# The Most Confusing SQL Functions

**Made clear and easy to understand**

Prepared by: Asfaw Gedamu

Do you find the following SQL functions confusing?

1. RANK vs DENSE_RANK  vs ROW_NUMBER

2. HAVING vs WHERE Clause

3. UNION vs UNION ALL

4. JOIN vs UNION

5. DELETE vs DROP vs TRUNCATE

6. CTE vs TEMP TABLE

7. SUBQUERIES vs CTE

8. ISNULL vs COALESCE

9. INTERSECT vs INNER JOIN

10. EXCEPT vs NOT IN

11. INNER JOIN vs LEFT JOIN vs RIGHT JOIN vs FULL JOIN

12. LAG() vs LEAD() Functions


You're not alone! Below is a concise guide to help visualize the key differences between these SQL commands and functions, complete with sample input data, SQL scripts, and expected outputs.


**1. RANK vs DENSE_RANK  vs ROW_NUMBER**

- RANK(): Leaves gaps after ties (1,1,3).  This function gives tied rows the same rank and then skips subsequent numbers.
- DENSE_RANK(): No gaps (1,1,2).  This function gives tied rows the same rank without skipping numbers.

- ROW_NUMBER assigns a unique, sequential number to each row.

👉 Use DENSE_RANK for leaderboards where gaps don't make sense!

Sample Table: players

```
1. | name  | score |
2. |-------|-------|
3. | Alice | 95    |
4. | Bob   | 95    |
5. | Carol | 90    |
```

SQL Query to Compare All Three Functions:

```sql
1. SELECT
2.     name,
3.     score,
4.     RANK () OVER (ORDER BY score DESC) AS rank,
5.     DENSE_RANK () OVER (ORDER BY score DESC) AS dense_rank,
6.     ROW_NUMBER () OVER (ORDER BY score DESC) AS row_num
7. FROM players;
8.
```

Expected Output Explanation:

RANK:

- Assigns the same rank to tied values but leaves gaps afterward.

Result:

```
1.    | name  | score | rank |
2.    |-------|-------|------|
3.    | Alice | 95    | 1    |
4.    | Bob   | 95    | 1    |
5.    | Carol | 90    | 3    |
6.
```

**DENSE_RANK:**

- Assigns the same rank to ties and does not leave gaps.

Result:

```
1.      | name  | score | dense_rank |
2.      |-------|-------|------------|
3.      | Alice | 95    | 1          |
4.      | Bob   | 95    | 1          |
5.      | Carol | 90    | 2          |
6.
```

**ROW_NUMBER:**

- Provides a unique sequential number to each row regardless of ties.

Result:

```
1.      | name  | score | row_num |
2.      |-------|-------|---------|
3.      | Alice | 95    | 1       |
4.      | Bob   | 95    | 2       |
5.      | Carol | 90    | 3       |
6.
```

**2. HAVING vs WHERE Clause**

- WHERE: Filters rows (before grouping).
- HAVING: Filters groups (after GROUP BY).

👉 "Show departments WHERE salary > 5000" vs "Show departments HAVING AVG(salary) > 5000."

Sample Table: employees

```
1. | department | salary|
2. |------------|-------|
```

```
 3. | HR          | 60000 |
 4. | IT          | 55000 |
 5. | HR          | 50000 |
 6. | IT          | 45000 |
 7. | IT          | 70000 |
 8. | Sales       | 40000 |
 9. | Sales       | 45000 |
10.
```

WHERE Clause (filters rows before GROUP BY):

```
1. SELECT department, COUNT() AS count
2. FROM employees
3. WHERE salary > 50000
4. GROUP BY department;
5.
```

Expected Output:

```
1. | department | count |
2. |------------|-------|
3. | HR         | 1     |
4. | IT         | 2     |
5.
```

HAVING Clause (filters groups after GROUP BY):

```
1. SELECT department, COUNT () AS emp_count
2. FROM employees
3. GROUP BY department
4. HAVING COUNT () > 1;
5.
```

Expected Output:

```
1. | department | emp_count |
2. |------------|-----------|
3. | HR         | 2         |
```

```
4. | IT          | 3            |
5. | Sales       | 2            |
6.
```

### 3. UNION vs UNION ALL

- UNION: "Merge and deduplicate."
- UNION ALL: "Merge and keep duplicates (way faster)."

