# Data Formatting Demo

## Version 1.0

**Author: Shuiliang (Leon) Wu**

## Table of Contents

# Introduction

This demo provides a guide to format provided `final_cleaned_diabetes.csv` raw data (cleaned) to meet `Requirements of Data Format` so that date can be recognized by the source code for federated learning.
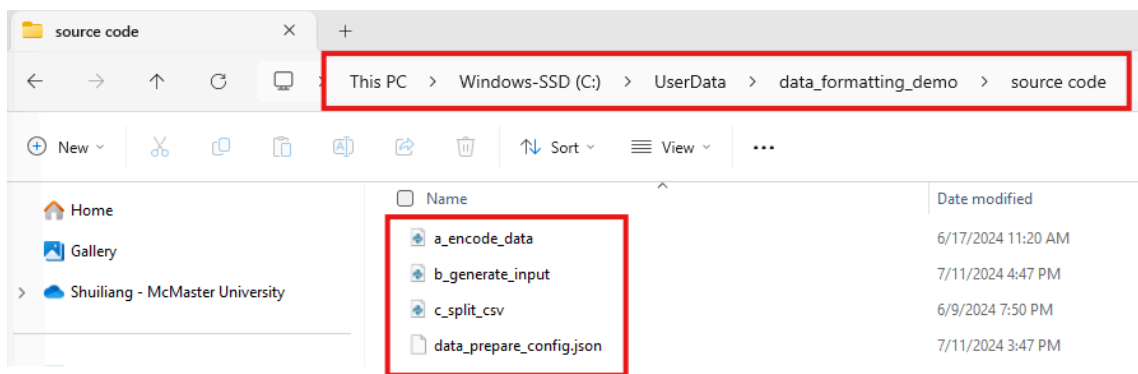
# What You Have Before Started

- Data Formatting Demo (this one)

- Source code in `data_formatting_demo` including four files:

    - `a_encode_data.py`

    - `b_generate_input.py`

    - `c_split_csv.py`

    - `data_prepare_config.json`

- Raw dataset `final_cleaned_diabetes.csv` (in `raw_data` folder)
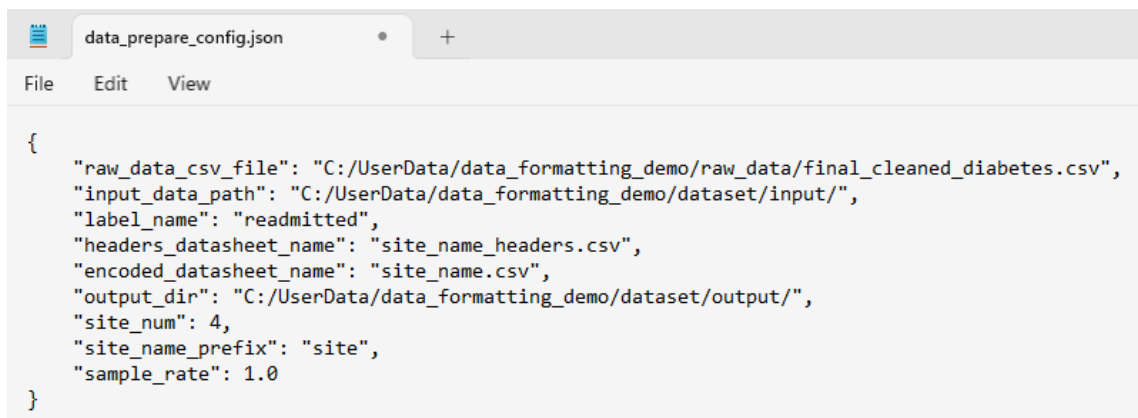
# Getting Started

This Demo is demonstrated using Windows 11 (operations in macOS are very similar).
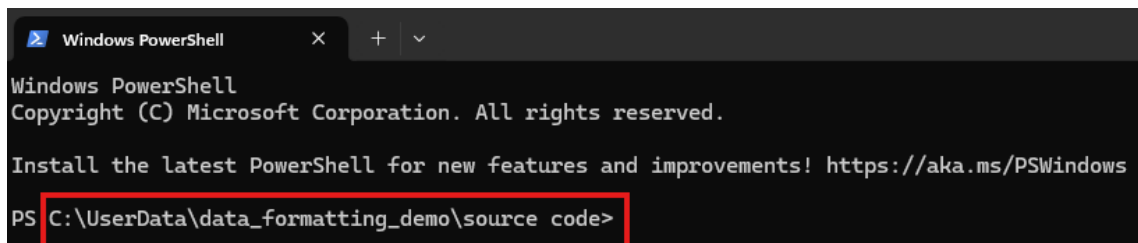
1. Open the `source code` folder in `data_formatting_demo`.

