

Air quality in Madrid

Ines Vilar

Final Project - Data Science - AllWomen

April 2020



- Home
- Compete
- Data
- Notebooks
- Discuss
- Courses
- More

Dataset

Air Quality in Madrid (2001-2018)

Different pollution levels in Madrid from 2001 to 2018

Decide Soluciones • updated 2 years ago (Version 5)

Data Tasks Kernels (29) Discussion (5) Activity Metadata Download (570 MB) New Notebook

Usability 7.9

License Other (specified in description)

Tags natural and physical sciences, time series, green living and environmental issues, environment, pollution

Context

In recent years, high levels of pollution in Madrid have forced the authorities to take measures to monitor and address this issue. One of the most effective and controversial measures was putting restrictions on the use of polluting vehicles in the city center.

To monitor the success of these measures Madrid's City Council Open Data website, makes publicly available daily and hourly historical data of the levels of contamination from different toxic particles, for a number of active stations across the whole city, from 2001 to real time.

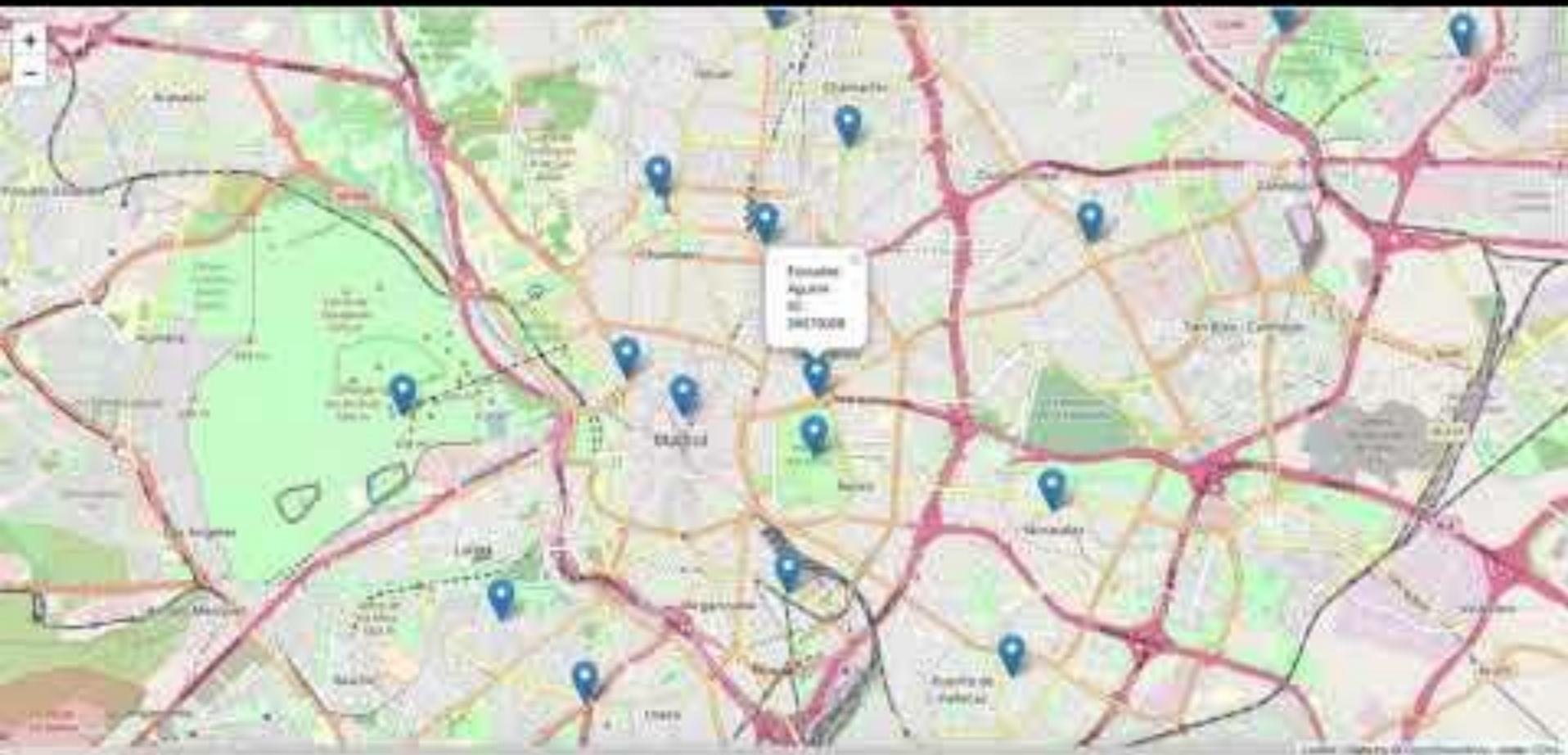
NO₂: nitrogen dioxide level measured in $\mu\text{g}/\text{m}^3$.
Long-term exposure is a cause of chronic lung diseases,
and are harmful for the vegetation.

O₃: ozone level measured in $\mu\text{g}/\text{m}^3$.
High levels can produce asthma, bronchitis or other
chronic pulmonary diseases in sensitive groups or
outdoor workers.

PM₁₀: particles smaller than 10 μm measured in $\mu\text{g}/\text{m}^3$.
Even though the cannot penetrate the alveolus, they can
still penetrate through the lungs and affect other organs.
Long term exposure can result in lung cancer and
cardiovascular complications.

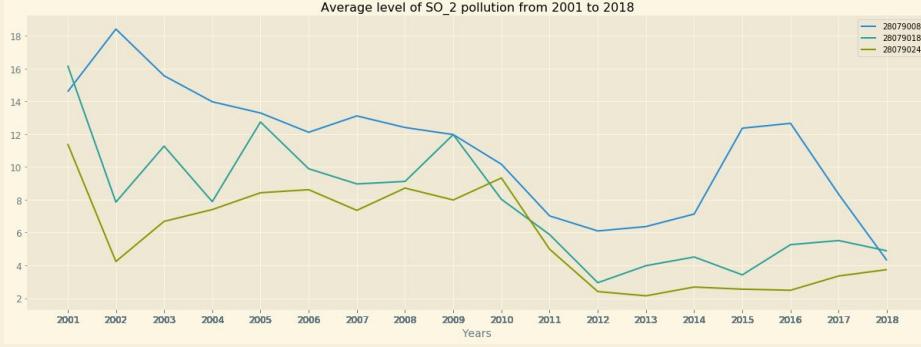
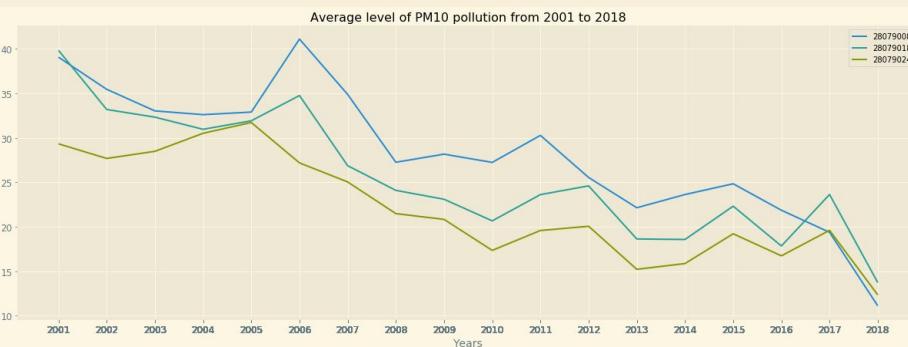
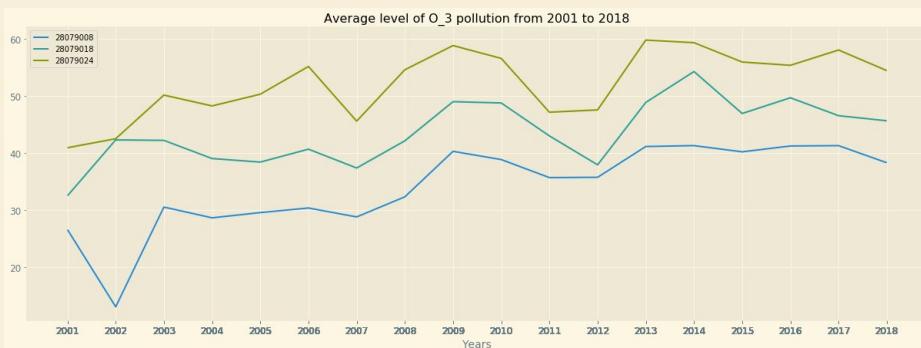
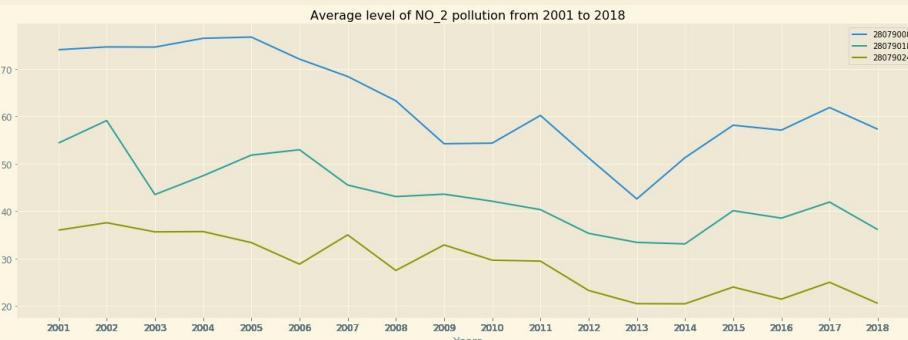
SO₂: sulphur dioxide level measured in $\mu\text{g}/\text{m}^3$.
High levels of sulphur dioxide can produce irritation in
the skin and membranes, and worsen asthma or heart
diseases in sensitive groups.



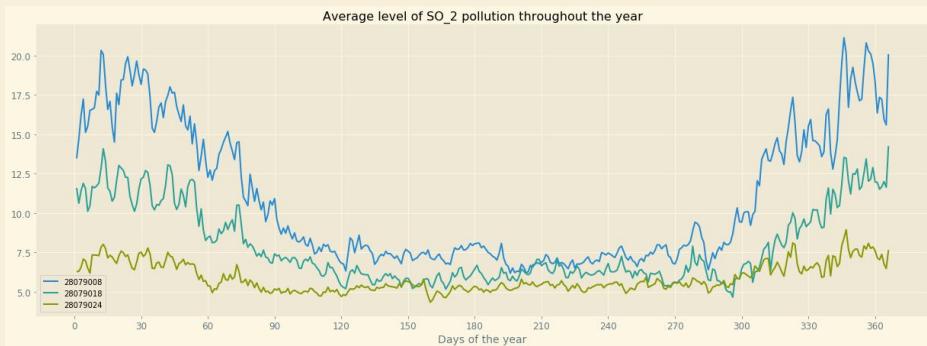
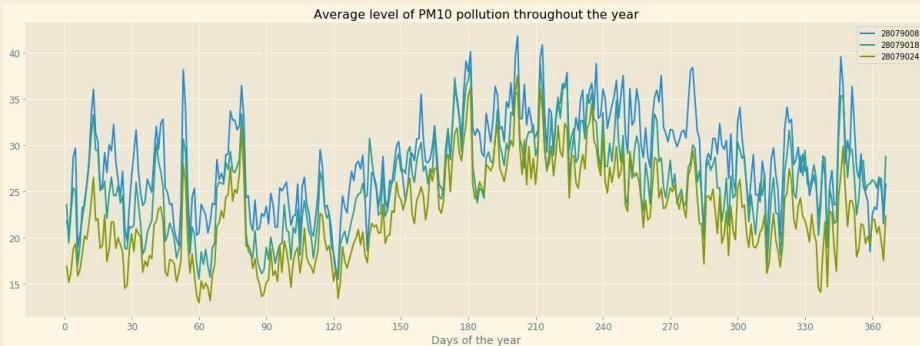
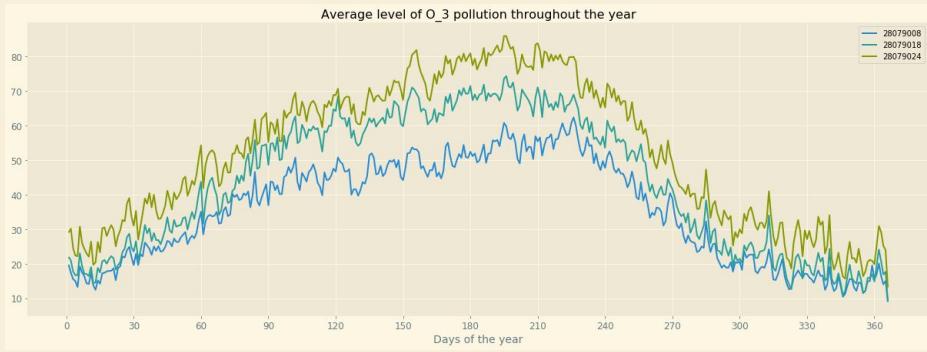
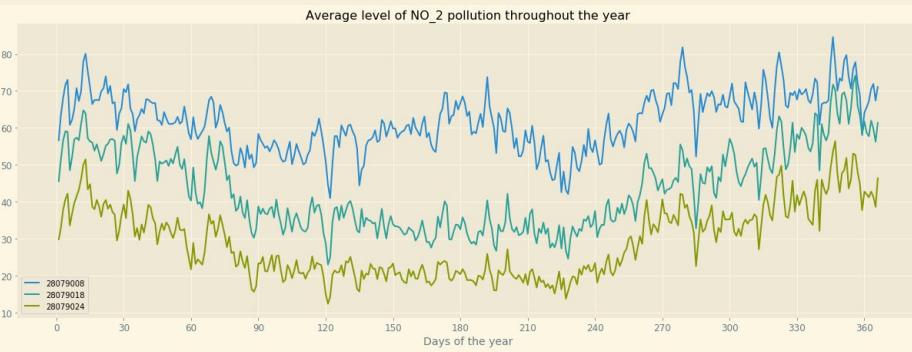


Map data © OpenStreetMap contributors

Let's start by looking at the 3 main components in any time series analysis: trend, seasonality and noise...

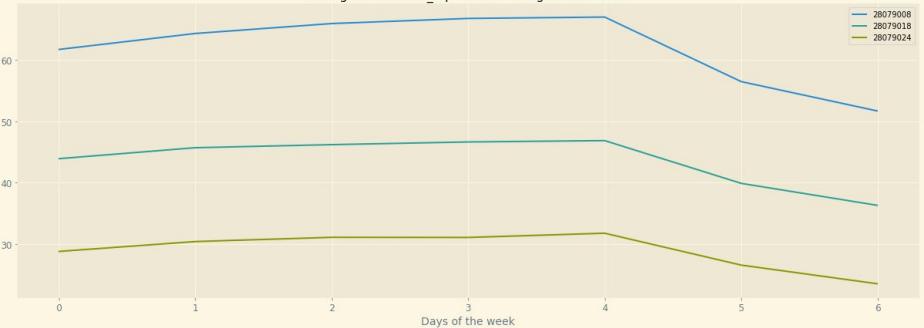


Do you see any significant patterns in the data throughout the year? What differences can you spot?

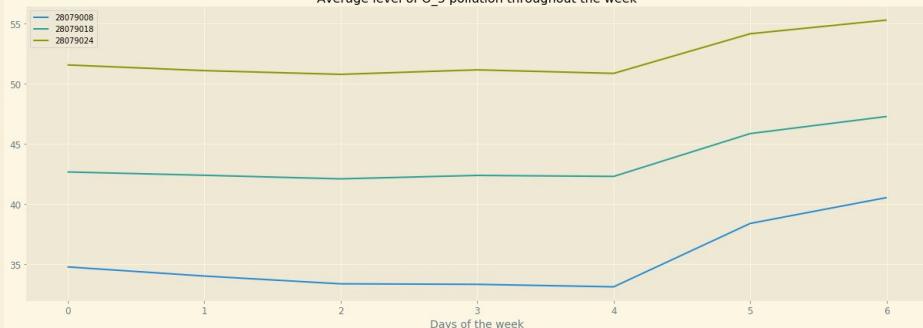


What about throughout the week? Can you see the difference between weekdays and weekends?

Average level of NO₂ pollution throughout the week



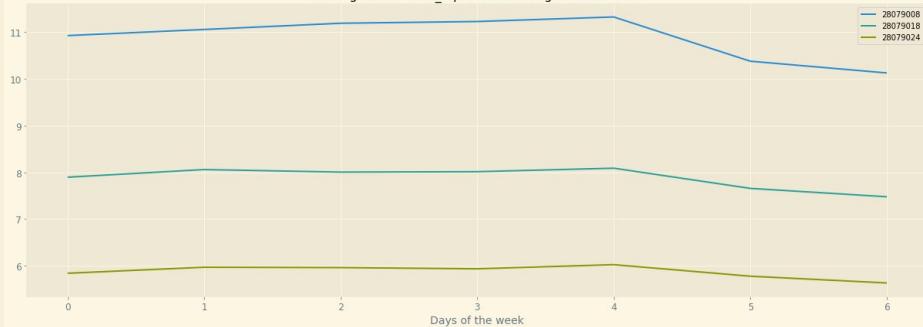
Average level of O₃ pollution throughout the week



Average level of PM10 pollution throughout the week

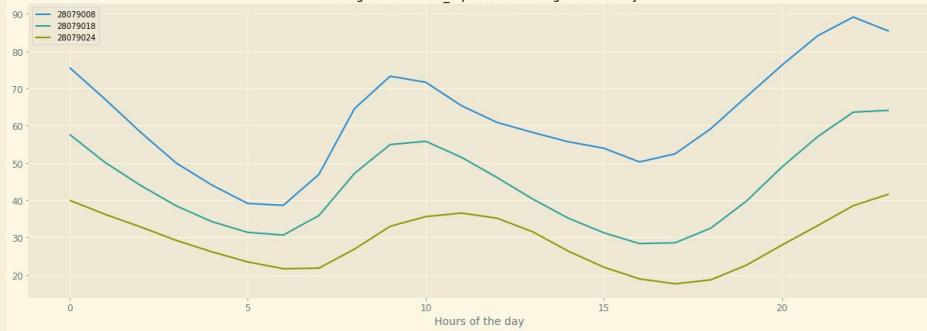


Average level of SO₂ pollution throughout the week

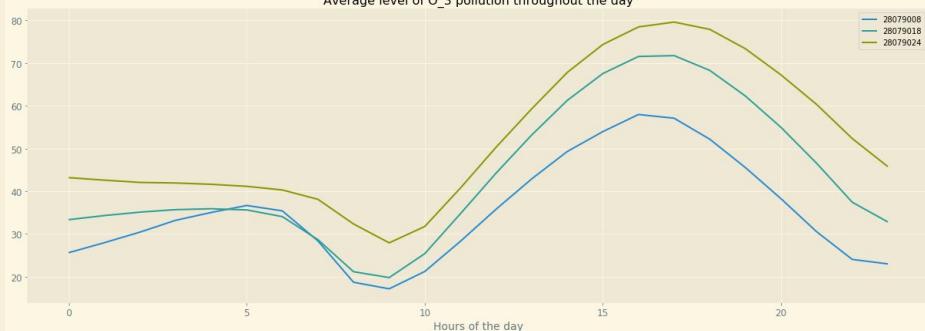


Can you spot the intra day patterns? What are the peak hours for pollution in the day?

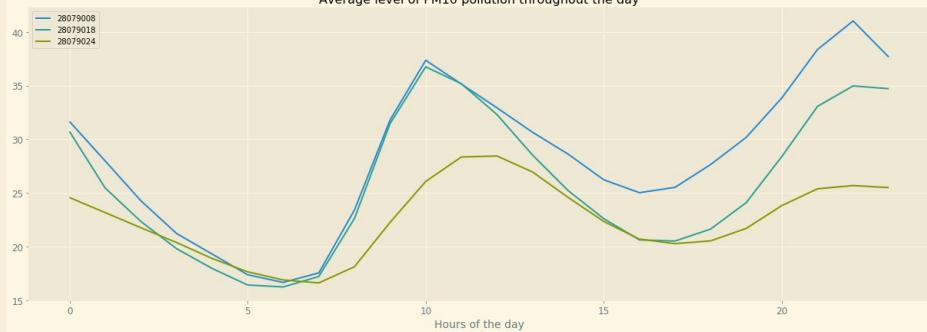
Average level of NO₂ pollution throughout the day



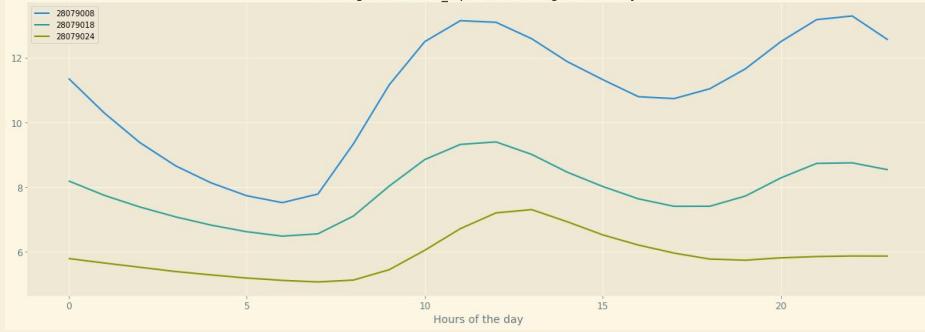
Average level of O₃ pollution throughout the day



Average level of PM10 pollution throughout the day



Average level of SO₂ pollution throughout the day







PREDICTIONS

Let's look at the rolling weekly predictions for one particle and one station, using an ARIMA model.

The particle is:

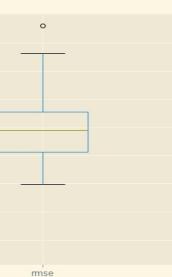
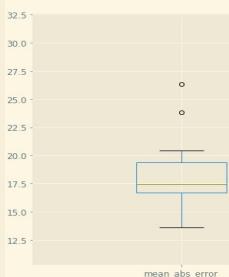
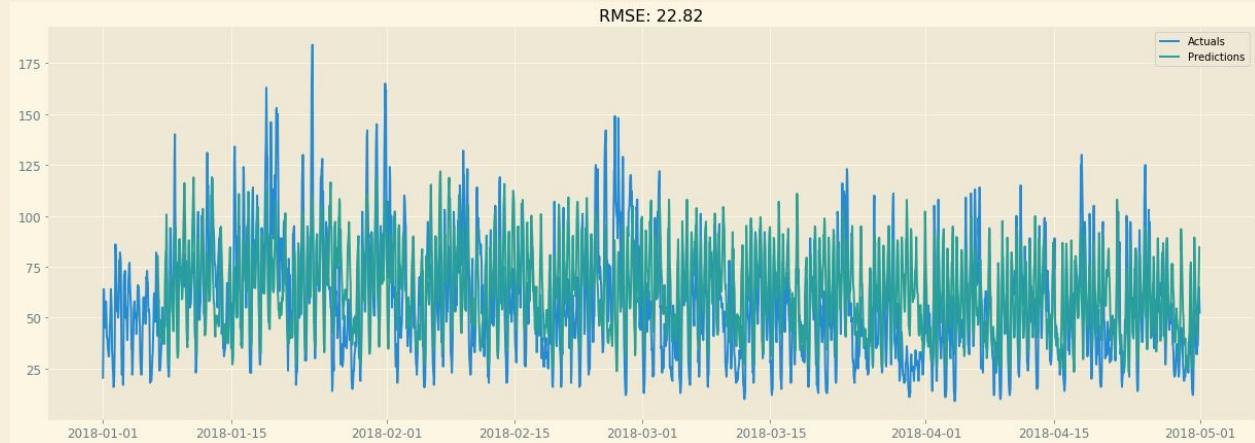
Nitrogen Dioxide (NO₂)

The location is:

Station ID: 28079008
Retiro/ Escuelas Aguirre

The model is:

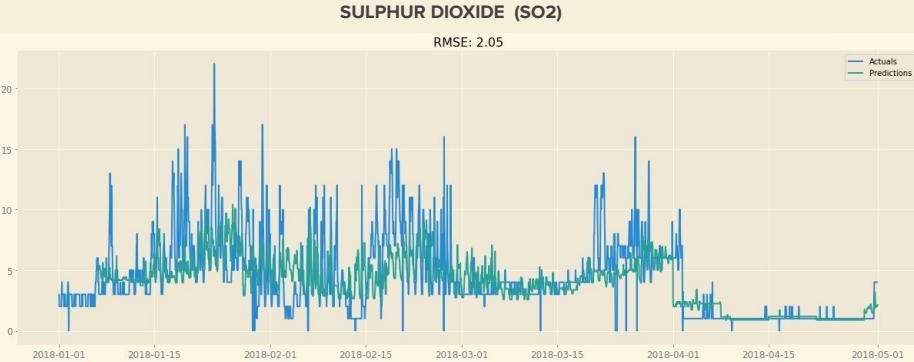
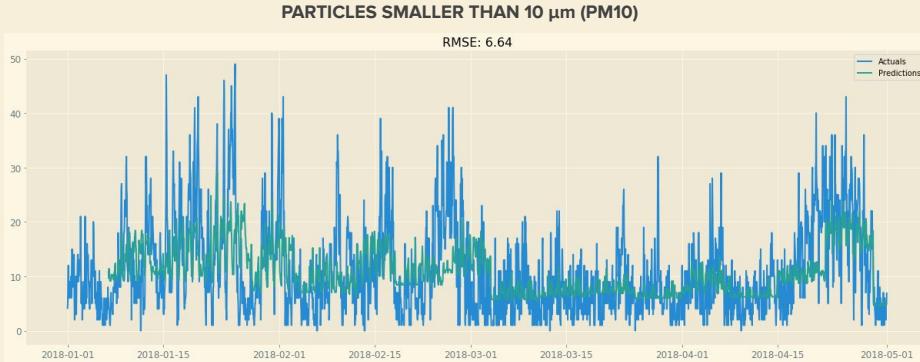
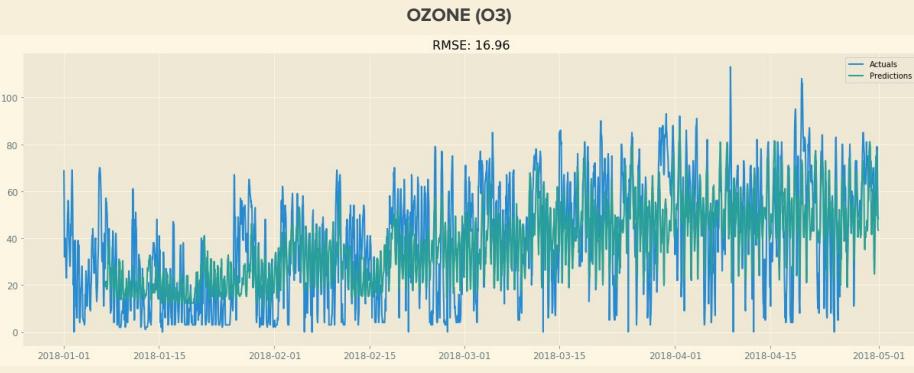
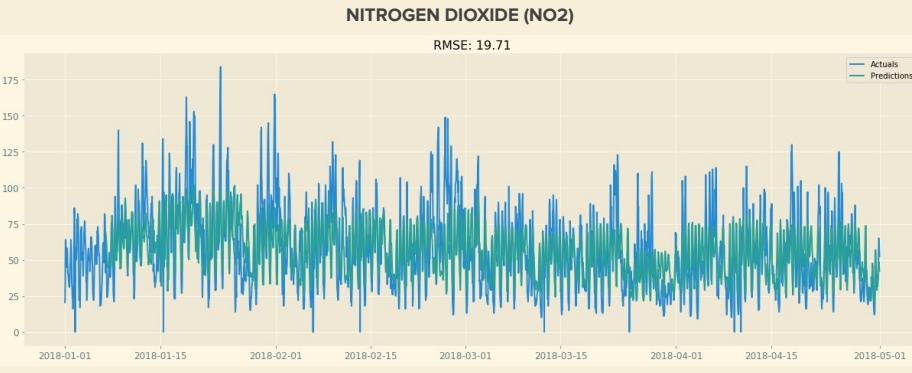
Auto ARIMA



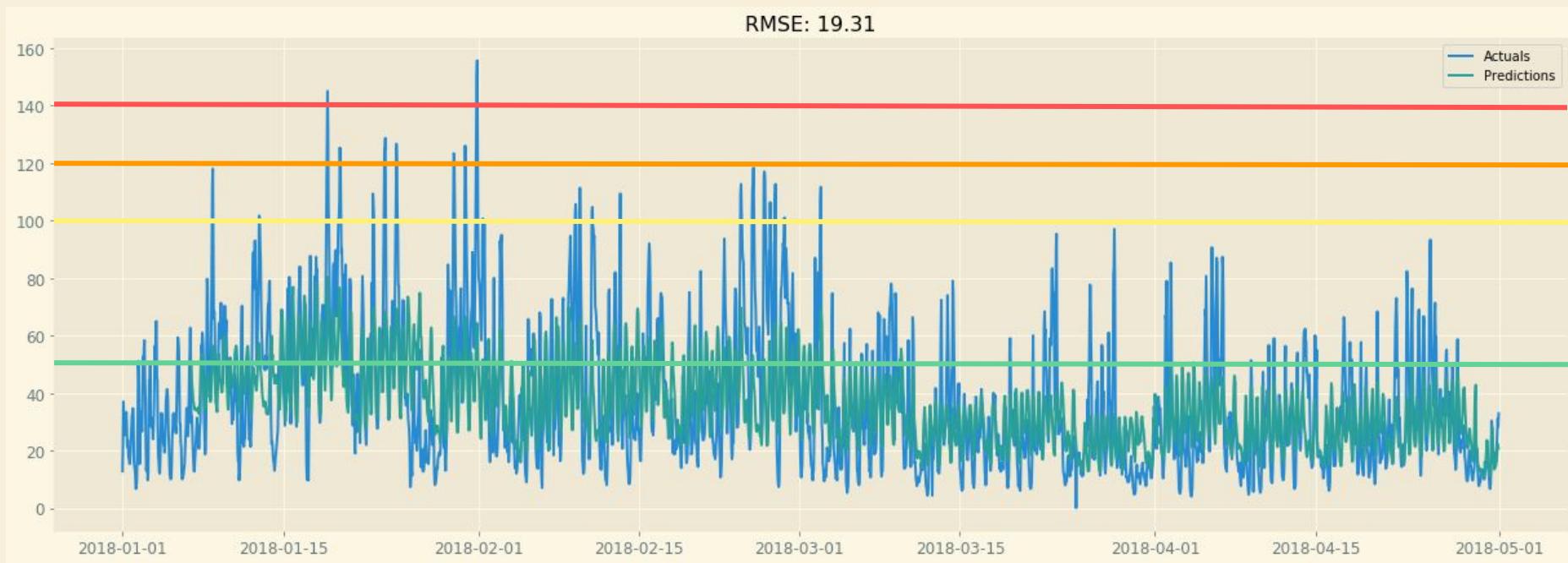
Now look at the rolling weekly predictions for the same particle and station, using an XGBoost model.



