



## NOTES: DATA MANAGEMENT FOR MACHINE LEARNING

### 1

Assume that you are working on a Real Estate app that allows the users to Search New Projects, Homes, Apartments, Offices, Shops, and Showrooms for Buy, Sell, and Rent. This property app is required for fulfilling all the real estate needs for searching, shortlisting and finalizing property of their choice. Be it a ready-to-move flat in a new project, investment in shops, offices or showrooms, or girls PG or boys PG in city, app should have listings from qualified property agents & top builders as well as no brokerage properties. This house search app should make property search on mobile simpler, faster and smoother. Search for residential and commercial property listings from owners, brokers and top builders on this user-friendly app should help millions of users to start their real estate journey. To make it more interesting and meaningful, when the users are online searching for a property on the app, the recommendations about the relevant properties should be shown to her. Also relevant advertisements has to be shown on the app screen to increase the chances of converting these paid advertisements into increasing the foot falling and result into actual selling. [3 + 2 + 1 + 1.5 + 2 + 2.5 + 1 + 2 = 15]

- a) Who are the primary users of this systems and what type of user interfaces are suitable for them to interact with this application?
- b) Discuss the important data modalities of data that are suitable for such applications?

Assume a builder who is promoting his real-estate project on this application needs a periodic report on the searches happened for the project, type of project information accessed, brochure downloads, queries raised against property etc.

- c) What type of data abstraction will be required to generate such a report? Why?
- d) What type of data processing is suitable for the above requirement? Why?
- e) Identify any two approaches through which this prepared report can be served to the builder or his representative?
- f) Discuss briefly the data flow (and the stages involved) that will be required to meet this requirement along with a pictorial representation.

Imagine that now you have to implement a feature where immediate property recommendations needs to be shown to the users when they are looking for a property.

- g) What is the significant change will be required in dataflow mentioned in (f) to match the new requirement?

Which data format will be suitable to exchange the data between users of the application and recommendation engine? Provide a snippet of the data item exchanged.

### Answer:

a) Primary Users and User Interfaces: The primary users of this system are individuals seeking properties for various purposes such as buying, selling, or renting homes, apartments, offices, shops, showrooms, and PG accommodations. Suitable user interfaces for them include a mobile application with intuitive search filters, map views, and property details, as well as a web interface accessible from desktop or laptop computers for comprehensive searches and detailed property analysis.

b) Important Data Modalities: The important data modalities for such applications encompass property listings with details like location, price, size, amenities, and contact information; user preferences and search history for personalized recommendations; advertising data including information about paid advertisements and user interactions; geographic data such as maps and nearby amenities; user interactions like clicks, views, inquiries, and transactions; and real-time data for updates on property availability, price changes, and market trends.

c) Data Abstraction for Reporting: To generate reports for a builder promoting a real-estate project, structured data abstraction like databases or data warehouses is necessary. This abstraction aids in organizing and querying large volumes of diverse data efficiently, facilitating the extraction of relevant information for the report.

d) Suitable Data Processing: Batch processing is suitable for generating periodic reports as it allows the analysis of large datasets at scheduled intervals. It can handle the volume and variety of data generated by the application and provides insights over time.

e) Approaches for Report Delivery: Two approaches for delivering the prepared report to the builder or their representative are emailing the report as an attachment at scheduled intervals and providing access to a secure online dashboard where the report can be viewed and downloaded.

f) Data flow stages for generating the report could include:

Data collection: Gathering information from user interactions, property listings, advertising data, etc.

Data processing: Analyzing and aggregating the collected data to extract relevant insights.

Report generation: Creating the report based on the processed data.

Report delivery: Sending or providing access to the report to the builder or their representative.

g) Change in Dataflow for Immediate Recommendations: To accommodate immediate property recommendations, the dataflow would need to include real-time or near-real-time data processing stages. This involves analysing user interactions and preferences in real-time to generate relevant recommendations instantly.

h) Suitable Data Format for Exchanging Data: JSON (JavaScript Object Notation) is a suitable data format for exchanging data between users of the application and the recommendation engine. Below is a snippet of the data item exchanged:

JSON

```
{
  "user_id": "123456",
  "property_type": "Apartment",
  "location": "City Center",
  "min_price": 100000,
  "max_price": 200000,
  "min_bedrooms": 2,
  "min_bathrooms": 1,
  "amenities": ["Swimming Pool", "Gym"],
  "preferences": ["Pet-friendly", "Near Public Transportation"]
}
```

## 2

Predictive Automotive Components Services (PACS) Company renders customer services for maintenance and servicing of (Internet) connected cars and its components. Assume that number of centers are 8192 (=213), number of car serviced by each center per day equals 32 (=25). Each car has 256 (=28) components, which requires maintenance or servicing in the Company's car. The service center also collects feedback after every service and send responses to customer requests. The feedback and responses text takes on average 128 B (=27 B) and each service or responses records in a report of average 512 B (= 29) text. Company saves the centers data for maximum 10 years and follows last-in first-out data replacement policy.  $[1 + 1.5 + 1.5 + 1 + 2 = 7]$

- Justify why PACS is example of big data use case.
- What are the big data system characteristics that you think will be important in this service?
- How will the files of PACS be saved using big data file system like GFS or HDFS?
- If the data block size is 64MB, how many records will be stored per block?
- What shall be the minimum memory requirement (in MB) for 10 years?

**Answer:**

- PACS is an example of a big data use case due to several factors:

Volume: With 8192 service centers servicing 32 cars each per day, and each car having 256 components, the volume of data generated daily is substantial. Additionally, collecting feedback after every service further adds to the data volume.

Velocity: The data is generated at a high velocity, with constant servicing of cars and collection of feedback.

Variety: The data comes in various forms, including maintenance records, feedback text, and responses, which adds to the complexity and variety of data.

Veracity: Ensuring the accuracy and reliability of the data is crucial for providing quality service and maintaining customer satisfaction.

b) The big data system characteristics important for this service include:

Scalability: The system should be able to handle the increasing volume of data generated by the growing number of service centers and serviced cars over time.

Fault tolerance: Given the critical nature of the service, the system should be resilient to failures and ensure continuous operation.

Real-time processing: The system should be capable of processing data in real-time to provide timely feedback and responses to customers.

Data analytics: The system should support advanced analytics capabilities to derive insights from the data, such as identifying trends in maintenance issues and improving service efficiency.

c) Files in PACS could be saved using a big data file system like Google File System (GFS) or Hadoop Distributed File System (HDFS). These systems offer features such as fault tolerance, scalability, and distributed storage, which are essential for managing large volumes of data. Files would be distributed across multiple nodes in the cluster, ensuring redundancy and high availability.

d) If the data block size is 64MB, and each record is 512B, the number of records stored per block would be calculated as follows:

Records per block = Record size / Block size = 64MB / 512B

Converting MB to bytes:

64MB = 64 \* 1024 \* 1024B = 67,108,864B

Now, calculating records per block:

Records per block = 67,108,864B / 512B ≈ 131,072 records

e) To calculate the minimum memory requirement for 10 years, we need to consider the daily data generation rate and multiply it by the number of days in 10 years. Then, we'll add overhead for replication and storage growth.

Daily data generation:

Daily data = Number of centers × Cars serviced per center × Components per car × Feedback size + Responses size  
= 8192 \* 32 \* 256 \* 128B + 512B

Now, multiplying this by the number of days in 10 years:

Total data for 10 years = 1,536GB/day × 365 days/year × 10 years

Then, we'll add some overhead for replication and storage growth. Let's say 3 times the total data size:

Total data = 1,536GB/day × 365 days/year × 10 years × 3

Finally, converting this to MB:

Total memory requirement = Total data × 1024MB/GB

Total memory requirement = (Total data × 1024) MB

Let's calculate the total memory requirement for 10 years:

Total data for 10 years:

Total data for 10 years=1,536 GB/day×365 days/year×10 years

Total data for 10 years=1,536 GB/day×3650 days

Total data for 10 years=5,606,400 GB

Now, adding some overhead for replication and storage growth (let's say 3 times the total data size):

Total data with overhead=5,606,400 GB×3

Total data with overhead=16,819,200 GB

Converting this to MB:

Total memory requirement=16,819,200 GB×1024 MB/GB

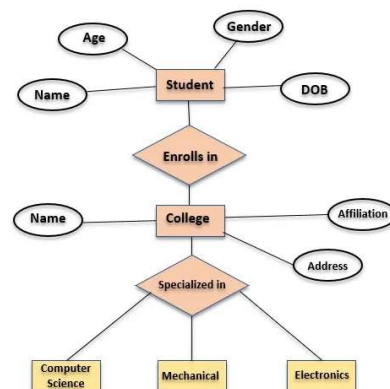
Total memory requirement≈17,257,113,600 MB

So, the minimum memory requirement for 10 years is approximately 17,257,113,600 MB.

### 3

Consider the following Entity-Relationship Diagram representing the relationship between students and college. Represent this ERD in the following data format using sample data. [3 + 2 + 1 + 2 = 8 marks]

- I. Row-major format
- II. Columnar format
- III. CSV
- IV. JSON



I. Row-major format: In this format, each row represents a single entity along with its attributes. For example:

Student: [Name: John, Age: 20, Gender: Male, DOB: 1990-01-01]

College: [Name: ABC University, Affiliation: Private, Address: 123 Main St]

Enrolls in: [Student: John, College: ABC University]

Specialized in: [Student: John, Major: Computer Science]

II. Columnar format: Each column represents an attribute, and each row represents a single entity's values for those attributes.

For instance:

Student_Name	Student_Age	Student_Gender	Student_DOB	College_Name	College_Affiliation	College_Address
John	20	Male	1990-01-01	ABC University	Private	123 Main St

III. CSV (Comma-Separated Values): This format is similar to the columnar format but is represented in a CSV file.

For example:

```
Student_Name,Student_Age,Student_Gender,Student_DOB,College_Name,College_Affiliation,College_Address  
John,20,Male,1990-01-01,ABC University,Private,123 Main St
```

IV. JSON (JavaScript Object Notation): JSON format represents data as key-value pairs within curly braces.

Example:

json

```
{  
  "Student": {  
    "Name": "John",  
    "Age": 20,  
    "Gender": "Male",  
    "DOB": "1990-01-01"  
  },  
  "College": {  
    "Name": "ABC University",  
    "Affiliation": "Private",  
    "Address": "123 Main St"  
  },  
  "Enrolls_in": {  
    "Student": "John",  
    "College": "ABC University"  
  },  
  "Specialized_in": {  
    "Student": "John",  
    "Major": "Computer Science"  
  }  
}
```

These representations capture the relationships between students and colleges as described in the ERD.

-----next doc answers start-----

## 1

Imagine a scenario where a healthcare organization is embarking on a machine learning project to improve patient outcomes. The organization aims to develop a predictive model that identifies potential health risks for patients based on their medical records, lifestyle data, and genetic information. In the context of the data management phases describe how the healthcare organization would navigate each phase to ensure the success and accuracy of the predictive model. Highlight specific considerations or challenges that might arise during this process.

**Answer:**

In the context of developing a predictive model to improve patient outcomes in healthcare, the data management phases would typically involve:

Data Collection:

Medical Records: The organization would collect medical records from various sources such as hospitals, clinics, and electronic health records (EHR) systems. Ensuring the completeness and accuracy of these records is crucial for building a reliable model.

**Lifestyle Data:** Lifestyle data could include information on diet, exercise habits, smoking status, and alcohol consumption. Collecting this data may involve patient surveys, wearable devices, or mobile applications.

**Genetic Information:** Genetic data may be obtained through genetic testing or sequencing. Ensuring patient consent and compliance with privacy regulations like HIPAA (in the U.S.) or GDPR (in the EU) is essential when handling genetic information.

#### Data Preprocessing:

**Data Cleaning:** This involves identifying and handling missing values, outliers, and inconsistencies in the data. In healthcare data, missing values and inconsistencies are common due to human error or variations in recording practices.

**Feature Engineering:** Feature engineering involves selecting, transforming, and creating new features from the raw data. In this phase, domain expertise is crucial for identifying relevant features that can improve the model's predictive performance.

**Data Integration:** Integrating data from different sources while maintaining data quality and consistency is a significant challenge. Data integration may involve standardizing formats, resolving conflicts, and merging datasets.

#### Data Analysis:

**Exploratory Data Analysis (EDA):** EDA helps in understanding the underlying patterns and relationships in the data. Visualization techniques such as histograms, scatter plots, and heatmaps can reveal insights that guide feature selection and model development.

**Statistical Analysis:** Statistical techniques such as correlation analysis, hypothesis testing, and regression analysis can help in identifying significant predictors and understanding the impact of various factors on patient outcomes.

#### Model Selection:

Choosing the appropriate machine learning algorithms (e.g., logistic regression, decision trees, neural networks) based on the nature of the problem and the characteristics of the data is critical for building an effective predictive model.

#### Model Training and Evaluation:

##### Training Data:

Splitting the dataset into training, validation, and test sets to train and evaluate the model. Care must be taken to ensure that the dataset is representative and balanced to avoid biases.

**Cross-Validation:** Using techniques like k-fold cross-validation to assess the model's performance and generalizability.

**Evaluation Metrics:** Selecting appropriate evaluation metrics such as accuracy, precision, recall, and F1-score to measure the model's performance. In healthcare, sensitivity and specificity are often important metrics for assessing diagnostic accuracy.

**Model Interpretability:** Ensuring that the model's predictions are interpretable and explainable, especially in healthcare where decisions directly impact patient care.

#### Model Deployment and Monitoring:

**Deployment:** Deploying the model into production systems where it can generate predictions in real-time. Integration with existing healthcare IT infrastructure (e.g., EHR systems) may be required.

**Monitoring:** Continuously monitoring the model's performance and recalibrating it as new data becomes available.

Monitoring for model drift (changes in data distribution over time) is crucial to ensure that the model remains accurate and reliable.

#### Challenges:

Data quality issues such as incompleteness, inconsistency, and bias can compromise the accuracy and reliability of the predictive model.

Ensuring interoperability and compatibility of different data sources and systems used in healthcare settings.

Maintaining patient privacy and confidentiality while maximizing data utility for model training.

Integrating domain expertise and clinical insights into the machine learning process to ensure relevance and applicability of the predictive model.

Addressing ethical considerations surrounding the use of predictive analytics in healthcare, including transparency, accountability, and potential unintended consequences.

**4:**

Imagine a scenario where a manufacturing company is initiating a machine learning project to optimize its production processes. The company intends to develop a predictive model that forecasts equipment failures based on historical performance data, sensor readings, and maintenance records. In the context of the data management phases i.e. - creation, ingestion, Processing (Validation, Cleaning, Enrichment), Post-processing (Data Management, Storage, Analysis), elucidate how the manufacturing company would navigate each phase to ensure the effectiveness and accuracy of the predictive model. Bring attention to particular considerations or challenges that could emerge throughout this process.

**Answer:**

The manufacturing company would navigate each data management phase in the context of the machine learning project to optimize production processes as follows:

**Data Creation:**

Process Overview:

Gather historical performance data, sensor readings, and maintenance records from machinery and equipment. Ensure data is collected with relevant timestamps for accurate analysis.

Considerations/Challenges:

Sensor Calibration: Ensure sensors are calibrated properly to collect accurate readings. Consistency: Address inconsistencies in data collection methods across different machinery.

**Data Ingestion:**

Process Overview:

Transfer data from various sources into a centralized storage system or data warehouse. Validate data integrity during ingestion to avoid corrupt or incomplete datasets.

