

## Table of contents

---

### Dataset

Collection

Pre-processing

Description

---

### Methodology

3 Ensemble Regressors

(RandomForest, XGBoost, LightGBM)

---

### Empirical Analysis & Result



## Data (Collection)



### Disaster Prevention Weather Data ( Seoul )

Source: Meteorological Data Open Portal  
( <https://data.kma.go.kr/cmmn/main.do>.  
accessed 10 Dec 2023. )

Data : Disaster Prevention Weather Data

Hourly Measurement data of disaster prevention weather by districts in Seoul

Period : 2016.1.1 ~ 2020.12.31

Variables : temperature, wind speed, wind direction, precipitation

\* Local air pressure, humidity, etc. had high multicollinearity with other variables, so only four variables above were selected



### Seoul Air Quality Data

Source: Air Korea  
( [https://www.airkorea.or.kr/web/sidoQualityCompare?itemCode=10008&pMENU\\_NO=102](https://www.airkorea.or.kr/web/sidoQualityCompare?itemCode=10008&pMENU_NO=102).  
accessed 10 Dec 2023. )

Data : Air Quality Data

Hourly measurement data of air quality monitored by districts of Seoul

Period : 2016.1.1 ~ 2020.12. 31

Variables : NO2 , PM10, PM2.5

\* NO2 was selected as a variable that affects the predictors: PM10 and PM2.5



### Seoul Greenspace data

Source: Seoul Basic Statistics Public Data Portal  
( <https://data.seoul.go.kr/dataList/368/S/2/datasetView.do>.  
accessed 10 Dec 2023. )

