

# Data Intake Report

**Name:** Data Analytics for Investment Decision- Assignment 2  
**Report date:** 02/03/2021  
**Internship Batch:** LISP01  
**Version:** 1.0  
**Data intake by:** Ajaegbu Ebuka Emmanuel  
**Data intake reviewer:** <intern who reviewed the report>  
**Data storage location:** <https://github.com/EEAjaegbu/Investment-Analytisc--Assignmet-2>

## DATA

### Customer ID Data File

<b>Total number of observations</b>	49171
<b>Total number of files</b>	1
<b>Total number of features</b>	5
<b>Base format of the file</b>	.txt
<b>Size of the data</b>	1304 KB

### Cabdata Data File

<b>Total number of observations</b>	359392
<b>Total number of files</b>	1
<b>Total number of features</b>	8
<b>Base format of the file</b>	.csv
<b>Size of the data</b>	24,437 KB

### City Data File

<b>Total number of observations</b>	20
<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 KB

### Transaction\_ID Data File

<b>Total number of observations</b>	440098
<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>	11,688 KB

**Proposed Approach:****Dedup Validation (Identification):**

The danger of duplicates in the data used for analysis is that it bias the result of the analysis. To identify duplicates entry or cases in the data, duplicated () method in pandas will be used to find duplicated rows while the drop\_duplicates () method in pandas will be used in removing duplicates entry in the data.