design, star schema model (or) snowflake schema model?

Star schema model, because the snowflake structure can reduce the effectiveness and more joins will be needed to execute a query.

20.Define concept hierarchy?

A concept hierarchy defines a sequence of mappings from a set of low-level concepts to higher-level concepts.

UNIT II

1. Define schema hierarchy?

A concept hierarchy that is a total (or) partial order among attributes in a database schema is called a schema hierarchy.

2. List out the OLAP operations in multidimensional data model?

- Roll-up
- Drill-down
- Slice and dice
- Pivot (or) rotate

3. What is roll-up operation?

The roll-up operation is also called drill-up operation which performs aggregation on a data cube either by climbing up a concept hierarchy for a dimension (or) by dimension reduction.

4. What is drill-down operation?

Drill-down is the reverse of roll-up operation. It navigates from less detailed data to more detailed data. Drill-down operation can be taken place by stepping down a concept hierarchy for a dimension.

5. What is slice operation?

The slice operation performs a selection on one dimension of the cube resulting in a sub cube.

6. What is dice operation?

The dice operation defines a sub cube by performing a selection on two (or) more dimensions.

7. What is pivot operation?

This is a visualization operation that rotates the data axes in an alternative presentation of the data.

8. List out the views in the design of a data warehouse?

- Top-down view
- Data source view
- Data warehouse view
- Business query view

9.What are the methods for developing large software systems?

- Waterfall method
- Spiral method

10.How the operation is performed in waterfall method?

The waterfall method performs a structured and systematic analysis at each step before proceeding to the next, which is like a waterfall falling from one step to the next.

11.List out the steps of the data warehouse design process?

- Choose a business process to model.
- Choose the grain of the business process
- Choose the dimensions that will apply to each fact table record.
- Choose the measures that will populate each fact table record.

12.Define ROLAP?

The ROLAP model is an extended relational DBMS that maps operations on multidimensional data to standard relational operations.

13.Define MOLAP?

The MOLAP model is a special purpose server that directly implements multidimensional data and operations.

14.Define HOLAP?

The hybrid OLAP approach combines ROLAP and MOLAP technology, benefiting from the greater scalability of ROLAP and the faster computation of MOLAP,(i.e.) a HOLAP server may allow large volumes of detail data to be stored in a relational database, while aggregations are kept in a separate MOLAP store.

15.What is enterprise warehouse?

An enterprise warehouse collects all the information's about subjects spanning the entire organization. It provides corporate-wide data integration, usually from one (or)more operational systems (or) external information providers. It contains detailed data as well as summarized data and can range in size from a few giga bytes to hundreds of giga bytes, tera bytes (or) beyond. An enterprise data warehouse may be implemented on traditional mainframes, UNIX super servers (or) parallel architecture platforms. It requires business modeling and may take years to design and build.

16.What is data mart?

Data mart is a database that contains a subset of data present in a data warehouse.Data marts are created to structure the data in a data warehouse according to issues such as hardware platforms and access control strategies. We can divide a data warehouse into data marts after the data warehouse has been created. Data marts are usually implemented on low-cost departmental servers that are UNIX (or) windows/NT based. The implementation cycle of the data mart is likely to be measured in weeks rather than months (or) years.

17.What are dependent and independent data marts?

Dependent data marts are sourced directly from enterprise data warehouses. Independent data marts are data captured from one (or) more operational systems (or)external information providers (or) data generated locally with in particular department(or) geographic area.

18.What is virtual warehouse?

A virtual warehouse is a set of views over operational databases. For efficient query processing, only some of the possible summary views may be materialized. A virtual warehouse is easy to build but requires excess capability on operational database servers.

19. Define indexing?

Indexing is a technique, which is used for efficient data retrieval (or) accessing data in a faster manner. When a table grows in volume, the indexes also increase in size requiring more storage.

20. Define metadata?

Metadata is used in data warehouse is used for describing data about data. (i.e.) meta data are the data that define warehouse objects. Metadata are created for the data names and definitions of the given warehouse.