Since the XGBoost model gave us the best results for NO₂, let's see how it performed for other particles.



Finally, let's compare the average of the 3 stations with the WHO thresholds for NO₂ contamination.



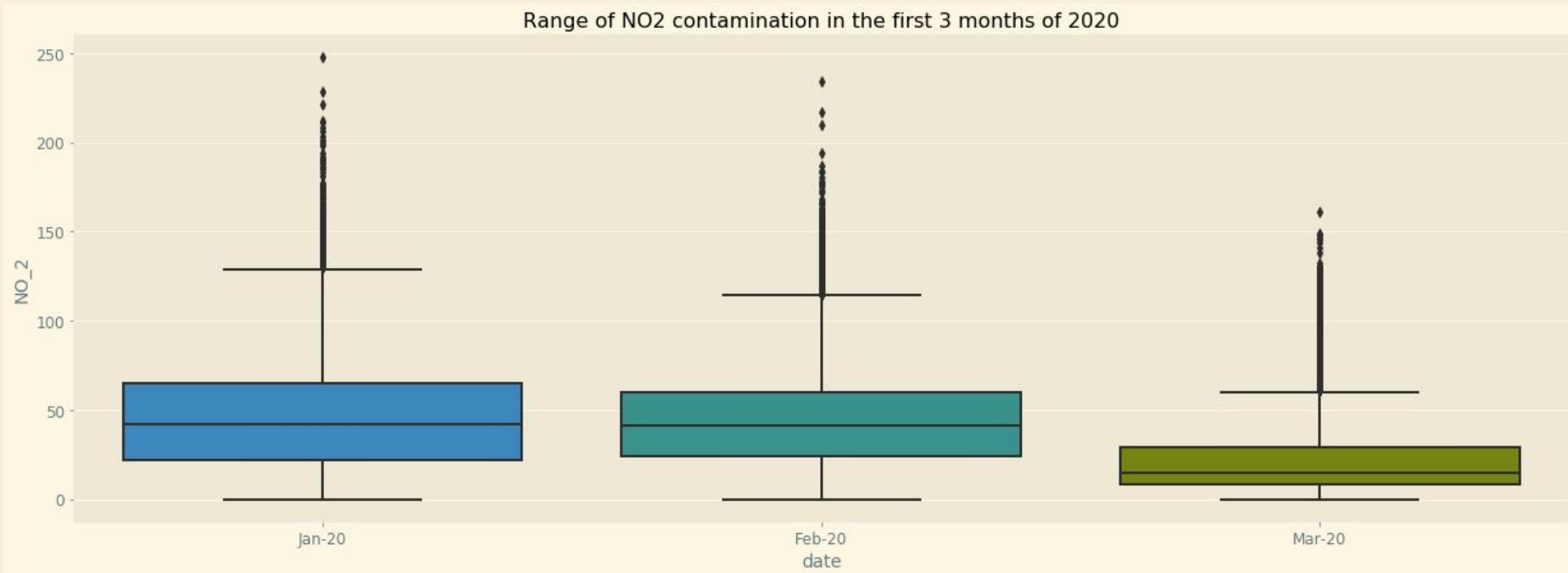
3 key take aways:

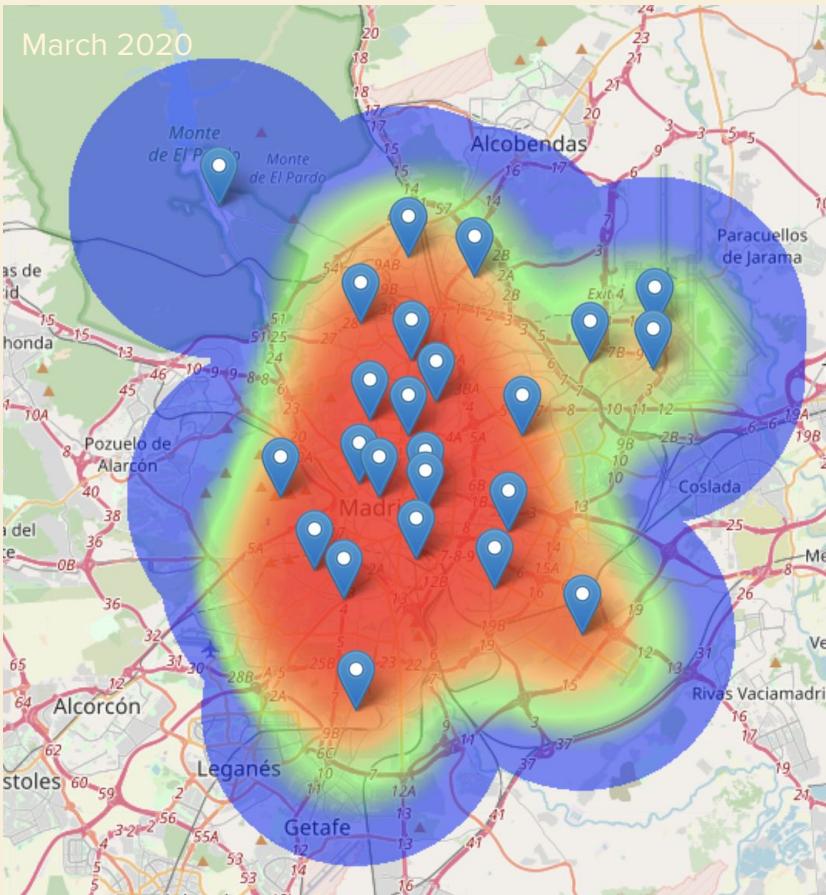
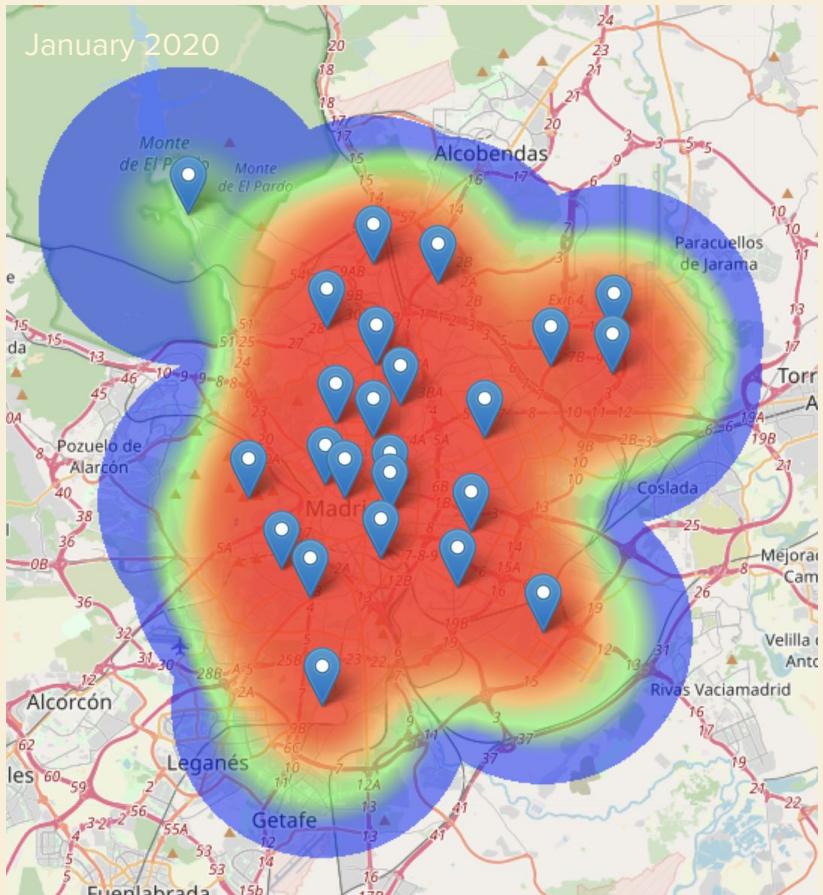
1. Most data is time series in its nature. Whether we want to predict trends in financial markets, sales of a product or levels of pollution, time is an important factor and should be factored into our models.
2. Both models perform better on stationary data. It's important to analyse any trends and patterns and transform the time series data if needed, so that it becomes stationary. Using non-stationary time series data produces unreliable results and leads to poor forecasting.
3. Any model has pros and cons. In this case, the XGB model had better overall performance while the ARIMA performed better at predicting dangerous levels of the particle. We can optimise our model by using an evaluation metric that penalises these bad predictions.

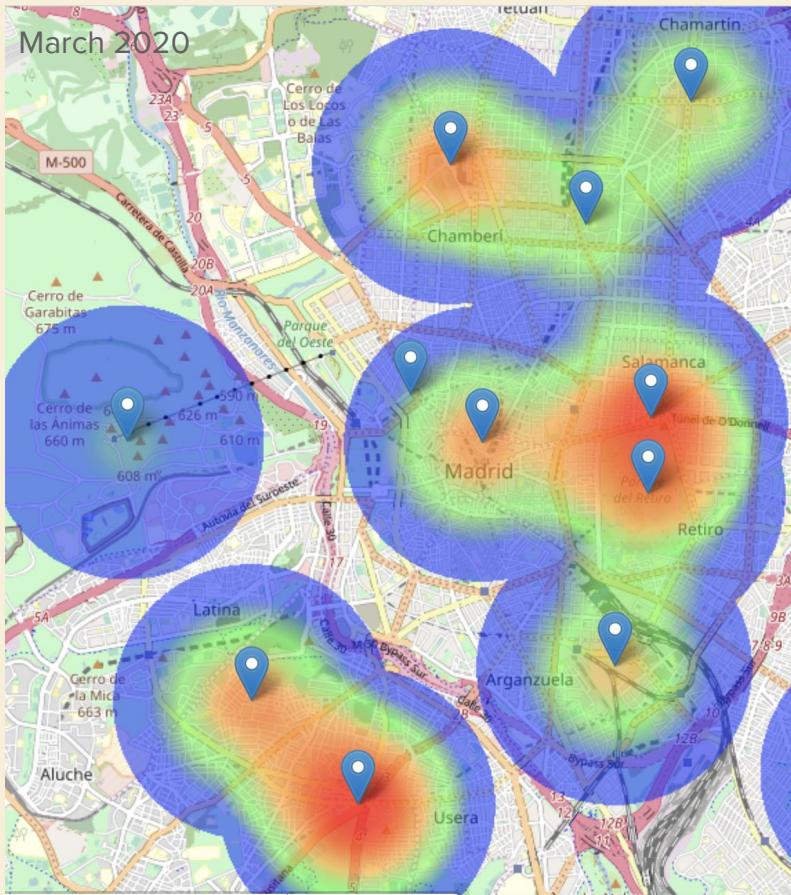
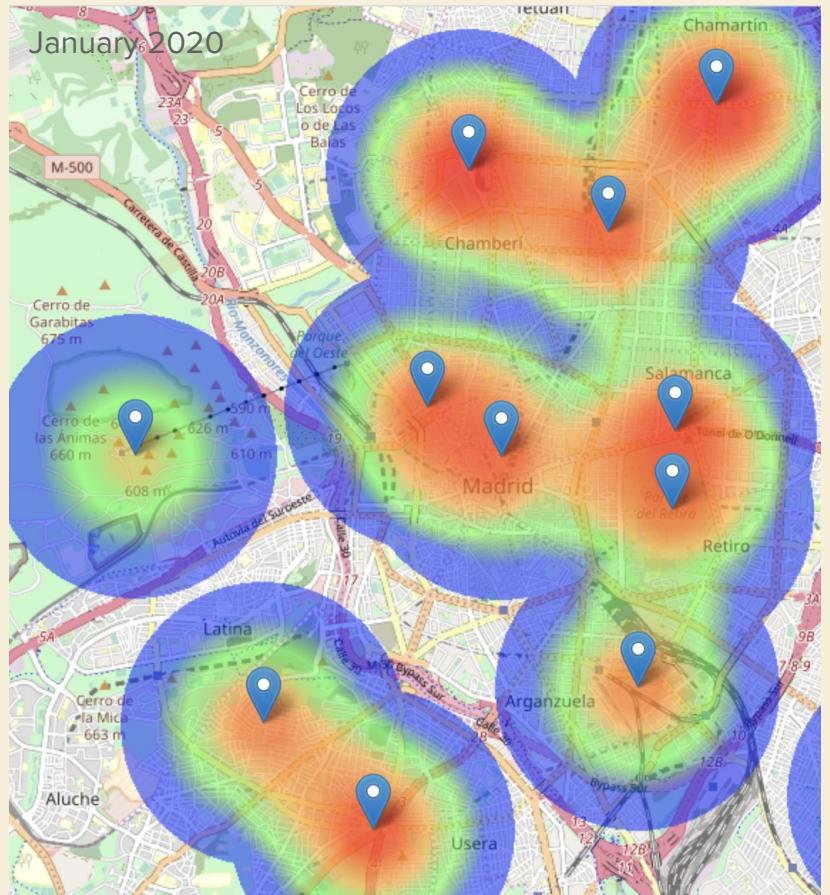
BONUS

30

Since the start of the Covid-19 crisis, levels of NO2 contamination have dropped significantly! 🌳🌍



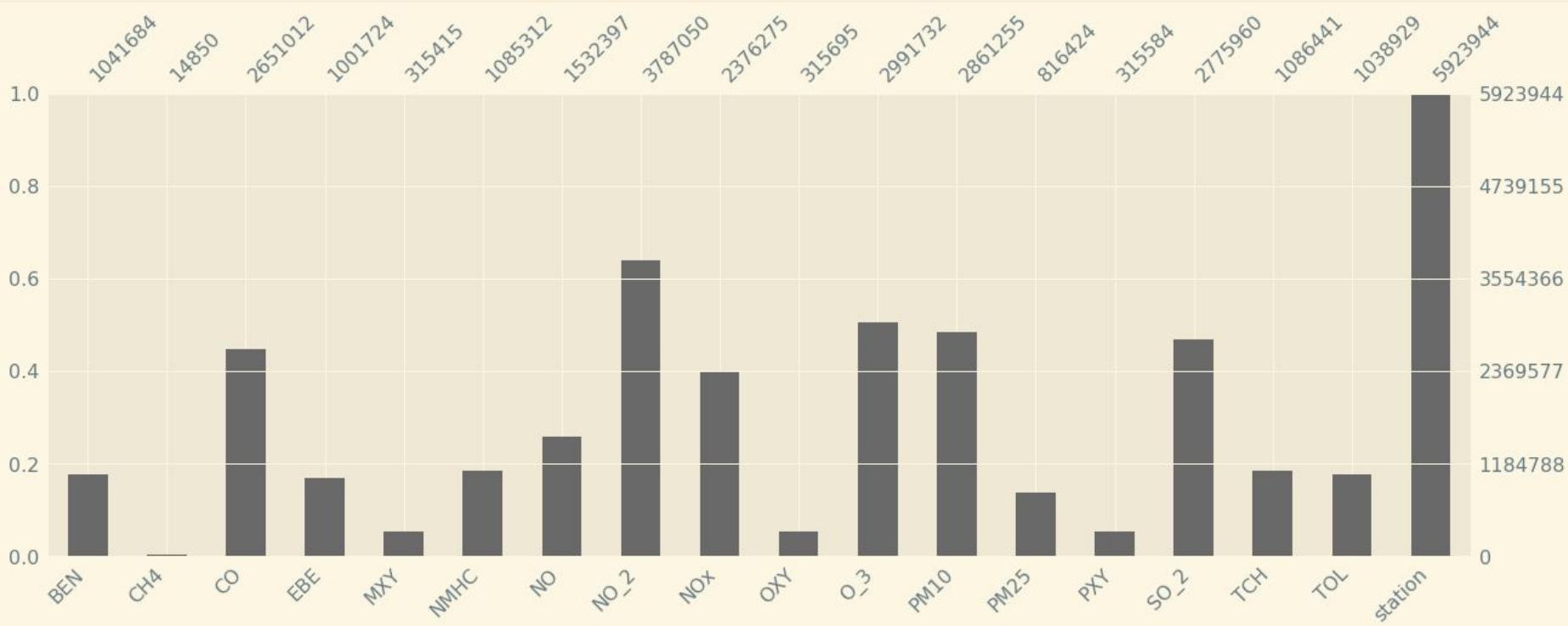




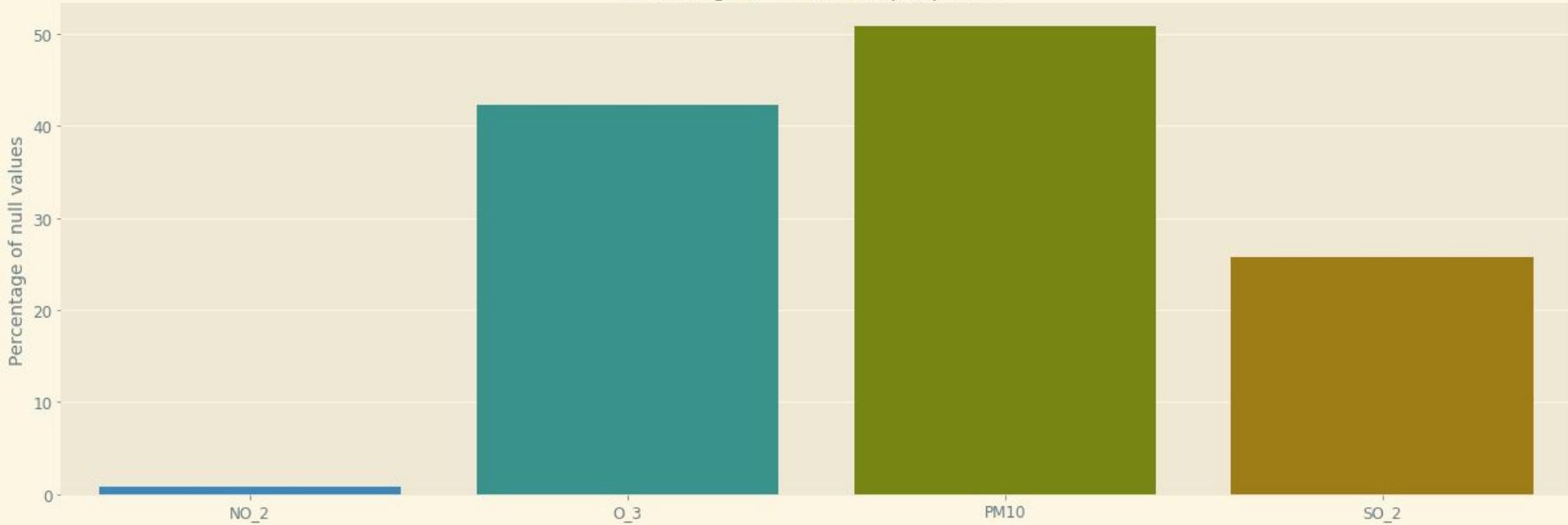
An aerial photograph of the Paris skyline during sunset or sunrise. The foreground shows the dense urban architecture of the 16th arrondissement, with its characteristic blue-tiled roofs. A wide, tree-lined boulevard cuts through the city, leading towards the horizon where the modern skyscrapers of La Défense are silhouetted against a bright sky. The Eiffel Tower is visible on the right side of the frame.

THANK YOU

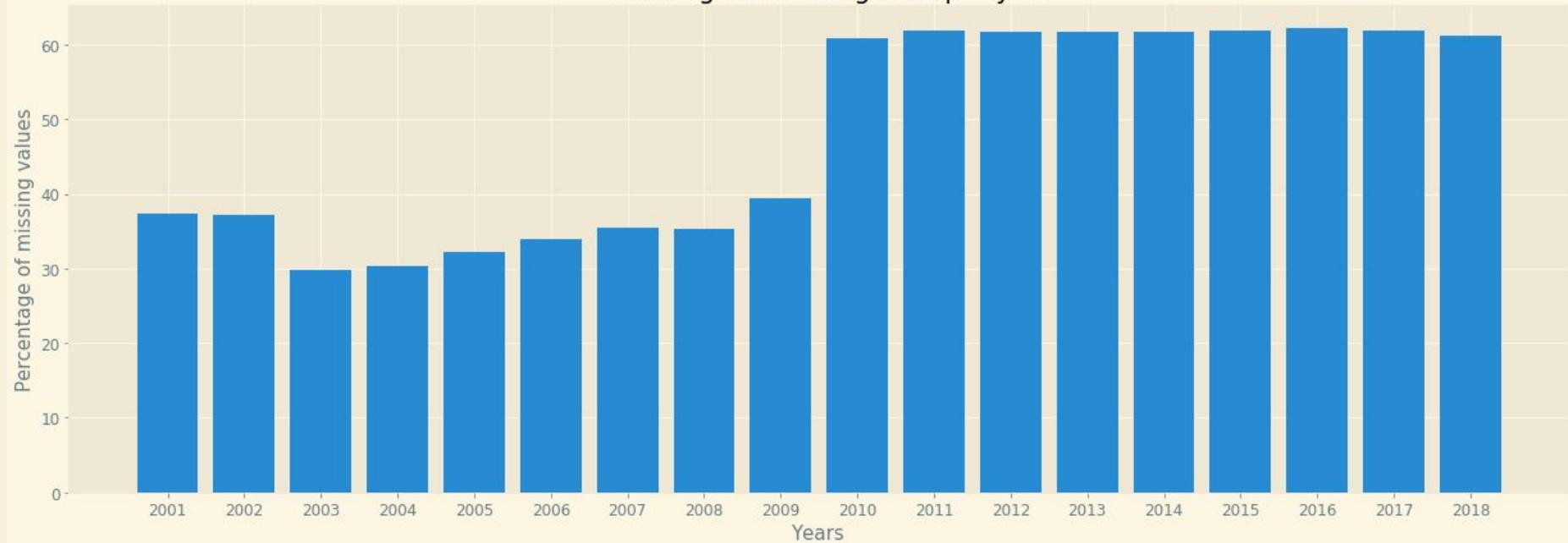
Missing Values



Percentage of null values per particle



Percentage of missing data per year



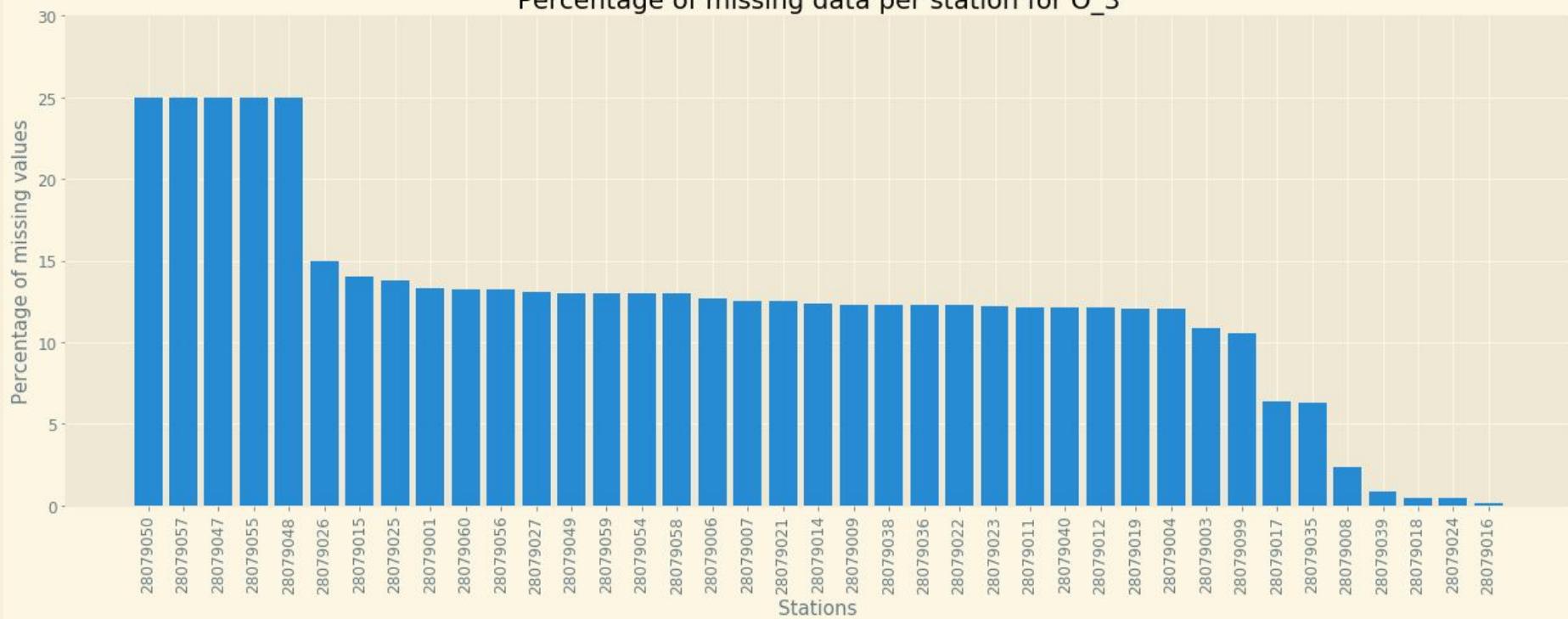
Percentage of missing data per station



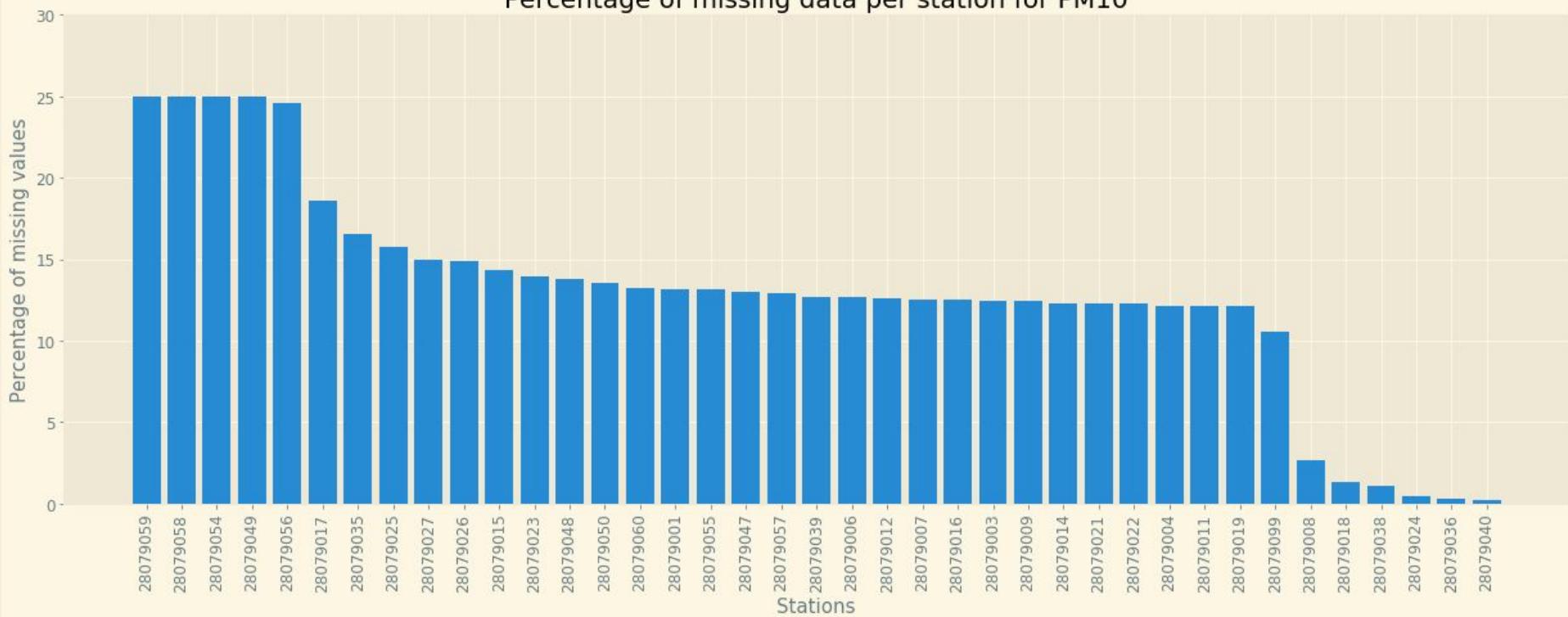
Percentage of missing data per station for NO_2



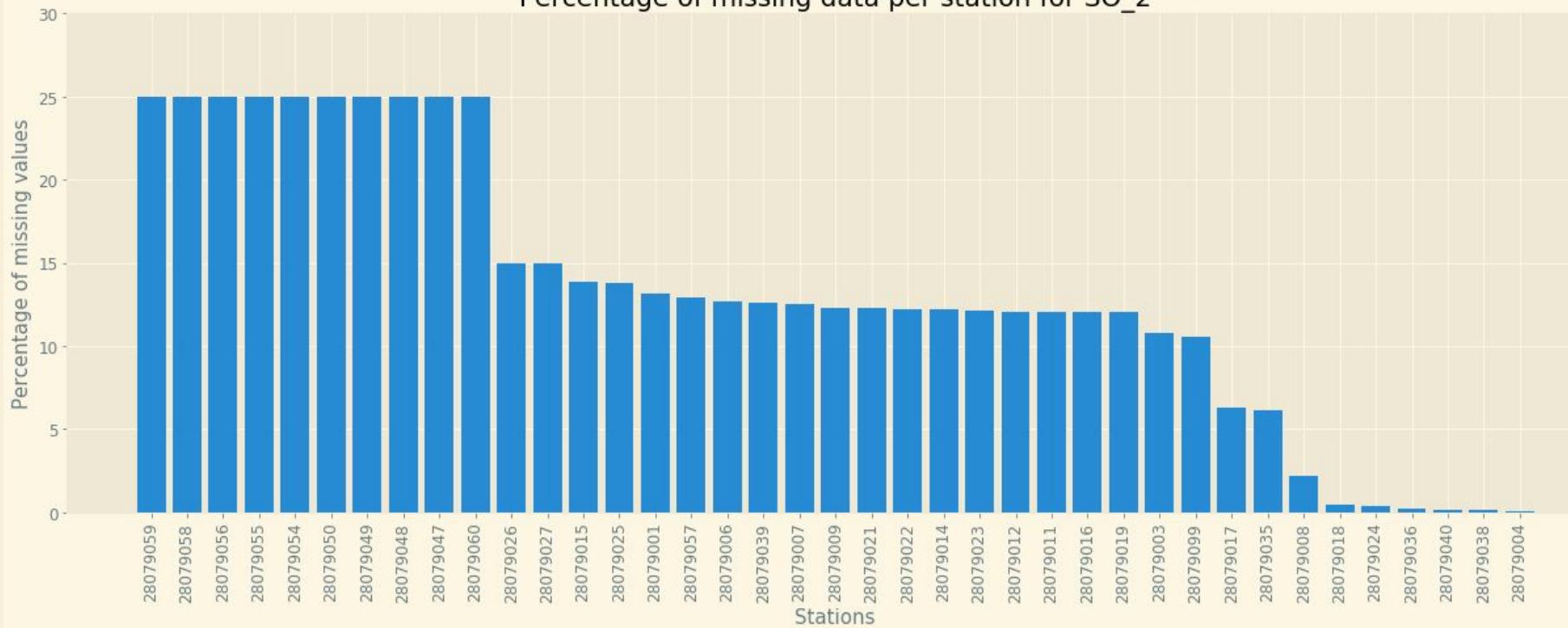
Percentage of missing data per station for O_3



Percentage of missing data per station for PM10

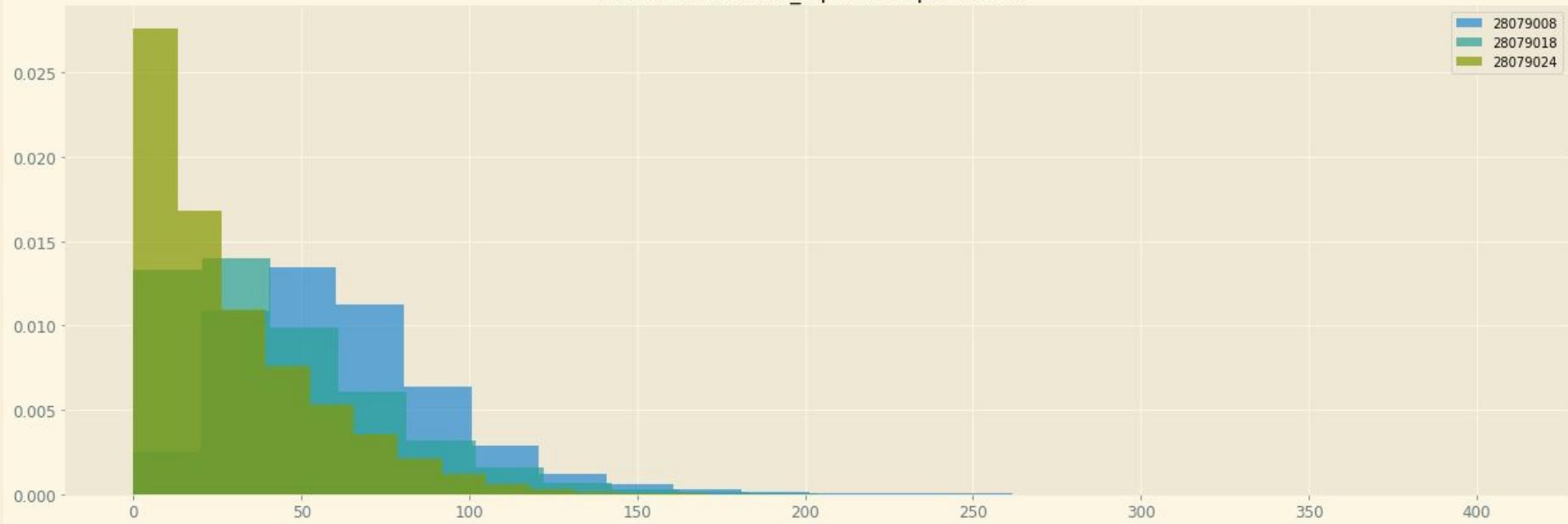


Percentage of missing data per station for SO_2

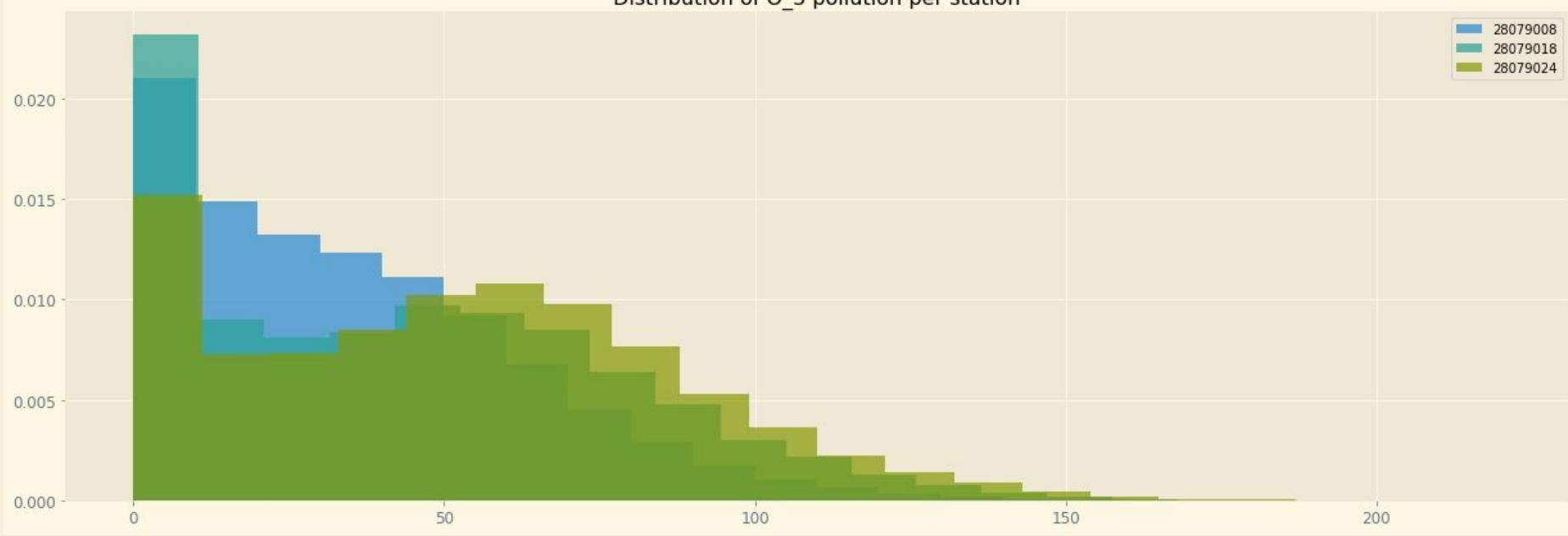


Distribution of the levels of
pollution for each particle

Distribution of NO₂ pollution per station



Distribution of O₃ pollution per station



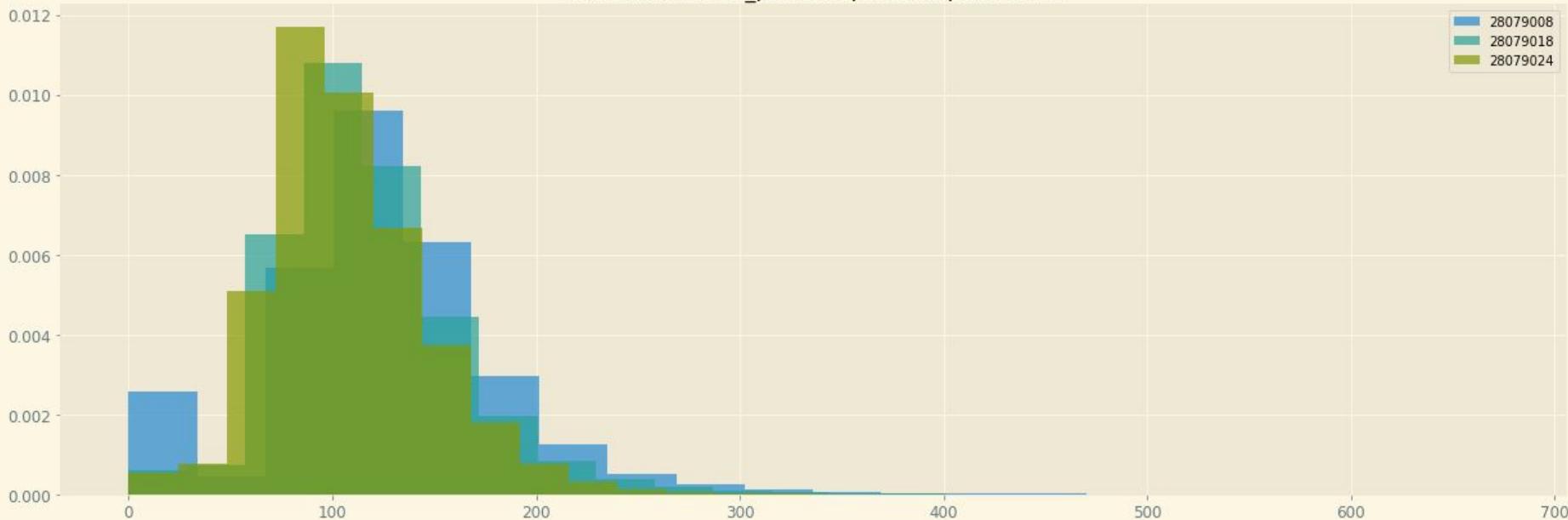
Distribution of PM10 pollution per station



Distribution of SO₂ pollution per station

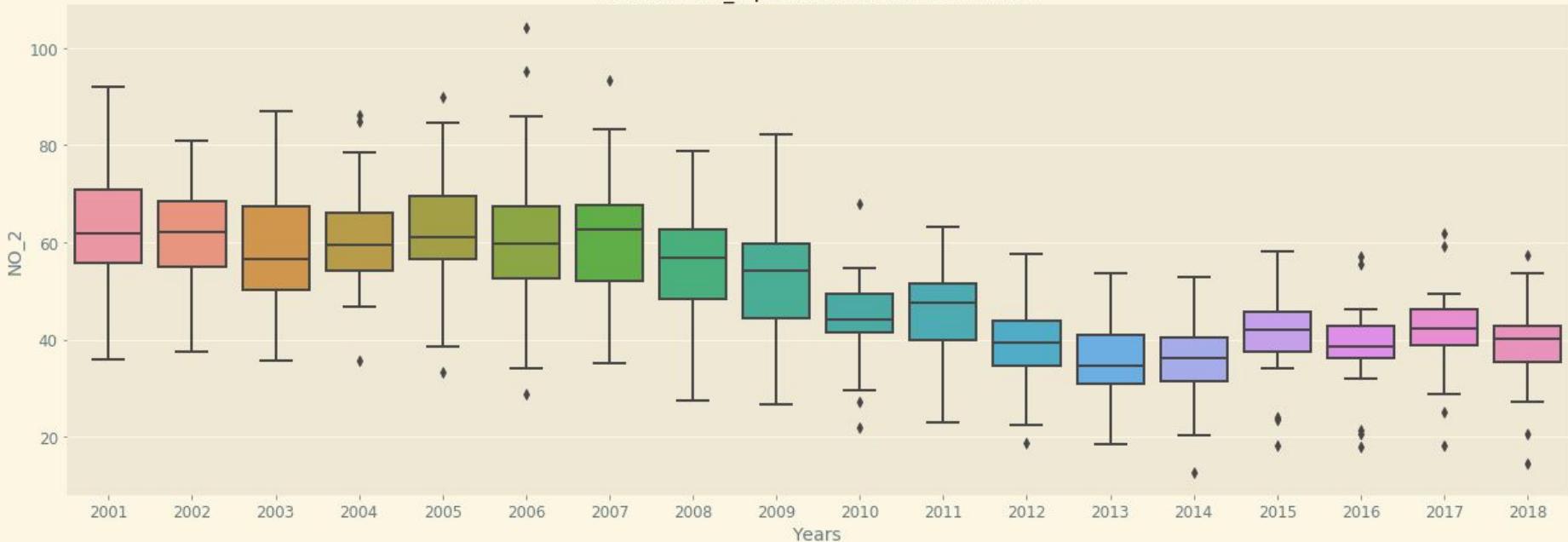


Distribution of all_particles pollution per station

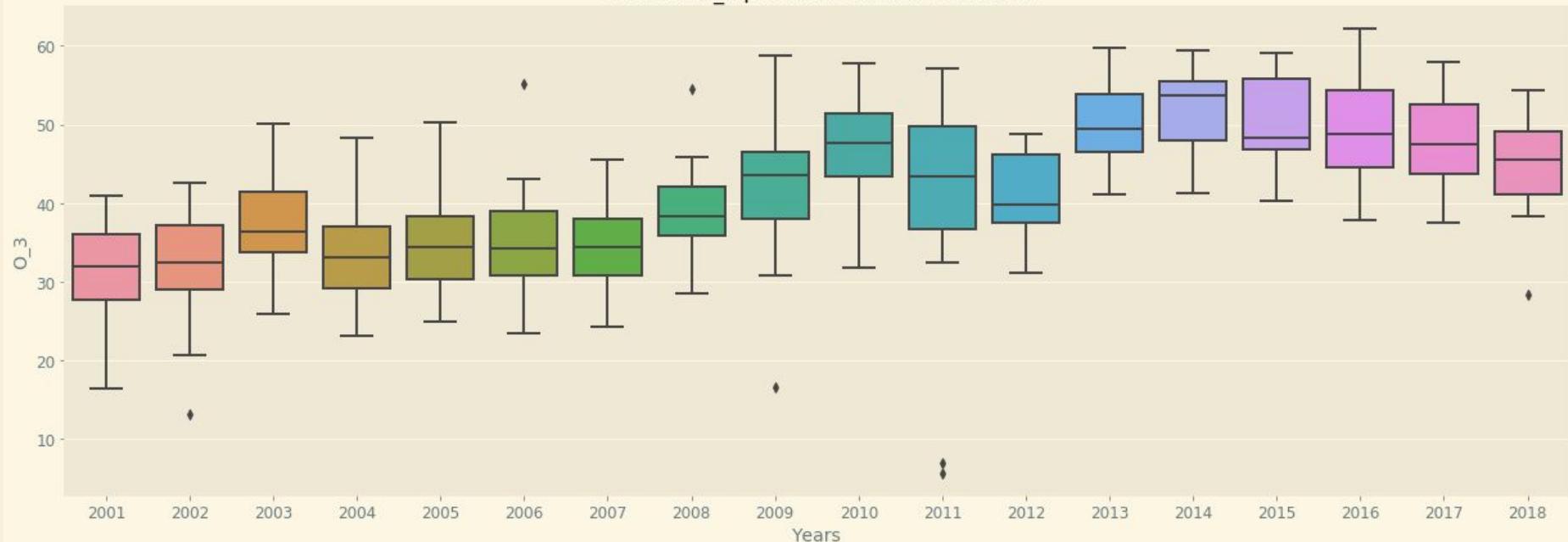


Evolution of the levels of
pollution throughout the years

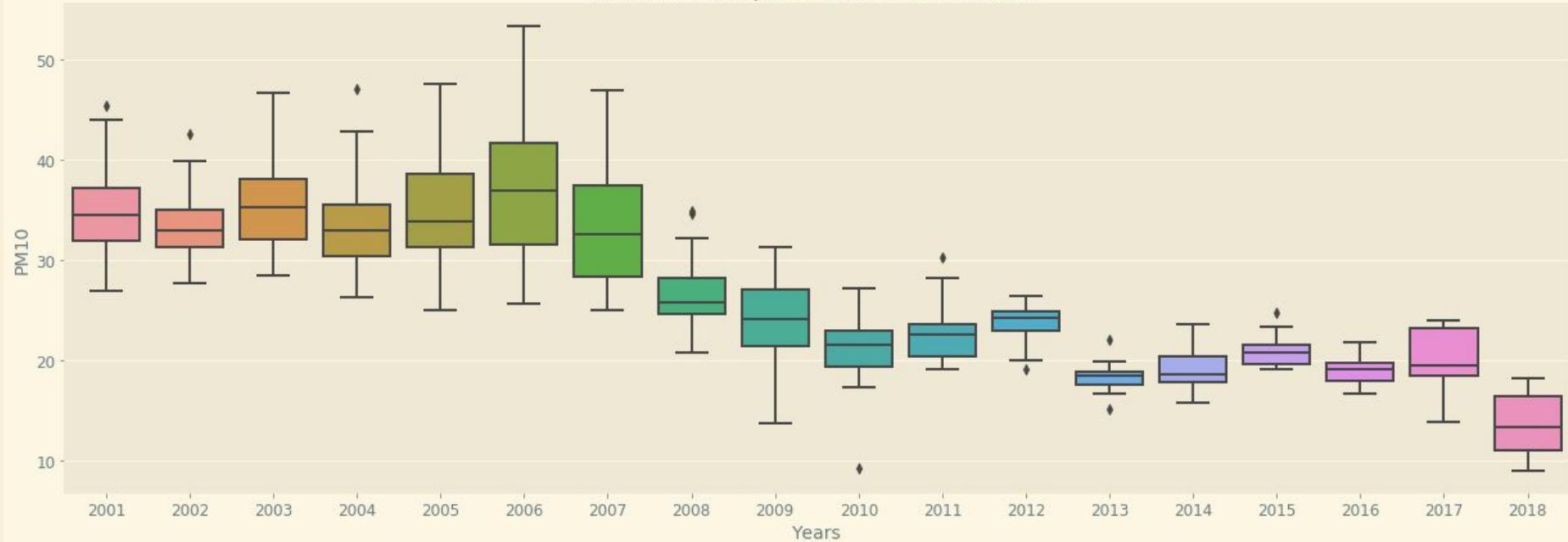
Levels of NO₂ pollution from 2001 to 2018



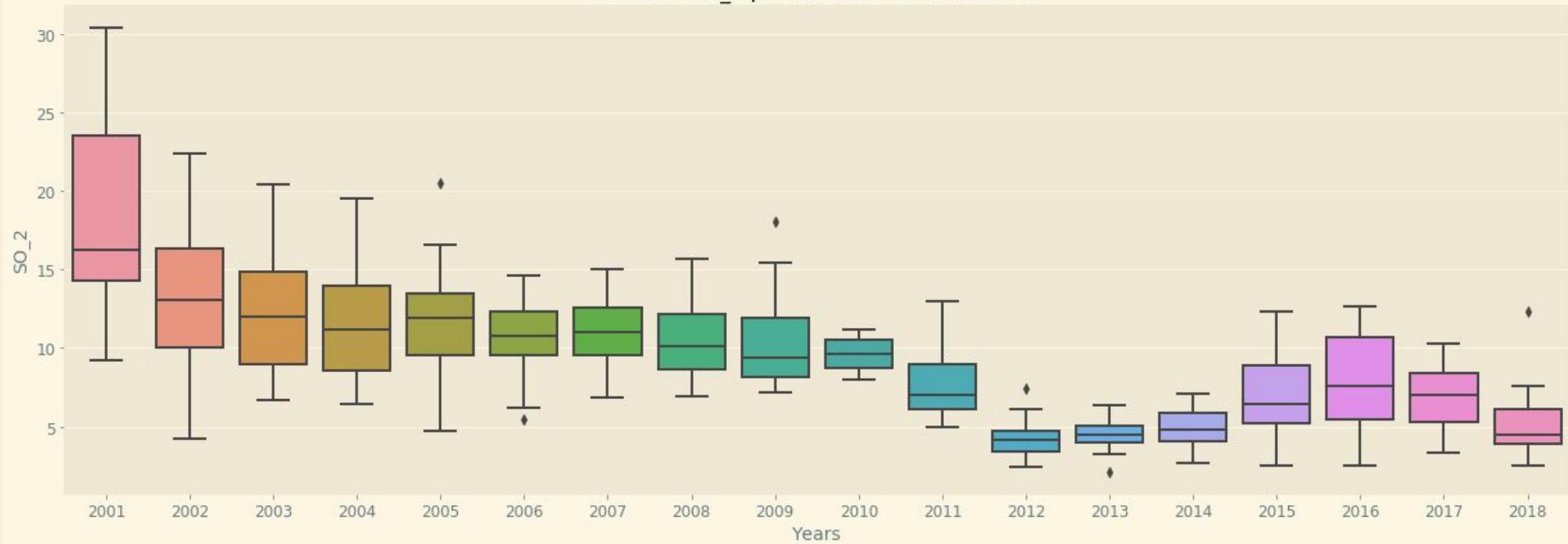
Levels of O₃ pollution from 2001 to 2018



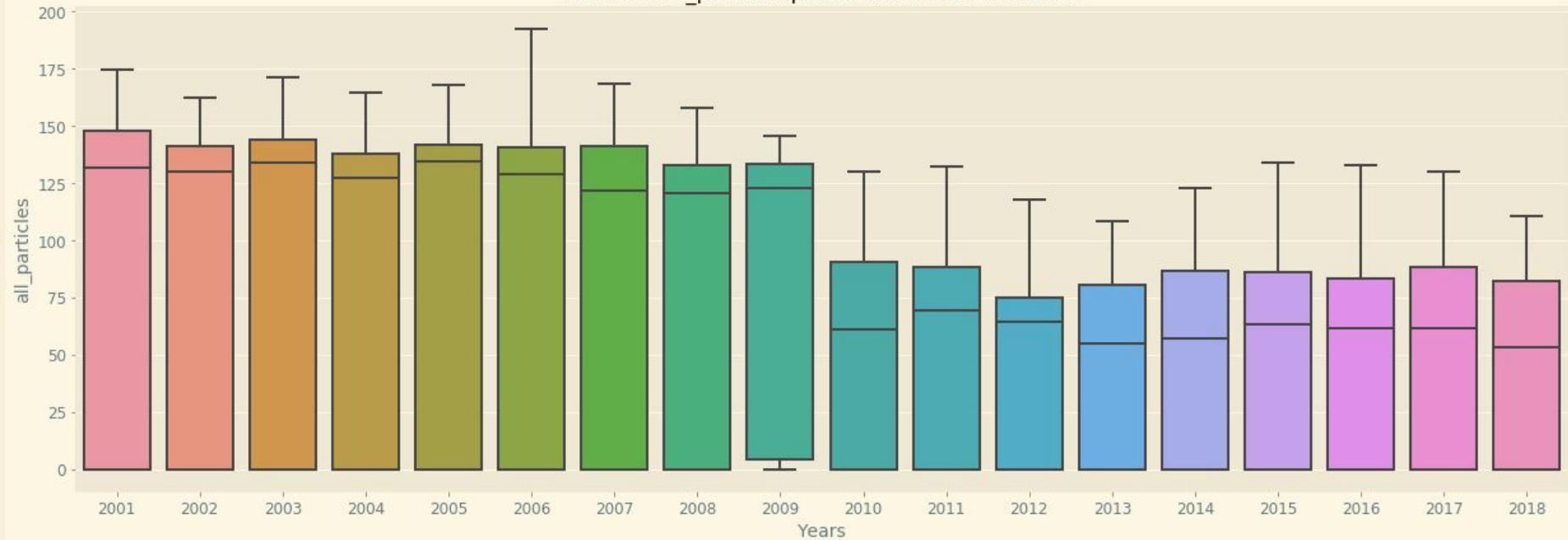
Levels of PM10 pollution from 2001 to 2018



Levels of SO₂ pollution from 2001 to 2018



Levels of all_particles pollution from 2001 to 2018

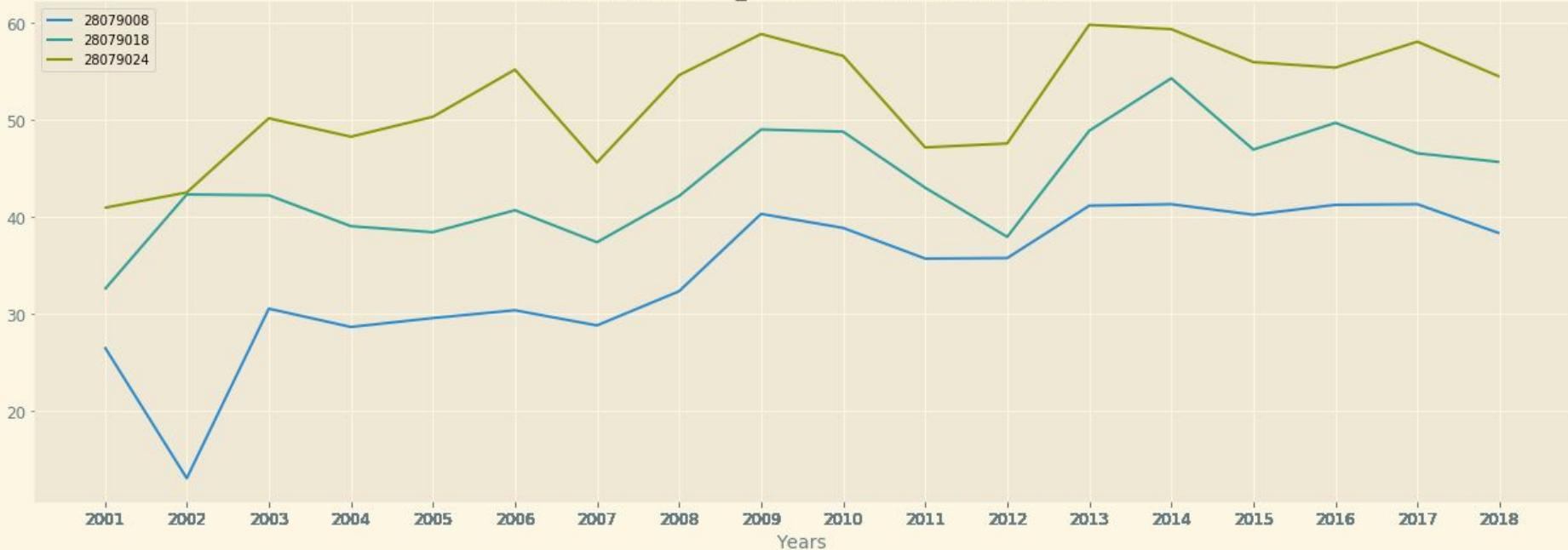


Average levels of pollution
throughout the years

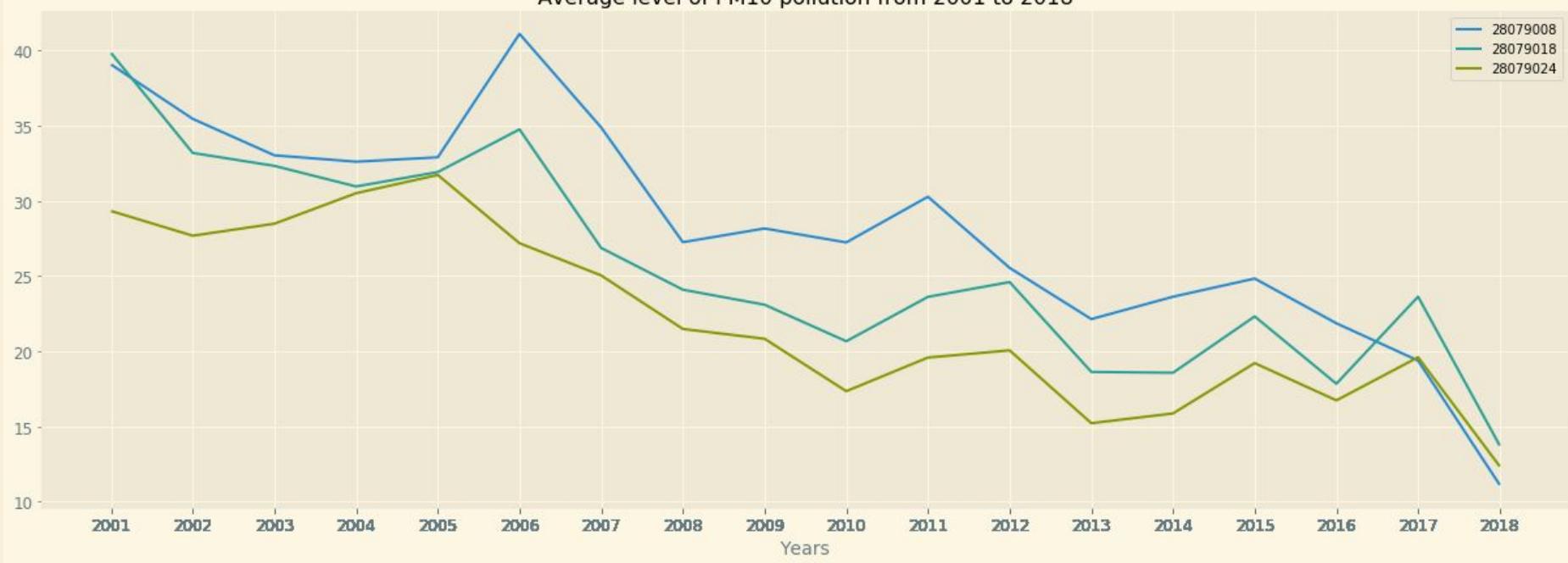
Average level of NO₂ pollution from 2001 to 2018