Data : Greenspace Data

Green space-related data for each district of Seoul prepared annually

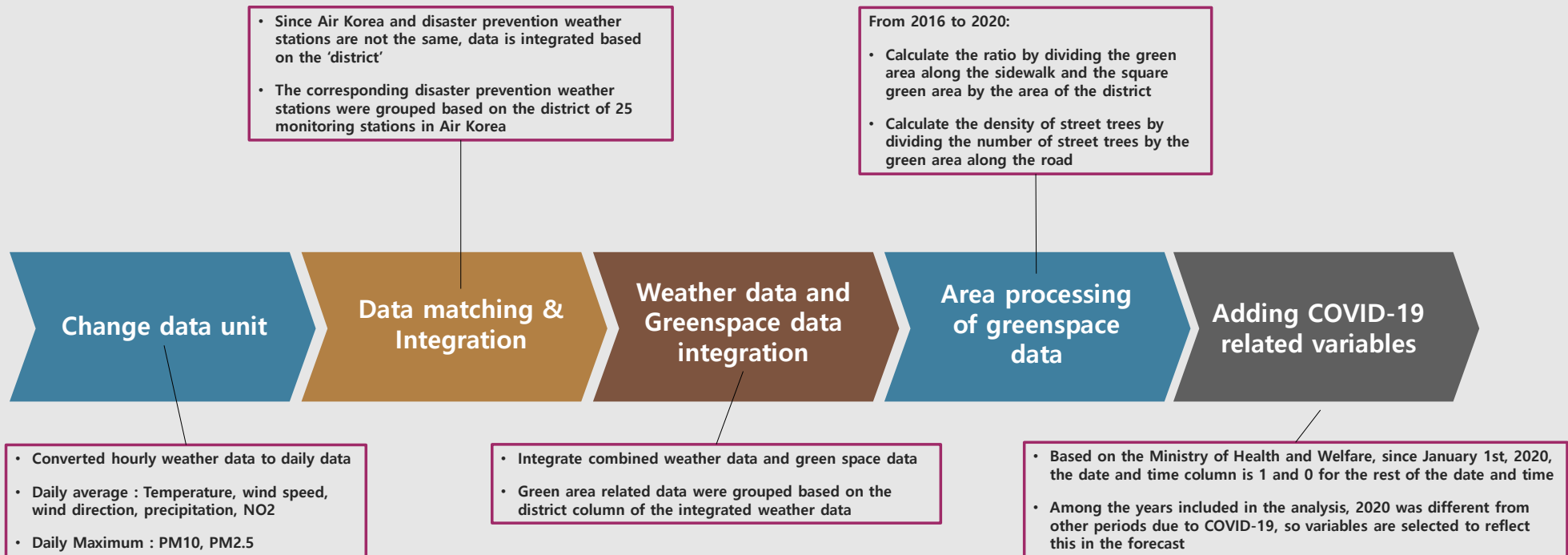
Period : 2016 ~ 2020

Variables: Green area by road, square green area, number of street trees, number of berry trees, number of zelkova trees, number of cherry trees, number of ginkgo trees

\* Four selected tree types were selected in order of the largest number of trees in Seoul's streets



## Data (Pre-processing)



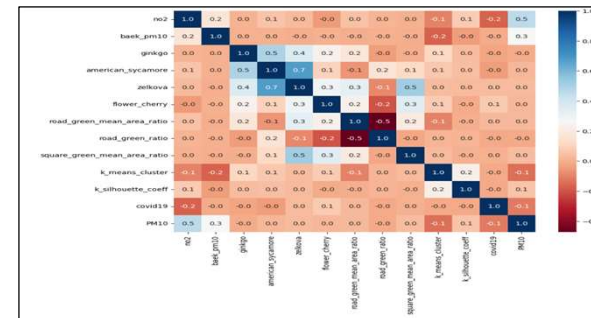


## Data (Description)

### Dataset

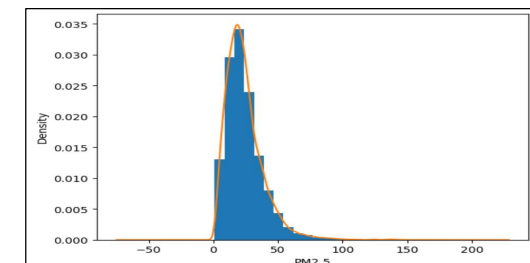
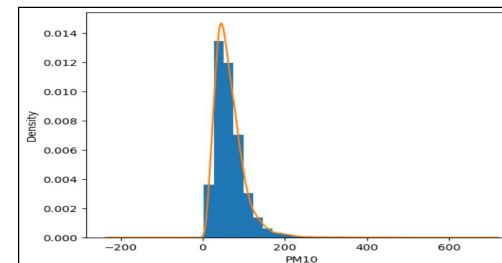
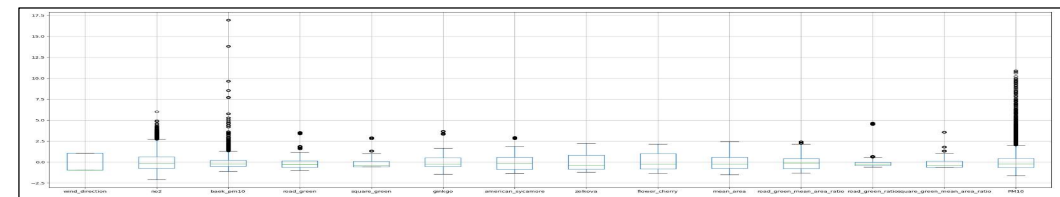
Variable	Unit	Source	Description
temperature	°C	Meteorological Data Open Portal	Converted regional hourly data to daily average data
wind_direction	°		0: 0° - 180° (easterly wind) 1: 180° - 360° (westerly wind)
wind_speed	m/s		Converted regional hourly data to daily average data
rainfall	0 or 1	Air Korea	0: Rainy X 1: Rainy
no2	µg/m <sup>3</sup>		Converted regional hourly data to daily average data
PM10	µg/m <sup>3</sup>		
PM2.5	µg/m <sup>3</sup>		
zelkova	trees / day	Seoul Basic Statistics Public Data Portal	Converted regional annual data to daily average data
ginkgo	trees / day		
flower_cherry	trees / day		
american_sycamore	trees / day	Seoul Basic Statistics Public Data Portal	Converted regional annual data to daily average data
road_green_mean_area_ratio	m <sup>2</sup> / m <sup>2</sup>		
road_green_ratio	(trees / day) / m <sup>2</sup>		
square_green_mean_area_ratio	m <sup>2</sup> / m <sup>2</sup>	Ministry of Health and Welfare	0: Before the COVID-19 outbreak ( ~2019.12) 1: After the COVID-19 outbreak (20.01~)
covid19	0 or 1		