👉 Use UNION ALL unless you explicitly need uniqueness.

Sample Data:

```
1. Table1 (names):
2. | name     |
3. |---------|
4. | Alice    |
5. | Bob      |
6. | Charlie |
7.
```

Table2 (names):

```
1. | name  |
2. |-------|
3. | Alice |
4. | David |
5.
```

UNION (removes duplicates):

```
1. SELECT name FROM Table1
2. UNION
3. SELECT name FROM Table2;
4.
```

Expected Output:

```
1. | name    |
2. |---------|
3. | Alice   |
4. | Bob     |
5. | Charlie |
6. | David   |
7.
```

UNION ALL (includes duplicates):

```
1. SELECT name FROM Table1
2. UNION ALL
3. SELECT name FROM Table2;
4.
```

Expected Output:

```
1. | name    |
2. |---------|
3. | Alice   |
4. | Bob     |
5. | Charlie |
6. | Alice   |
7. | David   |
8.
```

## 4. JOIN vs UNION

JOIN (combines columns from related rows):

Sample Tables:  employees

```
1. | emp_id | name  |
2. |--------|-------|
3. | 1      | Alice |
```

```
4. | 2       | Bob    |
5.
```

salaries

```
1. | emp_id | salary|
2. |--------|-------|
3. | 1      | 50000 |
4. | 2      | 60000 |
5. | 3      | 70000 |
6.
```

INNER JOIN Example:

```sql
1. SELECT a.emp_id, a.name, b.salary
2. FROM employees a
3. INNER JOIN salaries b ON a.emp_id = b.emp_id;
4.
```

Expected Output:

```
1. | emp_id | name   | salary|
2. |--------|-------|-------|
3. | 1      | Alice | 50000 |
4. | 2      | Bob   | 60000 |
5.
```

UNION (stacks similar rows):

Sample Queries:

```sql
1. SELECT name, 'employee' AS type FROM employees
2. UNION
3. SELECT name, 'manager' FROM (SELECT 2 AS emp_id, 'Bob' AS
name UNION SELECT 3, 'Carol');
4.
```

Expected Output:

```
1. | name  | type     |
2. |-------|----------|
3. | Alice | employee |
4. | Bob   | employee |
5. | Bob   | manager  |
6. | Carol | manager  |
7.
```

## 5. DELETE vs DROP vs TRUNCATE

- DELETE: "I'll remove these specific rows (and log every change)."
- TRUNCATE: "I'll wipe ALL rows (and reset the counter)."
- DROP: "I'll nuke the entire table (RIP)."

👉 Need speed? TRUNCATE. Need precision? DELETE.

Sample Table:  orders

```
1. | order_id | order_date  | amount |
2. |----------|-------------|--------|
3. | 1        | 2020-12-31  | 100    |
4. | 2        | 2021-01-15  | 150    |
5. | 3        | 2021-02-20  | 200    |
6.
```

DELETE (removes specific rows):

```
1. DELETE FROM orders
2. WHERE order_date < '2021-01-01';
3.
```

Expected Orders Table After DELETE:

```
1. | order_id | order_date  | amount |
2. |----------|-------------|--------|
3. | 2        | 2021-01-15  | 150    |
4. | 3        | 2021-02-20  | 200    |
5.
```

DROP (removes the table entirely):

```
1. DROP TABLE orders;
2.
```

TRUNCATE (removes all rows, retains structure):

```
1. TRUNCATE TABLE orders;
2.
```

Expected Orders Table After TRUNCATE:

```
1. | order_id | order_date | amount |
2. |----------|------------|--------|
3. | (empty)  |
4.
```

## 6. CTE vs TEMP TABLE

CTE (temporary result set within a single query):

- CTE: Disposable, single-query use.
- Temp Table: Reusable, session-persistent.