Average level of O₃ pollution from 2001 to 2018



Average level of PM10 pollution from 2001 to 2018



Average level of SO₂ pollution from 2001 to 2018

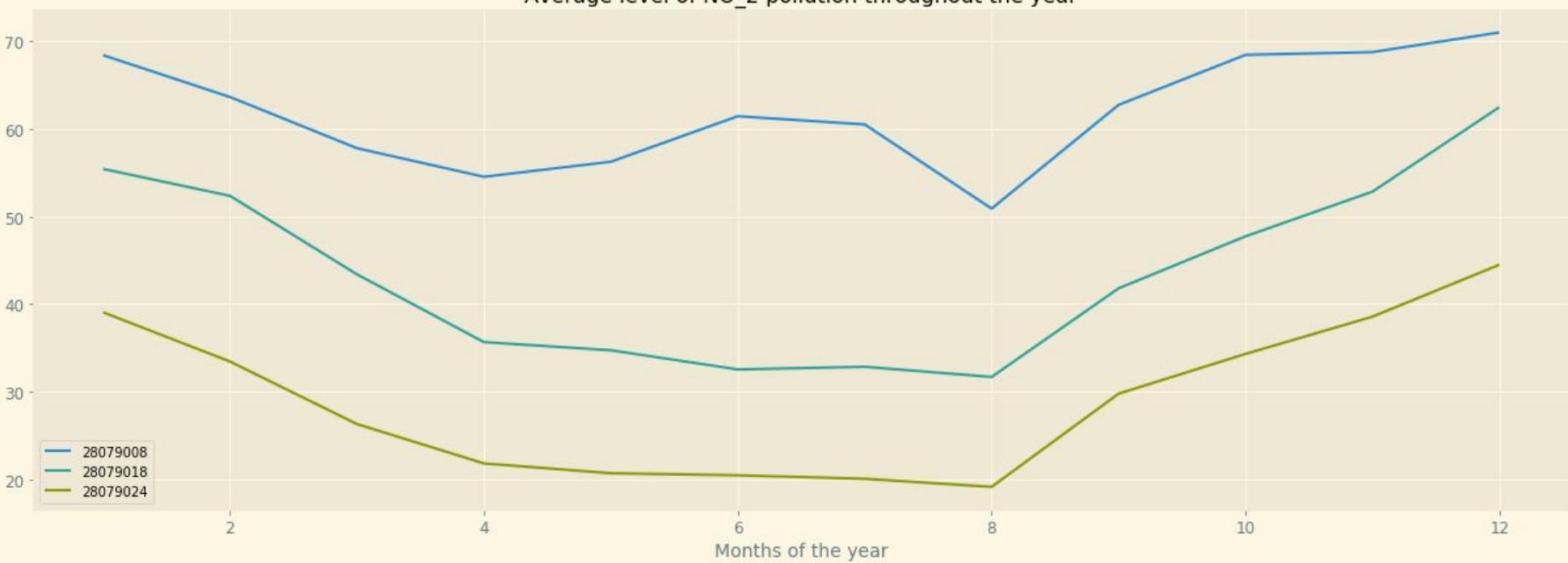


Average level of all_particles pollution from 2001 to 2018

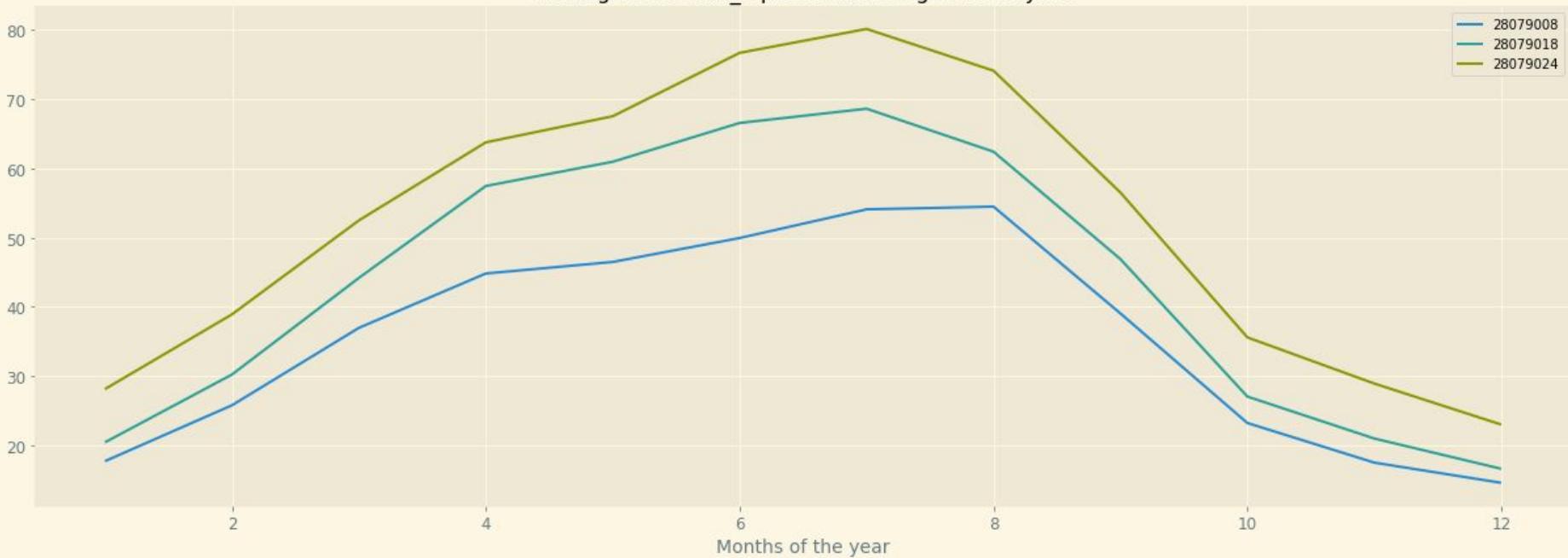


Average levels of pollution
throughout the months of the year

Average level of NO₂ pollution throughout the year



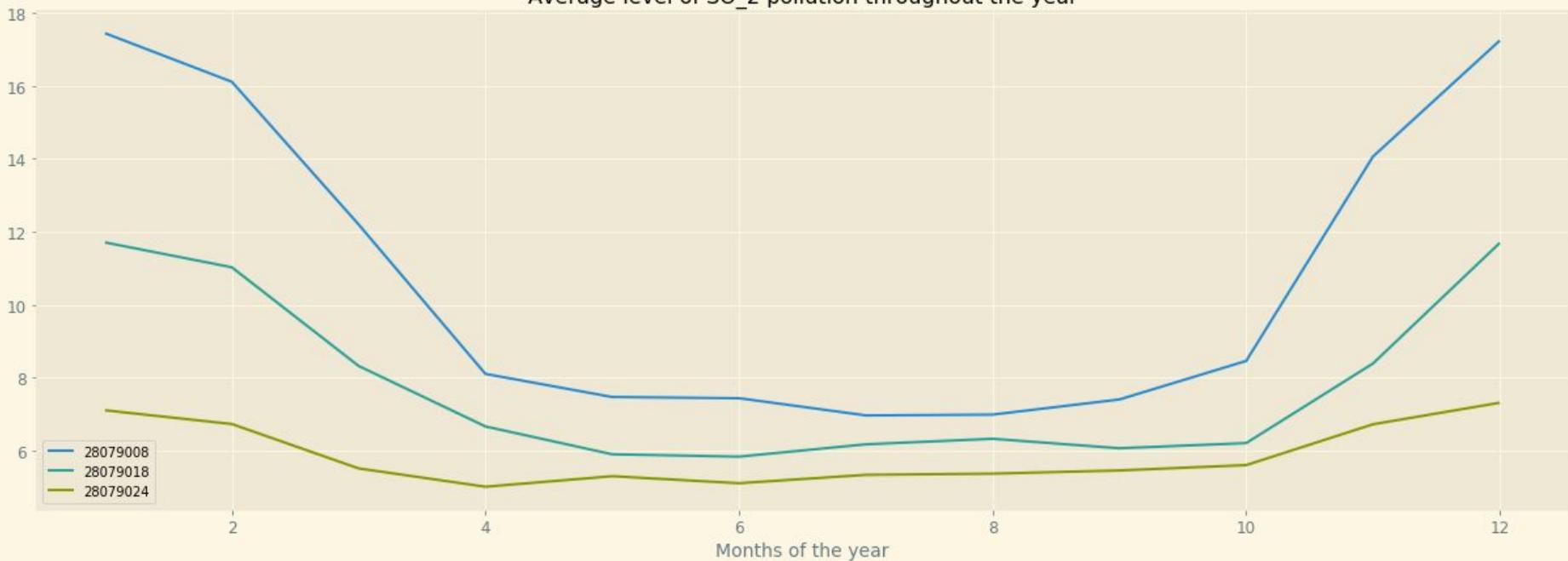
Average level of O₃ pollution throughout the year



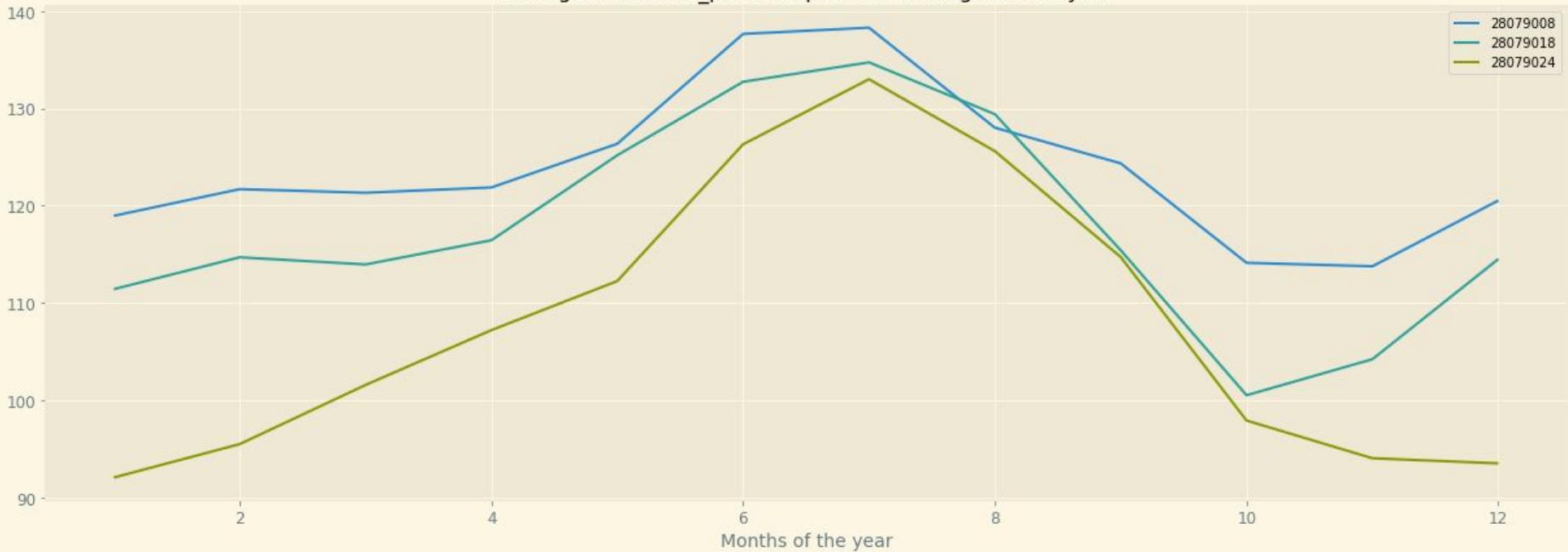
Average level of PM10 pollution throughout the year



Average level of SO₂ pollution throughout the year

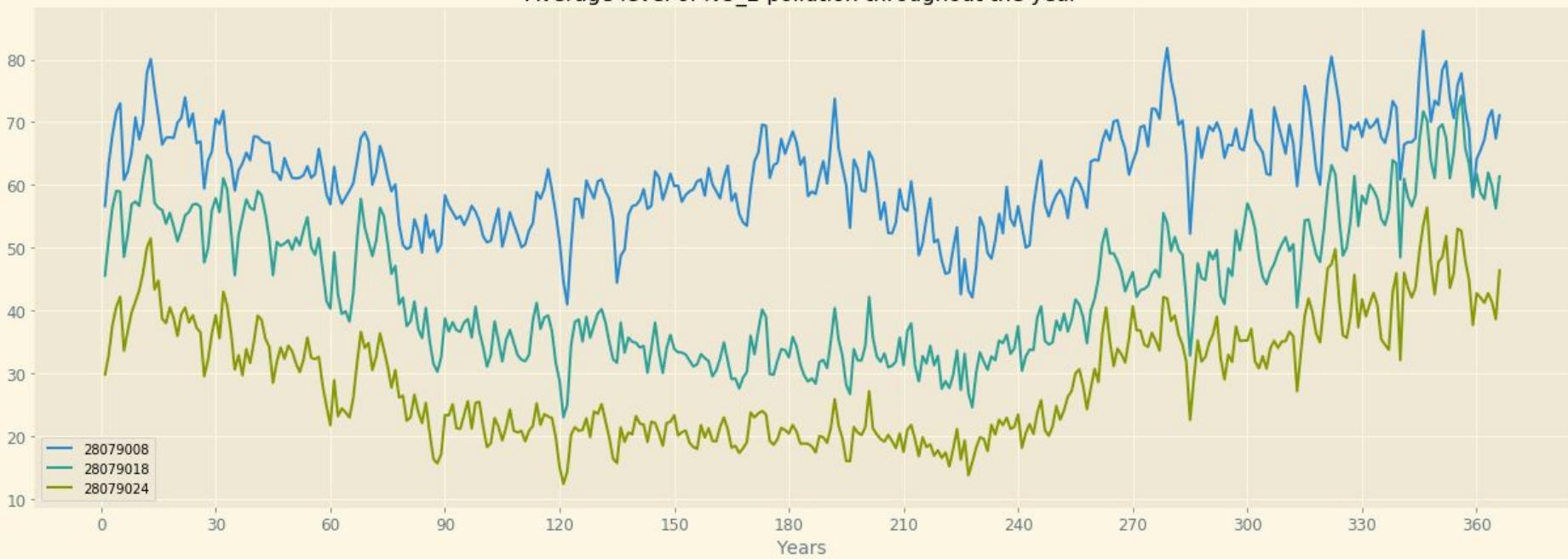


Average level of all_particles pollution throughout the year



Average levels of pollution
throughout the days of the year

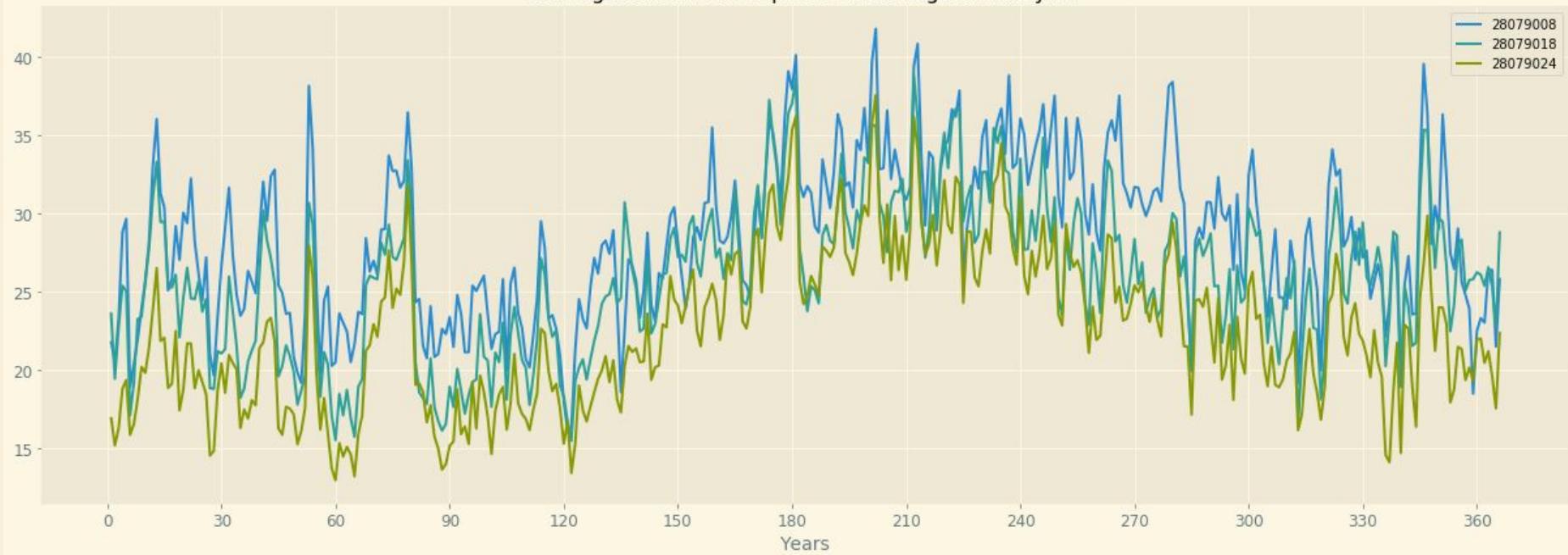
Average level of NO₂ pollution throughout the year



Average level of O₃ pollution throughout the year



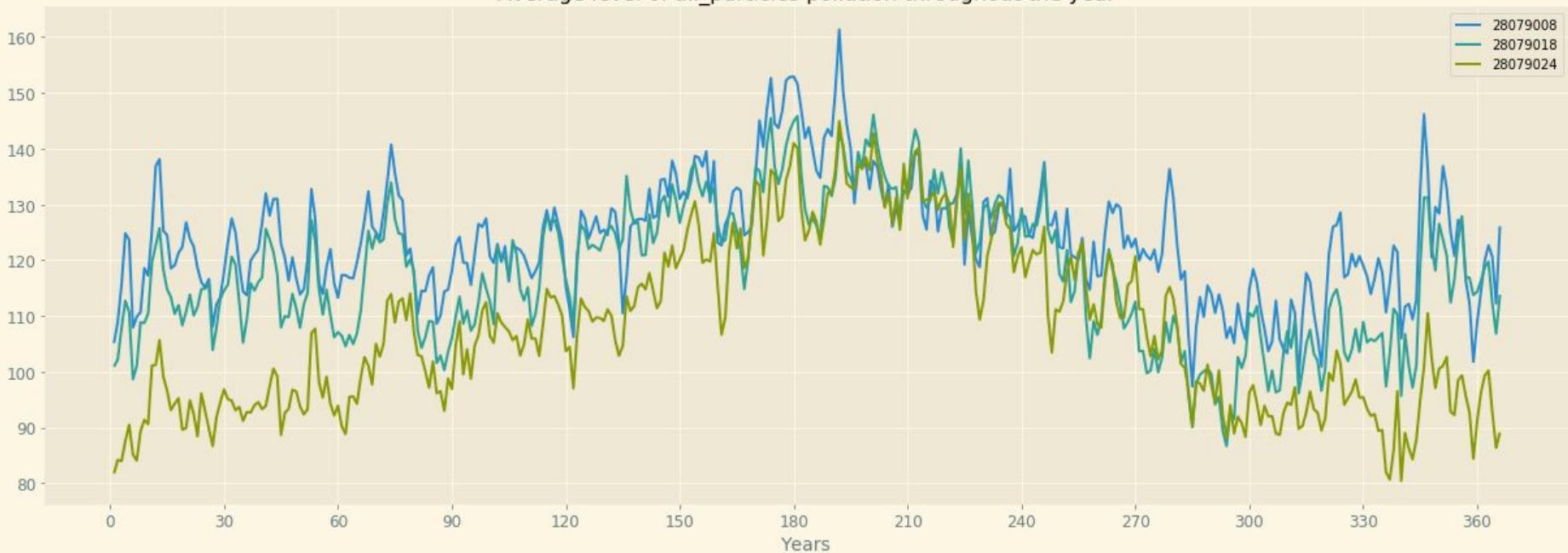
Average level of PM10 pollution throughout the year



Average level of SO₂ pollution throughout the year

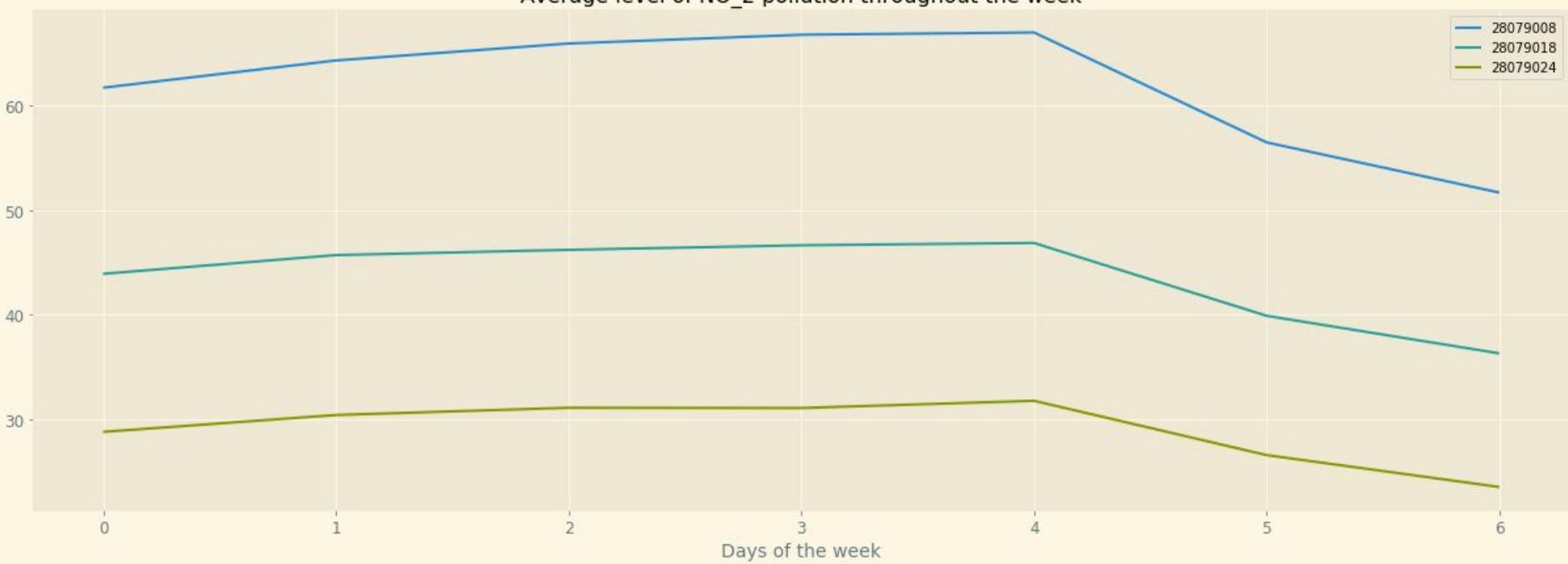


Average level of all_particles pollution throughout the year

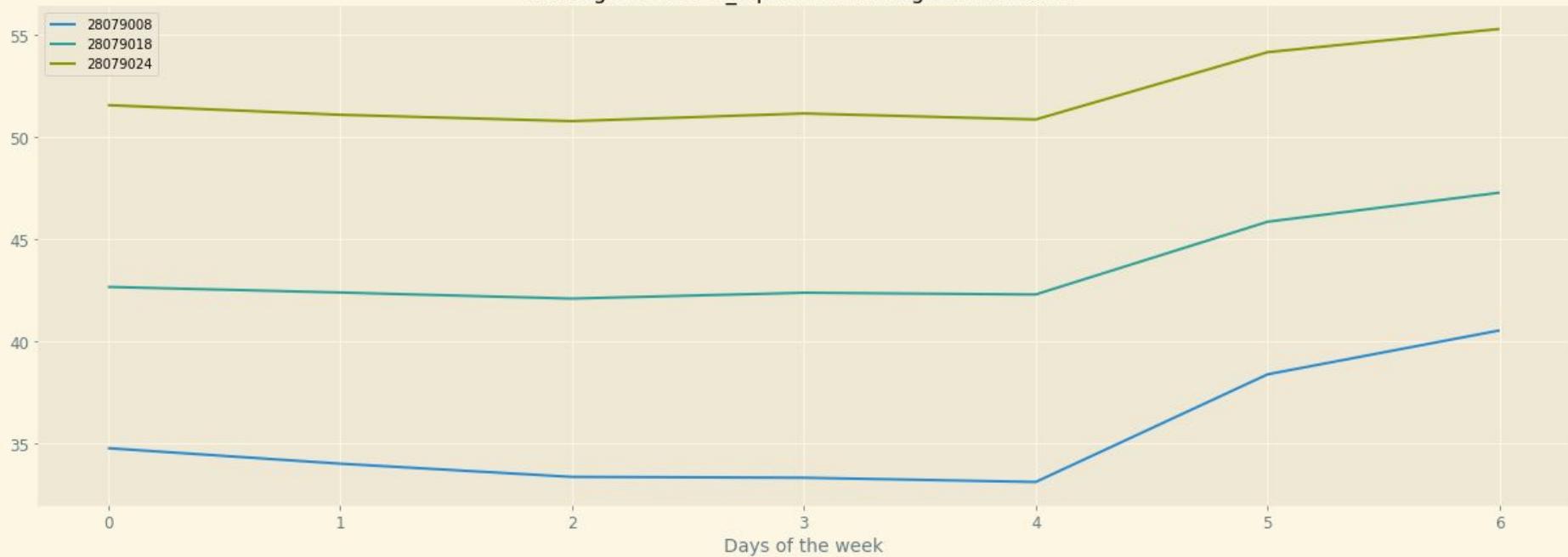


Average levels of pollution
throughout the week

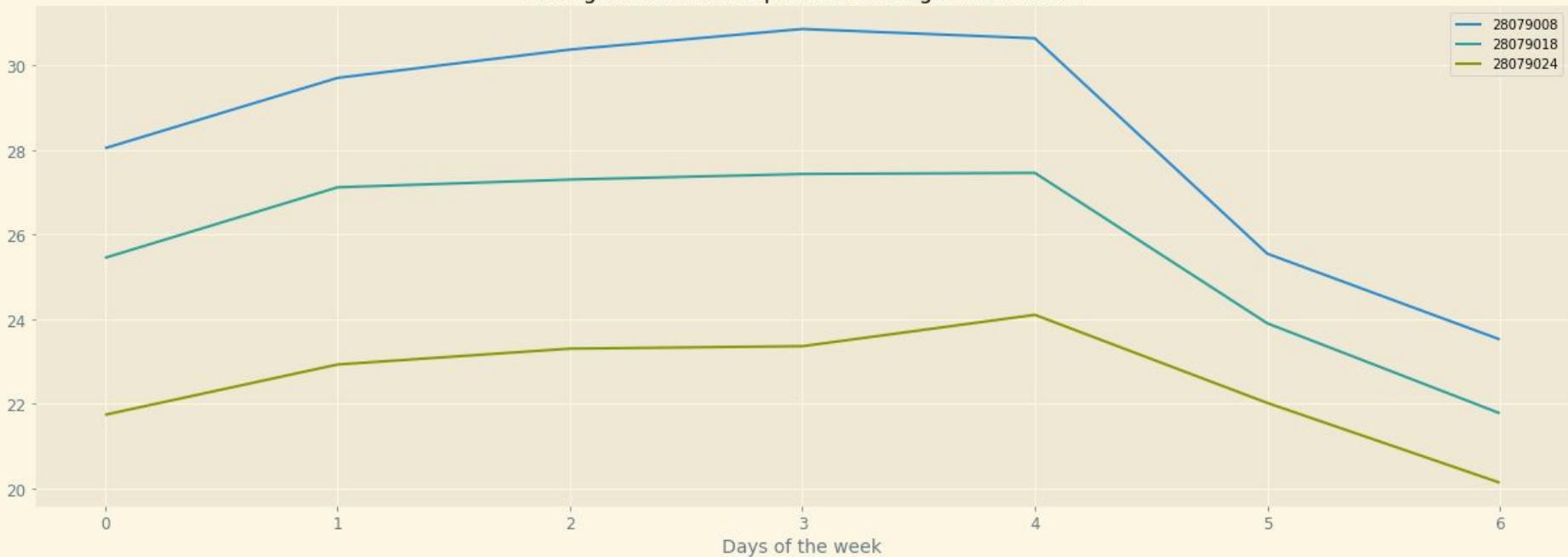
Average level of NO₂ pollution throughout the week



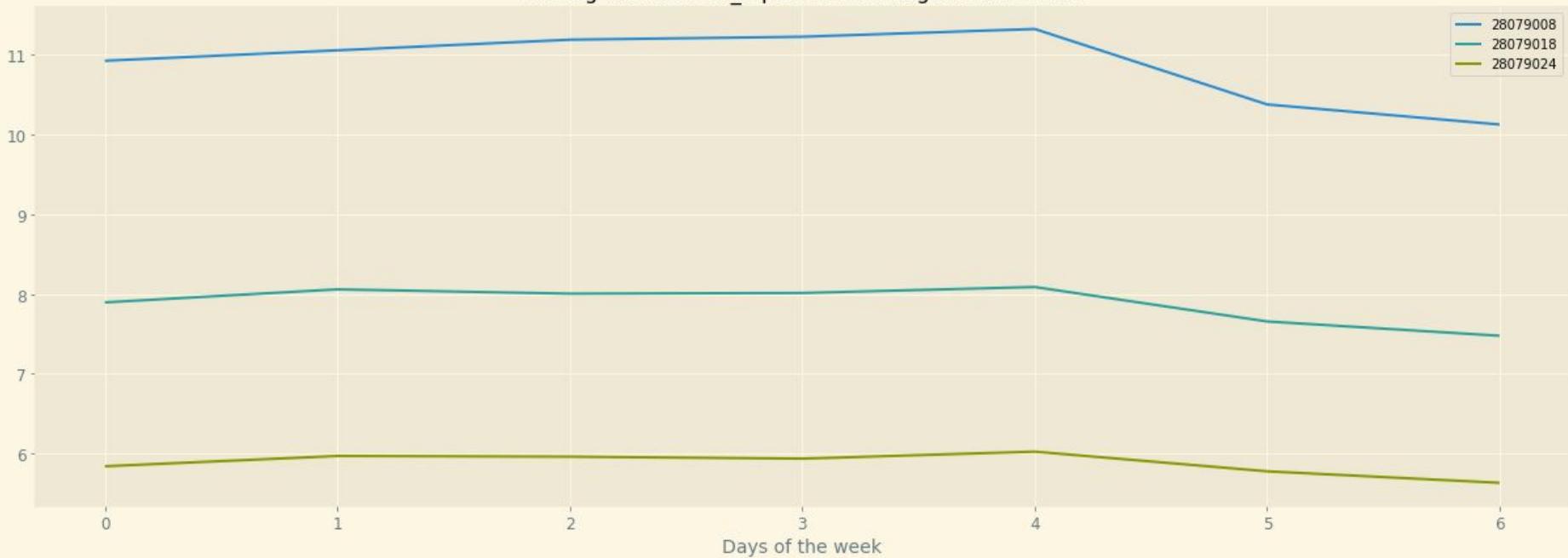
Average level of O₃ pollution throughout the week



