

Click to verify



[illegible]

when the number of tables and columns in your query increases, your readers won't have to track down which column is in which table. That a query might break if you join a table with an ambiguous column name (e.g., both tables include a field called Created_At. Note that field filters are incompatible with table aliases, so you'll need to remove aliases when connecting filter widgets to your Field Filters. SQL best practices for WHERE Filter with WHERE before HAVING Use a WHERE clause to filter superfluous rows, so you don't have to compute those values in the first place. Only after removing irrelevant rows, and after aggregating those rows and grouping them, should you include a HAVING clause to filter out aggregates. Avoid functions on columns in WHERE clauses Using a function on a column in a WHERE clause can really slow down your query, as the function makes the query non-sargable (i.e., it prevents the database from using an index to speed up the query). Instead of using the index to skip to the relevant rows, the function on the column forces the database to run the function on each row of the table. And remember, the concatenation operator || is also a function, so don't get fancy trying to concat strings to filter multiple columns. Prefer multiple conditions instead: Avoid SELECT hero, sidekick FROM superheroes WHERE hero || sidekick = 'BatmanRobin' Prefer SELECT hero, sidekick FROM superheroes WHERE hero = 'Batman' AND sidekick = 'Robin' Prefer = to LIKE This is not always the case. It's good to know that LIKE compares characters, and can be paired with wildcard operators like %, whereas the = operator compares strings and numbers for exact matches. The = can take advantage of indexed columns. This isn't the case with all databases, as LIKE can use indexes (if they exist for the field) as long as you avoid prefixing the search term with the wildcard operator, %. Which brings us to our next point: Avoid bookending wildcards in WHERE statements Using wildcards for searching can be expensive. Prefer adding wildcards to the end of strings. Prefixing a string with a wildcard can lead to a full table scan. Avoid SELECT column FROM table WHERE col LIKE "%wizar%" Prefer SELECT column FROM table WHERE col LIKE "wizar%" Prefer EXISTS to IN If you just need to verify the existence of a value in a table, prefer EXISTS to IN, as the EXISTS process exits as soon as it finds the search value, whereas IN will scan the entire table. IN should be used for finding values in lists. Similarly, prefer NOT EXISTS to NOT IN. SQL best practices for GROUP BY Order multiple groupings by descending cardinality Where possible, GROUP BY columns in order of descending cardinality. That is, group by columns with more unique values first (like IDs or phone numbers) before grouping by columns with fewer distinct values (like state or gender). SQL best practices for HAVING Only use HAVING for filtering aggregates And before HAVING, filter out values using a WHERE clause before aggregating and grouping those values. SQL best practices for SELECT SELECT columns, not stars Specify the columns you'd like to include in the results (though it's fine to use * when first exploring tables — just remember to LIMIT your results). SQL best practices for UNION Prefer UNION All to UNION If duplicates are not an issue, UNION ALL won't discard them, and since UNION ALL isn't tasked with removing duplicates, the query will be more efficient. SQL best practices for ORDER BY Sorting is expensive. If you must sort, make sure your subqueries are not needlessly sorting data. SQL best practices for INDEX This section is for the database admins in the crowd (and a topic too large to fit in this article). One of the most common things folks run into when experiencing performance issues in database queries is a lack of adequate indexing. Which columns you should index usually depends on the columns you're filtering by (i.e., which columns typically end up in your WHERE clauses). If you find that you're always filtering by a common set of columns, you should consider indexing those columns. Adding indexes Indexing foreign key columns and frequently queried columns can significantly decrease query times. Here's an example statement to create an index: CREATE INDEX product_title_index ON products (title) There are different types of indexes available, the most common index type uses a B-tree to speed up retrieval. Check out our article on making dashboards faster, and consult your database's documentation on how to create an index. Use partial indexes For particularly large datasets, or lopsided datasets, where certain value ranges appear more frequently, consider creating an index with a WHERE clause to limit the number of rows indexed. Partial indexes can also be useful for date ranges as well, for example if you want to index the past week of data only. Use composite indexes For columns that typically go together in queries (such as last_name, first_name), consider creating a composite index. The syntax is similar to creating a single index. For example: CREATE INDEX full_name_index ON customers (last_name, first_name) EXPLAIN Look for bottlenecks Some databases, like PostgreSQL, offer insight into the query plan based on your SQL code. Simply prefix your code with the keywords EXPLAIN ANALYZE. You can use these commands to check your query plans and look for bottlenecks, or to compare plans from one version of your query to another to see which version is more efficient. Here's an example query using the dvdrental sample database available for PostgreSQL. EXPLAIN ANALYZE SELECT title, release_year FROM film WHERE release_year > 2000; And the output: Seq Scan on film (cost=0.00..66.50 rows=1000 width=19) (actual time=0.008..0.311 rows=1000 loops=1) Filter: ((release_year)::integer > 2000) Planning Time: 0.062 ms Execution Time: 0.416 ms You'll see milliseconds required for planning time, execution time, as well as the cost, rows, width, times, loops, memory usage, and more. Reading these analyses is somewhat of an art, but you can use them to identify problem areas in your queries (such as nested loops, or columns that could benefit from indexing), as you refine them. Here's PostreSQL's documentation on using EXPLAIN. WITH Organize your queries with Common Table Expressions (CTE) Use the WITH clause to encapsulate logic in a common table expression (CTE). Here's an example of a query that looks for the products with the highest average revenue per unit sold in 2019, as well as max and min values. WITH product_orders AS (SELECT o.created_at AS order_date, p.title AS product_title, (o.subtotal / o.quantity) AS revenue_per_unit FROM orders AS o LEFT JOIN products AS p ON o.product_id = p.id -- Filter out orders placed by customer service for charging customers WHERE o.quantity > 0) SELECT product_title AS product, AVG(revenue_per_unit) AS avg_revenue_per_unit, MAX(revenue_per_unit) AS max_revenue_per_unit, MIN(revenue_per_unit) AS min_revenue_per_unit FROM product_orders WHERE order_date BETWEEN '2019-01-01' AND '2019-12-31' GROUP BY product ORDER BY avg_revenue_per_unit DESC The WITH clause makes the code readable, as the main query (what you're actually looking for) isn't interrupted by a long sub query. You can also use CTEs to make your SQL more readable if, for example, your database has fields that are awkwardly named, or that require a little bit of data munging to get the useful data. For example, CTEs can be useful when working with JSON fields. Here's an example of extracting and converting fields from a JSON blob of user events. WITH source_data AS (SELECT events->'data'-->'name' AS event_name, CAST(events->'data'-->'ts' AS timestamp) AS event_timestamp CAST(events->'data'-->'cust_id' AS int) AS customer_id FROM user_activity) SELECT event_name, event_timestamp, customer_id FROM source_data Alternatively, you could save a subquery as a Snippet: And yes, as you might expect, the Aerodynamic Leather Toucan fetches the highest average revenue per unit sold. SQL is amazing. But so is Metabase's Query Builder. You can compose queries using Metabase's graphical interface to join tables, filter and summarize data, create custom columns, and more. And with custom expressions, you can handle the vast majority of analytical use cases, without ever needing to reach for SQL. Questions composed using the Query Builder also benefit from automatic drill-through, which allows viewers of your charts to click through and explore the data, a feature not available to questions written in SQL. Glaring errors or omissions? There are libraries of books on SQL, so we're only scratching the surface here. You can share the secrets of your SQL sorcery with other Metabase users on our forum. CTEs are named sets of results that help keep your code organized. They allow you to reuse results in the same query, and perform multi-level aggregations. Next article