# SQL Data Preparation and Cleaning Report

**PROJECT TITLE:** STRAVA FITNESS DATA ANALYSIS

## OBJECTIVE

Cleaning, normalising, and standardising the raw statistics gathered from different fitness trackers (such as steps, sleep, calories, etc.) in order to get them ready for visual analytics in Power BI and Python was the aim of this SQL phase. I completed necessary data cleaning activities in the SQLite database to get the Fitbit datasets ready for analysis. This required adding many CSV files to the database, including hourlyCalories_merged, sleepDay_merged, and dailyActivity_merged. I handled missing data, made sure that column formats were consistent (particularly for date and time), and eliminated duplicate entries using GROUP BY and HAVING clauses. I converted datetime fields to normal date formats as necessary. Additionally, I used common keys like Id and ActivityDate to confirm the links across tables. These procedures made sure the data was accurate and dependable for further Power BI and Python analysis.

**Tools Used**: SQLite Workbench

## Data Files Handled:

- dailyActivity_merged

- heartrate_daily

- hourlyCalories_cleaned

- hourlyIntensities_cleaned

- hourlySteps_cleaned
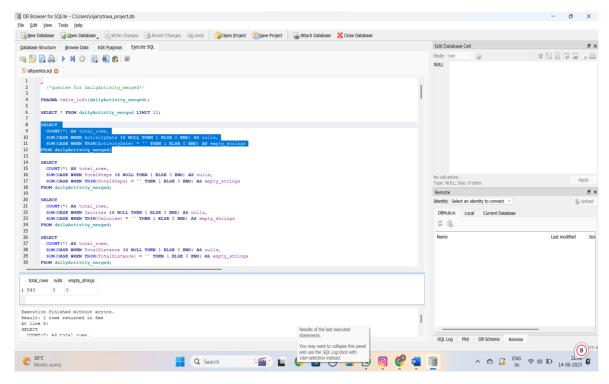
- sleepDay_cleaned

- weightLog_dates_fixed

We performed date format corrections, column standardisations, and ensured all datasets are free from NULLs and empty Strings. Also checked for duplicates and finally invalid values (like negative steps, calories, etc.)
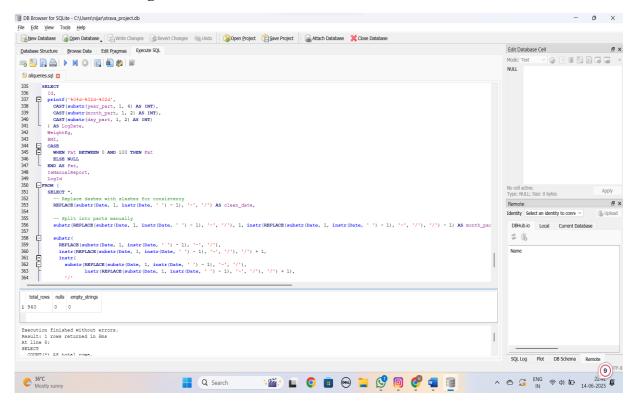
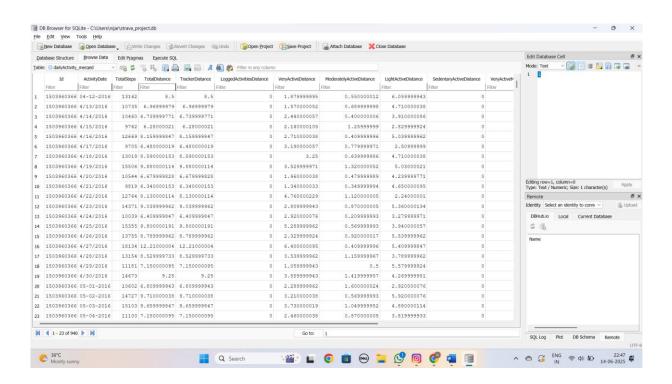# SQL QUERIES AND LOGIC

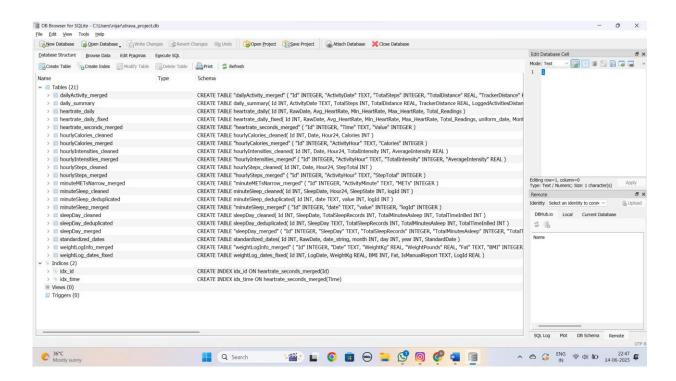## 1. Checking NULLS and Invalid values:



## 2. Checking Duplicates:

```
/* queires for heartrate_seconds_merged */

PRAGMA table_info(heartrate_seconds_merged);

SELECT * FROM heartrate_seconds_merged LIMIT 10;

SELECT
  SUM(CASE WHEN Id IS NULL THEN 1 ELSE 0 END) AS null_ids,
  SUM(CASE WHEN Time IS NULL THEN 1 ELSE 0 END) AS null_time,
  SUM(CASE WHEN Value IS NULL THEN 1 ELSE 0 END) AS null_value
FROM heartrate_seconds_merged;

SELECT *
FROM heartrate_seconds_merged
WHERE Value < 0;

-- for faster querying we are using indexing

CREATE INDEX idx_id ON heartrate_seconds(Id);
CREATE INDEX idx_time ON heartrate_seconds_merged(Time);

DROP TABLE IF EXISTS heartrate_daily;


CREATE TABLE heartrate_daily AS
SELECT
  Id,
  -- Just grab the first 10 characters if format is always MM/DD/YYYY or M/D/YYYY
  substr(Time, 1, instr(Time, ' ') - 1) AS RawDate,
```

### 3. Standardising Date Formats:



# Final Output and Observations

**After cleaning:**

- All datasets use either the mm-dd-yyy or mm/dd/yyy date format.

- NULL values and zero entries were removed.

- Column data types (like integers for steps and floats for weight) were validated.

**Challenges Faced:**

- Inconsistent date formats like yyyy-mm-dd 00:00 AM/PM required conditional handling.

- Some datasets lacked primary keys, so joins required caution.