SQL ML Interview Questions - Theory Questions

Question 1

What are the different types of JOIN operations in SQL?

Answer:

Theory

JOIN operations in SQL are used to combine rows from two or more tables based on a related column between them. The type of JOIN determines which rows from the tables will be included in the result set. Understanding joins is fundamental for creating feature sets for machine learning, as data is often spread across multiple normalized tables.

The primary types of JOINs are:

- 1. **INNER JOIN:** Returns only the rows where the join condition is met in both tables. It's the most common join type.
- 2. **LEFT JOIN (or LEFT OUTER JOIN):** Returns all rows from the left table and the matched rows from the right table. If there is no match, the columns from the right table will have NULL values.
- 3. **RIGHT JOIN (or RIGHT OUTER JOIN):** Returns all rows from the right table and the matched rows from the left table. If there is no match, the columns from the left table will have NULL values.
- 4. FULL OUTER JOIN: Returns all rows when there is a match in either the left or the right table. It combines the functionality of both LEFT and RIGHT JOINs. If there is no match, the missing side's columns will be NULL.
- 5. **CROSS JOIN:** Returns the Cartesian product of the two tables, i.e., it combines every row from the first table with every row from the second table. It's rarely used with a WHERE clause but can be useful for generating all possible combinations.
- 6. **SELF JOIN:** This is not a distinct join type but a regular join where a table is joined with itself. It's useful for querying hierarchical data or comparing rows within the same table.

Use Cases

- **INNER JOIN:** To combine customer data with their corresponding orders.
- **LEFT JOIN:** To find all customers, including those who have never placed an order. This is useful for identifying inactive users for a churn model.

- **RIGHT JOIN:** Less common, but could be used to find all products that have been ordered, including details of the customer who ordered them.
- FULL OUTER JOIN: To get a complete list of all customers and all orders, matching them where possible.
- **CROSS JOIN:** To generate a baseline feature set of all possible user-product pairs for a recommendation system.

Best Practices

- Always use explicit JOIN syntax (INNER JOIN, LEFT JOIN) instead of the old comma-based syntax in the FROM clause.
- Use meaningful table aliases to improve readability, especially in complex queries with multiple joins.
- Ensure the columns used in the ON clause are indexed to significantly improve query performance.

Pitfalls

- Accidental CROSS JOIN: Forgetting the ON clause can lead to a Cartesian product, which can be computationally massive and produce incorrect results.
- **Performance of OUTER JOINs:** LEFT, RIGHT, and FULL joins can be slower than INNER JOINs, especially on large tables, as they may need to scan more data.

Question 2

Explain the difference between WHERE and HAVING clauses.

Answer:

Theory

The fundamental difference between WHERE and HAVING lies in when they are applied during query execution. WHERE filters individual rows *before* any groupings are made, while HAVING filters groups of rows *after* they have been aggregated by the GROUP BY clause.

Order of Execution in a SQL Query:

- 1. FROM / JOIN
- 2. WHERE
- 3. GROUP BY
- 4. HAVING
- 5. SELECT
- 6. ORDER BY

7. LIMIT / OFFSET

Because of this order:

- The WHERE clause operates on row-level data. You cannot use aggregate functions (like SUM(), COUNT(), AVG()) in a WHERE clause.
- The HAVING clause operates on the results of aggregate functions. It is applied to the grouped data.

Code Example

```
Generated sql
-- Select customer_ids who have made more than 5 purchases
-- of over $10 each in the last year.

SELECT
customer_id,
COUNT(order_id) AS number_of_orders,
SUM(order_value) AS total_spent

FROM
orders

WHERE
order_date >= '2023-01-01' -- Filters individual rows BEFORE grouping
AND order_value > 10.00

GROUP BY
customer_id

HAVING
```

Explanation

1. **FROM orders**: The query starts with the orders table.

COUNT(order id) > 5; -- Filters groups AFTER grouping

- 2. WHERE order_date >= '2023-01-01' AND order_value > 10.00: SQL first filters the orders table, removing any rows representing orders that are old or have a low value. This reduces the amount of data that needs to be processed in the next step.
- 3. **GROUP BY customer_id**: The remaining rows are grouped by customer_id, so all orders from a single customer are put into one group.
- 4. **HAVING COUNT(order_id) > 5**: Now, the HAVING clause filters these *groups*. It keeps only the groups (customers) where the count of orders is greater than 5.
- 5. **SELECT ...**: Finally, the SELECT list is calculated for the remaining groups.

Best Practices & Optimization

• **Filter Early with WHERE**: Always use the WHERE clause to filter as many rows as possible before the GROUP BY phase. This is far more efficient because it reduces the

- number of rows that need to be aggregated, leading to lower memory usage and faster computation.
- Use HAVING only for filtering on aggregate functions. If a filter condition does not involve an aggregate, it belongs in the WHERE clause.

Question 3

Describe a subquery and its typical use case.

Answer:

Theory

A subquery, also known as an inner query or nested query, is a SELECT statement that is nested inside another SQL statement (the outer query). The result of the subquery is used by the outer query to complete its operation. Subqueries can be used in SELECT, FROM, WHERE, and HAVING clauses and can return a scalar value (one row, one column), a list of values (one column, multiple rows), or a derived table (multiple columns, multiple rows).

There are two main types:

- 1. **Non-Correlated Subquery:** The inner query is independent of the outer query. It executes once, and its result is used by the outer query.
- 2. **Correlated Subquery:** The inner query depends on the outer query for its values. It executes once for each row processed by the outer query, which can be inefficient.

Code Example

Generated sql

- -- Non-correlated subquery in a WHERE clause
- -- Find all employees who work in the 'Engineering' department.

SELECT employee_name
FROM employees
WHERE department_id IN (SELECT department_id FROM departments WHERE department_name = 'Engineering');

- -- Correlated subquery in a SELECT clause
- -- Find each employee's salary and the average salary of their department.

```
SELECT employee_name, salary,
```

(SELECT AVG(salary) FROM employees e2 WHERE e2.department_id = e1.department_id)
AS department_avg_salary
FROM
employees e1;

IGNORE_WHEN_COPYING_START content_copy download
Use code with caution. SQL
IGNORE_WHEN_COPYING_END

Explanation

- In the first example, the subquery SELECT department_id FROM departments WHERE department_name = 'Engineering' runs first and returns the ID for the Engineering department. The outer query then uses this ID to find all employees in that department.
- In the second example, the subquery is correlated. For each row e1 in the outer query, the inner query calculates the average salary for that specific employee's department (e2.department_id = e1.department_id).

Use Cases

- Filtering Data: Using IN, NOT IN, ANY, ALL in a WHERE clause.
- **Creating Derived Tables:** A subquery in the FROM clause can create a temporary, virtual table to be used in the main query.
- Calculating Derived Columns: A scalar subquery in the SELECT list can compute a
 value related to the current row.

Pitfalls and Optimization

- **Performance:** Correlated subqueries can be very slow because they execute repeatedly. Often, they can be rewritten using JOINs or window functions for better performance.
- **Readability:** Deeply nested subqueries can be difficult to read and maintain. Common Table Expressions (CTEs) are often a better alternative for structuring complex logic.

Question 4

Can you explain the use of indexes in databases and how they relate to Machine Learning?

Answer:

Theory

An index in a database is a data structure (most commonly a B-Tree) that improves the speed of data retrieval operations on a database table at the cost of additional writes and storage space. When you query a column with an index, the database can use the index to find the location of the requested rows directly, similar to using the index of a book to find a specific topic, rather than reading the entire book page by page (a full table scan).

How Indexes Relate to Machine Learning

Indexes are crucial for the efficiency of ML workflows that rely on data from SQL databases. Their primary impact is on the speed of data preparation and model deployment.

1. Faster Feature Engineering:

- ML feature engineering often involves complex queries with JOINs and WHERE clauses to combine data from multiple tables.
- Creating indexes on foreign keys (used in JOINs) and columns used for filtering (in WHERE clauses) can speed up the creation of your training dataset by orders of magnitude. For example, indexing customer_id and product_id would drastically accelerate joining sales, customers, and products tables.

2.

3. Efficient Data Sampling:

 When you need to create a representative sample from a massive dataset, indexes can help. For example, in stratified sampling, you might filter by a category (WHERE user_segment = 'premium'). If user_segment is indexed, finding all users in that segment is fast.

4.

5. Low-Latency Real-Time Predictions:

- o In a production environment, a deployed model often needs to fetch features for a single entity (e.g., a user visiting a website) to make a real-time prediction.
- A query like SELECT * FROM user_features WHERE user_id = ? must be extremely fast. Having a primary key index on user_id ensures this lookup is nearly instantaneous. Without an index, the database would have to scan the entire table, making real-time prediction impossible.

6.

Best Practices

- Index Strategically: Index columns that are frequently used in WHERE clauses, JOIN conditions, and ORDER BY clauses.
- Avoid Over-Indexing: Every index consumes storage and adds overhead to write operations (INSERT, UPDATE, DELETE), as the index must also be updated.
- **Use Composite Indexes:** For queries that filter on multiple columns, a composite index on those columns can be very effective.

Pitfalls

- Indexing Low-Cardinality Columns: Indexing a column with very few unique values (e.g., a boolean is_active column) is often not helpful, as a full table scan may be faster.
- **Ignoring Index Maintenance:** Indexes can become fragmented over time, requiring periodic maintenance (REINDEX, OPTIMIZE TABLE).

Question 5

Explain the importance of data normalization in SQL and how it affects Machine Learning models.

Answer:

Theory

It's crucial to distinguish between two different concepts of "normalization":

- 1. **Database Normalization (SQL):** This is a database design technique for organizing tables and columns to minimize data redundancy and improve data integrity. It involves following a set of rules called Normal Forms (1NF, 2NF, 3NF, etc.). The goal is to ensure that data is stored logically and efficiently, avoiding update, insertion, and deletion anomalies. A highly normalized database will have many small tables with relationships defined by foreign keys.
- 2. Feature Normalization/Standardization (Machine Learning): This is a data preprocessing technique used to change the scale of feature values. Common methods include Min-Max Scaling (scaling values to a [0, 1] range) and Standardization (rescaling to have a mean of 0 and a standard deviation of 1). This is essential for algorithms that are sensitive to the scale of input features, such as Gradient Descent-based models (Linear Regression, Neural Networks), SVMs, and distance-based algorithms like KNN.

How SQL Normalization Affects ML

SQL database normalization has a significant, indirect impact on the machine learning process.

Advantages:

- High Data Quality and Integrity: Normalization enforces consistency. For example, a
 customer's name is stored only once in a customers table, preventing inconsistencies
 that could arise if it were repeated in every orders record. High-quality, consistent data is
 the foundation of any good ML model.
- Reduced Storage: Minimized redundancy saves storage space.

Disadvantages (from an ML perspective):

- Complex Feature Engineering: To create a flat feature matrix for an ML model (where each row is an observation and each column is a feature), you must "denormalize" the data. This requires writing complex SQL queries with multiple JOIN operations to bring together related information from different tables.
- Query Performance Overhead: Executing these multi-table JOINs can be computationally expensive and slow, especially on large datasets. This can make the process of creating and iterating on training datasets time-consuming.

Best Practices

- OLTP vs. OLAP: In practice, production systems (OLTP Online Transaction Processing) are often highly normalized for data integrity. For analytics and machine learning (OLAP - Online Analytical Processing), data is often extracted, transformed, and loaded (ETL) into a denormalized structure like a data warehouse or data mart (e.g., a star schema). This denormalized view makes feature generation much faster and simpler.
- This means SQL normalization is critical at the data source, but the data is intentionally denormalized for the ML task.

Question 6

What are SQL Window Functions and how can they be used for Machine Learning feature engineering?

Answer:

Theory

SQL window functions perform calculations across a set of rows that are related to the current row. Unlike standard aggregate functions (SUM, COUNT), which collapse multiple rows into a single output row, window functions return a value for *every* row based on a "window" of related rows defined by the OVER() clause.

The OVER() clause has two key components:

- PARTITION BY: Divides the rows into partitions (groups). The window function is applied independently to each partition. This is similar to GROUP BY.
- ORDER BY: Orders the rows within each partition. This is crucial for functions that are order-sensitive, like LAG, LEAD, and running totals.

Use for Machine Learning Feature Engineering

Window functions are exceptionally powerful for creating sophisticated features, especially for time-series or sequential data.

1. Time-Series Features:

- LAG(column, n): Access data from n rows before the current row. Useful for creating features like "last month's sales" or "previous user session duration."
- LEAD(column, n): Access data from n rows after the current row. Useful for creating target variables (e.g., "will the user purchase in the next 7 days?").

2.

3. Ranking and Positional Features:

- o **ROW_NUMBER():** Assigns a unique number to each row within a partition.
- RANK() and DENSE_RANK(): Rank rows based on a value. Useful for features like "rank of a product's sales within its category."

4.

5. Cumulative and Rolling Aggregates:

 SUM(...) OVER (...) or AVG(...) OVER (...): Calculate running totals or moving averages. For example, "a customer's cumulative spending up to this point" or "30-day moving average of product sales."

6.

Code Example

Generated sql

-- Create features for a customer transaction dataset

SELECT

customer id,

transaction date,

transaction amount,

-- Feature 1: Customer's previous transaction amount

LAG(transaction_amount, 1, 0) OVER (PARTITION BY customer_id ORDER BY transaction_date) AS prev_transaction_amount,

-- Feature 2: Customer's running total of spending

SUM(transaction_amount) OVER (PARTITION BY customer_id ORDER BY transaction_date) AS running_total_spent,

-- Feature 3: Rank of this transaction's size for the customer

RANK() OVER (PARTITION BY customer_id ORDER BY transaction_amount DESC) AS transaction_rank

FROM

transactions;

IGNORE_WHEN_COPYING_START content_copy download Use code with caution. SQL

IGNORE_WHEN_COPYING_END

Explanation

- PARTITION BY customer_id: Ensures that all calculations are done independently for each customer.
- ORDER BY transaction_date: Orders each customer's transactions chronologically, which is essential for LAG and the running SUM.
- LAG(...): Gets the transaction amount from the previous row for that customer.
- SUM(...): Calculates the cumulative sum of transaction_amount up to the current row for that customer.
- RANK(): Ranks the transactions for each customer from largest to smallest.

Question 7

Explain how to discretize a continuous variable in SQL.

Answer:

Theory

Discretization, or binning, is the process of converting a continuous variable into a discrete, categorical variable by grouping values into a finite number of intervals or "bins". This technique can be useful for certain machine learning models (like Decision Trees or Naive Bayes) or to simplify a feature's representation.

There are two common strategies for discretization:

- 1. **Equal-Width Binning:** Divides the range of the variable into N bins of equal width. This is simple but can be sensitive to outliers.
- 2. **Equal-Frequency (Quantile) Binning:** Divides the data into N bins with approximately the same number of observations in each bin. This handles skewed data better.

Multiple Solution Approaches & Code Examples

Approach 1: Equal-Width Binning using CASE Statement

This is the most straightforward method for creating custom, fixed-width bins.

```
Generated sql
-- Discretize user age into predefined groups
SELECT
user_id,
age,
```

```
CASE
WHEN age >= 18 AND age <= 24 THEN '18-24'
WHEN age >= 25 AND age <= 34 THEN '25-34'
WHEN age >= 35 AND age <= 49 THEN '35-49'
ELSE '50+'
END AS age_group
FROM
users;

IGNORE_WHEN_COPYING_START
content_copy download
Use code with caution. SQL
IGNORE_WHEN_COPYING_END
```

Explanation: The CASE statement evaluates each user's age and assigns a corresponding string label for the age_group.

Approach 2: Equal-Frequency Binning using NTILE() Window Function

NTILE(N) is a window function that divides the rows in a partition into N ranked groups (or buckets) as evenly as possible.

```
Generated sql
-- Discretize user spending into 4 quartile groups (top 25%, etc.)

SELECT
user_id,
total_spending,
NTILE(4) OVER (ORDER BY total_spending) AS spending_quartile

FROM
user_summary;

IGNORE_WHEN_COPYING_START
content_copy download
Use code with caution. SQL
IGNORE WHEN COPYING END
```

Explanation:

- ORDER BY total_spending: The data must be ordered by the continuous variable you want to bin.
- NTILE(4): This divides the entire dataset into 4 groups. The users with the lowest spending get a spending_quartile of 1, the next group gets 2, and so on, up to the highest spenders in group 4. Each group will contain roughly 25% of the users.

Best Practices

- Choose the Right Method: Use CASE for domain-specific, meaningful bins (like standard age groups). Use NTILE when you want to create bins with equal numbers of observations, which is often better for model performance as it avoids creating very sparse categories.
- The number of bins is a hyperparameter. The choice can impact model performance, so it may require some experimentation.

Question 8

Explain how to perform binning of categorical variables in SQL for use in a Machine Learning model.

Answer:

Theory

Binning categorical variables typically refers to the process of handling high-cardinality features—categorical variables with a large number of unique values (e.g., zip codes, user IDs, job titles). High cardinality can cause problems for ML models, such as overfitting, increased memory usage, and poor performance.

The most common technique is to group infrequent or rare categories into a single, catch-all category, often named "Other". This reduces the number of distinct categories while preserving the information from the most frequent ones.

Code Example

The process involves two steps:

- 1. Identify the most frequent categories.
- 2. Use a CASE statement to re-label the infrequent ones.

This can be done efficiently using a Common Table Expression (CTE).

Generated sal

-- Bin a 'city' column, keeping the top 10 cities and grouping the rest as 'Other'.

```
WITH CityCounts AS (
-- Step 1: Find the top 10 most frequent cities
SELECT
city
FROM
users
```

```
GROUP BY
    city
  ORDER BY
    COUNT(*) DESC
  LIMIT 10
-- Step 2: Bin the cities based on whether they are in the top 10
SELECT
  user id,
  city AS original city,
  CASE
    WHEN city IN (SELECT city FROM CityCounts) THEN city
    ELSE 'Other'
  END AS city_binned
FROM
  users:
IGNORE WHEN COPYING START
content_copy download
Use code with caution. SQL
IGNORE_WHEN_COPYING_END
```

Explanation

- 1. **WITH CityCounts AS (...)**: This CTE first calculates the frequency of each city, orders them by count in descending order, and selects the top 10. The result of this CTE is a temporary list of the 10 most common cities.
- 2. **SELECT** ... **CASE** ...: The main query then goes through the users table. For each row, the CASE statement checks if the user's city is present in the CityCounts list.
 - If it is, the original city name is kept.
 - If it's not, the city is relabeled as 'Other'.

3.

Use Cases

- **Geographic Data:** Grouping rare zip codes, cities, or countries.
- User-Generated Content: Binning infrequent product tags or search terms.
- **Device Information:** Grouping rare mobile device models or browser versions.

Best Practices

• Threshold Selection: The threshold for what is considered "rare" (e.g., top 10, or any category with <1% frequency) is a crucial choice and can be tuned like a hyperparameter based on its effect on model performance.

 Domain Knowledge: Sometimes, rare categories might be highly predictive (e.g., a rare medical diagnosis). In such cases, blind binning could be harmful, and domain knowledge should guide the process.

Question 9

Describe SQL techniques to perform data sampling.

Answer:

Theory

Data sampling is the process of selecting a representative subset of data from a larger dataset. In machine learning, this is essential for:

- **Faster Model Development:** Training models on a smaller sample is much faster, allowing for quicker iteration and experimentation.
- Exploratory Data Analysis (EDA): Running queries on a sample is more efficient than on the full dataset.
- **Creating Train/Test Splits:** Generating random subsets for training, validation, and testing.

