# SQL Automated Data Cleaning

This documentation outlines a project with the goal of automating the data cleaning. The project aims to enhance data quality by leveraging stored procedure and event to remove duplicates, correct data inconsistencies and standardize data fields.

The project has two key objectives:

- To create a cleaned version of the data table.
- To automate the data cleaning process using stored procedure, event and trigger.

The project involves creating a new cleaned data table, developing a stored procedure for data cleaning, setting up automated events to schedule regular data cleaning and implementing triggers to clean data immediately after insertion.

In this project, the following key stages are undertaken:

- Creating the Cleaned Data Table: Defining and creating a new table for cleaned data.

*(The best practice is to have a copy of the raw data and then carry out the data cleaning process)*

```sql
1   USE bakery;
2
3   DELIMITER $$
4   DROP PROCEDURE IF EXISTS Copy_and_Clean_Data;
5   CREATE PROCEDURE Copy_and_Clean_Data()
6   BEGIN
7       -- CREATING The TABLE
8       CREATE TABLE IF NOT EXISTS `us_household_income_Cleaned` (
9           `row_id` int DEFAULT NULL,
10          `id` int DEFAULT NULL,
11          `State_Code` int DEFAULT NULL,
12          `State_Name` text,
13          `State_ab` text,
14          `County` text,
15          `City` text,
16          `Place` text,
17          `Type` text,
18          `Primary` text,
19          `Zip_Code` int DEFAULT NULL,
20          `Area_Code` int DEFAULT NULL,
21          `ALand` int DEFAULT NULL,
22          `AWater` int DEFAULT NULL,
23          `Lat` double DEFAULT NULL,
24          `Lon` double DEFAULT NULL,
25          `TimeStamp` TIMESTAMP DEFAULT NULL
26      ) ENGINE=InnoDB DEFAULT CHARSET=utf8mb4 COLLATE=utf8mb4_0900_ai_ci;
27
28      -- COPY DATA TO NEW TABLE
29      INSERT INTO us_household_income_Cleaned
30      SELECT *, CURRENT_TIMESTAMP
31      FROM bakery.us_household_income;
```

Adding a 'TimeStamp' column is important as it aids in debugging. With an automated process, data may change over time and 'TimeStamp' would make a significant difference when identifying and correcting errors.



- Procedure for Data Cleaning: Developing a stored procedure to copy data, remove duplicates and standardize fields.

```
35      -- 1. Remove Duplicates
36      DELETE FROM us_household_income_Cleaned
37      WHERE
38          row_id IN (
39              SELECT row_id
40      FROM (
41          SELECT row_id, id,
42              ROW_NUMBER() OVER (
43                  PARTITION BY id, `TimeStamp`
44                  ORDER BY id, `TimeStamp`) AS row_num
45          FROM
46              us_household_income_Cleaned
47      ) duplicates
48      WHERE
49          row_num > 1
50      );
51
52      -- 2. Standardization
53      UPDATE us_household_income_Cleaned
54      SET State_Name = 'Georgia'
55      WHERE State_Name = 'georia';
56
57      UPDATE us_household_income_Cleaned
58      SET County = UPPER(County);
59
60      UPDATE us_household_income_Cleaned
61      SET City = UPPER(City);
62
63      UPDATE us_household_income_Cleaned
64      SET Place = UPPER(Place);
65
66      UPDATE us_household_income_Cleaned
67      SET State_Name = UPPER(State_Name);
68
69      UPDATE us_household_income_Cleaned
70      SET `Type` = 'CDP'
71      WHERE `Type` = 'CPD';
72
73      UPDATE us_household_income_Cleaned
74      SET `Type` = 'Borough'
75      WHERE `Type` = 'Boroughs';
```

- **Testing the Procedure: Manually executing the procedure to validate its functionality.**



