

# Data Intake Report

Name: G2M: Insight for Cab Investment Firm

Report date: 07/14/2023

Internship Batch: LIMSUM23: 30

Version: 1.0

Data intake by: Susan Zhang

Data intake reviewer: N/A

Data storage location: [https://github.com/20szha/VC/blob/new-branch/Week%202/G2M\\_Insight\\_for\\_Cab\\_Investment\\_Firm.ipynb](https://github.com/20szha/VC/blob/new-branch/Week%202/G2M_Insight_for_Cab_Investment_Firm.ipynb)

## Tabular data details:

### Customer ID.csv

<b>Total number of observations</b>	49,171
<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>	1,027 KB

### Cab Data.csv

<b>Total number of observations</b>	359,392
<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>	20,663 KB

### City Data.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>	1 KB

### Transaction ID.csv

<b>Total number of observations</b>	440,098
<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>	8,788 KB

---

### After merging into Master Data.csv

<b>Total number of observations</b>	359,392
<b>Total number of files</b>	1 (merged from the top 4 files)

<b>Total number of features</b>	14
<b>Base format of the file</b>	.csv
<b>Size of the data</b>	47,439 KB

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

### **Proposed Approach:**

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

The following steps were taken to explore the data:

1. Read in all given .csv files
2. Cleaned up data if necessary
  - a. Converted excel time serial number into dates
3. Merge data into single master dataframe using the .merge() function
  - a. Merged on “City”, “Transaction ID”, and “Customer ID” to avoid duplicate records
  - b. Filled NA’s with ‘’
    - i. Original dataset did not have any anyways and those with no City, Customer ID, or Transaction ID matches should’ve been discarded
4. Explore the data and evaluate for better company to invest in
  - a. For example: sns.pairplot() function and .corr() function to visualize the relationship between each feature and the .groupby() function to compare the various features against one another

**Conclusion:** It appears that the Yellow Cab would be a better investment.