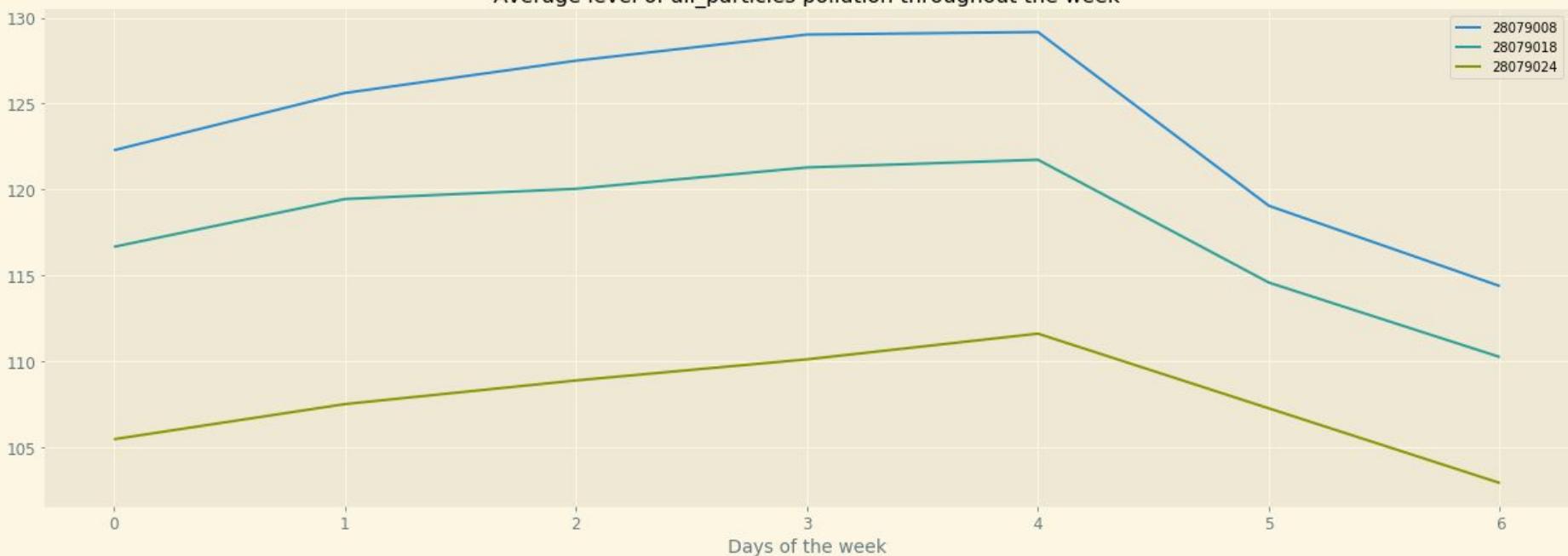
Average level of PM10 pollution throughout the week



Average level of SO₂ pollution throughout the week

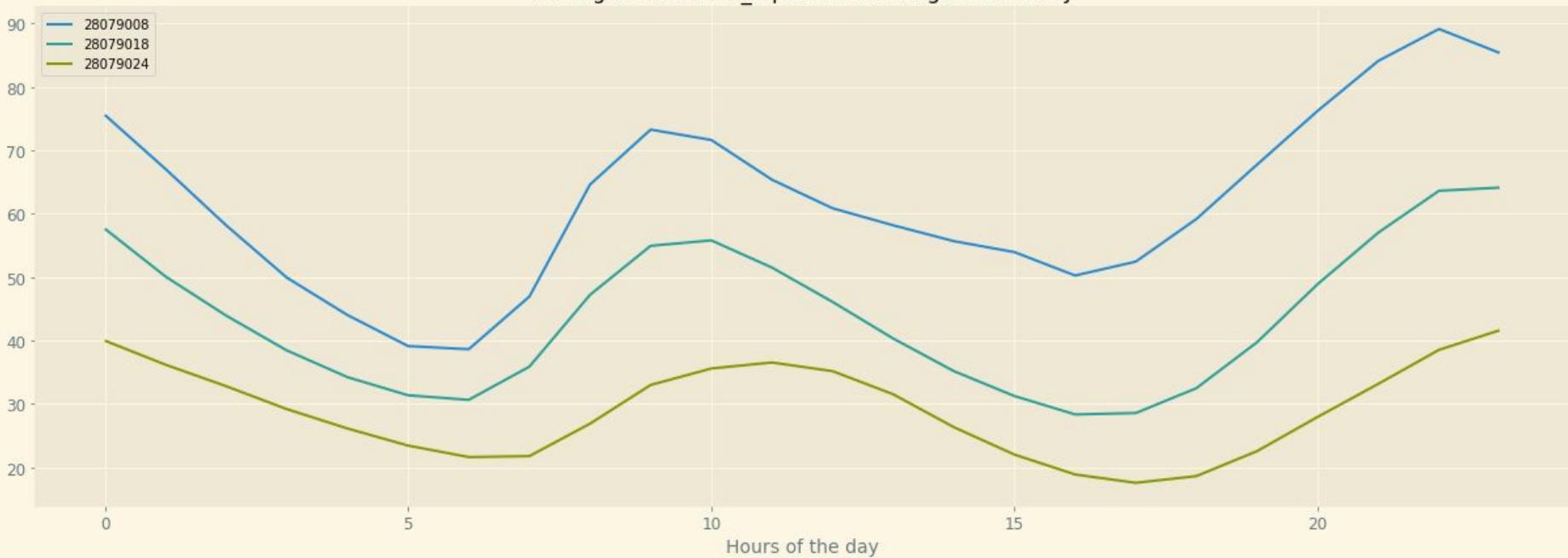


Average level of all_particles pollution throughout the week

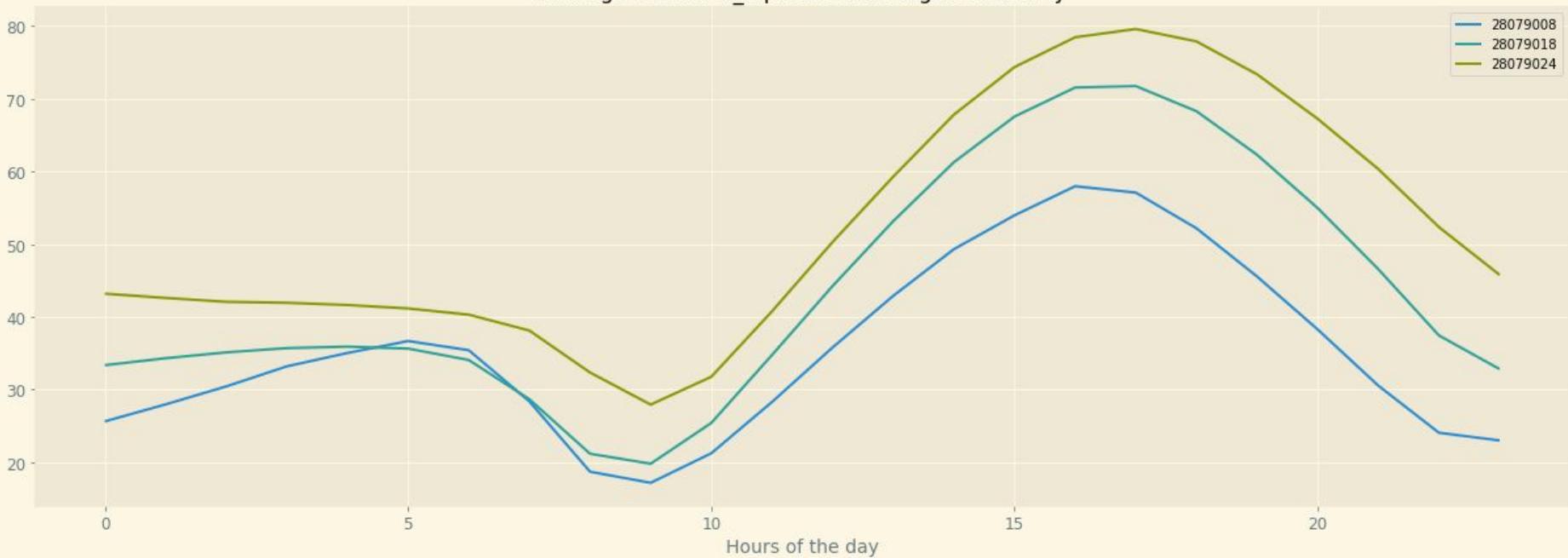


Average levels of pollution
throughout the day

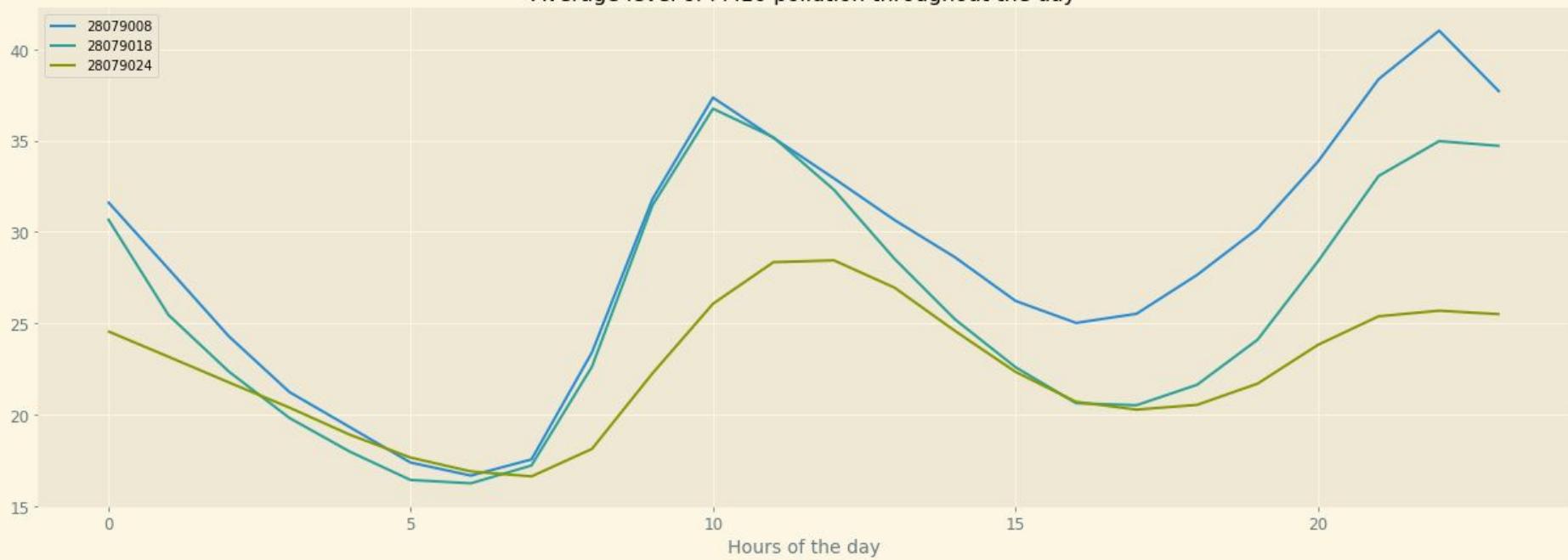
Average level of NO₂ pollution throughout the day



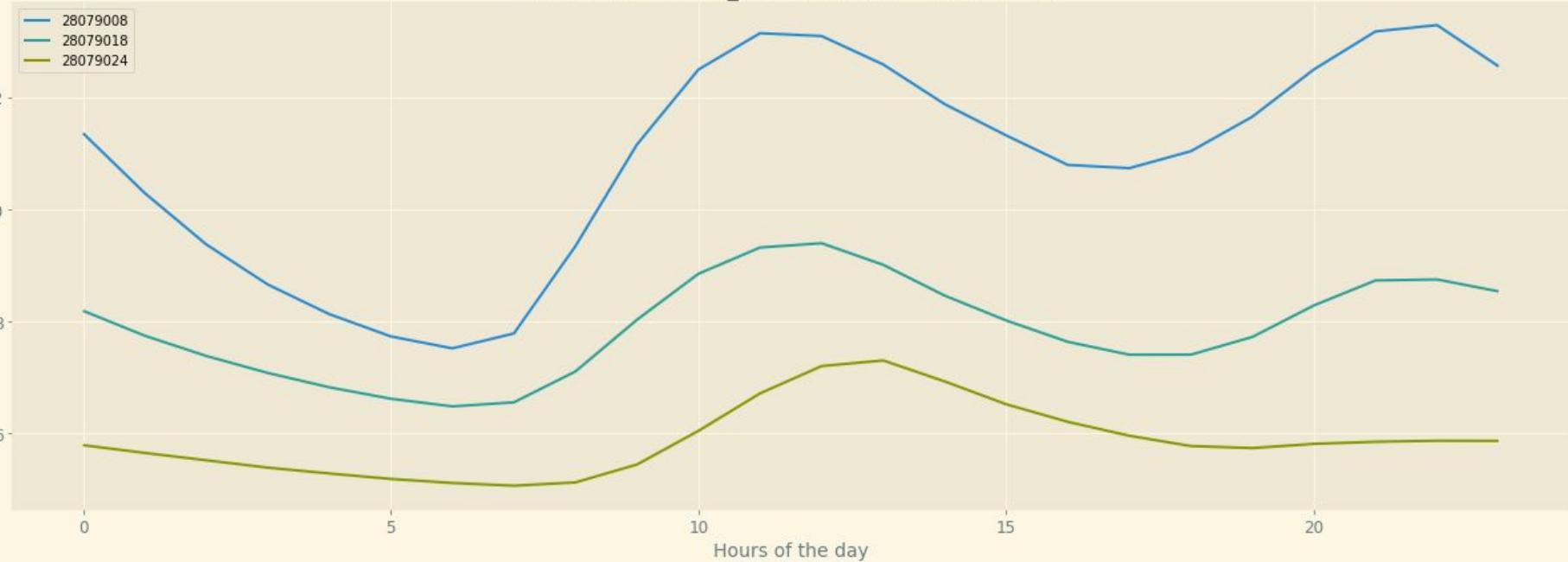
Average level of O₃ pollution throughout the day



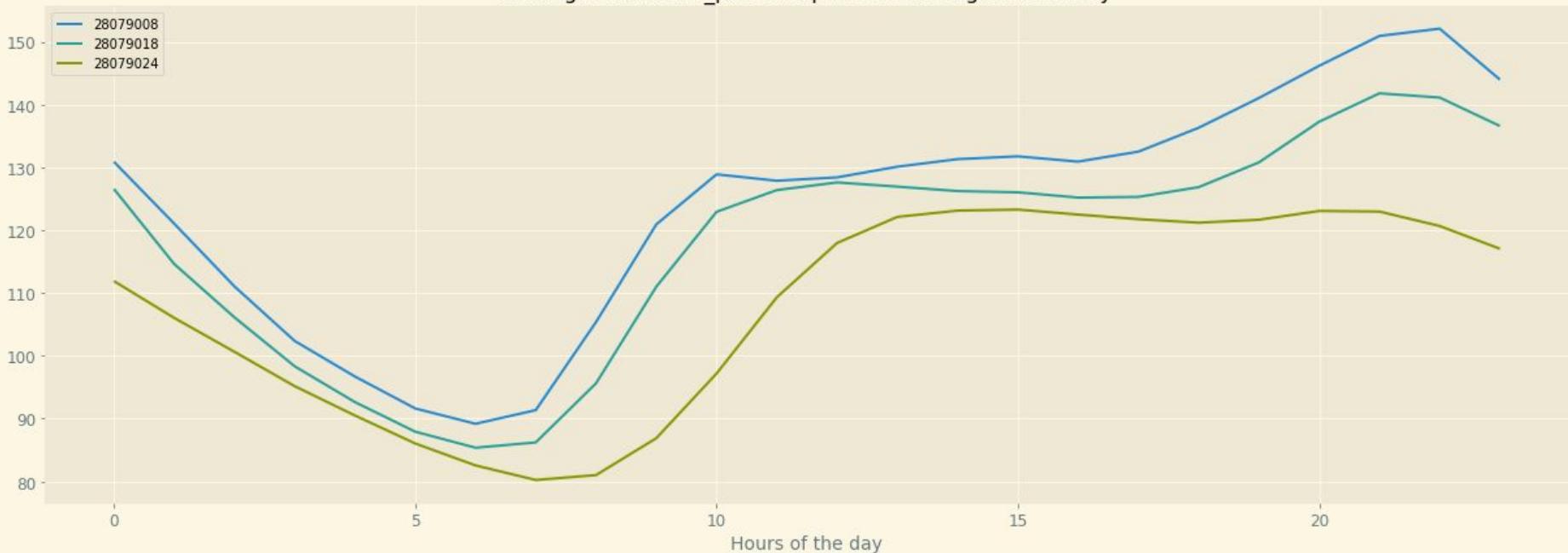
Average level of PM10 pollution throughout the day



Average level of SO₂ pollution throughout the day

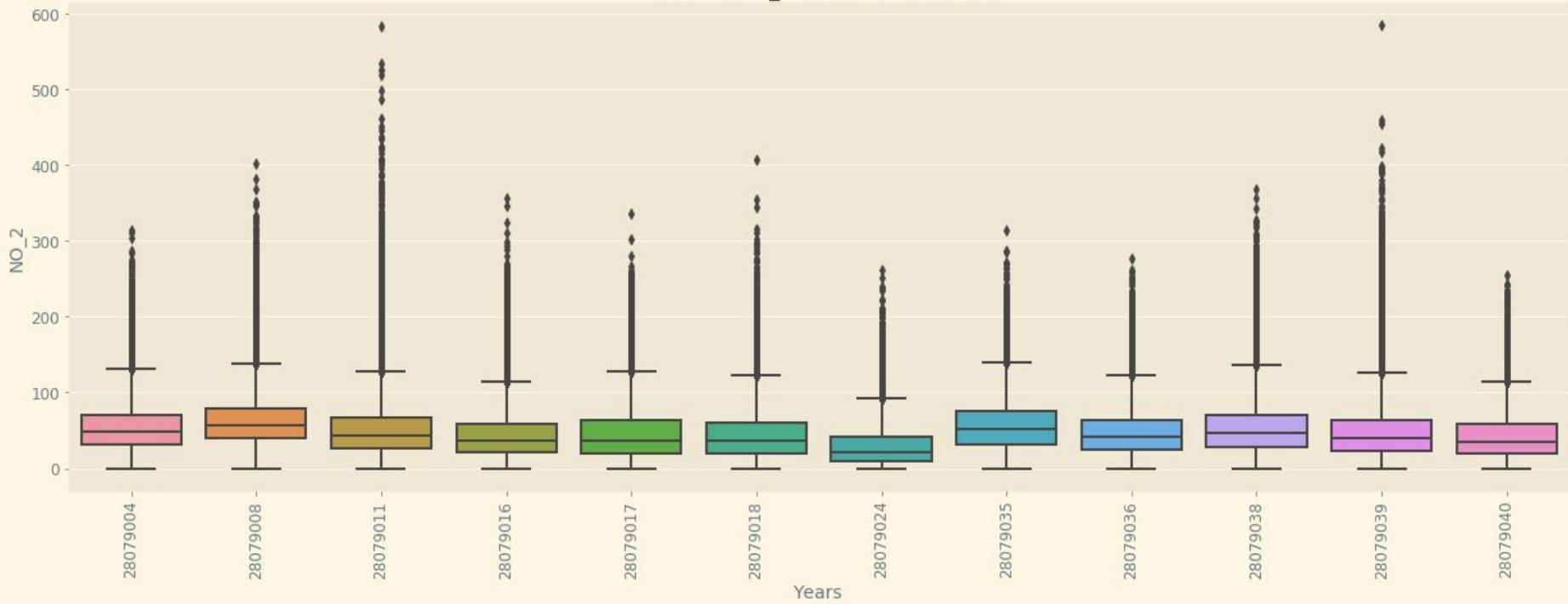


Average level of all_particles pollution throughout the day

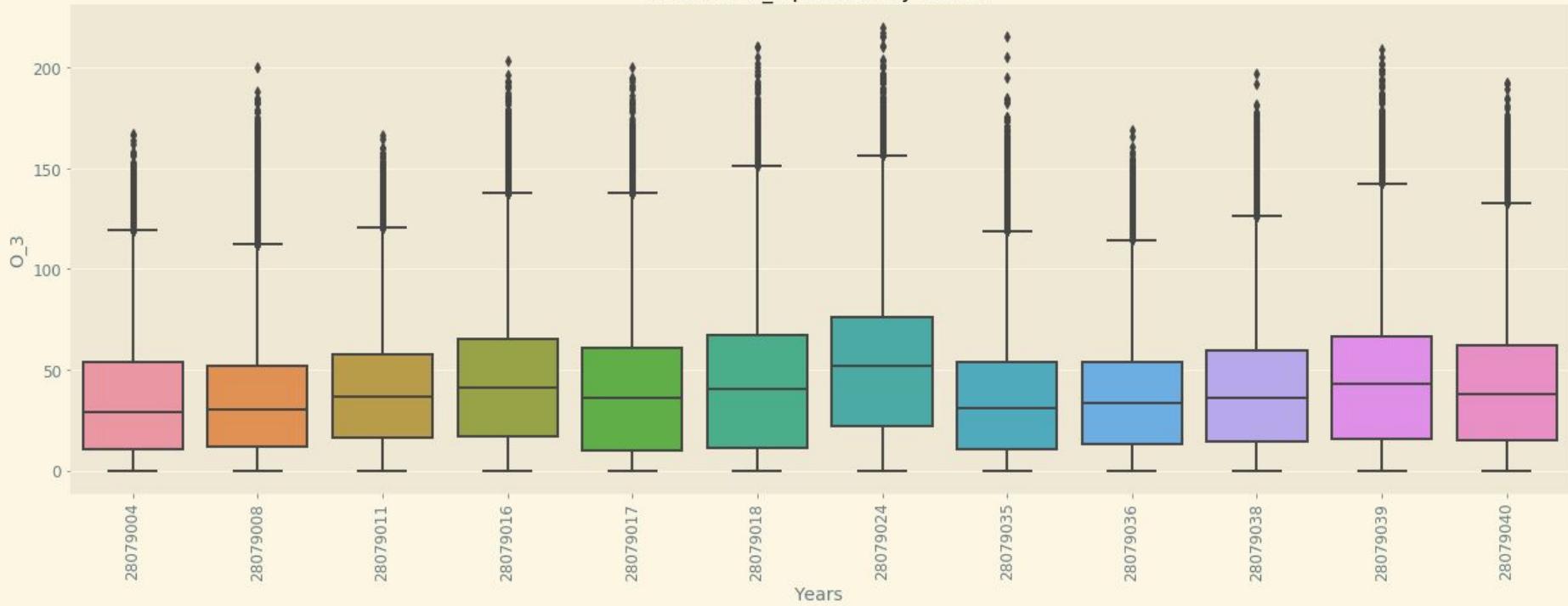


Levels of pollution across
different stations

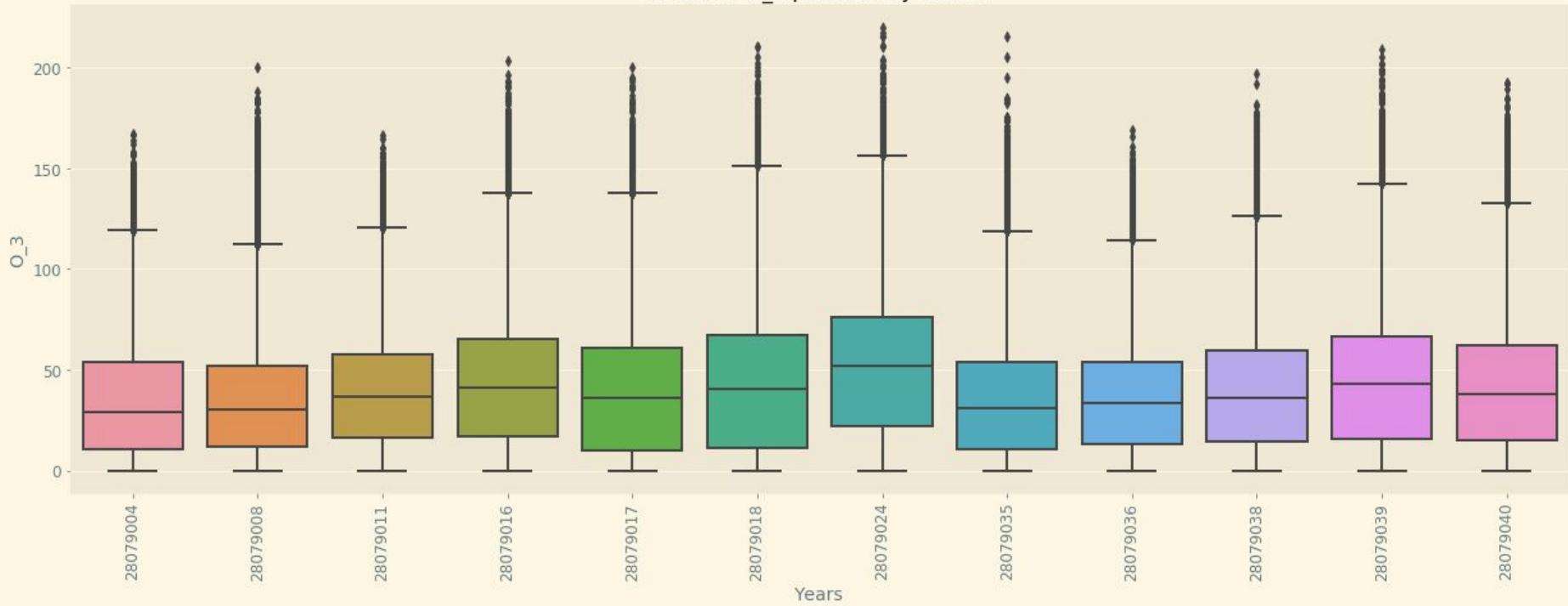
Levels of NO₂ pollution by station

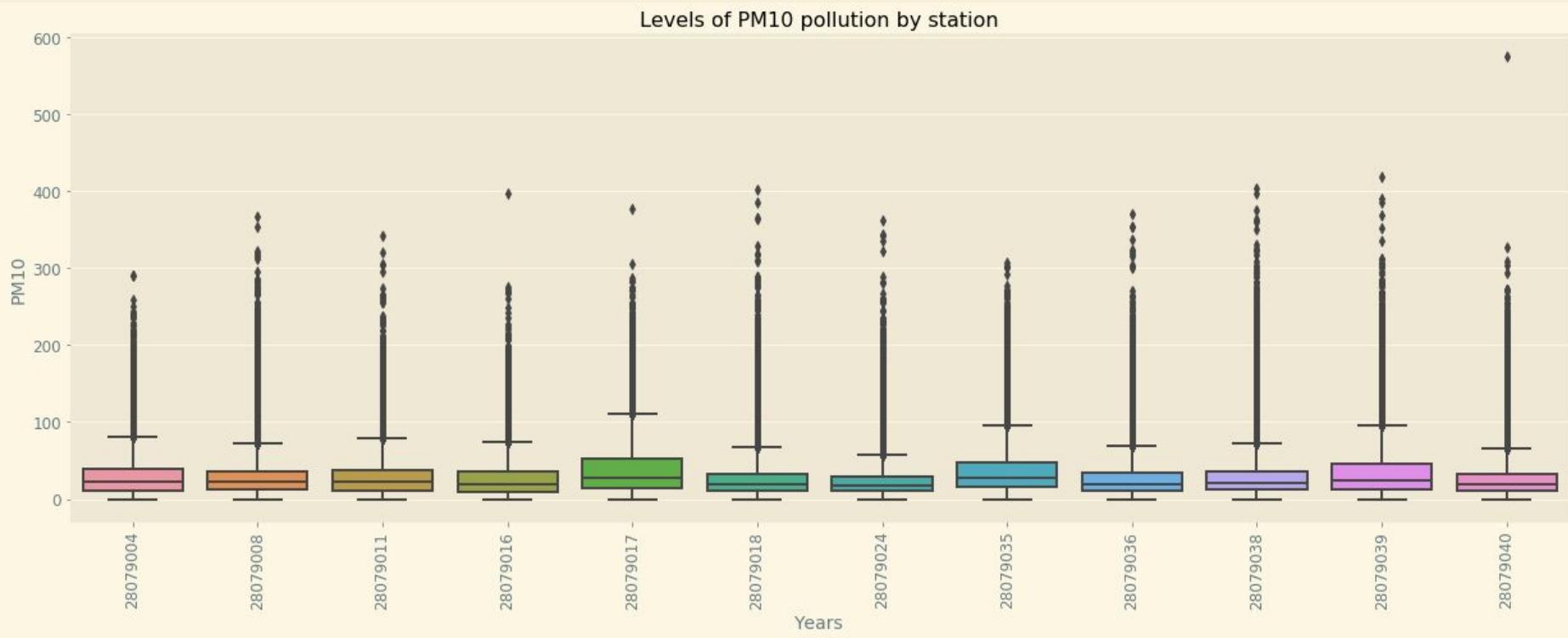


Levels of O₃ pollution by station

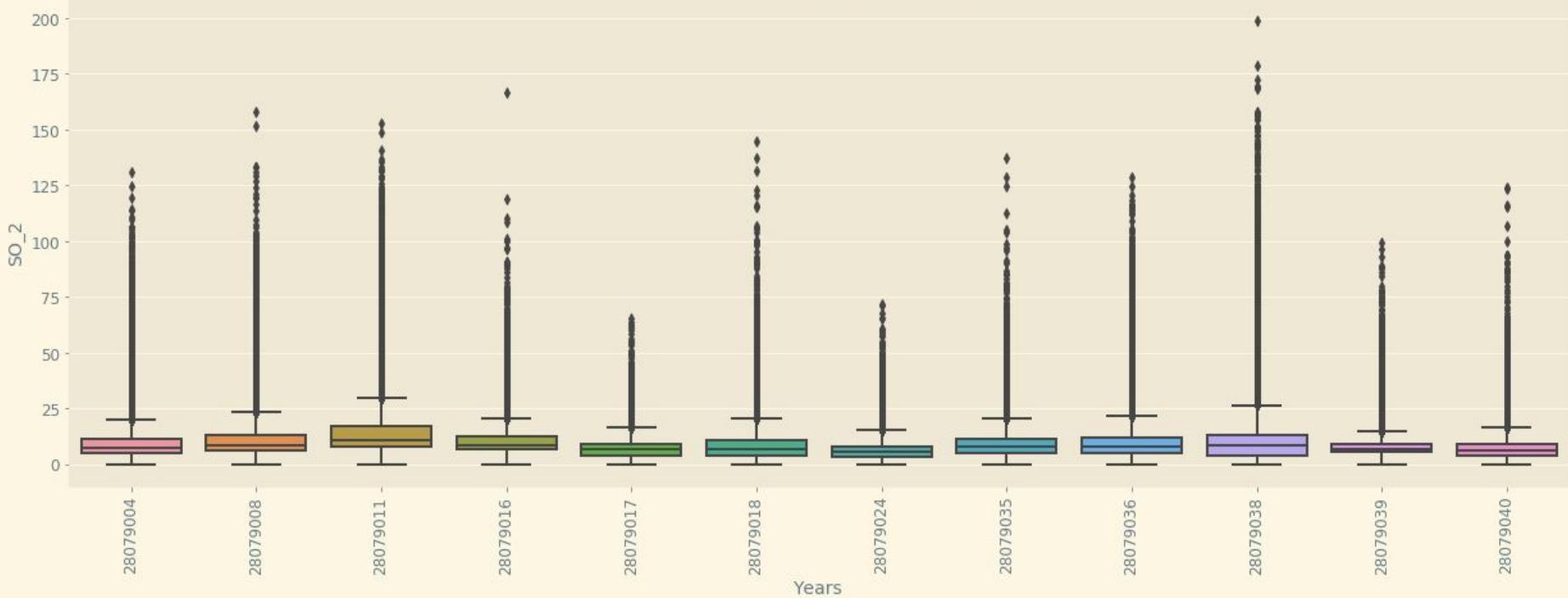


Levels of O₃ pollution by station

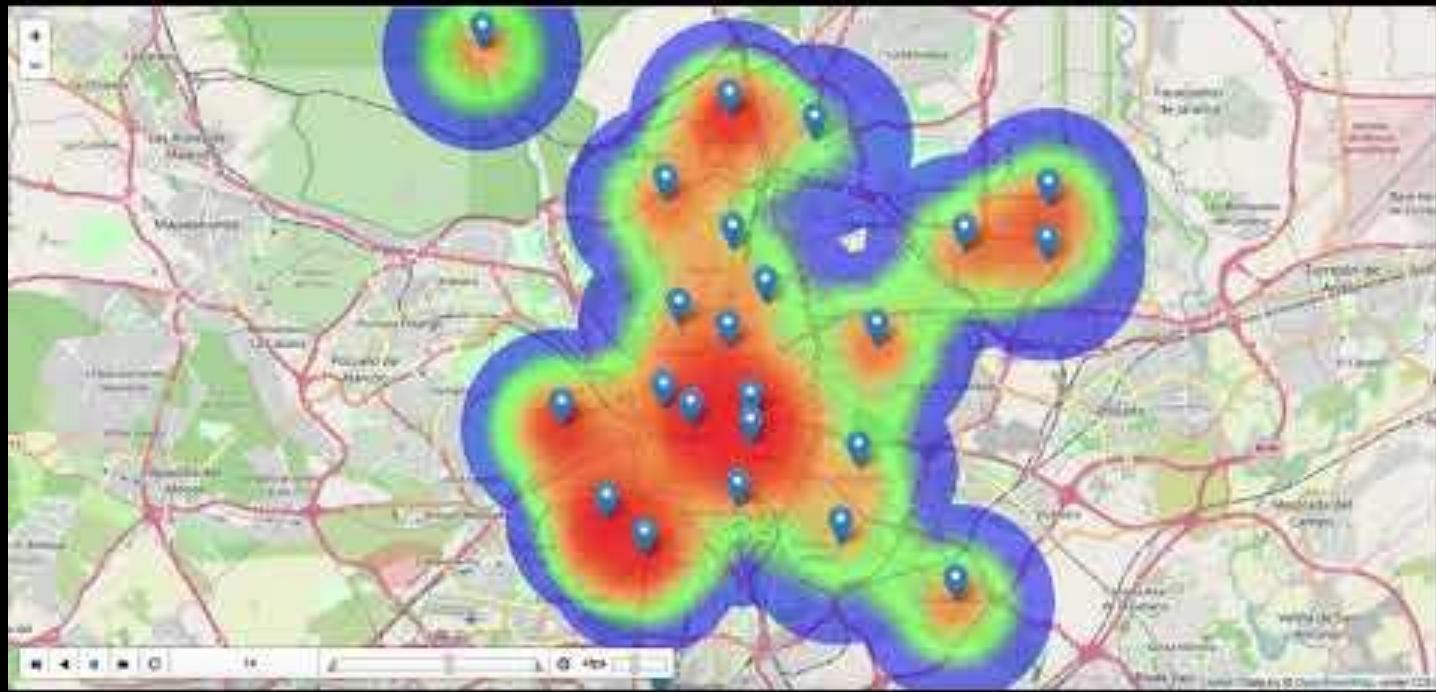




Levels of SO₂ pollution by station

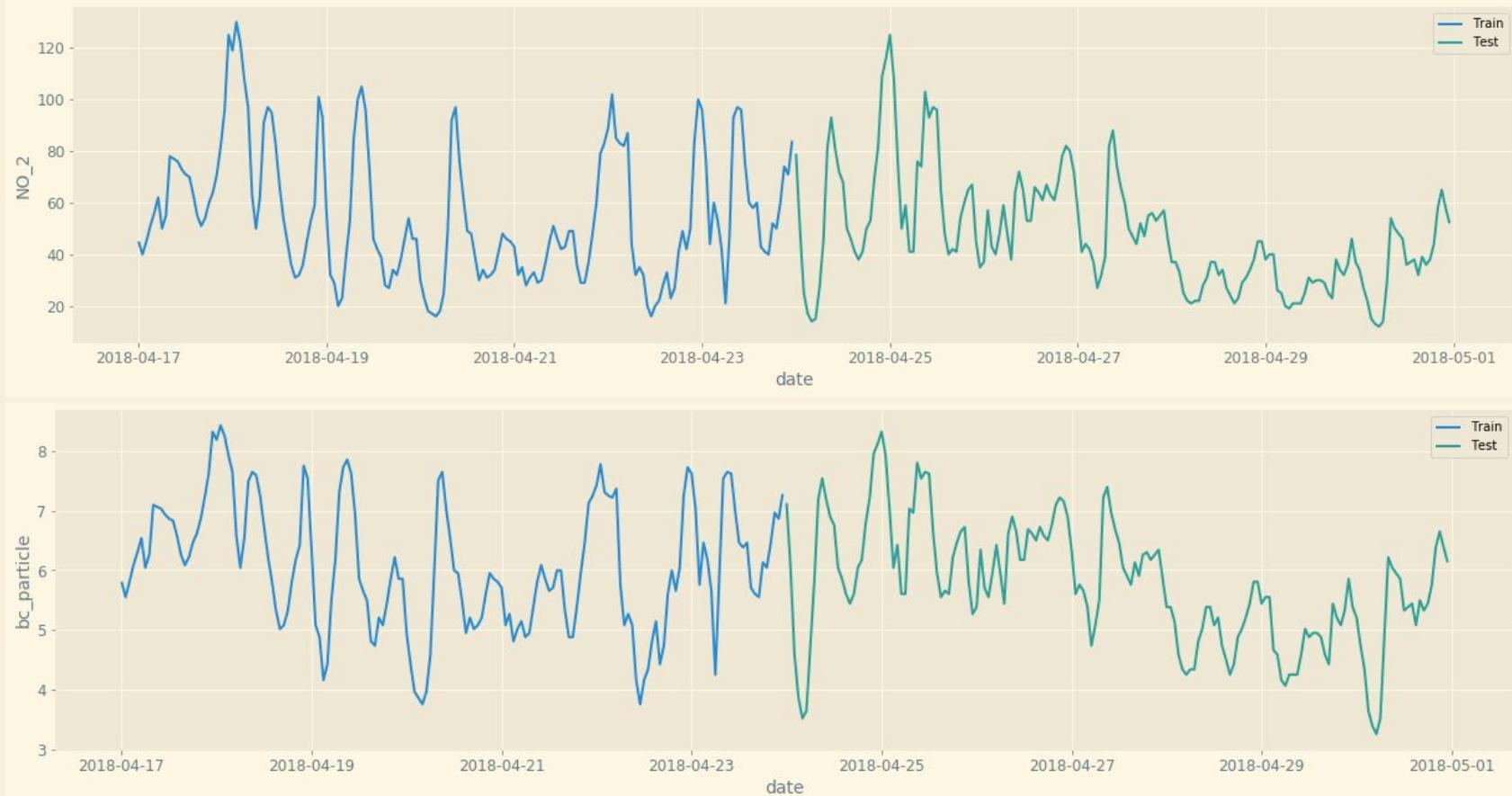




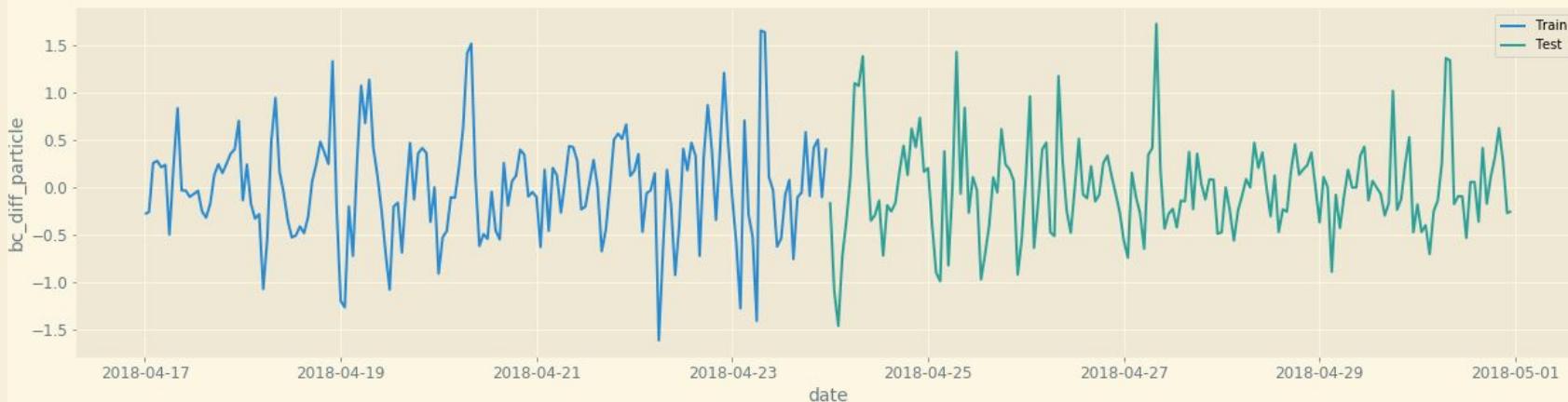
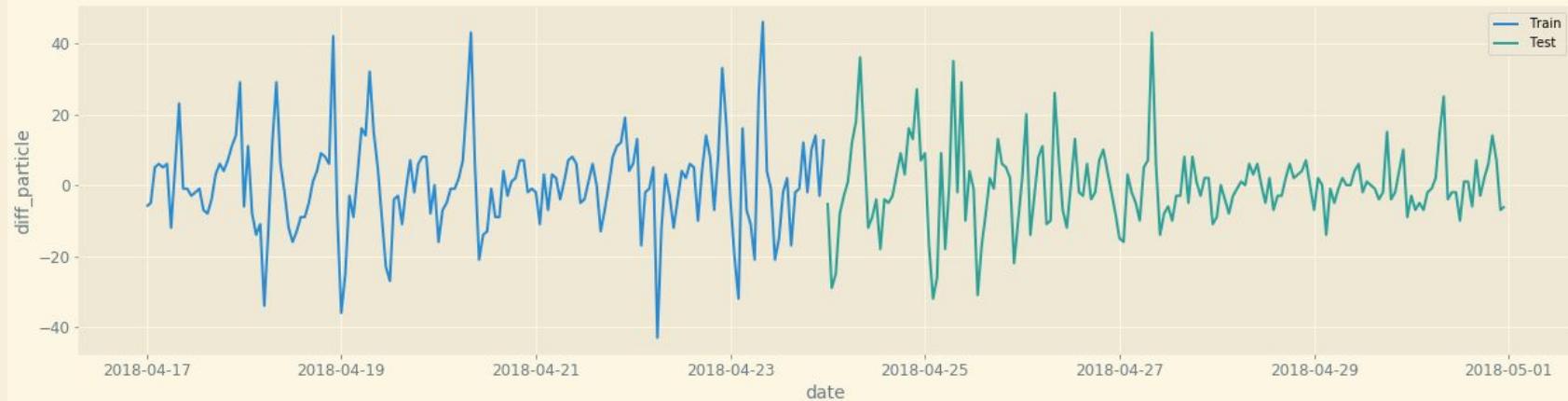


Target transformations

Box Cox Transformation



Difference & Box Cox Transformation

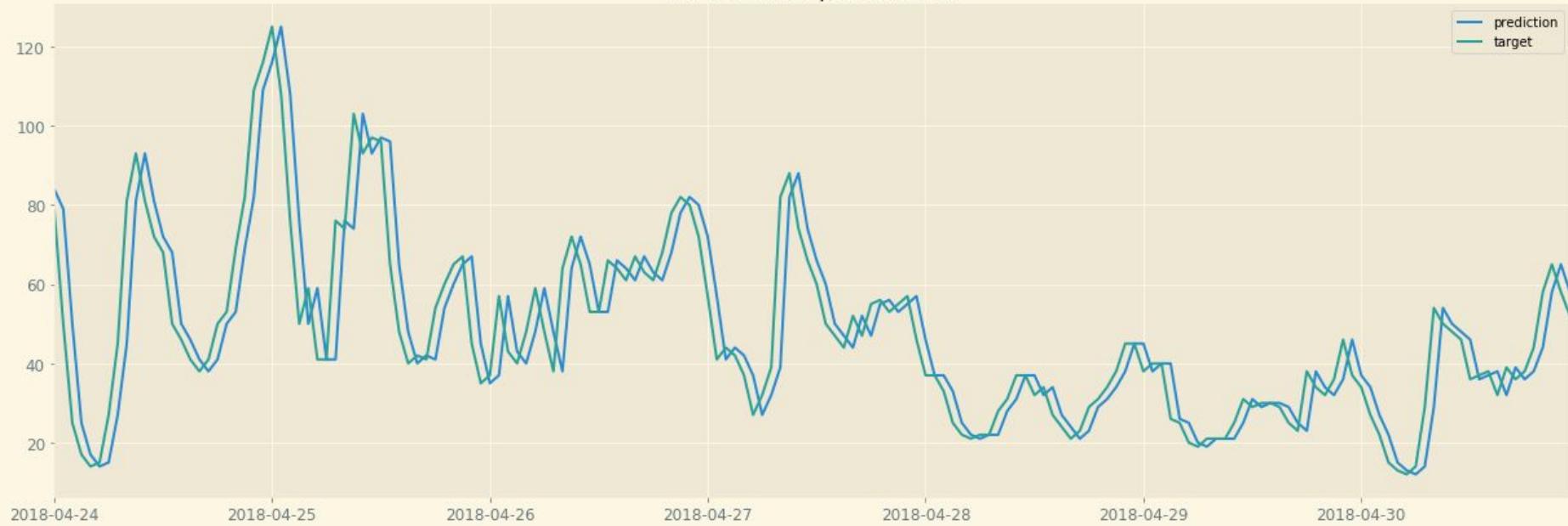


Basic Models

| models | r2_score | mean_abs_error | median_abs_error | mse | rmse |
|-----------------------------|-----------|----------------|------------------|------------|-----------|
| previous hour | 0.746614 | 7.773810 | 5.500000 | 123.452381 | 11.110913 |
| same hour, previous day | -0.073556 | 16.880952 | 13.000000 | 523.047619 | 22.870234 |
| same hour, previous week | -0.285305 | 19.583333 | 14.000000 | 626.214286 | 25.024274 |
| rolling average | 0.606051 | 9.854167 | 6.500000 | 191.936012 | 13.854097 |
| exponential rolling average | 0.584022 | 10.318985 | 7.176439 | 202.668927 | 14.236184 |

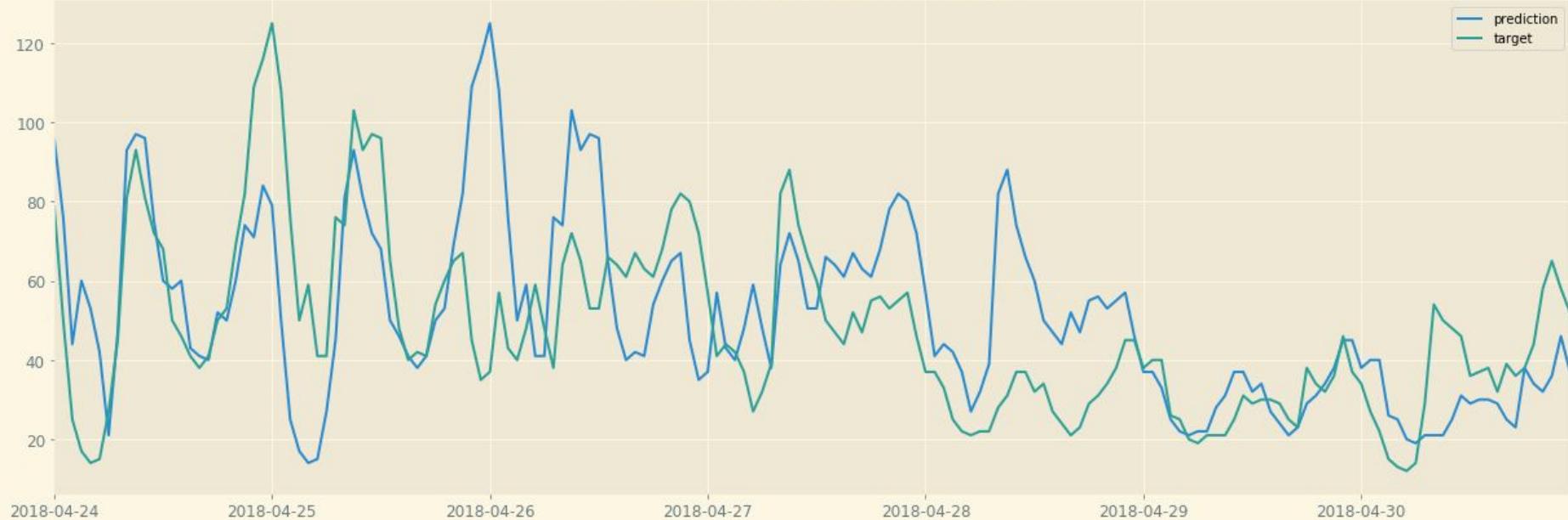
One hour shift

Predictions for previous hour



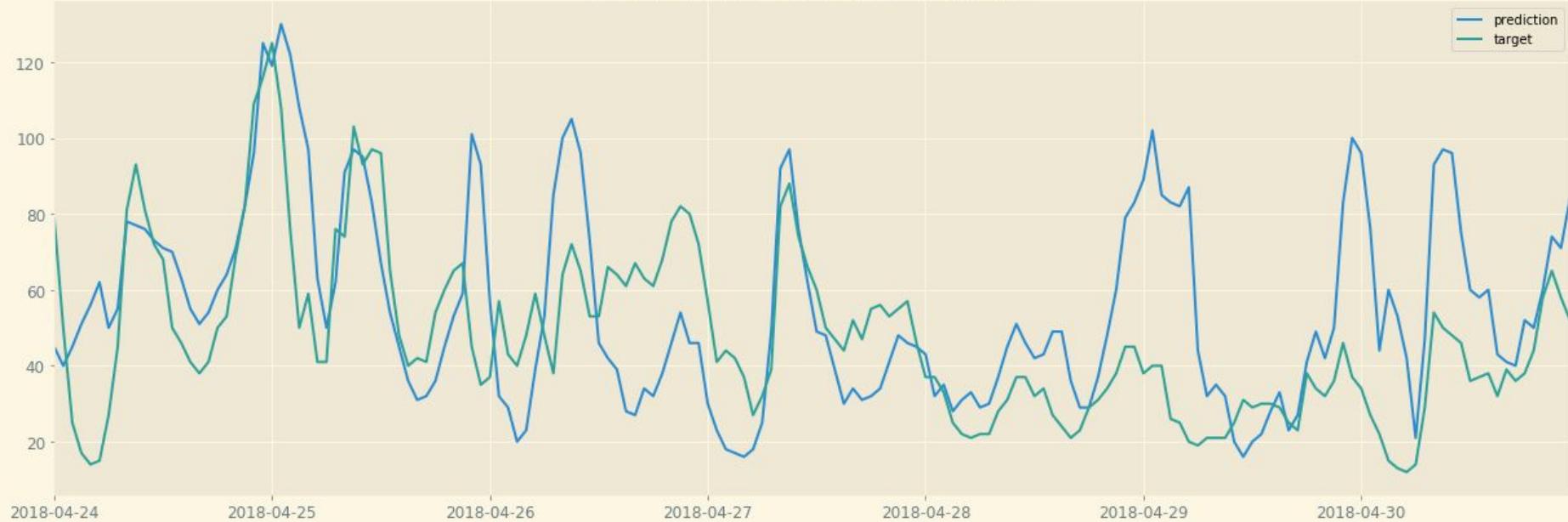
24 hour shift

Predictions for same hour, previous day



One week shift

Predictions for same hour, previous week



Rolling Average, window=2

Predictions for rolling average



Exponential Rolling Average, halflife=1

Predictions for exponential rolling average



Arima

| models | r2_score | mean_abs_error | median_abs_error | mse | rmse |
|--|-----------|----------------|------------------|-------------|-----------|
| ARIMA, order=(5, 0, 0), target: bc_particle | -0.039023 | 17.699984 | 15.000000 | 506.222909 | 22.499398 |
| ARIMA, order=(5, 1, 0), target: bc_particle | -1.761423 | 32.776925 | 34.945345 | 1345.394317 | 36.679617 |
| ARIMA, order=(5, 0, 0), target: bc_diff_particle | -2.836248 | 39.318706 | 41.298895 | 1869.060587 | 43.232633 |
| ARIMA, order=(0, 1, 5), target: bc_particle | -0.195759 | 20.287273 | 19.104218 | 582.586512 | 24.136829 |
| ARIMA, order=(0, 0, 5), target: bc_diff_particle | -2.769662 | 38.932398 | 40.999978 | 1836.619312 | 42.855797 |
| ARIMA, order=(5, 0, 5), target: bc_particle | 0.000398 | 17.431171 | 14.866814 | 487.016664 | 22.068454 |
| ARIMA, order=(5, 0, 5), target: bc_diff_particle | -2.791674 | 39.062975 | 41.013558 | 1847.343691 | 42.980736 |

Target: NO2 log | Station: 28079008 - Parque del Retiro/ Escuelas Aguirre | Model: AR

Order = (5, 0, 0)

RMSE: 22.50



Target: NO2 log | Station: 28079008 - Parque del Retiro/ Escuelas Aguirre | Model: AR

Order = (5, 1, 0)

RMSE: 36.68



Target: NO2 log diff | Station: 28079008 - Parque del Retiro/ Escuelas Aguirre | Model: AR

Order = (5, 0, 0)

RMSE: 43.23



Target: NO2 log | Station: 28079008 - Parque del Retiro/ Escuelas Aguirre | Model: MA

Order = (0, 1, 5)

RMSE: 24.14



Target: NO2 log diff | Station: 28079008 - Parque del Retiro/ Escuelas Aguirre | Model: MA

Order = (0, 0, 5)

RMSE: 42.86



Target: NO2 log | Station: 28079008 - Parque del Retiro/ Escuelas Aguirre | Model: ARMA

Order = (5, 0, 5)

RMSE: 22.07



Target: NO2 log diff | Station: 28079008 - Parque del Retiro/ Escuelas Aguirre | Model: ARMA

Order = (5, 0, 5)

RMSE: 42.98



Target: NO2 log | Station: 28079008 - Parque del Retiro/ Escuelas Aguirre | Model: ARIMA

Order = (5, 1, 5)

RMSE: 27.72



Model Tuning

| models | r2_score | mean_abs_error | median_abs_error | mse | rmse |
|---------------------------------|-----------|----------------|------------------|-------------|-----------|
| Auto ARIMA: bc_particle | -0.093783 | 19.169096 | 17.359020 | 532.902866 | 23.084689 |
| Auto ARIMA: bc_diff_particle | -2.769928 | 38.934043 | 40.996413 | 1836.748712 | 42.857306 |
| Auto SARIMA: bc_particle | -1.923727 | 31.099324 | 31.353438 | 1424.470578 | 37.742159 |

Target: NO2 log | Station: 28079008 - Parque del Retiro/ Escuelas Aguirre | Model: Auto ARIMA

Order = (5, 1, 5)

RMSE: 23.08



Target: NO2 log diff | Station: 28079008 - Parque del Retiro/ Escuelas Aguirre | Model: Auto ARIMA

Order = (5, 0, 5)

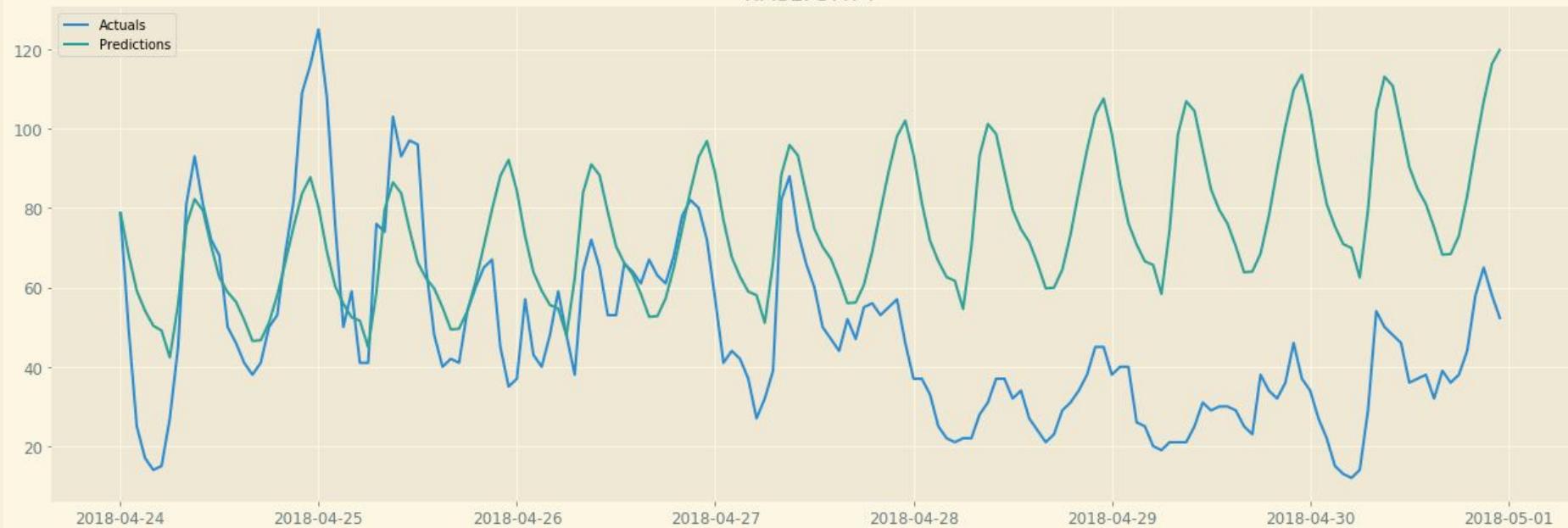
RMSE: 42.86



Target: NO2 log | Station: 28079008 - Parque del Retiro/ Escuelas Aguirre | Model: Auto SARIMA

Order = (2, 1, 2) | Seasonal order = (2, 0, 2, 12)

RMSE: 37.74



Rolling predictions with ARIMA

NO₂

Visualising the target variable



Testing Stationarity

Results of Dickey-Fuller Test:

Null Hypothesis: Unit Root Present
Test Statistic < Critical Value => Reject Null
P-Value =< Alpha(.05) => Reject Null

| | |
|--------------------|---------|
| Test Statistic | -19.065 |
| p-value | 0.000 |
| #Lags Used | 77.000 |
| Nr of Obs Used | 166386 |
| Critical Value 1% | 3.430 |
| Critical Value 5% | -2.862 |
| Critical Value 10% | -2.567 |

Results of KPSS Test:

Null Hypothesis: Data is Stationary
Test Statistic > Critical Value => Reject Null
P-Value =< Alpha(.05) => Reject Null

| | |
|---------------------|--------|
| Test Statistic | 8.344 |
| p-value | 0.010 |
| Lags Used | 77.000 |
| Critical Value 10% | 0.347 |
| Critical Value 5% | 0.463 |
| Critical Value 2.5% | 0.574 |
| Critical Value 1% | 0.739 |

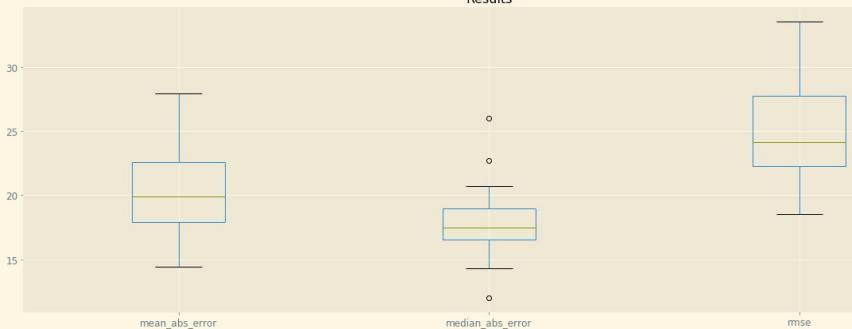


Station nr 28079008 - Parque del Retiro/ Escuelas Aguirre



Station nr 28079018 - San Isidro/Calle Farolillo

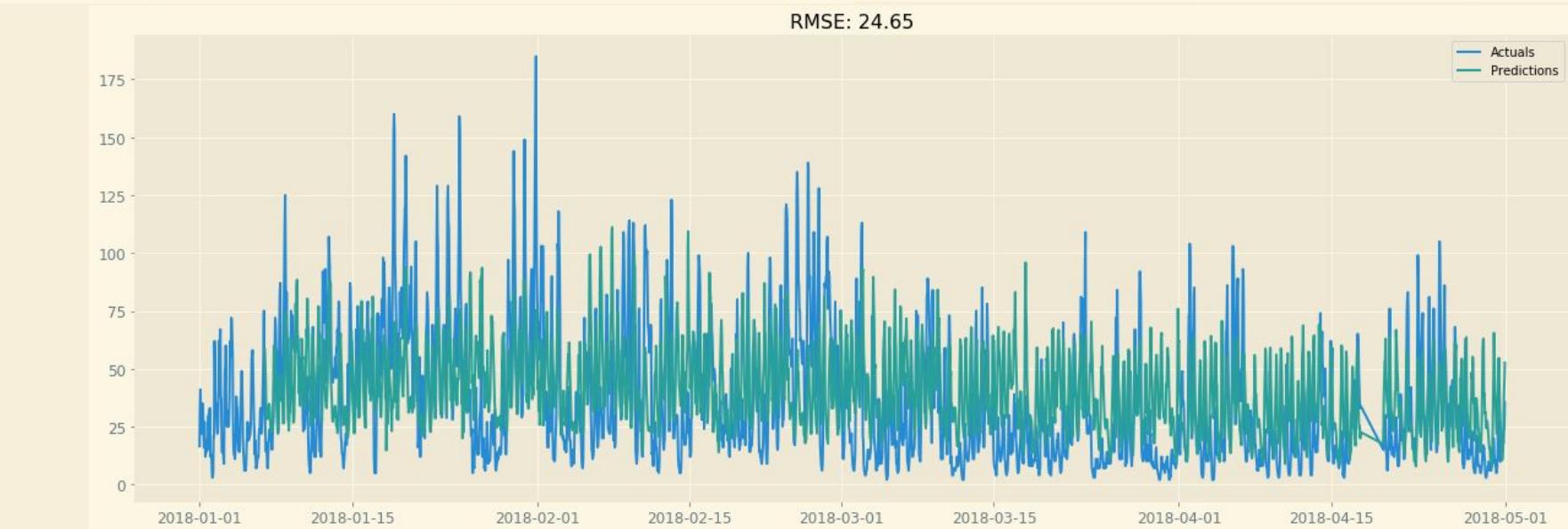
Results



RMSE: 24.65

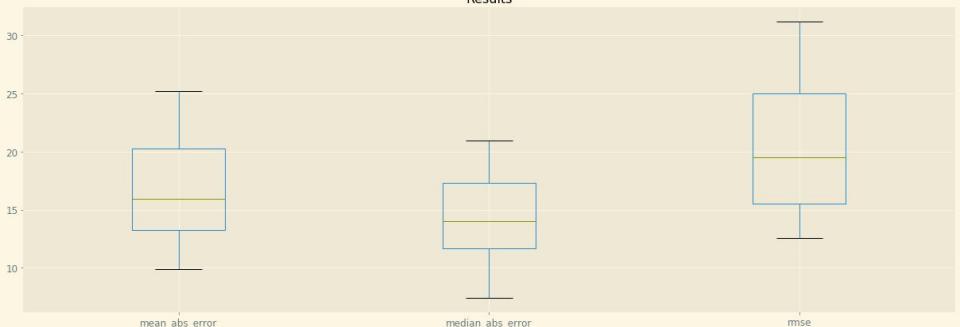


RMSE: 24.65



Station nr 28079024 - Casa de Campo

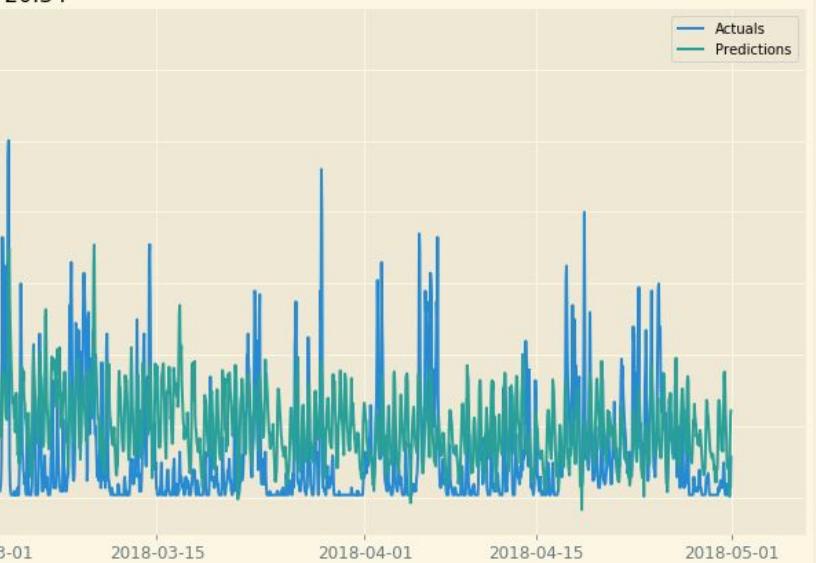
Results



RMSE: 20.34

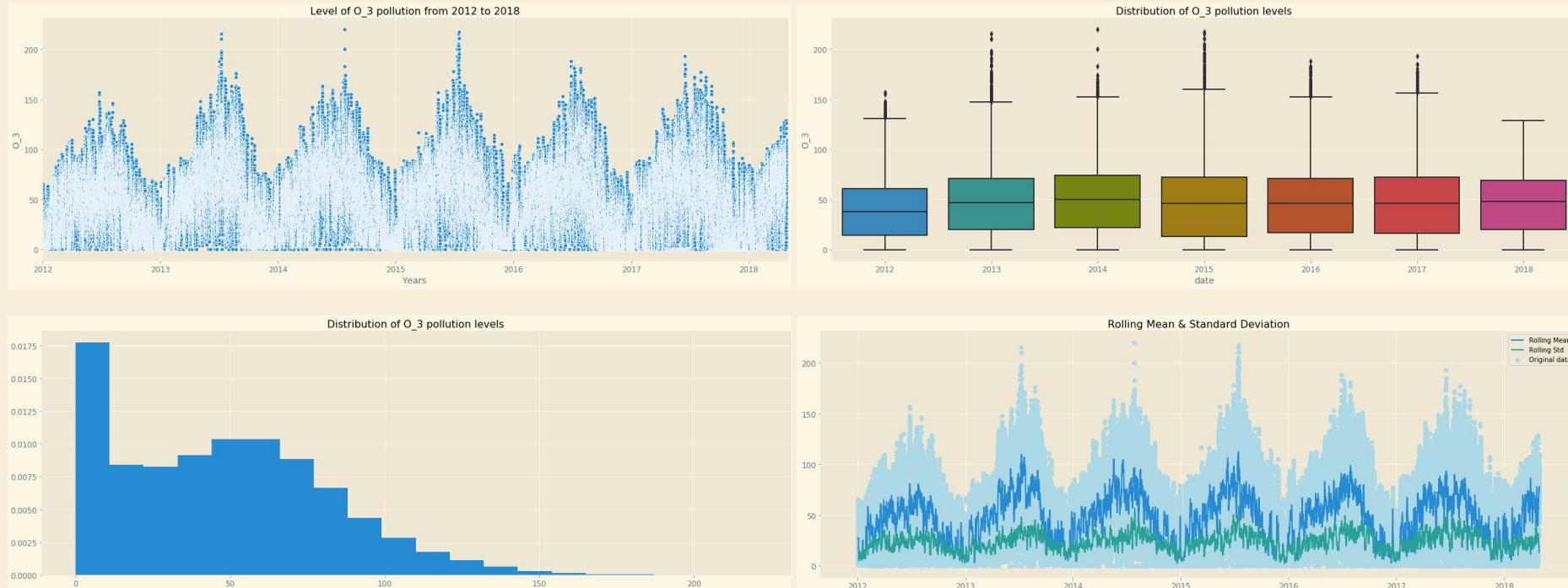


RMSE: 20.34



O3

Visualising the target variable



Testing Stationarity

Results of Dickey-Fuller Test:

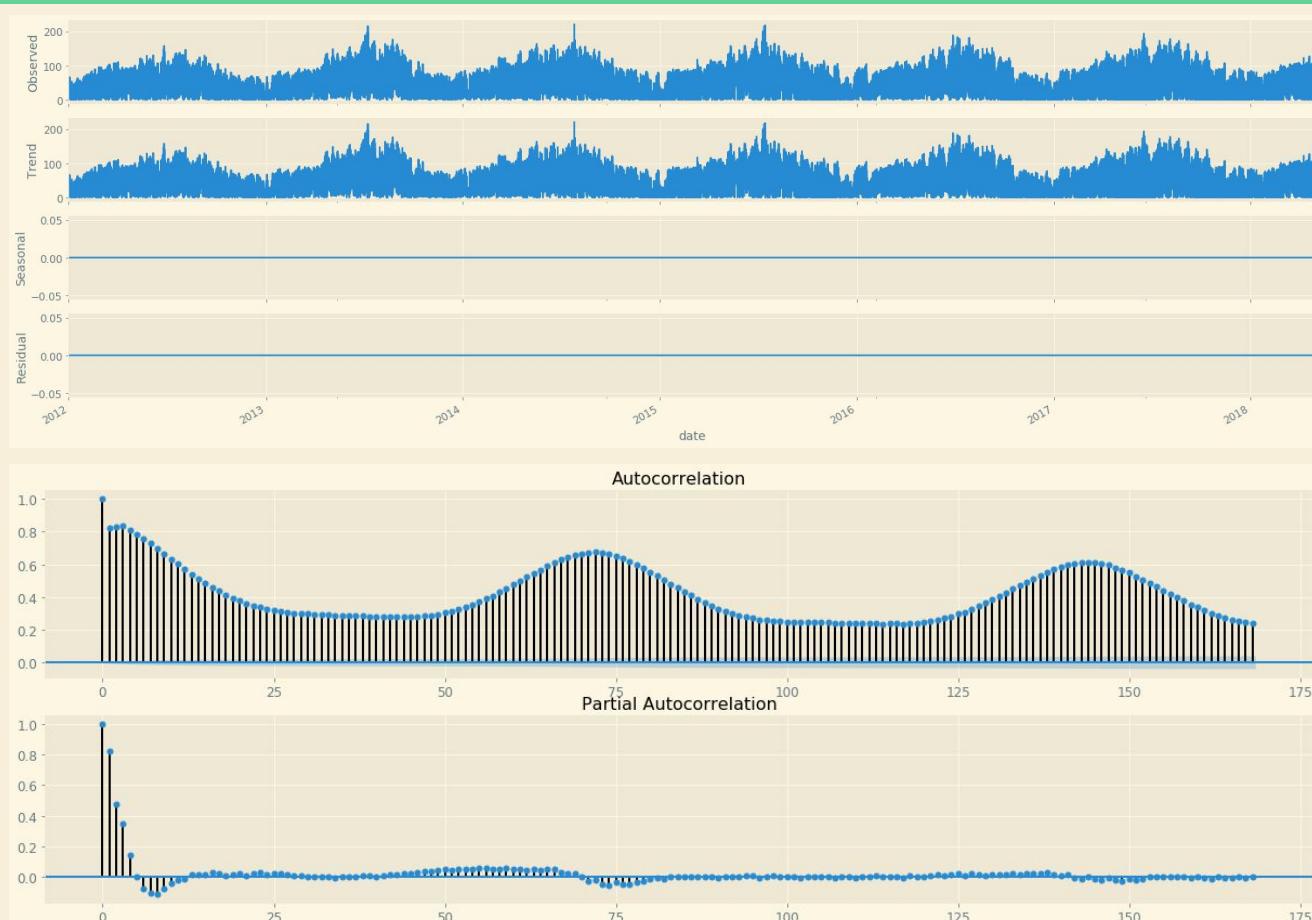
Null Hypothesis: Unit Root Present
Test Statistic < Critical Value => Reject Null
P-Value =< Alpha(.05) => Reject Null

| | |
|--------------------|------------|
| Test Statistic | -1.483e+01 |
| p-value | 1.891e-27 |
| #Lags Used | 7.700e+01 |
| Nr of Obs Used | 1.664e+05 |
| Critical Value 1% | -3.430 |
| Critical Value 5% | -2.862 |
| Critical Value 10% | -2.567 |

Results of KPSS Test:

Null Hypothesis: Data is Stationary
Test Statistic > Critical Value => Reject Null
P-Value =< Alpha(.05) => Reject Null

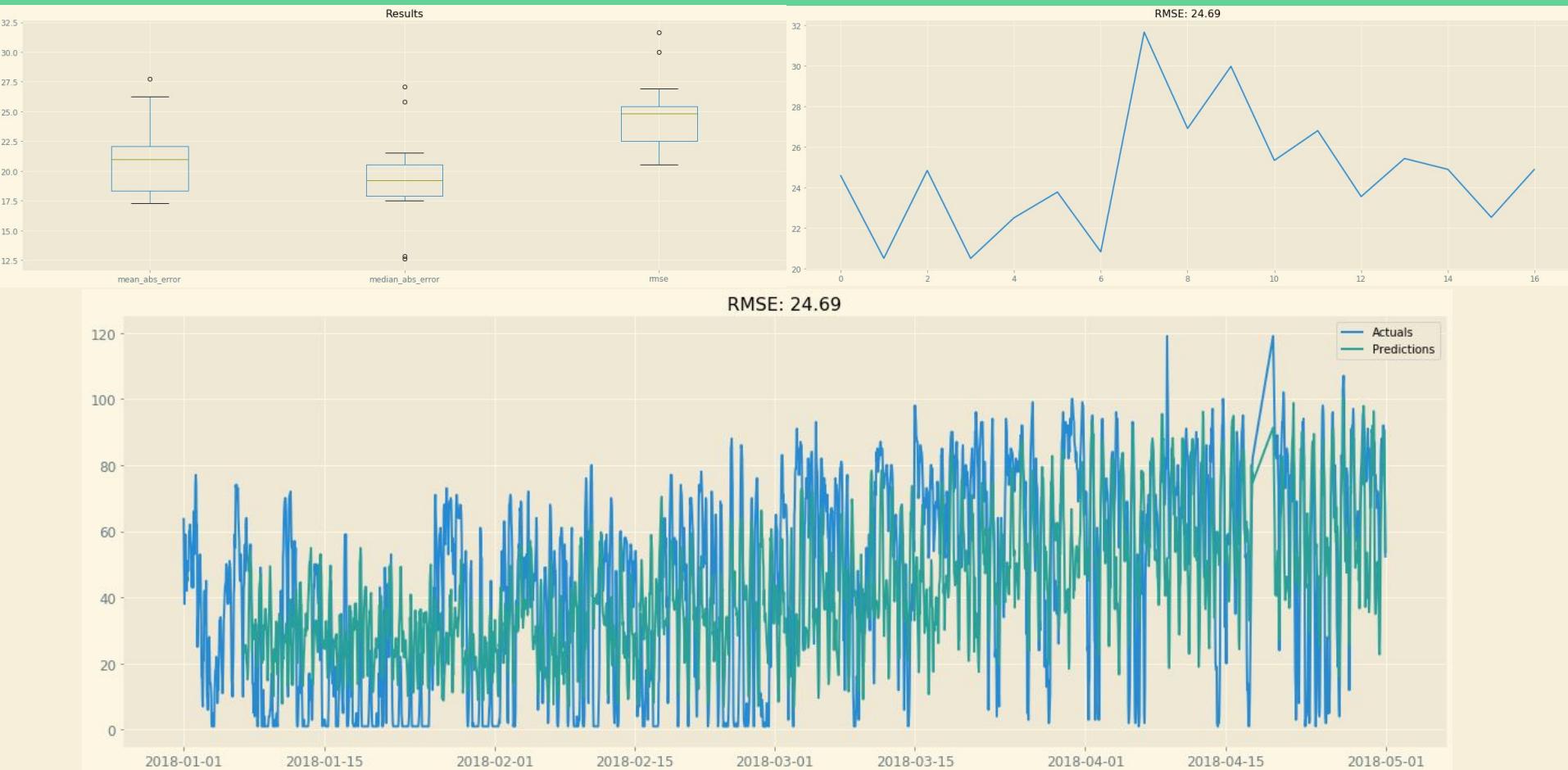
| | |
|---------------------|--------|
| Test Statistic | 1.562 |
| p-value | 0.010 |
| Lags Used | 77.000 |
| Critical Value 10% | 0.347 |
| Critical Value 5% | 0.463 |
| Critical Value 2.5% | 0.574 |
| Critical Value 1% | 0.739 |



Station nr 28079008 - Parque del Retiro/ Escuelas Aguirre



Station nr 28079018 - San Isidro/Calle Farolillo



Station nr 28079018 - San Isidro/Calle Farolillo



PM10

Visualising the target variable



Testing Stationarity

Results of Dickey-Fuller Test:

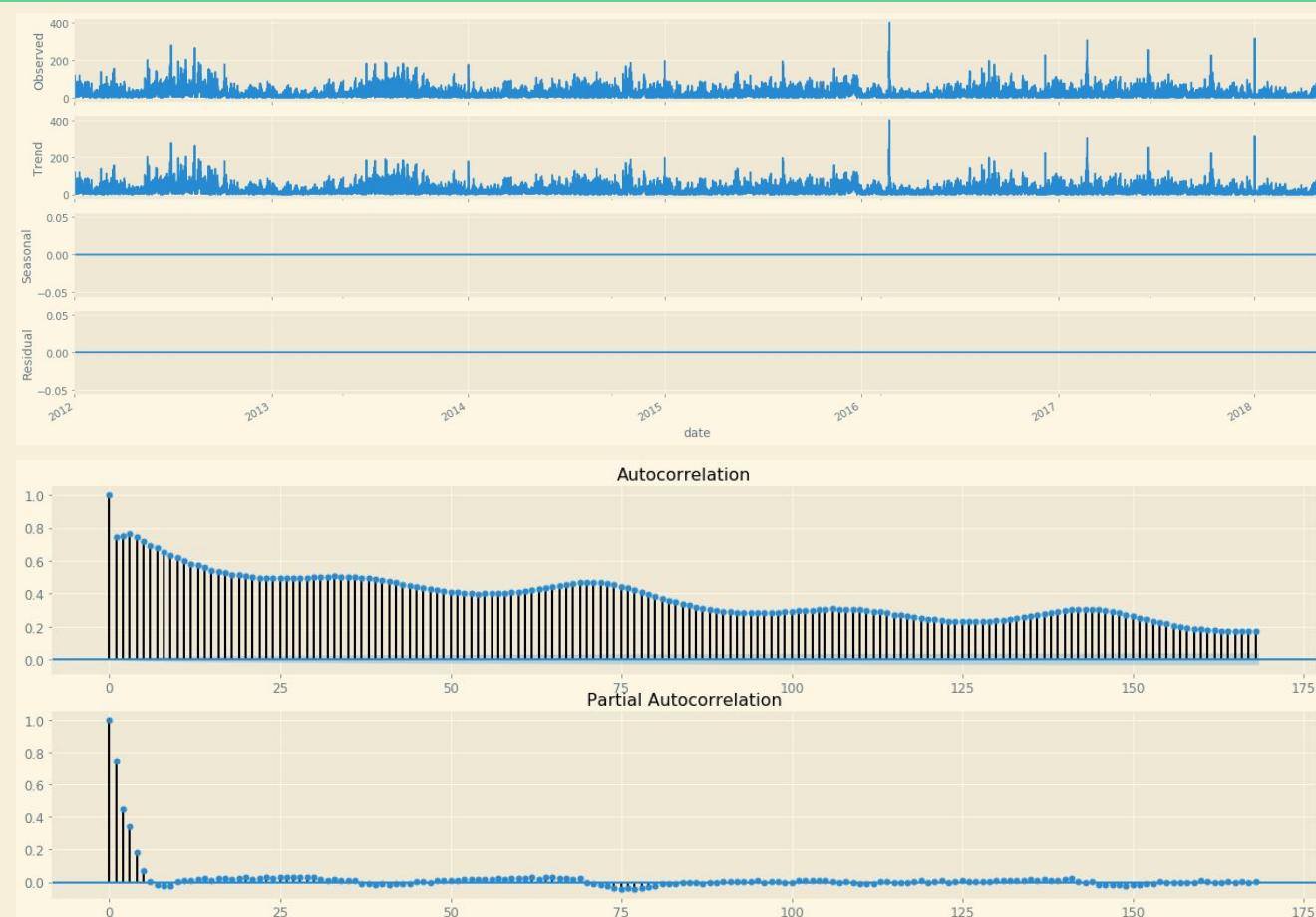
Null Hypothesis: Unit Root Present
Test Statistic < Critical Value => Reject Null
P-Value =< Alpha(.05) => Reject Null

| | |
|--------------------|---------|
| Test Statistic | -21.758 |
| p-value | 0.000 |
| #Lags Used | 77.000 |
| Nr of Obs Used | 166386 |
| Critical Value 1% | -3.430 |
| Critical Value 5% | -2.862 |
| Critical Value 10% | -2.567 |

Results of KPSS Test:

Null Hypothesis: Data is Stationary
Test Statistic > Critical Value => Reject Null
P-Value =< Alpha(.05) => Reject Null

| | |
|---------------------|--------|
| Test Statistic | 1.330 |
| p-value | 0.010 |
| Lags Used | 77.000 |
| Critical Value 10% | 0.347 |
| Critical Value 5% | 0.463 |
| Critical Value 2.5% | 0.574 |
| Critical Value 1% | 0.739 |



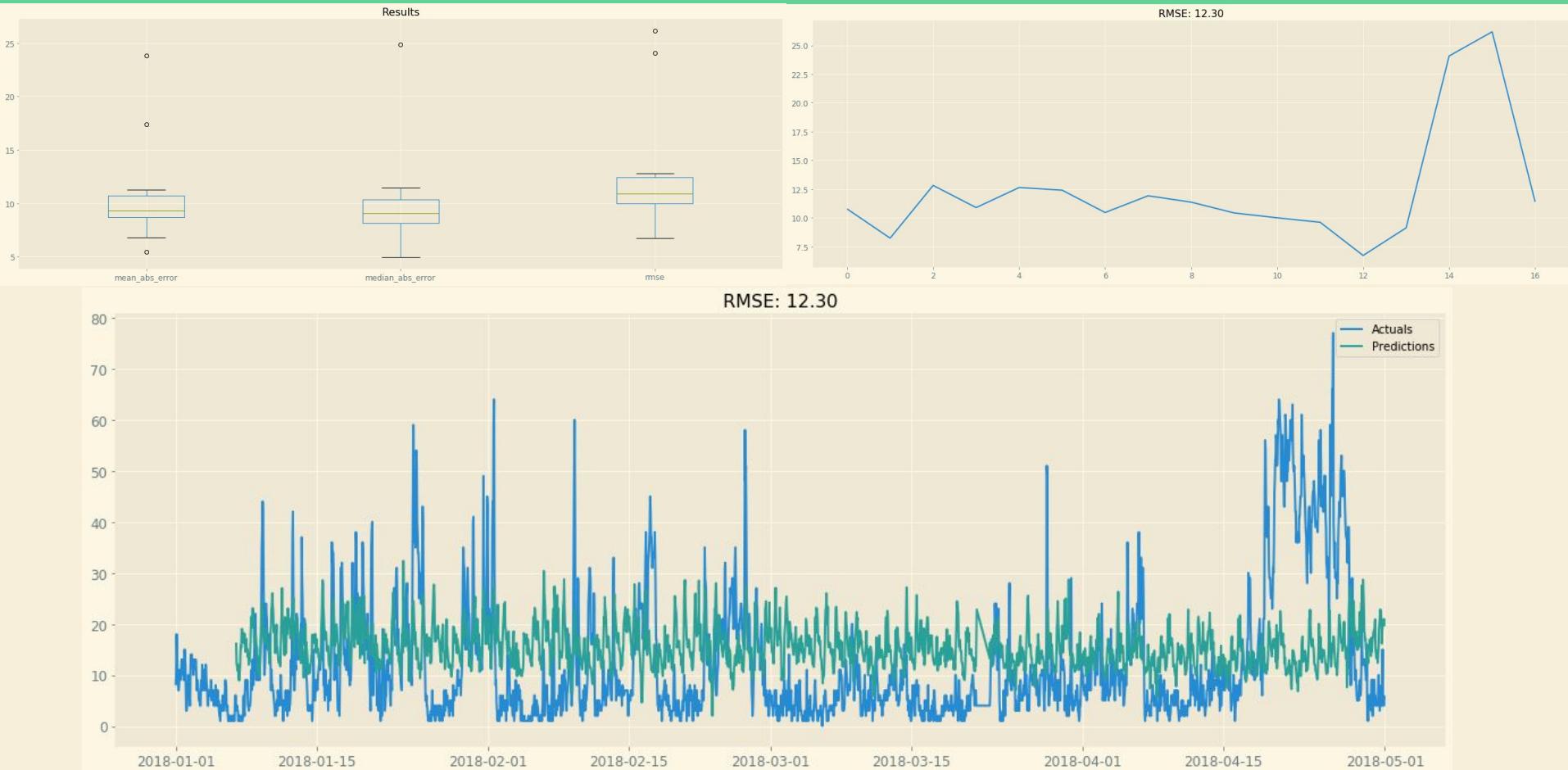
Station nr 28079008 - Parque del Retiro/ Escuelas Aguirre



Station nr 28079018 - San Isidro/Calle Farolillo

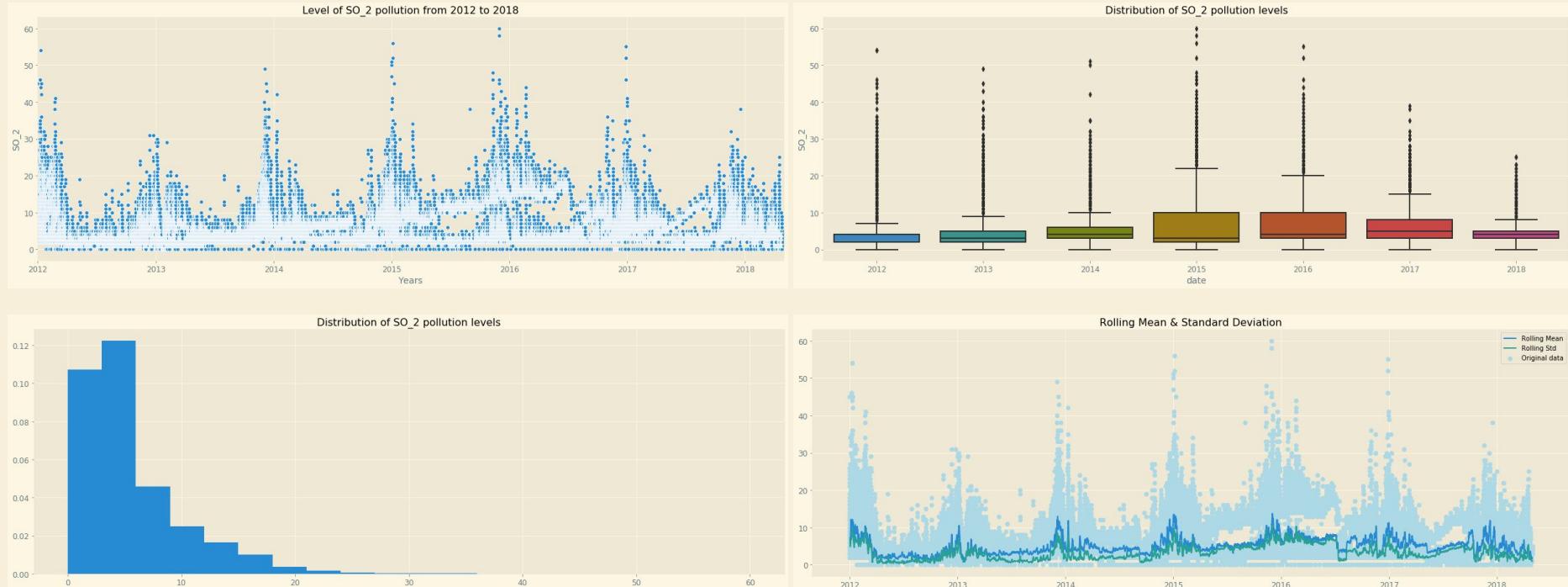


Station nr 28079018 - San Isidro/Calle Farolillo



SO₂

Visualising the target variable



Testing Stationarity

Results of Dickey-Fuller Test:

