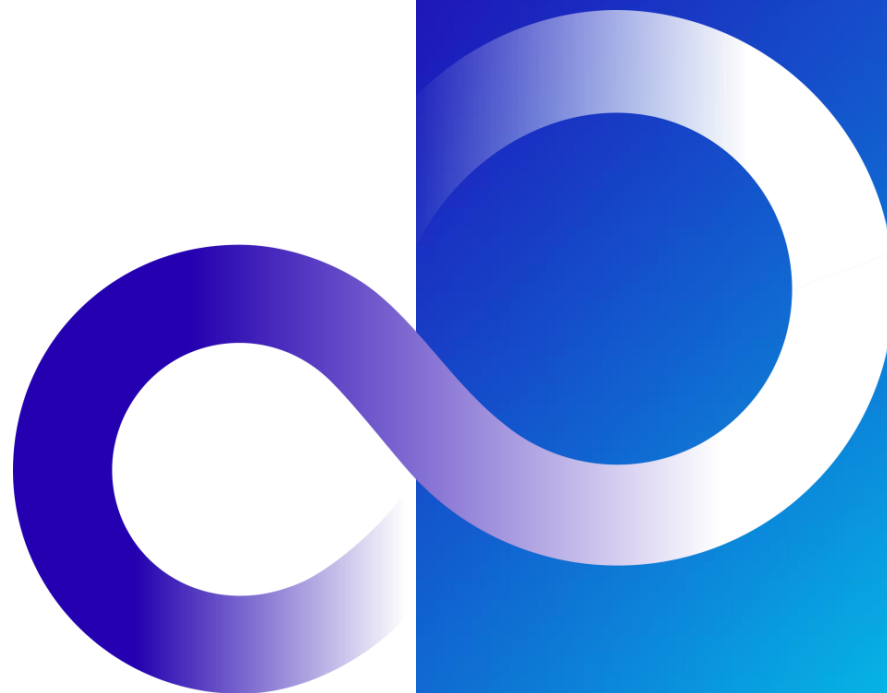# Fujitsu Kozuchi Auto Data Wrangling User Manual

Dec 20, 2023

FUJITSU Reserch

FUJITSU Limited

# Fujitsu Auto Data Wrangling



**Automatic pre-processing for tabular data using generative AI**

Reduce efforts of data preperation in AI application by automatic pre-processing for tabular data via data cleaning and data enrichment using generative AI.
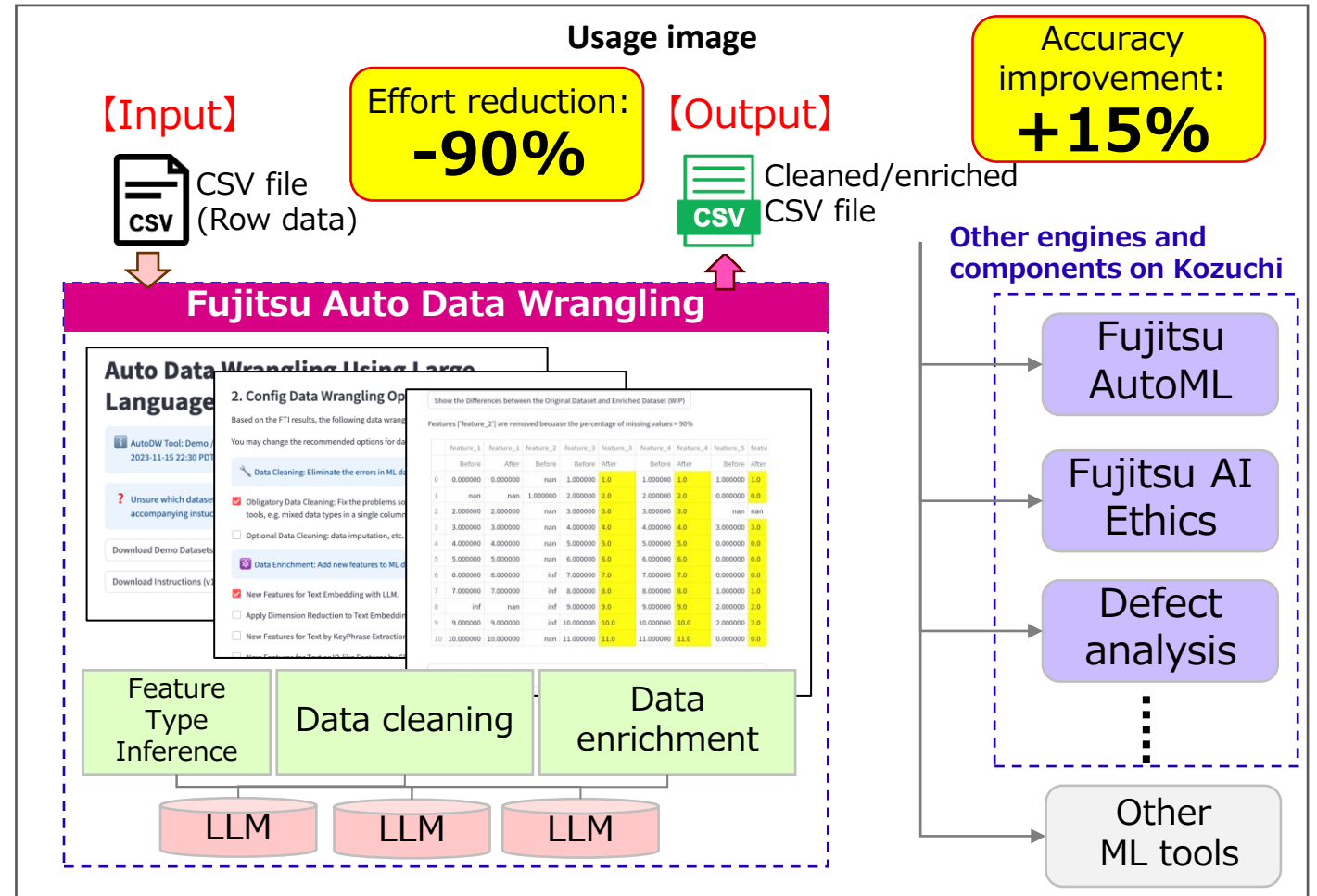
### Challenges

