

# Data Intake Report

Name: G2M insight for Cab Investment firm

Report date: 18 June, 2021

Internship Batch: LISUM01

Version: 1.0

Data intake by: Mohammad Shafiqul Islam

Data intake reviewer: Data Glacier

Data storage location: <https://github.com/msishuvo/G2M-insight-for-Cab-Investment-firm>

## Tabular data details:

<b>Total number of observations</b>	344953
<b>Total number of files</b>	4
<b>Total number of features</b>	21
<b>Base format of the file</b>	.csv
<b>Size of the data</b>	67.9 MB

## Cab data details:

<b>Total number of observations</b>	359392
<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>	19.2 MB

## City data details:

<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>	608 B

## Customer\_ID data details:

<b>Total number of observations</b>	49171
<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>	1.5 MB

**Transaction\_ID data details:**

<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>	10.1 MB

**Proposed Approach:**

- Duplicate data were removed by using the drop\_duplicate function from the pandas dataframe.
- The date had to be converted from an integer value to its original format.
- There were extra one-month data that contradicts the given time period. So they were removed.
- Outliers may present in the data, but due to the lack of other information, it is not considered.
- Apart from the original 15 features, some new columns such as profit, year, month, day, age group, weather are created using simple calculation.