# Question Bank ETE Data Mining Solution

**5 Marks Question**



**1. Demonstrate in detail about data mining steps in the process of knowledge discovery?**

Data mining is a process within the larger scope of knowledge discovery in databases (KDD). It involves extracting valuable insights and patterns from large datasets. The data mining process typically consists of the following steps:
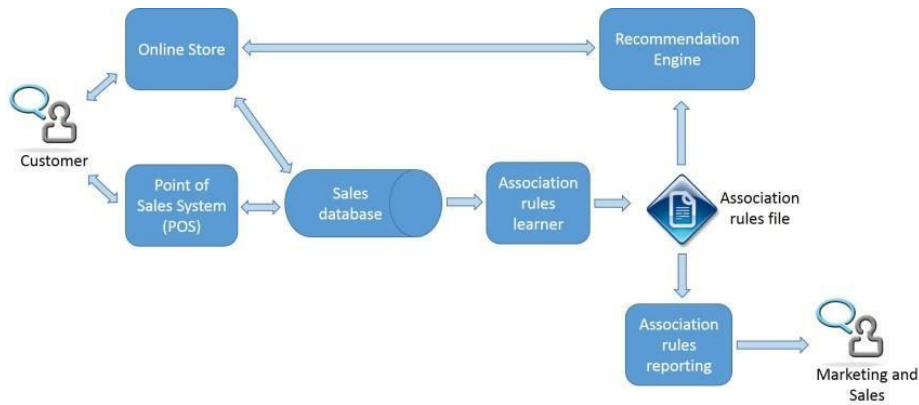
1. Problem Definition: Clearly define the objective of the data mining project.
2. Data Collection: Gather relevant data from various sources.
3. Data Cleaning: Remove noise, inconsistencies, and irrelevant information from the data.
4. Data Integration: Combine data from multiple sources into a unified dataset.
5. Data Selection: Identify and extract the subsets of data that are most relevant to the analysis.
6. Data Transformation: Convert data into a suitable format for analysis.
7. Data Mining: Apply various techniques and algorithms to extract patterns from the dataset.
8. Pattern Evaluation: Assess the quality and usefulness of the discovered patterns.
9. Knowledge Presentation: Communicate the findings to stakeholders through visualizations or reports.
10. Knowledge Utilization: Leverage the discovered knowledge to make informed decisions and drive outcomes.

**2. How will identify frequent itemset using Market Basket Analysis?**

To identify frequent itemsets using Market Basket Analysis, you can follow these steps:

1. Data Preparation: Gather transactional data that represents the items purchased by customers. Each transaction should contain a list of items bought together.

2. Encode the Data: Transform transactional data into a binary format, such as a matrix or itemsets.
3. Set Minimum Support Threshold: Determine the minimum frequency required for an itemset to be considered "frequent."
4. Generate Itemsets: Count the frequency of individual items in the dataset.
5. Prune Infrequent Itemsets: Remove itemsets that do not meet the minimum support threshold.
6. Generate Candidate Itemsets: Combine frequent itemsets to create larger itemsets.
7. Count Support for Candidate Itemsets: Scan data to count occurrences of each candidate itemset.
8. Prune Candidate Itemsets: Remove candidate itemsets that do not meet the minimum support threshold.
9. Repeat Steps 6-8: Iteratively generate candidate itemsets, count support, and prune.
10. Extract Frequent Itemsets: Collect all frequent itemsets obtained, representing items that occur together frequently.

## 3. Explain about Rare and Negative Pattern with suitable example.

Rare and negative patterns are two types of interesting patterns that can be discovered through data mining.

1. Rare Patterns: Rare patterns refer to infrequent or uncommon combinations of items or events in a dataset. These patterns can provide valuable insights and highlight unusual occurrences that deviate from the norm. They are often of interest in various domains such as fraud detection, anomaly detection, and rare disease diagnosis.

For example, in a retail dataset, a rare pattern could be a combination of items that are rarely purchased together, such as "diapers" and "wine." Discovering such a pattern could be valuable for targeted marketing campaigns or store layout optimization.

2. Negative Patterns: Negative patterns, also known as negative association rules or anti-patterns, represent associations between items that are significantly less likely to co-occur than expected by chance. These patterns highlight interesting negative relationships between items.