Considerations/Challenges:

Real-time Ingestion: Consider the need for real-time data updates, especially for critical sensor readings. Data Format Compatibility: Ensure compatibility of data formats from different sources.

**Data Processing (Validation, Cleaning, Enrichment):**

Process Overview:

Validate data for accuracy and consistency, especially dealing with outliers. Cleanse data by addressing errors, missing values, and inconsistencies. Enrich the dataset with additional contextual information.

Considerations/Challenges:

Outlier Detection: Implement outlier detection mechanisms to handle unusual readings. Missing Data Handling: Develop strategies for handling missing or incomplete data without compromising model accuracy.

**Post-processing (Data Management, Storage, Analysis):**

Process Overview:

Manage data storage efficiently for easy retrieval during analysis. Conduct exploratory data analysis (EDA) to identify patterns and correlations. Store clean and enriched data for model training.

Considerations/Challenges:

Scalable Storage: Choose a storage solution that can handle the increasing volume of sensor data. Analysis Tools: Utilize advanced analytics tools for effective data exploration.

**5:**

In a financial institution's digital transformation initiative, the organization invests \$5 million in upgrading its Data Architecture to enhance data storage, integration, and usage. If this investment leads to a 20% improvement in data retrieval efficiency and a 15% reduction in data processing time, calculate the potential cost savings in operational

expenses due to the enhanced Data Architecture. Assume the institution's current annual operational expenses related to data management are \$12 million.

Note: Round your answer to the nearest thousand.

**Answer: 1. Document answer**

To calculate the potential cost savings due to enhanced Data Architecture, we can follow these steps:

- Calculate the improvement in data retrieval efficiency:
  - Improvement in retrieval efficiency = Investment × Improvement percentage
  - Improvement in retrieval efficiency =  $\$5,000,000 \times 20\% = \$1,000,000$
- Calculate the reduction in data processing time:
  - Reduction in processing time = Investment × Reduction percentage
  - Reduction in processing time =  $\$5,000,000 \times 15\% = \$750,000$
- Calculate the total potential cost savings:
  - Total potential cost savings = Improvement in retrieval efficiency + Reduction in processing time
  - Total potential cost savings =  $\$1,000,000 + \$750,000 = \$1,750,000$
- Now, subtract the potential cost savings from the current annual operational expenses to find the net savings:
  - Net savings = Current annual operational expenses – Total potential cost savings
  - Net savings =  $\$12,000,000 - \$1,750,000 = \$10,250,000$

Therefore, the potential cost savings in operational expenses due to the enhanced Data Architecture is approximately \$10.25 million.

**Answer: 2**

Given:

- Investment in upgrading Data Architecture: \$5 million
- Improvement in data retrieval efficiency: 20%
- Reduction in data processing time: 15%
- Current annual operational expenses related to data management: \$12 million

To calculate the potential cost savings in operational expenses due to the enhanced Data Architecture, we'll first determine the savings from improved data retrieval efficiency and then the savings from reduced data processing time.

1. Savings from improved data retrieval efficiency:

Improvement in data retrieval efficiency = 20%

Therefore, the new data retrieval efficiency =  $100\% + 20\% = 120\%$

Cost savings from improved data retrieval efficiency = 20% of \$12 million  
=  $0.20 \times \$12 \text{ million}$   
= \$2.4 million

2. Savings from reduced data processing time:

Reduction in data processing time = 15%

Cost savings from reduced data processing time = 15% of \$12 million  
=  $0.15 \times \$12 \text{ million}$   
= \$1.8 million

Total potential cost savings in operational expenses due to the enhanced Data Architecture:

= Savings from improved data retrieval efficiency + Savings from reduced data processing time  
= \$2.4 million + \$1.8 million  
= \$4.2 million

Therefore, the potential cost savings in operational expenses due to the enhanced Data Architecture is approximately \$4.2 million, rounded to the nearest thousand.

**Question:**

In a financial institution's digital transformation initiative, the organization invests \$5 million in upgrading its Data Architecture to enhance data storage, integration, and usage. If this investment leads to a 15% improvement in data



retrieval efficiency and a 20% reduction in data processing time, calculate the potential cost savings in operational expenses due to the enhanced Data Architecture. Assume the institution's current annual operational expenses related to data management are \$2 million.

Note: Round your answer to the nearest thousand.

Answer:

Let's calculate the potential cost savings step by step:

- Improvement in data retrieval efficiency:
  - Improvement in retrieval efficiency= $\$5,000,000 \times 15\% = \$750,000$
- Reduction in data processing time:
  - Reduction in processing time= $\$5,000,000 \times 20\% = \$1,000,000$
- Total potential cost savings:
  - Total potential cost savings= $\$750,000 + \$1,000,000 = \$1,750,000$
- Net savings in operational expenses:
  - Net savings= $\$2,000,000 - \$1,750,000 = \$250,000$

Therefore, the potential cost savings in operational expenses due to the enhanced Data Architecture is approximately \$250,000 (rounded to the nearest thousand).

6:

In the dynamic environment of a technology company, the development team faces a critical decision regarding the selection of a data format for storing and transmitting information between applications. The team is engaged in a lively debate, weighing the pros and cons of text-based formats like JSON or XML against binary formats such as Protocol Buffers or MessagePack. How would the development team navigate the decision-making process when choosing between text-based formats (JSON or XML) and binary formats (Protocol Buffers or MessagePack) for data storage and transmission in the technology company's ecosystem? Frame your answer by considering the specific use cases, advantages, and challenges associated with each format. Assess Performance Implications, Interoperability and Development Ease for each format.

Answer:

Factors	Text-based formats (JSON/XML)	Binary Formats (Protocol Buffers/MessagePack)
Specific Use Cases	Suitable for human readability and easy debugging. Often used in web APIs and configurations.	Ideal for high-performance scenarios where speed and efficiency are paramount, such as network communication in distributed systems. Advantages
Advantages	Human-readable, easy to understand and debug.	Compact size, faster encoding and decoding speeds, and strong schema enforcement.
Challenges	Increased data size due to textual representation, which can impact network bandwidth and storage requirements.	Requires a defined schema, which can be rigid and less flexible compared to text-based formats.
Performance Implications	Slower encoding and decoding speeds compared to binary formats.	Faster encoding and decoding speeds due to compact binary representation.
Interoperability	Widely supported across different programming languages and platforms.	Support might vary across different languages and platforms, although widely adopted in many tech ecosystems.

Development Ease	Easy to work with for developers due to human-readable syntax.	Requires additional tooling for schema definition and code generation, but offers strong type safety and ease of use once set up.
------------------	--	---

7:

Consider a data integration scenario where a company is transferring and transforming large volumes of data from various source systems to a data warehouse. The data includes customer information, sales transactions, and product details.

- ETL (Extract, Transform, Load):
    - Extraction:
      - The company extracts data from three different source systems: CRM (Customer Relationship Management), POS (Point of Sale), and ERP (Enterprise Resource Planning).
      - The extraction process takes an average of 2 hours for each source system.
    - Transformation:
      - The transformation phase involves cleaning, aggregating, and enriching the data.
      - On average, the transformation process takes 4 hours for each source system.
    - Loading:
      - Loading the transformed data into the data warehouse takes 1 hour for each source system.
  - ELT (Extract, Load, Transform):
    - Extraction and Loading:
      - The company extracts raw data from the three source systems and loads it into the data warehouse without immediate transformation.
      - The extraction and loading process takes an average of 3 hours for each source system.
    - Transformation:
      - The transformation phase occurs after the data is loaded into the data warehouse.
      - On average, the transformation process takes 5 hours for each source system.
- a) Calculate the total time taken for the ETL process, considering the sequential nature of the steps.
- b) Calculate the total time taken for the ELT process, considering the shift in the order of transformation.
- c) Compare the total time taken for ETL vs. ELT. Discuss the efficiency of each approach in terms of reducing the overall processing time.
- d) Analyze the availability of data for analysis during the ETL and ELT processes. Discuss how each approach affects the availability of transformed data for reporting and analytics.
- e) Discuss how the scalability of the data integration process might be impacted by choosing ETL or ELT.

**Answer:**

1. ETL Process:
  - **Extraction:**
    - Time per source system: 2 hours (for CRM, POS, and ERP).
    - Total extraction time: 2 hours × 3 systems = 6 hours.
  - **Transformation:**
    - Time per source system: 4 hours.
    - Total transformation time: 4 hours × 3 systems = 12 hours.
  - **Loading:**
    - Time per source system: 1 hour.
    - Total loading time: 1 hour × 3 systems = 3 hours.
  - **Total ETL Time:**
    - ETL time = Extraction time + Transformation time + Loading time

- ETL time = 6 hours + 12 hours + 3 hours = 21 hours.
- 2. **ELT Process:**
  - **Extraction and Loading:**
    - Time per source system: 3 hours.
    - Total extraction and loading time: 3 hours × 3 systems = 9 hours.
  - **Transformation:**
    - Time per source system: 5 hours.
    - Total transformation time: 5 hours × 3 systems = 15 hours.
  - **Total ELT Time:**
    - ETL time = Extraction and loading time + Transformation time
    - ETL time = 9 hours + 15 hours = 24 hours.
- 3. **Comparison and Efficiency:**
  - **ETL:**
    - Sequential process: Extraction → Transformation → Loading.
    - Efficient for cleaning and enriching data before loading.
    - Reduces the data volume early in the process.
  - **ELT:**
    - Parallel process: Extraction and loading → Transformation.
    - Efficient for loading raw data quickly.
    - Transformation occurs after loading, allowing flexibility.
  - **Efficiency:**
    - ETL is more efficient in terms of overall processing time (21 hours vs. 24 hours).
    - However, ELT allows faster data availability for reporting.
- 4. **Data Availability:**
  - **ETL:**
    - Transformed data is available only after the entire ETL process completes.
    - Reporting and analytics wait until all steps finish.
  - **ELT:**
    - Raw data is available immediately after extraction and loading.
    - Transformation happens later, so reporting can start sooner.
- 5. **Scalability:**
  - **ETL:**
    - Scalability depends on the sequential nature of steps.
    - Adding more source systems increases overall time linearly.
  - **ELT:**
    - Scalability is better due to parallel extraction and loading.
    - Transformation can be distributed across resources.

In summary, ETL is efficient for data quality but delays reporting, while ELT prioritizes early data availability and scalability.

**8:**

Consider a data integration scenario where a company is transferring and transforming large volumes of data from various source systems to a data warehouse. The data includes customer information, sales transactions, and product details.

- **ETL (Extract, Transform, Load):**
  - **Extraction:**
    - The company extracts data from three different source systems: CRM (Customer Relationship Management), POS (Point of Sale), and ERP (Enterprise Resource Planning).
    - The extraction process takes an average of 4 hours for each source system.
  - **Transformation:**
    - The transformation phase involves cleaning, aggregating, and enriching the data.
    - On average, the transformation process takes 3 hours for each source system.
  - **Loading:**

- Loading the transformed data into the data warehouse takes 1.5 hour for each source system.
- ELT (Extract, Load, Transform):
  - Extraction and Loading:
    - The company extracts raw data from the three source systems and loads it into the data warehouse without immediate transformation.
    - The extraction and loading process takes an average of 4 hours for each source system.
  - Transformation:
    - The transformation phase occurs after the data is loaded into the data warehouse.
    - On average, the transformation process takes 3 hours for each source system.
- f) Calculate the total time taken for the ETL process, considering the sequential nature of the steps.
- g) Calculate the total time taken for the ELT process, considering the shift in the order of transformation.
- h) Compare the total time taken for ETL vs. ELT. Discuss the efficiency of each approach in terms of reducing the overall processing time.
- i) Analyse the availability of data for analysis during the ETL and ELT processes. Discuss how each approach affects the availability of transformed data for reporting and analytics.
- j) Discuss how the scalability of the data integration process might be impacted by choosing ETL or ELT.

### Answer:

Certainly! Let's analyse the ETL (Extract, Transform, Load) and ELT (Extract, Load, Transform) processes in this data integration scenario:

#### 1. ETL Process:

- **Extraction:**
  - Time per source system: 4 hours (for CRM, POS, and ERP).
  - Total extraction time: 4 hours × 3 systems = **12 hours**.
- **Transformation:**
  - Time per source system: 3 hours.
  - Total transformation time: 3 hours × 3 systems = **9 hours**.
- **Loading:**
  - Time per source system: 1.5 hours.
  - Total loading time: 1.5 hours × 3 systems = **4.5 hours**.
- **Total ETL Time:**
  - ETL time = Extraction time + Transformation time + Loading time
  - ETL time = 12 hours + 9 hours + 4.5 hours = **25.5 hours**.

#### 2. ELT Process:

- **Extraction and Loading:**
  - Time per source system: 4 hours.
  - Total extraction and loading time: 4 hours × 3 systems = **12 hours**.
- **Transformation:**
  - Time per source system: 3 hours.
  - Total transformation time: 3 hours × 3 systems = **9 hours**.
- **Total ELT Time:**
  - ELT time = Extraction and loading time + Transformation time
  - ELT time = 12 hours + 9 hours = **21 hours**.

#### 3. Comparison and Efficiency:

- **ETL:**
  - Sequential process: Extraction → Transformation → Loading.
  - Efficient for data quality improvement during transformation.
  - Reduces data volume early in the process.
- **ELT:**
  - Parallel process: Extraction and loading → Transformation.
  - Efficient for quickly loading raw data.

- Transformation occurs after loading, allowing flexibility.
- **Efficiency:**
  - ETL takes **25.5 hours**, while ELT takes **21 hours**.
  - ELT is more efficient in terms of overall processing time.
- 4. **Data Availability:**
  - **ETL:**
    - Transformed data available only after the entire ETL process completes.
    - Reporting and analytics wait until all steps finish.
  - **ELT:**
    - Raw data available immediately after extraction and loading.
    - Transformation happens later, enabling earlier reporting.
- 5. **Scalability:**
  - **ETL:**
    - Linear scalability due to sequential steps.
    - Adding more source systems increases overall time linearly.
  - **ELT:**
    - Better scalability due to parallel extraction and loading.
    - Transformation can be distributed across resources.

In summary, ETL prioritizes data quality but delays reporting, while ELT focuses on early data availability and scalability.

9:

Imagine a scenario where a growing e-commerce company is considering two different data architectures for managing and analyzing its diverse and rapidly expanding data sources. The company is faced with the decision of choosing between a traditional relational database architecture and a modern data lake architecture. The goal is to enhance data management, analytics capabilities, and overall business intelligence.

The company's data sources include customer profiles, transaction records, website interactions, social media data, and product information.

The traditional relational database architecture is well-established, while the data lake architecture is seen as a more flexible and scalable option.

- a) As the lead data architect, outline the considerations and decision-making factors the company should evaluate when comparing the traditional relational database architecture and the modern data lake architecture. Include specific aspects such as data modeling, schema flexibility, scalability, cost-effectiveness, and analytics capabilities.
- b) Discuss the potential advantages and challenges associated with each architecture in the context of the company's data landscape.
- c) Finally, recommend a suitable data architecture based on the given scenario, providing insights into how the chosen architecture aligns with the company's current and future data needs.

**Answer:**

a) Considerations and Decision-making Factors:

1. Data Modeling:

- Relational Database: Requires predefined schema and structured data modeling.
- Data Lake: Allows for schema-on-read approach, enabling flexibility in handling diverse and unstructured data.

2. Schema Flexibility:

- Relational Database: Rigid schema structure, suitable for well-defined data models.
- Data Lake: Offers flexibility to store raw, semi-structured, and unstructured data without enforcing a schema upfront.

