

# Data Intake Report

Name: G2M insight for Cab Investment firm

Report date: 16.10.2021

Internship Batch: LISUM04

Version: 1

Data intake by: A.M. Nimasha Chathurangani Attanayake

Data intake reviewer: J.M. Tharindu Chinthaka Jayaweera

Data storage location: <https://github.com/DataGlacier/DataSets>

<https://www.kaggle.com/donnetew/us-holiday-dates-2004-2021>

## 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>	20.1 MB

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>	1 KB

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>	1 MB

transaction\_id

<b>Total number of observations</b>	440098
<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>	8.58 MB

holiday\_data

<b>Total number of observations</b>	342
<b>Total number of files</b>	1
<b>Total number of features</b>	6
<b>Base format of the file</b>	csv
<b>Size of the data</b>	15.3 KB

**Proposed Approach:**

- Mention approach of dedup validation (identification)

Since the transaction IDs are unique to each and every trip, duplications can be checked based on transaction ID.

- Mention your assumptions (if you assume any other thing for data quality analysis)

According to the plotted graphs, there are outliers in the Price Charged feature. But since we do not have enough information on the components that made the Price Charged, it is not appropriate to treat it as an outlier.