**PROJECT SCREENSHOTS DOCUMENTATION**
**Topic:** Data Cleaning using Power BI & VS Code
**Internship:** Data Analytics / Data Science Internship
**Intern Name:** Harika Vavilapalli
**Tools Used:** Microsoft Power BI (Power Query), Visual Studio Code
**Description:** This document contains relevant screenshots showing the step-by-step data cleaning process performed in Power BI and automated data cleaning performed using Jypter Notebook in VS Code that are covered in both the Weeks (Week 1 & Week 2).

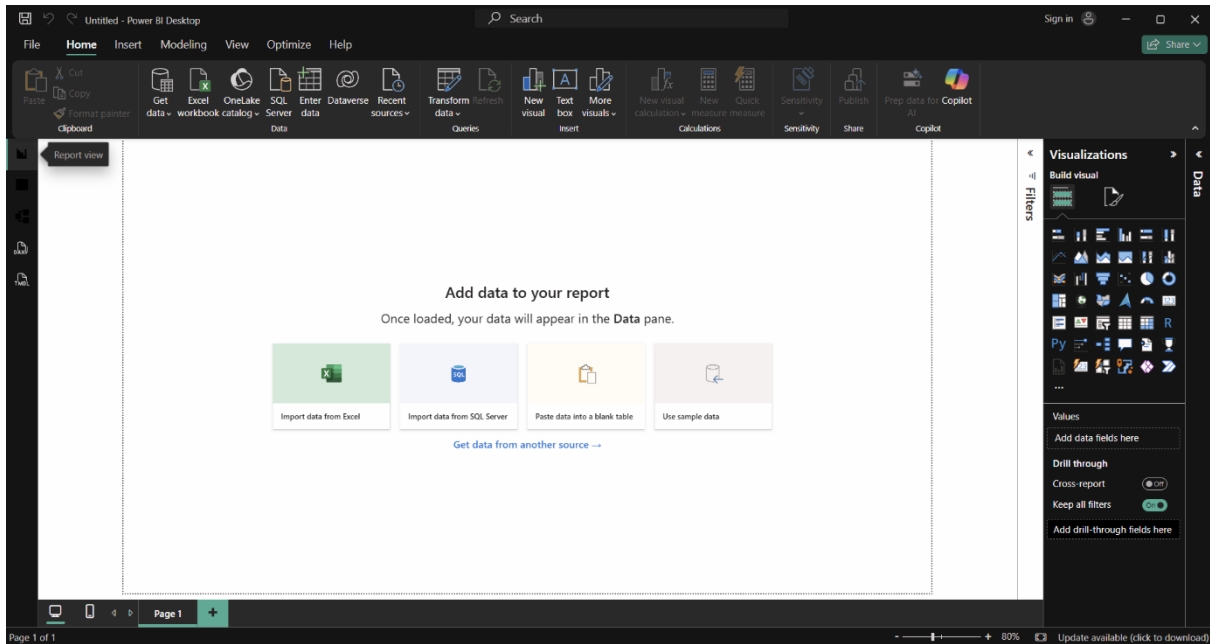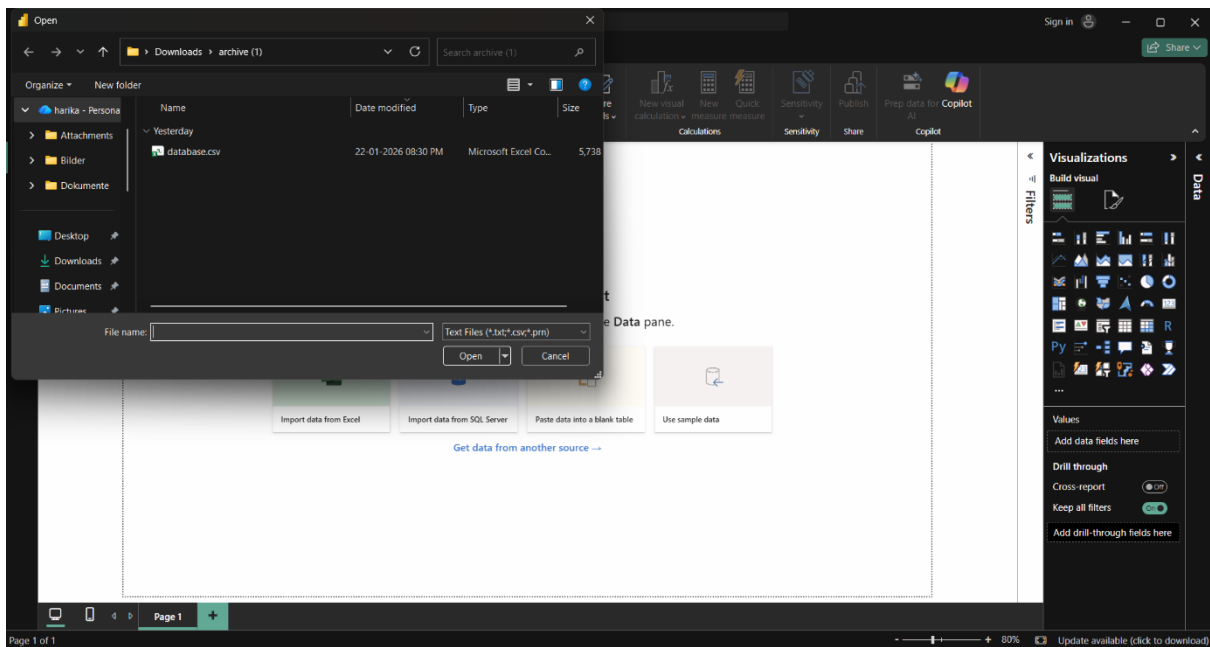# DATA CLEANING IN POWER BI