👉 CTEs for readability, Temp Tables for complex workflows.

Sample Table:  employees

```
1.  | emp_id | name  | salary |
2.  |--------|-------|--------|
3.  | 1      | Alice | 120000 |
4.  | 2      | Bob   | 90000  |
5.  | 3      | Carol | 110000 |
6.
```

CTE Example:

```
1.  WITH HighEarners AS (
2.    SELECT  FROM employees WHERE salary > 100000
3.  )
4.  SELECT  FROM HighEarners;
5.
```

Expected Output:

```
1.  | emp_id | name  | salary |
2.  |--------|-------|--------|
3.  | 1      | Alice | 120000 |
4.  | 3      | Carol | 110000 |
5.
```

TEMP TABLE (persists for session):

```
1.  CREATE TEMPORARY TABLE TempHighEarners AS
2.  SELECT  FROM employees WHERE salary > 100000;
3.
4.  SELECT  FROM TempHighEarners;
5.
```

Expected Output:

Same as the CTE output above.

### 7. SUBQUERIES vs CTE

Using SUBQUERIES:

Sample Tables:  employees

```
1. | emp_id | name  | dept_id |
2. |--------|-------|---------|
3. | 1      | Alice | 10      |
4. | 2      | Bob   | 20      |
5. | 3      | Carol | 10      |
6.
```

departments

```
1. | dept_id | location |
2. |---------|----------|
3. | 10      | NY       |
4. | 20      | LA       |
5.
```

Subquery Example:

```
1. SELECT name FROM employees
2. WHERE dept_id IN (SELECT dept_id FROM departments WHERE
location = 'NY');
3.
```

Expected Output:

```
1. | name  |
2. |-------|
3. | Alice |
4. | Carol |
5.
```

CTE Example:

```
1. WITH NYDepts AS (
2.    SELECT dept_id FROM departments WHERE location = 'NY'
3. )
4. SELECT name FROM employees WHERE dept_id IN (SELECT dept_id
FROM NYDepts);
5.
```

Expected Output:

Same as the subquery output.

## 8. ISNULL vs COALESCE

Sample Table:  contacts

```
1. | contact_id | phone    | mobile   |
2. |------------|---------|---------|
3. | 1          | NULL    | 123456   |
4. | 2          | 987654  | NULL     |
5. | 3          | NULL    | NULL     |
6.
```

ISNULL (SQL Server specific):

```
1. SELECT contact_id, ISNULL(phone, 'N/A') AS phone
2. FROM contacts;
3.
```

Expected Output:

```
1. | contact_id | phone  |
2. |------------|--------|
```

```
3. | 1          | N/A     |
4. | 2          | 987654  |
5. | 3          | N/A     |
6.
```

COALESCE (standard SQL):

```
1. SELECT contact_id, COALESCE(phone, mobile, 'N/A') AS
contact_number
2. FROM contacts;
3.
```

Expected Output:

```
1. | contact_id | contact_number |
2. |------------|----------------|
3. | 1          | 123456         |
4. | 2          | 987654         |
5. | 3          | N/A            |
6.
```

## 9. INTERSECT vs INNER JOIN

Sample Tables for INTERSECT:

Table1

```
1. | id |
2. |----|
3. | 1  |
4. | 2  |
```

```
5. | 3   |
6.
```

Table2

```
1. | id |
2. |----|
3. | 2  |
4. | 3  |
5. | 4  |
6.
```

INTERSECT Example:

```
1. SELECT id FROM Table1
2. INTERSECT
3. SELECT id FROM Table2;
4.
```

Expected Output:

```
1. | id |
2. |----|
3. | 2  |
4. | 3  |
5.
```

Sample Tables for INNER JOIN:

Table1

```
1.  | id | name |
2.  |----|------|
3.  | 1  | A    |
4.  | 2  | B    |
5.  | 3  | C    |
6.
```