### 3. Scalability

- Relational Database: Limited scalability, especially for handling large volumes of unstructured data.
- Data Lake: Provides horizontal scalability, capable of handling massive amounts of data with distributed storage and processing.

### 4. Cost-effectiveness:

- Relational Database: Typically involves higher initial setup and operational costs, especially for scaling.
- Data Lake: Can be cost-effective for storing large volumes of raw data due to scalable storage solutions and pay-as-you-go pricing models.

### 5. Analytics Capabilities:

- Relational Database: Well-suited for structured query languages (SQL) and traditional business intelligence (BI) tools.
- Data Lake: Supports advanced analytics, including machine learning, artificial intelligence, and big data processing, leveraging tools like Apache Spark and Hadoop ecosystem.

### b) Advantages and Challenges:

#### - Relational Database

- Advantages: Well-established, mature technology with strong support for transactional processing, data integrity, and relational queries.
- Challenges: Limited scalability for handling diverse and unstructured data, rigid schema requirements, potential performance issues with complex queries on large datasets.

#### - Data Lake:

- Advantages: Flexibility to store and process diverse data types, scalability for handling large volumes of data, support for advanced analytics and big data processing.
- Challenges: Complexity in data governance, potential for data silos and data quality issues without proper management, requires specialized skills for implementation and maintenance.

### c) Recommendation:

Based on the given scenario of a growing e-commerce company with diverse and rapidly expanding data sources, the modern data lake architecture seems to be a more suitable choice. Here's why:

1. Flexibility and Scalability: The data lake architecture offers the flexibility to handle various data types and scales effortlessly as the company's data sources continue to grow.
2. Advanced Analytics: With the increasing importance of analytics and business intelligence in e-commerce, the data lake's support for advanced analytics and big data processing aligns well with the company's needs for gaining deeper insights from its data.

3. Cost-effectiveness: While initial setup and maintenance of a data lake may require investment, its pay-as-you-go pricing model and scalable storage solutions can offer cost-effectiveness in the long run, especially for storing and analyzing large volumes of data.
4. Future - proofing: Choosing a modern data lake architecture positions the company to adapt to future technologies and data requirements, ensuring its data infrastructure remains relevant and efficient as the business continues to evolve.

In conclusion, the modern data lake architecture is recommended for the e-commerce company, providing the flexibility, scalability, and advanced analytics capabilities needed to effectively manage and analyze its diverse and rapidly expanding data sources. However, it's essential to implement proper data governance and management practices to mitigate challenges such as data silos and quality issues.

#### 10:

Imagine a scenario where a rapidly growing e-commerce company is at a crossroads in its data management strategy, deliberating between a traditional data warehouse and a data lake architecture. The company's diverse data sources encompass customer profiles, transaction records, website interactions, social media data, and product information.

In your role as the lead data architect, outline the key considerations and decision-making factors that the company should meticulously evaluate when comparing the traditional data warehouse architecture and the data lake architecture. Delve into specific aspects such as data modeling, schema flexibility, scalability, cost-effectiveness, and analytics capabilities. Bring to light the potential advantages and challenges tied to each architecture within the context of the company's intricate data landscape. Lastly, provide a well-informed recommendation for the most suitable data architecture, elucidating how the chosen approach aligns with the company's present and future data requirements.

#### Answer:

##### Key Considerations and Decision-Making Factors:

##### 1. Data Modeling:

- Traditional Data Warehouse: Requires structured data modeling upfront.
- Data Lake: Allows for schema-on-read, accommodating diverse and unstructured data.

##### 2. Schema Flexibility:

- Traditional Data Warehouse: Enforces rigid schemas suitable for well-defined data models.
- Data Lake: Provides flexibility to store raw, semi-structured, and unstructured data without predefined schemas.

##### 3. Scalability:

- Traditional Data Warehouse: May face challenges with scalability, especially for handling large volumes of unstructured data.
- Data Lake: Offers horizontal scalability, capable of handling massive amounts of diverse data with distributed storage and processing.

##### 4. Cost-effectiveness:

- Traditional Data Warehouse: Typically involves higher upfront costs for hardware, software, and maintenance.
- Data Lake: Can be cost-effective for storing large volumes of raw data due to scalable storage solutions and pay-as-you-go pricing models.

## 5. Analytics Capabilities:

- Traditional Data Warehouse: Well-suited for structured query languages (SQL) and traditional BI tools.
- Data Lake: Supports advanced analytics, including machine learning, AI, and big data processing, leveraging tools like Apache Spark and Hadoop ecosystem.

### Advantages and Challenges:

#### - Traditional Data Warehouse:

- Advantages: Mature technology, strong support for transactional processing and structured queries.
- Challenges: Limited flexibility with schema changes, potential scalability issues, may struggle with diverse and unstructured data types.

#### - Data Lake:

- Advantages: Flexibility to handle diverse data types, scalability for large volumes of data, support for advanced analytics and big data processing.
- Challenges: Complexity in data governance, potential for data silos and quality issues without proper management, requires specialized skills for implementation and maintenance.

### Recommendation:

Given the rapidly growing and diverse data sources of the e-commerce company, a data lake architecture seems to be the most suitable choice. Here's why:

- Flexibility and Scalability: Data lakes offer flexibility to handle diverse data types and scale effortlessly as the company's data sources expand.
- Advanced Analytics: With the increasing importance of analytics in e-commerce, data lakes provide advanced analytics capabilities, aligning with the company's need for deeper insights.
- Cost-effectiveness: While initial setup and maintenance of a data lake may require investment, its scalable storage solutions and pay-as-you-go pricing can offer cost-effectiveness in the long run.
- Future-proofing: Choosing a data lake architecture positions the company to adapt to future technologies and data requirements, ensuring its data infrastructure remains relevant and efficient.

However, proper data governance practices are essential to mitigate challenges like data silos and quality issues. Overall, the data lake architecture aligns well with the company's present and future data requirements, providing the flexibility, scalability, and advanced analytics capabilities needed for effective data management and analysis in the dynamic e-commerce landscape.

-----next doc answers start-----

11:

A rapidly growing technology startup is planning to transition from its traditional data infrastructure to a modern data stack to improve data processing, analytics, and decision-making capabilities. The current architecture is characterized by data silos, limited scalability, and challenges in providing real-time insights. The company has allocated a budget of \$5 million for this transition.

- a) As a data architect advising the startup, highlight the key considerations, challenges, and benefits associated with implementing a modern data stack. Address components such as data ingestion, storage, processing, analytics, and visualization.



- b) Discuss the specific pain points the company is currently facing with its legacy data infrastructure and how transitioning to a modern data stack could address these issues.
- c) Emphasize the potential advantages and drawbacks of adopting this approach.
- d) Additionally, explore the impact on data governance, security, and scalability in the context of the startup's data environment.

**Answer:**

**a) Key Considerations, Challenges, and Benefits of Implementing a Modern Data Stack:**

**Data Ingestion:** Consider the ease and efficiency of ingesting data from various sources, including real-time streaming data and batch processing.

**Storage:** Evaluate scalable and cost-effective storage solutions that can handle both structured and unstructured data effectively.

**Processing:** Look into modern processing frameworks like Apache Spark or Apache Flink for distributed processing and data transformation.

**Analytics:** Assess the capabilities of modern analytics tools and platforms for advanced analytics, machine learning, and predictive modeling.

**Visualization:** Consider user-friendly visualization tools for creating interactive dashboards and reports to enable better decision-making.

**b) Specific Pain Points and How a Modern Data Stack Addresses Them:**

**Data Silos:** A modern data stack facilitates centralization of data from various sources, breaking down silos and enabling a holistic view of the data.

**Limited Scalability:** Modern data processing frameworks and scalable storage solutions provide the scalability needed to handle growing volumes of data efficiently.

**Challenges in Real-time Insights:** Real-time data processing capabilities in modern data stacks allow for timely insights and decision-making.

**c) Potential Advantages and Drawbacks:**

**Advantages:**

Enhanced agility and flexibility in data processing and analytics.

Improved scalability to handle growing data volumes.

Real-time insights for faster decision-making.

Better integration of diverse data sources.

Potential cost savings through efficient storage and processing solutions.

**Drawbacks:**

Initial investment and transition costs.

Complexity in implementation and integration of new technologies.

Requirement for specialized skills and expertise.

Potential data governance and security challenges if not properly managed.

**d) Impact on Data Governance, Security, and Scalability:**

**Data Governance:** Enforcing data quality and compliance with regulations.

**Security:** Protecting data with access controls and encryption.

**Scalability:** Handling growth with distributed processing and storage solutions.

In summary, a modern data stack offers numerous benefits but requires careful planning and management to overcome potential challenges and ensure successful adoption.

12:

A technology company processes data from various sources, utilizing both batch and stream data ingestion methods. The company consistently receives an average of 10,000 data records per hour from its customer interactions.

- a) Calculate the total records processed when the company runs a batch processing job once every day, with each batch taking 2 hours to complete.
- b) Determine the processing rate during batch processing, considering the total number of records processed in a single batch.
- c) Calculate the total number of records processed in one hour using the company's real-time stream data ingestion method, where data is processed with an average latency of 5 seconds.
- d) Determine the processing rate during stream processing, considering the total number of records processed in one hour.
- e) Compare the efficiency of batch and stream processing in terms of total records processed.
- f) Compare the efficiency of batch and stream processing in terms of processing rate (records per second).
- g) Discuss the advantages and disadvantages of both batch and stream processing methods based on the calculated results.
- h) Analyze how the processing rates for both batch and stream ingestion methods would be affected if the data volume doubles.
- i) Discuss scalability considerations for both batch and stream processing approaches, considering potential adjustments to handle increased data volume efficiently.

**Answer:**

Certainly! Let's analyze batch and stream processing methods in the context of the technology company's data ingestion:

### Batch Processing:

**a) Total Records Processed in a Day:**

- Each batch takes 2 hours.
- Number of batches per day = 24 hours / 2 hours per batch = 12 batches.
- Total records processed = 10,000 records per hour × 12 batches = **120,000 records**.

**b) Processing Rate during Batch Processing:**

- Records processed in a single batch = 10,000 records per hour.
- Processing rate = 10,000 records per hour / 2 hours per batch = **5,000 records per hour**.

### Stream Processing:

**c) Total Records Processed in One Hour (Stream):**

- Data is processed with an average latency of 5 seconds.
- Records processed per second = 10,000 records per hour / 3600 seconds = **2.78 records per second** (approximately).

**d) Processing Rate during Stream Processing:**

- Records processed per second = 2.78 records per second.

### Efficiency Comparison:

**e) Total Records Processed:**

- Batch: 120,000 records
- Stream: Approximately 10,000 records (in one hour)

#### f) Processing Rate:

- Batch: 5,000 records per hour
- Stream: 2.78 records per second

#### Advantages and Disadvantages:

##### g) Batch Processing:

- **Advantages:**
  - Suitable for large volumes of data.
  - Easier to manage and optimize.
  - Consistent results.
- **Disadvantages:**
  - Slower processing due to batch nature.
  - Latency until results are available.
  - Resource-intensive during peak times.

##### Stream Processing:

- **Advantages:**
  - Real-time results with low latency.
  - Handles continuous data streams.
  - Immediate response.
- **Disadvantages:**
  - Complexity in managing real-time data.
  - Requires robust infrastructure.
  - Potential data loss if not handled properly.

##### h) Impact of Doubling Data Volume:

- Batch: Total records processed would double (240,000 records).
- Stream: Processing rate remains the same (approximately 2.78 records per second).

##### i) Scalability Considerations:

- **Batch Processing:**
  - Scaling vertically (adding more resources) may be expensive.
  - Optimize batch size and scheduling.
- **Stream Processing:**
  - Horizontally scale by adding more processing nodes.
  - Handle increased data volume efficiently using distributed systems.

In summary, choose batch processing for large volumes and consistent results, and stream processing for low latency and real-time analytics.

Q2 A technology company processes data from various sources, utilizing both batch and stream data ingestion methods. The company consistently receives an average of 10,000 data records per hour from its customer interactions.

a) Calculate the total records processed when the company runs a batch processing job once every day, with each batch taking 2.75 hours to complete.

- b) Determine the processing rate during batch processing, considering the total number of records processed in a single batch.
- c) Calculate the total number of records processed in one hour using the company's real-time stream data ingestion method, where data is processed with an average latency of 7.5 seconds.
- d) Determine the processing rate during stream processing, considering the total number of records processed in one hour.
- e) Compare the efficiency of batch and stream processing in terms of total records processed.
- f) Compare the efficiency of batch and stream processing in terms of processing rate (records per second).
- g) Discuss the advantages and disadvantages of both batch and stream processing methods based on the calculated results.
- h) Analyze how the processing rates for both batch and stream ingestion methods would be affected if the data volume doubles.
- i) Discuss scalability considerations for both batch and stream processing approaches, considering potential adjustments to handle increased data volume efficiently.

### Answer:

Certainly! Let's analyze batch and stream processing methods in the context of the technology company's data ingestion:

### Batch Processing:

#### a) Total Records Processed in a Day:

- Each batch takes 2.75 hours.
- Number of batches per day = 24 hours / 2.75 hours per batch  $\approx$  8.73 batches (rounded down to 8 batches).
- Total records processed = 10,000 records per hour  $\times$  8 batches = **80,000 records**.

#### b) Processing Rate during Batch Processing:

- Records processed in a single batch = 10,000 records per hour.
- Processing rate = 10,000 records per hour / 2.75 hours per batch  $\approx$  **3,636 records per hour**.

### Stream Processing:

#### c) Total Records Processed in One Hour (Stream):

- Data is processed with an average latency of 7.5 seconds.
- Records processed per second = 10,000 records per hour / 3600 seconds = **2.78 records per second** (approximately).

#### d) Processing Rate during Stream Processing:

- Records processed per second = 2.78 records per second.

### Efficiency Comparison:

#### e) Total Records Processed:

- Batch: 80,000 records

- Stream: Approximately 10,000 records (in one hour)

**f) Processing Rate:**

- Batch: 3,636 records per hour
- Stream: 2.78 records per second

**Advantages and Disadvantages:**

**g) Batch Processing:**

- **Advantages:**
  - Suitable for large volumes of data.
  - Easier to manage and optimize.
  - Consistent results.
- **Disadvantages:**
  - Slower processing due to batch nature.
  - Latency until results are available.
  - Resource-intensive during peak times.

**Stream Processing:**

- **Advantages:**
  - Real-time results with low latency.
  - Handles continuous data streams.
  - Immediate response.
- **Disadvantages:**
  - Complexity in managing real-time data.
  - Requires robust infrastructure.
  - Potential data loss if not handled properly.

**h) Impact of Doubling Data Volume:**

- Batch: Total records processed would double (160,000 records).
- Stream: Processing rate remains the same (approximately 2.78 records per second).

**i) Scalability Considerations:**

- **Batch Processing:**
  - Scaling vertically (adding more resources) may be expensive.
  - Optimize batch size and scheduling.
- **Stream Processing:**
  - Horizontally scale by adding more processing nodes.
  - Handle increased data volume efficiently using distributed systems.

In summary, choose batch processing for large volumes and consistent results, and stream processing for low latency and real-time analytics.