There are several sampling techniques that can be implemented in SQL.

Multiple Solution Approaches & Code Examples

Approach 1: Simple Random Sampling (Database-Specific)

Some SQL dialects have built-in, efficient sampling functions. This is the preferred method when available.

Generated sql

-- PostgreSQL, Oracle, DB2

SELECT * FROM large_table TABLESAMPLE BERNOULLI (10); -- Selects each row with a 10% probability

-- SQL Server

SELECT * FROM large_table TABLESAMPLE (10 PERCENT);

IGNORE_WHEN_COPYING_START content_copy download
Use code with caution. SQL
IGNORE_WHEN_COPYING_END

- BERNOULLI (Row-level sampling): Slower but more random.
- SYSTEM (Block-level sampling): Faster but less random, as it samples entire data pages.

Approach 2: ORDER BY RAND() (Universal but often inefficient)

This method works in most SQL databases but can be very slow on large tables because it requires assigning a random number to every row and then sorting the entire table.

```
Generated sql
-- Select 1000 random rows
SELECT *
FROM large_table
ORDER BY RAND() -- or RANDOM() in PostgreSQL/SQLite
LIMIT 1000;
IGNORE_WHEN_COPYING_START
content_copy download
Use code with caution. SQL
IGNORE_WHEN_COPYING_END
```

Pitfall: Do not use this on multi-million row tables in production. The performance degradation is significant.

Approach 3: Stratified Sampling using Window Functions

This technique ensures that subgroups within the data are represented proportionally in the sample. It's crucial for imbalanced datasets, such as fraud detection.

```
Generated sql
  -- Select 10% of users from each country for a sample dataset.
WITH RankedUsers AS (
  SELECT
    -- Assign a random rank to each user within their country
    ROW NUMBER() OVER (PARTITION BY country ORDER BY RAND()) as rn,
    -- Get the total count of users in each country
    COUNT(*) OVER (PARTITION BY country) as total in country
  FROM
    users
)
SELECT
FROM
  RankedUsers
WHERE
  -- Select the top 10% of rows from each country
```

rn <= total_in_country * 0.10;

IGNORE_WHEN_COPYING_START content_copy download
Use code with caution. SQL
IGNORE_WHEN_COPYING_END

Explanation:

- 1. The RankedUsers CTE partitions the data by country.
- ROW_NUMBER() OVER (PARTITION BY country ORDER BY RAND()) assigns a random rank to each user within their country partition.
- 3. The final WHERE clause selects the top N rows from each partition, effectively creating a 10% sample from each country.

Best Practices

- **Reproducibility:** If your SQL dialect supports it, set a random seed (SET SEED ...) before sampling to ensure you can get the same sample again.
- Understand the Method: Be aware of the performance trade-offs. Use TABLESAMPLE
 if available. Use window functions for stratification. Avoid ORDER BY RAND() on large
 datasets.

Question 10

How does SQL play a role in ML model deployment?

Answer:

Theory

SQL is a cornerstone of operationalizing machine learning models (MLOps), acting as the bridge between the data layer and the ML application layer. Its role in deployment is primarily centered around data retrieval for prediction, batch processing, and logging results.

Key Roles of SQL in ML Deployment

1. Real-Time Feature Retrieval:

- For online predictions (e.g., fraud detection, real-time recommendations), the model needs the latest features for a given entity (like a user or transaction).
- An application will execute a low-latency SQL query to fetch these features from a production database. The query is typically a simple, highly-optimized SELECT on an indexed key.

Example: SELECT age, last_login_days_ago, total_purchases FROM user_features WHERE user_id = '12345';

2.

3. Batch Scoring / Batch Predictions:

- Many ML applications run on a schedule (e.g., daily, hourly) to score a large number of entities at once.
- o A scheduled job (e.g., Airflow, cron) will execute a SQL script that:
 - a. Selects the new, unscored data.
 - b. Joins it with necessary feature tables.
 - c. Passes the data to the model for prediction (or sometimes, simple models can be implemented directly in SQL with CASE statements).
 - d. INSERTs or UPDATEs a results table with the predictions.
- **Example:** A nightly job to predict churn risk for all active customers.

4.

5. In-Database Machine Learning:

- Modern data warehouses (e.g., BigQuery ML, Snowflake, SQL Server ML Services) allow you to train and deploy models directly inside the database using SQL-like syntax.
- The entire prediction process can be encapsulated in a single SQL query.
- Conceptual Example: SELECT user_id, PREDICT(my_churn_model, *) FROM new users;
- o This simplifies the architecture by keeping computation close to the data.

6.

7. Logging and Monitoring:

- After a prediction is made, it's crucial to log the inputs (features), the output (prediction), and model metadata (e.g., model version).
- This is typically done with a simple INSERT statement into a predictions_log table.
- This log table is then used for monitoring model performance, detecting data drift, and for auditing and debugging purposes. SQL is then used again to query this log table to create monitoring dashboards.

8.

Best Practices

- **Optimize for Latency:** Feature retrieval queries for real-time use cases must be highly performant. Ensure all lookup keys are indexed.
- Version Control SQL Scripts: The SQL scripts used for feature engineering and batch scoring are part of your model's logic and should be version-controlled in Git along with the model code.

Question 11

What is the significance of in-database analytics for Machine Learning?

Answer:

Theory

In-database analytics (or in-database machine learning) is the approach of building, deploying, and running analytical and machine learning models directly within the database management system (DBMS) or data warehouse. This contrasts with the traditional approach of extracting data from the database, moving it to a separate specialized environment (like a Python/R server) for processing, and then loading the results back.

Significance and Key Advantages

1. Eliminates Data Movement:

- The biggest advantage is the reduction of data transfer. Moving terabytes of data out of a database is slow, expensive (in terms of network bandwidth and egress costs), and complex (ETL pipelines).
- By bringing the computation to the data, this bottleneck is eliminated, leading to much faster end-to-end workflows.

2.

3. Enhanced Security and Governance:

 Data never leaves the secure, managed environment of the database. This simplifies compliance with regulations like GDPR and HIPAA, as there are fewer systems to audit and secure. Access controls and data governance policies of the database are automatically inherited.

4.

5. Leverages Database Scalability and Parallelism:

Modern data warehouses are built for massive parallel processing (MPP).
 In-database ML leverages this power to train models on huge datasets much faster than would be possible on a single-node machine.

6.

7. Simplified Architecture and Reduced Latency:

- The MLOps pipeline becomes simpler with fewer moving parts. There is no need for a separate ML computation cluster.
- For real-time predictions, latency is reduced because there is no network hop to an external model-serving API.

8.

9. Democratization of ML:

 It allows data analysts and engineers who are proficient in SQL to build and deploy ML models without needing deep expertise in Python or other programming languages.

10.

Code Example

This is a conceptual example from a system like Google BigQuery ML.

```
Generated sql
  -- Step 1: Train a logistic regression model directly in SQL
CREATE OR REPLACE MODEL
 'my dataset.churn model'
OPTIONS(model type='LOGISTIC REG') AS
SELECT
 -- Label
will_churn,
 -- Features
 avg_session_duration,
 monthly_spend,
 device type
FROM
training data;
-- Step 2: Make predictions using the trained model
SELECT
 user id,
 predicted_will_churn
FROM
 ML.PREDICT(MODEL `my_dataset.churn_model`,
  (SELECT user_id, avg_session_duration, monthly_spend, device_type FROM new_users)
 );
IGNORE_WHEN_COPYING_START
content copy download
Use code with caution. SQL
IGNORE_WHEN_COPYING_END
```

Pitfalls and Trade-Offs

- Limited Algorithm Choice: Users are restricted to the ML algorithms implemented by the database vendor, which may not be as extensive or cutting-edge as those in open-source libraries like scikit-learn or PyTorch.
- Vendor Lock-in: The models and workflows are tied to a specific database technology, making migration difficult.
- Potential Cost: These advanced features may be available only in the more expensive tiers of a database service.

Question 12

Explain recursive SQL queries and how they can be used to prepare data for hierarchical Machine Learning algorithms.

Answer:

Theory

A recursive SQL query is a query that calls itself to traverse hierarchical or graph-like data structures. In standard SQL, this is implemented using a **Recursive Common Table Expression (CTE)**. A recursive CTE has a defined structure:

- 1. **Anchor Member:** A base query that runs once and returns the starting set of rows (e.g., the root of the hierarchy, like a CEO).
- 2. **UNION ALL:** This operator combines the results of the anchor and recursive members.
- 3. **Recursive Member:** A query that references the CTE itself. It joins the results from the previous step with the source table to find the next level of the hierarchy (e.g., the direct reports of the managers found in the previous step).
- 4. **Termination Condition:** The recursion stops when the recursive member returns no new rows.

Preparing Data for Hierarchical ML

Many real-world datasets have inherent hierarchies (organizational charts, product categories, geographic locations). Recursive queries are perfect for "flattening" this hierarchy or extracting features from it.

Use Cases:

- **Path Enumeration:** Generating a feature that shows the full path from the root to a node (e.g., Electronics > Audio > Headphones).
- Hierarchy Depth: Calculating how many levels deep a node is.
- Subtree Aggregation: Aggregating values up the hierarchy. For example, calculating
 the total sales for a product category by summing the sales of all its subcategories and
 products.

Code Example

Let's find the management chain for each employee in an employees table (employee_id, name, manager_id).

```
Generated sql
WITH RECURSIVE EmployeeHierarchy AS (
--- Anchor Member: Select employees with no manager (top level)
SELECT
employee_id,
name,
```

```
manager id,
    1 AS hierarchy_level,
    CAST(name AS VARCHAR(255)) AS path
  FROM
    employees
  WHERE
    manager id IS NULL
  UNION ALL
  -- Recursive Member: Join employees to their managers
  SELECT
    e.employee id,
    e.name,
    e.manager id,
    eh.hierarchy_level + 1,
    eh.path || ' -> ' || e.name
  FROM
    employees e
  INNER JOIN
    EmployeeHierarchy eh ON e.manager id = eh.employee id
SELECT * FROM EmployeeHierarchy;
IGNORE_WHEN_COPYING_START
content copy download
Use code with caution. SQL
IGNORE_WHEN_COPYING_END
```

Explanation

- 1. **Anchor Member:** Selects the top-level employees (CEO/founders) where manager_id is NULL. It initializes hierarchy_level to 1 and the path to the employee's name.
- 2. **Recursive Member:** Joins the employees table back to the EmployeeHierarchy CTE. It finds all employees e whose manager is someone already in the hierarchy eh.
- 3. **Feature Generation:** For each new employee found, it calculates two new features:
 - hierarchy_level: Incremented by 1 from their manager's level.
 - o path: Appends their own name to their manager's path.

4.

5. The process repeats until it reaches the employees at the bottom of the hierarchy who manage no one.

These generated features (hierarchy_level, path) can then be used in an ML model.

Pitfalls

• Infinite Loops: If the data contains cycles (e.g., A is a manager of B, and B is a manager of A), the recursion will never terminate. It's crucial to ensure the hierarchical data is a proper directed acyclic graph (DAG) or add a depth limit to the query to prevent this.

Question 13

Describe how graph-based features can be generated from SQL data.

Answer:

Theory

Many relational datasets can be interpreted as graphs, where tables of entities represent **nodes** (e.g., customers, products) and transaction tables represent **edges** (e.g., an orders table links customers to products). Even without a dedicated graph database, SQL can be used to compute powerful graph-based features that capture the connectivity and topology of the data.

The typical SQL representation of a graph is an edge list: a table with source_node and target node columns.

Generating Graph Features with SQL

- 1. **Node Degree (Centrality):** This is the most basic feature, representing the number of connections a node has.
 - Code: A simple GROUP BY and COUNT.
 - SELECT user_id, COUNT(*) AS num_friends FROM social_graph GROUP BY user id;
 - **Use Case:** In a social network, a user with a high degree is highly connected. In fraud detection, a user connected to many fraudulent accounts is suspicious.

2.

- 3. **Neighborhood Features (Ego Graph Features):** These are aggregated features of a node's immediate neighbors.
 - Code: Requires a JOIN to access the neighbors' attributes.
 - SELECT g.user_id, AVG(u.age) AS avg_friend_age FROM social_graph g JOIN users u ON g.friend id = u.user id GROUP BY g.user id;
 - Use Case: The average age of a user's friends could be a predictive feature for targeted advertising.

4.

5. **Triangle Count & Local Clustering Coefficient:** Measures how interconnected a node's neighbors are with each other. A high clustering coefficient suggests a tightly-knit community.

 Code: This is more complex and requires multiple self-joins on the edge table to find paths of length 3 that form a triangle (A->B, B->C, and A->C).

Generated sql

-- Conceptual query for finding triangles SELECT a.source, a.target, c.target FROM edges a JOIN edges b ON a.target = b.source JOIN edges c ON b.target = c.source WHERE a.source = c.target;

6.

```
IGNORE_WHEN_COPYING_START content_copy download
Use code with caution. SQL
IGNORE_WHEN_COPYING_END
```

 Use Case: Identifying communities in social networks or collusive behavior in financial networks.

7.

- 8. **Shared Neighbors:** The number of neighbors two nodes have in common.
 - o **Code:** Requires joining the edge table with itself.
 - SELECT a.user_id, b.user_id, COUNT(*) AS shared_friends FROM social_graph a JOIN social_graph b ON a.friend_id = b.friend_id WHERE a.user_id < b.user_id GROUP BY a.user_id, b.user_id;
 - Use Case: A core feature in recommendation systems (link prediction). If two
 users share many friends, they are likely to know each other.

9.

Optimization and Best Practices

- **Performance:** Graph queries, especially those involving multiple self-joins (like triangle counting), can be extremely slow on large tables.
- **Pre-computation:** It is a common practice to pre-calculate these graph features using a scheduled SQL script and store them in a dedicated user_graph_features table. This avoids re-computing them on the fly for model training or prediction.

Question 14

What are SQL Common Table Expressions (CTEs) and how can they be used for feature generation?

Answer:

Theory

A Common Table Expression (CTE) is a named, temporary result set that you can reference within a SELECT, INSERT, UPDATE, or DELETE statement. It is defined using the WITH clause and exists only for the duration of that single query.

CTEs are primarily used for:

- Improving Readability: They break down long, complex queries into logical, sequential steps, making the code easier to read, understand, and debug compared to deeply nested subqueries.
- Modularity: Each CTE can represent a distinct transformation step.
- Recursion: CTEs are the only way to write recursive queries in standard SQL.

Using CTEs for Feature Generation

CTEs provide an elegant and organized way to build complex features that require multiple steps of aggregation and transformation. You can "chain" CTEs, where each subsequent CTE builds upon the previous one.

Code Example

Let's generate customer features like total spending, number of orders, and average order value.

```
Generated sql
  WITH
-- Step 1: Calculate total spending and order count for each customer
CustomerSpending AS (
  SELECT
    customer id,
    COUNT(order id) AS num orders,
    SUM(order value) AS total spending
  FROM
    orders
  GROUP BY
    customer id
),
-- Step 2: Join with customer data and calculate additional features
CustomerFeatures AS (
  SELECT
    c.customer_id,
    c.signup_date,
    cs.num orders,
    cs.total_spending,
    -- Calculate average order value using results from the first CTE
```

```
cs.total_spending / cs.num_orders AS avg_order_value
FROM
    customers c
JOIN
    CustomerSpending cs ON c.customer_id = cs.customer_id
)
-- Final selection of the generated feature set
SELECT * FROM CustomerFeatures;

IGNORE_WHEN_COPYING_START
content_copy download
Use code with caution. SQL
IGNORE_WHEN_COPYING_END
```

Explanation

- 1. **CustomerSpending CTE:** This first logical unit calculates two basic aggregate features: the total number of orders and total spending for each customer.
- 2. **CustomerFeatures CTE:** This second CTE takes the results from CustomerSpending and joins them with the main customers table. It then calculates a new feature, avg order value, by combining two columns from the previous CTE.
- 3. **Final SELECT:** The final query is clean and simple; it just selects all the columns from the final CTE, which now contains the complete feature set.

This structure is much clearer than a single, large query with multiple subqueries or joins.

CTE vs. Subquery vs. Temporary Table

- **Readability:** CTEs are generally more readable than nested subqueries.
- **Scope:** CTEs are scoped to a single statement. Temporary tables persist for the entire session and must be explicitly dropped.
- **Optimization:** A standard CTE is like a macro that the optimizer expands. It doesn't automatically cache results if referenced multiple times (though some advanced systems can). A temporary table is explicitly materialized, which can be better if you need to reuse the intermediate result many times.

Question 15

Explain the role of partitioning in large-scale SQL databases.

Answer:

Theory

Table partitioning is a database optimization technique used to divide a very large table into smaller, more manageable segments called partitions, while still treating it as a single logical table. The database system manages these partitions automatically based on a defined **partition key** and **partitioning scheme**.

The most common partitioning scheme is **horizontal partitioning**, where rows are distributed into different partitions based on the value of one or more columns (the partition key). A common choice for the partition key is a date or timestamp column.

Role and Benefits in Machine Learning Contexts

1. Massive Query Performance Improvement via Partition Pruning:

- This is the most significant benefit. When a query includes a filter (WHERE clause) on the partition key, the database optimizer can perform partition pruning. This means it knows it only needs to scan the data in the relevant partition(s) and can completely ignore all others.
- ML Use Case: When creating a training set, you often filter by a time range (e.g., "use data from the last 2 years"). If the table is partitioned by month, the query will only read the 24 relevant monthly partitions instead of scanning the entire multi-terabyte table, resulting in a dramatic speedup.

2.

3. Efficient Data Management and Maintenance:

- Operations can be performed at the partition level. For example, you can
 efficiently delete old data by dropping an entire partition (ALTER TABLE ... DROP
 PARTITION), which is an almost instantaneous metadata operation, far faster
 than a DELETE statement that has to scan and log every row.
- ML Use Case: This is crucial for managing data for model retraining. For a model that retrains on a 1-year sliding window, you can easily drop the oldest month's partition and add the newest one, simplifying data lifecycle management.

4.

5. Improved Manageability:

 Smaller partitions can be individually backed up, restored, or indexed, making database administration tasks more manageable on very large tables.

6.

Code Example (Conceptual)

First, you define a partitioned table.

Generated sql
-- Conceptual DDL for a partitioned table
CREATE TABLE sales (
sale_id INT,
product_id INT,
sale_amount DECIMAL(10, 2),

```
sale_date DATE
) PARTITION BY RANGE (sale_date);
```

-- Partitions are created for specific ranges

CREATE TABLE sales_2022_q4 PARTITION OF sales FOR VALUES FROM ('2022-10-01') TO ('2023-01-01');

CREATE TABLE sales_2023_q1 PARTITION OF sales FOR VALUES FROM ('2023-01-01') TO ('2023-04-01');

...

IGNORE_WHEN_COPYING_START content_copy download
Use code with caution. SQL
IGNORE_WHEN_COPYING_END

Then, a query on the partition key will benefit from pruning.

Generated sql

-- This query will ONLY scan the sales_2023_q1 partition.