Fig 1: Power BI Canvas.



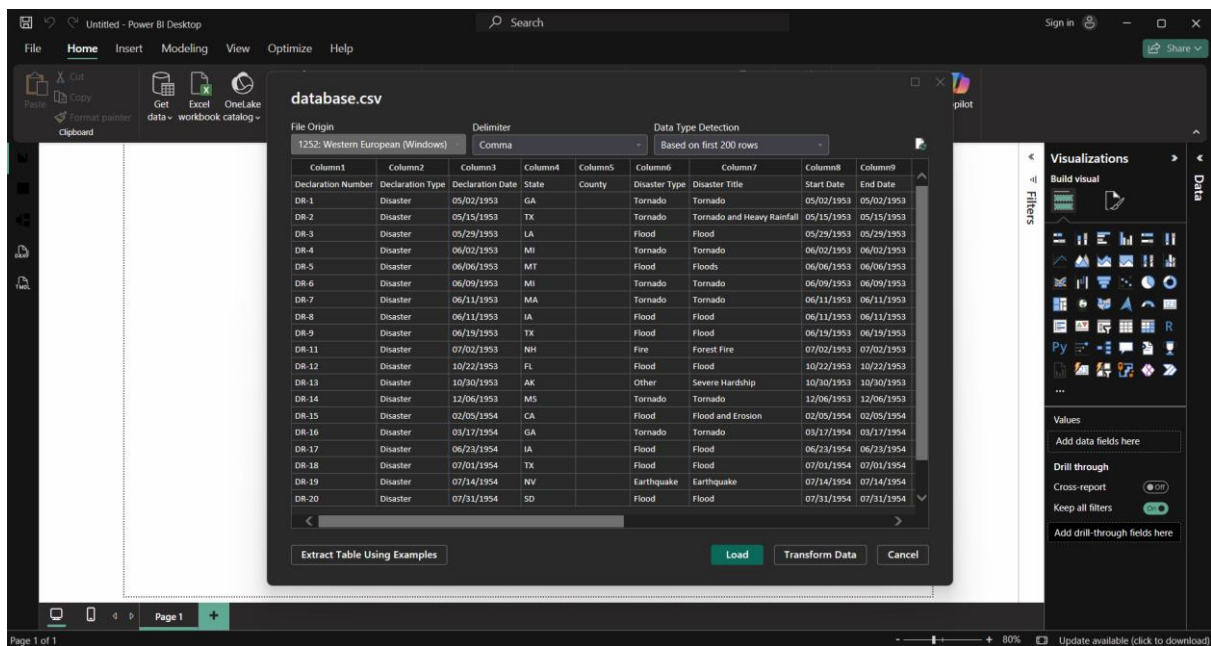Fig 2: Importing Dataset Using **Get Data**.