For example, in a movie recommendation system, a negative pattern could be the association between the genres "Action" and "Romance," indicating that customers who like action movies are less likely to prefer romance movies. Discovering negative patterns can be useful in personalization, cross-selling strategies, or understanding customer preferences.

## 4. Compare OLTP and OLAP.

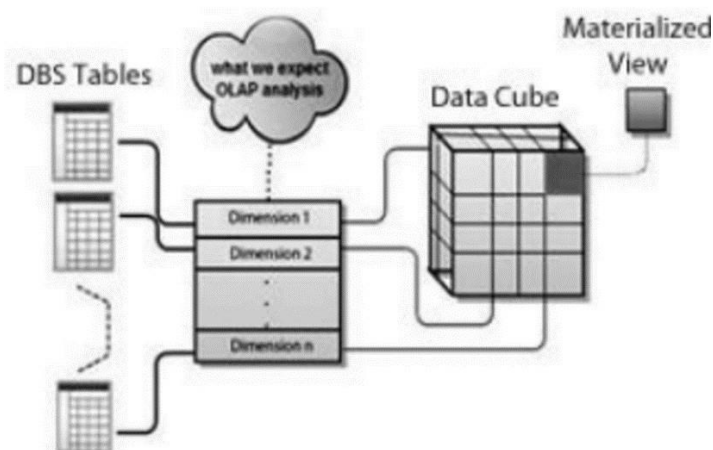|  | OLTP | OLAP |
| --- | --- | --- |
| Purpose | Supports real-time transactional operations. | Supports complex, analytical queries and reporting. |
| Database Design | Designed for efficient data modification (insert, update, delete). | Designed for efficient data retrieval and analysis. |

| | | |
|---|---|---|
| Data Granularity | Deals with detailed and current transactional data. | Deals with aggregated and historical data. |
| Data Structure | Normalized relational database structure. | Dimensional or multi-dimensional database structure. |
| Query Types | Simple, short and frequent queries. | Complex, ad-hoc and long-running queries. |
| User Interaction | Used by operational staff to perform day-to-day transactions. | Used by managers and analysts for decision-making. |
| Performance | Emphasizes quick response times for individual transactions. | Emphasizes query performance for complex analysis. |

**5. Discuss about data cube materialization**

Data cube materialization is the process of precomputing and storing summarized data in a multidimensional structure called a data cube. It involves aggregating and storing data at different levels of granularity to enable faster query response times and facilitate complex analysis. By precalculating aggregations, data cube materialization reduces the need for repetitive calculations during query execution, improving overall performance. However, it also requires additional storage space to store the precomputed results. Materialized data cubes are commonly used in OLAP (Online Analytical Processing) systems to provide efficient and interactive analysis capabilities for decision support and data exploration tasks.

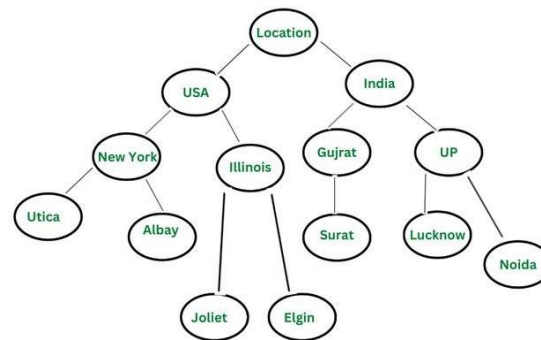The process of data cube materialization typically involves the following steps:

1. Cube Definition: Identify the dimensions and measures to be included in the data cube. Dimensions define the perspectives from which the data can be analyzed, while measures represent the quantitative values of interest.
2. Cube Generation: Perform the necessary calculations and aggregations to generate the summary data for each combination of dimensions and measures. This involves aggregating and summarizing the source data based on the defined dimensions.
3. Cube Storage: Store the precomputed data cubes in a specialized data structure optimized for efficient query processing. This may involve storing the cube in a multidimensional array, a relational table, or other formats depending on the underlying database system.
4. Cube Maintenance: Periodically update and refresh the data cubes to ensure they remain up to date with the source data. This involves handling incremental updates or modifications to the underlying data while maintaining the integrity of the precomputed cubes.