## Table2

```
1.  | id | value |
2.  |----|-------|
3.  | 2  | X     |
4.  | 3  | Y     |
5.  | 4  | Z     |
6.
```

## INNER JOIN Example:

```
1.  SELECT a.id, a.name, b.value
2.  FROM Table1 a
3.  INNER JOIN Table2 b ON a.id = b.id;
4.
```

## Expected Output:

```
1.  | id | name | value |
2.  |----|------|-------|
3.  | 2  | B    | X     |
4.  | 3  | C    | Y     |
5.
```

## 10. EXCEPT vs NOT IN

Sample Tables:

Table1

```
1. | id |
2. |----|
3. | 1  |
4. | 2  |
5. | 3  |
6. | 4  |
7.
```

Table2

```
1. | id |
2. |----|
3. | 2  |
4. | 4  |
5.
```

EXCEPT Example:

```
1. SELECT id FROM Table1
2. EXCEPT
3. SELECT id FROM Table2;
4.
```

Expected Output:

```
1. | id |
2. |----|
3. | 1  |
4. | 3  |
5.
```

NOT IN Example:

```
1. SELECT id FROM Table1
2. WHERE id NOT IN (SELECT id FROM Table2);
3.
```

Expected Output:

```
1. | id |
2. |----|
3. | 1  |
4. | 3  |
5.
```

**11. INNER JOIN vs LEFT JOIN vs RIGHT JOIN vs FULL JOIN**

Definitions:

- INNER JOIN: Returns rows when there is a match between the tables.
- LEFT JOIN: Returns all rows from the left table and matching rows from the right table. Unmatched rows in the right table return NULL.
- RIGHT JOIN: Returns all rows from the right table and matching rows from the left table. Unmatched rows in the left table return NULL.
- FULL JOIN: Returns rows when there is a match in either the left or the right table, filling unmatched rows with NULL.

Sample Tables:

Customers Table

```
| CustomerID | Name    | City  |
|------------|---------|-------|
| 1          | Alice   | Paris |
| 2          | Bob     | Tokyo |
| 3          | Charlie | Delhi |
```

Orders Table:

```
| OrderID | CustomerID | Product |
|---------|------------|---------|
| 101     | 1          | Laptop  |
| 102     | 2          | Phone   |
| 103     | 4          | Camera  |
```

SQL Queries and Outputs:

INNER JOIN

```sql
SELECT Customers.Name, Orders.Product
FROM Customers
INNER JOIN Orders ON Customers.CustomerID = Orders.CustomerID;
```

Output:

```
| Name  | Product |
|-------|---------|
| Alice | Laptop  |
| Bob   | Phone   |
```

LEFT JOIN

```sql
SELECT Customers.Name, Orders.Product
FROM Customers
LEFT JOIN Orders ON Customers.CustomerID = Orders.CustomerID;
```

Output:

```
| Name    | Product |
|---------|---------|
| Alice   | Laptop  |
| Bob     | Phone   |
| Charlie | NULL    |
```

RIGHT JOIN

```sql
SELECT Customers.Name, Orders.Product
FROM Customers
RIGHT JOIN Orders ON Customers.CustomerID = Orders.CustomerID;
```

Output:

```
| Name  | Product |
|-------|---------|
| Alice | Laptop  |
| Bob   | Phone   |
| NULL  | Camera  |
```

FULL JOIN

```sql
SELECT Customers.Name, Orders.Product
FROM Customers
FULL JOIN Orders ON Customers.CustomerID = Orders.CustomerID;
```

Output:

```
| Name    | Product |
|---------|---------|
| Alice   | Laptop  |
| Bob     | Phone   |
| Charlie | NULL    |
| NULL    | Camera  |
```

## 12. LAG() vs LEAD() Functions

Definitions:

- LAG(): Accesses data from the previous row in the result set.

- LEAD(): Accesses data from the next row in the result set.