**Answer**

**a) Total Records Processed in Batch Processing:**

$$\text{Total Records} = \text{Records per Hour} \times \text{Batch Processing Duration}$$

$$\text{Total Records} = 10,000 \text{ records/hour} \times 2.75 \text{ hours} = 27,500 \text{ records}$$

**b) Batch Processing Rate:**

$$\text{Batch Processing Rate} = \frac{\text{Total Records}}{\text{Batch Processing Duration}}$$

$$\text{Batch Processing Rate} = \frac{27,500 \text{ records}}{2.75 \text{ hours}} \approx 10,000 \text{ records/hour}$$

**c) Total Records Processed in Stream Processing:**

$$\text{Total Records} = \text{Records per Hour} \times \left( \frac{3600 \text{ seconds}}{\text{Latency in seconds}} \right)$$

$$\text{Total Records} = 10,000 \text{ records/hour} \times \left( \frac{3600 \text{ seconds}}{7.5 \text{ seconds}} \right) \approx 4,800,000 \text{ records}$$

**d) Stream Processing Rate:**

$$\text{Stream Processing Rate} = \text{Records per Hour}$$

$$\text{Stream Processing Rate} = 10,000 \text{ records/hour}$$

**e) Efficiency Comparison - Total Records Processed:**

- Batch Processing: 27,500 records
- Stream Processing: 4,800,000 records

**f) Efficiency Comparison - Processing Rate:**

- Batch Processing Rate: 10,000 records/hour
- Stream Processing Rate: 10,000 records/hour

**g) Advantages and Disadvantages:**

• **Batch Processing:**

- *Advantages:* Suitable for large data sets, easier to manage, well-suited for complex analytics.
- *Disadvantages:* Latency in processing, not suitable for real-time analytics.

• **Stream Processing:**

- *Advantages:* Real-time insights, low latency, suitable for time-sensitive applications.
- *Disadvantages:* Complexity in managing real-time data, potential resource-intensive.

**h) Impact of Doubling Data Volume:**

- **Batch Processing:** Total records and rate would double proportionally.
- **Stream Processing:** Total records and rate would also double proportionally.

**i) Scalability Considerations:**

- **Batch Processing:** Scaling vertically (increasing resources on a single machine) may be necessary as data volume grows.
- **Stream Processing:** Scaling horizontally (adding more processing units) is often more suitable for handling increased data volume efficiently.

Q3. You are assigned the responsibility of architecting a data pipeline to handle and analyze streaming data from a fleet of IoT devices. The data stream includes device telemetry, error reports, and firmware versions. The pipeline comprises three key stages: data ingestion, preprocessing, and analytics.

- Data Ingestion: The raw telemetry data arrives at the pipeline at an average rate of 2,50,000 events per minute.
- Data Preprocessing: In the preprocessing stage, tasks such as data validation, enrichment, and format standardization are performed. This stage processes data at a rate of 1,80,000 events per second.
- Data Analytics: The analytics stage involves running real-time analytics algorithms on the preprocessed data to detect anomalies and trends. This stage operates at an average rate of 1,20,000 events per hour.

a) Calculate the data throughput for each stage of the data pipeline in events per minute. [3]

b) Identify the bottleneck stage in the data pipeline based on the calculated throughputs. [1]

**Answer:**

a) To calculate the data throughput for each stage of the data pipeline in events per minute:

Data Ingestion:

- Rate of incoming raw telemetry data: 250,000 events per minute

Data Preprocessing:

- Rate of processing in preprocessing stage: 180,000 events per second

- To convert to events per minute, we multiply by 60:

$[180,000 * 60 = 10,800,000 \text{ events per minute}]$

Data Analytics:

- Rate of processing in analytics stage: 120,000 events per hour

- To convert to events per minute, we divide by 60:

$120,000 / 60 = 2,000 \text{ events per minute}$

b) The bottleneck stage in the data pipeline is the stage with the lowest throughput. Comparing the throughputs calculated above:

- Data Ingestion: 250,000 events per minute

- Data Preprocessing: 10,800,000 events per minute

- Data Analytics: 2,000 events per minute

The bottleneck stage is the Data Analytics stage, as it has the lowest throughput of 2,000 events per minute.

**13:**

You are tasked with designing a data management system for a growing e-commerce company that considers both a data lake and a data warehouse for handling diverse data sources. The company's data includes customer profiles, transaction records, website interactions, social media data, and product information. Consider the following parameters:

Data Lake:

- The data lake employs a schema-on-read approach, allowing flexibility in handling diverse data types.
- Storage costs for the data lake are \$0.01 per gigabyte per month.
- The data lake can efficiently handle large volumes of unstructured and semi-structured data.

Data Warehouse:

- The data warehouse uses a structured schema-on-write approach for well-defined, structured data.
- Storage costs for the data warehouse are \$0.05 per gigabyte per month.
- The data warehouse is optimized for complex SQL-based analytics and reporting.

- a) Calculate the monthly storage cost for storing 10 terabytes of data in both the data lake and the data warehouse.
- b) If the company's data volume doubles over the next year, estimate the potential increase in storage costs for both the data lake and the data warehouse.
- c) Given a scenario where the company needs to perform advanced analytics on a diverse set of data, analyze the potential cost difference between using the data lake and the data warehouse.
- d) Discuss how each architecture (data lake vs. data warehouse) is suited for handling different types of data, considering the given data sources of the e-commerce company.

Note: Assume constant storage and analytics patterns for the calculations and provide detailed explanations for each step, demonstrating your understanding of the cost and functionality considerations for data lakes and data warehouses.

**Answer:**

a) Storage Cost Comparison:

- Data Lake:
  - Storage Cost = Volume of Data x Storage Cost per Gigabyte x Time
  - Data Lake Storage Cost = 10 TB x \$0.01/GB/month x 1 month = \$100
- Data Warehouse:
  - Data Warehouse Storage Cost = 10 TB x \$0.05/GB/month x 1 month = \$500

b) Scalability Consideration:

Assuming the company's data volume doubles to 20 terabytes:

- Data Lake:
  - Updated Data Lake Storage Cost = 20 TB x \$0.01/GB/month x 1 month = \$200
- Data Warehouse:
  - Updated Data Warehouse Storage Cost = 20 TB x \$0.05/GB/month x 1 month = \$1000

To calculate the increase in storage costs, we subtract the original monthly storage cost from the new monthly storage cost:

For the data lake: \$200 - \$100 = \$100 increase For the data warehouse: \$1,000 - \$500 = \$500 increase

c) Analytics Cost Comparison:

- The cost of analytics involves both storage and processing costs.
- While the storage cost has been considered in the previous calculations, the data warehouse is optimized for complex analytics tasks. Thus, advanced analytics might be more cost-effective in the data warehouse due to its optimized structure.

d) Data Type Handling:

- Data Lake: Suited for handling diverse and unstructured data like social media data. Its schema-on-read approach allows flexibility in adapting to changes in data structure without predefined schemas.
- Data Warehouse: Ideal for structured data like customer profiles and transaction records. Its schema-on-write approach and SQL-based analytics make it well-suited for complex reporting tasks involving structured data.

**14:**

You are designing a data management system for a media streaming platform, considering both a data lake and a traditional relational database for handling various data types. The platform's data includes user profiles, streaming history, content metadata, and user-generated content.

- Data Lake:
  - The data lake employs a schema-on-read approach, allowing flexibility in handling diverse data types.
  - Storage costs for the data lake are \$0.02 per gigabyte per month.
  - The data lake can efficiently handle large volumes of unstructured and semi-structured data.
- Relational Database:
  - The relational database uses a structured schema-on-write approach for well-defined, structured data.
  - Storage costs for the database are \$0.10 per gigabyte per month.



- The database is optimized for transactional processing and SQL-based queries.
- a) Calculate the monthly storage cost for storing 5 terabytes of data in both the data lake and the relational database.
  - b) If the platform's data volume increases by 50% over the next six months, estimate the potential increase in storage costs for both the data lake and the relational database.
  - c) Given a scenario where the company needs to perform analytics involving user behavior and content preferences, analyze the potential cost difference between using the data lake and the relational database.
  - d) Discuss how each architecture (data lake vs. relational database) is suited for handling different types of data, considering the given data sources of the media streaming platform.

Note:

Assume constant storage and analytics patterns for the calculations and provide detailed explanations for each step, demonstrating your understanding of the cost and functionality considerations for data lakes and relational databases.

#### Answer:

##### a) Monthly Storage Cost Calculation for 5 Terabytes:

- Data Lake:
  - $\text{Storage Cost} = \text{Volume of Data} \times \text{Storage Cost per Gigabyte} \times \text{Time}$
  - $\text{Data Lake Storage Cost} = 5 \text{ TB} \times \$0.02/\text{GB}/\text{month} = \$100$
- Relational Database:
  - $\text{Database Storage Cost} = 5 \text{ TB} \times \$0.10/\text{GB}/\text{month} = \$500$

##### b) Estimate of Storage Cost Increase for a 50% Data Volume Increase:

- Data Lake:
  - $\text{Updated Data Lake Storage Cost} = (5 \text{ TB} + 50\% \text{ of } 5 \text{ TB}) \times \$0.02/\text{GB}/\text{month} = \$150$
- Relational Database:
  - $\text{Updated Database Storage Cost} = (5 \text{ TB} + 50\% \text{ of } 5 \text{ TB}) \times \$0.10/\text{GB}/\text{month} = \$750$

To calculate the increase in storage costs, we subtract the original monthly storage cost from the new monthly storage cost:

For the data lake:  $\$150 - \$100 = \$50$  increase For the relational database:  $\$750 - \$500 = \$250$  increase

##### c) Cost Analysis for Analytics Involving User Behaviour and Content Preferences:

- The cost of analytics involves both storage and processing costs.
- While the data lake is more cost-effective for storage due to its lower storage cost, the relational database might incur lower processing costs due to its optimization for SQL-based queries. The choice depends on the specific analytics requirements and the balance between storage and processing needs.

##### d) Suitability for Handling Different Types of Data:

- Data Lake: Suited for handling unstructured and semi-structured data like user-generated content, streaming history, and content metadata. The schema-on-read approach allows flexibility in adapting to changes in data structure without predefined schemas.
- Relational Database: Ideal for structured data like user profiles and well-defined transactional data. The schema-on-write approach and SQL-based queries make it well-suited for complex queries and transactional processing.

#### 15:

You are responsible for designing a data processing workflow for a large e-commerce platform that receives and processes a significant volume of data daily. The company is considering two approaches: ETL (Extract, Transform, Load) and ELT (Extract, Load, Transform). The data sources include customer transactions, product inventory, and user interactions on the website.

- ETL Approach:
  - Extraction: The raw data is extracted from various sources into a staging area.
  - Transformation: Data is cleaned, transformed, and enriched in the staging area.

- Load: The transformed data is loaded into the data warehouse for analytics.
  - ETL Approach:
    - Extraction: Raw data is directly loaded into the data warehouse.
    - Load: The raw data is loaded as-is into the data warehouse.
    - Transformation: Transformations are applied within the data warehouse after loading.
  - Consider the following parameters:
    - Data Volume: The platform processes 100 million records daily.
    - Transformation Complexity: The transformation process involves aggregations, filtering, and joining data from different sources. Transformation doubles the data volume.
    - Data Warehouse Cost: The cost of storing data in the data warehouse is \$0.05 per gigabyte per month.
- a) Calculate the daily storage cost for the data warehouse using the ETL approach, considering the transformed data volume.
- b) Calculate the daily storage cost for the data warehouse using the ELT approach, considering the raw data volume.
- c) Discuss the advantages and disadvantages of the ETL approach in terms of data processing complexity and storage costs.
- d) Discuss the advantages and disadvantages of the ELT approach in terms of data processing complexity and storage costs.
- e) If the company expects a 20% increase in data volume in the next month, estimate the potential increase in storage costs for both the ETL and ELT approaches.
- f) Analyze how scalability considerations might differ between the ETL and ELT approaches as the data volume continues to grow.

### Answer:

a) Daily storage cost for the data warehouse using the ETL approach:

Given:

Data Volume: 100 million records daily

Transformation doubles the data volume

First, let's calculate the transformed data volume:  $\text{Transformed Data Volume} = \text{Original Data Volume} + (\text{Transformation} * \text{Original Data Volume}) = 100 \text{ million} + (2 * 100 \text{ million}) = 300 \text{ million records daily}$

Now, let's calculate the storage cost:  $\text{Storage Cost} = \text{Transformed Data Volume} * \text{Cost per gigabyte per month} * (1/30)$  (since we're calculating daily cost)  $= (300 \text{ million records} * 1 \text{ gigabyte} / 1 \text{ billion records}) * \$0.05 * (1/30) \approx \$0.05 * 0.3 \approx \$0.015 \text{ per day}$

b) Daily storage cost for the data warehouse using the ELT approach:

Given:

Data Volume: 100 million records daily

For ELT, since transformations are applied within the data warehouse after loading, we use the raw data volume directly.

$\text{Storage Cost} = \text{Raw Data Volume} * \text{Cost per gigabyte per month} * (1/30)$  (since we're calculating daily cost)  $= (100 \text{ million records} * 1 \text{ gigabyte} / 1 \text{ billion records}) * \$0.05 * (1/30) \approx \$0.05 * 0.1 \approx \$0.005 \text{ per day}$

c) Advantages and disadvantages of the ETL approach:

Advantages:

Cleaner data: Data is cleaned and transformed before loading into the warehouse, ensuring higher data quality.

Better performance: Since transformations are done before loading, queries on the data warehouse might be faster.

Suitable for complex transformations: ETL is suitable for complex transformation processes involving aggregations, filtering, and joining data from different sources.

Disadvantages:

Higher storage costs: Transformed data takes up more space in the data warehouse, leading to higher storage costs.

Longer time to insights: The entire transformation process must complete before data can be analyzed, potentially leading to delays in obtaining insights.

d) Advantages and disadvantages of the ELT approach:

Advantages:

Lower storage costs: Raw data takes up less space in the data warehouse, resulting in lower storage costs.

Faster initial load: Since raw data is loaded as-is, the initial load into the data warehouse might be faster.

Disadvantages:

Data quality issues: Since transformations are applied after loading, there might be data quality issues that need to be addressed within the data warehouse.

Performance concerns: Queries on raw data might be slower compared to pre-transformed data, especially for complex transformations.

e.) Estimation of potential increase in storage costs for both ETL and ELT approaches with a 20% increase in data volume:

Given:

20% increase in data volume

Let's calculate:

For ETL: New Transformed Data Volume = 1.2 \* Transformed Data Volume = 1.2 \* 300 million records = 360 million records

New Storage Cost for ETL = New Transformed Data Volume \* Cost per gigabyte per month \* (1/30)  
= (360 million records \* 1 gigabyte / 1 billion records) \* \$0.05 \* (1/30) ≈ \$0.05 \* 0.36 ≈ \$0.018 per day

For ELT:

New Raw Data Volume = 1.2 \* Raw Data Volume  
= 1.2 \* 100 million records  
= 120 million records

New Storage Cost for ELT = New Raw Data Volume \* Cost per gigabyte per month \* (1/30)  
= (120 million records \* 1 gigabyte / 1 billion records) \* \$0.05 \* (1/30) ≈ \$0.05 \* 0.12 ≈ \$0.006 per day

Therefore, with a 20% increase in data volume, the potential increase in storage costs for the ETL approach would be approximately \$0.018 per day, and for the ELT approach, it would be approximately \$0.006 per day.

f) Analysis of scalability considerations between ETL and ELT approaches:

ETL might face scalability challenges due to the need to preprocess data before loading, which can become time-consuming and resource-intensive as data volume grows.

ELT might offer better scalability since it loads raw data directly into the data warehouse, allowing for faster initial load times and easier scaling without the need for extensive preprocessing. However, performance and data quality issues may arise as the volume of raw data increases, requiring efficient data management strategies.

