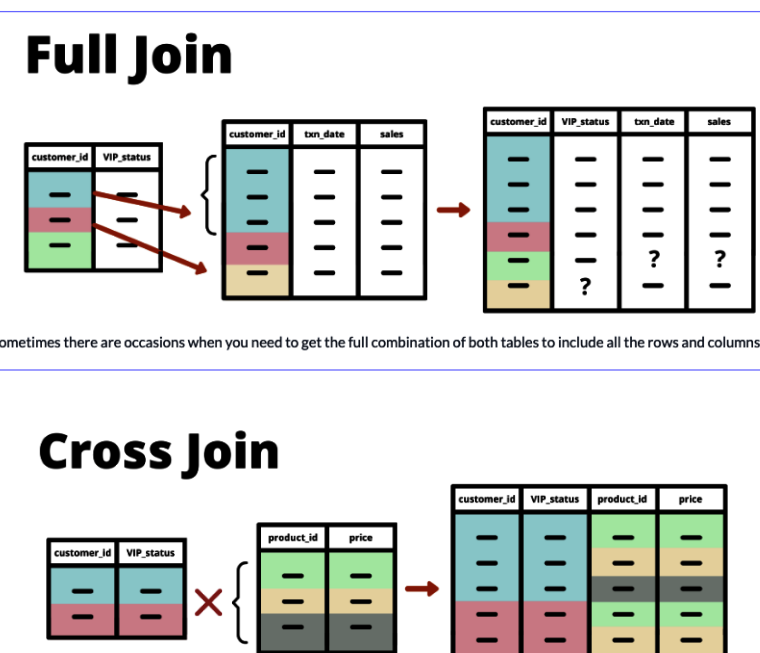


→ Temporary table + CTE

DROP TABLE IF EXISTS table\_name;  
CREATE TEMP TABLE table\_name AS  $\begin{pmatrix} \text{WITH } \langle \text{cte\_names } (\text{col}_1, \text{col}_2, \dots, \text{col}_n) \text{ AS } (\text{VALUES} \begin{pmatrix} \begin{pmatrix} \text{row 1} \\ \vdots \\ \text{row n} \end{pmatrix} \end{pmatrix} \end{pmatrix}$

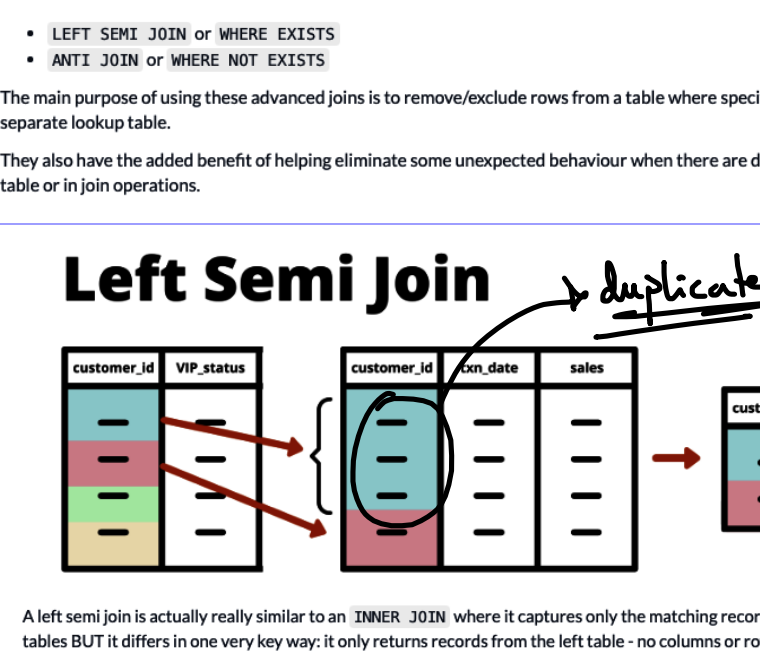
select \* from cte\_name;

Full Join



Sometimes there are occasions when you need to get the full combination of both tables to include all the rows and columns.

Cross Join



A cross join creates a full combination of all the rows in both tables that are being joined. This type of join is also referred to as a Cartesian join which is named after the original mathematician's contribution - the Cartesian product of two sets.

→ FROM t1 CROSS JOIN t2 / FROM t1, t2

- LEFT SEMI JOIN or WHERE EXISTS
- ANTI JOIN or WHERE NOT EXISTS

The main purpose of using these advanced joins is to remove/include rows from a table where specific column values based off a reported table/join.