**6. Explain the Role of Concept Hierarchies.**

The role of concept hierarchies in data analysis and decision-making is to provide a structured representation of data at different levels of abstraction. Concept hierarchies organize data attributes or dimensions into hierarchical structures that capture the relationships between different levels of detail or generalization. Here are five key roles of concept hierarchies:

1. Data Summarization: Concept hierarchies provide an aggregated view of data at different levels, allowing for a concise understanding of patterns and trends.
2. Data Generalization: Concept hierarchies enable the movement from specific details to more general concepts, aiding in the identification of high-level insights and patterns.
3. Data Drill-Down: Concept hierarchies support exploring detailed data by navigating from higher-level concepts to specific information, uncovering hidden patterns and anomalies.
4. Data Integration: Concept hierarchies facilitate the integration of data from various sources or levels of detail, enabling comparison and alignment for multidimensional analysis.
5. Data Navigation and Exploration: Concept hierarchies offer a navigational structure for users to interactively explore complex datasets, helping them discover insights and make informed decisions.



**Concept Hierarchy for Dimension Location**

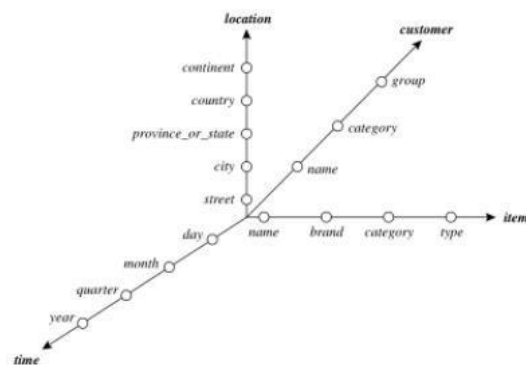**7. What are the steps in designing the data warehouse?**

Designing a data warehouse involves several key steps to ensure that the resulting data structure supports efficient and effective data analysis. Here are the necessary steps in designing a data warehouse:

1. Requirement Gathering: Understand business needs and data sources.
2. Data Source Analysis: Analyze source systems and identify relevant data entities and quality issues.
3. Data Modeling: Develop a conceptual model to represent data structure and relationships.
4. Schema Design: Design physical schema for efficient querying and data retrieval.
5. Extract, Transform, Load (ETL): Implement processes to extract, transform, and load data into the warehouse.
6. Metadata Design: Define metadata structures to document and manage warehouse information.
7. Performance Optimization: Optimize design for efficient data retrieval and analysis.
8. Security and Access Control: Establish security measures and access controls.
9. Testing and Validation: Thoroughly test and validate data warehouse design and ETL processes.

| 10. Deployment and Maintenance: Deploy warehouse, monitor performance, and maintain/update as needed. |
| --- |

## 8. Explain Starnet Query Model

In a starnet model, radial lines emanate from a central point. Each radial line represents a dimension of the data, and the hierarchy of that dimension is represented along the line. Based on the starnet query model, the availability of OLAP operations such as Roll-up and Drill-down can be known. Figure shows an example of the starnet query model. The model has four dimensions which are time, location, customer, and item. To query the data, the data can be rolled up along the time dimension from day to month. The data can also be drilled down along the location dimension from country to city.



### 9. Briefly explain a recommended approach for data warehouse development.

A recommended approach for data warehouse development is to follow a systematic and iterative process that involves the following steps:

1. Requirement Gathering: Understand the business objectives, data needs, and reporting requirements of the organization.
2. Data Source Analysis: Analyze the structure, content, and quality of the source systems' data. Identify the relevant entities, attributes, and relationships that need to be captured in the data warehouse.
3. Data Modeling: Develop a conceptual data model that represents the dimensions, hierarchies, and measures of the data warehouse.
4. Schema Design: Translate the conceptual data model into a physical schema design. This involves creating tables, defining primary and foreign keys, and establishing relationships between tables
5. ETL Development: Design and implement the Extract, Transform, Load (ETL) processes that will extract data from the source systems, transform it into the desired format, and load it into the data warehouse.
6. Testing and Validation: Thoroughly test and validate the data warehouse to ensure data accuracy, integrity, and adherence to business rules.
7. Deployment and User Training: Deploy the data warehouse into a production environment, ensuring proper infrastructure and security measures are in place.
8. Ongoing Maintenance and Enhancement: Establish a regular maintenance and monitoring process for the data warehouse