# EYouth X DEPI Tech Challenge

# Data Cleaning and Modeling Report

**Presented to:** EYouth

**Date of Submission:** 28/2/2025

**Submitted by:** Mohammed Yasir Abduljawad

# Contents

# Objectives of Data Cleaning

The primary objective of this data cleaning process is to ensure the accuracy, consistency, and reliability of the **Supplier Quality Dataset**. By addressing common data quality issues, this process aims to enhance the validity of subsequent analysis and reporting. The specific goals of data cleaning include:

- **Handling Missing Data:** Identifying and addressing any missing values in key columns to prevent gaps in analysis.

- **Correcting Inconsistencies:** Standardizing naming conventions, data formats, and measurement units to ensure uniformity across records.

- **Removing Duplicates:** Eliminating redundant records to avoid double counting and biased results.

- **Resolving Errors:** Identifying and correcting incorrect values, such as negative defect counts or unrealistic downtime durations.

- **Ensuring Data Integrity:** Verifying relationships between different columns, such as supplier IDs matching the correct supplier names.

By conducting a thorough data cleaning process, the dataset will be optimized for accurate supplier quality analysis, ensuring meaningful insights for decision-making.

# Data Sources

The dataset used in this analysis originates from supplier quality monitoring systems and consists of records related to **defect rates, vendor, plant, and material details.**

The dataset is provided in **Google sheets** and exported in **Excel format (Supplier Quality Analysis Dataset.xlsx)** to import it to the Power Query easily, structured into multiple fields capturing essential supplier quality details needed for analysis.

# Cleaning steps

## Template

**The steps will be in this template for each table:**

**Table name:**

**Notes:**

**Changes made:**

**Data type changes:**

Transformed all sheets except description sheet.

## Category table



*Figure 1: Category Table Before*

There are two columns **Category, Sub-Category, and Sub-Category ID.** The category columns and sub-category columns are the same. There are no missing data, duplicates or errors.

**Changes made:**

- Sub-Category column is removed
- Sub-Category ID column header is changed to Category ID

**Data type changes:**

- **Category ID column** is changed from **Whole number** to **Text**



*Figure 2: Category Table After*

# Defect table



*Figure 3: Defect Table Before*

- The defect column has **305 defects with 272 distinct and 245 unique defects**. This indicates that there is a problem with duplicates, deleting them will drop **17%** of the records which will lead to serious problems so we could clarify that instead of deleting them. Also there are a lot of data that is quite similar like **"Wrong Size" and "Wrong Size", "Warped" and "Warping"** which cannot be deleted for the same reason. Although deleting is not the only solution available but all solutions would be impractical in such situations.

**Data type changes:**

- **Defect ID column** is changed from **Whole number** to **Text**



*Figure 4: Defect Table After*

## Defect type table



*Figure 5: Defect type Table Before*

There are no missing data, duplicates or errors. The IDs only might seem odd but no problem with that.

**Data Changes:**

- **Defect Type ID column** is changed from **Whole number** to **Text**



*Figure 6: Defect type Table After*

## Material table



*Figure 7: Material Table Before*

There are no missing data, duplicates or errors.

**Data changes:**

- **Material Type ID column** is changed from **Whole number** to **Text**



*Figure 8: Material Table After*

# Plant Table



Figure 9: Plant Table Before

There are no missing data, duplicates or errors. It will be better to separate state abbreviations from the cities which can be used later to apply a hierarchy.

**Changes made:**

- Replaced **Cincinnati OH** with **Cincinnati, OH**.
- Split Plant column to two Plant.1 and Plant.2 using delimiters ",".
- Trimmed Plant.2 column
- Rename columns plant.1 = Plant City and plant.2 = Plant State
- Replaced **Wi** with **WI**

**Data changes:**

- **Material Type ID column** is changed from **Whole number** to **Text**.



Figure 10: Plant Table After

## Vendor Table



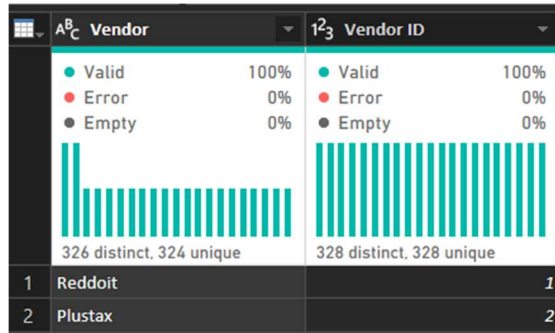*Figure 11: Vendor Table Before*

There are no missing data or errors. But there are two duplicate **roundphase (80 & 113)** and **Quotefix (125 & 144).** But deleting them may cause problems as they are recorded in the report so we could **clarify** that **instead of deleting** them.

**Data changes:**

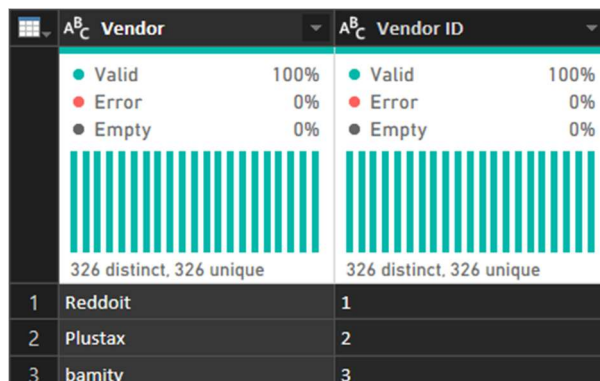- **Vendor ID column** is changed from **Whole number** to **Text.**



*Figure 12: Vendor Table After*
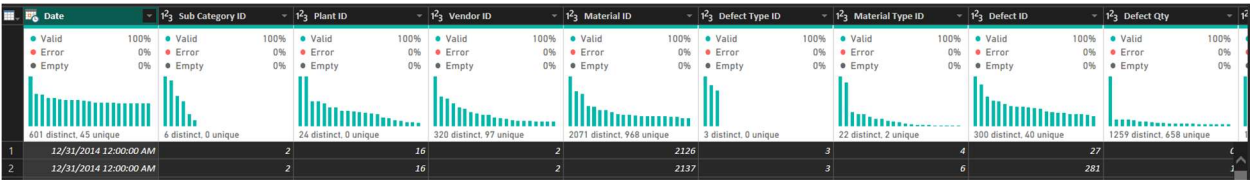
# Defect Report Table:



*Figure 13: Defect Report Table Before*

There are no missing data or errors. There is 6145 record, 193 of which are duplicates. The precious similarity in numbers may suggest a problem with the record entry that caused this entry. Anyways this should be sent back to the data owner but let's assume the presence of the duplicates. After removing them the records will be 5952.

## Changes made:

- Removed duplicate rows.
- Renamed the Sub Category ID to Category ID to fit the change done on the Category Table.

## Data Changes:

- **Defect Type ID column** is changed from **Whole number** to **Text**
- **Category ID column** is changed from **Whole number** to **Text**
- **Material Type ID column** is changed from **Whole number** to **Text**
- **Plant ID column** is changed from **Whole number** to **Text**
- **Vendor ID column** is changed from **Whole number** to **Text**
- **Material ID column** is changed from **Whole number** to **Text**
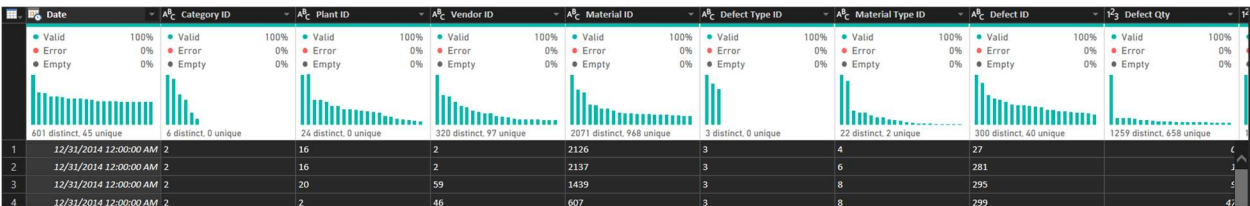- **Defect ID column** is changed from **Whole number** to **Text**



