

# Data Intake Report

Name: G2M insight for Cab Investment firm

Report date: 04.07.2022

Internship Batch:LISUM 10:30

Version: 1.0

Data intake by:Prajesh Tejani

Data intake reviewer:<intern who reviewed the report>

Data storage location: <https://www.kaggle.com/datasets/gopalkalpnde/nyc-tlc-data>

## Tabular data details:

<b>Total number of observations</b>	7821717
<b>Total number of files</b>	1
<b>Total number of features</b>	19
<b>Base format of the file</b>	yellow_tripdata_2015-01.csv
<b>Size of the data</b>	2.0 GB

## Proposed Approach:

- Firstly I Read data using different method such as dask,modin and pandas
- Validate column(Remove White\_space and convert name in small Letters)
- Validation with YAML file
- convert data into .txt.gz format
- Location:<https://drive.google.com/file/d/1bgywre2DHTTrVyPDhN28oSIVKDOqcKH5/view?usp=sharing>