SELECT * FROM sales WHERE sale date >= '2023-02-01' AND sale date < '2023-03-01';

IGNORE_WHEN_COPYING_START content_copy download
Use code with caution. SQL
IGNORE_WHEN_COPYING_END

Best Practices

- Choose a partition key that is used in almost all performance-critical queries. A date column is a very common and effective choice.
- Avoid creating too many small partitions, as this can add overhead to the query planner.

Question 16

Describe how you could use SQL to report the performance metrics of a Machine Learning model.

Answer:

Theory

After an ML model is deployed, its predictions are typically logged to a database table for monitoring and analysis. This log table usually contains the entity identifier, the features used (or

a reference to them), the model's prediction, the predicted probability, and, once it becomes available, the actual outcome (the ground truth). SQL is the perfect tool for querying this log table to calculate a wide range of performance metrics.

```
Assume we have a prediction_logs table like this: prediction_logs(prediction_id, model_version, entity_id, predicted_label, true_label, predicted_probability)
```

Calculating Metrics with SQL

1. Accuracy: The proportion of correct predictions.

```
Generated sql
SELECT
model_version,
AVG(CASE WHEN predicted_label = true_label THEN 1.0 ELSE 0.0 END) AS accuracy
FROM
prediction_logs
GROUP BY
model_version;

IGNORE_WHEN_COPYING_START
content_copy download
Use code with caution. SQL
IGNORE_WHEN_COPYING_END
```

Explanation: CASE WHEN ... creates a column of 1s (correct) and 0s (incorrect). AVG() on this column gives the accuracy.

2. Confusion Matrix: A breakdown of true positives, true negatives, false positives, and false negatives.

```
Generated sql
SELECT
true_label,
predicted_label,
COUNT(*) as count
FROM
prediction_logs
WHERE
model_version = 'v2.1'
GROUP BY
true_label, predicted_label;
```

-- Or more pivot-table like:

```
SELECT
  true_label,
  COUNT(CASE WHEN predicted label = 1 THEN 1 END) AS predicted as 1,
  COUNT(CASE WHEN predicted label = 0 THEN 1 END) AS predicted as 0
FROM
  prediction logs
GROUP BY
  true_label;
IGNORE WHEN COPYING START
content copy download
Use code with caution. SQL
IGNORE WHEN COPYING END
3. Precision, Recall, and F1-Score: These can be calculated from the confusion matrix values.
Generated sql
  WITH Metrics AS (
  SELECT
    SUM(CASE WHEN predicted label = 1 AND true label = 1 THEN 1 ELSE 0 END) AS
true positives,
    SUM(CASE WHEN predicted label = 1 AND true label = 0 THEN 1 ELSE 0 END) AS
false positives,
    SUM(CASE WHEN predicted label = 0 AND true label = 1 THEN 1 ELSE 0 END) AS
false negatives
  FROM prediction_logs
SELECT
  true_positives * 1.0 / (true_positives + false_positives) AS precision,
  true positives * 1.0 / (true positives + false negatives) AS recall
FROM
  Metrics;
IGNORE_WHEN_COPYING_START
content copy download
Use code with caution. SQL
IGNORE_WHEN_COPYING_END
4. Data for ROC/AUC Curve: To plot an ROC curve, you need the True Positive Rate (TPR)
```

and False Positive Rate (FPR) at various probability thresholds. This can be generated using

window functions.

-- Generate points needed to plot an ROC curve

Generated sql

SELECT

predicted probability,

SUM(CASE WHEN true_label = 1 THEN 1 ELSE 0 END) OVER (ORDER BY predicted_probability DESC) * 1.0 / SUM(CASE WHEN true_label = 1 THEN 1 ELSE 0 END) OVER() AS tpr,

SUM(CASE WHEN true_label = 0 THEN 1 ELSE 0 END) OVER (ORDER BY predicted_probability DESC) * 1.0 / SUM(CASE WHEN true_label = 0 THEN 1 ELSE 0 END) OVER() AS fpr FROM prediction logs;

IGNORE_WHEN_COPYING_START content_copy download
Use code with caution. SQL
IGNORE_WHEN_COPYING_END

Use Cases

- Building automated performance monitoring dashboards (e.g., in Metabase, Tableau, Power BI) that are powered by these SQL queries.
- Comparing the performance of different model versions in an A/B test.
- Segmenting performance by different data slices (e.g., "What is the model's accuracy for users in Germany?").

Question 17

Describe how you would version control the datasets used for building Machine Learning models in SQL.

Answer:

Theory

Versioning datasets is a critical component of reproducible machine learning. It ensures that you can always trace a model back to the exact data on which it was trained. While Git is excellent for versioning code, it's not designed for versioning large datasets. In an SQL environment, data versioning requires different strategies.

Multiple Solution Approaches

1. Snapshot Tables / Full Copies:

- Method: The most straightforward approach. When a dataset is ready for training, you
 create a complete, immutable copy of it in a new table with a version identifier in its
 name.
- Example: CREATE TABLE training_set_churn_v3_20231115 AS SELECT ... FROM ...;
- Pros: Conceptually simple, explicit, and guarantees perfect reproducibility.
- **Cons:** Extremely high storage cost due to data duplication. Not feasible for very large datasets.

2. Versioning the Generation Code (Code-as-Version):

- Method: Instead of versioning the data itself, you version the SQL script (generate_features.sql) that produces the training dataset. This script is stored in a version control system like Git, alongside the model training code.
- **Pros:** Very lightweight and integrates perfectly with standard developer workflows (Git).
- Cons: Does not guarantee perfect reproducibility if the underlying source data can change (i.e., rows are updated or deleted). It only works reliably if the source data is immutable or append-only.

3. Versioning via Timestamps (Slowly Changing Dimensions):

- Method: Design source tables to capture history. For example, use a Type 2 Slowly
 Changing Dimension (SCD) approach, where tables have valid_from and valid_to
 timestamp columns. To reconstruct a dataset, the feature generation query includes a
 WHERE clause to select data AS OF a specific point in time.
- Example Query: SELECT * FROM user_profiles WHERE '2023-10-01' BETWEEN valid from AND valid to;
- **Pros:** Storage-efficient, provides a full historical view of the data.
- **Cons:** Significantly increases the complexity of queries and requires careful data pipeline design to maintain the history tables.

4. Using a Feature Store:

- Method: A feature store is a centralized system designed specifically to manage ML features. Features are versioned, stored, and can be retrieved for training (as a point-in-time correct dataset) or serving (for a single entity). The feature store itself is often built on top of a database.
- **Pros:** The gold standard for mature MLOps. Solves versioning, reproducibility, and the online/offline skew problem.
- **Cons:** Requires significant engineering effort to build or cost to buy a commercial solution.

Recommended Best Practice

A practical and effective hybrid approach is to **version the feature generation SQL script in Git** (Approach 2) and, for critical production models, **also create a snapshot table** (Approach

1) of the exact training data used. This provides both lightweight development versioning and ironclad production reproducibility.

Question 18

What is Data Lineage, and how can you track it using SQL?

Answer:

Theory

Data lineage provides a complete history of data's journey. It tracks the data's origin, what happens to it over time, and where it moves. It answers critical questions like:

- **Source:** Where did this data come from?
- **Transformation:** What calculations or business logic were applied to create this feature?
- **Destination:** Which reports, dashboards, and ML models consume this data?

In the context of machine learning, data lineage is essential for reproducibility, debugging, governance, and impact analysis.

Tracking Data Lineage with SQL

Tracking lineage manually or automatically with SQL-based tools is challenging but possible through several methods:

1. Manual Documentation and Code Parsing:

- The most basic method is to manually read SQL scripts and document dependencies. If you see CREATE TABLE features AS SELECT col1, col2 FROM source_data;, you know the features table depends on source_data.
- This can be partially automated with scripts that parse SQL code to identify source and target tables in FROM and JOIN clauses.

2.

3. Leveraging Database Metadata (Information Schema):

- SQL databases maintain a set of system tables, often called the **Information Schema**, which contains metadata about all the objects in the database (tables, columns, views, etc.).
- You can query the information schema to discover relationships. For example, querying INFORMATION_SCHEMA.VIEW_TABLE_USAGE can tell you which tables a specific view depends on. This helps track lineage through views.

4.

5. Using Logging and Audit Trails:

 Many databases can log all queries that are executed. By parsing these query logs, you can reconstruct a history of data movement and transformation. For example, you can see that Job_X ran an INSERT INTO table_B SELECT ... FROM table_A, which establishes a lineage link from A to B.

6.

7. Adopting Data Build Tools (like dbt):

- This is the modern, preferred approach. Tools like dbt (Data Build Tool) allow you
 to define your data transformations as a series of SELECT statements. Dbt then
 automatically infers the dependencies between all your models (tables/views) by
 parsing the ref() functions in your SQL code.
- It uses these dependencies to build a Directed Acyclic Graph (DAG) of your entire data pipeline, which is the data lineage. This lineage graph can be visualized, making it easy to understand complex data flows.

8.

Importance in Machine Learning

- **Reproducibility:** To reproduce a model, you must reproduce its feature data. Lineage tells you exactly how that feature data was created.
- **Debugging and Root Cause Analysis:** If a model's performance suddenly drops, lineage allows you to trace its input features back to the source. You might discover a problem in an upstream data pipeline that corrupted the data.
- **Impact Analysis:** Before changing or deprecating a source table, you can use lineage to see exactly which downstream ML models and dashboards will be affected.
- **Governance and Compliance:** For regulations like GDPR, lineage is required to track the use of personally identifiable information (PII) throughout your systems.

SQL ML Interview Questions - Coding Questions

Question 1

Write a SQL query that joins two tables and retrieves only the rows with matching keys.

Answer:

Theory

To retrieve only the rows with matching keys from two tables, you use an INNER JOIN. This is the most common type of join. It creates a new result table by combining column values of two tables (table1 and table2) based upon the join-predicate. The query scans table1 and for each row it finds, it scans table2 for all matching rows based on the ON condition.

Code Example

Let's assume we have a customers table and an orders table. We want to find all orders placed by customers who exist in our customers table.

Sample Data:

```
customers table:
| customer_id | name |
|------|
| 1 | Alice |
| 2 | Bob |
| 3 | Charlie |

orders table:
| order_id | customer_id | amount |
|------|
| 101 | 1 | 50 |
| 102 | 2 | 75 |
| 103 | 1 | 120 |
| 104 | 4 | 30 | (No matching customer)
```

SQL Query:

```
Generated sql
SELECT
c.customer_id,
c.name,
o.order_id,
o.amount
FROM
customers AS c
INNER JOIN
orders AS o ON c.customer_id = o.customer_id;
```

Result:

```
| customer_id | name | order_id | amount |
|------|------|------|
| 1 | Alice | 101 | 50 |
| 2 | Bob | 102 | 75 |
| 1 | Alice | 103 | 120 |
```

Explanation

1. **SELECT c.customer_id**, ...: Specifies the columns we want to see in the final result.

- 2. **FROM customers AS c**: Defines the first (left) table and gives it a short alias c for readability.
- 3. **INNER JOIN orders AS o**: Specifies that we are joining with the orders table (aliased as o) and that we only want rows that have a match in both tables.
- 4. **ON c.customer_id = o.customer_id**: This is the join condition. It tells the database to match rows where the customer_id value is the same in both the customers and orders tables. The row for order 104 is excluded because its customer_id of 4 does not exist in the customers table.

Best Practices

- **Use Aliases**: Always use table aliases (like c and o) to make your query shorter and more readable, especially when joining multiple tables.
- Index Join Keys: For optimal performance, the columns used in the ON clause (customer_id in this case) should be indexed. This allows the database to find matches much faster than scanning the entire table.

Question 2

Create a SQL query to pivot a table transforming rows into columns.

Answer:

Theory

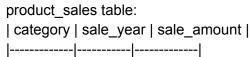
Pivoting is a data transformation technique that rotates data from a state of rows to a state of columns. It's commonly used to create summary reports or to prepare a "wide" format dataset for machine learning models where each distinct category value becomes a new feature column.

While some SQL dialects have a specific PIVOT operator (like SQL Server or Oracle), the most portable and universal method is to use a combination of an aggregate function and CASE statements.

Code Example

Let's say we have a product_sales table recording sales amount by product category for different years. We want to transform this into a report where each year is a column.

Sample Data:



```
| Electronics | 2022 | 5000 |
| Books | 2022 | 1500 |
| Electronics | 2023 | 6200 |
| Books | 2023 | 1800 |
```

SQL Query (Standard SQL):

```
Generated sql
  SELECT
  category.
  SUM(CASE WHEN sale_year = 2022 THEN sale_amount ELSE 0 END) AS sales_2022,
  SUM(CASE WHEN sale year = 2023 THEN sale amount ELSE 0 END) AS sales 2023
FROM
  product_sales
GROUP BY
  category;
IGNORE WHEN_COPYING_START
content_copy download
Use code with caution. SQL
IGNORE_WHEN_COPYING_END
```

Result:

```
| category | sales 2022 | sales 2023 |
|-----|
| Books | 1500 | 1800 |
| Electronics | 5000 | 6200 |
```

Explanation

- 1. **GROUP BY category**: This clause groups all rows for the same category (Electronics, Books) together. The final result will have one row for each unique category.
- SUM(CASE WHEN ...): This is the core of the pivot logic. It's an aggregate function that works on each group defined by GROUP BY.
 - CASE WHEN sale_year = 2022 THEN sale_amount ELSE 0 END: For the sales 2022 column, this expression checks each row within a group. If the year is 2022, it takes the sale amount; otherwise, it takes 0.
 - SUM(...): The SUM function then adds up these values. For the 'Electronics' group, it calculates SUM(5000, 0) which equals 5000. For the 'Books' group, it's SUM(1500, 0).
 - The same logic applies to the sales 2023 column.

3.

Pitfalls

Static Columns: This approach requires you to know the values you want to pivot into columns (e.g., 2022, 2023) ahead of time. If a new year 2024 is added to the data, you must manually add a new SUM(CASE ...) expression to the query. This is a limitation of static SQL.

Question 3

Write a SQL query to calculate moving averages.

Answer:

Theory

A moving average (or rolling average) is a common time-series analysis technique used to smooth out short-term fluctuations and highlight longer-term trends or cycles. It is calculated by averaging a subset of data points over a specific window of time. In SQL, this is achieved efficiently using window functions.

Code Example

Let's calculate a 7-day moving average of daily revenue for a specific product.

Sample Data:

```
daily_revenue table:

| report_date | product_id | revenue |

|------|-----|-----|

| 2023-11-01 | 101 | 150 |

| 2023-11-02 | 101 | 160 |

| ... | ... | ... |

| 2023-11-07 | 101 | 180 |

| 2023-11-08 | 101 | 190 |
```

SQL Query:

```
Generated sql
SELECT
report_date,
product_id,
revenue,
AVG(revenue) OVER (
PARTITION BY product_id
ORDER BY report_date
ROWS BETWEEN 6 PRECEDING AND CURRENT ROW
```

```
) AS "7_day_moving_avg"
FROM
daily_revenue;

IGNORE_WHEN_COPYING_START
content_copy download
Use code with caution. SQL
IGNORE_WHEN_COPYING_END
```

Explanation

- AVG(revenue) OVER (...): This specifies that we want to compute an average using a window function.
- **PARTITION BY product_id**: This clause divides the data into separate groups for each product_id. The moving average calculation will restart for each product, which is crucial if you have data for multiple products in the same table.
- ORDER BY report_date: Within each partition, the rows are ordered by report_date.
 This is essential for any time-series calculation to ensure the window is applied chronologically.
- ROWS BETWEEN 6 PRECEDING AND CURRENT ROW: This is the frame clause that
 defines the "window" for the average. It tells the database to include the current row and
 the six rows that came before it (totaling 7 rows) in the calculation. For dates with
 missing data, RANGE is often a better choice.

Use Cases

- Feature Engineering: Create smoothed features for time-series forecasting models to reduce noise.
- Financial Analysis: Track trends in stock prices or company revenues.
- Business Intelligence: Create smoother trend lines in performance dashboards.

Question 4

How can you create lagged features in SQL?

Answer:

Theory

A lagged feature is a feature that contains a value from a previous time step. For a given observation at time t, a lagged feature would be the value of a variable at time t-n, where n is the lag period. These features are fundamental for sequential and time-series models as they

capture temporal dependencies. The LAG() window function in SQL is designed specifically for this purpose.

Code Example

Let's create a feature that shows each customer's previous purchase amount.

```
Sample Data:
```

```
purchases table:
| customer_id | purchase_date | amount |
|-----|
| 1 | 2023-01-10 | 50 |
| 1 | 2023-02-15 | 75 |
| 2 | 2023-03-05 | 100 |
| 1 | 2023-04-01 | 60 |
SQL Query:
Generated sql
  SELECT
  customer_id,
  purchase_date,
  amount.
  LAG(amount, 1, 0) OVER (
    PARTITION BY customer id
    ORDER BY purchase date
  ) AS previous_purchase_amount
FROM
  purchases;
IGNORE WHEN COPYING START
content_copy download
Use code with caution. SQL
IGNORE_WHEN_COPYING_END
Result:
| customer_id | purchase_date | amount | previous_purchase_amount |
|-----|
| 1 | 2023-01-10 | 50 | 0 |
| 1 | 2023-02-15 | 75 | 50 |
| 1 | 2023-04-01 | 60 | 75 |
```

Explanation

| 2 | 2023-03-05 | 100 | 0 |

• LAG(amount, 1, 0): This is the core function call.

- o amount: The column whose previous value we want to retrieve.
- 1: The lag offset. 1 means we want the value from the immediately preceding row.
- 0: The default value to use if there is no preceding row (i.e., for the first purchase of each customer). Using a default avoids NULLs.

•

- **OVER (...)**: This defines the window for the LAG function.
 - PARTITION BY customer_id: This is critical. It ensures that we only look at previous purchases for the same customer. The calculation is isolated within each customer's data.
 - ORDER BY purchase_date: This orders each customer's purchases chronologically, defining what "previous" means.

•

Use Cases

- **Predictive Modeling**: Create features like "time since last event," "previous page visited," or "last month's sales."
- **Change Detection**: Calculate the difference between the current value and the previous value to detect significant changes.

Question 5

Describe how to compute a ratio feature within groups using SQL.

Answer:

Theory

Calculating a ratio within a group helps normalize a value against its local context. For example, instead of using a raw salary value, a feature representing an employee's salary as a percentage of their department's total salary can be more predictive. This type of feature is easily created using SQL window functions, which can compute a group-level aggregate and make it available on every row of that group.

Code Example

Let's calculate a feature that represents each product's sales as a ratio of its category's total sales.

Sample Data:

product_sales table:
| product_name | category | sales |

```
|------|
| Laptop A | Electronics | 5000 |
| Laptop B | Electronics | 3000 |
| Book X | Books | 1500 |
| Book Y | Books | 500 |
```

SQL Query:

```
Generated sql
SELECT
product_name,
category,
sales,
sales * 1.0 / SUM(sales) OVER (PARTITION BY category) AS ratio_of_category_sales
FROM
product_sales;

IGNORE_WHEN_COPYING_START
content_copy download
Use code with caution. SQL
IGNORE_WHEN_COPYING_END
```

Result:

Explanation

- 1. SUM(sales) OVER (PARTITION BY category): This is the key window function.
 - PARTITION BY category divides the data into partitions, one for 'Electronics' and one for 'Books'.
 - SUM(sales) is calculated for each partition independently. For the 'Electronics' partition, the sum is 5000 + 3000 = 8000. For the 'Books' partition, the sum is 1500 + 500 = 2000.
 - Crucially, this partition-level sum (8000 or 2000) is appended to every row within its respective partition.