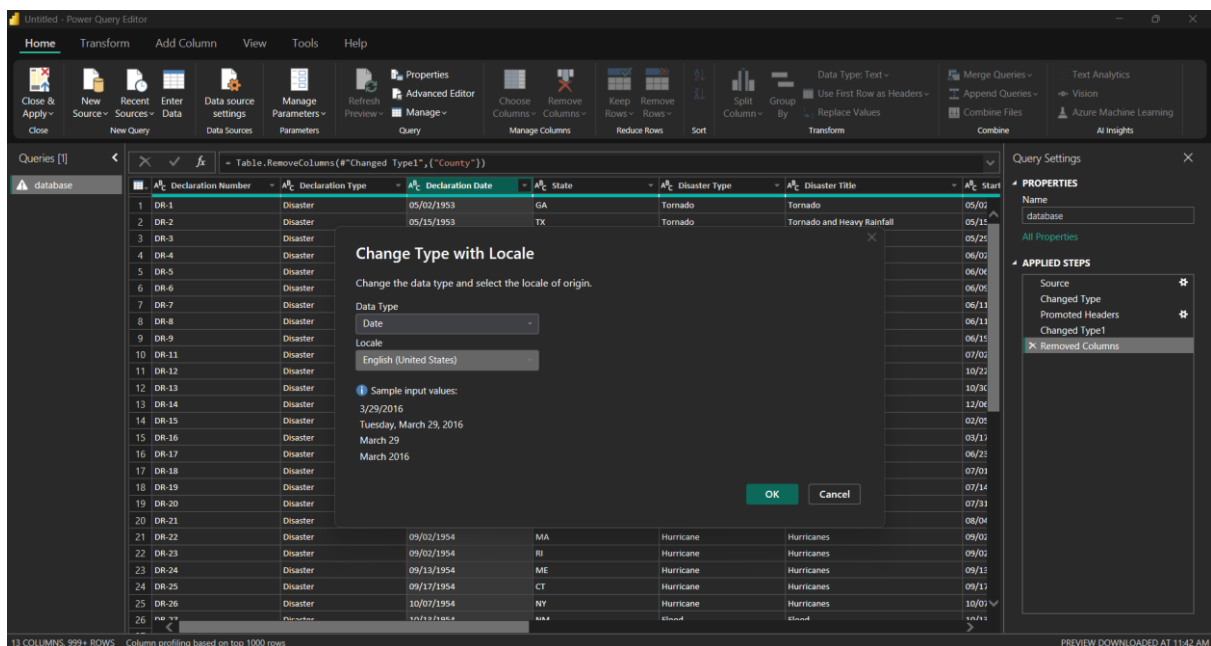Fig 3: Transforming the Data Using **Transform Data**.



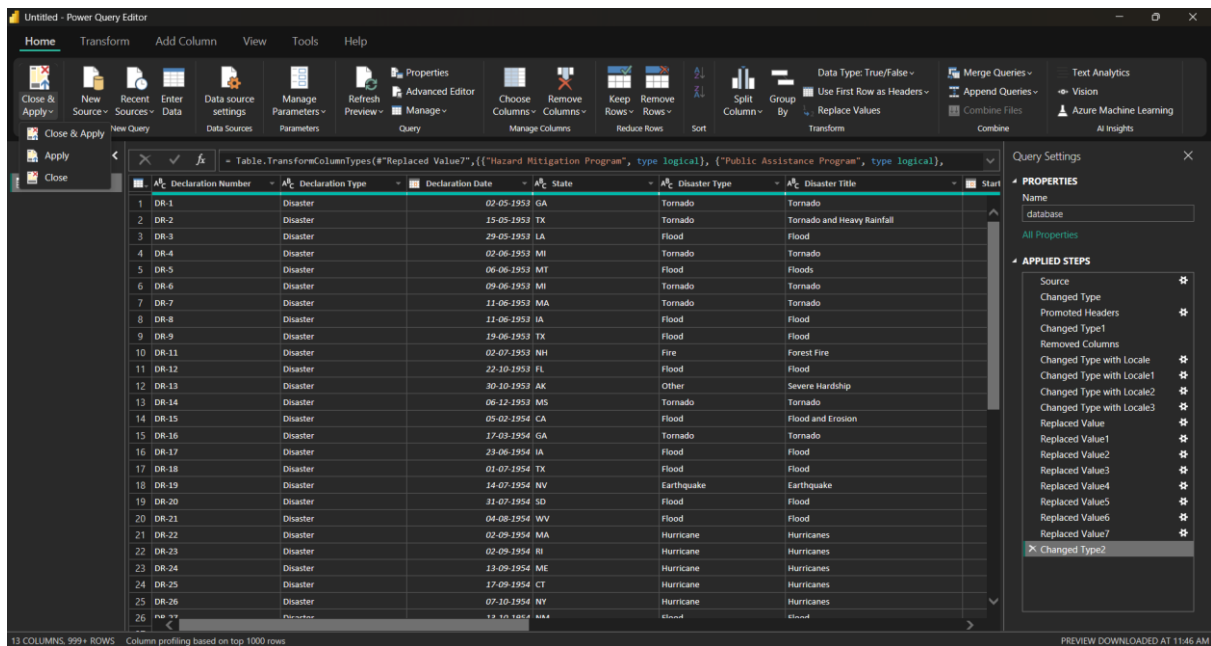Fig 4: Changing The Data Type of the column Using **Locale**.