Sample Table:

Sales Table:

```
| SaleID | Product | SaleAmount | SaleDate    |
|--------|---------|------------|-------------|
| 1      | Laptop  | 1000       | 2023-01-01 |
| 2      | Phone   | 500        | 2023-01-02 |
| 3      | Tablet  | 700        | 2023-01-03 |
```

SQL Queries and Outputs:

LAG() Example

```sql
SELECT Product, SaleAmount, LAG(SaleAmount) OVER (ORDER BY
SaleDate) AS PreviousSale
FROM Sales;
```

Output:

```
| Product | SaleAmount | PreviousSale |
|---------|------------|--------------|
| Laptop  | 1000       | NULL         |
| Phone   | 500        | 1000         |
| Tablet  | 700        | 500          |
```

LEAD() Example

```sql
SELECT Product, SaleAmount, LEAD(SaleAmount) OVER (ORDER BY
SaleDate) AS NextSale
FROM Sales;
```

Output:

```
| Product | SaleAmount | NextSale |
|---------|------------|----------|
| Laptop  | 1000       | 500      |
| Phone   | 500        | 700      |
| Tablet  | 700        | NULL     |
```

That's a great follow-up! Even though LAG() and LEAD() may return NULL in some cases (especially at the boundaries of data), they are still incredibly valuable for many use cases.

### Why Use LAG() and LEAD()?

1. Comparing Sequential Rows:

   In datasets where order matters (like time-series data), you often need to compare the current row with a previous or next row to analyze trends.

   E.g., "Has the sale increased or decreased compared to the previous day?"

2. Trend Analysis:

   They are essential for calculating moving averages, growth rates, or detecting anomalies.

3. Handling Gaps in Data:

   Even when NULL  is returned, it's a helpful indicator that there is no relevant value available (like missing events or boundary conditions).

4. Flexible Calculation:

   You can handle NULL gracefully using functions like COALESCE() to replace them with default values (like 0 or the current row's value).

 Practical Demonstration

## Sales Data Example

```
| SaleDate   | Product | SaleAmount |
|-----------|---------|------------|
| 2023-01-01 | Laptop  | 1000       |
| 2023-01-02 | Phone   | 800        |
| 2023-01-03 | Tablet  | 1200       |
```

## SQL Query

```sql
SELECT
    SaleDate, Product, SaleAmount,
    LAG(SaleAmount) OVER (ORDER BY SaleDate) AS PreviousSale,
    LEAD(SaleAmount) OVER (ORDER BY SaleDate) AS NextSale,
    SaleAmount - LAG(SaleAmount) OVER (ORDER BY SaleDate) AS SaleChange
FROM Sales;
```

## Result:

```
|SaleDate    |Product  |SaleAmount| PreviousSale| NextSale| SaleChange|
|-----------|---------|----------|------------------
| 2023-01-01 | Laptop  | 1000     | NULL        | 800     | NULL       |
| 2023-01-02 | Phone   | 800      | 1000        | 1200    | -200       |
| 2023-01-03 | Tablet  | 1200     | 800         | NULL    | 400        |
```

## Why the NULLs Are Useful

- In row 1, NULL for PreviousSale clearly shows it's the first row with no earlier reference.
- In row 3, NULL for NextSale indicates it's the last row, so there's no subsequent sale to compare.
- The SaleChange column provides meaningful insights for the rows that have valid comparisons.

Even though NULLs appear at boundaries, they are indicators of valid structural limits in your dataset. Handling them with logic (like COALESCE() or conditional checks) ensures these functions remain powerful for insights. Would you like suggestions on how to visualize or analyze this further?

**What are moving averages, growth rates, or detecting anomalies?**

**1. Moving Averages**

A moving average is a calculation that takes the average of a fixed number of recent data points. It's useful for smoothing out short-term fluctuations and identifying long-term trends.