### Check correlation coefficient & multicollinearity



	Feature	VIFscore
0	no2	6.478875
1	baek_pm10	2.646961
2	ginkgo	6.543249
3	american_sycamore	9.483502
4	zelkova	10.977438
5	flower_cherry	4.279494
6	road_green_mean_area_ratio	4.461421
7	road_green_ratio	1.725992
8	square_green_mean_area_ratio	2.463852
9	covid19	1.266677
10	PM10	4.151445

### Check Box Plot and Distribution



### Train/Test Dataset classification

```
1 X_train, X_test, y_train, y_test = train_test_split(X, Y, test_size=0.2, random_state=2021, stratify=Y)
```

Train/Test 8:2 separation, apply 'stratify = y' to set y to be more balanced



## Methodology

### Modeling – Regression Model

#### ● Linear Regression

- p-value

Assumes a straight-line relationship between input and output, finds the best-fit line to make predictions.

#### ● RandomForest Regressor

- Ensemble
- feature importance

- SHAP value  
(Shapely Additive Explanations)

Builds multiple decision trees and averages their predictions for better accuracy and resistance of overfitting.

#### ● XGBoost Regressor

- Ensemble
- feature importance

- SHAP value

Combines weak predictive models (usually trees) to create a strong model, optimizing for both speed and accuracy.

#### ● LightGBM Regressor

- Ensemble
- feature importance

- SHAP value

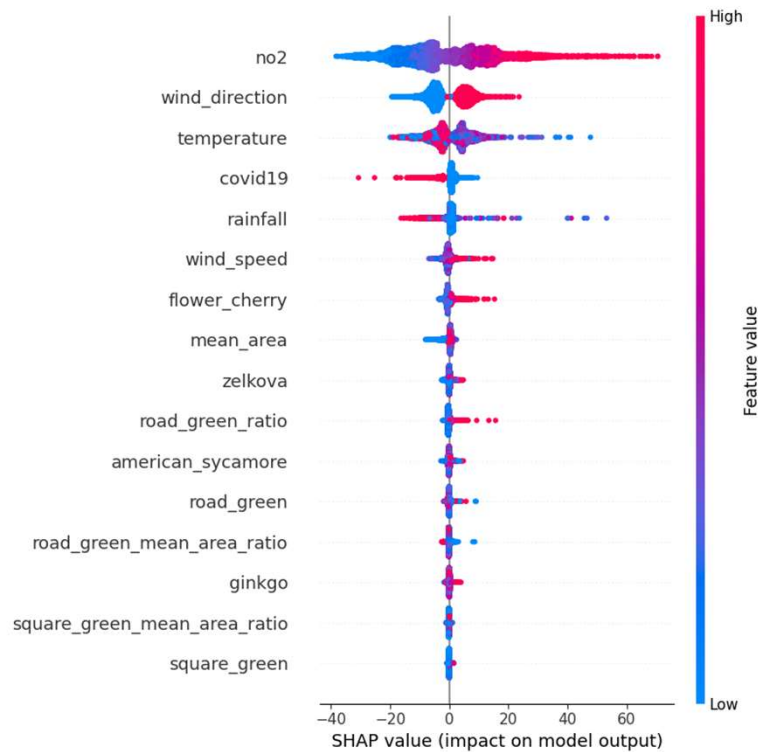
Similarly with XGBoost, it uses gradient boosting with a focus on efficiency and speed, especially for large datasets.



