# Data Management Deliverables for Personal Expense Tracker Project
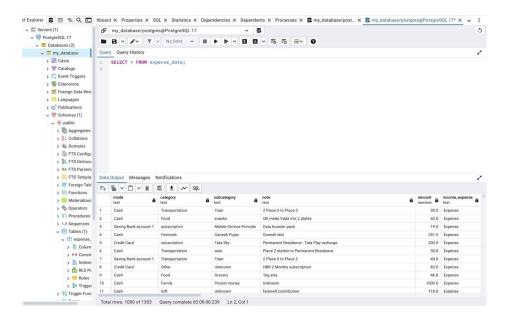
## 1. Cleaned and Validated Datasets

- **Initial Data**: The original dataset, `Daily Household Transactions.csv`, contained raw transaction records. This file included various transaction types (e.g., expenses, income) and metadata fields such as transaction category, subcategory, notes, amount, and mode.
- **Cleaning Process**: During data cleaning, issues such as missing values, incorrect data types, duplicates, and inconsistent categorization were identified and addressed.
  - **Duplicates**: Duplicate transactions were removed.
  - **Missing Values**: Missing values in essential columns (e.g., `amount`, `category`) were either filled or the rows were removed if essential information was missing.
  - **Data Type Conversion**: The `amount` field was converted to numeric format, and date fields were standardized to a consistent format.
  - **Categorization**: Standardized expense categories and subcategories for consistency.
- **Cleaned Dataset**: The cleaned data, `updated_data.csv`, is stored in PostgreSQL under the table `expense_data`, ready for analysis.
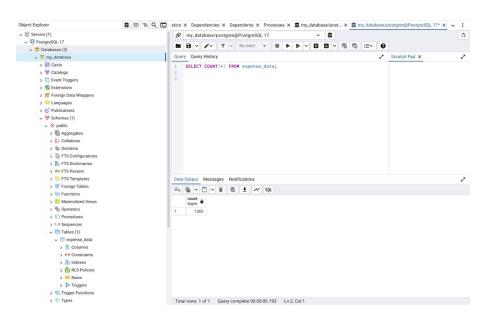
## 2. Database Management Reports

- **Database Structure**: The cleaned dataset was imported into PostgreSQL for efficient data management and querying. The database, `my_database`, contains the `expense_data` table, structured with the following fields:
  - `mode`: Transaction method (e.g., Cash, Credit Card).
  - `category`: General category of the transaction (e.g., Food, Transportation).
  - `subcategory`: Specifics within the main category (e.g., Snacks under Food).
  - `note`: Additional description or notes about the transaction.
  - `amount`: Numerical value of the transaction.
  - `income_expense`: Indicator of whether the transaction is an income or expense.
  - `currency`: The currency of the transaction amount (e.g., INR).

**Queries Performed**:

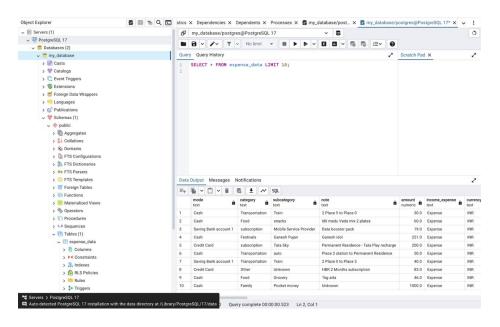- **Full Table Query**: `SELECT * FROM expense_data;` - Displays all records in the `expense_data` table for review.
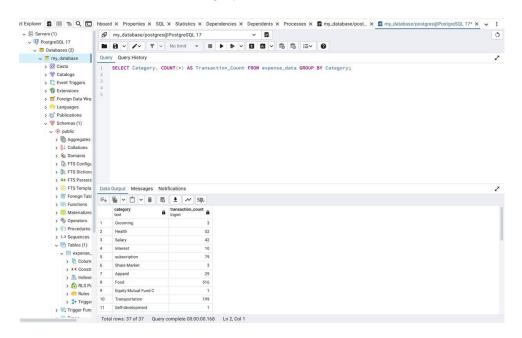
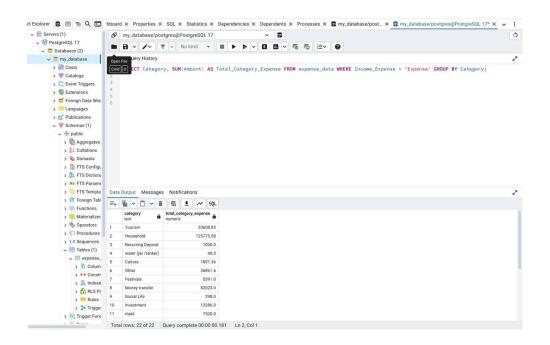- **Record Count**: `SELECT COUNT(*) FROM expense_data;` - Confirms the total number of transactions in the dataset.



- **Sample Query**: `SELECT * FROM expense_data LIMIT 10;` - Fetches a sample of 10 records for quick inspection

- **Category-Based Transaction Count**: `SELECT category, COUNT(*) AS transaction_count FROM expense_data GROUP BY category;` - Provides a count of transactions per category.



- **Category-Based Total Expense**: `SELECT category, SUM(amount) AS total_category_expense FROM expense_data WHERE income_expense = 'Expense' GROUP BY category;` - Sums up expenses per category

## 3. Data Quality Assessments

- **Data Completeness**: Verified that all essential columns (`amount`, `category`, `income_expense`) contain data with minimal missing values after cleaning.
- **Data Consistency**: Ensured that categories and subcategories are standardized, avoiding inconsistencies in spelling or naming conventions.
- **Accuracy Checks**: Cross-validated transaction amounts and categorization based on descriptions in the `note` field.
- **Integrity in PostgreSQL**: Ensured that the data uploaded to PostgreSQL maintained integrity by using constraints (e.g., numeric validation on `amount`) and setting up data types accordingly.

## 4. Data Collection Methodologies

- **Data Sources**: Transaction data was collected from user inputs and uploads, encompassing various daily transactions.
- **Data Entry Standards**: Established guidelines for entering data to ensure consistency. Categories were predefined to minimize manual entry errors, and descriptive notes were encouraged for clarity.
- **Automation Potential**: Data was loaded into PostgreSQL, allowing for automated data updates and additional processing if needed for future transactions.