- **Automating the Cleaning Process: Setting up a MySQL event to run the cleaning procedure periodically and creating triggers for real-time data cleaning.**

```sql
101
102      -- Create Event
103      DROP EVENT run_data_cleaning;
104      CREATE EVENT run_data_cleaning
105          ON SCHEDULE EVERY 30 DAY
106          DO CALL Copy_and_Clean_Data();
107
108      -- CREATE TRIGGER
109      DELIMITER $$
110      CREATE TRIGGER Transfer_clean_data
111          AFTER INSERT ON bakery.us_household_income
112          FOR EACH ROW
113      BEGIN
114          CALL Copy_and_Clean_Data();
115      END $$
116      DELIMITER ;
117
118      INSERT INTO bakery.us_household_income
119      (`row_id`,`id`,`State_Code`,`State_Name`,`State_ab`,`County`,`City`,`Place`,`Type`,`Primary`,`Zip_Code`,`Area_Code`,`ALand`,`AWater`,`Lat`,`Lon`)
120      VALUES
121      (121671,37025904,37,'North Carolina','NC','Alamance County','Charlotte','Alamance','Track','Track',28215,980,24011255,98062070,35.2661197,-80.6865346);
122
123
```
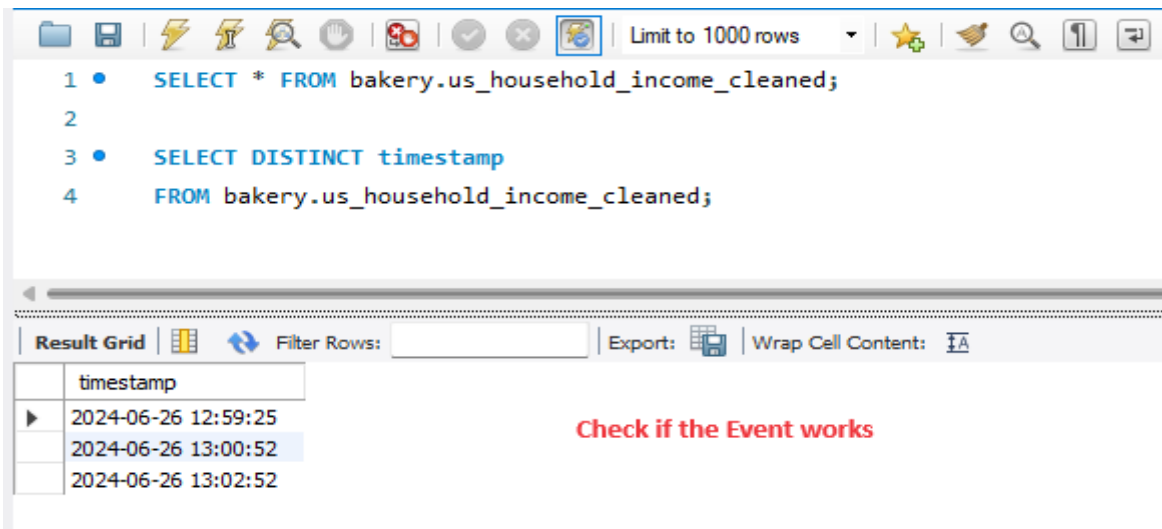
- We might encounter the error of getting the table empty when running the event. That's when 'TimeStamp' comes in. It is unique to the time when the stored procedure runs. So, by including 'TimeStamp' when querying for duplicates, only the timestamped data will be removed when the event runs.

```sql
35    -- 1. Remove Duplicates
36    DELETE FROM us_household_income_Cleaned
37    WHERE
38        row_id IN (
39            SELECT row_id
40        FROM (
41            SELECT row_id, id,
42                ROW_NUMBER() OVER (
43                    PARTITION BY id, `TimeStamp`
44                    ORDER BY id, `TimeStamp`) AS row_num
45            FROM
46                us_household_income_Cleaned
47        ) duplicates
48        WHERE
49            row_num > 1
50    );
```

- Verification and Validation: Performing queries to check data integrity and consistency in the cleaned data table.

*(Initially, I tested with a 2-minute scheduled event, but later changed it to 30 days.)*

```sql
1 ● SELECT * FROM bakery.us_household_income_cleaned;
2
3 ● SELECT DISTINCT timestamp
4   FROM bakery.us_household_income_cleaned;
```

| timestamp |
|---|
| 2024-06-26 12:59:25 |
| 2024-06-26 13:00:52 |
| 2024-06-26 13:02:52 |

**Check if the Event works**

**TIPS**

When dealing with large datasets, it's essential to optimize SQL queries and processes to ensure efficient performance. Here are a few techniques to consider:

INDEXES:

- Create indexes on columns that are frequently used in WHERE clauses, joins, and sorting operations to speed up query performance.

Common Table Expressions (CTEs):

- Use CTEs for better readability and to break down complex queries into simpler parts.

Batch Processing:

- For large data migrations or updates, consider processing data in smaller batches to avoid long-running transactions and reduce locking issues.