

# Data Intake Report

Name: G2M insight for Cab Investment firm

Report date: 11/11/2024

Internship Batch: LISUM39

Version:1.0

Data intake by: Vinicius Brun

Data intake reviewer:

Data storage location: <https://github.com/vinbrun/data-glacier-virtual-internship/tree/main/week-2/DataSets>

## Tabular data details: Cab\_Data.csv

Total number of observations	359393
Total number of files	1
Total number of features	7
Base format of the file	.csv
Size of the data	20.1 MB

## Tabular data details: City.csv

Total number of observations	20
Total number of files	1
Total number of features	3
Base format of the file	.csv
Size of the data	0.8 KB

## Tabular data details: Customer\_ID.csv

Total number of observations	49171
Total number of files	1
Total number of features	4
Base format of the file	.csv
Size of the data	1.0 MB

## Tabular data details: Transaction\_ID.csv

Total number of observations	440098
Total number of files	1
Total number of features	3
Base format of the file	.csv
Size of the data	8.6 MB

**Proposed Approach:**

- There appears to be no missing values in the existing columns.
- Hypothesis: investigate why there are more rows in Transaction IDs than in Cab\_Data.
- Deduplication approach:
  - Identify exact duplicates based on the column Transaction\_ID in the file Transaction\_ID.csv (unique identifier).
  - Identify exact duplicates based on the column Customer\_ID in the file Customer\_ID.csv (unique identifier).
  - Identify exact duplicates based on the column Transaction\_ID in the file Cab\_Data.csv.
- Data validation approach:
  - Join Cab\_Data.csv with Transaction\_ID.csv on Transaction\_ID to ensure that each trip in Cab\_Data.csv has a corresponding transaction in Transaction\_ID.csv.
  - Join Transaction\_ID.csv with Customer\_ID.csv on Customer\_ID to ensure that transaction in Transaction\_ID.csv has a corresponding customer in Customer\_ID.csv