

# Data Intake Report

Name: Cab Investment Firm Analysis  
Report date: 03/15/2025  
Internship Batch: LISUM43  
Version: 1.0  
Data intake by: Duong Duc  
Data intake reviewer:  
Data storage location: GitHub Repository

## Tabular data details:

### Cab\_Data.csv

<b>Total number of observations</b>	359392
<b>Total number of files</b>	4
<b>Total number of features</b>	7
<b>Base format of the file</b>	.csv
<b>Size of the data</b>	59.92 MB

### City.csv

<b>Total number of observations</b>	20
<b>Total number of files</b>	4
<b>Total number of features</b>	3
<b>Base format of the file</b>	.csv
<b>Size of the data</b>	0.00MB

### Customer\_ID.csv

<b>Total number of observations</b>	49171
<b>Total number of files</b>	4
<b>Total number of features</b>	4
<b>Base format of the file</b>	.csv
<b>Size of the data</b>	4.03 MB

### Transaction\_ID.csv

<b>Total number of observations</b>	440098
<b>Total number of files</b>	4
<b>Total number of features</b>	3
<b>Base format of the file</b>	.csv
<b>Size of the data</b>	32.32 MB

**Proposed Approach:**

- Deduplication and Validation:
  - Duplicates will be identified and removed based on Transaction ID and Customer ID.
  - Columns with missing values will be removed.
  - Outliers in Price\_Charged and Cost\_of\_Trip will be flagged for further reviewed.
- Assumptions:
  - Profit is calculated only using Priced\_Charged and Cost\_of\_Trip to calculate profit.
  - Because there is no data on trip duration, we cannot pinpoint outliers in the Price\_Charged feature.

**Note: Convert this doc in pdf and provide the link of pdf file in your dashboard.  
Please do not forget to remove this section while converting the file into pdf.**