- Applying machine learning to tabular data in the field requires a lot of effort of pre-processing, such as data cleaning and data enrichment
- Conventional machine learning cannot handle a variety of text features as they are, which limits their accuracy

### Solutions

- **Automatic data cleaning** based on **Feature Type Inference** using LLMs
- **Automatic data enrichment** that creates new features by analyzing existing text features using LLMs
- **Achieving both automation and scalability** by using different LLMs depending on data wrangling processes

### Strengths

- Effort reduction of data preparation (90% reduction compared with manual coding of data wrangling)
- Automatic data enrichment for text features not addressed by other companies' data wrangling tools
- Improving accuracy of ML 15%+ (in case of using Fujitsu AutoML)

**Usage image**

【Input】 → CSV file (Row data)

Effort reduction: **-90%**

【Output】 → Cleaned/enriched CSV file

Accuracy improvement: **+15%**

**Fujitsu Auto Data Wrangling**

Feature Type Inference | Data cleaning | Data enrichment

LLM | LLM | LLM

**Other engines and components on Kozuchi**

- Fujitsu AutoML
- Fujitsu AI Ethics
- Defect analysis
- Other ML tools

# What you can and cannot do?

- **What you can do :** <span style="color:red">Preprocessing Before Machine Learning</span>
  - Data Cleaning
    - Format the input data to avoid errors when applying it to machine learning tools such as Fujitsu AutoML
  - Data Enrichment
    - Analyze the feature columns of input data and add new feature columns that contribute to improving analysis efficiency

- **What it can't do: Machine learning**
  - For machine learning processing such as data learning, classification, and prediction, please use existing machine learning tools such as Fujitsu AutoML.

# Two ways to use

- ## Interactive Demo
  - You can experience the operation of the Fujitsu Auto Data Wrangling using sample data provided by us.

- ## Proof of Concept (PoC)
  - You can use the Fujitsu Auto Data Wrangling with your own data.

