

# Data Intake Report

Name: Data Analyst: Cross selling recommendation - Group Project

Report date:

Internship Batch: LISUM21

Version: 1.0

Data intake by: Bogdan-Remus Pintilie

Data intake reviewer: Raka Prasetya Nugraha, Parwinder Singh

Data storage location:

<https://drive.google.com/file/d/1OCOZGjoL14tgKZtCl4Q9cifuL8xwIAo6/view?usp=sharing>

<https://drive.google.com/file/d/1wCwF2D6lzxh50rBFWXwW330ByiHhVK0w/view?usp=sharing>

## Tabular data details:

<b>Total number of observations</b>	929615
<b>Total number of files</b>	1 (Test.csv)
<b>Total number of features</b>	24
<b>Base format of the file</b>	.csv
<b>Size of the data</b>	105.2 MB

<b>Total number of observations</b>	13647309
<b>Total number of files</b>	1 (Train.csv)
<b>Total number of features</b>	48
<b>Base format of the file</b>	.csv
<b>Size of the data</b>	2.14 GB

## Proposed Approach:

1. Independence: The assumption of independence suggests that the data points being analyzed are not influenced by each other. Independence allows for unbiased analysis and statistical tests that rely on the assumption of independence.

2. Validity: The assumption of validity implies that the data represents what it claims to represent. Valid data is relevant and appropriate for the analysis objectives, and it accurately captures the intended information or measurements.

3. Integrity: The assumption of data integrity suggests that the data has not been tampered with or altered inappropriately. Data integrity ensures that the data retains its original state and remains reliable throughout the analysis process.

\*4. Duplicate identification: The code below shows that the key field contains duplicates.

\*5. We improve data relevance by omitting the following fields since they do not provide insights and help address the problem statement: conyuemp,