UNIT 3

1. Define Data mining.

It refers to extracting or “mining” knowledge from large amount of data. Data mining is a process of discovering interesting knowledge from large amounts of data stored either, in database, data warehouse, or other information repositories

2. Give some alternative terms for data mining.

- Knowledge mining
- Knowledge extraction
- Data/pattern analysis.
- Data Archaeology
- Data dredging

3. What is KDD.

KDD-Knowledge Discovery in Databases.

4. What are the steps involved in KDD process.

- Data cleaning
- Data Mining
- Pattern Evaluation
- Knowledge Presentation
- Data Integration

- Data Selection
- Data Transformation

5.What is the use of the knowledge base?

Knowledge base is domain knowledge that is used to guide search or evaluate the Interestingness of resulting pattern. Such knowledge can include concept hierarchies used to organize attribute /attribute values in to different levels of abstraction.

6 What is the purpose of Data mining Technique?

It provides a way to use various data mining tasks.

7.Define Predictive model.

It is used to predict the values of data by making use of known results from a different set of sample data.

8.Define descriptive model

- It is used to determine the patterns and relationships in a sample data. Data mining tasks that belongs to descriptive model:
- Clustering
- Summarization
- Association rules
- Sequence discovery

9. Define the term summarization

The summarization of a large chunk of data contained in a web page or a document.

Summarization = characterization=generalization

10. List out the advanced database systems.

- Extended-relational databases
- Object-oriented databases
- Deductive databases
- Spatial databases
- Temporal databases
- Multimedia databases
- Active databases

- Scientific databases
- Knowledge databases

11. Define cluster analysis

Cluster analyses data objects without consulting a known class label. The class labels are not present in the training data simply because they are not known to begin with.

12. Describe challenges to data mining regarding data mining methodology and user interaction issues.

- Mining different kinds of knowledge in databases
- Interactive mining of knowledge at multiple levels of abstraction
- Incorporation of background knowledge
- Data mining query languages and ad hoc data mining
- Presentation and visualization of data mining results
- Handling noisy or incomplete data
- Pattern evaluation

13. Describe challenges to data mining regarding performance issues.

- Efficiency and scalability of data mining algorithms
- Parallel, distributed, and incremental mining algorithms

14. Describe issues relating to the diversity of database types.

- Handling of relational and complex types of data
- Mining information from heterogeneous databases and global information systems

15. What is meant by pattern?

Pattern represents knowledge if it is easily understood by humans; valid on test data with some degree of certainty; and potentially useful, novel, or validates a hunch

about which the user was curious. Measures of pattern interestingness, either objective or subjective, can be used to guide the discovery process.

16. How is a data warehouse different from a database?

Data warehouse is a repository of multiple heterogeneous data sources, organized

under a unified schema at a single site in order to facilitate management decision-making. Database consists of a collection of interrelated data.

17 Define Association Rule Mining.

Association rule mining searches for interesting relationships among items in a given data set

18. When we can say the association rules are interesting?

Association rules are considered interesting if they satisfy both a minimum support threshold and a minimum confidence threshold. Users or domain experts can set such thresholds.

19. Define support and confidence in Association rule mining.

Support S is the percentage of transactions in D that contain $A \cup B$.

Confidence c is the percentage of transactions in D containing A that also contain B .

Support $(A \Rightarrow B) = P(A \cup B)$

Confidence $(A \Rightarrow B) = P(B/A)$

20. How are association rules mined from large databases?

- I step: Find all frequent item sets:
- II step: Generate strong association rules from frequent item sets

21. Describe the different classifications of Association rule mining.

- Based on types of values handled in the Rule
 - i. Boolean association rule
 - ii. Quantitative association rule
- Based on the dimensions of data involved
 - i. Single dimensional association rule
 - ii. Multidimensional association rule
- Based on the levels of abstraction involved
 - i. Multilevel association rule
 - ii. Single level association rule
- Based on various extensions
 - i. Correlation analysis
 - ii. Mining max patterns

UNIT 4

1 What is the purpose of Apriori Algorithm?

Apriori algorithm is an influential algorithm for mining frequent item sets for Boolean association rules. The name of the algorithm is based on the fact that the algorithm uses prior knowledge of frequent item set properties.