2.

- 3. **sales** * **1.0** / ...: In the outer SELECT, we take the individual row's sales value and divide it by the partition-level sum we just calculated.
 - Multiplying by 1.0 is a common trick to force floating-point division to get a decimal result, avoiding integer division which would result in 0 in many cases.

Pitfalls and Optimization

- **Division by Zero**: If a group's total sum could be zero, the query would fail. To prevent this, wrap the denominator in NULLIF: NULLIF(SUM(sales) OVER (...), 0).
- Performance: On databases that support it, defining an index on the PARTITION BY key (category in this case) can improve performance.

Question 6

Write a SQL query that identifies and removes duplicate records from a dataset.

Answer:

Theory

Identifying and removing duplicate records is a common data cleaning task. A record is considered a duplicate if another row exists with the same values across a specific set of columns (the "uniqueness key"). The most robust and flexible way to handle this in SQL is to use the ROW_NUMBER() window function within a Common Table Expression (CTE) to assign a unique rank to each record within a duplicate group, and then delete all but one.

Code Example

Let's assume we have a contacts table where there are duplicate entries for the same email address. We want to keep only the most recently added record for each email.

Sample Data:

Multiple Solution Approaches

Approach 1: DELETE using a CTE (Recommended)

```
Generated sql
WITH NumberedContacts AS (
SELECT
```

```
contact_id,
ROW_NUMBER() OVER(
PARTITION BY email
ORDER BY created_at DESC
) AS rn
FROM
contacts
)
DELETE FROM contacts
WHERE contact_id IN (
SELECT contact_id FROM NumberedContacts WHERE rn > 1
);
IGNORE_WHEN_COPYING_START
content_copy download
Use code with caution. SQL
IGNORE_WHEN_COPYING_END
```

Explanation:

- 1. WITH NumberedContacts AS (...): We define a CTE to first identify the duplicates.
- 2. ROW_NUMBER() OVER(...):
 - o PARTITION BY email: This groups all rows with the same email together.
 - ORDER BY created_at DESC: Within each group, it orders the records by their creation date in descending order (newest first).
 - ROW_NUMBER() then assigns a sequential number to each row in the group.
 The newest record gets rn = 1, the second newest gets rn = 2, and so on.

3.

- 4. **DELETE FROM contacts** ...: The outer query deletes rows from the original contacts table.
- 5. **WHERE contact_id IN (...)**: The WHERE clause specifies which rows to delete by selecting the contact_ids from our CTE where the row number (rn) is greater than 1. These are the duplicates we want to remove.

Approach 2: Create a New Table

This approach avoids a DELETE statement, which can be safer.

```
Generated sql
-- Step 1: Create a new table with only the unique records

CREATE TABLE contacts_deduped AS

WITH NumberedContacts AS (
SELECT
*,
ROW NUMBER() OVER(PARTITION BY email ORDER BY created at DESC) AS rn
```

```
FROM contacts
)
SELECT contact_id, email, created_at FROM NumberedContacts WHERE rn = 1;
-- Step 2: (Optional) Rename tables to replace the original
-- RENAME TABLE contacts TO contacts_old;
-- RENAME TABLE contacts_deduped TO contacts;

IGNORE_WHEN_COPYING_START content_copy download
Use code with caution. SQL
IGNORE WHEN COPYING END
```

Best Practices

- Always Backup: Before running a destructive DELETE operation on a production table, always create a backup.
- **Define Duplicates Clearly**: The PARTITION BY clause is critical. It defines what combination of columns makes a record a duplicate.
- **Define Which Record to Keep**: The ORDER BY clause inside OVER() is equally important. It determines which of the duplicate records will be assigned rn = 1 and therefore be kept.

Question 7

In SQL, how would you format strings or concatenate columns for text-based Machine Learning features?

Answer:

Theory

Concatenating (joining) multiple string columns or formatting them into a single text feature is a common preprocessing step for text-based ML models. For example, combining a product's title, description, and category into one text block can provide a richer input for NLP models. SQL provides standard functions like CONCAT or operators like || for this purpose.

Code Example

Let's create a combined text feature from a product's name, brand, and description.

Sample Data:

```
products table:
| product_id | name | brand | description |
|------|------|------|
| 101 | Smart Watch | TechCo| A watch that tracks fitness. |
| 102 | Running Shoe| FitGear| NULL |
```

SQL Query:

Note: CONCAT_WS (Concatenate With Separator) is available in PostgreSQL, MySQL, and SQL Server. For other databases, you might chain CONCAT or || operators.

Result:

```
| product_id | combined_feature | |------| | | 101 | Name: Smart Watch | Brand: TechCo | Desc: A watch that tracks fitness. | | 102 | Name: Running Shoe | Brand: FitGear | Desc: N/A |
```

Explanation

- || Operator: A standard way to concatenate strings. 'Name: ' || name results in 'Name: Smart Watch'.
- **COALESCE**(**description**, 'N/A'): This function is crucial for handling NULL values. It returns the first non-NULL value in its argument list. If description is NULL, it returns the string 'N/A', preventing the entire concatenated string from becoming NULL.
- CONCAT_WS(' | ', ...): This function joins multiple strings together using a specified separator (in this case, ' | '). It is often preferred because it automatically ignores NULL

arguments (after they've been handled by COALESCE if needed) and makes the query cleaner than chaining many || operators.

Use Cases

- **NLP Feature Engineering**: Creating a single document from multiple text fields before feeding it into a TF-IDF vectorizer or a language model.
- **Categorical Features**: Combining multiple low-cardinality categorical features (e.g., country and device type) into a single interaction feature ('US-Mobile').
- Search: Creating a searchable text field that aggregates information from multiple columns.

Question 8

Write a SQL stored procedure that calls a Machine Learning scoring function.

Answer:

Theory

A stored procedure is a pre-compiled set of one or more SQL statements that are stored in the database. In the context of MLOps, a stored procedure can be used to create a reusable, controlled, and secure endpoint for running a deployed ML model. This allows applications to get predictions without having to implement the feature retrieval logic themselves.

The exact syntax for calling an ML model varies significantly across database platforms (e.g., SQL Server's sp_execute_external_script, Snowflake's external functions, BigQuery ML's ML.PREDICT). The following example uses a generic, conceptual syntax to illustrate the logic.

Code Example (Conceptual - T-SQL/SQL Server like)

This procedure takes a customer_id, retrieves the latest features for that customer, and calls a hypothetical prediction function to score their churn risk.

Generated sal

CREATE OR ALTER PROCEDURE dbo.GetChurnPrediction

@CustomerID INT,

@PredictedChurnProba FLOAT OUTPUT

AS

BEGIN

-- Enforce security and prevent SQL injection SET NOCOUNT ON;

```
-- Step 1: Declare variables to hold the features
  DECLARE @TenureMonths INT;
  DECLARE @MonthlySpend DECIMAL(10, 2);
  -- Step 2: Retrieve the latest features for the given customer
  SELECT TOP 1
    @TenureMonths = DATEDIFF(month, signup date, GETDATE()),
    @MonthlySpend = avg_monthly_spend
  FROM
    customer features
  WHERE
    customer_id = @CustomerID;
  -- Step 3: Call the hypothetical machine learning scoring function
  -- The actual call syntax is database-specific.
  -- This example assumes a function that takes features and returns a score.
  EXECUTE PREDICT(
    MODEL = 'ChurnModel v2',
    INPUT_DATA = { 'tenure': @TenureMonths, 'spend': @MonthlySpend },
    OUTPUT SCORE = @PredictedChurnProba OUTPUT
  );
END;
GO
-- How to call the procedure
DECLARE @ChurnScore FLOAT;
EXEC dbo.GetChurnPrediction @CustomerID = 12345, @PredictedChurnProba =
@ChurnScore OUTPUT;
SELECT @ChurnScore AS ChurnProbability;
IGNORE WHEN COPYING START
content copy download
Use code with caution. SQL
IGNORE_WHEN_COPYING_END
```

Explanation

- 1. **CREATE OR ALTER PROCEDURE** ...: Defines the procedure name (GetChurnPrediction) and its parameters. @CustomerID is an input, and @PredictedChurnProba is an output parameter to return the result.
- Feature Retrieval: A standard SELECT query fetches the necessary features for the model from a pre-built customer_features table. This logic is now centralized and hidden from the client application.

- 3. **Model Scoring**: The EXECUTE PREDICT(...) line is a conceptual representation of calling the model. It passes the retrieved features as input and captures the model's output score in the @PredictedChurnProba variable.
- 4. **OUTPUT Parameter**: The procedure returns the prediction score via the OUTPUT parameter, which the calling application can then use.

Best Practices

- Security: Use stored procedures to grant limited, indirect access to data. Applications
 can EXECUTE the procedure without having direct SELECT permissions on the
 underlying feature tables.
- **Performance**: Feature retrieval queries within the procedure must be highly optimized, typically using an index on the lookup key (customer_id).
- Modularity: Encapsulating the scoring logic makes model updates easier. You can
 update the procedure to call a new model version (ChurnModel_v3) without changing
 any client application code.

Question 9

How would you construct a complex SQL query to extract time series features for a Machine Learning model?

Answer:

Theory

Extracting a rich set of time series features often requires multiple stages of transformation and aggregation. A complex query for this task is best constructed using Common Table Expressions (CTEs) to break the logic into clean, readable, and sequential steps. We can generate features like rolling averages, lagged values, and period-over-period growth rates in a single, powerful query.

Code Example

Let's build a feature set for a sales forecasting model. For each product and each month, we want to calculate:

- 1. Total monthly sales.
- 2. Sales from the previous month (lagged feature).
- 3. A 3-month rolling average of sales.
- 4. Month-over-month sales growth percentage.

Generated sql

```
WITH MonthlySales AS (
  -- Step 1: Aggregate daily sales into monthly totals for each product.
  SELECT
    product id,
    DATE TRUNC('month', sale date)::date AS sale month,
    SUM(sale amount) AS total monthly sales
  FROM
    sales data
  GROUP BY
    product id, sale month
),
WindowFeatures AS (
  -- Step 2: Use window functions to create time-series features.
  SELECT
    product id,
    sale_month,
    total_monthly_sales,
    -- Feature: Sales from the previous month
    LAG(total_monthly_sales, 1, 0) OVER (PARTITION BY product_id ORDER BY
sale month) AS previous month sales,
    -- Feature: 3-month rolling average
    AVG(total_monthly_sales) OVER (
      PARTITION BY product id ORDER BY sale month
      ROWS BETWEEN 2 PRECEDING AND CURRENT ROW
    ) AS "3_month_rolling_avg"
  FROM
    MonthlySales
)
-- Step 3: Calculate final features like growth rate.
SELECT
  product id,
  sale month,
  total_monthly_sales,
  previous month sales,
  "3 month rolling avg",
  -- Feature: Month-over-month growth rate
  (total_monthly_sales - previous_month_sales) * 1.0 / NULLIF(previous_month_sales, 0) AS
"mom growth rate"
FROM
  WindowFeatures
ORDER BY
  product_id, sale_month;
IGNORE_WHEN_COPYING_START
```

content_copy download
Use code with caution. SQL
IGNORE WHEN COPYING END

Explanation

- MonthlySales CTE: The first step is to roll up the raw, daily data into the desired time grain (monthly). This GROUP BY simplifies the dataset for the next stage.
 DATE TRUNC is a standard function for this.
- 2. **WindowFeatures CTE**: This is where the core feature engineering happens. It takes the monthly sales data and applies window functions:
 - LAG() calculates the previous month sales.
 - AVG() with a ROWS BETWEEN frame calculates the 3 month rolling avg.
 - All window functions are partitioned by product_id to ensure calculations are done independently for each product.

3.

4. **Final SELECT**: The final query uses the features generated in the previous CTE to calculate a derived feature: the month-over-month growth rate. Using NULLIF prevents division-by-zero errors.

Best Practices

- Modularity with CTEs: Breaking down the logic into CTEs makes the query vastly easier to debug and maintain. You can test each CTE independently.
- **Performance**: Ensure the source table is indexed on sale_date and product_id to speed up the initial aggregation and partitioning.
- Data Granularity: Always start by aggregating your data to the correct time series interval (e.g., daily, weekly, monthly) before applying window functions.

Question 10

Discuss ways to implement regular expressions in SQL for natural language processing tasks.

Answer:

Theory

Regular expressions (regex) are a powerful tool for pattern matching in text, making them very useful for cleaning and feature engineering in NLP tasks directly within SQL. Most modern SQL dialects provide functions to leverage regex. The common functions are:

- Matching/Searching (REGEXP_LIKE, ~): Checks if a string contains a pattern. Returns a boolean.
- Extracting (REGEXP_SUBSTR, REGEXP_MATCHES): Extracts the part of a string that matches a pattern.
- Replacing (REGEXP_REPLACE): Replaces parts of a string that match a pattern with another string.

Code Examples (using PostgreSQL syntax)

Use Case 1: Feature Engineering - Check for a specific keyword

Create a binary feature indicating if a product review mentions "delivery".

```
Generated sql
SELECT
review_text,
(review_text ~* 'delivery') AS mentions_delivery -- ~* for case-insensitive match
FROM product_reviews;

IGNORE_WHEN_COPYING_START
content_copy download
Use code with caution. SQL
IGNORE_WHEN_COPYING_END
```

Use Case 2: Data Cleaning - Remove all numbers from text

Clean text data before tokenization.

IGNORE WHEN COPYING END

```
Generated sql
SELECT
review_text,
REGEXP_REPLACE(review_text, '\d+', ", 'g') AS text_without_numbers
-- '\d+' matches one or more digits. 'g' flag means replace all occurrences.
FROM product_reviews;

IGNORE_WHEN_COPYING_START
content_copy download
Use code with caution. SQL
```

Use Case 3: Information Extraction - Extract email addresses

Extract structured information from unstructured text.

```
Generated sql
SELECT
comment_text,
(REGEXP MATCHES(comment text, '\S+@\S+\.\S+'))[1] AS extracted email
```

-- \S+ matches one or more non-whitespace characters. FROM user_comments;

IGNORE_WHEN_COPYING_START content_copy download
Use code with caution. SQL
IGNORE_WHEN_COPYING_END

Explanation

- **~* Operator**: A PostgreSQL shorthand for case-insensitive regex matching. The query (review text **~*** 'delivery') returns true if "delivery", "Delivery", etc., is found.
- **REGEXP_REPLACE(string, pattern, replacement, flags)**: Finds all substrings matching pattern and replaces them. The 'g' (global) flag is crucial for replacing all instances, not just the first one.
- **REGEXP_MATCHES(string, pattern)**: Returns a set of matching substrings. Accessing the first element [1] gives the first full match.

Debugging and Performance

- **Complexity**: Regex can be hard to write and debug. It's best to test patterns on a small sample of data or using an online regex tester before running them on a large table.
- Performance: Regex operations, especially on long text columns, can be computationally intensive. If a regex-based feature is needed frequently, it's often a best practice to pre-compute it and store it in a dedicated column rather than calculating it on the fly in every query.

Question 11

Write a SQL script to identify and replace missing values with the column mean.

Answer:

Theory

Imputing missing values (NULLs) is a standard data preprocessing step. Replacing them with the column's mean is a common and simple strategy. The most efficient way to do this in SQL without multiple queries is to use window functions. A window function can calculate the mean of the entire column and make that value available on every row, allowing you to replace NULLs in a single pass.

Code Example

Let's impute missing age values in a users table with the average age of all users.

Sample Data:

```
users table:

| user_id | age |

|-----|

| 1 | 30 |

| 2 | 25 |

| 3 | NULL|

| 4 | 35 |

| 5 | NULL|
```

Multiple Solution Approaches

Approach 1: SELECT with Window Function (for creating a new view/table)

This is non-destructive and ideal for feature generation.

```
Generated sql
SELECT
user_id,
COALESCE(
age,
AVG(age) OVER ()
) AS imputed_age
FROM
users;

IGNORE_WHEN_COPYING_START
content_copy download
Use code with caution. SQL
IGNORE_WHEN_COPYING_END
```

Explanation:

- 1. **AVG(age) OVER ()**: This window function calculates the average of the age column. Because the OVER() clause is empty, the window is the entire table. The average (30+25+35)/3 = 30 is calculated and made available on every row.
- 2. **COALESCE(age, ...)**: This function returns the first non-NULL value it finds.
 - o If age is NOT NULL (e.g., for user 1), it returns the original age (30).
 - If age IS NULL (e.g., for user 3), it returns the second argument, which is the overall average age (30).

3.

Approach 2: UPDATE statement (to modify the table in place)

This is a destructive operation.

```
Generated sql
-- Using a CTE for clarity
WITH Stats AS (
SELECT AVG(age) as mean_age FROM users
)
UPDATE users
SET age = (SELECT mean_age FROM Stats)
WHERE age IS NULL;
IGNORE_WHEN_COPYING_START
content_copy download
Use code with caution. SQL
IGNORE_WHEN_COPYING_END
```

Explanation: This approach first calculates the mean in a CTE (or subquery) and then uses that value in an UPDATE statement to fill in the NULLs. This requires two scans of the table (one to calculate the average, one to update).

Best Practices

- **Grouped Imputation**: Mean imputation is often more powerful when done per group. For example, imputing the age with the average age of users *in the same city*. This is easily done by adding a PARTITION BY clause: AVG(age) OVER (PARTITION BY city).
- Mean vs. Median: The mean is sensitive to outliers. If the data is skewed, imputing with the median can be a more robust choice. The median can be calculated using PERCENTILE_CONT(0.5) WITHIN GROUP (ORDER BY age) OVER ().

Question 12

Create a SQL query that normalizes a column (scales between 0 and 1).

Answer:

Theory

Min-Max normalization is a feature scaling technique that rescales a numeric feature to a fixed range, typically [0, 1]. This is important for machine learning algorithms that are sensitive to the scale of input features, such as neural networks, SVMs, or any model using gradient descent.

```
The formula for Min-Max normalization is:
normalized_value = (value - min_value) / (max_value - min_value)
```

In SQL, this can be implemented efficiently using window functions to find the min and max values of the column in a single pass.

Code Example

Let's normalize the monthly_spend column in a customers table.

```
Sample Data:
customers table:
| customer_id | monthly_spend |
|-----|
| 1 | 50 |
| 2 | 200 |
| 3 | 20 | (Min)
| 4 | 500 | (Max)
SQL Query:
Generated sql
  WITH MinMax AS (
  SELECT
    MIN(monthly spend) AS min spend,
    MAX(monthly_spend) AS max_spend
  FROM
    customers
)
SELECT
  customer id,
  monthly_spend,
  -- Apply the normalization formula
  (monthly_spend - min_spend) * 1.0 / NULLIF((max_spend - min_spend), 0) AS
normalized_spend
FROM
  customers, MinMax;
IGNORE WHEN COPYING START
content_copy download
Use code with caution. SQL
IGNORE_WHEN_COPYING_END
A more modern approach using window functions avoids the cross join:
Generated sql
  SELECT
  customer id,
```

```
monthly_spend,
(monthly_spend - MIN(monthly_spend) OVER ()) * 1.0 /
NULLIF((MAX(monthly_spend) OVER () - MIN(monthly_spend) OVER ()), 0) AS
normalized_spend
FROM
customers;

IGNORE_WHEN_COPYING_START
content_copy download
Use code with caution. SQL
IGNORE WHEN COPYING END
```

Explanation

- 1. MIN(monthly_spend) OVER () and MAX(monthly_spend) OVER (): These window functions calculate the minimum (20) and maximum (500) values of the monthly_spend column across the entire table. The empty OVER() clause specifies the entire dataset as the window. These global min/max values are then available on every single row.
- 2. **(value min)** / **(max min)**: The query then applies the normalization formula directly for each row.
- 3. * **1.0**: This ensures floating-point division.
- 4. **NULLIF((max min), 0)**: This is a crucial safety check. If all values in the column are the same, max min will be 0, causing a division-by-zero error. NULLIF returns NULL if the two arguments are equal, which makes the result of the division NULL instead of an error.