*Figure 14: Defect Report Table After*

# Data Model

## Objective of Data Model

The objective of data modeling is to organize the dataset into a clear and structured format for accurate analysis. This involves linking the **defect report table** with **related tables** such as category, vendor, material, plant, and defect to ensure consistency. **ID formats are standardized,** and a calendar table is created to manage dates properly. These steps help in tracking defects, evaluating supplier. With a well-structured **Star Schema**, the data becomes easier to analyze in Power BI, leading to better decision-making.
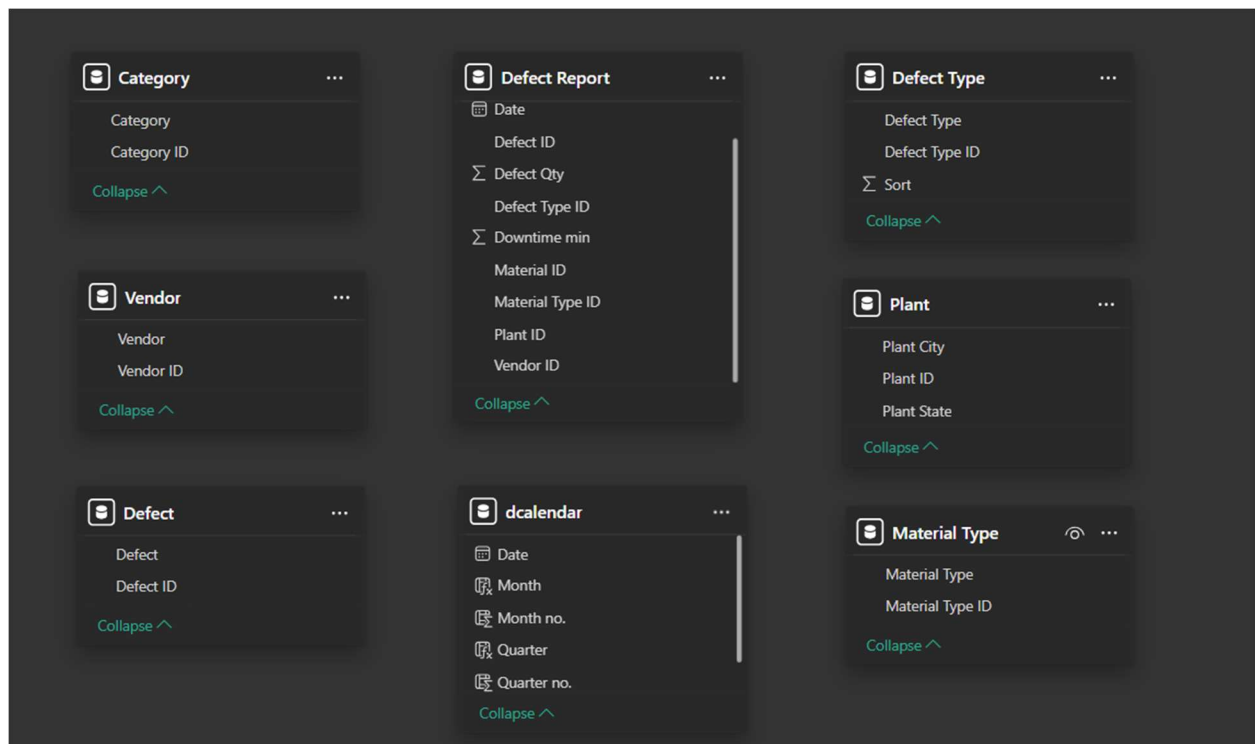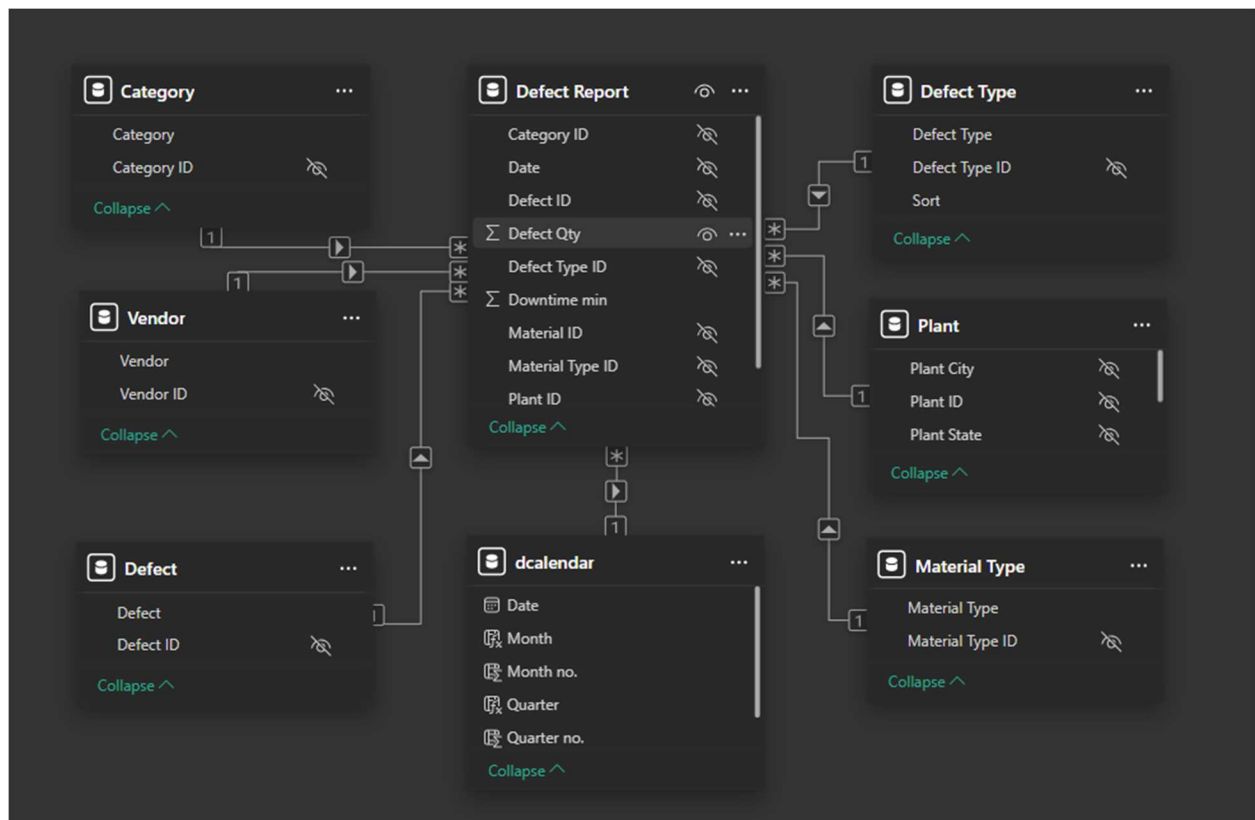
## Model



*Figure 15: Before Data Modeling*

*Figure 16: Data Modelled based on Star Schema*