

# Data Intake Report

Name: G2M Insight for Cab Investment Firm

Report date: July 20, 2025

Internship Batch: LISUM47

Version:<1.0>

Data intake by: Paula McCree-Bailey

Data intake reviewer:<intern who reviewed the report>

Data storage location: [https://github.com/pmb-7684/DataGlacier\\_Internship/tree/main/Module2](https://github.com/pmb-7684/DataGlacier_Internship/tree/main/Module2)

## Tabular data details:

	<b>Cab_Data.csv</b>
<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>	20,633 KB

	<b>City.csv</b>
<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>	1 KB

	<b>Customer_ID.csv</b>
<b>Total number of observations</b>	49,171
<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>	1,027 KB

	<b>Transaction_ID.csv</b>
<b>Total number of observations</b>	440,098
<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>	8,788 KB

	Result_df.csv
<b>Total number of observations</b>	359,392
<b>Total number of files</b>	1
<b>Total number of features</b>	24
<b>Base format of the file</b>	.csv
<b>Size of the data</b>	8,788 KB

### Proposed Approach:

- Dedup validation (identification)  
To check for duplicates in the dataset, ``duplicates = result_df[result_df.duplicated()]`` which scans for rows that are exact duplicates of previous ones. Also used ``duplicates = result_df[result_df.duplicated(subset=key_cols, keep=False)]`` to make sure all duplicates are captured. There are no duplicate observations in our dataset.
- Assumptions:
  - The “Price Change” feature contains a lot of outliers. We are provided with the duration or rate. The decision was to retain the outliers.
  - “User feature represents the number of people taking a cab.
  - “Profit” feature is created as revenue (Price\_charges) – cost of goods (Cost\_of\_trip)