

# Data Intake Report

Name: <G2M insight for Cab Investment firm>

Report date: <10/11/2022>

Internship Batch:<LISUM14>

Version:<1.0>

Data intake by:<Zhuoming Yang>

Data intake reviewer:<>

Data storage location: <GitHub >

## Tabular data details:

<b>Total number of observations</b>	<358,392 of rows>
<b>Total number of files</b>	<5>
<b>Total number of features</b>	<60>
<b>Base format of the file</b>	<.csv >
<b>Size of the data</b>	<102.7MB>

## Proposed Approach:

- **Data: Merged\_data.csv, Cab\_Data.csv, City.csv, Customer\_ID.csv, Transaction\_ID.csv, StormEvent\_2016, StormEvent\_2017, StormEvent\_2018:**
- **Cab\_Data.csv** – this file includes details of transaction for 2 cab companies.
- **Customer\_ID.csv** – this is a mapping table that contains a unique identifier which links the customer's demographic details.
- **Transaction\_ID.csv** – this is a mapping table that contains transaction to customer mapping and payment mode.
- **City.csv** – this file contains list of US cities, their population and number of cab users.
- **StormEvent\_2016\_2017\_2019** – this file contains records of yearly storm information US
  - Customer\_data['customer\_id'] = transaction\_data['customer\_id']
  - Transaction\_data['transactoin\_id']= cab\_data['transaction\_id']
  - Cab\_data['city'] = city\_data['city']
  - Join every data together by their foreign key -> Merged\_data
  - **Important Feature:**
    - **Location:** State, City, Population
    - **Customer information:** Income(usd/month), gender
    - **Trip information:** KM\_travelled, Price\_Charged, Profit, Loss
- **Assumptions:**
  - Null Values
  - Duplicates
  - Date Format need to be converted into a better format for analyzing
  - Data Types need to be change