

# Data Intake Report

Name: Cab Industry Analysis

Report date:13/02/2024

Internship Batch:LISUM30

Version:<1.0>

Data intake by:Omar Mahmoud Hamdan

Data intake reviewer:Data Glacier

Data storage location: [omarhamdaan/Internship\\_EDA\\_W2 \(github.com\)](https://github.com/omarhamdaan/Internship_EDA_W2)

## Tabular data details:

Transaction\_ID

<b>Total number of observations</b>	44098
<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>	8.5MB

## Tabular data details:

Customer\_ID

<b>Total number of observations</b>	49171
<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>	1KB

#### Tabular data details:

City

<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>	20MB

#### Tabular data details:

Cab\_Data

<b>Total number of observations</b>	359392
<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>	1027KB

#### **Deduplication Validation (Identification):**

1. Utilized Python's pandas library to identify and remove duplicate entries based on key identifiers such as Transaction ID and Customer ID.
2. Assured unique entries for accurate analysis.

#### **Assumptions for Data Quality Analysis:**

1. Assumed that each entry in the Transaction and Cab Data files represents a unique ride.
2. The 'Date of Travel' fields were considered accurate and were converted to a datetime format for time series analysis.
3. Population and Users data in the City.csv file were assumed to be the latest and relevant to the period of the cab data.