

Data Intake Report

Name: Cloud and API deployment

Report date: 08/11/2023

Internship Batch: LISUM23

Version:1.0

Data intake by: Madoka Fujii

Data intake reviewer:

Data storage location:

Tabular data details: housing.csv

Total number of observations	89885
Total number of files	1
Total number of features	5
Base format of the file	.CSV
Size of the data	6.7MB

Proposed Approach:

- Find a large size of CSV
- Try to read with Pandas, Pandas with chunks, Dask, Dask with chunks, Ray and check the speed to complete
- Conduct basic Validations (check missing data, white spaces from the col name, and data type)
- Generate YAML and util.py
- Validate with YAML file
- Create a gz format of CSV with pipe separated text file (|)
- Summary the total number of rows, total number of columns, and the file size