## Empirical analysis / Results (RF)

### Results Interpretation 1 – RandomForest Regressor

#### SHAP Value Check



#### SHAP Value Interpretation

- \* It is shown that the more red dot concentrated toward +(positive) and blue dot toward -(negative), the more the variable affects to y value.
- \* The higher the index, the higher fine dust had a positive correlation.

#### 1) Representative variables which shows positive correlation

NO2 / wind\_direction / wind\_speed / flower\_cherry

#### 2) Representative variables which shows negative correlation

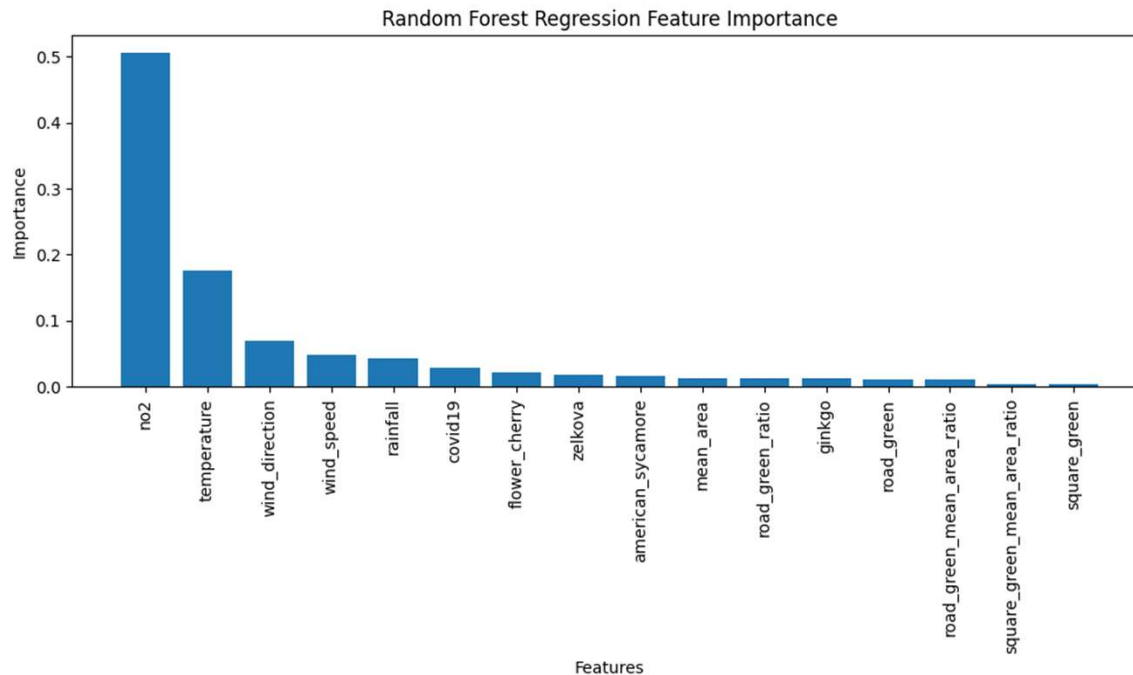
temperature / covid19 / rainfall



## Empirical analysis / Results (RF)

### Results Interpretation 1 – RandomForest Regressor

#### Importance of features



#### Important Features

NO2 > temperature > wind\_direction

(0.1761) (0.0692) (0.04778)

\* Air pollution and atmospheric conditions seems important for prediction.

