

# Data Intake Report

Name: G2M Insights for Cab Investment Firm

Report date: 22-10-2023

Internship Batch: LISUM26

Version:<1.0>

Data intake by: Mohammad Shehzar Khan

Data intake reviewer:<intern who reviewed the report>

Data storage location: <location URL eg: github, cloud>

## 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>	19.2 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>	608 Bytes

<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>	1.5 MB

<b>Name of the file</b>	Transaction_ID.csv
<b>Total number of observations</b>	440,097
<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:**

- Mention approach of dedup validation (identification)
- Mention your assumptions (if you assume any other thing for data quality analysis)

The approach has been done in the following manner:

- Data Exploration & Understanding
- Feature Engineering
- Data Merging
- Data Visualization
- Data Analysis
- Recommendations

Assumptions:

- Outliers are present in Price\_Charged feature but due to unavailability of trip duration details ,we are not treating this as outlier.
- Profit of rides are calculated keeping other factors constant and only Price\_Charged and Cost\_of\_Trip features used to calculate profit.
- Users feature of city dataset is treated as number of cab users in the city. We have assumed that this can be other cab users as well(including Yellow and Pink cab).
- When we looked at the Cost of the Trip for each ride, there were some rides that were unusually costly as compared to all other rides. Possible reasons for such expensive rides could be either Premium rides, or Overnight rides, or Long distance rides. We have called these people as premium customers.
- We have assumed that a trip upto 8km is short distance trip, a trip from 8km to 30km is Medium distance trip, and a trip above 30km is a long distance trip.
- We have assumed that any customer who uses the same cab service atleast every 4 months on average is a loyal customer to that cab service. Here, we have data for 36 months approx., so any customer who has 9 rides with a particular company is a loyal customer to that company.