

# Data Intake Report

Name: G2M Insight for Cab Investment Firm

Report date: October 14<sup>th</sup>, 2022

Internship Batch: LISUM14

Version: 1.0

Data intake by: Nolan Piloza-Hibbit

Data intake reviewer: Nolan Piloza-Hibbit

Data storage location: [https://github.com/ngpiloza/DG\\_Internship\\_WK2](https://github.com/ngpiloza/DG_Internship_WK2)

## Cab data details:

<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 details:

<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

## Customer\_ID data details:

<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

## Transaction\_ID data details:

<b>Total number of observations</b>	44,098
<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

**Proposed Approach:**

- Duplicate removal can be accomplished using basic Excel functions and built-in tools. The chosen method will be the **Data > Remove Duplicates** tool.
- Datetime changes can be made using pandas built-in functions to display the proper date.
- Some data types will need to be changed to work with numeric only functions (i.e. changing 3,450 str to an int).