#### Unimportant Features

road\_green > square\_green  
(includes mean area ratio)

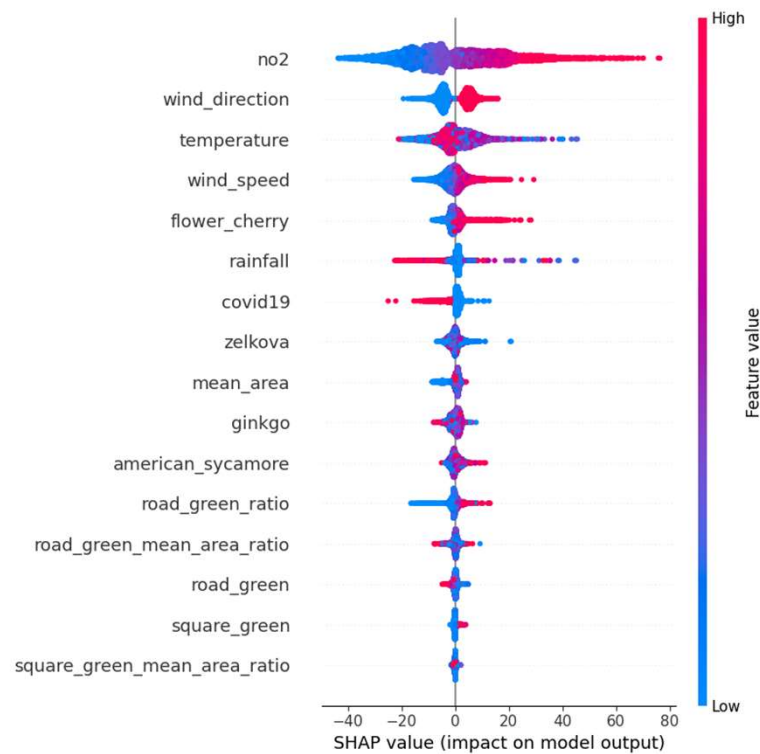
(0.0038) (0.0036)

\* Greenspace data did not made strong relationship between fine dust in this model.

## Empirical analysis / Results (XGB)

### Results Interpretation 2 – XGB Regressor

#### SHAP Value Check



#### SHAP Value Interpretation

- \* It is shown that the more red dot concentrated toward +(positive) and blue dot toward -(negative), the more the variable affects to y value.
- \* The higher the index, the higher fine dust had a positive correlation.