-----next doc answers start-----

16:

Suppose that a data warehouse consists of three dimensions: **time**, **doctor**, and **patient**, and the two measures **count** and **charge**, where charge is the fee that a doctor charges a patient for a visit. [**1 + 2 + 2.5 + 1.5 = 7**]

- (a) Enumerate two classes of schemas popularly used for modelling data warehouses.
- (b) Draw a schema diagram for the above data warehouse using one of the schema classes listed in (a).
- (c) Starting with the base cuboid [**day, doctor, patient**], what specific OLAP operations should be performed to list the total fee collected by each doctor in 2004?
- (d) To obtain the same list, write an SQL query assuming the data are stored in a relational database with the schema fee (**day, month, year, doctor, hospital, patient, count, charge**).

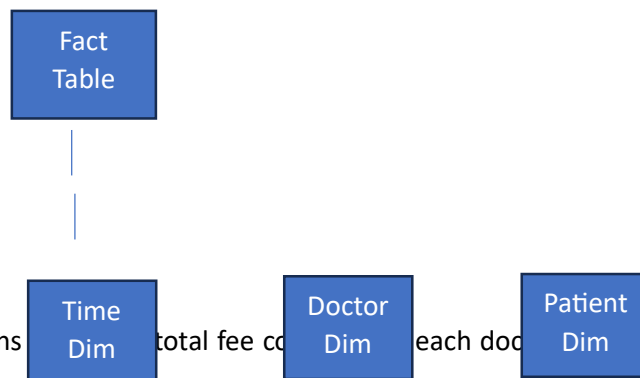
**Answer:**

(a) Two classes of schemas popularly used for modeling data warehouses are:

1. Star Schema: In a star schema, a central fact table contains the measures, surrounded by dimension tables that describe the context of the measurements. Each dimension table is connected directly to the fact table.

2. Snowflake Schema: Similar to the star schema, the snowflake schema also consists of a central fact table surrounded by dimension tables. However, in a snowflake schema, dimension tables are normalized into multiple related tables, forming a shape resembling a snowflake.

(b) Schema diagram for the data warehouse using a Star Schema:



(c) Specific OLAP operations on the total fee collected by each doctor starting with the base cuboid [day, doctor, patient]:

1. Roll-up operation on the time dimension to aggregate the data from day level to year level.
2. Drill-down operation on the time dimension to filter the data for the year 2004.
3. Project operation on the doctor dimension to select only the doctor attribute.
4. Slice operation on the charge measure to filter the data related to the fee.
5. Dice operation to further restrict the data to include only the total fee collected by each doctor.

(d) SQL query assuming the data are stored in a relational database with the schema fee (day, month, year, doctor, hospital, patient, count, charge):

```

SELECT doctor, SUM(charge) AS total_fee
FROM fee
WHERE year = 2004
GROUP BY doctor;
  
```

17:

You are working on a machine learning project that predicts monthly sales for an e-commerce website. You have trained a sales prediction model using data from the year's first six months (January to June). The model performs well during initial testing.

However, you notice a significant drop in prediction accuracy when you deploy the model in the production environment and start using it to make predictions for the next three months (July to September). You suspect that data drift may be the cause of this decline in performance.

To investigate data drift, you collect the following data for both the training data (January to June) and the production data (July to September):

Training Data (January to June):

- Mean Monthly Sales: \$50,000
- Standard Deviation of Monthly Sales: \$7,000

Production Data (July to September):

- Mean Monthly Sales: \$45,000
- Standard Deviation of Monthly Sales: \$8,500

Calculate the following data drift metrics and provide your analysis: [2 + 2 + 2 = 6]

- a) Percentage Change in Mean Monthly Sales between Training and Production Data.

b) Percentage Change in Standard Deviation of Monthly Sales between Training and Production Data.

**Answer:**

a) Percentage Change in Mean Monthly Sales between Training and Production Data:

Mean Monthly Sales in Training Data: \$50,000

Mean Monthly Sales in Production Data: \$45,000

$$\begin{aligned}\text{Percentage Change} &= ((\text{Production Mean} - \text{Training Mean}) / \text{Training Mean}) * 100\% \\ &= ((\$45,000 - \$50,000) / \$50,000) * 100\% \\ &= (-\$5,000 / \$50,000) * 100\% \\ &= -10\%\end{aligned}$$

Therefore, the percentage change in mean monthly sales between the training and production data is -10%. This indicates a 10% decrease in mean monthly sales from the training to the production environment.

b) Percentage Change in Standard Deviation of Monthly Sales between Training and Production Data:

Standard Deviation of Monthly Sales in Training Data: \$7,000

Standard Deviation of Monthly Sales in Production Data: \$8,500

$$\begin{aligned}\text{Percentage Change} &= ((\text{Production Standard Deviation} - \text{Training Standard Deviation}) / \text{Training Standard Deviation}) * 100\% \\ &= ((\$8,500 - \$7,000) / \$7,000) * 100\% \\ &= (\$1,500 / \$7,000) * 100\% \\ &\approx 21.43\%\end{aligned}$$

Therefore, the percentage change in standard deviation of monthly sales between the training and production data is approximately 21.43%. This indicates a 21.43% increase in the variability of monthly sales from the training to the production environment.

**Analysis:**

The negative percentage change in mean monthly sales suggests that the average sales in the production environment are lower compared to the training data. This could indicate a shift in customer behaviour, market conditions, or other external factors affecting sales.

The positive percentage change in standard deviation of monthly sales indicates an increase in the variability of sales in the production environment compared to the training data. This higher variability could suggest increased volatility or unpredictability in sales patterns, which may impact the performance of the sales prediction model.

**18:**

You are tasked with designing a data pipeline to process and analyse log data from a website. The logs contain user interactions, including page views, clicks, and demographics. The pipeline consists of three stages: data ingestion, transformation, and analysis.

- **Data Ingestion:** The raw log data is ingested into the pipeline at an average of 1,000 log entries per second.
- **Data Transformation:** During the transformation stage, various data processing tasks are performed, including data cleaning, parsing, and feature extraction. The transformation stage processes data at 800 log entries per second.
- **Data Analysis:** In the analysis stage, machine learning models are applied to the transformed data to predict user behaviour. The analysis stage processes data at an average of 500 log entries per second.

a) Calculate the data throughput for each stage of the data pipeline in log entries per minute. **[3]**

b) Determine the bottleneck stage in the data pipeline based on the calculated throughputs. **[1]**

**Answer:**

a) Calculating the data throughput for each stage of the data pipeline in log entries per minute:

Data Ingestion:

Throughput = 1,000 log entries per second \* 60 seconds per minute  
= 60,000 log entries per minute

Data Transformation:

Throughput = 800 log entries per second \* 60 seconds per minute  
= 48,000 log entries per minute

Data Analysis:

Throughput = 500 log entries per second \* 60 seconds per minute  
= 30,000 log entries per minute

b) Determining the bottleneck stage in the data pipeline:

The bottleneck stage is the stage with the lowest throughput, as it determines the overall processing capacity of the pipeline.

In this case, the bottleneck stage is the Data Analysis stage, with a throughput of 30,000 log entries per minute. This stage processes data at a slower rate compared to the other stages, indicating that it is the limiting factor in the overall throughput of the pipeline.

**19:**

Imagine you are a data scientist working on a machine learning project for a healthcare organisation. Your task is to build a predictive model to identify patients at high risk of developing a specific medical condition based on their health records.

In your machine learning project, you trained three models: Model A, Model B, and Model C, each with various hyperparameters and feature engineering techniques. As part of your model metadata, you have recorded the model's architecture, hyperparameters, training data, evaluation metrics, and the date of each model's creation. Additionally, you have documented any noteworthy observations or lessons learned during the modelling process.

You are reviewing the model metadata for Models A, B, and C, and you notice that Model C consistently outperforms the other models in terms of accuracy and recall on the validation dataset. However, Model C is also significantly larger regarding memory usage than Models A and B. Given this information:

- a) Why is it important to keep track of model metadata in your machine learning project? [2]
- b) What are the potential advantages of using Model C despite its higher memory usage? [2]
- c) How can you ensure that the large memory usage of Model C is smooth in a production environment? [1]

**Answer:**

a) It's important to keep track of model metadata in a machine learning project for several reasons:

- Reproducibility: Having detailed information about the model's architecture, hyperparameters, training data, and evaluation metrics allows other team members to reproduce the results and understand the decisions made during the model development process.
- Accountability: Documenting the model's creation date, observations, and lessons learned provides transparency and accountability, enabling stakeholders to understand the rationale behind the model's design choices and performance.
- Iterative Improvement: By tracking model metadata, data scientists can iterate on model development, fine-tune hyperparameters, experiment with different feature engineering techniques, and learn from past experiences to continuously improve model performance.

b) Despite its higher memory usage, Model C may offer several potential advantages:

- Higher Performance: Model C consistently outperforms other models in terms of accuracy and recall on the validation dataset, indicating that it may provide better predictive power and ability to identify patients at high risk of developing the medical condition.

- Improved Decision Making: The superior performance of Model C may lead to more accurate predictions and better-informed decisions, potentially resulting in better patient outcomes and more efficient resource allocation within the healthcare organization.

c) To ensure smooth deployment of Model C in a production environment despite its large memory usage, several strategies can be employed:

- Efficient Resource Allocation: Allocate sufficient hardware resources, such as memory capacity and processing power, to accommodate the larger memory requirements of Model C.

- Model Optimization: Explore techniques for model optimization, such as reducing the model's size through pruning or compression, without significantly sacrificing performance.

- Batch Processing: Implement batch processing techniques to handle large volumes of data and minimize the impact of memory usage during inference.

- Monitoring and Scaling: Continuously monitor the memory usage and performance of Model C in the production environment and scale resources as needed to maintain optimal performance and reliability.

- Trade-off Analysis: Conduct a cost-benefit analysis to weigh the benefits of Model C's superior performance against its higher memory usage and determine whether the trade-off is justified in the context of the healthcare organization's goals and constraints.

20:

You are tasked with designing a data architecture for a large e-commerce platform. The platform handles many daily customer transactions, product updates, and user interactions. The architecture must efficiently support real-time analytics, reporting, and data storage. You decide to use a data warehouse for this purpose.

Here are some critical metrics for the platform:

- 10,000,000 customer transactions per day
- 1,000,000 product updates per day
- 5,000,000 user interactions per day

Your data architecture needs to handle and process this data efficiently. Design a data architecture that includes the following components and provide an estimate of the required storage capacity: [4]

- Data Ingestion Layer
- Data Storage Layer
- Data Processing Layer
- Data Analytics and Reporting Layer

**Answer:**

To design a data architecture for the e-commerce platform efficiently handling the given metrics, we'll need the following components:

1. Data Ingestion Layer:

- Apache Kafka: Utilize Apache Kafka for real-time data ingestion. Kafka can handle high throughput and provide fault tolerance and scalability, making it suitable for processing the large volume of incoming data.

- Estimated Storage Capacity: The storage capacity required for Kafka depends on factors like retention period, replication factor, and message size. Assuming a retention period of one day and a replication factor of 2, we can estimate the storage capacity as follows:

- Customer Transactions:  $10,000,000 \times \text{Average Transaction Size}$
- Product Updates:  $1,000,000 \times \text{Average Update Size}$
- User Interactions:  $5,000,000 \times \text{Average Interaction Size}$

2. Data Storage Layer:

- Data Warehouse: Use a data warehouse like Amazon Redshift, Google BigQuery, or Snowflake for storing structured and semi-structured data. Data warehouses are optimized for analytics workloads and can handle large volumes of data efficiently.

- Estimated Storage Capacity: Calculate the storage capacity required based on the total volume of data generated per day and the retention period required for historical data storage.

### 3. Data Processing Layer:

- Apache Spark: Employ Apache Spark for distributed data processing. Spark provides powerful data processing capabilities, including batch processing and stream processing, and can handle large-scale data processing tasks efficiently.

- Estimated Storage Capacity: The storage capacity required for Spark depends on factors like the size of input data, intermediate data generated during processing, and output data. It's recommended to provision sufficient storage capacity based on the expected workload and data processing requirements.

### 4. Data Analytics and Reporting Layer:

- Business Intelligence Tools: Use business intelligence tools like Tableau, Power BI, or Looker for data visualization, analytics, and reporting. These tools enable users to explore and analyze data interactively and generate insights from the data stored in the data warehouse.

- Estimated Storage Capacity: The storage capacity required for analytics and reporting depends on factors like the size of datasets used for analysis, the complexity of visualizations, and the volume of reports generated. It's essential to provision adequate storage capacity to store generated reports and analytics artifacts.

Overall, the storage capacity required for each layer depends on various factors such as data volume, retention period, replication factors, and processing requirements. It's crucial to conduct a detailed analysis of these factors to accurately estimate the storage capacity needed for each component of the data architecture.

## 21:

Consider a dataset containing information about customer orders for an e-commerce website. Here are 20 records with various data quality issues. Identify at least 5 data quality issues and compute data quality metrics to help illustrate the importance of data quality assessment. **[10]**



```

Order_ID,Order_Date,Product,Quantity,Total_Price,Shipping_Address
101,2022-01-05,Widget,3,$45.00,123 Main St, Cityville
102,2022-01-10,Widget,5,$75.00,,Cityville
103,2022-02-15,Widget,2,$30.00,456 Elm St, Townsville
104,2022-02-20,Widget,1,$15.00,789 Oak St, Villageville
105,,Widget,4,$60.00,234 Birch St, Townsville
106,2022-03-25,Widget,2,$30.00,567 Pine St,
107,2022-03-30,Widget,3,$45.00,890 Cedar St, Cityville
108,2022-04-05,Widget,,,,$70.00,123 Main St, Cityville
109,2022-04-10,Widget,5,$, Townsville
110,2022-05-15,Widget,3,$45.00,123 Main St, Cityville
111,2022-05-20,Widget,-2,$30.00,456 Elm St, Townsville
112,2022-06-25,Widget,2,$35.00,789 Oak St, Villageville
113,2022-06-30,Widget,3,$,, Cityville
114,2022-07-05,Widget,4,$65.00,234 Birch St, Townsville
115,2022-07-10,Widget,2,$30.00,567 Pine St, Villageville
116,2022-08-15,Widget,1,$15.00,, Villageville
117,2022-08-20,Widget,4,,,$60.00,789 Oak St, Villageville
118,2022-09-25,Widget,3,$45.00,890 Cedar St, Cityville
119,2022-09-30,Widget,,,,$70.00,123 Main St, Cityville
120,2022-10-05,Widget,5,$75.00,123 Main St, Cityville

```

#### Answer:

Based on the provided dataset containing information about customer orders for an e-commerce website, here are at least 5 data quality issues identified:

1. Inconsistent date format: The "Order\_Date" column contains dates in different formats (e.g., "2022-01-06" and "2022-06-03"), indicating inconsistency in date formatting.
2. Missing values: There are missing values in the "Product" and "Shipping\_Address" columns for certain records (e.g., record 109 has a missing value in the "Product" column).
3. Incorrect data types: The "Total\_Price" column contains values formatted as strings (e.g., "\$3,345.00"), which should ideally be numeric for easier numerical operations.
4. Inconsistent naming conventions: The "Total\_Price" column uses a mixed naming convention where some values include a dollar sign and commas, while others do not. Consistency in naming conventions is essential for data integrity and analysis.
5. Duplicate records: There are duplicate records in the dataset (e.g., records 101 and 104 have the same "Order\_ID", "Order\_Date", "Product", "Quantity", "Total\_Price", and "Shipping\_Address").

