

# Data Intake Report

Name: G2M insight for Cab Investment firm

Report date: 14<sup>th</sup>-Aug-23

Internship Batch: LISUM24

Version:1.0

Data intake by: Anusha Asim

Data intake reviewer: N/A

Data storage location: <https://github.com/DataGlacier/DataSets/blob/main/README.md>

## Tabular data details:

File 1: Cab\_Data.csv

<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>	20.10 MB

File 2: City.csv

<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

File 3: Customer\_ID.csv

<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>	1.00 MB

File 4: Transaction\_ID.csv

<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:

- Implement deduplication validation to ensure data accuracy using Pandas `drop_duplicates()` method.
- One-hot encode categorical variables before conducting statistical tests.
- Conduct exploratory data analysis (EDA) to understand data distribution and relationships.
- Perform statistical tests for hypothesis validation.
- Utilize data visualization techniques to present insights effectively.
- Assumptions:
  - The variables in the dataset follow a normal distribution for statistical tests.
  - The assumption of homoscedasticity holds, indicating equal error variances in statistical analyses.
  - Data preprocessing includes removal of missing values and outliers.
  - Dataset is complete and representative of the specified time period.