Use Cases

- Preprocessing for ML: Essential for distance-based algorithms (like KNN), gradient-based optimization (Linear Regression, Neural Networks), and SVMs.
- Data Visualization: Scaling features to a common range can make them comparable on a single chart.

Question 13

Generate a feature that is a count over a rolling time window using SQL.

Answer:

Theory

A rolling count is a time-series feature that counts the number of events that occurred within a specific, preceding time window relative to each data point. For example, "how many

transactions did this user make in the last 7 days?". This type of feature is extremely valuable for capturing recent activity and momentum. It is implemented in SQL using the COUNT window function with a RANGE or ROWS based frame clause.

Code Example

Let's create a feature for each user transaction that counts how many transactions that same user made in the preceding 30 days.

Sample Data:

```
transactions table:

| transaction_id | user_id | transaction_time |

|------|

| 1 | 101 | 2023-10-15 10:00:00 |

| 2 | 101 | 2023-11-01 11:00:00 |

| 3 | 102 | 2023-11-05 12:00:00 |

| 4 | 101 | 2023-11-20 13:00:00 |
```

SQL Query (using PostgreSQL/Oracle interval syntax):

```
Generated sql
  SELECT
  transaction id,
  user_id,
  transaction time,
  COUNT(transaction id) OVER (
    PARTITION BY user_id
    ORDER BY transaction time
    RANGE BETWEEN INTERVAL '30' DAY PRECEDING AND CURRENT ROW
  ) - 1 AS count_transactions_last_30d
FROM
  transactions;
IGNORE WHEN COPYING START
content_copy download
Use code with caution. SQL
IGNORE_WHEN_COPYING_END
```

Note: The syntax for date/time intervals can vary. In SQL Server, you might use ROWS UNBOUNDED PRECEDING with a filter on the date, or another method.

Explanation

- COUNT(transaction_id) OVER (...): This specifies a windowed count.
- PARTITION BY user_id: Ensures the count is performed independently for each user.

- **ORDER BY transaction_time**: Orders the data chronologically, which is essential for any time-window calculation.
- RANGE BETWEEN INTERVAL '30' DAY PRECEDING AND CURRENT ROW: This is
 the frame clause that defines the rolling window. For each row, it includes all other rows
 for that user from the 30 days leading up to and including the current row's
 transaction_time. RANGE is generally preferred over ROWS for time-based windows as
 it correctly handles irregular time gaps.
- -1: The count includes the current row itself. To get the count of *previous* events within the window, we subtract 1.

Use Cases

- **Fraud Detection**: A sudden spike in the number of transactions in the last hour can be a strong fraudulent signal.
- **User Engagement**: Tracking the number of logins or actions in the last 7 days as a feature for a churn model.
- Credit Risk: Counting the number of late payments in the last year.

Question 14

Code an SQL function that categorizes continuous variables into bins.

Answer:

Theory

A User-Defined Function (UDF) is a reusable code block that can be called from within SQL queries. Creating a UDF to handle binning (discretization) is a great way to encapsulate complex business logic and ensure consistency across multiple queries and reports. The function takes a continuous value as input and returns a categorical string label representing the bin.

Code Example

Let's create a SQL function called BinUserAge that categorizes a user's age into predefined groups.

SQL Function Definition (PostgreSQL syntax):

Generated sql
CREATE OR REPLACE FUNCTION BinUserAge(age INT)
RETURNS VARCHAR AS \$\$
BEGIN

```
RETURN CASE
WHEN age IS NULL THEN 'Unknown'
WHEN age BETWEEN 18 AND 24 THEN '18-24'
WHEN age BETWEEN 25 AND 34 THEN '25-34'
WHEN age BETWEEN 35 AND 49 THEN '35-49'
WHEN age >= 50 THEN '50+'
ELSE 'Under 18'
END;
END;
END;
S$ LANGUAGE plpgsql;

IGNORE_WHEN_COPYING_START
content_copy download
Use code with caution. SQL
IGNORE_WHEN_COPYING_END
```

How to Use the Function in a Query:

```
Generated sql
SELECT
user_id,
age,
BinUserAge(age) AS age_group
FROM
users;

IGNORE_WHEN_COPYING_START
content_copy download
Use code with caution. SQL
IGNORE WHEN COPYING END
```

Explanation

- 1. **CREATE OR REPLACE FUNCTION BinUserAge(age INT)**: This defines a function named BinUserAge that accepts one integer argument, age.
- 2. **RETURNS VARCHAR**: Specifies that the function will return a string (the bin label).
- 3. **\$\$... \$\$ LANGUAGE plpgsql**: This is the standard syntax in PostgreSQL for defining the function body. The language is plpgsql.
- 4. **BEGIN** ... **END**: This block contains the function's logic.
- 5. **RETURN CASE** ... **END**: A standard CASE statement implements the binning logic. It evaluates the input age and returns the corresponding string. It's crucial to handle NULL values and edge cases gracefully.

Use Cases & Trade-offs

- Reusability & Consistency: The primary benefit. The same binning logic can be applied in dozens of queries without repeating the CASE statement, preventing errors and inconsistencies.
- **Readability**: Queries become cleaner. SELECT BinUserAge(age) is much more readable than a long, nested CASE statement.
- Performance Considerations: In some database systems, calling a UDF for every row
 can be slower than using an inline CASE statement because it may prevent the query
 optimizer from creating the most efficient execution plan. For very large-scale data
 processing, this trade-off between readability and performance should be considered.

Question 15

Implement a SQL solution to compute the TF-IDF score for text data.

Answer:

Theory

TF-IDF (Term Frequency-Inverse Document Frequency) is a numerical statistic that reflects how important a word is to a document in a collection or corpus. It's a very advanced task to implement in pure SQL, but it's a great demonstration of complex data manipulation using CTEs.

The calculation is broken down into three parts:

- 1. **Term Frequency (TF)**: How often a word appears in a document, normalized by the total number of words in that document. TF = (count of term t in doc d) / (total words in doc d)
- 2. **Inverse Document Frequency (IDF)**: A measure of how rare a word is across all documents. IDF = log(total number of docs / number of docs containing term t)
- 3. **TF-IDF**: The product of the two. TF-IDF = TF * IDF

Code Example

This is a complex query using multiple CTEs. We'll start with a table of tokenized words. (In a real-world scenario, the tokenization step itself would likely be done outside of SQL).

Sample Data:

```
doc_words table (pre-tokenized):
| doc_id | word |
|-----|
| 1 | sql |
| 1 | is |
| 1 | cool |
```

```
|2|sql|
| 2 | is |
|2|powerful|
SQL Query:
Generated sql
  WITH
-- Step 1: Calculate word counts per document
WordCounts AS (
  SELECT doc_id, word, COUNT(*) AS term_count
  FROM doc words
  GROUP BY doc id, word
),
-- Step 2: Calculate total words per document
DocStats AS (
  SELECT doc_id, COUNT(*) AS total_words
  FROM doc words
  GROUP BY doc_id
-- Step 3: Calculate Term Frequency (TF)
TermFrequency AS (
  SELECT
    wc.doc id,
    wc.word,
    wc.term_count * 1.0 / ds.total_words AS tf
  FROM WordCounts wc
  JOIN DocStats ds ON wc.doc_id = ds.doc_id
-- Step 4: Calculate Inverse Document Frequency (IDF)
InverseDocFrequency AS (
  SELECT
    word,
    LOG( (SELECT COUNT(DISTINCT doc_id) FROM doc_words) * 1.0 / COUNT(DISTINCT
doc_id) ) AS idf
  FROM doc_words
  GROUP BY word
-- Final Step: Calculate TF-IDF
SELECT
  tf.doc id,
  tf.word,
  tf.tf.
  idf.idf,
  tf.tf * idf.idf AS tf_idf
```

FROM TermFrequency tf

JOIN InverseDocFrequency idf ON tf.word = idf.word

ORDER BY doc_id, tf_idf DESC;

IGNORE_WHEN_COPYING_START content_copy download
Use code with caution. SQL
IGNORE_WHEN_COPYING_END

Explanation

- 1. WordCounts CTE: Counts how many times each word appears in each document.
- 2. DocStats CTE: Counts the total number of words in each document.
- 3. **TermFrequency CTE**: Joins the first two CTEs to calculate TF using the formula term count / total words.
- 4. **InverseDocFrequency CTE**: This is the most complex part.
 - GROUP BY word ensures we calculate IDF for each unique word.
 - The LOG() function calculates the IDF.
 - (SELECT COUNT(DISTINCT doc_id) FROM doc_words) is a scalar subquery to get the total number of documents in the corpus.
 - COUNT(DISTINCT doc_id) within the main query counts how many documents each specific word appears in.

5.

6. **Final SELECT**: Joins the TermFrequency and InverseDocFrequency CTEs on word to get both scores, then multiplies them to get the final tf idf.

Best Practices and Optimization

- **Materialization**: This query is extremely computationally intensive. It should not be run on the fly. The standard practice is to run it as a batch job and materialize the results into a dedicated feature_tfidf table for later use.
- **Tokenization**: The initial step of breaking text into words (tokenization) and handling things like stop words and stemming is usually better handled in a programming language like Python before the data is loaded into the doc_words table.

SQL ML Interview Questions - Coding Questions

Question 1

Write a SQL query that joins two tables and retrieves only the rows with matching keys.

Answer:

Theory

To retrieve only the rows with matching keys from two tables, you use an INNER JOIN. This is the most common type of join. It creates a new result table by combining column values of two tables (table1 and table2) based upon the join-predicate. The query scans table1 and for each row it finds, it scans table2 for all matching rows based on the ON condition. This effectively filters out any rows from either table that do not have a corresponding match in the other.

Code Example

Scenario: We have a customers table and an orders table. We want a list of customers who have placed orders, along with their order details.

Sample Data:

SQL Query:

```
Generated sql
SELECT
c.customer_id,
c.name,
o.order_id,
o.amount
FROM
customers AS c
INNER JOIN
orders AS o ON c.customer_id = o.customer_id;
```

Result:

```
| customer_id | name | order_id | amount |
|------|
| 1 | Alice | 101 | 50 |
| 1 | Alice | 103 | 120 |
| 2 | Bob | 102 | 75 |
```

Explanation

- 1. **SELECT c.customer_id, c.name, o.order_id, o.amount**: This specifies the columns we want to see in the final result. Using aliases (c. and o.) is a best practice to specify which table each column comes from.
- 2. **FROM customers AS c**: This defines the first (left) table and gives it a short alias c for readability.
- INNER JOIN orders AS o: This specifies that we are joining with the orders table (aliased as o) and that we only want rows that have a match in both tables. JOIN is shorthand for INNER JOIN.
- 4. **ON c.customer_id = o.customer_id**: This is the join condition. It tells the database to match rows where the customer_id value is the same in both the customers and orders tables. The row for order 104 is excluded because its customer_id of 4 does not exist in the customers table. Similarly, 'Charlie' is excluded because he has no orders.

Best Practices

- **Use Aliases**: Always use table aliases (like c and o) to make your query shorter and more readable, especially when joining multiple tables.
- Index Join Keys: For optimal performance, the columns used in the ON clause (customer_id in this case) should be indexed. An index on the foreign key column (orders.customer_id) is particularly important.

Question 2

Create a SQL query to pivot a table transforming rows into columns.

Answer:

Theory

Pivoting is a data transformation technique that rotates data from a "long" format (multiple rows for different attributes) to a "wide" format (attributes as distinct columns). It's commonly used to create summary reports or to prepare a dataset for machine learning models where each distinct category value becomes a new feature column.

While some SQL dialects have a specific PIVOT operator (like SQL Server or Oracle), the most portable and universal method is to use a combination of an aggregate function (like SUM or MAX) and CASE statements, grouped by the desired row identifier.

Code Example

Scenario: We have a product sales table recording sales amount by product category for different years. We want to transform this into a report where each year is a column.

Sample Data:

```
product sales table:
| category | sale_year | sale_amount |
|-----|
| Electronics | 2022 | 5000 |
| Books | 2022 | 1500 |
| Electronics | 2023 | 6200 |
| Books | 2023 | 1800 |
```

SQL Query (Standard SQL):

```
Generated sql
  SELECT
  category,
  SUM(CASE WHEN sale_year = 2022 THEN sale_amount ELSE 0 END) AS sales_2022,
  SUM(CASE WHEN sale year = 2023 THEN sale amount ELSE 0 END) AS sales 2023
FROM
  product_sales
GROUP BY
  category;
IGNORE WHEN COPYING START
content_copy download
Use code with caution. SQL
IGNORE WHEN COPYING END
```

Result:

```
| category | sales_2022 | sales_2023 |
|-----|
| Books | 1500 | 1800 |
| Electronics | 5000 | 6200 |
```

Explanation

1. **GROUP BY category**: This clause groups all rows for the same category (Electronics, Books) together. The final result will have one row for each unique category.

- 2. **SUM(CASE WHEN ...)**: This is the core of the pivot logic. It's an aggregate function that works on each group defined by GROUP BY.
 - CASE WHEN sale_year = 2022 THEN sale_amount ELSE 0 END: For the sales_2022 column, this expression checks each row within a group. If the year is 2022, it takes the sale amount; otherwise, it takes 0.
 - SUM(...): The SUM function then adds up these values. For the 'Electronics' group, it calculates SUM(5000, 0) which equals 5000. For the 'Books' group, it's SUM(1500, 0).
 - The same logic applies to the sales_2023 column.

3.

Pitfalls

• Static Columns: This approach requires you to know the values you want to pivot into columns (e.g., 2022, 2023) ahead of time. If a new year 2024 is added to the data, you must manually add a new SUM(CASE ...) expression to the query. This is a limitation of static SQL. Dynamic SQL is required to handle an unknown number of pivot columns.

Question 3

Write a SQL query to calculate moving averages.

Answer:

Theory

A moving average (or rolling average) is a common time-series analysis technique used to smooth out short-term fluctuations and highlight longer-term trends or cycles. It is calculated by averaging a subset of data points over a specific window of time. In SQL, this is achieved efficiently and elegantly using window functions, specifically AVG() combined with an ORDER BY and a frame clause (ROWS or RANGE).

Code Example

Scenario: Calculate a 7-day moving average of daily revenue for a specific product.

Sample Data:

```
daily_revenue table:
| report_date | product_id | revenue |
|------|-----|------|
| 2023-11-01 | 101 | 150 |
| 2023-11-02 | 101 | 160 |
| ... | ... | ... |
```

```
| 2023-11-07 | 101 | 180 |
| 2023-11-08 | 101 | 190 |
```

SQL Query:

```
Generated sql
  SELECT
  report date,
  product id,
  revenue.
  AVG(revenue) OVER (
    PARTITION BY product id
    ORDER BY report date
    ROWS BETWEEN 6 PRECEDING AND CURRENT ROW
  ) AS "7 day moving avg"
FROM
  daily_revenue;
IGNORE_WHEN_COPYING_START
content copy download
Use code with caution. SQL
IGNORE_WHEN_COPYING_END
```

Explanation

- AVG(revenue) OVER (...): This specifies that we want to compute an average using a window function.
- **PARTITION BY product_id**: This clause divides the data into separate groups for each product_id. The moving average calculation will restart for each product, which is crucial if you have data for multiple products in the same table.
- ORDER BY report_date: Within each partition, the rows are ordered by report_date.
 This is essential for any time-series calculation to ensure the window is applied chronologically.
- ROWS BETWEEN 6 PRECEDING AND CURRENT ROW: This is the frame clause that defines the "window" for the average. It tells the database to include the current row and the six rows that came before it (totaling 7 rows) in the calculation.

Use Cases

- Feature Engineering: Create smoothed features for time-series forecasting models to reduce noise.
- Financial Analysis: Track trends in stock prices or company revenues.
- Business Intelligence: Create smoother trend lines in performance dashboards.
- Anomaly Detection: Compare the current value to its moving average to spot unusual deviations.

Question 4

How can you create lagged features in SQL?

Answer:

Theory

A lagged feature is a feature that contains a value from a previous time step. For a given observation at time t, a lagged feature would be the value of a variable at time t-n, where n is the lag period. These features are fundamental for sequential and time-series models as they capture temporal dependencies (e.g., momentum, seasonality). The LAG() window function in SQL is designed specifically for this purpose.

Code Example

Scenario: Create a feature that shows each customer's previous purchase amount.

Sample Data:

```
purchases table:

| customer_id | purchase_date | amount |

|-----|

| 1 | 2023-01-10 | 50 |

| 1 | 2023-02-15 | 75 |

| 2 | 2023-03-05 | 100 |

| 1 | 2023-04-01 | 60 |
```

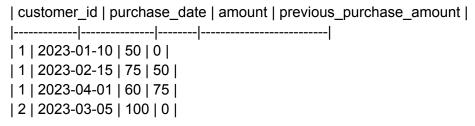
SQL Query:

```
Generated sql
SELECT
customer_id,
purchase_date,
amount,
LAG(amount, 1, 0) OVER (
PARTITION BY customer_id
ORDER BY purchase_date
) AS previous_purchase_amount
FROM
purchases;

IGNORE_WHEN_COPYING_START
content_copy download
```

Use code <u>with caution</u>. SQL IGNORE WHEN COPYING END

Result:



Explanation

- LAG(amount, 1, 0): This is the core function call.
 - o amount: The column whose previous value we want to retrieve.
 - 1: The lag offset. 1 means we want the value from the immediately preceding row. A value of 2 would get the value from two rows ago.
 - 0: The default value to use if there is no preceding row (i.e., for the first purchase of each customer). Using a default avoids NULLs, which is often desirable for ML models.

•

- **OVER (...)**: This defines the window for the LAG function.
 - PARTITION BY customer_id: This is critical. It ensures that we only look at previous purchases for the same customer. The calculation is isolated within each customer's data.
 - ORDER BY purchase_date: This orders each customer's purchases chronologically, defining what "previous" means.

•

Use Cases

- **Predictive Modeling**: Create features like "time since last event," "previous page visited," or "last month's sales."
- Change Detection: Calculate the difference between the current value and the previous value (amount LAG(amount, 1, amount) OVER (...)) to detect significant changes.

Question 5

Describe how to compute a ratio feature within groups using SQL.

Answer:

Theory

Calculating a ratio within a group helps normalize a value against its local context, often creating a more powerful predictive feature than the raw value alone. For example, an employee's salary is a number, but their salary as a percentage of their department's total salary indicates their relative importance within that department. This type of feature is easily created using SQL window functions, which can compute a group-level aggregate and make it available on every row of that group, allowing for row-level calculations with group-level context.

Code Example

Scenario: Calculate a feature that represents each product's sales as a ratio of its category's total sales.