#### 1) Representative variables which shows positive correlation

NO2 / wind\_direction / wind\_speed / flower\_cherry

#### 2) Representative variables which shows negative correlation

temperature / rainfall / covid19

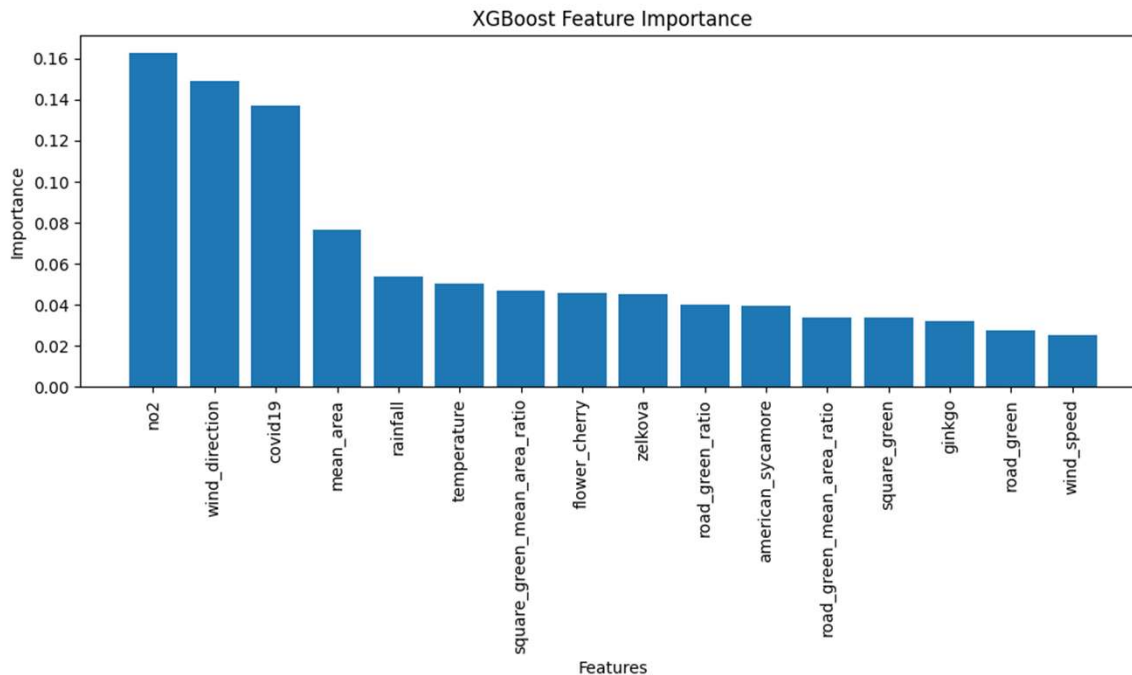




## Empirical analysis / Results (XGB)

### Results Interpretation 2 – XGB Regressor

#### Importance of features



#### Important Features

NO2 > wind\_direction > covid19

(0.1488)      (0.1372)      (0.0766)

\* Quite similar with RandomForest, newly revealed that the COVID-19 variable was also highly relevant with fine dust.

#### Unimportant Features

road\_green > wind\_speed

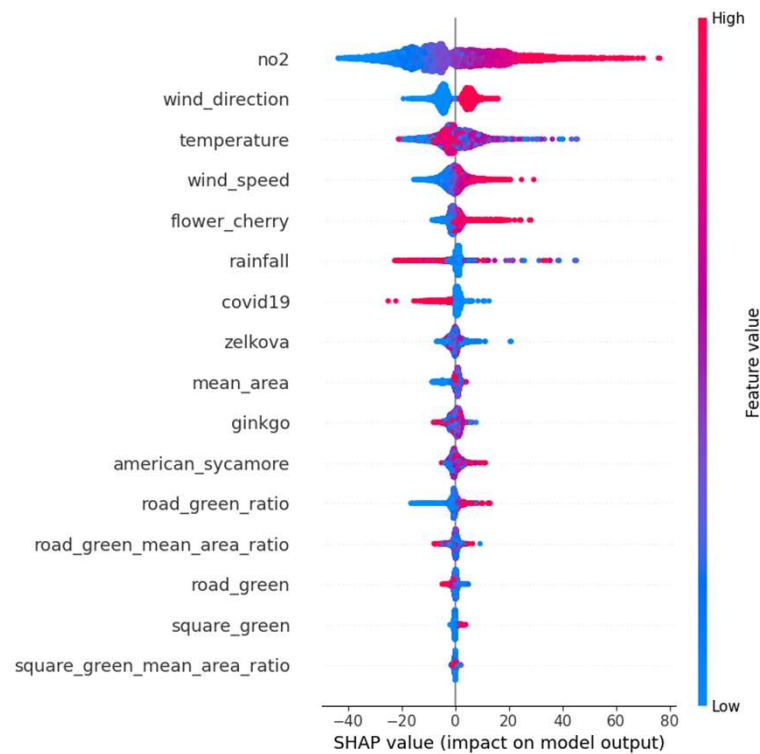
(0.0278)      (0.0250)

\* Turns out that speed of wind is not significantly related unlike the wind direction (shown in RF model)

## Empirical analysis / Results (LGBM)

### Results Interpretation 3 – LGBM Regressor

#### SHAP Value Check



#### SHAP Value Interpretation

- \* It is shown that the more red dot concentrated toward +(positive) and blue dot toward -(negative), the more the variable affects to y value.
- \* The higher the index, the higher fine dust had a positive correlation.

