

<p>## 1. What is data mining? Discuss the significance of data mining?</p> <p>*Data mining* is the process of discovering patterns, correlations, and anomalies within large datasets to predict outcomes and extract useful information. It involves techniques from machine learning, statistics, and database systems.</p> <p>*Significance of data mining*:</p> <ul style="list-style-type: none">- Helps businesses make data-driven decisions- Enables prediction of future trends and behaviors- Identifies hidden patterns in large datasets- Improves customer relationship management- Enhances fraud detection and risk management- Supports scientific discovery and research- Optimizes marketing campaigns and sales strategies	<p>## 13. List some evaluation metrics used in Classification.</p> <p>Evaluation metrics for classification:</p> <ol style="list-style-type: none">1. *Accuracy*: $(TP+TN)/(TP+TN+FP+FN)$2. *Precision*: $TP/(TP+FP)$3. *Recall/Sensitivity*: $TP/(TP+FN)$4. *F1-Score*: $2(Precision*Recall)/(Precision+Recall)$	<ul style="list-style-type: none">- Simple structure, easy to understand- Optimized for query performance- May have redundant data in dimensions <p>2. *Snowflake Schema*:</p> <ul style="list-style-type: none">- Normalized version of star schema- Dimension tables are normalized into multiple related tables- Reduces redundancy but more complex queries <p>3. *Galaxy Schema (Fact Constellation)*:</p> <ul style="list-style-type: none">- Multiple fact tables sharing dimension tables- Used for complex analysis across multiple business processes <p>4. *Data Vault Modeling*:</p> <ul style="list-style-type: none">- Hybrid approach combining 3NF and dimensional modeling- Consists of hubs, links, and satellites- Designed for flexibility and historical tracking <p>5. *Dimensional Modeling*:</p> <ul style="list-style-type: none">- Focuses on business processes (facts) and context (dimensions)- Uses conformed dimensions for consistency- Optimized for analytical queries
<p>## 2. Discuss the steps involved in KDD?</p> <p>*Knowledge Discovery in Databases (KDD)* involves these steps:</p> <ol style="list-style-type: none">1. *Data Cleaning*: Remove noise and inconsistent data2. *Data Integration*: Combine data from multiple sources3. *Data Selection*: Retrieve relevant data from the database4. *Data Transformation*: Convert data into appropriate forms for mining5. *Data Mining*: Apply intelligent methods to extract patterns6. *Pattern Evaluation*: Identify truly interesting patterns7. *Knowledge Presentation*: Present discovered knowledge to users	<p>## 20. Define Data Warehouse. Explain its characteristics.</p> <p>*Data Warehouse*: A subject-oriented, integrated, time-variant, and non-volatile collection of data in support of management's decision-making process.</p> <p>*Characteristics*:</p> <ol style="list-style-type: none">1. *Subject-Oriented*: Organized around major subjects (e.g., customers, products)2. *Integrated*: Consistent naming, encoding, and formatting across sources3. *Time-Variant*: Data is stored with a time dimension for historical analysis4. *Non-Volatile*: Data is read-only and not updated or deleted5. *Summarized*: Contains aggregated data rather than detailed transactions6. *Large Volume*: Typically contains terabytes of historical data7. *Optimized for Analysis*: Designed for complex queries rather than transactions8. *Diverse Sources*: Integrates data from multiple operational systems9. *Decision Support*: Supports business intelligence and analytics	<p>## 23. What is OLAP? Explain its characteristics.</p> <p>*OLAP (Online Analytical Processing)*: A category of software tools that provides analysis of data stored in a database, enabling complex analytical calculations, trend analysis, and sophisticated data modeling.</p> <p>*Characteristics*:</p> <ol style="list-style-type: none">1. *Multidimensional View*: Data is organized in cubes with multiple dimensions2. *Fast Query Performance*: Optimized for complex analytical queries3. *Interactive Analysis*: Users can drill down, roll up, slice and dice data4. *Aggregated Data*: Works with summarized rather than detailed data5. *Historical Perspective*: Focuses on trends over time6. *Calculated Metrics*: Supports complex calculations and KPIs7. *What-if Analysis*: Enables scenario modeling and forecasting8. *Consistent Reporting*: Provides single version of truth