Sample Data:

```
product_sales table:
| product_name | category | sales |
|-----|
| Laptop A | Electronics | 5000 |
| Laptop B | Electronics | 3000 |
| Book X | Books | 1500 |
| Book Y | Books | 500 |
```

SQL Query:

```
Generated sql
SELECT
product_name,
category,
sales,
sales * 1.0 / SUM(sales) OVER (PARTITION BY category) AS ratio_of_category_sales
FROM
product_sales;

IGNORE_WHEN_COPYING_START
content_copy download
Use code with caution. SQL
IGNORE_WHEN_COPYING_END
```

Result:

Explanation

- 1. SUM(sales) OVER (PARTITION BY category): This is the key window function.
 - PARTITION BY category divides the data into partitions, one for 'Electronics' and one for 'Books'.
 - SUM(sales) is calculated for each partition independently. For the 'Electronics' partition, the sum is 5000 + 3000 = 8000. For the 'Books' partition, the sum is 1500 + 500 = 2000.
 - Crucially, this partition-level sum (8000 or 2000) is appended to every row within its respective partition.

2.

- 3. **sales** * **1.0** / ...: In the outer SELECT, we take the individual row's sales value and divide it by the partition-level sum we just calculated.
 - Multiplying by 1.0 is a common trick to force floating-point division to get a decimal result, avoiding integer division which would result in 0 in many cases.

4.

Pitfalls and Optimization

- Division by Zero: If a group's total sum could be zero, the query would fail. To prevent this, wrap the denominator in NULLIF: sales * 1.0 / NULLIF(SUM(sales) OVER (PARTITION BY category), 0).
- Performance: On databases that support it, defining an index on the PARTITION BY key (category in this case) can improve performance.

Question 6

Write a SQL query that identifies and removes duplicate records from a dataset.

Answer:

Theory

Identifying and removing duplicate records is a common data cleaning task. A record is considered a duplicate if another row exists with the same values across a specific set of columns (the "uniqueness key"). The most robust and flexible way to handle this in SQL is to use the ROW_NUMBER() window function within a Common Table Expression (CTE). This allows you to assign a unique rank to each record within a duplicate group, and then selectively identify or delete all but one of them.

Code Example

Scenario: We have a contacts table where there are duplicate entries for the same email address. We want to keep only the most recently added record for each email.

Sample Data:

```
contacts table:
| contact_id | email | created_at |
|------|
| 1 | alice@example.com | 2023-11-01 10:00:00 |
| 2 | bob@example.com | 2023-11-02 11:00:00 |
| 3 | alice@example.com | 2023-11-03 12:00:00 | (Duplicate of contact 1)
| 4 | carol@example.com | 2023-11-04 13:00:00 |
```

Multiple Solution Approaches

Approach 1: Identify Duplicates (Non-destructive)

First, it's safest to write a query that just identifies the duplicates.

```
Generated sql
  WITH NumberedContacts AS (
  SELECT
    contact id,
    email,
    created at,
    ROW NUMBER() OVER(
      PARTITION BY email
      ORDER BY created at DESC
    ) AS rn
  FROM
    contacts
SELECT * FROM NumberedContacts WHERE rn > 1;
IGNORE WHEN COPYING START
content copy download
Use code with caution. SQL
IGNORE_WHEN_COPYING_END
```

Approach 2: DELETE using a CTE (Recommended for Removal)

This query will remove the identified duplicates from the table.

```
Generated sql
WITH NumberedContacts AS (
SELECT
contact_id,
ROW_NUMBER() OVER(
```

```
PARTITION BY email
ORDER BY created_at DESC
) AS rn
FROM
contacts
)
DELETE FROM contacts
WHERE contact_id IN (
SELECT contact_id FROM NumberedContacts WHERE rn > 1
);
IGNORE_WHEN_COPYING_START
content_copy download
Use code with caution. SQL
IGNORE_WHEN_COPYING_END
```

Explanation

- 1. WITH NumberedContacts AS (...): We define a CTE to first identify the duplicates.
- 2. ROW_NUMBER() OVER(...):
 - PARTITION BY email: This is the uniqueness key. It groups all rows with the same email together.
 - ORDER BY created_at DESC: Within each group, it orders the records by their creation date in descending order (newest first). This determines which record we keep.
 - ROW_NUMBER() then assigns a sequential number (rn) to each row in the group. The newest record gets rn = 1, the second newest gets rn = 2, and so on.
 Any row with rn > 1 is a duplicate we want to remove.

3.

- 4. **DELETE FROM contacts** ...: The outer query deletes rows from the original contacts table.
- 5. **WHERE contact_id IN (...)**: The WHERE clause specifies which rows to delete by selecting the contact_ids from our CTE where the row number (rn) is greater than 1.

Best Practices

- Always Backup: Before running a destructive DELETE operation on a production table, always create a backup or run the select query first to verify what will be deleted.
- **Define Duplicates Clearly**: The PARTITION BY clause is critical. It defines what combination of columns makes a record a duplicate.
- **Define Which Record to Keep**: The ORDER BY clause inside OVER() is equally important. It determines which of the duplicate records will be assigned rn = 1 and therefore be kept.

Question 7

In SQL, how would you format strings or concatenate columns for text-based Machine Learning features?

Answer:

Theory

Concatenating (joining) multiple string columns or formatting them into a single text feature is a common preprocessing step for text-based ML models. For example, combining a product's title, description, and category into one text block can provide a richer input for NLP models. SQL provides standard functions like CONCAT or operators like || for this purpose. It is also crucial to handle NULL values gracefully to avoid losing information.

Code Example

Scenario: Create a combined text feature from a product's name, brand, and description to be used in an NLP model.

Sample Data:

```
products table:
| product id | name | brand | description |
|-----|
| 101 | Smart Watch | TechCo| A watch that tracks fitness. |
| 102 | Running Shoe| FitGear| NULL |
```

SQL Query:

```
Generated sal
  SELECT
  product id,
  -- Using CONCAT WS to handle separators and NULLs gracefully
  CONCAT_WS(' | ',
    'Name: ' || name,
    'Brand: ' | brand,
    'Desc: ' || COALESCE(description, 'N/A')
  ) AS combined feature
FROM
  products;
IGNORE_WHEN_COPYING_START
content copy download
Use code with caution. SQL
IGNORE_WHEN_COPYING_END
```

Note: CONCAT_WS (Concatenate With Separator) is available in PostgreSQL, MySQL, and modern SQL Server. For other databases, you might chain CONCAT or || operators.

Result:

```
| product_id | combined_feature | | |
|---|---|---|---|
| 101 | Name: Smart Watch | Brand: TechCo | Desc: A watch that tracks fitness. |
| 102 | Name: Running Shoe | Brand: FitGear | Desc: N/A |
```

Explanation

- || Operator: A standard way to concatenate strings (in PostgreSQL/Oracle; + in SQL Server). 'Name: ' || name results in 'Name: Smart Watch'.
- **COALESCE**(**description**, 'N/A'): This function is crucial for handling NULL values. It returns the first non-NULL value in its argument list. If description is NULL, it returns the string 'N/A', preventing the entire concatenated string from becoming NULL.
- **CONCAT_WS(' | ', ...)**: This function joins multiple strings together using a specified separator (in this case, ' | '). It is often preferred because it automatically ignores NULL arguments and makes the query cleaner than chaining many || operators.

Use Cases

- **NLP Feature Engineering**: Creating a single document from multiple text fields before feeding it into a TF-IDF vectorizer or a language model.
- **Categorical Features**: Combining multiple low-cardinality categorical features (e.g., country and device type) into a single interaction feature ('US-Mobile').
- **Search**: Creating a searchable text field that aggregates information from multiple columns.

Question 8

Write a SQL stored procedure that calls a Machine Learning scoring function.

Answer:

Theory

A stored procedure is a pre-compiled set of one or more SQL statements that are stored in the database. In the context of MLOps, a stored procedure is an excellent way to operationalize a machine learning model. It encapsulates the complex logic of feature retrieval and model invocation into a single, reusable database object. This creates a secure and performant API for applications to get predictions without needing to know the underlying data schema or feature engineering logic.

The exact syntax for calling an ML model varies significantly across database platforms (e.g., SQL Server's sp_execute_external_script, Snowflake's external functions, BigQuery ML's ML.PREDICT). The following example uses a generic, conceptual syntax to illustrate the universal logic.

Code Example (Conceptual - T-SQL/SQL Server like)

Scenario: A procedure takes a customer_id, retrieves the latest features for that customer, and calls a hypothetical prediction function to score their churn risk.

```
Generated sql
  CREATE OR ALTER PROCEDURE dbo.GetChurnPrediction
  @CustomerID INT,
  @PredictedChurnProba FLOAT OUTPUT
AS
BEGIN
  -- Enforce security and prevent SQL injection side effects
  SET NOCOUNT ON;
  -- Step 1: Declare variables to hold the features
  DECLARE @TenureMonths INT;
  DECLARE @MonthlySpend DECIMAL(10, 2);
  DECLARE @HasActiveSupportContract BIT;
  -- Step 2: Retrieve the latest features for the given customer
  -- This query might be complex in a real scenario
  SELECT TOP 1
    @TenureMonths = DATEDIFF(month, signup_date, GETDATE()),
    @MonthlySpend = avg_monthly_spend,
    @HasActiveSupportContract = has_active_contract
  FROM
    customer_features
  WHERE
    customer id = @CustomerID;
  -- Step 3: Call the hypothetical machine learning scoring function
  -- The actual call syntax is database-specific.
  -- This example assumes a native PREDICT function.
  SELECT @PredictedChurnProba = probability
  FROM PREDICT(MODEL = 'ChurnModel_v2',
         DATA = (SELECT @TenureMonths as tenure,
                  @MonthlySpend as spend,
                  @HasActiveSupportContract as support_contract)
         );
END:
```

-- How to call the procedure from an application
DECLARE @ChurnScore FLOAT;
EXEC dbo.GetChurnPrediction @CustomerID = 12345, @PredictedChurnProba = @ChurnScore OUTPUT;
SELECT @ChurnScore AS ChurnProbability;

IGNORE_WHEN_COPYING_START content_copy download
Use code with caution. SQL
IGNORE_WHEN_COPYING_END

Explanation

- 1. **CREATE OR ALTER PROCEDURE** ...: Defines the procedure name (GetChurnPrediction) and its parameters. @CustomerID is an input, and @PredictedChurnProba is an OUTPUT parameter to return the result.
- Feature Retrieval: A standard SELECT query fetches the necessary features for the model from a pre-built customer_features table. This logic is now centralized and hidden from the client application.
- 3. **Model Scoring**: The PREDICT(...) line is a conceptual representation of calling the model. It passes the retrieved features as input and captures the model's output score in the @PredictedChurnProba variable.
- 4. **OUTPUT Parameter**: The procedure returns the prediction score via the OUTPUT parameter, which the calling application can then use.

Best Practices

- Security: Use stored procedures to grant limited, indirect access to data. Applications
 can EXECUTE the procedure without having direct SELECT permissions on the
 underlying feature tables.
- **Performance**: Feature retrieval queries within the procedure must be highly optimized, typically using an index on the lookup key (customer_id).
- Modularity: Encapsulating the scoring logic makes model updates easier. You can
 update the procedure to call a new model version (ChurnModel_v3) without changing
 any client application code.

Question 9

How would you construct a complex SQL query to extract time series features for a Machine Learning model?

Answer:

Theory

Extracting a rich set of time series features often requires multiple stages of transformation and aggregation. A complex query for this task is best constructed using Common Table Expressions (CTEs) to break the logic into clean, readable, and sequential steps. This modular approach allows you to build layers of features, where each CTE uses the results of the previous one. We can generate features like rolling averages, lagged values, and period-over-period growth rates in a single, powerful query.

Code Example

Scenario: We need to build a feature set for a sales forecasting model. For each product and each month, we want to calculate:

- 1. Total monthly sales.
- Sales from the previous month (lagged feature).
- 3. A 3-month rolling average of sales.
- 4. Month-over-month sales growth percentage.

```
Generated sql
  WITH MonthlySales AS (
  -- Step 1: Aggregate daily sales into monthly totals for each product.
  -- This sets the grain of our time series.
  SELECT
    product id,
    DATE TRUNC('month', sale date)::date AS sale month,
    SUM(sale amount) AS total monthly sales
  FROM
    sales_data
  GROUP BY
    product id, sale month
WindowFeatures AS (
  -- Step 2: Use window functions to create time-series features.
  SELECT
    product id,
    sale month,
    total_monthly_sales,
    -- Feature: Sales from the previous month
    LAG(total_monthly_sales, 1, 0) OVER (PARTITION BY product_id ORDER BY
sale_month) AS previous_month_sales,
    -- Feature: 3-month rolling average
    AVG(total_monthly_sales) OVER (
```

```
PARTITION BY product id ORDER BY sale month
      ROWS BETWEEN 2 PRECEDING AND CURRENT ROW
    ) AS "3 month rolling avg"
  FROM
    MonthlySales
-- Step 3: Calculate final derived features like growth rate.
SELECT
  product id,
  sale month,
  total monthly sales,
  previous month sales,
  "3 month rolling avg",
  -- Feature: Month-over-month growth rate
  (total monthly sales - previous month sales) * 1.0 / NULLIF(previous month sales, 0) AS
"mom growth rate"
FROM
  WindowFeatures
ORDER BY
  product id, sale month;
IGNORE_WHEN_COPYING_START
content copy download
Use code with caution. SQL
IGNORE_WHEN_COPYING_END
```

Explanation

- MonthlySales CTE: The first step is to roll up the raw, daily data into the desired time grain (monthly). This GROUP BY simplifies the dataset for the next stage.
 DATE TRUNC is a standard function for this.
- 2. **WindowFeatures CTE**: This is where the core feature engineering happens. It takes the monthly sales data and applies window functions:
 - LAG() calculates the previous month sales.
 - o AVG() with a ROWS BETWEEN frame calculates the 3_month_rolling_avg.
 - All window functions are partitioned by product_id to ensure calculations are done independently for each product.

3.

4. **Final SELECT**: The final query uses the features generated in the previous CTE to calculate a derived feature: the month-over-month growth rate. Using NULLIF prevents division-by-zero errors.

Best Practices

- Modularity with CTEs: Breaking down the logic into CTEs makes the query vastly easier to debug and maintain. You can test each CTE independently by running SELECT * FROM <CTE name>.
- **Performance**: Ensure the source table is indexed on sale_date and product_id to speed up the initial aggregation and partitioning.
- **Data Granularity**: Always start by aggregating your data to the correct time series interval (e.g., daily, weekly, monthly) before applying window functions.

Question 10

Discuss ways to implement regular expressions in SQL for natural language processing tasks.

Answer:

Theory

Regular expressions (regex) are a powerful tool for pattern matching in text, making them very useful for cleaning and feature engineering in NLP tasks directly within SQL. Most modern SQL dialects provide functions to leverage regex. This allows for sophisticated text manipulation without having to extract the data to an external environment. The common functions are:

- Matching/Searching (REGEXP_LIKE, ~, RLIKE): Checks if a string contains a pattern. Returns a boolean.
- Extracting (REGEXP_SUBSTR, REGEXP_MATCHES): Extracts the part of a string that matches a pattern.
- Replacing (REGEXP_REPLACE): Replaces parts of a string that match a pattern with another string.

Code Examples (using PostgreSQL syntax for richness)

Use Case 1: Feature Engineering - Check for a specific keyword

Create a binary feature indicating if a product review mentions "delivery" or "shipping".

```
Generated sql
SELECT
review_text,
(review_text ~* '(delivery|shipping)') AS mentions_logistics
-- ~* for case-insensitive match. | acts as OR.
FROM product_reviews;

IGNORE_WHEN_COPYING_START
content_copy download
```

```
Use code with caution. SQL IGNORE WHEN COPYING END
```

Use Case 2: Data Cleaning - Remove URLs and non-alphanumeric characters Clean text data before tokenization.

```
Generated sql
SELECT
review_text,
REGEXP_REPLACE(
REGEXP_REPLACE(review_text, 'https?://\S+', ", 'g'), -- Remove URLs
'[^a-zA-Z0-9\s]', ", 'g' -- Remove non-alphanumeric chars except space
) AS cleaned_text
FROM product_reviews;

IGNORE_WHEN_COPYING_START
content_copy download
Use code with caution. SQL
IGNORE_WHEN_COPYING_END
```

Use Case 3: Information Extraction - Extract hashtags from tweets

Extract structured information from unstructured text.

```
Generated sql
SELECT
tweet_text,
REGEXP_MATCHES(tweet_text, '#(\w+)', 'g') AS hashtags
-- \w+ matches one or more word characters. 'g' finds all matches.
FROM tweets;

IGNORE_WHEN_COPYING_START
content_copy download
Use code with caution. SQL
IGNORE_WHEN_COPYING_END
```

Explanation

- **~* Operator**: A PostgreSQL shorthand for case-insensitive regex matching. The query (review_text ~* '(delivery|shipping)') returns true if any of the keywords are found.
- REGEXP_REPLACE(string, pattern, replacement, flags): Finds all substrings matching pattern and replaces them. The 'g' (global) flag is crucial for replacing all instances, not just the first one. Here, we nest them to perform multiple cleaning steps.
- **REGEXP_MATCHES(string, pattern, 'g')**: Returns a set of all matching substrings. This is powerful for extracting all occurrences of a pattern, like all hashtags in a tweet.

Debugging and Performance

- **Complexity**: Regex can be hard to write and debug. It's best to test patterns on a small sample of data or using an online regex tester before running them on a large table.
- Performance: Regex operations, especially on long text columns, can be computationally intensive and may prevent the use of standard indexes. If a regex-based feature is needed frequently, it's often a best practice to pre-compute it and store it in a dedicated column rather than calculating it on the fly in every query.

Question 11

Write a SQL script to identify and replace missing values with the column mean.

Answer:

Theory

Imputing missing values (NULLs) is a standard data preprocessing step. Replacing them with the column's mean is a common and simple strategy for numerical data. The most efficient way to do this in SQL without multiple queries is to use window functions. A window function can calculate the mean of the entire column (or a group) and make that value available on every row, allowing you to identify and replace NULLs in a single pass.

Code Example

Scenario: Impute missing age values in a users table with the average age of all users.

Sample Data:

```
users table:

| user_id | age | city |

|------|

| 1 | 30 | NYC |

| 2 | 25 | London |

| 3 | NULL| NYC |

| 4 | 35 | London |

| 5 | NULL| NYC |

| 6 | 45 | London |
```

Multiple Solution Approaches

Approach 1: SELECT with Global Mean (for creating a new view/table)

This is non-destructive and ideal for feature generation.

```
Generated sql
SELECT
user_id,
COALESCE(
age,
AVG(CAST(age AS FLOAT)) OVER ()
) AS imputed_age
FROM
users;

IGNORE_WHEN_COPYING_START
content_copy download
Use code with caution. SQL
IGNORE_WHEN_COPYING_END
```

Approach 2: SELECT with Grouped Mean (More Advanced)

This is often more effective, as it imputes with a more relevant mean.

```
Generated sql
SELECT
user_id,
age,
city,
COALESCE(
age,
AVG(CAST(age AS FLOAT)) OVER (PARTITION BY city)
) AS imputed_age_by_city
FROM
users;

IGNORE_WHEN_COPYING_START
content_copy download
Use code with caution. SQL
IGNORE_WHEN_COPYING_END
```

Approach 3: UPDATE statement (to modify the table in place)

This is a destructive operation and should be used with caution.

```
Generated sql
    UPDATE users
SET age = sub.mean_age
FROM (SELECT AVG(age) as mean_age FROM users) AS sub
WHERE users.age IS NULL;
IGNORE_WHEN_COPYING_START
```

content_copy download
Use code with caution. SQL
IGNORE WHEN COPYING END

Explanation (for Approach 2)

- 1. AVG(CAST(age AS FLOAT)) OVER (PARTITION BY city): This window function calculates the average age.
 - PARTITION BY city divides the data into groups (NYC, London). The average is calculated independently for each city.
 - The average for that user's city is then made available on every row within that city's partition.

