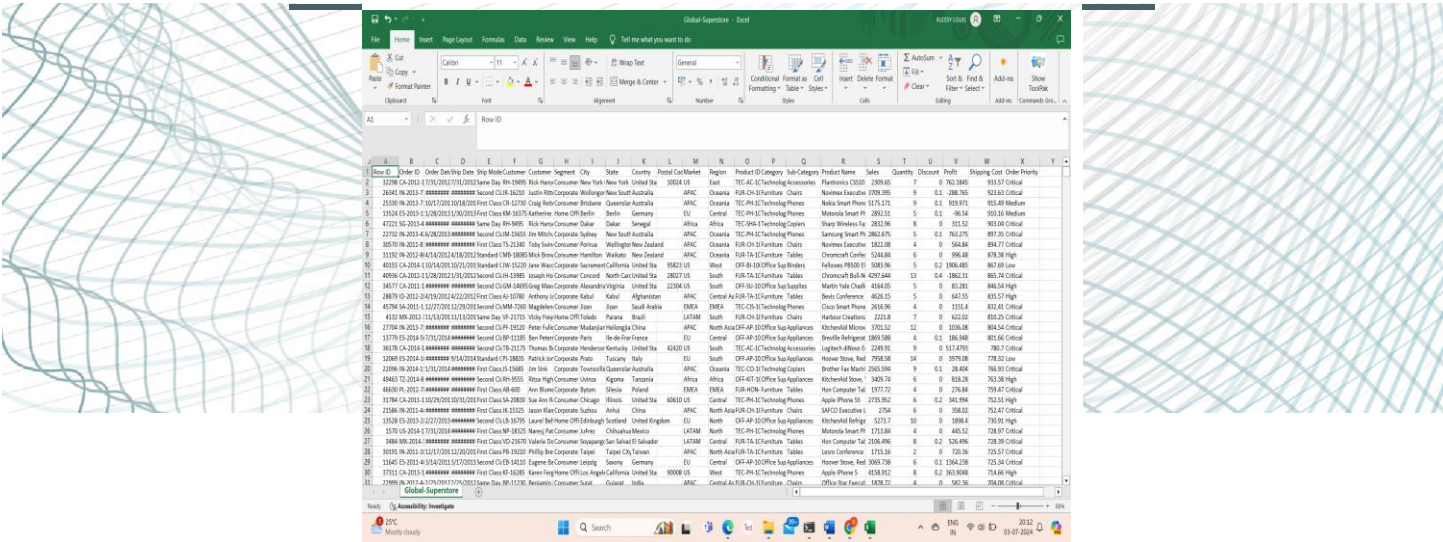# OpenRefine Task Report



**Blessy Louis**

Date

03.07.2024

Blessylouis2002@gmail.com

# INTRODUCTION

OpenRefine is a powerful open-source tool for working with messy data. It allows users to clean, transform, and enrich datasets by providing an intuitive interface for performing a wide variety of data wrangling tasks. OpenRefine supports operations such as data cleaning, clustering, normalization, and transformation. It can handle large datasets and is capable of importing data from various formats, including CSV, Excel, JSON, and XML. With its ability to undo and redo steps, users can experiment with different data cleaning strategies and track changes easily. OpenRefine is particularly useful for data scientists and analysts who need to prepare data for analysis or visualization.
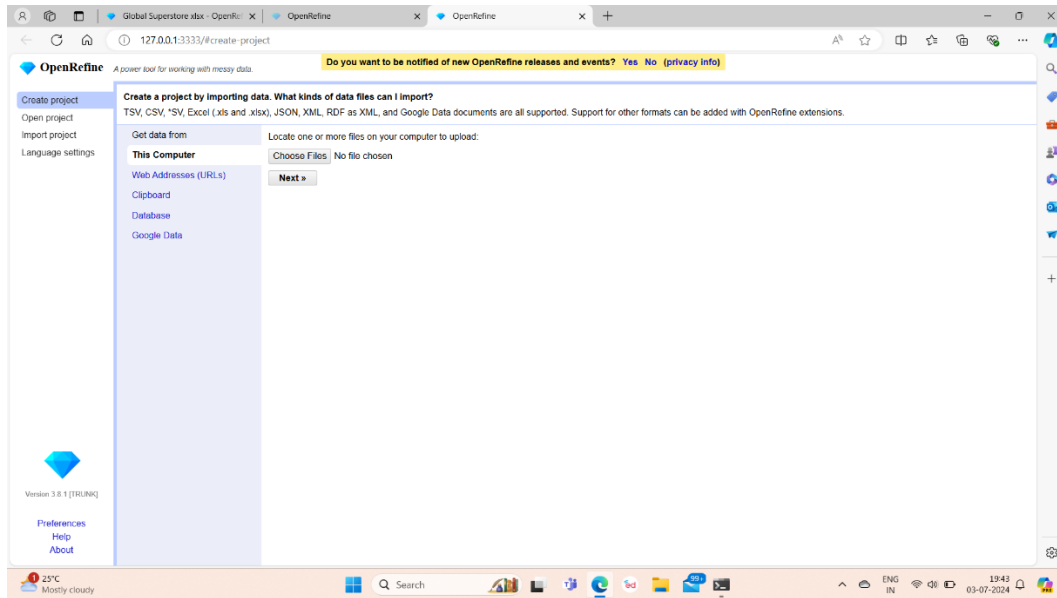
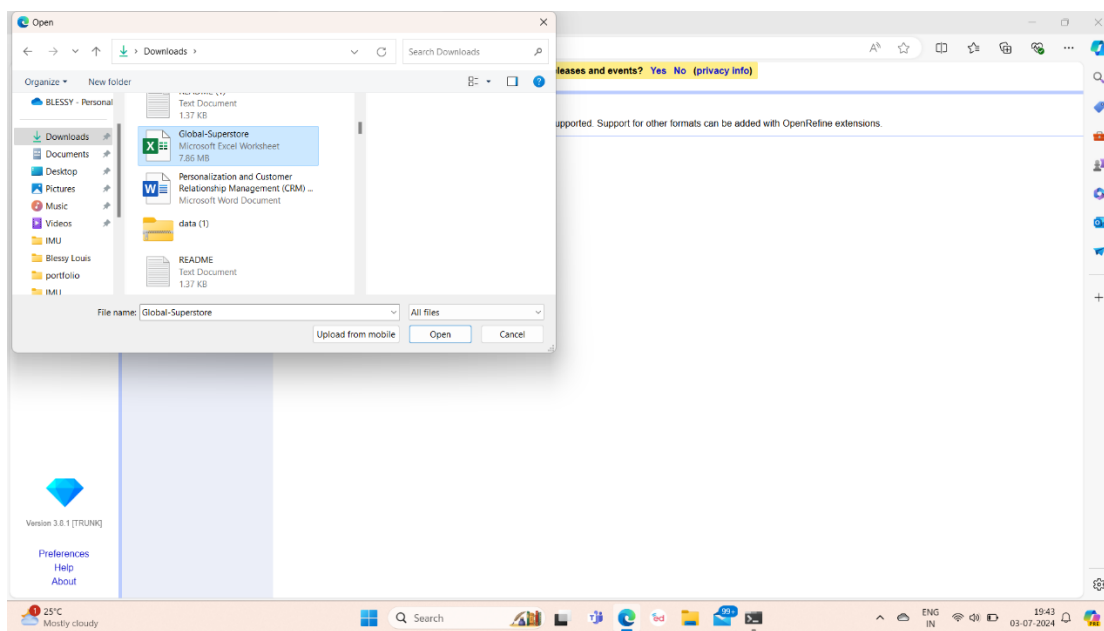<div style="text-align: center;">

# THE DATASET

</div>

## Data Description

This dataset contains detailed transactional information regarding orders, spanning various attributes. Each record is uniquely identified by a Row ID and includes an Order ID to specify individual transactions. The dataset captures the timeline of each order with Order Date and Ship Date, alongside the shipping preferences indicated by the Ship Mode. Customer details are meticulously recorded, including Customer ID, Customer Name, Segment (which likely categorizes the type of customer), City, State, Country, and Postal Code. Additionally, the dataset distinguishes the Market and Region of each transaction. Product details are comprehensive, featuring Product ID, Category, Sub-Category, and Product Name, reflecting a diverse range of items. Financial aspects are captured through fields like Sales, Quantity, Discount, Profit, and Shipping Cost. Finally, the Order Priority field suggests the urgency assigned to each order, potentially influencing the shipping process. This rich dataset is ideal for analyzing sales performance, customer behavior, and logistical efficiency across different geographical and market segments.

## Dataset Upload



After the OpenRefine application is installed, we open the application after extracting the zipped file, we encounter an interface like above.



We select the dataset file from the local files

The dataset is uploaded successfully to the application.

# Data Formatting





Formatting the order date column to date format

Similarly converting ship date to date format.



Performing text facet: