

# Data Intake Report

Name: EDA Cab Companies

Report date: 6<sup>th</sup> of March 2023

Internship Batch: LISUM19

Version: N/A

Data intake by: Mario Rodrigues Peres

Data intake reviewer: N/A

Data storage location: GitHub

## Tabular data details:

### Dataset1 - 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>	21.2 MB		
<b>Data type (unique values) of features</b>	Transaction ID	int64	(359392)
	Date of Travel	int64	(1095)
	Company	object	(2)
	City	object	(19)
	KM Travelled	float64	(874)
	Price Charged	float64	(99176)
	Cost of Trip	float64	(16291)

### Dataset2 - 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>	759 B		

<b>Data type (unique values) of features</b>	City	object	(20)
	Population	object	(20)
	Users	object	(20)

### Dataset3 - 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>	1.1 MB		
<b>Data type (unique values) of features</b>	Customer ID	int64	(49171)
	Gender	object	(2)
	Age	int64	(48)
	Income (USD/Month)	int64	(23341)

### Dataset4 - 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>	9 MB		
<b>Data type (unique values) of features</b>	Transaction ID	int64	(440098)
	Customer ID	int64	(49171)
	Payment_Mode	object	(2)

### Proposed Approach:

- There is no duplicated registers in any of the provided tables (checked using duplicated function of pandas library in python). Tables 1, 3 and 4 have unique IDs, whereas table 2 presents unique cities. Therefore, all registers are unique in all datasets.
- Data was merged into a single table and, basically, descriptive statistics was used to assess the dataset.