2. Define anti-monotone property.

If a set cannot pass a test, all of its supersets will fail the same test as well.

3. How to generate association rules from frequent item sets?

Association rules can be generated as follows

For each frequent item set l , generate all non empty subsets of l .

For every non empty subsets s of l , output the rule " $S \Rightarrow (l-s)$ " if

Support count(l) = min_conf,

Support_count(s)

Where min_conf is the minimum confidence threshold.

4. Give few techniques to improve the efficiency of Apriori algorithm.

- Hash based technique
- Transaction Reduction
- Portioning
- Sampling
- Dynamic item counting

5. What are the things suffering the performance of Apriori candidate generation technique.

- Need to generate a huge number of candidate sets
- Need to repeatedly scan the scan the database and check a large set of candidates by pattern matching

6. Describe the method of generating frequent item sets without candidate generation.

Frequent-pattern growth(or FP Growth) adopts divide-and-conquer strategy.

Steps:

Compress the database representing frequent items into a frequent pattern tree or FP tree

Divide the compressed database into a set of conditional database

Mine each conditional database separately

7. Mention few approaches to mining Multilevel Association Rules

- Uniform minimum support for all levels(or uniform support)
- Using reduced minimum support at lower levels(or reduced support)
- Level-by-level independent
- Level-cross filtering by single item
- Level-cross filtering by k-item set

8. What are multidimensional association rules?

Association rules that involve two or more dimensions or predicates

- Interdimension association rule: Multidimensional association rule with no repeated predicate or dimension
- Hybrid-dimension association rule: Multidimensional association rule with multiple occurrences of some predicates or dimensions.

9. Define constraint-Based Association Mining.

Mining is performed under the guidance of various kinds of constraints provided by the user.

The constraints include the following

- Knowledge type constraints
- Data constraints
- Dimension/level constraints
- Interestingness constraints
- Rule constraints.

10. Define the concept of classification.

Two step process

- A model is built describing a predefined set of data classes or concepts.
- The model is constructed by analyzing database tuples described by attributes. The model is used for classification.

11 What is Decision tree?

A decision tree is a flow chart like tree structures, where each internal node denotes a test on an attribute, each branch represents an outcome of the test, and leaf nodes represent classes or class distributions. The top most in a tree is the root node.

12. What is Attribute Selection Measure?

The information Gain measure is used to select the test attribute at each node in the decision tree. Such a measure is referred to as an attribute selection measure or a measure of the goodness of split.

13. Describe Tree pruning methods.

When a decision tree is built, many of the branches will reflect anomalies in the training data due to noise or outlier. Tree pruning methods address this problem of over fitting the data.

Approaches:

- Pre pruning
- Post pruning

14. Define Pre Pruning

A tree is pruned by halting its construction early. Upon halting, the node becomes a leaf. The leaf may hold the most frequent class among the subset samples.

15. Define Post Pruning.

Post pruning removes branches from a “Fully grown” tree. A tree node is pruned by removing its branches.

Eg: Cost Complexity Algorithm

16. What is meant by Pattern?

Pattern represents the knowledge.

17. Define the concept of prediction.

Prediction can be viewed as the construction and use of a model to assess the class of an unlabeled sample or to assess the value or value ranges of an attribute that a given sample is likely to have.

18 What is the use of Regression?

Regression can be used to solve the classification problems but it can also be used for applications such as forecasting. Regression can be performed using many different types of techniques; in actually regression takes a set of data and fits the data to a formula.