They also have the added benefit of helping eliminate some unexpected behaviour when there are duplicates in the "target" table or in join operations.

Left Semi Join



A left semi join is actually really similar to an INNER JOIN, where it captures only the matching records between left and right tables BUT it filters out rows from the original table that do not have a matching record in the right table.

→ Target table acts as a filter

Syntax

SELECT FROM <base-table> WHERE EXISTS (SELECT 1 FROM <target-table> WHERE ...)

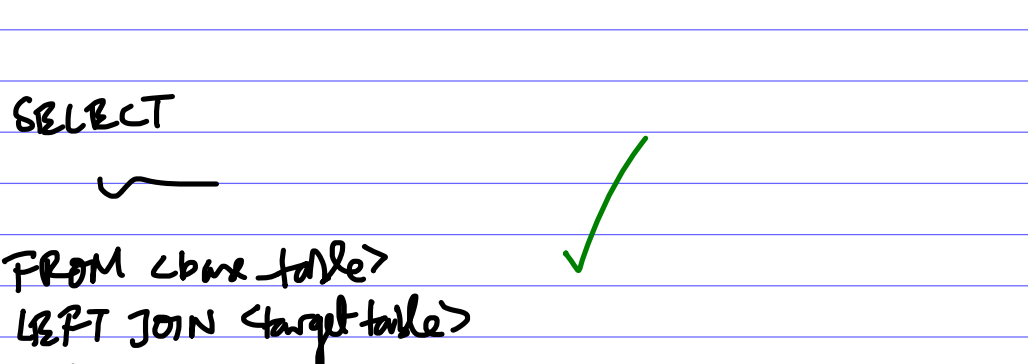
ORDER OF TABLES IMPORTANT

SELECT DISTINCT base-table, FROM <base-table> LEFT JOIN <target-table> ON WHERE <target-table> IS NOT NULL;

INEFFICIENT WHEN WE WANT ONLY RESULTS FROM LHS TABLE.

SELECT DISTINCT base-table, FROM <base-table> WHERE JOIN <target-table> ON

Anti Join



An ANTI JOIN is the opposite of a LEFT SEMI JOIN, where only records from the left which do not appear on the right are returned. It is expressed as a WHERE NOT EXISTS clause in PostgreSQL, but similar to the LEFT SEMI JOIN some SQL flavour support a direct ANTI JOIN syntax also.

Syntax

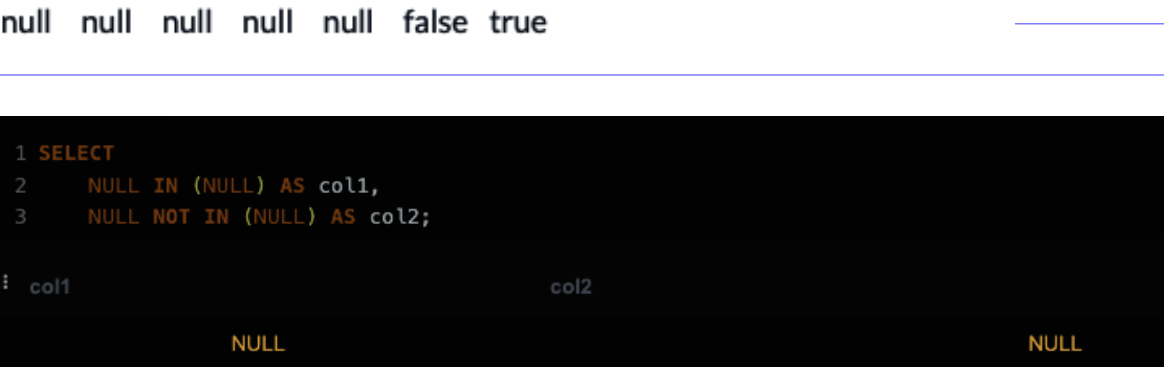
SELECT FROM <base-table> WHERE NOT EXISTS (SELECT 1 FROM <target-table> WHERE ...)

OR

ANTI JOIN <target-table> ON

SELECT FROM <base-table> LEFT JOIN <target-table> ON WHERE <target-table> IS NULL;

\* If we are joining tables which have NULL in the JOINING column THE RESULTS DON'T ACT AS EXPECTED.



col1 col2 col3 col4 col5 col6 col7  
null null null null null false true

```
1 SELECT  
2 NULL IN (NULL) AS col1,  
3 NULL NOT IN (NULL) AS col2;  
1 col1 col2  
NULL NULL
```

To keep things simple and less prone to issues - try to always stick with LEFT SEMI JOIN and ANTI JOIN instead of IN and NOT IN, where you do not know what exact values will appear for the IN subquery.