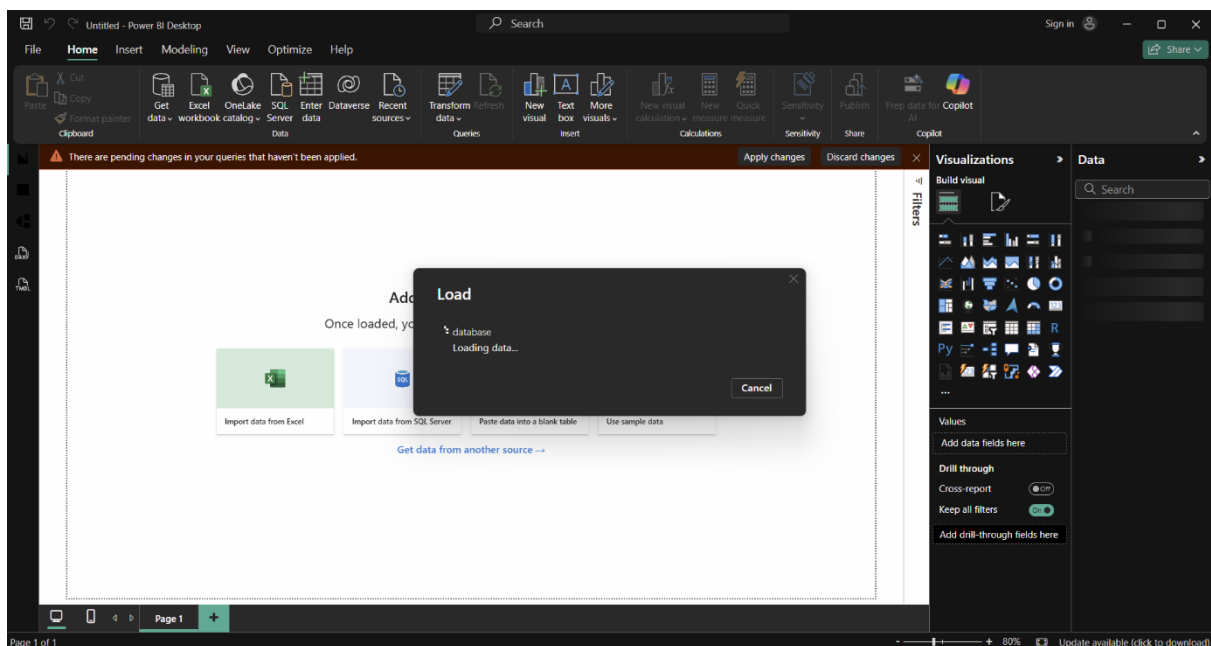Fig 5: Selecting **Close & Apply** after data cleaning.



Fig 6: Power BI automatically Loads the Cleaned Data

Fig 7: Cleaned Data in **Table View**.



