

# Data Intake Report

Name: G2M Case Study

Report date: 07/05/2022

Internship Batch: LISUM11

Version: 1.0

Data intake by: Anthony Ghimpu

Data intake reviewer: (Individual Assignment)

Data storage location: Computer Disk

## Tabular data details:

<b>Total number of observations</b>	359,392
<b>Total number of files</b>	1
<b>Total number of features</b>	7
<b>Base format of the file</b>	.csv
<b>Size of the data</b>	20.1 mb

<b>Total number of observations</b>	20
<b>Total number of files</b>	1
<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>Total number of observations</b>	49172
<b>Total number of files</b>	1
<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>Total number of observations</b>	440098
<b>Total number of files</b>	1
<b>Total number of features</b>	3
<b>Base format of the file</b>	.csv
<b>Size of the data</b>	8.58 MB

## Proposed Approach:

- Mention approach of dedup validation (identification)
  - Every observation will be checked for duplication in an iterative manner. If any duplicates arise, they will be dropped from the dataset.
- Mention your assumptions (if you assume any other thing for data quality analysis)
  - There are no outliers within the dataset. (Data is not skewed)
  - All missing data will result in dropping of observations.