19 What are the requirements of cluster analysis?

The basic requirements of cluster analysis are

- Dealing with different types of attributes.
- Dealing with noisy data.
- Constraints on clustering.
- Dealing with arbitrary shapes.
- High dimensionality
- Ordering of input data
- Interpretability and usability
- Determining input parameter and
- Scalability

20.What are the different types of data used for cluster analysis?

The different types of data used for cluster analysis are interval scaled, binary, nominal, ordinal and ratio scaled data.

UNIT V

1. Define Clustering?

Clustering is a process of grouping the physical or conceptual data object into clusters.

2. What do you mean by Cluster Analysis?

A cluster analysis is the process of analyzing the various clusters to organize the different objects into meaningful and descriptive objects.

3. What are the fields in which clustering techniques are used?

- Clustering is used in biology to develop new plants and animal taxonomies.
- Clustering is used in business to enable marketers to develop new distinct groups of their customers and characterize the customer group on basis of purchasing.
- Clustering is used in the identification of groups of automobiles Insurance policy customer.
- Clustering is used in the identification of groups of house in a city on the basis of house type, their cost and geographical location.

- Clustering is used to classify the document on the web for information discovery.

4. What are the requirements of cluster analysis?

The basic requirements of cluster analysis are

- Dealing with different types of attributes.
- Dealing with noisy data.
- Constraints on clustering.
- Dealing with arbitrary shapes.
- High dimensionality
- Ordering of input data
- Interpretability and usability
- Determining input parameter and
- Scalability

5. What are the different types of data used for cluster analysis?

The different types of data used for cluster analysis are interval scaled, binary, nominal, ordinal and ratio scaled data.

6. What are interval scaled variables?

Interval scaled variables are continuous measurements of linear scale. For Example , height and weight, weather temperature or coordinates for any cluster. These measurements can be calculated using Euclidean distance or Minkowski distance.

7. Define Binary variables? And what are the two types of binary variables?

Binary variables are understood by two states 0 and 1, when state is 0, variable is absent and when state is 1, variable is present. There are two types of binary variables, symmetric and asymmetric binary variables. Symmetric variables are those variables that have same state values and weights. Asymmetric variables are those variables that have not same state values and weights.

8. Define nominal, ordinal and ratio scaled variables?

A nominal variable is a generalization of the binary variable. Nominal variable has more than two states, For example, a nominal variable, color consists of four states, red, green, yellow, or black.

In Nominal variables the total number of states is N and it is denoted by letters, symbols or integers.

An ordinal variable also has more than two states but all these states are ordered in a meaningful sequence.

A ratio scaled variable makes positive measurements on a non-linear scale, such as exponential scale, using the formula

$$Ae^{Bt} \text{ or } Ae^{-Bt}$$

Where A and B are constants.

9. What do u mean by partitioning method?

In partitioning method a partitioning algorithm arranges all the objects into various partitions, where the total number of partitions is less than the total number of objects. Here each partition represents a cluster. The two types of partitioning method are k-means and k-medoids.

10. Define CLARA and CLARANS?

Clustering in LARge Applications is called as CLARA. The efficiency of CLARA depends upon the size of the representative data set. CLARA does not work properly if any representative data set from the selected representative data sets does not find best k-medoids.

To recover this drawback a new algorithm, Clustering Large Applications based upon RANdomized search (CLARANS) is introduced. The CLARANS works like CLARA, the only difference between CLARA and CLARANS is the clustering process that is done after selecting the representative data sets.

11. What is Hierarchical method?

Hierarchical method groups all the objects into a tree of clusters that are arranged in a hierarchical order. This method works on bottom-up or top-down approaches.

12. Differentiate Agglomerative and Divisive Hierarchical Clustering?

Agglomerative Hierarchical clustering method works on the bottom-up approach.

In Agglomerative hierarchical method, each object creates its own clusters. The single Clusters are merged to make larger clusters and the process of merging continues until all the singular clusters are merged into one big cluster that consists of all the objects.