# Access to the environment

- Preparation
  - Download "Environment_Access_Manual" from the [access manual](access manual)
  - **PoC**
    - Refer to "How to Connect with Azure VPN Gateway" section
  - **Interactive Demo**
    - Refer to "Bastion connection with Remote Desktop" section
- Access
  - http://10.0.0.139:8550/ from a web browser
  - Confirmation is successful if the start screen, as displayed on the right, appears

# 2. How to use Fujitsu Auto Data Wrangling webapp

FUJITSU

# 2.1 Getting Started with Sample Datasets (Optional)

**FUJITSU**



## Auto Data Wrangling Using Large Language Models

ℹ️ AutoDW Tool: Demo / Test Edition. Current in Development. Version v1.1.11. Last Update on 2023-11-15 22:30 PDT

❓ Unsure which datasets to begin with? No worries! Download our demo datasets and accompanying instuctions below

[Download Demo Datasets] ⬅ Download sample demo datasets

[Download Instructions (v1.1)] ⬅ Download user instructions

- If you are unsure about which datasets to begin with, you may download some sample demo datasets and check the effects of auto data wrangling first

- You may also download the user manual / instructions

# Sample Demo Dataset #1



- **titanic.csv**
  - Common dataset for demo purpose
  - Target column: Survived
  - ML task: Classification

# Sample Demo Dataset #2

FUJITSU



- consolidated_coin_data.csv
  - A representative data to show the effect of FTI
    - The numbers look numerical values, but they are strings (or numerical numbers embedded in strings)
  - Target column: Currency
  - ML task: Classification

- **dirty_data.csv**
  - A representative data to show the effect of data cleaning
    - A small dataset full of different type of errors (e.g. mixed data type in a single column, NaN value decoded as "?", etc.)
  - Target column: target
  - ML task: Classification

# Sample Demo Dataset #4



- **NYC_Airbnb_2019(changed).csv**
  - A representative data to show the effect of FTI, data cleaning, data enrichment
    - Include feature columns of text type, datetime type, and numeric type
    - Embedded by editing dirty_data.csv similar error
    - Created by modifying Airbnb listings and metrics in NYC, NY, USA (2019)©Airbnb Licensed under CC BY 4.0
  - Target column: Price
  - ML task:Regression

The following pages are explained using this data. Use this file in the demo video.

# 2.2 Specify Dataset & Problem Setting (1)

## 1. Upload Dataset (CSV) for Data Wrangling

Please upload your dataset for ML applications.

Drag and drop file here
Limit 1GB per file • CSV

Browse files

- Two approaches to upload a dataset (CSV) for data wrangling
  - Drag and drop a csv file
  - Click "Brower files" button and select a csv file

# 2.2 Specify Dataset & Problem Setting (2)



NY_Airbnb_2019(changed).csv  0.8MB

| | id | name | neighbourhood | room_type | minim |
|---|---|---|---|---|---|
| 0 | 2,539 | Clean & quiet apt home by the park | Kensington | Private room | 1 |
| 1 | 2,595 | Skylit Midtown Castle | Midtown | Entire home/apt | 1 |
| 2 | 3,647 | None | Harlem | Private | None |
| 3 | 3,831 | Cozy Entire Floor of Brownstone | Clinton Hill | Entire home/apt | 1 |
| 4 | 5,022 | Entire Apt: Spacious Studio/Loft by central park | None | Entire home/apt | None |
| 5 | 5,099 | Large Cozy 1 BR Apartment In Midtown East | Murray Hill | Entire home/apt | 3 |
| 6 | 5,121 | BlissArtsSpace! | Bedford-Stuyvesant | Private room | 45 |
| 7 | 5,178 | Large Furnished Room Near B'way | Hell's Kitchen | Private room | None |
| 8 | 5,203 | Cozy Clean Guest Room - Family Apt | Upper West Side | None | 2 |
| 9 | 5,238 | Cute & Cozy Lower East Side 1 bdrm | Chinatown | Entire home/apt | none |

- A preview of the uploaded dataset is displayed.

- Select target variables (columns you want to predict) in "Target Columns"

- Select ML tasks, the two options are:
  - Classification: for ML tasks that predict discrete class labels
  - Regression: for ML tasks that predict a continuous quantity

- After the Target columns and ML task are specified, feature type inference (FTI) is automatically executed for the uploaded dataset, which predicts the feature type for each columns of the dataset.