2.

- 3. **COALESCE(age, ...)**: This function returns the first non-NULL value it finds.
 - If age is NOT NULL, it returns the original age.
 - If age IS NULL, it returns the second argument, which is the average age for that user's city.

4.

Best Practices

- **Grouped Imputation**: As shown in Approach 2, mean imputation is often more powerful when done per group.
- Mean vs. Median: The mean is sensitive to outliers. If the data is skewed, imputing with
 the median can be a more robust choice. The median can be calculated using
 PERCENTILE_CONT(0.5) WITHIN GROUP (ORDER BY age) OVER (PARTITION BY
 ...) in most modern SQL dialects.

Question 12

Create a SQL guery that normalizes a column (scales between 0 and 1).

Answer:

Theory

Min-Max normalization is a feature scaling technique that rescales a numeric feature to a fixed range, typically [0, 1]. This is a critical preprocessing step for many machine learning algorithms that are sensitive to the scale of input features, such as neural networks, SVMs, and any distance-based or gradient-based model.

The formula for Min-Max normalization is: normalized value = (value - min value) / (max value - min value) In SQL, this can be implemented efficiently using window functions to find the minimum and maximum values of the column in a single pass without requiring joins or subqueries.

Code Example

Scenario: Normalize the monthly_spend column in a customers table.

```
Sample Data:
customers table:
| customer_id | monthly_spend |
|------|
| 1 | 50 |
| 2 | 200 |
| 3 | 20 | (Min)
| 4 | 500 | (Max)

SQL Query:

Generated sql
```

```
Generated sql
SELECT
customer_id,
monthly_spend,
(monthly_spend - MIN(monthly_spend) OVER ()) * 1.0 /
NULLIF((MAX(monthly_spend) OVER () - MIN(monthly_spend) OVER ()), 0) AS
normalized_spend
FROM
customers;
```

IGNORE_WHEN_COPYING_START content_copy download
Use code with caution. SQL
IGNORE_WHEN_COPYING_END

Result:

Explanation

1. MIN(monthly_spend) OVER () and MAX(monthly_spend) OVER (): These window functions calculate the minimum (20) and maximum (500) values of the monthly_spend

- column across the entire table. The empty OVER() clause specifies the entire dataset as the window. These global min/max values are then available on every single row.
- 2. **(value min)** / **(max min)**: The query then applies the normalization formula directly for each row.
- 3. * **1.0**: This ensures floating-point division, preventing integer division which might truncate the result.
- 4. **NULLIF((max min), 0)**: This is a crucial safety check. If all values in the column are the same, max min will be 0, causing a division-by-zero error. NULLIF returns NULL if its two arguments are equal, which makes the result of the division NULL instead of an error.

Use Cases

- Preprocessing for ML: Essential for distance-based algorithms (like KNN), gradient-based optimization (Linear Regression, Neural Networks), and SVMs.
- **Data Visualization**: Scaling features to a common range can make them comparable on a single chart.

Question 13

Generate a feature that is a count over a rolling time window using SQL.

Answer:

Theory

A rolling count is a time-series feature that counts the number of events that occurred within a specific, preceding time window relative to each data point. For example, "how many transactions did this user make in the last 7 days?". This type of feature is extremely valuable for capturing recent activity, frequency, and momentum. It is implemented in SQL using the COUNT window function with a RANGE or ROWS based frame clause that defines the lookback period.

Code Example

Scenario: For each user transaction, create a feature that counts how many transactions that same user made in the preceding 30 days.

Sample Data:

```
transactions table:
| transaction_id | user_id | transaction_time |
|------|-----|
| 1 | 101 | 2023-10-15 10:00:00 |
```

```
| 2 | 101 | 2023-11-01 11:00:00 |
| 3 | 102 | 2023-11-05 12:00:00 |
| 4 | 101 | 2023-11-20 13:00:00 |
```

SQL Query (using PostgreSQL/Oracle interval syntax):

```
Generated sql
  SELECT
  transaction id,
  user id,
  transaction_time,
  COUNT(transaction id) OVER (
    PARTITION BY user id
    ORDER BY transaction time
    RANGE BETWEEN INTERVAL '30' DAY PRECEDING AND CURRENT ROW
  ) - 1 AS count_transactions_last_30d
FROM
  transactions:
IGNORE WHEN COPYING START
content copy download
Use code with caution. SQL
IGNORE WHEN COPYING END
```

Note: The syntax for date/time intervals can vary. In SQL Server, you might use CAST(transaction_time as bigint) and RANGE BETWEEN 2592000 PRECEDING AND CURRENT ROW for seconds.

Explanation

- **COUNT(transaction_id) OVER (...)**: This specifies a windowed count.
- PARTITION BY user_id: Ensures the count is performed independently for each user.
- **ORDER BY transaction_time**: Orders the data chronologically, which is essential for any time-window calculation.
- RANGE BETWEEN INTERVAL '30' DAY PRECEDING AND CURRENT ROW: This is
 the frame clause that defines the rolling window. For each row, it includes all other rows
 for that user whose transaction_time is within the 30 days leading up to and including the
 current row's transaction_time. RANGE is generally preferred over ROWS for
 time-based windows as it correctly handles irregular time gaps.
- -1: The count includes the current row itself. To get the count of *previous* events within the window, we subtract 1.

Use Cases

- **Fraud Detection**: A sudden spike in the number of transactions in the last hour can be a strong fraudulent signal.
- **User Engagement**: Tracking the number of logins or actions in the last 7 days as a feature for a churn model.
- Credit Risk: Counting the number of late payments in the last year.

Question 14

Code an SQL function that categorizes continuous variables into bins.

Answer:

Theory

A User-Defined Function (UDF) is a reusable code block that can be called from within SQL queries. Creating a UDF to handle binning (discretization) is an excellent way to encapsulate complex business logic and ensure consistency across multiple queries and reports. The function takes a continuous value as input and returns a categorical string label representing the bin. This promotes code reuse and maintainability.

Code Example

Scenario: Create a SQL function called BinUserAge that categorizes a user's age into predefined groups.

SQL Function Definition (PostgreSQL syntax):

```
Generated sql
CREATE OR REPLACE FUNCTION BinUserAge(age INT)
RETURNS VARCHAR AS $$
BEGIN
RETURN CASE
WHEN age IS NULL THEN 'Unknown'
WHEN age BETWEEN 18 AND 24 THEN '18-24'
WHEN age BETWEEN 25 AND 34 THEN '25-34'
WHEN age BETWEEN 35 AND 49 THEN '35-49'
WHEN age >= 50 THEN '50+'
ELSE 'Under 18'
END;
END;
$$ LANGUAGE plpgsql;

IGNORE_WHEN_COPYING_START
```

```
content_copy download
Use code with caution. SQL
IGNORE WHEN COPYING END
```

How to Use the Function in a Query:

```
Generated sql
SELECT
user_id,
age,
BinUserAge(age) AS age_group
FROM
users;

IGNORE_WHEN_COPYING_START
content_copy download
Use code with caution. SQL
IGNORE_WHEN_COPYING_END
```

Explanation

- 1. **CREATE OR REPLACE FUNCTION BinUserAge(age INT)**: This defines a function named BinUserAge that accepts one integer argument, age.
- 2. **RETURNS VARCHAR**: Specifies that the function will return a string (the bin label).
- 3. **\$\$... \$\$ LANGUAGE plpgsql**: This is the standard syntax in PostgreSQL for defining the function body. The language is plpgsql (Procedural Language/PostgreSQL).
- 4. **BEGIN ... END**: This block contains the function's logic.
- 5. **RETURN CASE ... END**: A standard CASE statement implements the binning logic. It evaluates the input age and returns the corresponding string. It's crucial to handle NULL values and edge cases (like age < 18) gracefully.

Use Cases & Trade-offs

- Reusability & Consistency: The primary benefit. The same binning logic can be applied in dozens of queries without repeating the CASE statement, preventing errors and inconsistencies when business rules change.
- **Readability**: Queries become cleaner and more semantic. SELECT BinUserAge(age) is much more readable than a long, nested CASE statement.
- Performance Considerations: In some database systems, calling a UDF for every row
 can be slower than using an inline CASE statement because it may prevent the query
 optimizer from creating the most efficient execution plan (context switching). For very
 large-scale data processing, this trade-off between readability and performance should
 be considered and tested.

Question 15

Implement a SQL solution to compute the TF-IDF score for text data.

Answer:

Theory

TF-IDF (Term Frequency-Inverse Document Frequency) is a numerical statistic that reflects how important a word is to a document in a collection or corpus. It is a very advanced task to implement in pure SQL, but it serves as an excellent demonstration of complex data manipulation using chained Common Table Expressions (CTEs).

The calculation is broken down into three parts:

- 1. **Term Frequency (TF)**: How often a word appears in a document, normalized by the total number of words in that document. TF = (count of term t in doc d) / (total words in doc d)
- 2. **Inverse Document Frequency (IDF)**: A measure of how rare a word is across all documents. It dampens the weight of common words (like 'the', 'is'). IDF = log(total number of docs / number of docs containing term t)
- 3. **TF-IDF**: The product of the two, which gives a weight to each word in each document. TF-IDF = TF * IDF

Code Example

Scenario: We start with a table of tokenized words (a real-world pipeline would first tokenize raw text). We want to compute the TF-IDF score for each word in each document.

Sample Data:

```
doc_words table (pre-tokenized):
| doc_id | word |
|------|
| 1 | sql |
| 1 | is |
| 1 | cool |
| 2 | sql |
| 2 | powerful |
```

SQL Query:

```
Generated sql
WITH
-- Step 1: Calculate raw word counts per document
WordCounts AS (
SELECT doc_id, word, COUNT(*) AS term_count
```

```
FROM doc words
  GROUP BY doc_id, word
),
-- Step 2: Calculate total words per document to normalize TF
DocStats AS (
  SELECT doc id, COUNT(*) AS total words
  FROM doc words
  GROUP BY doc id
),
-- Step 3: Calculate Term Frequency (TF)
TermFrequency AS (
  SELECT
    wc.doc id,
    wc.word,
    wc.term count * 1.0 / ds.total words AS tf
  FROM WordCounts wc
  JOIN DocStats ds ON wc.doc_id = ds.doc_id
),
-- Step 4: Calculate Inverse Document Frequency (IDF)
InverseDocFrequency AS (
  SELECT
    word,
    LOG( (SELECT COUNT(DISTINCT doc_id) FROM doc_words) * 1.0 / COUNT(DISTINCT
doc id)) AS idf
  FROM doc_words
  GROUP BY word
-- Final Step: Calculate TF-IDF by joining the TF and IDF results
SELECT
  tf.doc_id,
  tf.word.
  tf.tf,
  idf.idf,
  tf.tf * idf.idf AS tf_idf
FROM TermFrequency tf
JOIN InverseDocFrequency idf ON tf.word = idf.word
ORDER BY doc_id, tf_idf DESC;
IGNORE_WHEN_COPYING_START
content copy download
Use code with caution. SQL
IGNORE_WHEN_COPYING_END
```

Explanation

- WordCounts CTE: Counts how many times each word appears in each document.
- 2. **DocStats CTE**: Counts the total number of words in each document.
- 3. **TermFrequency CTE**: Joins the first two CTEs to calculate TF using the formula term_count / total_words.
- 4. InverseDocFrequency CTE: This is the most complex part.
 - GROUP BY word ensures we calculate IDF for each unique word.
 - o The LOG() function calculates the IDF.
 - (SELECT COUNT(DISTINCT doc_id) FROM doc_words) is a scalar subquery to get the total number of documents in the corpus.
 - COUNT(DISTINCT doc_id) within the main query counts how many documents each specific word appears in.

5.

6. **Final SELECT**: Joins the TermFrequency and InverseDocFrequency CTEs on word to get both scores, then multiplies them to get the final tf_idf.

Best Practices and Optimization

- Materialization: This query is extremely computationally intensive. It should be run as a
 batch job and its results stored in a dedicated feature_tfidf table for later use by ML
 models.
- Tokenization: The initial step of breaking text into words (tokenization), as well as removing stop words and stemming, is usually better handled in a programming language like Python before the data is loaded into the doc_words table.

SQL ML Interview Questions - Scenario-Based Questions

Question 1

How would you write a SQL query to select distinct values from a column?

Answer:

Theory

To retrieve a unique list of values from a column in a table, you use the DISTINCT keyword in your SELECT statement. When DISTINCT is specified, the database engine processes the query results to remove all duplicate rows before returning the final set. This is a fundamental operation in exploratory data analysis (EDA) for understanding the cardinality and distribution of categorical features.

Multiple Solution Approaches

Approach 1: Using DISTINCT (Recommended and Standard)

This is the most direct and readable way to get unique values.

Generated sql

-- Select the unique list of product categories from the sales table. SELECT DISTINCT category FROM sales:

Approach 2: Using GROUP BY

This approach achieves the same result. The GROUP BY clause groups all rows with the same value in the specified column into a single summary row. When used without an aggregate function, its output is a list of the unique values in that column.

Generated sql

Select the unique list of product categories using GROUP BY.
 SELECT category
 FROM sales
 GROUP BY category;

IGNORE_WHEN_COPYING_START content_copy download
Use code with caution. SQL
IGNORE WHEN COPYING END

Explanation

- DISTINCT: The DISTINCT keyword is applied to all columns in the SELECT list.
 SELECT DISTINCT col1, col2 will return unique combinations of col1 and col2. In this case, SELECT DISTINCT category simply finds and returns each unique category name once.
- GROUP BY: The GROUP BY category clause tells the database to collect all rows with
 the same category value and condense them into a single representative row. Since we
 are not applying any aggregate functions like COUNT() or SUM(), the effect is simply
 returning the list of unique groups, which are the unique category names.

Performance Analysis and Trade-offs

- DISTINCT vs. GROUP BY: In most modern SQL database optimizers, there is no significant performance difference between SELECT DISTINCT and GROUP BY for the simple task of retrieving a unique list of values from a single column. Both operations typically require the database to perform a sort or use a hash table on the result set to identify and eliminate duplicates.
- **Readability**: SELECT DISTINCT is generally preferred for this specific task because its intent is clearer and more explicit. You are asking for "distinct values." GROUP BY is

- more semantically tied to aggregation, so using it without an aggregate function can be slightly less intuitive for other developers reading the code.
- **Indexing**: The performance of both queries can be significantly improved if the column you are selecting from (category in this case) is indexed. An index allows the database to retrieve the sorted list of unique values directly from the index structure without having to scan the entire table and then sort it.

Use Cases

- **Feature Engineering**: Discovering all possible values of a categorical feature to decide on an encoding strategy (e.g., one-hot encoding).
- **Data Exploration (EDA)**: Quickly understanding the variety of data in a column, such as finding all unique countries, user segments, or error codes.
- **Populating UI Elements**: Getting a unique list of items to populate a dropdown menu in a user interface.

Question 2

How would you optimize a SQL query that seems to be running slowly?

Answer:

Theory

Optimizing a slow SQL query is a systematic process of identifying and resolving performance bottlenecks. The goal is to reduce the amount of work the database has to do, such as minimizing disk I/O, CPU usage, and data transfer. My approach would follow these steps, starting with the highest-impact and most common solutions.

Step-by-Step Optimization Strategy

- 1. Analyze the Execution Plan (EXPLAIN plan):
 - This is the first and most critical step. The execution plan is the roadmap that the database optimizer has chosen to run the query. I would use EXPLAIN (or EXPLAIN ANALYZE in PostgreSQL, SHOWPLAN in SQL Server) to understand:
 - Full Table Scans: Is the database reading an entire table when it could be using an index? This is often the biggest bottleneck.
 - **Join Algorithms**: What type of join is being used (e.g., Nested Loop, Hash Join, Merge Join)? Is it appropriate for the data size?
 - Row Estimates vs. Actual Rows: Are the database's statistics up-to-date? If the estimated number of rows is vastly different from the actual number, the optimizer might choose a poor plan.

2.

3. Implement an Indexing Strategy:

- Missing or improper indexes are the most common cause of slow queries. Based on the EXPLAIN plan, I would:
 - Index WHERE clause columns: Any column used for filtering is a prime candidate for an index.
 - Index JOIN key columns: The foreign key columns used in ON clauses must be indexed to avoid slow nested loop joins on large tables.
 - Consider Composite Indexes: If the query filters on multiple columns, a composite (multi-column) index can be highly effective. The order of columns in the index should match the order in the query.

0

4.

5. Rewrite the Query Logic:

- o Sometimes the guery itself is written inefficiently. I would look for:
 - **SELECT***: Avoid using SELECT*, especially in joins. Explicitly list only the columns you need. This reduces data transfer and can sometimes allow for index-only scans.
 - Correlated Subqueries: These can be very slow as they execute once for each row of the outer query. I would rewrite them using JOINs or Common Table Expressions (CTEs).
 - WHERE vs. HAVING: Use WHERE to filter rows *before* aggregation and HAVING only to filter *after* aggregation. Filtering early reduces the amount of data that needs to be grouped.
 - UNION vs. UNION ALL: If I don't need to remove duplicates, UNION ALL is much faster because it avoids the costly sort/hash operation that UNION performs.
 - LIKE with leading wildcards: A query like WHERE name LIKE '%john' cannot use a standard B-tree index. If possible, avoid leading wildcards or use full-text indexing.

0

6.

7. Data Structure and Database-Level Optimizations:

- If query rewriting and indexing aren't enough, I would look at the underlying data structure:
 - Partitioning: For very large tables (especially time-series data), partitioning the table by a key (e.g., sale_date) allows the database to scan only the relevant partitions (partition pruning), drastically improving performance.
 - Materialized Views: For complex and frequent aggregation queries (e.g., for a dashboard), I would pre-calculate and store the results in a materialized view. The dashboard then queries this much smaller, pre-aggregated view.

■ **Updating Statistics**: Ensure database statistics are up-to-date (ANALYZE in PostgreSQL) so the query optimizer can make informed decisions.

0

8.

Code Example (Rewriting a Correlated Subquery)

Slow Query (Correlated Subquery):

```
Generated sql
-- Get the latest order for each customer

SELECT
c.customer_id,
c.name,
(SELECT MAX(order_date) FROM orders o WHERE o.customer_id = c.customer_id) AS

last_order_date

FROM
customers c;

IGNORE_WHEN_COPYING_START

content_copy download

Use code with caution. SQL

IGNORE WHEN COPYING END
```

Optimized Query (Using JOIN and GROUP BY):

```
Generated sql
  SELECT
  c.customer_id,
  c.name,
  o.last_order_date
FROM
  customers c
LEFT JOIN (
  SELECT customer id, MAX(order date) AS last order date
  FROM orders
  GROUP BY customer id
) AS o ON c.customer_id = o.customer_id;
IGNORE_WHEN_COPYING_START
content_copy download
Use code with caution. SQL
IGNORE_WHEN_COPYING_END
```

The optimized version calculates all the max dates in one pass over the orders table, which is far more efficient.