#### 1) Representative variables showing a positive correlation

NO2 / wind\_direction / wind\_speed / flower\_cherry

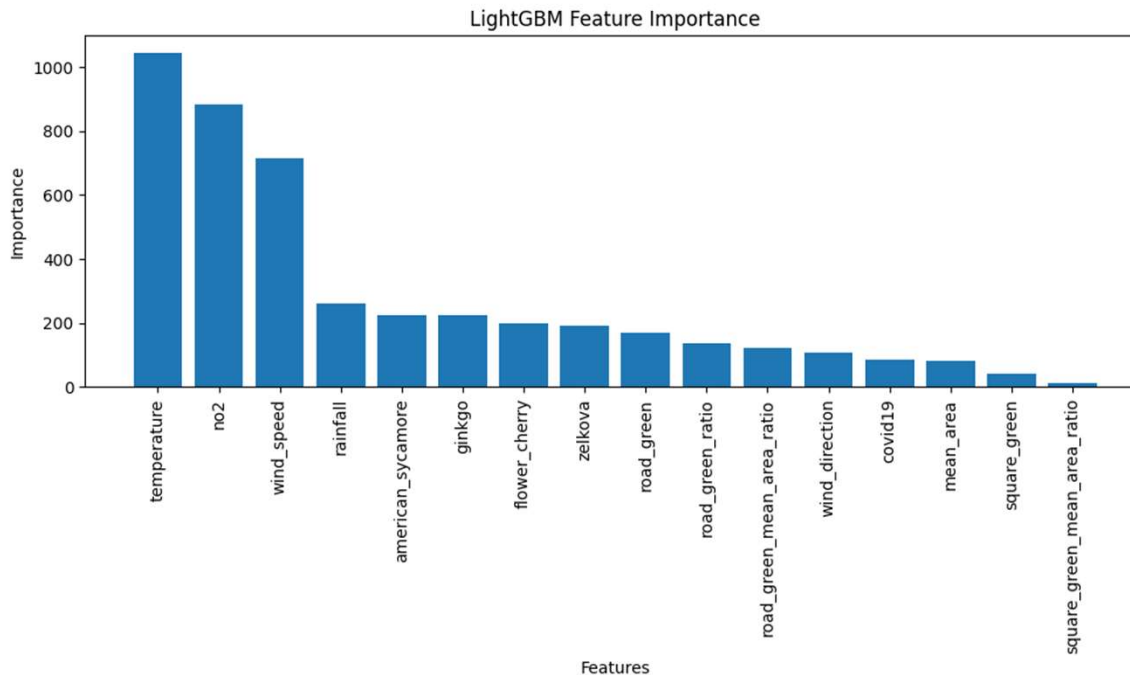
#### 2) Representative variables showing a negative correlation

temperature / covid19 / rainfall

## Empirical analysis / Results (LGBM)

### Results Interpretation 3 – LGBM Regressor

#### Importance of features



#### Important Features

temperature > NO2 > wind\_speed  
(884) (714) (263)

\* Wind speed has risen again as an important variable, which seems to require further consideration.

\*In the SHAP analysis, wind\_speed always showed a high correlation: considerable as important one

#### Unimportant Features

Square\_green > Square\_green\_mean\_area\_ratio  
(42) (13)

\* Same result with RF model, in which the greenspace variables were not important as weather variables



## Conclusion

Model	MSE	Important Variables
RandomForest	994.927	NO2 / Temperature / Wind direction
XGBoost	923.103	NO2 / Wind direction / COVID-19
LightGBM	921.886	Temperature / NO2 / Wind speed

- By integrating 3 types of data, we were able to analyze the variables that affect fine dust.
- 1) NO<sub>2</sub> is the most influential variable that shows a positive correlation with fine dust. The remaining positive correlation variables are wind\_direction, wind\_speed.
- 2) Temperature is the most influential variable that shows negative correlation with fine dust. The remaining negative correlation variable is COVID-19.
- 3) Weather-related variables mainly affected the results, but some regional green-related variables such as flower\_cherry also had a significant positive correlation.
- If we conduct analysis with more data sources, it would help to improve climate predictions.