

# Data Intake Report

Name: G2M Insight for Cab Investment

Internship Batch: LISUM10: 30

Version: 1.0

Data intake by:

Fatemeh Bagheri

Data intake reviewer:

Data storage location: <https://github.com/DataGlacier/DataSets>

## Tabular data details:

Cab\_Data.csv

<b>Total number of observations</b>	359392
<b>Total number of files</b>	
<b>Total number of features</b>	7
<b>Base format of the file</b>	.csv
<b>Size of the data</b>	20.2 MB

City.csv

<b>Total number of observations</b>	20
<b>Total number of files</b>	
<b>Total number of features</b>	3
<b>Base format of the file</b>	.csv
<b>Size of the data</b>	759 B

Customer\_ID.csv

<b>Total number of observations</b>	49171
<b>Total number of files</b>	
<b>Total number of features</b>	4
<b>Base format of the file</b>	.csv
<b>Size of the data</b>	1 MB

Transaction\_ID.csv

<b>Total number of observations</b>	440098
<b>Total number of files</b>	
<b>Total number of features</b>	3
<b>Base format of the file</b>	.csv
<b>Size of the data</b>	20.2 MB

**Note: Replicate same table with file name if you have more than one file.**

**Proposed Approach:**

- Load all data as Pandas dataframe (python or Jupiter)
- This data assumes there are no duplicates in Transaction ID (Cab\_Data and Transaction\_ID), and in Customer ID (Customer\_Data), City (City), so index these columns when loading each dataframe
- Cut the same rows and columns (include all the necessary information)
- Cab\_ID and Transaction\_ID merged into one Pandas dataframe using the join function
- Added additional columns to Cab\_Data: Profit Generated, Trips, Yellow Cab Trips, Pink Cab Trips, Share of Yellow Cab Trips, Share of Pink Cab Trips – to make summaries easier to access
- Added additional columns to Customer\_ID: Age Group, Income Bracket – for making a breakdown of users by demographic characteristic
- Correct some rows with nan or unreadable text
- Used the groupby function to show summaries of data, sometimes merging with other frames, such as breakdown of users by city, and the City dataframe
- Used Pandas Libraries to make several plots including, time series and k means