

# Data Intake Report

Name: G2M Case Study  
Report date: 04.10.2021  
Internship Batch: LISUM04  
Version: 1.0  
Data intake by: Tolga Yaz  
Data intake reviewer:  
Data storage location:

## Tabular data details:

<b>Name of the File</b>	Cab_Data.csv
<b>Total number of observations</b>	359 392
<b>Total number of features</b>	7
<b>Base format of the file</b>	.csv
<b>Size of the data</b>	20.10 MB

<b>Name of the File</b>	Customer_ID.csv
<b>Total number of observations</b>	49 171
<b>Total number of features</b>	4
<b>Base format of the file</b>	.csv
<b>Size of the data</b>	1.00 MB

<b>Name of the File</b>	Transaction_ID.csv
<b>Total number of observations</b>	440 098
<b>Total number of features</b>	3
<b>Base format of the file</b>	.csv
<b>Size of the data</b>	8.58 MB

<b>Name of the File</b>	City.csv
<b>Total number of observations</b>	20
<b>Total number of features</b>	3
<b>Base format of the file</b>	.csv
<b>Size of the data</b>	4.00 kB

<b>Name of the File</b>	dfMain.csv
<b>Total number of observations</b>	359 392
<b>Total number of features</b>	15
<b>Base format of the file</b>	.csv
<b>Size of the data</b>	37.30 MB

**Proposed Approach:**

- The data was transformed into Python Pandas Dataframes to be processed.
- The dataframes were observed carefully to identify the nulls but there was no null.
- .duplicated() method of Python Pandas Library was used to identify duplicates. No necessary duplicates were found to drop.
- The four .csv files that were provided were turned into one main file (dfMain.csv) after some necessary steps like dropping the unnecessary columns and adding some necessary columns. The changes can be observed in the Jupyter Notebook provided.