Question 3

How would you merge multiple result sets in SQL without duplicates?

Answer:

Theory

To merge the results of two or more SELECT statements into a single result set while automatically removing any duplicate rows, you use the UNION operator. The UNION operator combines the rows from each query and then performs a deduplication step, which is equivalent to applying DISTINCT to the entire combined result set.

For a UNION to work, the SELECT statements being combined must adhere to two rules:

- 1. They must have the same number of columns.
- 2. The corresponding columns must have compatible data types.

Code Example

Scenario: We have two tables, employees_us and employees_eu. An employee might exist in both tables if they have transferred. We want a single, unique list of all employee IDs across both regions.

Sample Data:

```
employees_us table:
| employee_id | name |
|------|
| 101 | Alice |
| 102 | Bob |

employees_eu table:
| employee_id | name |
|------|
| 102 | Bob | (Duplicate)
| 103 | Charlie |
```

SQL Query using UNION:

```
Generated sql SELECT employee_id, name FROM employees_us
```

UNION

SELECT employee_id, name FROM employees_eu;

```
IGNORE_WHEN_COPYING_START content_copy download
Use code with caution. SQL
IGNORE_WHEN_COPYING_END
```

Result:

```
| employee_id | name |
|------|
| 101 | Alice |
| 102 | Bob | (Appears only once)
| 103 | Charlie |
```

Explanation

- 1. **SELECT employee_id, name FROM employees_us**: The first query retrieves all employees from the US table.
- 2. **UNION**: This operator signals to the database that the results of the next query should be appended.
- 3. **SELECT employee_id, name FROM employees_eu**: The second query retrieves all employees from the EU table.
- 4. **Deduplication**: The database engine internally combines the two result sets and then filters out any rows that are identical. In this case, the row (102, 'Bob') exists in both sets, so UNION ensures it only appears once in the final output.

UNION vs. UNION ALL (Performance and Trade-offs)

It's crucial to understand the difference between UNION and its counterpart, UNION ALL.

- UNION: Removes duplicate rows. This operation is computationally expensive because it requires the database to sort or hash the entire combined result set to find the duplicates.
- **UNION ALL**: Does **not** remove duplicate rows. It simply appends the results. This is **much faster** and uses fewer resources.

Best Practice:

If you know for certain that the combined result sets will not have duplicates, or if duplicates are acceptable for your use case, **always use UNION ALL**. Only use UNION when you specifically require the deduplication.

For example, if we wanted to see every single employment record, even if it's a duplicate, we would use:

Generated sql
SELECT employee_id, name FROM employees_us
UNION ALL
SELECT employee_id, name FROM employees_eu;

IGNORE_WHEN_COPYING_START content_copy download
Use code with caution. SQL
IGNORE_WHEN_COPYING_END

This would return four rows, with (102, 'Bob') appearing twice.

Question 4

How would you handle very large datasets in SQL for Machine Learning purposes?

Answer:

Theory

Handling very large datasets (VLDBs) in SQL for machine learning requires a strategic approach focused on minimizing data movement and performing computations as efficiently as possible. The goal is to avoid pulling the entire dataset into memory. My strategy would involve a combination of sampling for development, in-database processing for feature engineering, and leveraging database-native optimizations.

Step-by-Step Strategy

1. Develop on a Sample:

- I would never work with the full dataset during initial development and experimentation. I would create a representative sample of the data.
- Stratified Sampling: For classification tasks with imbalanced classes (like fraud detection), I would use stratified sampling to ensure the sample preserves the original class distribution.
- SQL Implementation: Use database-specific sampling functions
 (TABLESAMPLE in PostgreSQL/SQL Server) or window functions
 (ROW_NUMBER() OVER (PARTITION BY ... ORDER BY RAND())) for stratified sampling.

2.

3. Perform Feature Engineering In-Database:

 Instead of extracting raw data and then processing it in Python/R, I would leverage the power of the SQL engine to do the heavy lifting. The database is optimized for large-scale data transformations.

- Complex SQL: Use CTEs and window functions to compute sophisticated features like rolling averages, lagged values, and group-based ratios directly in the database.
- Materialize Features: The final feature engineering query would run as a scheduled job (e.g., using Airflow or dbt), and its results would be stored in a new, clean feature table (or materialized view). The ML model training script then reads from this much smaller, processed table.

4.

5. Leverage Database Architecture and Optimizations:

- Partitioning: If the data is time-series based (e.g., logs, transactions), the table must be partitioned by date. My feature engineering queries would always include a WHERE clause on the partition key (e.g., WHERE sale_date > '2022-01-01') to enable partition pruning, which dramatically reduces the amount of data scanned.
- Columnar Storage: Modern data warehouses (like BigQuery, Snowflake, Redshift) use columnar storage. This is highly efficient for analytical queries because the query only needs to read the data for the columns it actually needs, not the entire row. My queries would avoid SELECT * to take full advantage of this.

6.

7. Consider In-Database Machine Learning:

- For certain tasks, I would consider using in-database ML platforms (e.g., BigQuery ML, Snowflake ML). These tools allow you to train and get predictions from models using SQL-like syntax without ever moving the data out of the warehouse.
- This is the ultimate solution for avoiding data movement, enhancing security, and leveraging the database's parallel processing capabilities.

8.

Code Example (Using Sampling and Partition Filtering)

This conceptual query shows how to generate features from a large, partitioned sales table.

Generated sql

SELECT

- -- Create a feature set using 10% of the data from the last year
- -- This assumes the 'sales' table is partitioned by 'sale_date'

```
CREATE TABLE ml_features_sample AS
WITH SampledData AS (
    SELECT *
    FROM sales TABLESAMPLE BERNOULLI (10) -- Use database sampling
    WHERE sale_date >= '2023-01-01' AND sale_date < '2024-01-01' -- Use partition pruning
)
-- Perform feature engineering only on the small, relevant sample
```

```
customer_id,
    COUNT(DISTINCT product_id) AS distinct_products_bought,
    AVG(sale_amount) AS avg_sale_amount,
    -- ... more features ...
FROM
    SampledData
GROUP BY
    customer_id;

IGNORE_WHEN_COPYING_START
content_copy download
Use code with caution. SQL
IGNORE WHEN COPYING END
```

Question 5

Discuss how you would design a system to regularly feed a Machine Learning model with SQL data.

Answer:

Theory

Designing a system to regularly feed an ML model with data from SQL requires building a robust, automated data pipeline. The architecture should be reliable, scalable, and maintainable. I would design the system by decoupling the different stages of the process: data extraction, feature computation, and model consumption (for both training and inference). I would use a workflow orchestration tool to manage the schedule and dependencies.

The system design depends on the model's use case: **batch prediction** (scoring many entities on a schedule) or **real-time prediction** (scoring a single entity on demand).

System Design for Batch Predictions (e.g., Daily Churn Prediction)

This is the most common pattern.

- 1. **Data Source:** The production OLTP database (e.g., PostgreSQL, MySQL) containing raw, normalized application data (customers, orders, etc.).
- 2. Workflow Orchestrator (e.g., Airflow, Prefect, Dagster):
 - This is the brain of the system. It manages the schedule (e.g., run daily at 1 AM) and dependencies between tasks.

3.

4. ETL/ELT Pipeline (using SQL and a tool like dbt):

- Task 1 (Extract & Load): A scheduled task copies new or updated raw data from the OLTP database into a centralized data warehouse (e.g., BigQuery, Snowflake, Redshift).
- Task 2 (Transform): A separate task, often managed by dbt, runs a series of SQL models (scripts) against the data warehouse. These SQL queries perform the heavy lifting of feature engineering: joining tables, aggregating data, calculating window functions, etc. The output is a clean, wide feature table where each row represents an entity (e.g., a customer) and each column is a feature. This process is version-controlled in Git.

5.

- 6. Model Training & Scoring Environment (e.g., Python on a VM or Kubernetes):
 - Task 3 (Scoring): The orchestrator triggers a Python script. This script:
 - a. Connects to the data warehouse and reads the latest data from the feature table.
 - b. Loads the trained ML model file (e.g., from an S3 bucket or model registry).
 - c. Makes predictions for each row.
 - d. Writes the predictions (e.g., customer_id, churn_probability, model_version) back into a dedicated predictions table in the data warehouse.

7.

8. **Downstream Consumption:** Other services or BI tools (like Tableau) can now use the predictions table to display results or take action.

System Design for Real-Time Predictions (e.g., Live Fraud Detection)

This requires a low-latency architecture.

- 1. **Online Feature Store:** In addition to the offline feature table in the data warehouse, we maintain an online feature store (e.g., Redis, DynamoDB). A streaming process or the batch pipeline keeps this online store updated with the latest features.
- 2. Model Serving API (e.g., Flask/FastAPI in a Docker container):
 - An application sends a request to this API with an entity ID (e.g., transaction id).

3.

- 4. Low-Latency SQL Query:
 - The API executes a highly optimized, simple SQL query against the online feature store or a read-replica of the production database to fetch the features for that *single* entity. The query must use an index for near-instantaneous lookup.
 - Example: SELECT * FROM real time features WHERE transaction id = ?;

5.

6. **Prediction and Response:** The model makes a prediction, and the API returns it to the calling application in milliseconds.

Question 6

How would you extract and prepare a confusion matrix for a classification problem using SQL?

Answer:

Theory

A confusion matrix is a table that summarizes the performance of a classification model. For a binary classification problem, it shows the counts of True Positives (TP), True Negatives (TN), False Positives (FP), and False Negatives (FN). We can generate these values directly from a predictions table in SQL using conditional aggregation. The predictions table should contain the model's prediction and the actual ground truth label.

The strategy is to pivot the data, grouping by the true label and then counting the predicted labels for each group using CASE statements.

Code Example

Scenario: We have a prediction_logs table that stores the results of a binary churn model (1 = churn, 0 = no churn). We want to generate a confusion matrix to evaluate its performance.

Sample Data:

SQL Query:

```
Generated sql
SELECT
true_label,
COUNT(CASE WHEN predicted_label = 1 THEN 1 END) AS predicted_as_churn,
COUNT(CASE WHEN predicted_label = 0 THEN 1 END) AS predicted_as_no_churn
FROM
prediction_logs
GROUP BY
true_label
ORDER BY
true_label DESC;
```

IGNORE_WHEN_COPYING_START content_copy download
Use code with caution. SQL
IGNORE_WHEN_COPYING_END

Result (The Confusion Matrix):

	true_label (Actual) predicted_as_churn (Predicted 1) predicted_as_no_churn (Predicted 0)
ı	
ı	
	1 (Is Churn) 2 (TP) 1 (FN)
I	0 (Is Not Churn) 1 (FP) 2 (TN)

Explanation

- 1. **GROUP BY true_label**: This groups the data into two buckets: one for all customers who actually churned (true_label = 1) and one for all who did not (true_label = 0). The query will produce one row for each of these buckets.
- 2. **COUNT(CASE WHEN predicted_label = 1 THEN 1 END)**: This is the conditional aggregation for the "Predicted as Churn" column.
 - For each row in a group, the CASE statement checks if the predicted_label is 1. If
 it is, it returns a value of 1; otherwise, it returns NULL.
 - The COUNT() function only counts non-NULL values. So, this expression effectively counts how many times the model predicted 1 within that group.

3.

- 4. **COUNT(CASE WHEN predicted_label = 0 THEN 1 END)**: This follows the same logic to count how many times the model predicted 0 within each group.
- 5. **ORDER BY true_label DESC**: This is for formatting, ensuring the "Actual Positive" (1) row appears before the "Actual Negative" (0) row, which is the standard convention for a confusion matrix.

Use Cases

- Model Performance Dashboards: This SQL query can directly power a visualization of a confusion matrix in a BI tool like Tableau or Looker, allowing for continuous performance monitoring.
- Calculating Metrics: Once you have the TP, FP, TN, and FN counts from this query, you
 can easily calculate other key metrics like Precision, Recall, and F1-score in subsequent
 SQL queries or in your BI tool.

Question 7

How would you log and track predictions made by a Machine Learning model within a SQL environment?

Answer:

Theory

Logging and tracking ML predictions is a cornerstone of MLOps and responsible AI. It enables model monitoring, debugging, auditing, and performance analysis. I would design a dedicated logging table in the SQL database with a well-defined schema to capture all relevant information about each prediction event. The application serving the model would be responsible for writing a new record to this table every time it makes a prediction.

Step 1: Design the Prediction Log Table Schema

I would create a table, for example model_predictions_log, with a schema designed for comprehensive analysis.

```
Generated sql
  CREATE TABLE model predictions log (
                   BIGINT PRIMARY KEY GENERATED ALWAYS AS IDENTITY, -- Unique
  prediction id
ID for each prediction
  model version
                    VARCHAR(50) NOT NULL,
                                                            -- Version of the model used
(e.g., 'churn-v2.1.3')
                 VARCHAR(255) NOT NULL,
  entity id
                                                        -- ID of the item being scored
(e.g., customer id, transaction id)
  prediction_timestamp TIMESTAMP WITH TIME ZONE DEFAULT CURRENT_TIMESTAMP,
-- When the prediction was made
  predicted label
                    VARCHAR(100),
                                                      -- The model's output label (e.g.,
'churn', 'not_churn')
  predicted probability FLOAT,
                                                   -- The model's confidence score (e.g.,
0.87)
  feature vector
                   JSONB,
                                                  -- (Optional but powerful) A JSON
object of the features used for the prediction
  ground_truth_label VARCHAR(100) NULL,
                                                           -- The actual outcome, to be
updated later
                                                                     -- When the
  ground truth updated at TIMESTAMP WITH TIME ZONE NULL
ground truth was recorded
);
-- For performance on a large log table, partitioning is essential.
-- PARTITION BY RANGE (prediction_timestamp);
IGNORE WHEN COPYING START
content copy download
Use code with caution. SQL
IGNORE_WHEN_COPYING_END
```

Explanation of Schema Design

- prediction_id: A unique primary key for every log entry.
- model_version: Crucial for comparing the performance of different model versions.
- entity_id: Links the prediction back to the business entity.
- prediction_timestamp: Essential for time-based analysis and monitoring for concept drift.
- predicted label and predicted probability: The core outputs of the model.
- feature_vector (JSONB): This is incredibly useful for debugging. If a model makes a strange prediction, you can inspect the exact features it used. JSONB is an indexed binary JSON format in PostgreSQL.
- **ground_truth_label**: This column is initially NULL. A separate batch process would run later (e.g., after 30 days for a churn model) to backfill this column with the actual outcome, enabling performance calculation.

Step 2: Implement the Logging Logic

The application code that calls the model would execute an INSERT statement after every prediction.

Conceptual Python code:

```
Generated python
  def get_prediction(customer_id, features):
  # model.predict(...) returns label and probability
  label, proba = model.predict(features)
  model_version = "churn-v2.1.3"
  # Log the prediction to the database
  log_query = """
  INSERT INTO model predictions log
    (model_version, entity_id, predicted_label, predicted_probability, feature_vector)
  VALUES (%s, %s, %s, %s, %s);
  db_cursor.execute(log_query, (model_version, customer_id, label, proba,
json.dumps(features)))
  db connection.commit()
  return label, proba
IGNORE_WHEN_COPYING_START
content copy download
Use code with caution. Python
IGNORE_WHEN_COPYING_END
```

Step 3: Use SQL to Analyze the Logs

With the data logged, I can now use SQL to monitor and analyze model performance.

Example Analysis Query (Model Accuracy Over Time):

```
Generated sql
  SELECT
  DATE TRUNC('day', prediction timestamp)::date AS prediction day,
  model version,
  COUNT(*) AS total predictions,
  AVG(CASE WHEN predicted_label = ground_truth_label THEN 1.0 ELSE 0.0 END) AS
accuracy
FROM
  model_predictions_log
WHERE
  ground truth label IS NOT NULL -- Only include predictions where we know the outcome
GROUP BY
  prediction day, model version
ORDER BY
  prediction day DESC;
IGNORE WHEN COPYING START
content copy download
Use code with caution. SQL
IGNORE WHEN COPYING END
```

Question 8

Discuss how to manage the entire lifecycle of a Machine Learning model using SQL tools.

Answer:

Theory

Managing the entire ML lifecycle using SQL-centric tools involves adopting a "data-centric MLOps" approach. This philosophy places the database and SQL at the heart of the workflow, from data preparation to monitoring. This approach is particularly powerful in organizations where data teams are strong in SQL and it promotes reproducibility, governance, and collaboration. Tools like **dbt (Data Build Tool)** are central to this modern workflow, complemented by features in modern data warehouses.

Here is how I would manage each stage of the lifecycle:

1. Data Preparation and Feature Engineering

- Tool: SQL itself, orchestrated by dbt.
- Process: Instead of Python scripts, the entire feature engineering logic is defined as a series of SQL SELECT statements (models) in dbt. Dbt manages the dependencies between these models, creating a Directed Acyclic Graph (DAG) of the entire transformation pipeline.

• Lifecycle Management:

- Version Control: All dbt SQL models are version-controlled in Git. This means our entire feature logic is auditable and reproducible.
- Testing: Dbt allows us to write data quality tests (e.g., not_null, unique) directly in YAML files, ensuring feature data integrity.
- Documentation: Dbt automatically generates documentation and a lineage graph, showing exactly how each feature is derived.

2. Model Training

• **Tool: In-database ML** (e.g., BigQuery ML, Snowflake ML) or a Python script that reads from the final dbt feature table.

Process:

- In-Database: I can use a SQL command like CREATE MODEL to train a model directly on the feature table created by dbt. The model itself becomes a first-class object within the database.
- Traditional: A Python script connects to the data warehouse, runs SELECT *
 FROM final_features_table, and trains the model. The key is that the data source
 is a clean, versioned, and tested table produced by the dbt pipeline.

3. Model Deployment and Serving

• Tool: SQL Stored Procedures, UDFs, or in-database PREDICT functions.

Process:

- The trained model (whether in-database or external) is exposed via a SQL interface.
- Batch Serving: A dbt model can call the PREDICT function to score a batch of new customers and materialize the results in a predictions table.
- Real-time Serving: A stored procedure can be created that takes an entity_id, looks up features, calls the PREDICT function, and returns the score. This provides a secure and stable API endpoint for applications.

4. Model Monitoring

• **Tool: SQL queries** on log tables, visualized in a **BI tool** (e.g., Tableau, Looker, Metabase).

Process:

- As described in the previous question, all predictions are logged to a dedicated table.
- SQL queries are written to calculate performance metrics (accuracy, precision, recall using the confusion matrix technique) and data drift (comparing statistical properties of features over time).
- These SQL queries power automated dashboards, providing a near real-time view of model health.

•

Summary of the SQL-Centric Lifecycle

Lifecycle Stage Primary SQL-Centric Tool/Technique

Data Prep dbt for version-controlled SQL transformations, CTEs, Window Functions

Training In-database CREATE MODEL (BigQuery ML) / Reading from dbt models

Deployment Stored Procedures, UDFs, PREDICT() functions

Monitoring SQL queries on prediction logs, BI dashboards

Governance Git (for SQL code), dbt docs (for data lineage)

This approach creates a highly transparent, reproducible, and robust ML lifecycle that is deeply integrated with the organization's core data platform and accessible to a wide range of data professionals.