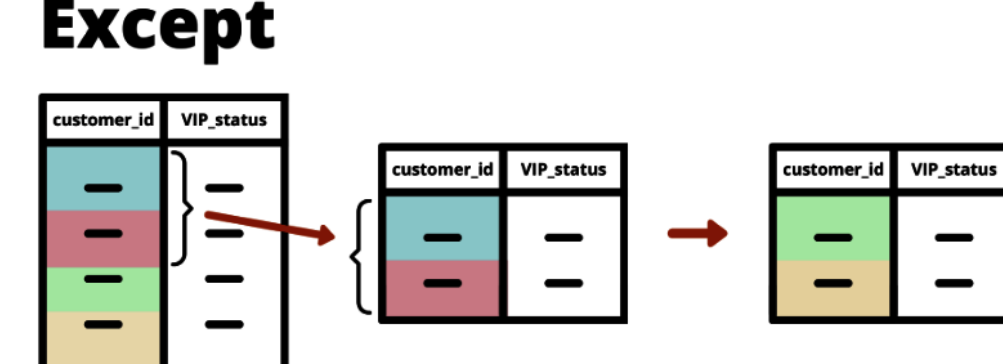
The behaviour when null values are present in the subquery for the filtering is not always stable and may lead to incorrect values!

Union

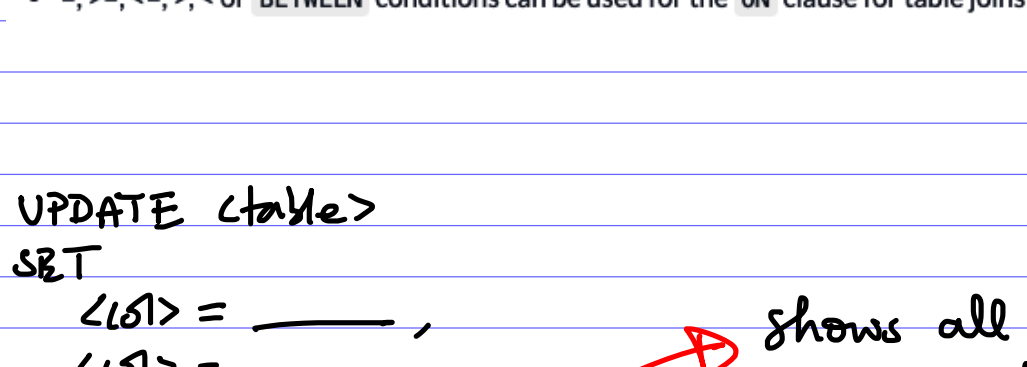


Simply put - UNION will be the union between two sets or SELECT results.

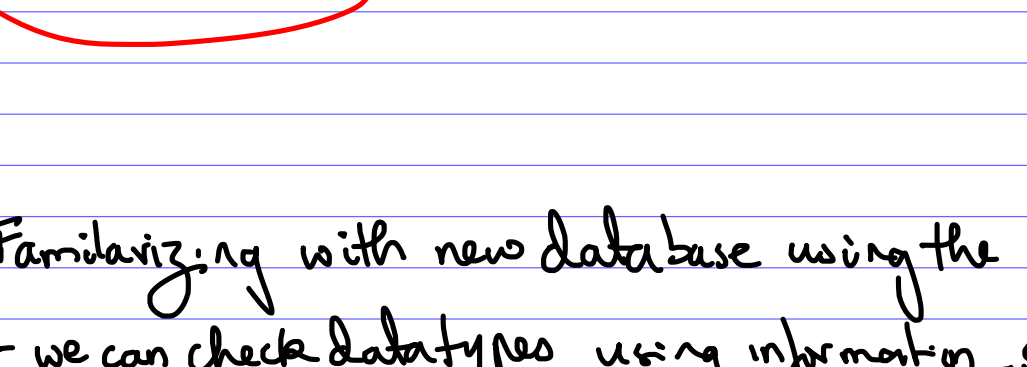
Union All



Intersect



Except



• =, >=, <=, >, < or BETWEEN conditions can be used for the ON clause for table joins

→ UPDATE <table>

SET <col> = ...  
WHERE ...  
RETURNING \*

shows all updated rows as query output

→ Familiarizing with new database using the information-scheme  
- we can check datatypes using information-scheme.columns table

select table-name, column-name, data-type  
from information-scheme.columns  
where table-name in ( ... );

High lvl overview of the dataset we will be working with

select schema-name, table-name, column-name  
from information-scheme.columns  
where column-name ILIKE '...';

→ Checking for missing values in multiple columns?

select \* from <table>  
where (col1 + col2 + ... + coln) is NULL;

Propagation of nulls  
1 + null = null

\* When dealing with time series based data & we encounter missing values, use the day before NEVER the day after.