Fig 8: Cleaned Data In **Model View**.
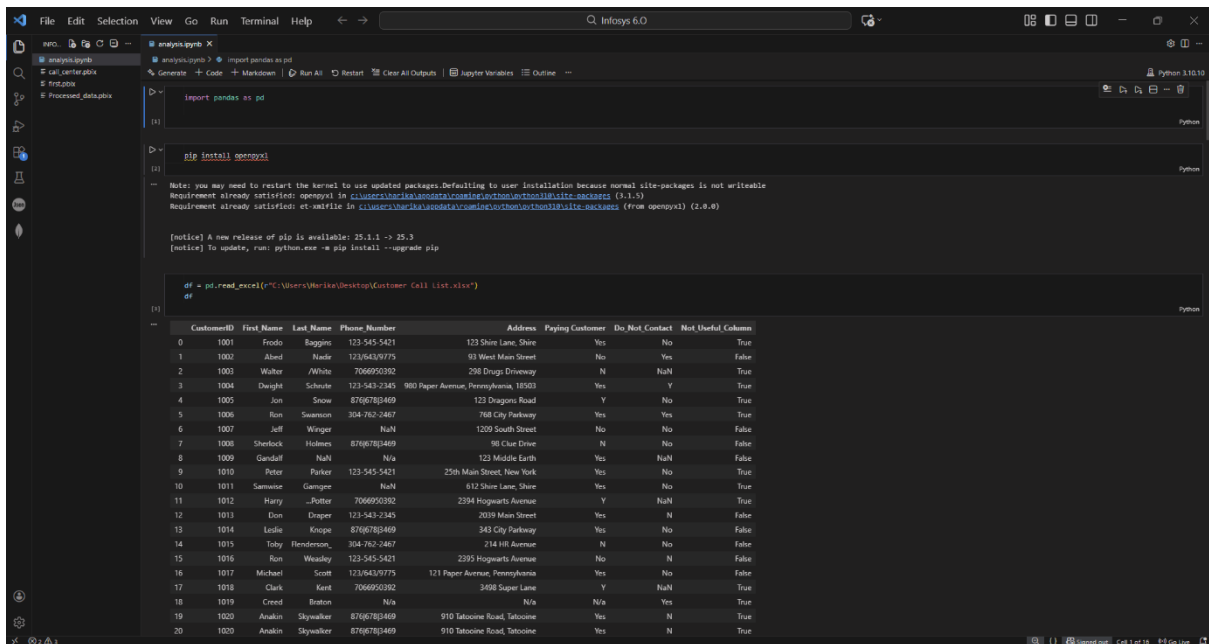
# DATA CLEANING USING VS CODE (JYPTER NOTEBOOK)



Fig 1: Creating a Jypter file using ipynb extension and Loading the Data.
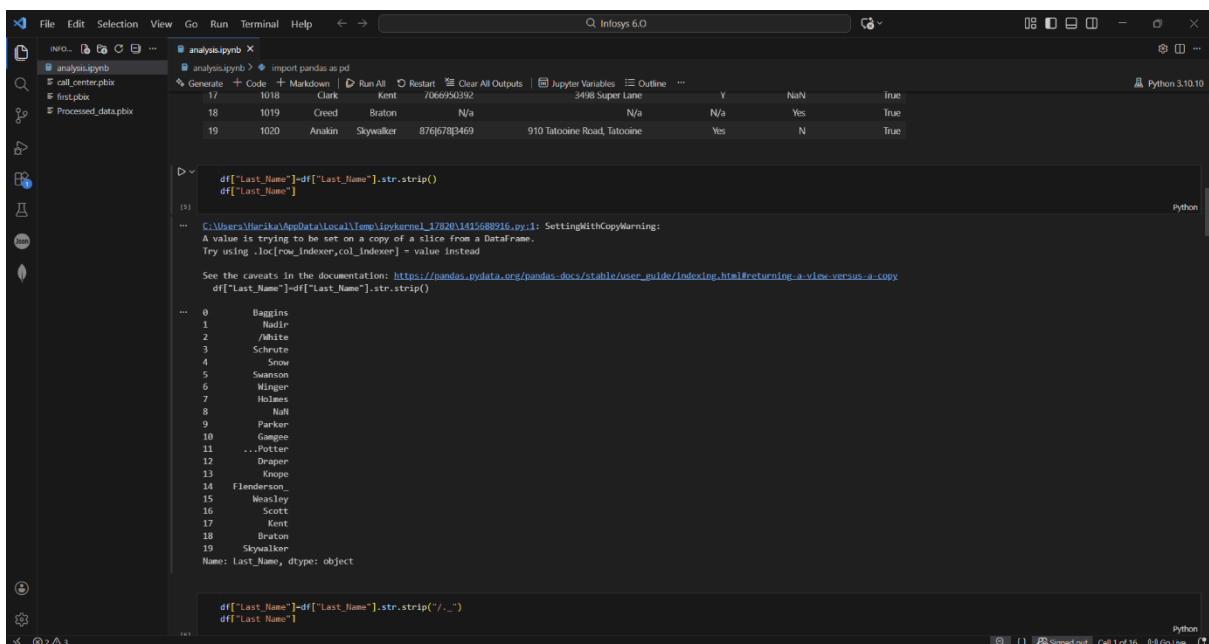


