# All things PostgreSQL from MinervaDB

Follow



Photo by Birger Strahl on Unsplash

# How to Improve PostgreSQL Efficiency: Removing Duplicate Rows Easily

## Boosting PostgreSQL Performance: Simple Steps to Eliminate Duplicate Data

Shiv Iyer

## Impact of Duplicate Rows on PostgreSQL Performance

Duplicate rows in a PostgreSQL database can significantly impact its
performance and efficiency in several ways:

1. **Increased Disk Space Usage**: Duplicates consume additional disk space
   unnecessarily, which can lead to increased storage costs and potentially
   degrade disk I/O performance as the database grows larger than needed.

2. **Reduced Query Performance**: More rows mean more data for PostgreSQL
   to scan through during queries, which can slow down query execution
   times, especially for full table scans or when indexes are not used
   effectively.

3. **Complications in Indexing**: Indexes on columns with many duplicate values
   are less effective. While duplicates in non-unique indexes are not directly
   harmful to the index's functionality, they still increase the size of the index,
   potentially reducing cache hit ratios and slowing down index scans.

4. **Data Integrity Issues**: In many cases, duplicates can lead to data integrity
   problems, making it difficult to ensure accurate data analysis, reporting, and
   decision-making based on the data.

5. **Increased Load on Maintenance Tasks**: Routine database maintenance
   tasks, such as vacuuming, indexing, and backups, can take longer to
   complete because there's simply more data to process.

## How to Eliminate Duplicate Records in PostgreSQL

Eliminating duplicate records involves identifying them and then deciding on a
strategy to remove or consolidate them. Here are general steps to remove
duplicates while keeping one instance of each duplicated set:

### 1. Using the `DISTINCT` Clause

For non-persistent remov~~al~~ ~~duplication of unique~~ ~~rows~~, use the `DISTINCT`
clause in your queries. T~~his~~ ~~~~ from the table but can
be used for reporting or data retrieval purposes.

COPY

```
SELECT DISTINCT column1, column2, ...
FROM my_table;
```

## 2. Deleting Duplicates While Keeping One Copy

If you need to remove duplicates from a table and keep one row of each duplicate set, one method is to use a CTE (Common Table Expression) with the row_number() window function:

COPY

```
WITH cte AS (
  SELECT
    column1, column2,
    row_number() OVER (PARTITION BY column_to_deduplicate ORDER BY id) AS rn
  FROM
    my_table
)
DELETE FROM cte
WHERE rn > 1;
```

◄ ███████████████████████████████ ►

This query assigns a row number to each row within a partition of duplicated records, ordered by some unique identifier ( id in this example). It then deletes all but the first row of each set of duplicates.

## 3. Using Temporary Tables to Remove Duplicates

Another approach involves creating a temporary table to hold the distinct rows, deleting the original data, and then repopulating the original table with the de-duplicated data:

COPY

```
CREATE TEMPORARY TABLE
SELECT DISTINCT ON (column_to_deduplicate) *
```

```
FROM my_table;


DELETE FROM my_table;


INSERT INTO my_table
SELECT * FROM temp_table;


DROP TABLE temp_table;
```

## Preventing Duplicates

- **Constraints**: Use unique constraints or primary keys to prevent the insertion of duplicate rows:

COPY

```
ALTER TABLE my_table ADD CONSTRAINT unique_constraint_name UNIQUE (colu
```

◄ ━━━━━━━━━━━━━━━━━━━━━ ►

- **Upsert**: Use the ON CONFLICT clause to perform an upsert (update or insert), which can prevent duplicates during data insertion:

COPY

```
INSERT INTO my_table (column1, column2)
VALUES ('value1', 'value2')
ON CONFLICT (column1) DO UPDATE SET column2 = EXCLUDED.column2;
```

Maintaining a database free of unnecessary duplicates is crucial for optimizing performance and ensuring data integrity. Regular monitoring and cleanup, combined with constraints to prevent duplicates, can help maintain the database's efficiency and accuracy.

## pg_test_fsync for Fsync Method Selection in PostgreSQL

Enhance PostgreSQL performance by leveraging pg_test_fsync to assess and...

minervadb.xyz

## PostgreSQL Transaction Logs: Implementation & Optimization

Discover how PostgreSQL WAL ensures data integrity and durability. Optimiz...

minervadb.xyz

**Boosting PostgreSQL Performance: Best Practices for Efficient Caching and Faster Query Response**

Learn how to boost PostgreSQL performance by configuring efficient caching…

minervadb.xyz

**Visualizing RLS Policy Checks in PostgreSQL Query Plans: Optimizing Performance and Enhancing Data Security**

Learn to visualize RLS policy checks in PostgreSQL query plans to optimize…

minervadb.xyz

# Subscribe to our newsletter

Read articles from **All things PostgreSQL from MinervaDB** directly inside your inbox. Subscribe to the newsletter, and don't miss out.

Enter your email address          SUBSCRIBE

SQL    PostgreSQL    MySQL    Databases    Data Science    Oracle

SQL Server    Devops    Cloud    Artificial Intelligence

Machine Learning    Analytics for Ecommerce    big data

## MORE ARTICLES

Shiv Iyer

Shiv Iyer

## How to Use Cumulative Aggregation in PostgreSQL: A Step-by-Step Retail Use Case

Cumulative Aggregation in PostgreSQL Cumulative aggregation, also referred to as running totals or c...

## Learn PostgreSQL Embedded JOINs: Examples and Techniques

In PostgreSQL, embedding JOINs within JOINs refers to the practice of chaining multiple JOIN operati...

Shiv Iyer

## Easy Steps to Transfer Data from Amazon EMR to Redshift

To load data from Amazon EMR (Elastic MapReduce) to Amazon Redshift, you need to follow a series of ...

©2025 All things PostgreSQL from MinervaDB

Archive  ·  Privacy policy  ·  Terms