# 2.3 Feature Type Inference (2)



The ML task is: Regression

Feature Type Inference (FTI) Results:

[i] Optional: User can change the FTI results by clicking the cells of Feature Type column

| Column Name | Column Type | Feature Type |
|---|---|---|
| price | Target | Numeric |
| id | Feature | ID |
| name | Feature | Sentence |
| neighbourhood | Feature | Categorical |
| room_type | Feature | Categorical |
| minimum_nights | Feature | Categorical |
| last_review | Feature | Datetime |

- After the FTI prediction is completed, a table with three columns is displayed:
  - Column Name: shows the name of each column in the uploaded dataset
  - Column Type: shows a column is a target or a feature
  - Feature Type: shows the predicted feature type for each column

- Feature Types
  - **Numeric**: quantitative data, such as 1, 2, 3, ⋯
  - **Categorical**: a variable with a set number of groups, such as male, female
  - **Datetime**: dates and times, such as 11-23-2022, 15:20PM, etc.
  - **Sentence**: a set of words, such as "auto data wrangling tools are useful", etc.
  - **URL**: Uniform Resource Locator, such as https://www.fujitsu.com/global/about/research/
  - **Embed**: A string with numerical value embedded, such as $1,000
  - **List**: A sequence of several variables, grouped together, such as [A, B, C], [1, 2, 3], etc.
  - **ID**: An identity column, such as the index
  - **Unit**: a determinate quantity as a standard of measurement, such as 100m, 60kg, etc.
  - **Sign**: A string with symbols such as >=, <= with numerical values, examples are >50, <=100, etc.
  - **Range**: A string shows a range of numerical values, such as 60-100

# Edit FTI Results (Optional)

- The "Feature Type" column in the FTI results is editable
- If user feels some of the FTI results is not satisfied, you can click the cell to change the predicted FTI type
- The following processing will be based on the feature type that user selected

# 2.4 Config data wrangling options (1)



- Based on the FTI results, certain data wrangling options are suggested automatically.
- Nonetheless, users have the flexibility to modify these options according to their requirements.
- Two major data wrangling options:
  - Data Cleaning: eliminate the errors in the dataset
  - Data enrichment: add new features to the dataset

# 2.4 Config data wrangling options (2)
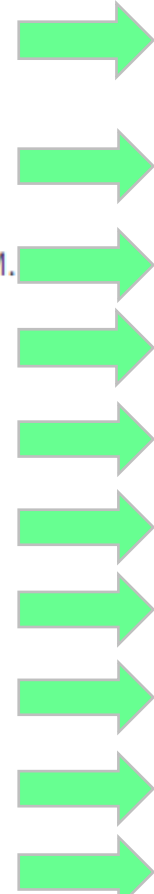
🔧 Data Cleaning: Eliminate the errors in ML datasets

☑ Obligatory Data Cleaning: Fix the problems so that the dataset can be handled by common AutoML tools, e.g. mixed data types in a single column

☐ Optional Data Cleaning: data imputation, etc.

- Data cleaning has two options:
  - **Obligatory data cleaning** is essential to address issues in the dataset, enabling seamless compatibility with popular AutoML tools. This process involves checking for potential errors, such as decoding datasets when necessary, refining headers, eliminating irrelevant features, addressing NaN cells in the target, handling infinite values, cleaning columns with mixed data types, and cleaning text columns, etc.
  - **Optional data cleaning** serves as an additional layer on the obligatory data cleaning results to ensure the dataset is fully prepared for AutoML tools. This step may involve operations such as data imputation and encoding the target column. These operations are considered optional, as many common AutoML tools are capable of performing them as well.

# 2.4 Config data wrangling options (3)

Data Enrichment: Add new features to ML datasets

☑ New Features for Text Embedding with LLM.

☐ Apply Dimension Reduction to Text Embedding

☐ New Features for Text by KeyPhrase Extraction with LLM.

☐ New Features for Text or ID-like Features by Clustering with LLM.

☐ New Features generated from List.

☑ New Features generated from Datatime.

☐ New Features generated from URL.

☐ New Features generated from Embedded Numbers.

☐ New Features generated from Number Ranges.

☐ New Features generated from Unit Features.
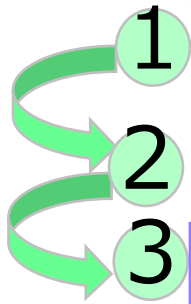
☐ New Features generated from Inequality Sign.

● **Data enrichment based on FTI**

Add new text embedding features using LLM for the "Sentence" features (more details in the next page)

Add new keyphrase features for the "Sentence" features using LLM

Add new clustering features for the "Sentence" features using LLM

Add new features for the "List" features

Add new features for the "Datetime" features

Add new features for the "URL" features

Add new features for the "Embed" features

Add new features for the "Range" features

Add new features for the "Unit" features

Add new features for the "Sign" features

Data Enrichment: Add new features to ML datasets

① New Features for Text Embedding with LLM.

