

# Readme

There are mainly three notebooks in the 'src' directory. They are listed and described below:

## 1. Data\_Preprocessing.ipynb

This notebook is mainly for data preprocessing. It will generate the preprocessed data files in two directories named 'cleaned\_data' and 'imputed\_data'.

- The 'cleaned\_data' will include the data files with NaN value for the 35 air quality stations.
- The 'imputed\_data' directory includes the One-Hot encoded and data imputed files for final prediction.

## 2. XGBoost.ipynb

This notebook will generate the final 'submission.csv' file with the model XGBoost.

## 3. LSTM.ipynb

This notebook will use the LSTM model for data analysis and simple prediction. Since the result of LSTM is not better than XGBoost. So this notebook will not make the final prediction for the 35 air quality stations.

Meanwhile, there is a directory name 'MSBD5002PROJECT\_data' in the 'src' directory. There are three files in this directory:

- **airQuality\_station.csv**: includes the station location information for 35 air quality stations
- **Beijing\_grid\_weather\_station.csv**: includes the location information for the grid weather stations
- **observed\_weather\_stations.csv**: include the location information for the observed\_weather stations