Use Case: Smoothing daily sales data over a 3-day period

Sample Data:

```
| SaleDate   | SaleAmount |
|------------|------------|
| 2023-01-01 | 100        |
| 2023-01-02 | 200        |
| 2023-01-03 | 300        |
| 2023-01-04 | 400        |
| 2023-01-05 | 500        |
```

SQL Query for 3-Day Moving Average

```sql
SELECT
    SaleDate,
    SaleAmount,
    AVG(SaleAmount) OVER (ORDER BY SaleDate ROWS BETWEEN 2
PRECEDING AND CURRENT ROW) AS MovingAvg3Days
FROM Sales;
```

Result:

```
| SaleDate   | SaleAmount | MovingAvg3Days |
|------------|------------|----------------|
| 2023-01-01 | 100        | 100            |
| 2023-01-02 | 200        | 150            |
| 2023-01-03 | 300        | 200            |
| 2023-01-04 | 400        | 300            |
| 2023-01-05 | 500        | 400            |
```

Explanation: The average smooths out the sharp changes and shows a more stable trend.

## 2. Growth Rates

A growth rate measures the percentage change from one period to the next.

Use Case: Calculating month-over-month sales growth

Sample Data:

```
| Month      | SaleAmount |
|------------|------------|
| 2023-01    | 1000       |
| 2023-02    | 1500       |
| 2023-03    | 2000       |
```

SQL Query for Growth Rate

```sql
SELECT
    Month,
    SaleAmount,
    (SaleAmount - LAG(SaleAmount) OVER (ORDER BY Month)) * 100.0
/ LAG(SaleAmount) OVER (ORDER BY Month) AS GrowthRate
FROM Sales;
```

Result:

```
| Month     | SaleAmount | GrowthRate |
|-----------|------------|------------|
| 2023-01   | 1000       | NULL       |
| 2023-02   | 1500       | 50.00      |
| 2023-03   | 2000       | 33.33      |
```

Explanation: In February, sales grew by 50% compared to January, and in March, they grew by 33.33% compared to February.

## 3. Detecting Anomalies

An anomaly is an unusual or unexpected data point that differs significantly from the rest.

Use Case: Detecting days when sales are significantly higher than the moving average

Sample Data:

```
| SaleDate   | SaleAmount |
|------------|------------|
| 2023-01-01 | 100        |
| 2023-01-02 | 200        |
| 2023-01-03 | 3000       |   -- Anomaly
| 2023-01-04 | 400        |
| 2023-01-05 | 500        |
```

SQL Query to Detect Anomalies

```sql
SELECT
    SaleDate,
    SaleAmount,
    AVG(SaleAmount) OVER (ORDER BY SaleDate ROWS BETWEEN 2
PRECEDING AND CURRENT ROW) AS MovingAvg3Days,
    CASE
        WHEN SaleAmount > 2 * AVG(SaleAmount) OVER (ORDER BY
SaleDate ROWS BETWEEN 2 PRECEDING AND CURRENT ROW)
        THEN 'Anomaly'
        ELSE 'Normal'
    END AS AnomalyStatus
FROM Sales;
```

Result:

```
| SaleDate    | SaleAmount | MovingAvg3Days | AnomalyStatus |
|-------------|------------|----------------|---------------|
| 2023-01-01  | 100        | 100            | Normal        |
| 2023-01-02  | 200        | 150            | Normal        |
| 2023-01-03  | 3000       | 1100           | Anomaly       |
| 2023-01-04  | 400        | 1200           | Normal        |
| 2023-01-05  | 500        | 1400           | Normal        |
```

Explanation: The sales on `2023-01-03` are much higher than the moving average, flagging it as an anomaly.

Summary

- Moving Averages: Smooth out fluctuations for better trend analysis.
- Growth Rates: Show percentage changes between periods for performance tracking.
- Anomaly Detection: Helps identify outliers for further investigation.
- LAG() vs LEAD(): Work within a single table to provide relative row data for analytical queries.
- **LAG()** looks backward (to the previous row) from the current row. Therefore, for the first row, there's no previous row, which results in NULL.
- **LEAD()** looks forward (to the next row) from the current row. For the last row, there is no next row, which results in NULL.
- The reason **LAG()** and **LEAD()** can return NULL is due to their access pattern in the result set.

This guide, complete with sample input and output data, should help novice users clearly understand the differences between these SQL commands and functions. Enjoy exploring and practicing these concepts!