2. Right-click `data_prepare_config.json`, select `Edit in Notepad` then update the parameters accordingly.



```json
{
    "raw_data_csv_file": "C:/UserData/data_formatting_demo/raw_data/final_cleaned_diabetes.csv",
    "input_data_path": "C:/UserData/data_formatting_demo/dataset/input/",
    "label_name": "readmitted",
    "headers_datasheet_name": "site_name_headers.csv",
    "encoded_datasheet_name": "site_name.csv",
    "output_dir": "C:/UserData/data_formatting_demo/dataset/output/",
    "site_num": 4,
    "site_name_prefix": "site",
    "sample_rate": 1.0
}
```
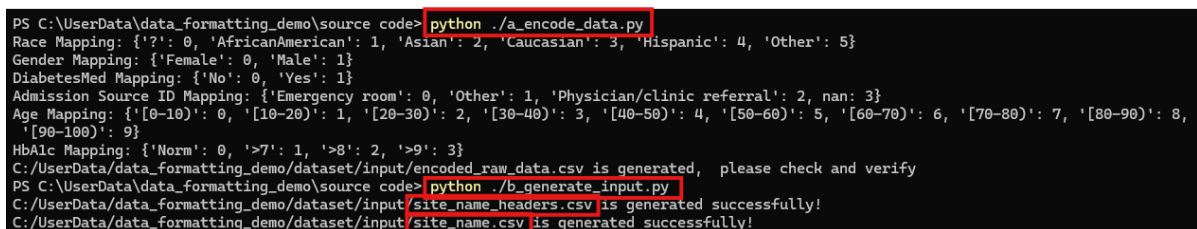
3. Right-click at the blank space in `source code` folder, select `Open in Terminal`. The terminal shall be pop-up as below:
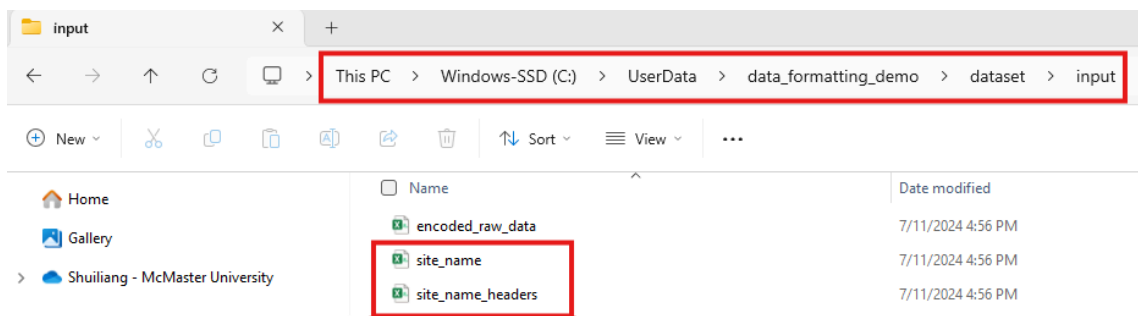


4. Formatting the data by running below in terminal one-by-one:

```
python ./a_encode_data.py
python ./b_generate_input.py
```



5. There shall be `site_name.csv` and `site_name_headers.csv` generated in `input` folder under `dataset` folder. Replace `site_name` with your actual site name, and they are ready to be used.

6. Since there will be four sites in total in the demo of `NVIDIA FLARE User Guide for Project Manager` and `NVIDIA FLARE User Guide for Site Admin`, the formatted dataset is split into four datasets with equal amount of data in each datasets by running:

```
python ./c_split_csv.py
```

```
PS C:\UserData\data_formatting_demo\source code> python .\c_split_csv.py
site1= start_index=0 end_index=4254
File copied to C:/UserData/data_formatting_demo/dataset/output/site1.csv
site2= start_index=4254 end_index=8508
File copied to C:/UserData/data_formatting_demo/dataset/output/site2.csv
site3= start_index=8508 end_index=12762
File copied to C:/UserData/data_formatting_demo/dataset/output/site3.csv
site4= start_index=12762 end_index=17018
File copied to C:/UserData/data_formatting_demo/dataset/output/site4.csv
File copied to C:/UserData/data_formatting_demo/dataset/output/site1_header.csv
File copied to C:/UserData/data_formatting_demo/dataset/output/site2_header.csv
File copied to C:/UserData/data_formatting_demo/dataset/output/site3_header.csv
File copied to C:/UserData/data_formatting_demo/dataset/output/site4_header.csv
```

7. There shall be four datasets generated in `output` folder under `dataset` folder. Replace `site#` with the actual site name, and they are ready to be used for demo.