② Apply Dimension Reduction to Text Embedding

② appears when ① is checked

Select input type for reduced dimension:

● Number Input
○ Slider

③ appears when ② is checked

Enter the reduced dimension

7

Press Enter to apply

- **Text embedding configuration**
  - You can apply dimension reduction to text embedding by checking ②
  - You can configure the number of dimensions of text embedding (from 1 to 768) by using either a number input box or slider, as shown in ③

**FUJITSU**

📙 Actions to be Performed & Explanations

Obligatory Data Cleaning will be conducted. The Data cleaning module will check for possible dataset errors and fix them, including but not limited to: decode datasets as necessary, clean headers, remove irrelevant features, drop NaN cells in target, process and replace infinite values, handle columns with mixed data types, encode the target column for machine learning compatibility, text column cleaning, etc.

LLM Embedding will be conducted

Embedding dimension is:768

Because columns ['last_review', 'reviews_per_month'] are detected as Datatime features, new features will be generated from these Datatime features, for example, MM/DD/YYYY => MM, DD, YYYY

Columns ['price'] are string instead of numerical values. To better use the features, the string will be converted into numerical values, or the numerical values embedded in the string will be extracted.

Because columns ['price'] are detected as Embedded Number features, new features will be generated from these Embedded Number features, for example, $1,000 => 1000

- The data wrangling actions are summarized, explained, and are displayed in UI.

## 3. Conduct Data Wrangling

Press the button to start data wrangling

Start Data Wrangling

- After the data wrangling options are configured, click the button "Start Data Wrangling" to start the data cleaning and enrichment processing based on the configurations

**FUJITSU**



- After the data wrangling processing is done, the cleaned & enriched dataset is displayed

**FUJITSU**

**Descriptive statistics**

| | 57 | last_review_Year | last_review_Month | last_review_Day | last_review_WeekDay | last_review_Hour |
|---|---|---|---|---|---|---|
| count | 56 | 8,660 | 8,660 | 8,660 | 8,660 | 8,660 |
| unique | ne | None | None | None | None | None |
| top | ne | None | None | None | None | None |
| freq | ne | None | None | None | None | None |
| mean | 37 | 2,017.7411 | 6.3028 | 15.9072 | 3.1374 | 0 |
| std | 35 | 1.6403 | 2.6837 | 9.8054 | 2.2032 | 0 |
| min | 78 | 2,008 | 1 | 1 | 0 | 0 |
| 25% | 16 | 2,016 | 5 | 6 | 1 | 0 |
| 50% | 39 | 2,019 | 6 | 17 | 3 | 0 |
| 75% | 53 | 2,019 | 8 | 24 | 5 | 0 |

- The descriptive statistics for the cleaned & enriched dataset is also displayed

Click to Download Enriched Dataset

Click to Download ML Task Specification in JSON

Show the Differences between the Original Dataset and Enriched Dataset

Click to View EDA Report for the Enriched Dataset

- Click the button "Click to Download Enriched Dataset" to download the cleaned & enriched dataset to your local machine. The downloaded dataset is in CSV format and can be further used for other applications such as AutoML.

FUJITSU



○Click the button "Click to Download ML Task Specification" to download the ML task specification for the cleaned & enriched dataset to your local machine. The ML task specification is in JSON format and can be further used for other applications such as AutoML.

