

# Data Intake Report

Name: G2M insight for Cab Investment firm

Report date: 12th June 2021

Internship Batch: LISUM01

Version:1.0

Data intake by: Samuel Alejandro Cueva Lozano

Data intake reviewer: Samuel Alejandro Cueva Lozano

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

## Tabular data details:

**Name of the file:** Cab\_Data.csv

<b>Total number of observations</b>	359393
<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>	21.2MB

**Name of the file:** Customer\_ID.csv

<b>Total number of observations</b>	49172
<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>	1.1MB

**Name of the file:** Transaction\_ID.csv

<b>Total number of observations</b>	440099
<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>	9.0MB

**Name of the file:** City.csv

<b>Total number of observations</b>	21
<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

## Proposed Approach:

### Duplicate values

To consider two or more records as duplicates these must have the same values in the most of the attributes (excluding unique identifier) or very similar values in the case that they are strings (e.g. 'Samuel Cueva Lozano' = 'Samy Cueva Loz' and the others attributes be the same).

A record in the data set will be considered duplicate if there is a transaction with the same kilometers traveled, the same city, the same company and the same day.

### Join tables

The table *Cab\_Data.csv* will be joined with the table *Customer\_ID.csv* by the table *Transaction\_ID.csv* since the latter has the *Transaction ID* (from *Cab\_Data.csv*) and Customer ID (from *Customer\_ID.csv*) attributes, then the table *City.csv* will also join them.

### Null values

The *Empty values*, *NaN values* and values like "?" will be treated as **Null values**

Fields with **too many null values** would be removed or filled according to the following:

- Fields like *Age*, *Gender*, *Income*, *Company* and *Users* would be filled because they have missing values that depends on their hypothetical values or depend on other attributes.
- The rest of fields would be removed because they have missing values due to bad configuration, issues with data collection, or untraceable random reasons.

After cleaning at the field level, rows with null values would be removed.

### Outlier detection

Outliers will be removed from numerical fields so that they don't negatively affect the analysis.

A fixed threshold would be used for the *Age*, *Population* and *Users* attribute to avoid inconsistent data, then the IQR Score will be used to filter out the outliers in all attributes.

### Data Transformation

The field *Date of Travel* from the table *Cab\_Data.csv* could be transformed to a more readable format for better understanding.