Divisive Hierarchical clustering method works on the top-down approach. In this method all the objects are arranged within a big singular cluster and the large cluster is continuously divided into smaller clusters until each cluster has a single object.

13. What is CURE?

Clustering Using Representatives is called as CURE. The clustering algorithms generally work on spherical and similar size clusters. CURE overcomes the problem of spherical and similar size cluster and is more robust with respect to outliers.

14. Define Chameleon method?

Chameleon is another hierarchical clustering method that uses dynamic modeling. Chameleon is introduced to recover the drawbacks of CURE method. In this method two clusters are merged, if the interconnectivity between two clusters is greater than the interconnectivity between the objects within a cluster.

15. Define Density based method?

Density based method deals with arbitrary shaped clusters. In density-based method, clusters are formed on the basis of the region where the density of the objects is high.

16. What is a DBSCAN?

Density Based Spatial Clustering of Application Noise is called as DBSCAN. DBSCAN is a density based clustering method that converts the high-density objects regions into clusters with arbitrary shapes and sizes. DBSCAN defines the cluster as a maximal set of density connected points.

17. What do you mean by Grid Based Method?

In this method objects are represented by the multi resolution grid data structure. All the objects are quantized into a finite number of cells and the collection of cells build the grid structure of objects. The clustering operations are performed on that grid structure. This method is widely used because its processing time is very fast and that is independent of number of objects.

18. What is a STING?

Statistical Information Grid is called as STING; it is a grid based multi resolution clustering method. In STING method, all the objects are contained into rectangular cells, these cells are kept into various levels of resolutions and these levels are arranged in a hierarchical structure.

19. Define Wave Cluster?

It is a grid based multi resolution clustering method. In this method all the objects

are represented by a multidimensional grid structure and a wavelet transformation is applied for finding the dense region. Each grid cell contains the information of the group of objects that map into a cell. A wavelet transformation is a process of signaling that produces the signal of various frequency sub bands.

20. What is Model based method?

For optimizing a fit between a given data set and a mathematical model based methods are used. This method uses an assumption that the data are distributed by probability distributions. There are two basic approaches in this method that are

1. Statistical Approach
2. Neural Network Approach.

21. What is the use of Regression?

Regression can be used to solve the classification problems but it can also be used for applications such as forecasting. Regression can be performed using many different types of techniques; in actually regression takes a set of data and fits the data to a formula.

22. What are the reasons for not using the linear regression model to estimate the output data?

There are many reasons for that, One is that the data do not fit a linear model, It is possible however that the data generally do actually represent a linear model, but the

linear model generated is poor because noise or outliers exist in the data.

Noise is erroneous data and outliers are data values that are exceptions to the usual and expected data.

23. What are the two approaches used by regression to perform classification?

Regression can be used to perform classification using the following approaches

1. Division: The data are divided into regions based on class.
2. Prediction: Formulas are generated to predict the output class value.

24. What do you mean by logistic regression?

Instead of fitting a data into a straight line logistic regression uses a logistic curve.

The formula for the univariate logistic curve is

$$P = \frac{e^{(C_0 + C_1 X_1)}}{1 + e^{(C_0 + C_1 X_1)}}$$

The logistic curve gives a value between 0 and 1 so it can be interpreted as the probability of class membership.

25. What is Time Series Analysis?

A time series is a set of attribute values over a period of time. Time Series Analysis may be viewed as finding patterns in the data and predicting future values.

26. What are the various detected patterns?

Detected patterns may include:

• Trends : It may be viewed as systematic non-repetitive changes to the values over time.

- “ Cycles : The observed behavior is cyclic.
- “ Seasonal : The detected patterns may be based on time of year or month or day.
- “ Outliers : To assist in pattern detection , techniques may be needed to remove or reduce the impact of outliers.

27. What is Smoothing?

Smoothing is an approach that is used to remove the nonsystematic behaviors found in time series. It usually takes the form of finding moving averages of attribute values. It is used to filter out noise and outliers.