DO NOT USE FUTURE DATA

**CALCULATION (INPUTS) OVER PARTITION BY ORDER BY FRAME CLAUSE**

**PARTITION BY**

- Splits rows of dataset into groups or window frames based on unique combination of values
- When unused, calculations will be performed on the entire dataset

**ORDER BY**

- Sorts rows within window frames by column or expression
- When unused, window frame rows are processed in a random order

**FRAME CLAUSE**

- Used to include/exclude rows within each window frame when performing calculations
- Default frame clause is RANGE BETWEEN UNBOUNDED PRECEDING AND CURRENT ROW

\* Window functions are not allowed in HAVING  
\* we can only use the outputs from the GROUP BY query in the HAVING clause.

1 Use a Random Function in the 'WHERE' Clause

In many SQL dialects, you can use a random number function directly in the 'where' clause to filter rows probabilistically. This method can work well if you want an approximate sample size or random selection based on a specific probability.

PostgreSQL

In PostgreSQL, you can use the RANDOM() function directly in the WHERE clause.

```
SELECT *  
FROM my_table  
WHERE RANDOM() < 0.1;
```

• Explanation: RANDOM() generates a random floating-point number between 0 and 1 for each row. The condition RANDOM() < 0.1 will approximately select 10% of the rows, providing a random sample without sorting the entire dataset.

MySQL

In MySQL, the equivalent function is RAND().

```
SELECT *  
FROM my_table  
WHERE RAND() < 0.1;
```

• Explanation: Similar to PostgreSQL, RAND() generates a random number between 0 and 1 for each row. The condition RAND() < 0.1 will approximately select 10% of the rows.

\* PERCENT\_RANK() vs CUME\_DIST() vs NTILE()

\* ORDER BY NULLS FIRST NULLS LAST

→ WINDOW clause simplification: if we find ourselves repeating a certain window multiple times in the same query, we can use the WINDOW clause after the final FROM section of the SQL query & refer to a simple alias in the window function call within the select expressions in the body of the SQL query.

select sum(<col>) over w  
from <table>  
window  
w as ( ... ),  
w2 as ( ... );

→ sum(<col>) over (order by ...)

↳ cumulative sum

→ Default FRAME: RANGE BETWEEN UNBOUNDED PRECEDING AND CURRENT ROW

**FRAME MODE BETWEEN START FRAME AND END FRAME FRAME EXCLUSION**

**START FRAME**

- Defines where window frame begins relative to current row
- UNBOUNDED or OFFSET
- PRECEDING or FOLLOWING
- OFFSET can be used with integer or date inputs in RANGE or ROWS mode
- When unused, defaults to UNBOUNDED PRECEDING and starts from the first sorted row per group

**END FRAME**

- Defines where window frame ends relative to current row
- UNBOUNDED or OFFSET
- PRECEDING or FOLLOWING
- OFFSET can be used with integer or date inputs in RANGE or ROWS mode
- When unused, defaults to CURRENT ROW and BETWEEN must not be used before the start frame

**FRAME EXCLUSION**

- Used to include/exclude rows inside the frame
- CURRENT ROW excludes only the current row and keeps everything else
- TIES removes duplicate rows with the same value as current row
- GROUP removes all rows with equal value as the current row including itself
- If unused, defaults to NO OTHERS and no rows are excluded

actual output using default behaviour

Experimentation with Recursive table sum 'n' natural numbers

PRACTICE MORE

→ FRAME MODES

ROWS	1	1	2	6	9	20	20	25
RANGE	8	8	7	3	11	11	16	
GROUPS	3	3	2	1	1	1	2	

UNBOUNDED, PRECEDING, FOLLOWING, CURRENT ROW

FRAME EXCLUSION: CURRENT ROW, TIES, GROUP, NO OTHERS (DEFAULT), PRACTICAL USE CASES??