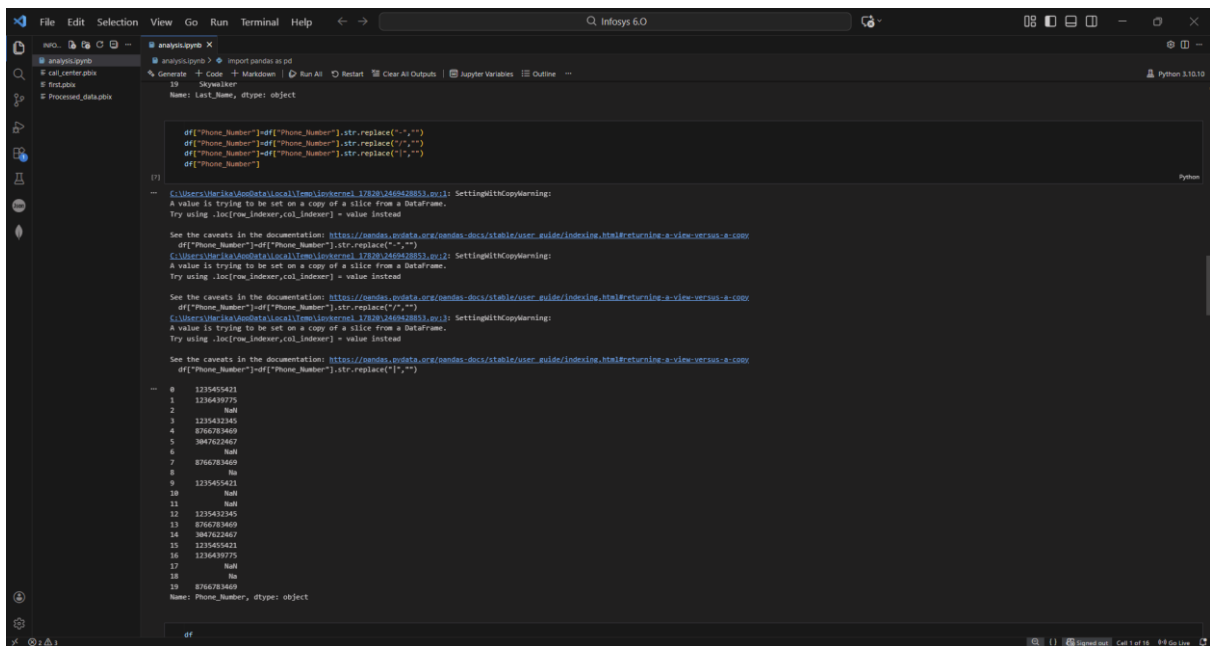Fig 2: Removing the Unwanted Elements in the column.