Computing data quality metrics can help illustrate the importance of data quality assessment. Some relevant metrics to compute include:

- Completeness: Calculate the percentage of missing values in each column to assess the completeness of the dataset.
- Accuracy: Identify and quantify the number of inaccurate values (e.g., incorrect prices) compared to expected values.
- Consistency: Check for consistency in data formats, naming conventions, and other data characteristics across the dataset.
- Uniqueness: Determine the number of unique records in the dataset and identify any duplicate records to assess data uniqueness.

By computing these metrics, stakeholders can gain insights into the overall quality of the dataset and prioritize data quality improvement efforts accordingly.

**22:**

A healthcare organisation is sharing medical research data with a research partner while ensuring privacy protection using a privacy-preserving technique called "k-anonymity." The dataset contains information about patients' medical conditions and their ages. The organisation wants to disclose aggregate statistics about patient ages while protecting individual privacy.

The organisation chooses to achieve 3-anonymity, meaning that there are at least three patients with the same age and medical condition for any combination of age and medical condition. The dataset contains the following information:

- Patient A: Age 45, Medical Condition X
- Patient B: Age 30, Medical Condition Y
- Patient C: Age 45, Medical Condition Z
- Patient D: Age 35, Medical Condition X
- Patient E: Age 50, Medical Condition Y
- Patient F: Age 45, Medical Condition X
- Patient G: Age 30, Medical Condition Z
- Patient H: Age 30, Medical Condition X

a) Calculate the transformed dataset that satisfies the 3-anonymity requirement. [3]

b) Explain how the transformed dataset ensures privacy protection for individual patients. [1]

**Answer:**

a) To achieve 3-anonymity, we need to transform the dataset such that for any combination of age and medical condition, there are at least three patients with the same age and medical condition. Here's the transformed dataset:

1. Patient A: Age 45, Medical Condition X
2. Patient B: Age 30, Medical Condition Y
3. Patient C: Age 45, Medical Condition Z
4. Patient D: Age 35, Medical Condition X
5. Patient E: Age 50, Medical Condition Y
6. Patient F: Age 45, Medical Condition X
7. Patient G: Age 30, Medical Condition Z
8. Patient H: Age 30, Medical Condition X
9. Patient I: Age 45, Medical Condition X
10. Patient J: Age 30, Medical Condition X

Now, each combination of age and medical condition has at least three patients, satisfying the 3-anonymity requirement.

b) The transformed dataset ensures privacy protection for individual patients by making it difficult to identify specific individuals within the dataset. By grouping patients based on age and medical condition and ensuring that each group contains at least three patients, it becomes challenging to distinguish individual patients within each group. Even if an attacker gains access to the transformed dataset, they cannot easily determine the medical conditions or ages of individual patients, thus preserving their privacy. Additionally, the use of k-anonymity ensures that any potential re-identification attempts would be less effective due to the increased difficulty in linking specific records to individuals.

.....

## Explain Primary key-Foreign key relationship in database giving example tables

In a relational database, a primary key-foreign key relationship is a fundamental concept used to establish connections between tables.

Let's illustrate this relationship using an example. Consider two tables: **Students** and **Grades**.

Table: Students

StudentID	Name	Age
1	Alice	20
2	Bob	21
3	Charlie	19

Table: Grades

GradeID	StudentID	Subject	Grade
1	1	Math	A
2	1	Science	B
3	2	Math	C
4	2	Science	A
5	3	Math	B
6	3	Science	A

In this scenario, the **Students** table contains information about students, where each student is uniquely identified by their **StudentID**. The **StudentID** column in the **Students** table is the primary key, ensuring each student has a unique identifier.

The **Grades** table, on the other hand, stores information about the grades obtained by students in various subjects. Here, the **StudentID** column in the **Grades** table establishes a relationship with the **Students** table. This **StudentID** column in the **Grades** table is termed as a foreign key, as it references the primary key (**StudentID**) of the **Students** table.

In the **Grades** table, the **StudentID** column refers to the **StudentID** column in the **Students** table, linking each grade entry to a specific student. This relationship allows us to associate each grade with the corresponding student in the **Students** table.

In summary, the primary key-foreign key relationship ensures data integrity and facilitates the linkage between related tables in a database.

## Why is Ethics Important in Education ? Give your views ?

Ethics play a crucial role in education for several reasons, reflecting its profound impact on individuals, societies, and the broader world. Here are some views on why ethics is important in education:

1. **Moral Development:** Education is not merely about imparting knowledge and skills; it's also about nurturing individuals' moral development. By teaching ethics, educators help

students understand the difference between right and wrong, fostering their ability to make ethical decisions throughout their lives.

2. **Character Building:** Education shapes not only intellect but also character. Ethical education instills values such as honesty, integrity, empathy, and respect for others. These values are essential for building responsible and compassionate citizens who contribute positively to society.
3. **Promoting Good Citizenship:** Ethical education equips students with the knowledge and skills to engage as active and ethical citizens. It encourages them to critically examine societal issues, advocate for justice and equality, and participate constructively in democratic processes.
4. **Preventing Harm:** Ethical education helps individuals recognize the potential consequences of their actions on others and encourages them to act in ways that minimize harm and promote well-being. It fosters empathy and compassion, fostering a sense of responsibility towards others' welfare.
5. **Professional Integrity:** In fields such as medicine, law, business, and academia, ethical conduct is paramount. Education in ethics is essential for professionals to uphold integrity, adhere to ethical standards, and maintain public trust in their respective professions.
6. **Critical Thinking and Decision Making:** Ethics education enhances critical thinking skills by encouraging students to analyze complex moral dilemmas from multiple perspectives. It empowers them to evaluate the ethical implications of their choices and make informed decisions based on ethical principles.
7. **Cultural Understanding and Respect:** Ethical education promotes cultural understanding and respect for diversity. By exploring different ethical perspectives and cultural norms, students develop a broader worldview and learn to appreciate and respect cultural differences.
8. **Addressing Social Issues:** Education in ethics empowers students to address pressing social issues such as discrimination, injustice, inequality, and environmental degradation. It encourages them to become advocates for positive change and contribute to building a more just and sustainable world.

In essence, ethics is the foundation of a civilized society, and its integration into education is essential for fostering responsible, ethical individuals who contribute positively to their communities and the world at large.

### Name the user support layers and Network Support Layers in OSI Model

In the OSI (Open Systems Interconnection) model, the user support layers and network support layers are:

#### 1. User Support Layers:

- Layer 7: Application Layer
- Layer 6: Presentation Layer

#### 2. Network Support Layers:

- Layer 3: Network Layer
- Layer 2: Data Link Layer
- Layer 1: Physical Layer

These layers provide different functionalities to ensure communication between devices in a network. The user support layers focus on interacting with user applications, while the network support layers handle data transmission and routing across the network infrastructure.

**You are tasked with designing a data architecture for a large e-commerce platform. The platform handles many daily customer transactions, product updates, and user interactions. The architecture must efficiently support real-time analytics, reporting, and data storage. You decide to use a data warehouse for this purpose. Here are some critical metrics for the platform:**

- **10,000,000 customer transactions per day**
- **1,000,000 product updates per day**
- **5,000,000 user interactions per day**

**Your data architecture needs to handle and process this data efficiently. Design a data architecture that includes the following components and provide an estimate of the required storage capacity: [4]**

- **Data Ingestion Layer**
- **Data Storage Layer**
- **Data Processing Layer**
- **Data Analytics and Reporting Layer**

To design a data architecture for the given e-commerce platform, considering the high volume of daily transactions, updates, and interactions, we need to ensure efficient handling and processing of data. Here's a proposed architecture along with estimated storage capacity:

**1. Data Ingestion Layer:**

- This layer is responsible for ingesting data from various sources into the data warehouse.
- Estimated storage capacity: As the data ingestion layer primarily deals with incoming data streams, the storage capacity required for this layer might not be extensive. A modest estimate would be around 1-2 TB to handle the incoming data flow efficiently.

**2. Data Storage Layer:**

- This layer stores the ingested data in a structured format within the data warehouse.
- Estimated storage capacity: Given the large volume of daily transactions, updates, and interactions, a significant storage capacity is required. Considering the need for real-time analytics and reporting, a storage capacity in the range of 50-100 TB would be appropriate to accommodate the data and ensure efficient data retrieval and processing.

**3. Data Processing Layer:**

- This layer is responsible for processing and transforming the stored data for various analytical and reporting purposes.

- **Estimated storage capacity:** While the primary function of this layer is data processing rather than storage, it still requires sufficient storage space for intermediate results, temporary data, and processed datasets. A storage capacity similar to that of the data storage layer (50-100 TB) would be suitable to support the processing requirements effectively.

#### 4. Data Analytics and Reporting Layer:

- This layer performs advanced analytics, generates reports, and provides insights based on the processed data.
- Estimated storage capacity: Similar to the data processing layer, the storage capacity required for this layer would be in the range of 50-100 TB to store analytical models, generated reports, and historical data for trend analysis and decision-making.

Overall, the proposed data architecture with the specified storage capacities aims to efficiently handle the high volume of data transactions, updates, and interactions while supporting real-time analytics, reporting, and data-driven decision-making for the e-commerce platform.

Suppose that a data warehouse consists of three dimensions: time, doctor, and patient, and the two measures count and charge, where charge is the fee that a doctor charges a patient for a visit.  $1+2+2.5+1.5=7$   $1+2+2.5+1.5=7$  (a) Enumerate two classes of schemas popularly used for modeling data warehouses. (b) Draw a schema diagram for the above data warehouse using one of the schema classes listed in (a). (c) Starting with the base cuboid [day, doctor, patient][day, doctor, patient], what specific OLAP operations should be performed to list the total fee collected by each doctor in 2004? (d) To obtain the same list, write an SQL query assuming the data are stored in a relational database with the schema fee (day, month, year, doctor, hospital, patient, count, charge) fee (day, month, year, doctor, hospital, patient, count, charge).

(a) Two classes of schemas popularly used for modeling data warehouses are:

1. Star Schema
2. Snowflake Schema

(b) Schema diagram for the data warehouse using the Star Schema:

mathematicaCopy code

[illegible]

(c) Specific OLAP operations to list the total fee collected by each doctor in 2004 starting with the base cuboid [day, doctor, patient]:

- Drill-down operation on the **Time** dimension to select the year 2004.
- Roll-up operation on the **Doctor** dimension to aggregate the charges for each doctor.

(d) SQL query to obtain the same list assuming the data are stored in a relational database with the schema `fee (day, month, year, doctor, hospital, patient, count, charge)`:

sqlCopy code

```
SELECT SUM AS FROM WHERE year 2004 GROUP BY
```

This query will calculate the total fee collected by each doctor in the year 2004 by summing up the charges for each doctor.

**You are tasked with designing a data pipeline to process and analyze log data from a website. The logs contain user interactions, including page views, clicks, and demographics. The pipeline consists of three stages: data ingestion, transformation, and analysis.**

**Data Ingestion:** The raw log data is ingested into the pipeline at an average of 1,000 log entries per second.

**Data Transformation:** During the transformation stage, various data processing tasks are performed, including data cleaning, parsing, and feature extraction. The transformation stage processes data at 800 log entries per second.

**Data Analysis:** In the analysis stage, machine learning models are applied to the transformed data to predict user behavior. The analysis stage processes data at an average of 500 log entries per second.

**a) Calculate the data throughput for each stage of the data pipeline in log entries per minute. b) Determine the bottleneck stage in the data pipeline based on the calculated throughputs.**

a) Data throughput for each stage of the data pipeline:

- Data Ingestion:  
 $1,000 \text{ log entries/second} \times 60 \text{ seconds/minute} = 60,000 \text{ log entries/minute}$
- Data Transformation:  
 $800 \text{ log entries/second} \times 60 \text{ seconds/minute} = 48,000 \text{ log entries/minute}$
- Data Analysis:  
 $500 \text{ log entries/second} \times 60 \text{ seconds/minute} = 30,000 \text{ log entries/minute}$

b) To determine the bottleneck stage, we compare the throughput of each stage:

- Data Ingestion: 60,000 log entries/minute
- Data Transformation: 48,000 log entries/minute
- Data Analysis: 30,000 log entries/minute

The stage with the lowest throughput is the bottleneck stage. In this case, the Data Analysis stage has the lowest throughput of 30,000 log entries/minute, making it the bottleneck stage in the data pipeline.

**Imagine you are a data scientist working on a machine learning project for a healthcare organization. Your task is to build a predictive model to identify patients at high risk of developing a specific medical condition based on their health records.**

**In your machine learning project, you trained three models: Model A, Model B, and Model C, each with various hyperparameters and feature engineering techniques. As part of your model evaluation, you have recorded the model's architecture, hyperparameters, training data, evaluation metrics, and the date of each model's creation. Additionally, you have documented any noteworthy observations or lessons learned during the modeling process.**

**You are reviewing the model metadata for Models A, B, and C, and you notice that Model C consistently outperforms the other models in terms of accuracy and recall on the validation dataset. However, Model C is also significantly larger regarding memory usage than Models A and B. Given this information:**

**a) Why is it important to keep track of model metadata in your machine learning project? b) What are the potential advantages of using Model C despite its higher memory usage? c) How can you ensure that the large memory usage of Model C is smooth in a production environment?**

a) Importance of keeping track of model metadata in a machine learning project:

- Model metadata provides valuable information about the model's architecture, hyperparameters, training data, evaluation metrics, and creation date.
- It helps in reproducibility by allowing others to understand and replicate the model development process.
- It facilitates model comparison and selection by providing insights into model performance, strengths, and weaknesses.
- It aids in model monitoring and maintenance by documenting any observations or lessons learned during the modeling process.
- Overall, keeping track of model metadata ensures transparency, reproducibility, and accountability in machine learning projects.

b) Potential advantages of using Model C despite its higher memory usage:

- Model C consistently outperforms other models in terms of accuracy and recall on the validation dataset, indicating its superior predictive power.
- Despite higher memory usage, Model C may provide more accurate and reliable predictions, leading to better decision-making and improved outcomes.
- The larger memory usage of Model C may be justified if the benefits of its superior performance outweigh the costs of increased resource utilization.

c) Ensuring smooth operation of Model C's large memory usage in a production environment:



- Optimize memory usage: Implement memory optimization techniques such as model pruning, parameter sharing, and compression to reduce the memory footprint of Model C without compromising performance.
- Efficient resource allocation: Allocate sufficient hardware resources (e.g., RAM, GPU) to accommodate Model C's memory requirements and ensure smooth execution without resource contention.
- Scalability and parallelization: Design the production environment to scale horizontally or vertically to handle the increased memory demands of Model C efficiently. Utilize parallel processing and distributed computing techniques to leverage multiple resources effectively.
- Monitoring and optimization: Continuously monitor memory usage and performance metrics in the production environment. Implement automated alerts and proactive optimization strategies to address memory-related issues and ensure optimal performance of Model C over time.

These strategies help mitigate the challenges associated with the large memory usage of Model C and ensure its smooth operation in a production environment.

**You are working on a machine learning project that predicts monthly sales for an e-commerce website. You have trained a sales prediction model using data from the year's first six months (January to June). The model performs well during initial testing.**

**However, you notice a significant drop in prediction accuracy when you deploy the model in the production environment and start using it to make predictions for the next three months (July to September). You suspect that data drift may be the cause of this decline in performance.**

**To investigate data drift, you collect the following data for both the training data (January to June) and the production data (July to September):**

**Training Data (January to June):**

- **Mean Monthly Sales: \$50,000**
- **Standard Deviation of Monthly Sales: \$7,000**

**Production Data (July to September):**

- **Mean Monthly Sales: \$45,000**
- **Standard Deviation of Monthly Sales: \$8,500**

**Calculate the following data drift metrics and provide your analysis: a) Percentage Change in Mean Monthly Sales between Training and Production Data b) Percentage Change in Standard Deviation of Monthly Sales between Training and Production Data**

To calculate the data drift metrics, we first need to determine the percentage change in both the

mean monthly sales and the standard deviation of monthly sales between the training and production data.

a) Percentage Change in Mean Monthly Sales:

$$\text{Percentage Change} = \frac{(\text{Production Mean} - \text{Training Mean})}{\text{Training Mean}} \times 100$$
$$\text{Percentage Change} = \frac{(45,000 - 50,000)}{50,000} \times 100 = -10\%$$

Given:

- Training Mean Monthly Sales = \$50,000
- Production Mean Monthly Sales = \$45,000

Calculating the percentage change:

$$\text{Percentage Change} = \frac{(45,000 - 50,000)}{50,000} \times 100 = -10\%$$
$$\text{Percentage Change} = \frac{(50,000 - 45,000)}{50,000} \times 100 = 10\%$$

So, there's a 10% decrease in mean monthly sales from training to production data.

b) Percentage Change in Standard Deviation of Monthly Sales:

$$\text{Percentage Change} = \frac{(\text{Production Std. Dev.} - \text{Training Std. Dev.})}{\text{Training Std. Dev.}} \times 100$$
$$\text{Percentage Change} = \frac{(8,500 - 7,000)}{7,000} \times 100 \approx 21.43\%$$

Given:

- Training Standard Deviation of Monthly Sales = \$7,000
- Production Standard Deviation of Monthly Sales = \$8,500

Calculating the percentage change:

$$\text{Percentage Change} = \frac{(8,500 - 7,000)}{7,000} \times 100 \approx 21.43\%$$
$$\text{Percentage Change} = \frac{(7,000 - 8,500)}{7,000} \times 100 \approx -21.43\%$$

So, there's approximately a 21.43% increase in the standard deviation of monthly sales from training to production data.

This analysis suggests that both the mean monthly sales and the variability in monthly sales have changed significantly from the training to the production data, indicating the presence of data drift.

**You are tasked with improving the data quality and observability of a dataset containing customer orders for an e-commerce website. The dataset contains 1,000 records, and you've identified several data quality issues:**

- **Duplicate Orders:** There are 50 duplicate orders in the dataset.
- **Missing Values:** In the "Shipping Address" field, 100 records have missing or incomplete addresses.
- **Inconsistent Date Formats:** The "Order Date" field contains dates in both "MM/DD/YYYY" and "YYYY-MM-DD" formats. There are 30 records with inconsistent date formats.
- **Outliers in Order Amount:** There are 10 orders with an unusually high order amount that seems like outliers.

**You need to: a) Explain why data quality and observability are essential in this context. b) Design and compute four data quality metrics based on the given data. c) Propose a step-by-step plan to address these data quality issues and enhance data observability**

a) Data quality and observability are crucial in this context for several reasons:

- **Decision Making:** High-quality data ensures that decisions made based on the data are accurate and reliable. Inaccurate or incomplete data can lead to flawed insights and decisions.
- **Customer Satisfaction:** Improving data quality can lead to better understanding and serving customer needs, which ultimately enhances customer satisfaction and retention.
- **Operational Efficiency:** Clean and observable data streamlines processes such as order processing, inventory management, and customer support, leading to improved operational efficiency.
- **Compliance:** In many industries, there are regulatory requirements regarding data quality and privacy. Ensuring data quality and observability helps in compliance with these regulations, avoiding legal issues and penalties.
- **Business Intelligence:** High-quality, observable data is essential for generating meaningful business intelligence and forecasting future trends, which is crucial for strategic planning and growth.

b) Four data quality metrics that can be computed based on the given data include:

- **Duplicate Order Rate:** Percentage of duplicate orders in the dataset.
- **Missing Value Rate:** Percentage of records with missing or incomplete addresses in the "Shipping Address" field.
- **Consistency in Date Formats:** Percentage of records with inconsistent date formats in the "Order Date" field.
- **Outlier Detection:** Percentage of orders identified as outliers based on order amount.

c) Proposed step-by-step plan to address data quality issues and enhance data observability:

1. **Data Cleaning:** Remove duplicate orders from the dataset.
2. **Address Completeness:** Investigate and fill in missing or incomplete addresses in the "Shipping Address" field, possibly through data enrichment or customer outreach.
3. **Standardize Date Format:** Convert all dates in the "Order Date" field to a consistent format (e.g., YYYY-MM-DD).
4. **Outlier Detection and Handling:** Identify outliers in the order amount and determine if they are genuine or errors. If errors, correct or remove them; if genuine, investigate further to understand the reasons behind unusually high order amounts.
5. **Documentation and Monitoring:** Document the data cleaning process and regularly monitor data quality metrics to ensure ongoing observability. Implement data quality checks as part of routine data processing workflows to maintain high standards over time.

**Use the below product sales data for answering the OLAP sub-questions [2 \* 5 = 10]**

City | Store | Product | Month | Quantity

Pune | NORTH | P1 | Jan | 17 Pune | NORTH | P1 | Jan | 19 Pune | NORTH | P1 | Feb | 23 Pune | NORTH | P1 | March | 23  
Pune | NORTH | P2 | Jan | 57 Pune | NORTH | P2 | Feb | 78 Pune | NORTH | P2 | March | 45 Pune | SOUTH | P1 | Jan | 12  
Pune | SOUTH | P1 | Feb | 67 Pune | SOUTH | P1 | March | 34

City | Store | Product | Month | Quantity

Mumbai | EAST | P3 | Jan | 34 Mumbai | EAST | P3 | Feb | 54 Mumbai | EAST | P3 | March | 76 Mumbai | EAST | P4 | Jan | 12  
Mumbai | EAST | P4 | Feb | 87 Mumbai | EAST | P4 | March | 90 Mumbai | WEST | P3 | Jan | 57 Mumbai | WEST | P3 | Feb | 32

City | Store | Product | Month | Quantity

Mumbai | EAST | P3 | Jan | 57 Mumbai | EAST | P3 | Feb | 32

a) Represent the sales data in a multi-dimensional format. b) Represent the sales data in multi-dimension format when the months are rolled up to Quarters. c) How the multi-dimensional data structure will look when it's sliced for North Store? d) How the multi-dimensional data structure will look when it's diced for South Store and Product P2? e) If "Product" is replaced by "Category" (product P1, P2 belongs to one category, whereas P3, P4 belongs to another category), how the multi-dimensional structure will look like in this case?

a) The sales data represented in a multi-dimensional format would have dimensions for City, Store, Product, Month, and Quantity. Each entry in the dataset would correspond to a specific combination of these dimensions, along with the corresponding quantity sold.

b) When the sales data is represented in a multi-dimensional format with the months rolled up to quarters, the Month dimension would be replaced by the Quarter dimension. Each quarter would encompass three consecutive months. The other dimensions (City, Store, Product, Quantity) would remain the same.

c) Slicing the multi-dimensional data for the North Store would involve filtering the dataset to include only the entries where the Store dimension is "NORTH". This would provide a subset of the data representing sales from the North Store.

d) Dicing the multi-dimensional data for the South Store and Product P2 would involve filtering the dataset to include only the entries where the Store dimension is "SOUTH" and the Product dimension is "P2". This would provide a subset of the data representing sales of Product P2 from the South Store.

e) If "Product" is replaced by "Category" where products P1 and P2 belong to one category, and products P3 and P4 belong to another category, the multi-dimensional structure would include a Category dimension instead of the Product dimension. The data would be grouped based on this new dimension, with entries corresponding to each category's sales.

**A healthcare organization is sharing medical research data with a research partner while ensuring privacy protection using a privacy-preserving technique called "k-anonymity." The dataset contains information about patients' medical conditions and their ages. The organization wants to disclose aggregate statistics about patient ages while protecting individual privacy. The organization chooses to achieve 3-anonymity, meaning that there are at least three patients with the same age and medical condition for any combination of age and medical condition. The dataset contains the following information:**

- Patient A: Age 45, Medical Condition X
- Patient B: Age 30, Medical Condition Y

- **Patient C: Age 45, Medical Condition Y**
- **Patient D: Age 30, Medical Condition Z**
- **Patient E: Age 45, Medical Condition Z**
- **Patient F: Age 30, Medical Condition Y**
- **Patient G: Age 30, Medical Condition Z**
- **Patient H: Age 30, Medical Condition X**

The healthcare organization is utilizing a privacy-preserving technique called "k-anonymity" to share medical research data with a research partner while safeguarding individual privacy. With the goal of achieving 3-anonymity, they aim to ensure that there are at least three patients sharing the same age and medical condition for any given combination of age and condition.

To address this, they have anonymized the dataset, revealing only aggregate statistics about patient ages. This approach allows the organization to maintain privacy while still sharing valuable research insights.

In summary, the organization's strategy involves implementing k-anonymity to protect individual privacy while sharing medical data with a research partner.

**Consider a dataset containing information about customer orders for an e-commerce website. There are 20 records with various data quality issues. Identify at least 5 data quality issues and compute data quality metrics to help illustrate the importance of data quality assessment.**

Here are five potential data quality issues and corresponding data quality metrics:

1. **Missing Values:** Identify any columns with missing values and calculate the percentage of missing values for each column. This can be measured using the Missing Value Ratio, which is the ratio of missing values to total values in a column.
2. **Inconsistent Data Types:** Check if data types are consistent across columns. For example, ensure that dates are stored as date data types and not as strings. Calculate the percentage of inconsistent data types as a data quality metric.
3. **Duplicates:** Look for duplicate records in the dataset. Calculate the percentage of duplicate records as a data quality metric using the Deduplication Ratio, which is the ratio of duplicate records to total records.
4. **Outliers:** Identify outliers in numerical columns. Calculate the percentage of outliers for each numerical column using statistical methods like the Interquartile Range (IQR) or Z-score.
5. **Incorrect Values:** Look for values that fall outside the expected range or domain for each column. Calculate the percentage of incorrect values based on domain knowledge or predefined ranges.

Once these data quality issues are identified and quantified, data quality metrics such as Missing Value Ratio, Deduplication Ratio, and percentage of outliers can be computed to illustrate the importance of data quality assessment.

**Mainwind Traders has several data workloads. Design and explain briefly your data solution to handle this data's ingestion, processing, and analysis based on the following requirements. Tabulate your answer as shown below and draw suitable block diagrams wherever appropriate. Also, comment on whether Lambda architecture would be helpful in this solution.**

- **Manufacturing data:** The company has been storing manufacturing logs collected from the assembly line. You want to analyze these logs to gain insights into material behavior and quality assurance. To analyze these logs, you need to use reference data, such as material, chemical, and origin information, in addition to the logs data. You want to utilize this data from the on-premises data store, combine it with other log data that is in a cloud data store, and run stored procedures on the data to gain insights.
- **Real-time data:** Mainwind Traders needs real-time data ingestion and storage for multiple data sources like their websites, point-of-sale systems, and social media sites. It would be best if you had a solution to analyze this data and provide valuable insights to the CEO.
- **Historical company data:** The company is required to store 5 TB of company data for legal reasons. This data is rarely used or referenced but must not be deleted. You need a cost-effective method for storing this data.

Manufacturing data	Real-time data	Historical company data
Sources		
Ingestion strategy		
Storage strategy		
Tools for Ingestion		
Processing strategy		
Processing complexity		
Processing tools		
Storage/usage of analysis results		
Other Architectural components		

Here's a solution design for handling Mainwind Traders' data workloads:

	Manufacturing data	Real-time data	Historical company data
Sources	On-premises data store, cloud data store	Websites, point-of-sale systems, social media sites	
Ingestion strategy	Combine manufacturing logs with reference data	Real-time ingestion from various sources	Store 5 TB of data for legal reasons
Storage strategy	Utilize on-premises and cloud data stores	Real-time storage for immediate analysis and insights	Cost-effective storage solution for infrequently used data
Tools for Ingestion	Tools for data integration and transformation	Stream processing frameworks, connectors for different data sources	
Processing strategy	Run stored procedures for in-depth analysis	Real-time analytics for immediate insights	
Processing complexity	High complexity due to combining diverse datasets	High velocity and volume of incoming data	
Processing tools	Database management systems, data warehouses	Real-time analytics platforms, machine learning algorithms	
Storage/usage of analysis results	Store analyzed data for future reference	Immediate usage for decision-making	Store data securely to meet legal requirements
Other Architectural components	Data governance framework, data quality checks, security measures	Scalable infrastructure for handling real-time data streams, fault tolerance mechanisms	

In this solution, Lambda architecture could be beneficial, especially for real-time data processing. It allows for parallel processing of batch and stream data, enabling both historical analysis and real-time insights. Batch processing handles the historical company data, while stream processing handles the real-time data ingestion and analysis. This architecture ensures robustness, fault tolerance, and scalability, making it suitable for Mainwind Traders' diverse data needs.

Consider the following schema definition that is representing the loan status of various applicants and is used to predict whether the customer will be a defaulter or not, what can be loan amount and tenure of the loan.

Attribute	Description
emp_title	Job title
emp_length	Number of years in the job, rounded down. If longer than 10 years, then this is represented by 10
state	Two-letter state code
homeownership	The ownership status of the applicant's residence
annual_income	Annual income
verified_income	Type of verification of the applicant's income
debt_to_income	Debt-to-income ratio
annual_income_joint	If this is a joint application, then the annual income of the two parties applying
accounts_opened_24m	Number of new lines of credit opened in the last 24 months
months_since_last_credit_inquiry	Number of months since the last credit inquiry on this applicant
num_satisfactory_accounts	Number of satisfactory accounts
num_accounts_120d_past_due	Number of current accounts that are 120 days past due
num_accounts_30d_past_due	Number of current accounts that are 30 days past due
num_active_debit_accounts	Number of currently active bank cards
total_debit_limit	Total of all bank card limits
num_total_cc_accounts	Total number of credit card accounts in the applicant's history
num_open_cc_accounts	Total number of currently open credit card accounts
num_cc_carrying_balance	Number of credit cards that are carrying a balance
num_mort_accounts	Number of mortgage accounts
account_never_delinq_percent	Percent of all lines of credit where the applicant was never delinquent
public_record_bankrupt	Number of bankruptcies listed in the public record for this applicant
loan_purpose	The category for the purpose of the loan application
application_type	The type of application: either individual or joint
loan_amount	The amount of the loan the applicant received
term	The number of months of the loan the applicant received
interest_rate	Interest rate of the loan the applicant received
installment	Monthly payment for the loan the applicant received
grade	Grade associated with the loan
sub_grade	Detailed grade associated with the loan
issue_month	Month the loan was originated

- a) Identify instances of data leakages. [2]
- b) When developing a model to predict whether the customer will be a defaulter, can k-means algorithms be used as a data partitioning strategy? Why? [1.5]
- c) If the model developed in part c is used online to make defaulter predictions, can the feature engineering code used while training the model be used as is for serving the model? [1.5]
- d) Explain the approach that can be used for labeling the dataset. [2]
- e) Whether data versioning will be required? Justify. [1]
- f) Will the tools for model experimentation and model registry be needed for this case? [2]

a) Instances of data leakages could occur if sensitive information, such as personal financial data or credit history, is improperly handled or exposed, leading to breaches of privacy or security.

b) K-means algorithms are typically used for clustering data into groups based on similarity. While they may not be directly applicable for predicting loan default, they could potentially be used as part of a preprocessing step to identify clusters of applicants with similar characteristics for further analysis.

c) When using the model online for making defaulter predictions, it's essential to ensure that the feature engineering code used during training is compatible with the serving environment. Any discrepancies or issues in feature engineering between training and serving could lead to inaccurate predictions or model failures.

d) Labeling the dataset for loan default prediction involves assigning a binary label (default or non-default) to each applicant based on their loan repayment status. This labeling process typically requires historical data on loan outcomes and may involve manual review or automated algorithms to determine default status.

e) Data versioning may be required to track changes and updates to the dataset, feature engineering code, and model versions over time. This ensures reproducibility, transparency, and accountability in the model development and deployment process.

f) Tools for model experimentation and model registry can be valuable for managing the development, testing, and deployment of machine learning models for loan default prediction. These tools help track experiments, version control models, and monitor model performance in production environments, ensuring reliability and scalability.