<p>## 3. Discuss various types of data used in Data Mining?</p> <p>Types of data used in data mining:</p> <ol style="list-style-type: none">1. *Structured Data*: Organized in tables (e.g., relational databases)2. *Unstructured Data*: No predefined format (e.g., text, images)3. *Semi-structured Data*: Neither raw nor strictly structured (e.g., XML, JSON)4. *Time-series Data*: Data points indexed in time order5. *Spatial Data*: Related to geographic locations6. *Multimedia Data*: Images, audio, video7. *Stream Data*: Continuous, real-time data flows8. *Graph Data*: Network or web data	<p>## 21. Describe ETL process in detail.</p> <p>*ETL (Extract, Transform, Load)* Process:</p> <ol style="list-style-type: none">1. *Extract*:<ul style="list-style-type: none">- Gather data from various source systems (databases, files, APIs)- Handle different data formats (structured, semi-structured, unstructured)- May involve full extraction or incremental extraction- Resolve connectivity and access issues2. *Transform*:<ul style="list-style-type: none">- *Cleaning*: Handle missing values, correct errors, standardize formats- *Integration*: Resolve naming conflicts, unify measurement units- *Aggregation*: Summarize data at desired levels- *Derivation*: Calculate new fields or metrics- *Filtering*: Select relevant data- *Sorting*: Order data appropriately- *Validation*: Ensure data quality and consistency3. *Load*:<ul style="list-style-type: none">- *Initial Load*: First-time population of data warehouse- *Incremental Load*: Periodic updates with new/changed data- *Full Refresh*: Complete reload of data- Index creation for performance- Constraints and integrity checks <p>Additional ETL components:</p> <ul style="list-style-type: none">- Scheduling and automation- Error handling and logging- Metadata management- Data quality monitoring- Performance optimization	<p>## 24. What is Data Cube? Explain the types of OLAP.</p> <p>*Data Cube*: A multidimensional structure that allows data to be modeled and viewed in multiple dimensions (e.g., product, geography, time).</p> <p>*Types of OLAP*:</p> <ol style="list-style-type: none">1. *MOLAP (Multidimensional OLAP)*:<ul style="list-style-type: none">- Stores data in optimized multidimensional arrays- Fast query performance- Limited in handling large volumes of data2. *ROLAP (Relational OLAP)*:<ul style="list-style-type: none">- Uses relational database with star/snowflake schema- Handles large data volumes- Slower than MOLAP for complex queries3. *HOLAP (Hybrid OLAP)*:<ul style="list-style-type: none">- Combines MOLAP and ROLAP approaches- Stores aggregations in MOLAP, detailed data in ROLAP- Balances performance and scalability4. *DOLAP (Desktop OLAP)*:<ul style="list-style-type: none">- Client-side OLAP tools- Works with extracted subsets of data- Good for personal analysis
<p>## 10. Why is managing inconsistencies important in data mining?</p> <p>Managing inconsistencies is important because:</p> <ol style="list-style-type: none">1. *Accuracy*: Ensures reliable results and predictions2. *Decision Making*: Prevents incorrect business decisions3. *Model Performance*: Improves quality of mining models4. *Data Integration*: Necessary when combining multiple sources5. *Efficiency*: Reduces processing time and resources6. *Reproducibility*: Enables consistent results across analyses7. *Compliance*: Meets regulatory requirements for data quality8. *Customer Trust*: Maintains confidence in data-driven insights	<p>## 22. What is data warehouse design? Discuss some modeling types.</p> <p>*Data Warehouse Design*: The process of defining the architecture, components, and structure of a data warehouse to meet business requirements.</p> <p>*Modeling Types*:</p> <ol style="list-style-type: none">1. *Star Schema*:<ul style="list-style-type: none">- Central fact table surrounded by dimension tables	