

# Data Intake Report

Name: G2M insight for Cab Investment firm

Report date: 20<sup>th</sup> September 2022

Internship Batch: LISUM13: 30

Version: 1.0

Data intake by: Rania Tarek Fleifel

Data intake reviewer:

Data storage location: [G2M Cab datasets](#)

## Tabular data details:

### 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>	20.1 MB
<b>Unique identifier feature</b>	Transaction ID
<b>Dupe validation</b>	No duplicates

### 2) 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 MB
<b>Unique identifier feature</b>	Customer ID
<b>Dupe validation</b>	No duplicates

### 3) 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
<b>Unique identifier feature</b>	Transaction ID
<b>Dupe validation</b>	No duplicates

### 4) 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>	759 bytes

<b>Unique identifier feature</b>	City
<b>Dupe validation</b>	No duplicates

*Variables after data cleaning and new features are: (Total:19)*

- 'trans\_id'--> ride's unique identifier
- 'company'--> yellow\_cab or pink\_cab
- 'city'--> ride's location
- 'kms'--> ride's distance covered in kms
- 'c\_paytype'--> ride's payment method cash or credit
- 'costumer\_id' --> rider's unique identifier
- 'gender'--> rider's gender
- 'c\_age'--> age of rider
- 'c\_income'--> income USD/month of rider
- 'population'--> city's population
- 'users'--> city's cab users
- 'traveldate'--> ride's date %d-%M-%Y
- 'cost\_km'--> ride's 'cost/kms'
- 'charged\_km'--> ride's 'charged/kms'
- 'weekendnot'--> whether date of ride is a weekend (Sat,Sun) or not
- 'holidaynot'--> whether date of ride is a US holiday (\*\*<https://www.timeanddate.com/holidays/us/?hol=25>) or not
- 'agegrp'--> rider's age group (youth or adult or older)
- 'year'--> extracted from traveldate
- 'month'--> extracted from traveldate

### **Deduplication Approach:**

- For each file, find the number of unique values for each feature, the feature with the maximum number of unique features is the “Unique identifier” of the file.
- **If** the unique identifier feature has number of unique values equal the total number of observations, there’s no duplicated observations (entirely repeated rows).
- **Else**, keep only one of the duplicated rows & drop the rest.

### **Assumptions:**

- If a feature doesn’t have neither zeros or empty values, there’s no missing data
- If a feature has values of different types, there’s faulty data
- Merge transaction\_ID.csv and costumer\_ID.csv on “Customer\_ID” feature.
- Data where the company choice is not provided is not helpful in making final decision(e.g. costumer IDs that don’t have any transaction IDs within our desired time frame)
- Data provided is for a time interval more than the interval I am looking into for decision making
- The market un-utilized by both companies is up-for-grabs (i.e. there are no other competitors)
- The cost of a ride is based on kms only, the rest of the charged amount is the profit per ride
- XYZ is interested in the company with biggest growth potential, not current market share.