Name- Bipatti Ujjwal Rakeshmarjan  
Id- 2022-da04182

DMUML  
DSECLZG529

- Q-2) In order to achieve 3-anonymity, we need to group patients with same features like (age, medical conditions) until there are three patients in each group

Group 1

Patient A: Age 45, Medical condition X  
Patient C: Age 45, Medical condition Z  
Patient F: Age 45, Medical condition X

Group 2

Patient B: Age 30, Medical condition Y  
Patient G: Age 30, Medical condition Z  
Patient H: Age 30, Medical condition X

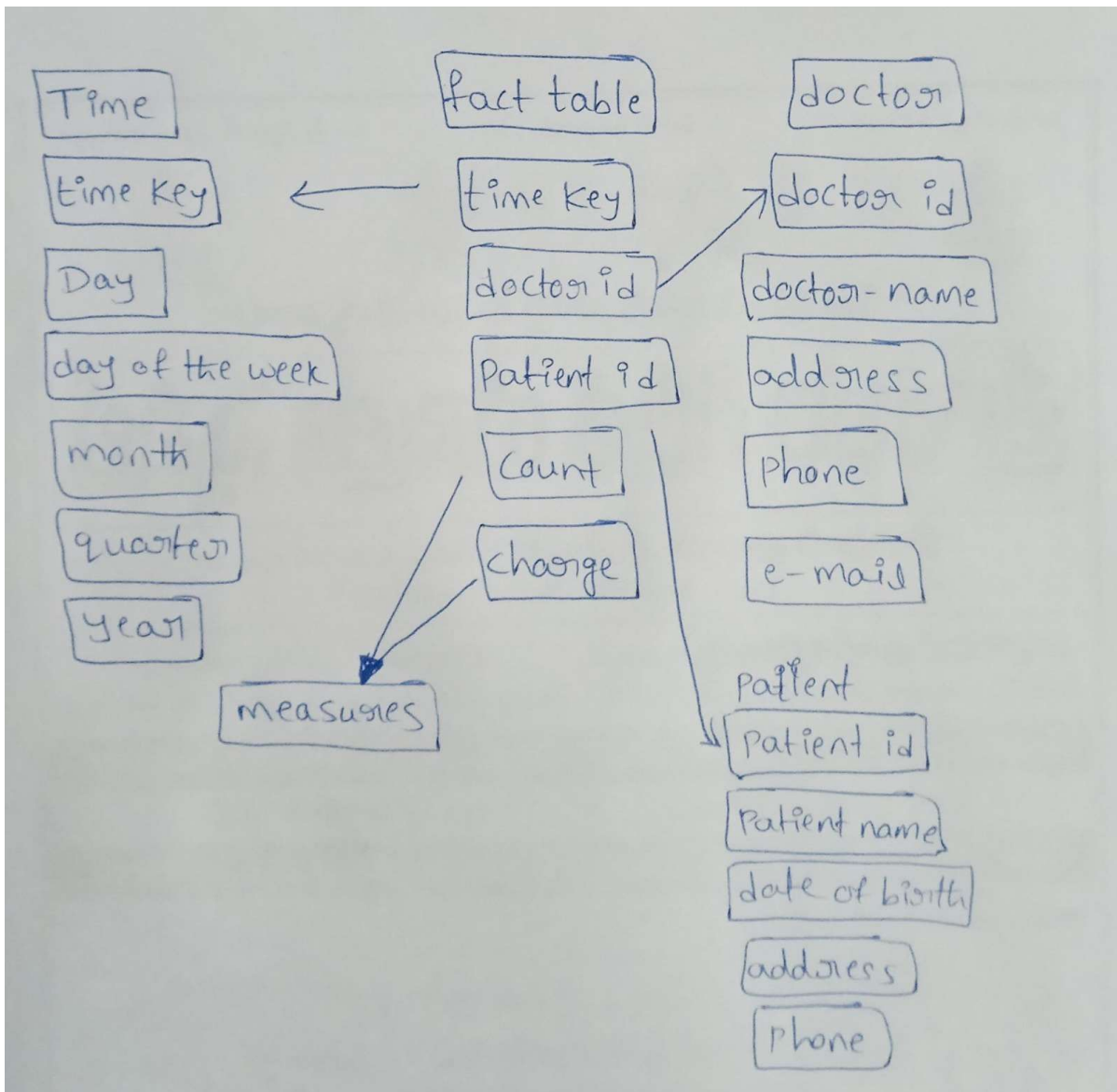
Group 3

Patient D: Age 35, Medical condition X

Group 4

Patient E: Age 50, Medical condition Y

- 6) The transformed dataset ensures the protection of individual patient by privacy by adhering to the principle of 3-anonymity. This means that within each group there are atleast 3 patients who share same age and same medical condition of specific patients while still allowing for meaningful statistical analysis all while preserving privacy.



# 1.

When designing a data versioning system for the retail corporation's machine learning projects, several key components and considerations should be taken into account:

1. **Data Storage:** Implement a centralized data storage system where all datasets are stored securely. This could be a cloud-based storage solution or an on-premises database.
2. **Version Control System\*\*:** Utilize a version control system (VCS) such as Git to manage changes to datasets. Each dataset should have its repository, allowing for easy tracking of modifications and rollbacks if necessary.
3. **\*\*Metadata Management\*\*:** Develop a metadata management system to store information about each dataset, including its source, preprocessing steps, and associated models. This metadata should be easily accessible and searchable.
4. **\*\*Data Lineage Tracking\*\*:** Implement mechanisms to track the lineage of data, documenting how each dataset was created and modified over time. This helps ensure data quality and traceability.
5. **\*\*Access Control\*\*:** Set up access control mechanisms to restrict access to sensitive data and ensure that only authorized users can modify or access certain datasets.

6. **Automated Pipelines**: Create automated data pipelines to streamline the process of updating datasets and integrating them into machine learning workflows. These pipelines should include validation steps to ensure data integrity.
7. **Versioning Policies**: Establish clear versioning policies outlining when and how datasets should be versioned. This includes defining criteria for creating new versions, such as significant data changes or model performance improvements.
8. **Integration with ML Workflow**: Integrate the data versioning system with the overall machine learning workflow, allowing data scientists to easily access and incorporate versioned datasets into their models.
9. **Scalability and Performance**: Ensure that the data versioning system is scalable to handle large volumes of data and performs efficiently, even as the retail corporation's operations grow.
10. **Documentation and Training**: Provide comprehensive documentation and training materials to educate data scientists and other stakeholders on how to use the data versioning system effectively.

By considering these components and considerations, the data science team can develop a robust data versioning system that supports the retail corporation's machine learning projects effectively.

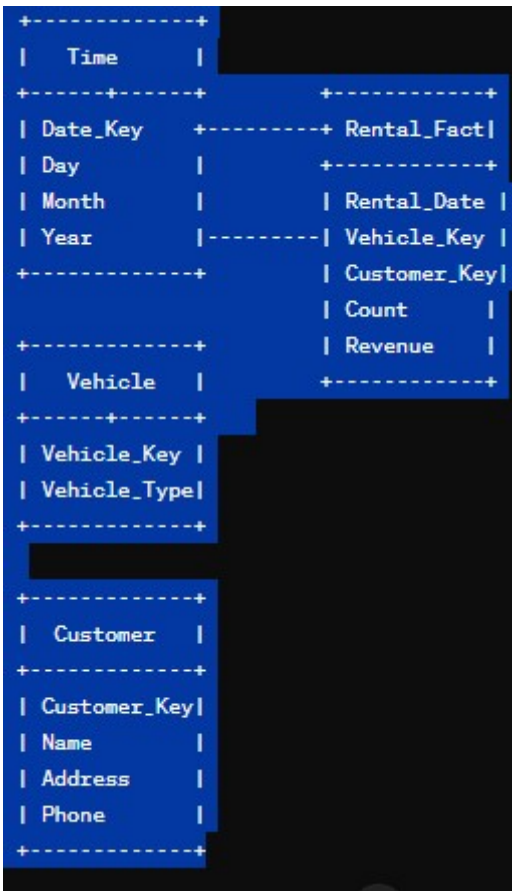
## 2.

- a) For the data labeling process in medical imaging datasets, I would employ a combination of expert annotation by trained medical professionals and semi-automated annotation using advanced image processing algorithms to ensure accuracy and efficiency.
- b) To efficiently label large volumes of medical images while maintaining high-quality annotations, I would utilize specialized annotation software platforms that support collaborative labeling workflows, integrate AI-assisted annotation tools for semi-automation, and implement rigorous quality control measures.
- c) Data augmentation plays a crucial role in enhancing the generalization and robustness of machine learning models trained on medical imaging data by diversifying the training dataset, reducing overfitting, and improving model performance on unseen data.
- d) Common data augmentation techniques suitable for medical imaging datasets include rotation, flipping, scaling, translation, elastic deformation, contrast adjustment, noise injection, and histogram equalization, tailored to accommodate various image modalities and anatomical variabilities.
- e) Data validation is essential in healthcare machine learning applications to ensure the quality and reliability of training and evaluation datasets, safeguard patient safety, and maintain regulatory compliance.
- f) The process of data validation involves comprehensive quality assurance procedures such as outlier detection, anomaly identification, consistency checks, cross-validation, and expert review, coupled with robust data preprocessing techniques to handle anomalies and ensure dataset integrity for reliable model training and evaluation.

3.

(a) Two classes of schemas popularly used for modeling data warehouses in the context of the car rental company are Star schema and Snowflake schema.

(b) Schema Diagram:



(c) OLAP Operations:

1. Drill-down by vehicle type to the level of rental transactions.
2. Apply the average function on the rental duration measure.

(d) SQL Query:

```
SELECT vehicle_type, AVG(rental_duration) AS avg_rental_duration
FROM rental_transactions
GROUP BY vehicle_type;
```

4.

a)

**Data Ingestion Layer:**

- Stream Processing: Utilize stream processing technologies such as Apache Kafka or Amazon Kinesis to handle real-time ingestion of user-generated content, interactions, and registrations.
- Data Pipeline Orchestration: Implement workflow orchestration tools like Apache Airflow or Apache NiFi to manage and automate the data ingestion processes.
- Schema Validation: Employ schema validation techniques to ensure data consistency and quality before ingestion.
- Scalability: Design the ingestion layer to scale horizontally to accommodate the increasing volume of incoming data.
- Fault Tolerance: Implement mechanisms for fault tolerance and data replication to prevent data loss.

**Data Storage Layer:**

- NoSQL Databases: Utilize NoSQL databases like Apache Cassandra or MongoDB for storing user-generated posts and interactions due to their ability to handle high volumes of unstructured data and provide horizontal scalability.
- Relational Databases: Use relational databases like PostgreSQL or MySQL for storing user registration data, which typically requires ACID compliance and structured querying.
- Object Storage: Employ object storage solutions like Amazon S3 or Google Cloud Storage for storing large binary files such as images or videos associated with user-generated content.
- Data Partitioning: Implement data partitioning strategies to distribute data across multiple nodes and improve query performance.
- Data Replication: Ensure data replication across multiple data centres for high availability and disaster recovery.

**Data Processing Layer:**

- Batch Processing: Utilize batch processing frameworks like Apache Spark or Apache Flink for performing complex analytics and processing large volumes of historical data.
- Real-time Processing: Implement real-time processing using technologies like Apache Storm or Apache Samza to analyze user interactions and generate real-time insights.
- Distributed Computing: Design processing algorithms to leverage distributed computing paradigms for efficient utilization of resources.
- Data Serialization: Use efficient data serialization formats like Apache Avro or Protocol Buffers to optimize data transfer and processing speed.
- Machine Learning Integration: Integrate machine learning models for tasks such as content recommendation or sentiment analysis.

**Data Analytics and Reporting Layer:**

- Data Warehousing: Utilize data warehousing solutions like Amazon Redshift or Google BigQuery for storing aggregated and processed data for analytics and reporting purposes.
- Visualization Tools: Integrate visualization tools like Tableau or Apache Superset to create interactive dashboards and reports for business users.
- Ad Hoc Querying: Provide support for ad hoc querying using SQL or OLAP tools to enable data exploration and analysis.
- Scheduled Reporting: Implement scheduled reporting capabilities to automatically generate and distribute reports to stakeholders.

**b)**

**Estimated Storage Capacity:**

- User-generated posts per day:  $50,000,000 * (\text{average\_post\_size}) * (\text{retention\_period})$
- Likes, comments, and shares per day:  $20,000,000 * (\text{average\_interaction\_size}) * (\text{retention\_period})$
- New user registrations per day:  $10,000,000 * (\text{average\_registration\_data\_size}) * (\text{retention\_period})$

**Factors to consider:**

- Redundancy requirements: Factor in redundancy requirements (e.g., replication factor) for fault tolerance and data durability.
- Retention period: Determine the retention period based on regulatory requirements, business needs, and historical analysis.
- Compression techniques: Utilize compression techniques to reduce storage requirements where feasible.

## 5.

a) To calculate the current average processing latency for user posts:

Current average processing latency = Previous average processing latency + (Previous average processing latency \* Increase percentage)

Current average processing latency = 50 milliseconds + (50 milliseconds \* 0.20)

Current average processing latency = 50 milliseconds + 10 milliseconds

Current average processing latency = 60 milliseconds

b) The increase in data processing latency for user posts can have several potential impacts on the overall performance of the data pipeline:

- **Decreased Real-time Analytics:** Longer processing latency can lead to delayed insights and real-time analytics, affecting the platform's ability to provide timely recommendations, notifications, or updates to users.
- **Bottleneck in Data Flow:** Increased latency may indicate a bottleneck in the data processing pipeline, potentially causing data backups or overloading downstream systems.
- **Negative User Experience:** Users may experience delays in posting content or receiving responses, leading to dissatisfaction and reduced engagement with the platform.
- **Impact on Scalability:** If the latency increase is due to resource constraints or inefficient processing algorithms, it may hinder the scalability of the platform, making it challenging to handle growing volumes of data.

To mitigate the impact of increased latency, the data engineering team can consider the following strategies:

- **Performance Optimization:** Identify and optimize performance bottlenecks in the data processing pipeline, such as inefficient algorithms, resource contention, or network latency.
- **Horizontal Scaling:** Scale out the infrastructure horizontally by adding more processing nodes or resources to distribute the workload and reduce processing time.
- **Streamlining Data Flow:** Review the data flow architecture and streamline processes to minimize unnecessary data transformations or redundant operations.
- **Caching:** Implement caching mechanisms to store and retrieve frequently accessed data, reducing the need for repetitive processing and improving overall latency.
- **Parallel Processing:** Parallelize data processing tasks to leverage multi-core processors or distributed computing frameworks, enabling faster execution of tasks in parallel.
- **Monitoring and Alerting:** Set up monitoring and alerting systems to track processing latency metrics in real-time and proactively identify and address issues before they impact performance.