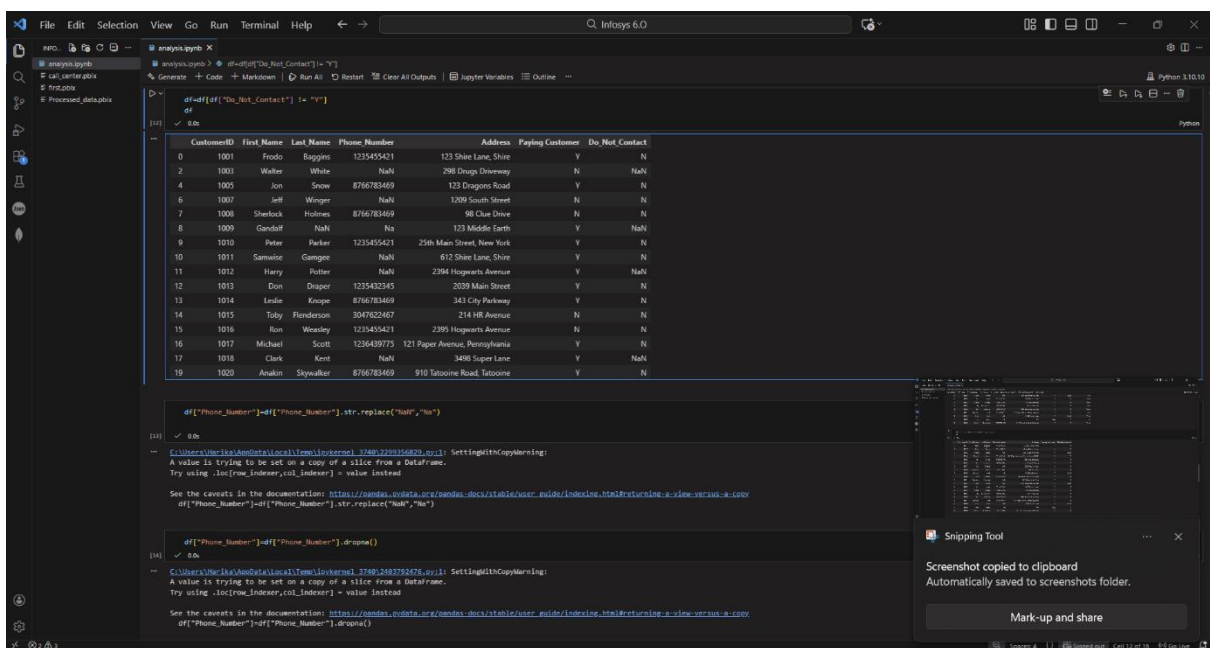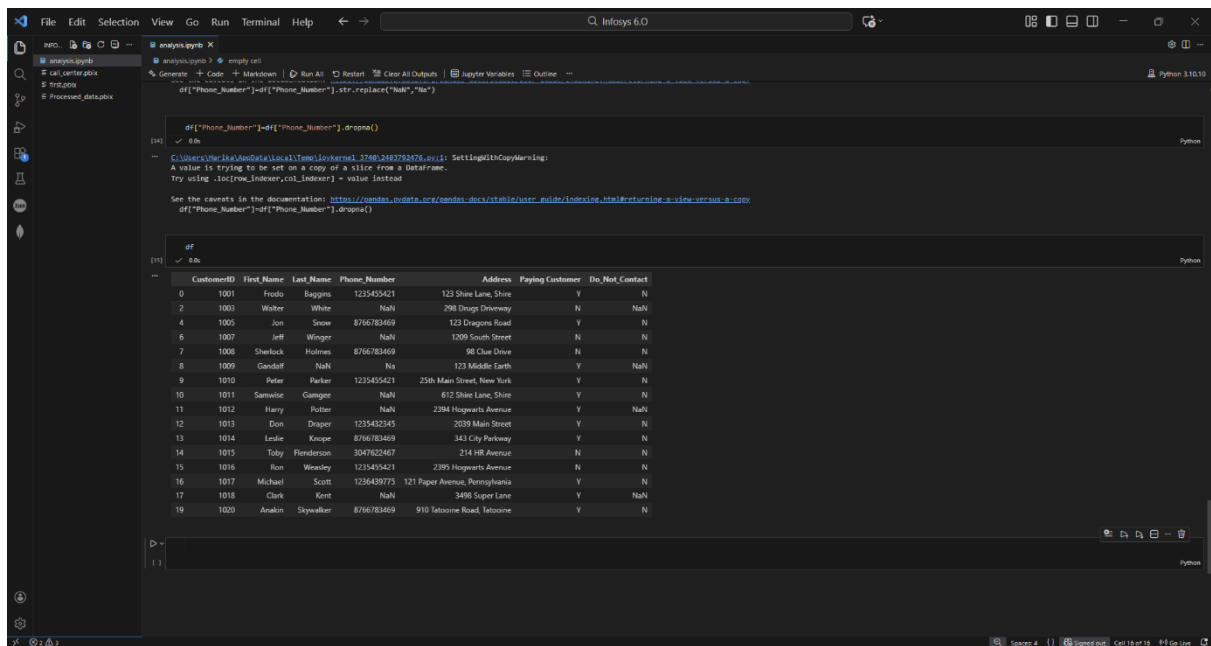Fig 3: Replacing the values with suitable data types.



Fig 4 : Dropping the irrelevant or unused elements in the Data.

Fig 5: Loading the Cleaned Dataset.