- Click the above button to
  - Show the differences (highlighted in yellow color) between the original dataset and enrich dataset
    - By default, only 10 rows are displayed, but you can configure it by clicking "select the number of rows to show"
  - Explain the reason if columns are dropped
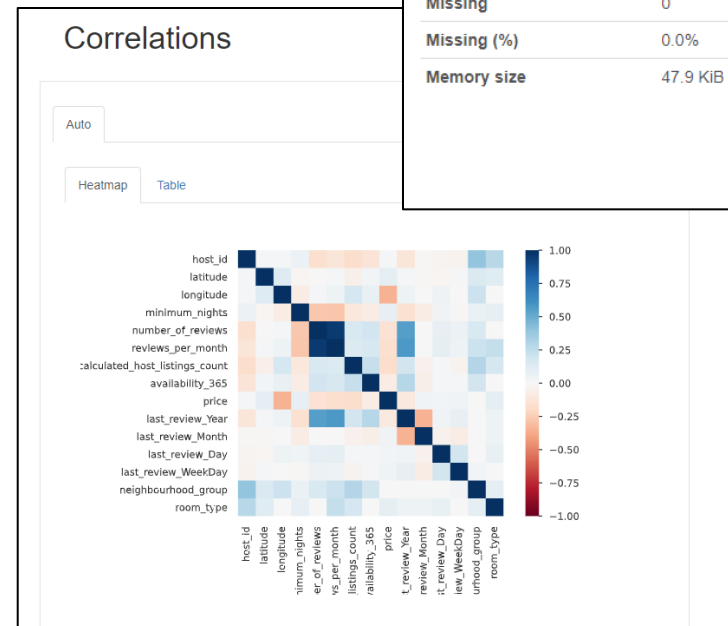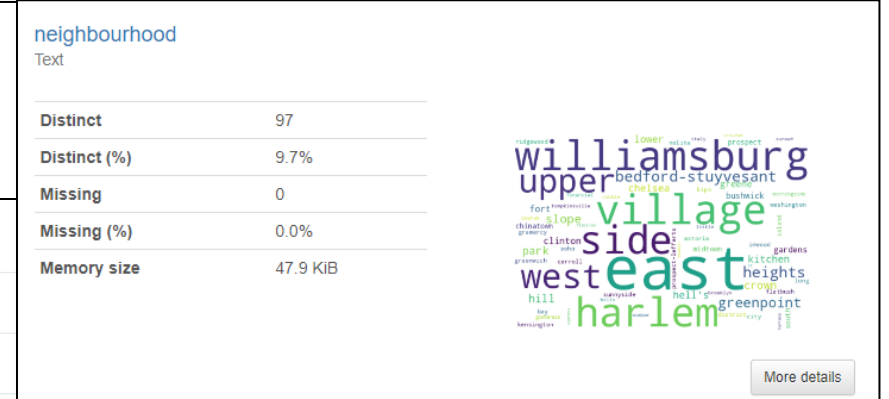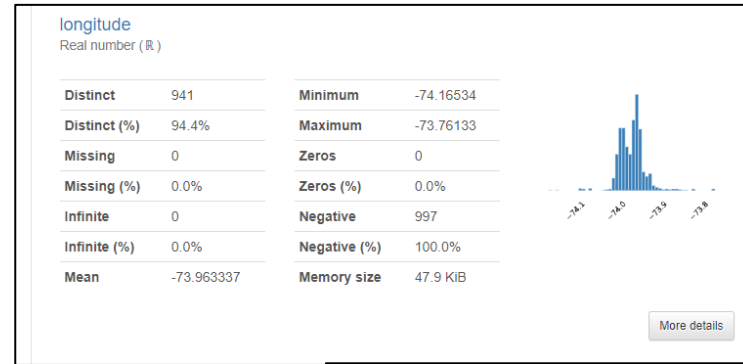
30

FUJITSU

Click to Download Enriched Dataset

Click to Download ML Task Specification in JSON

Show the Differences between the Original Dataset and Enriched Dataset

Click to View EDA Report for the Enriched Dataset

- Click the button to generate the EDA report for the cleaned & enriched dataset
  - Note: When there are too many columns in the enriched dataset (e.g. if the LLM text embedding is applied), the generation of EDA report could be slow due to the interaction visualizations between every column pairs

**longitude**
Real number ( ℝ )

| | | | |
|---|---|---|---|
| Distinct | 941 | Minimum | -74.16534 |
| Distinct (%) | 94.4% | Maximum | -73.76133 |
| Missing | 0 | Zeros | 0 |
| Missing (%) | 0.0% | Zeros (%) | 0.0% |
| Infinite | 0 | Negative | 997 |
| Infinite (%) | 0.0% | Negative (%) | 100.0% |
| Mean | -73.963337 | Memory size | 47.9 KiB |

More details

**neighbourhood**
Text

| | |
|---|---|
| Distinct | 97 |
| Distinct (%) | 9.7% |
| Missing | 0 |
| Missing (%) | 0.0% |
| Memory size | 47.9 KiB |

More details

**Correlations**

Auto

Heatmap    Table

# Try a Different Dataset



- If you complete the data wrangling for a dataset and want to try a different dataset, please click here to start the data wrangling for a new dataset.

# Contact us (for customers)

- Please contact your Fujitsu representative.

Thank you

FUJITSU

© 2023 Fujitsu Limited