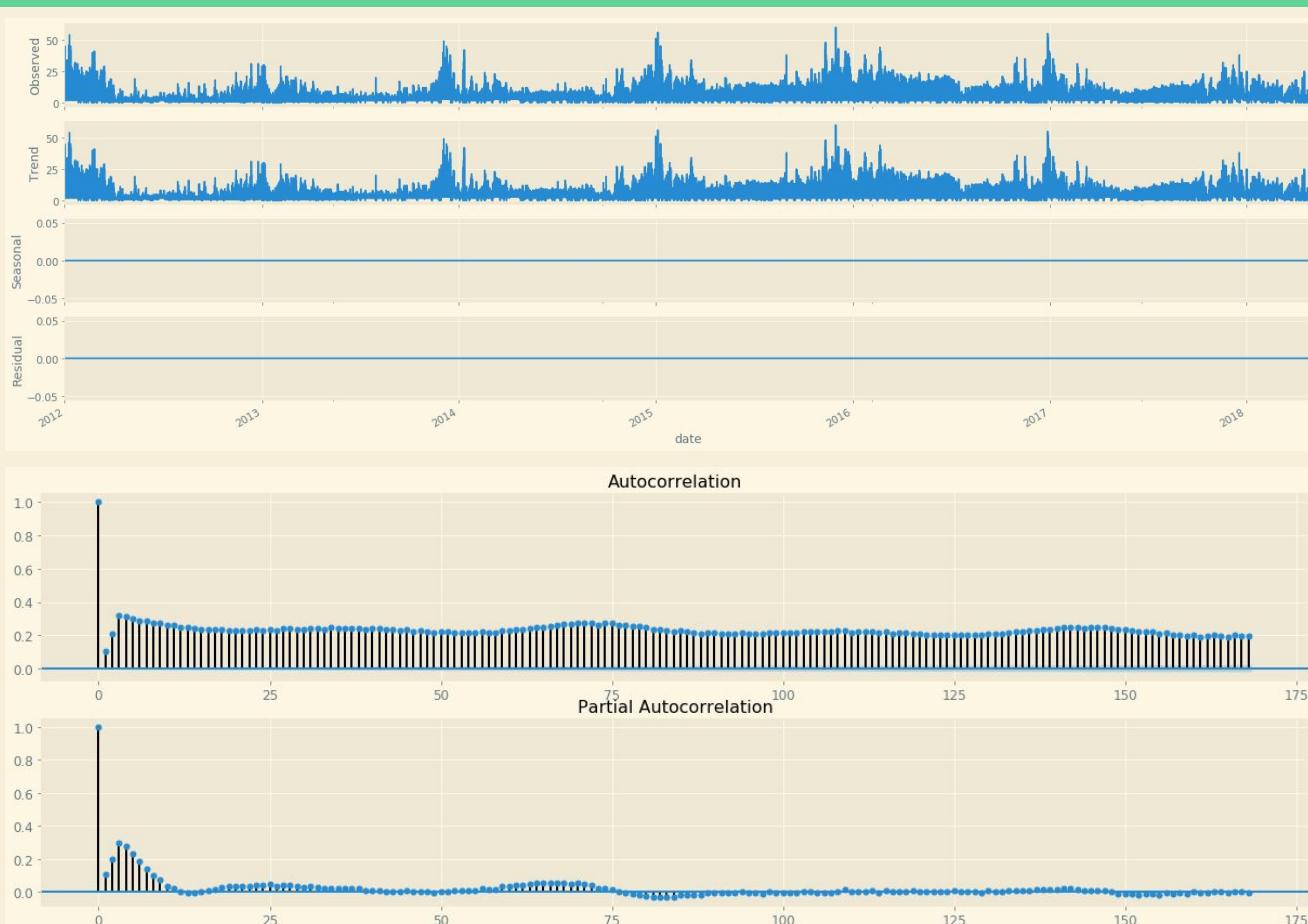
Null Hypothesis: Unit Root Present
Test Statistic < Critical Value => Reject Null
P-Value =< Alpha(.05) => Reject Null

| | |
|--------------------|-----------|
| Test Statistic | -8.382 |
| p-value | 2.490e-13 |
| #Lags Used | 7.700e+01 |
| Nr of Obs Used | 1.664e+05 |
| Critical Value 1% | -3.430 |
| Critical Value 5% | -2.862 |
| Critical Value 10% | -2.567 |

Results of KPSS Test:

Null Hypothesis: Data is Stationary
Test Statistic > Critical Value => Reject Null
P-Value =< Alpha(.05) => Reject Null

| | |
|---------------------|--------|
| Test Statistic | 32.792 |
| p-value | 0.010 |
| Lags Used | 77.000 |
| Critical Value 10% | 0.347 |
| Critical Value 5% | 0.463 |
| Critical Value 2.5% | 0.574 |
| Critical Value 1% | 0.739 |



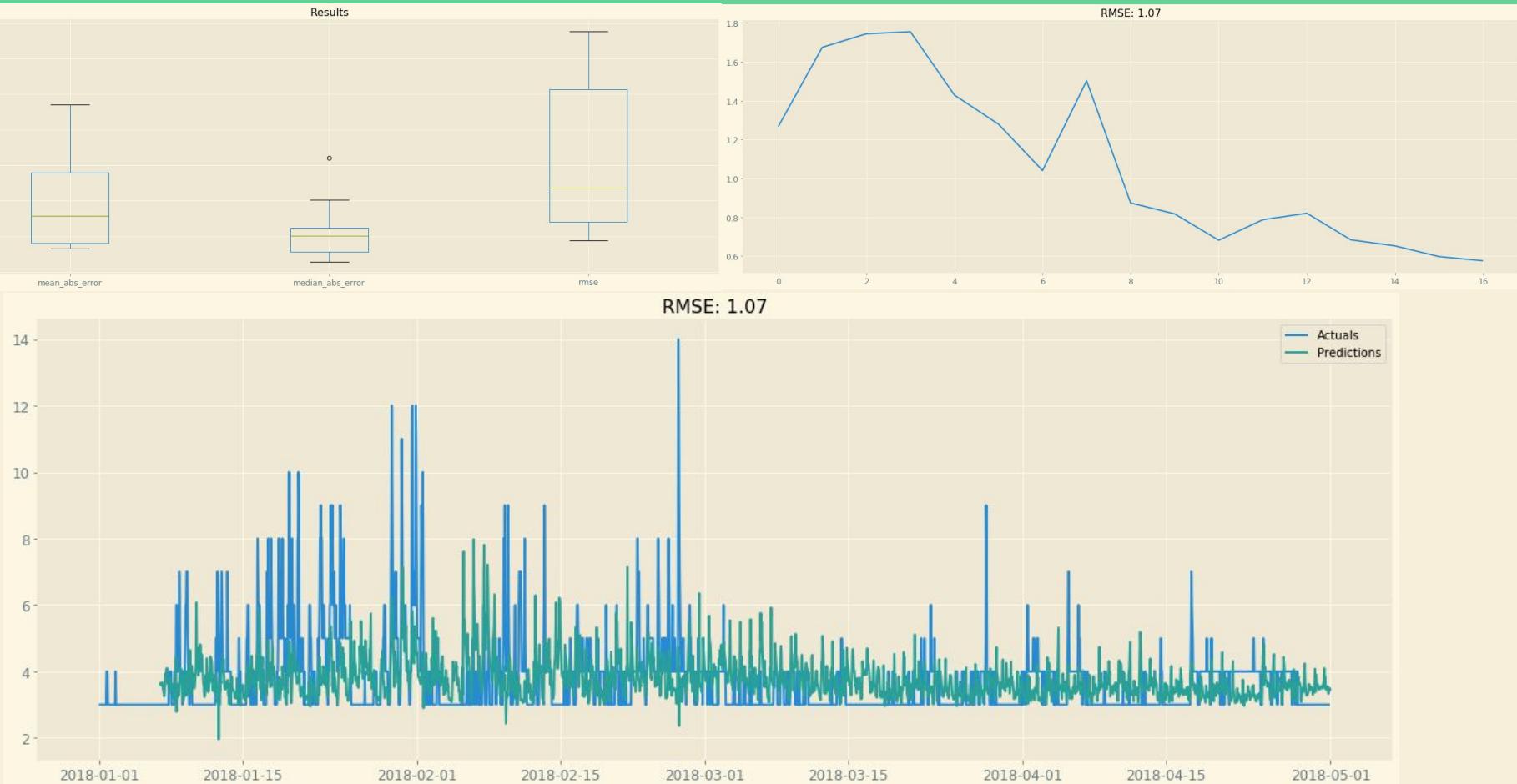
Station nr 28079008 - Parque del Retiro/ Escuelas Aguirre



Station nr 28079018 - San Isidro/Calle Farolillo



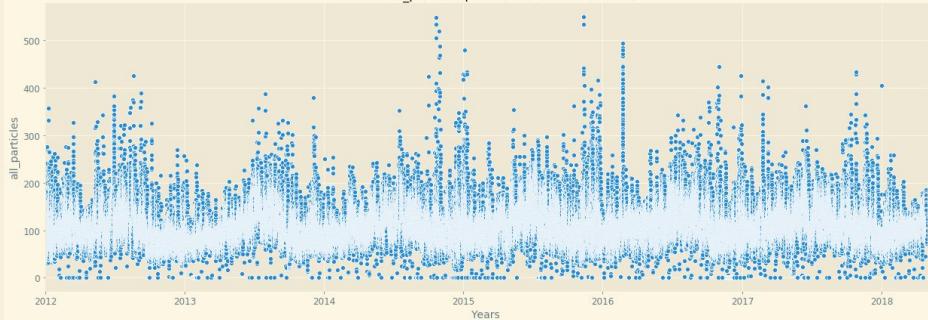
Station nr 28079018 - San Isidro/Calle Farolillo



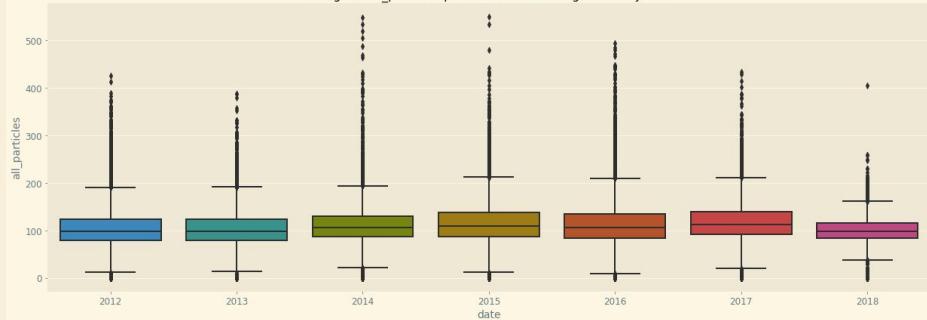
All Particles

Visualising the target variable

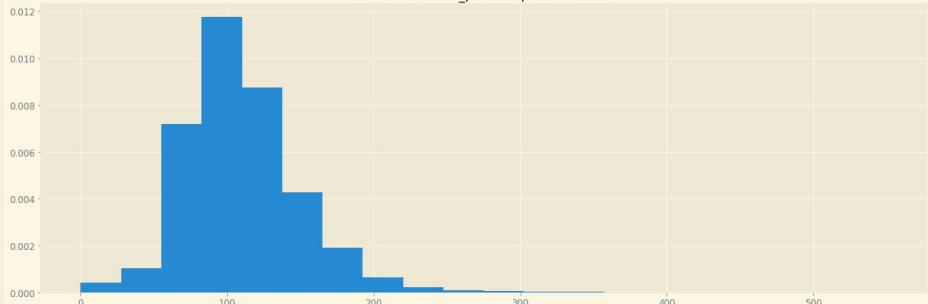
Level of all_particles pollution from 2012 to 2018



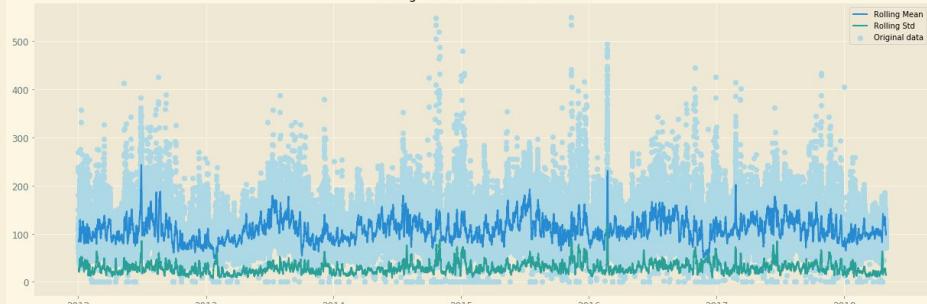
Range of all_particles pollution levels throughout the years



Distribution of all_particles pollution levels



Rolling Mean & Standard Deviation



Testing Stationarity

Results of Dickey-Fuller Test:

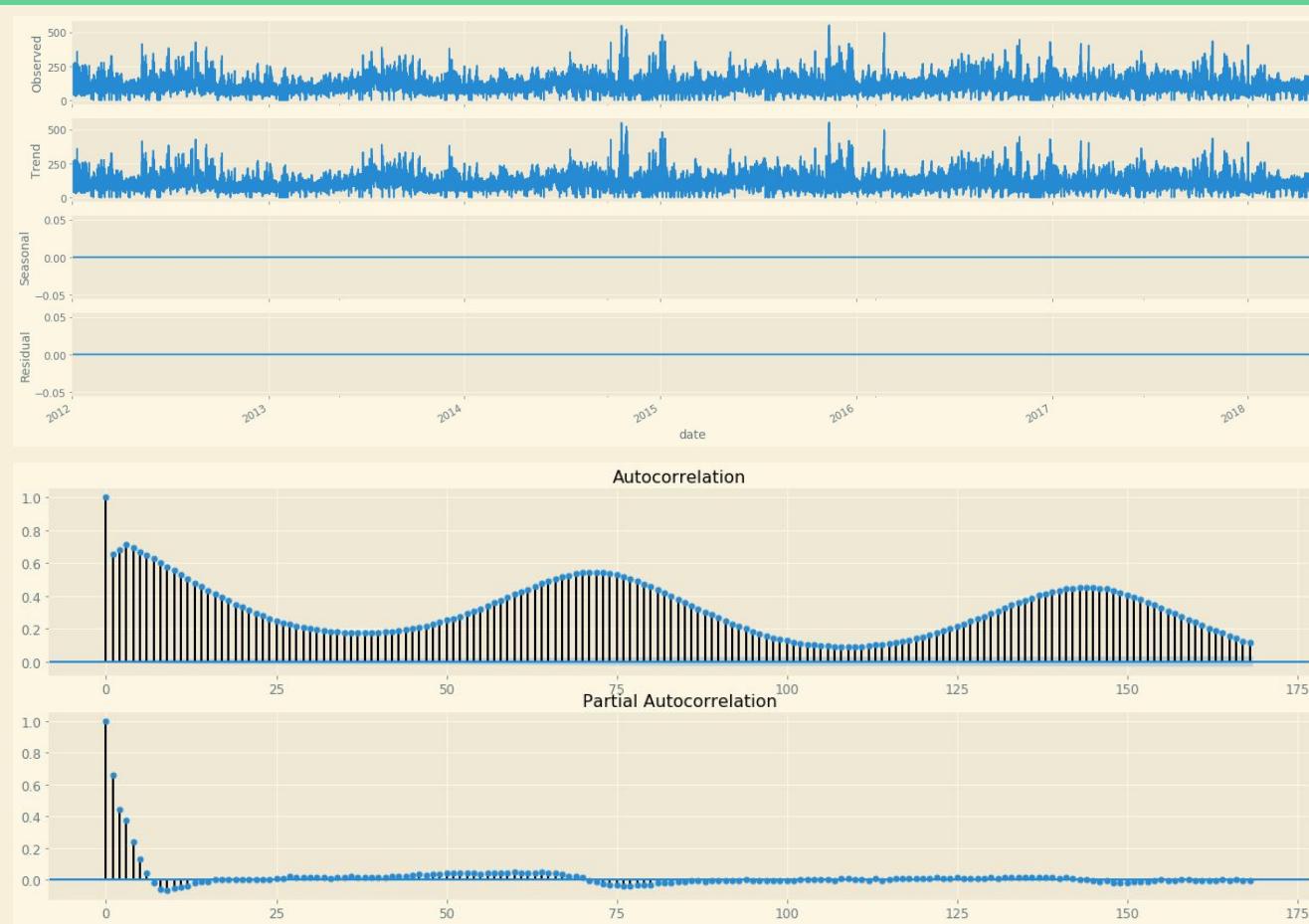
Null Hypothesis: Unit Root Present
Test Statistic < Critical Value => Reject Null
P-Value =< Alpha(.05) => Reject Null

| | |
|--------------------|-------------|
| Test Statistic | -1.568 e+01 |
| p-value | 1.464 e-28 |
| #Lags Used | 7.700 e+01 |
| Nr of Obs Used | 1.664 e+05 |
| Critical Value 1% | -3.430 |
| Critical Value 5% | -2.862 |
| Critical Value 10% | -2.567 |

Results of KPSS Test:

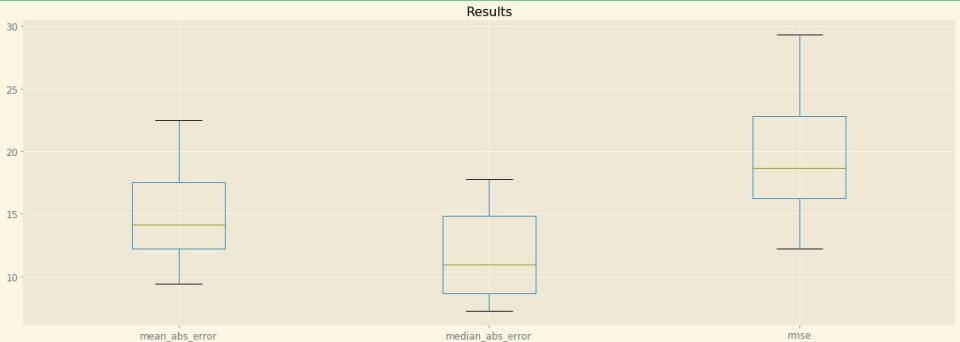
Null Hypothesis: Data is Stationary
Test Statistic > Critical Value => Reject Null
P-Value =< Alpha(.05) => Reject Null

| | |
|---------------------|--------|
| Test Statistic | 5.715 |
| p-value | 0.010 |
| Lags Used | 77.000 |
| Critical Value 10% | 0.347 |
| Critical Value 5% | 0.463 |
| Critical Value 2.5% | 0.574 |
| Critical Value 1% | 0.739 |



Station nr 28079008 - Parque del Retiro/ Escuelas Aguirre

Results



RMSE: 19.99



Station nr 28079018 - San Isidro/Calle Farolillo



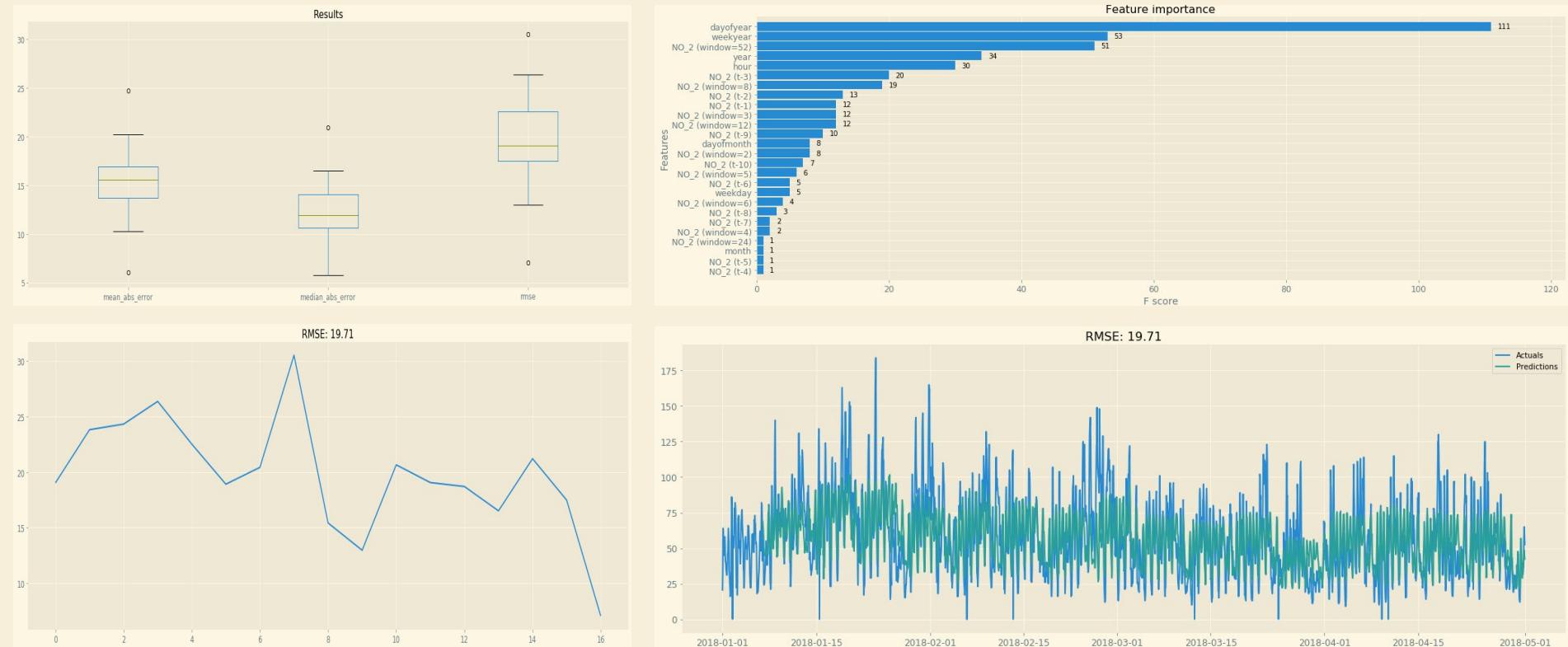
Station nr 28079024 - Casa de Campo



Rolling predictions with XGB

NO₂

Station nr 28079008 - Parque del Retiro/ Escuelas Aguirre

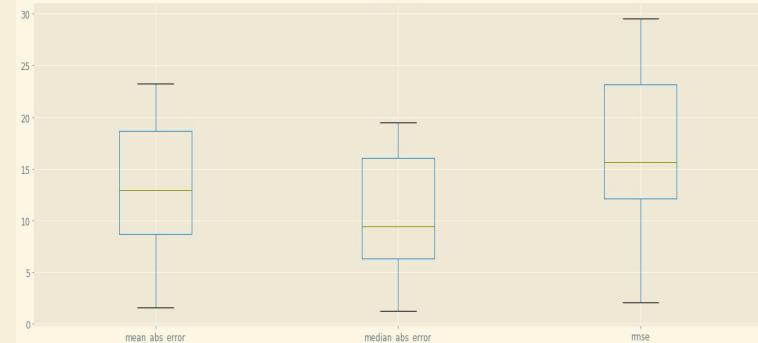


Station nr 28079018 - San Isidro/Calle Farolillo



Station nr 28079024 - Casa de Campo

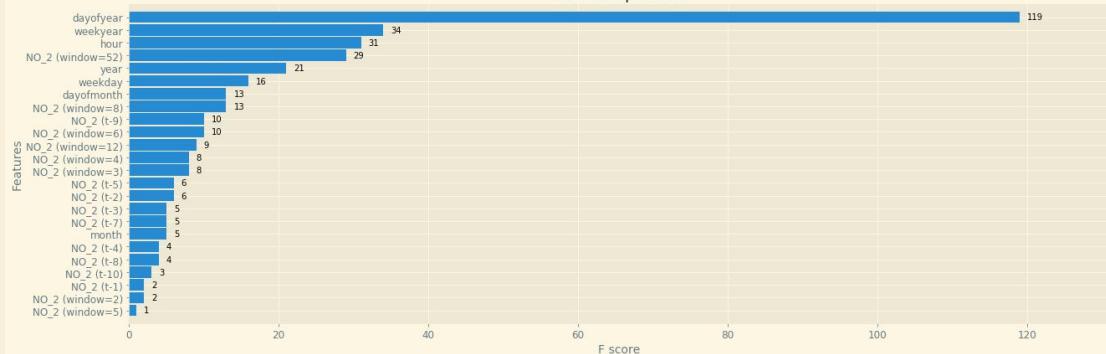
Results



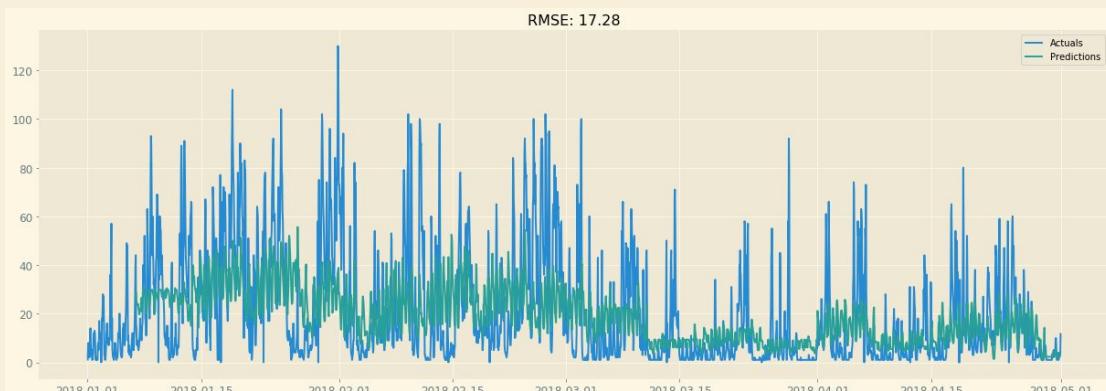
RMSE: 17.28



Feature importance

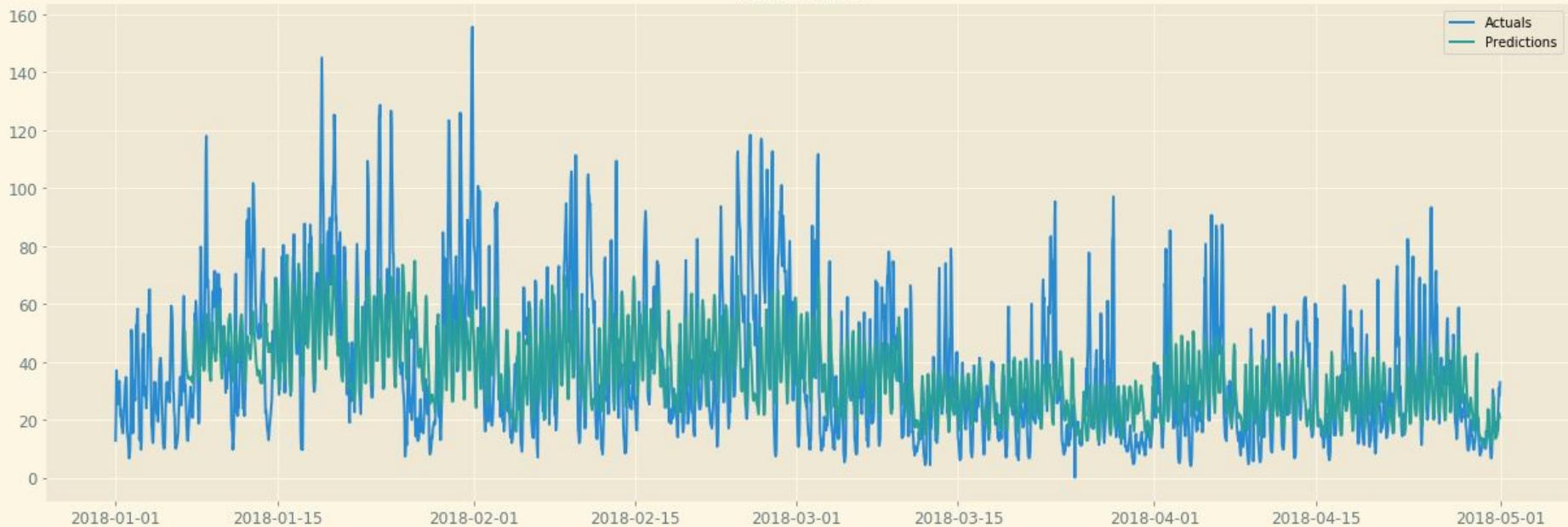


RMSE: 17.28



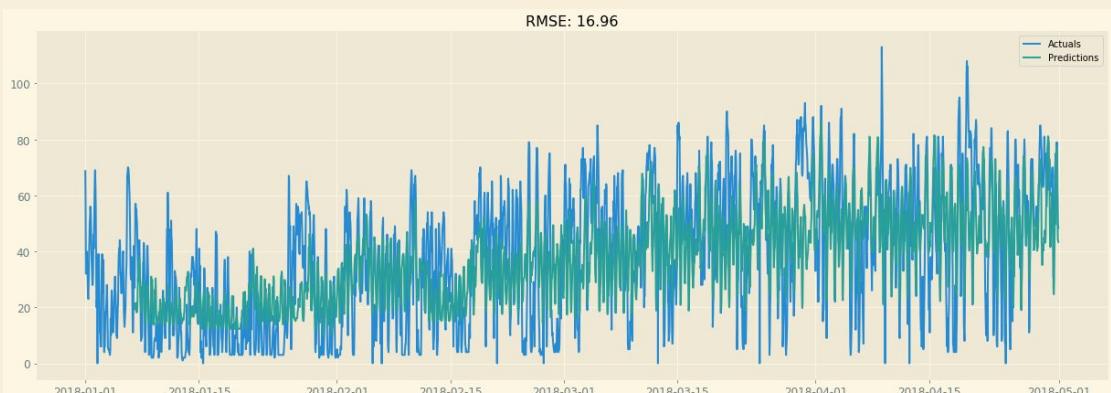
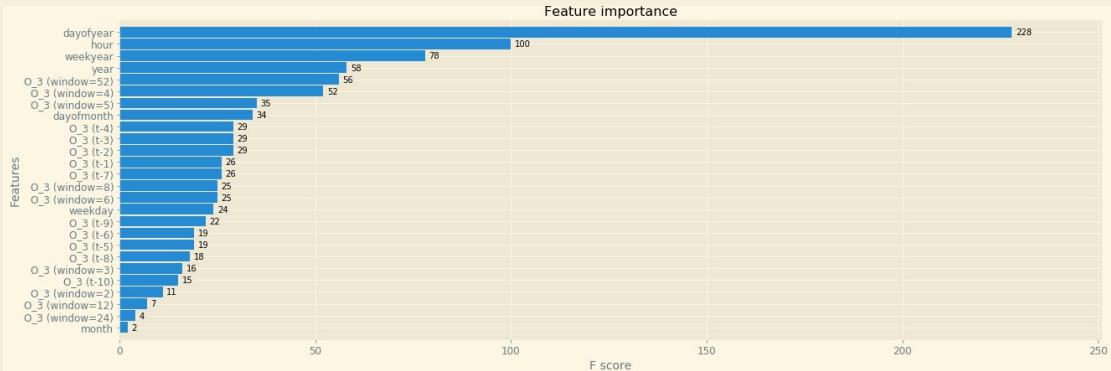
Average of the 3 stations

RMSE: 19.31

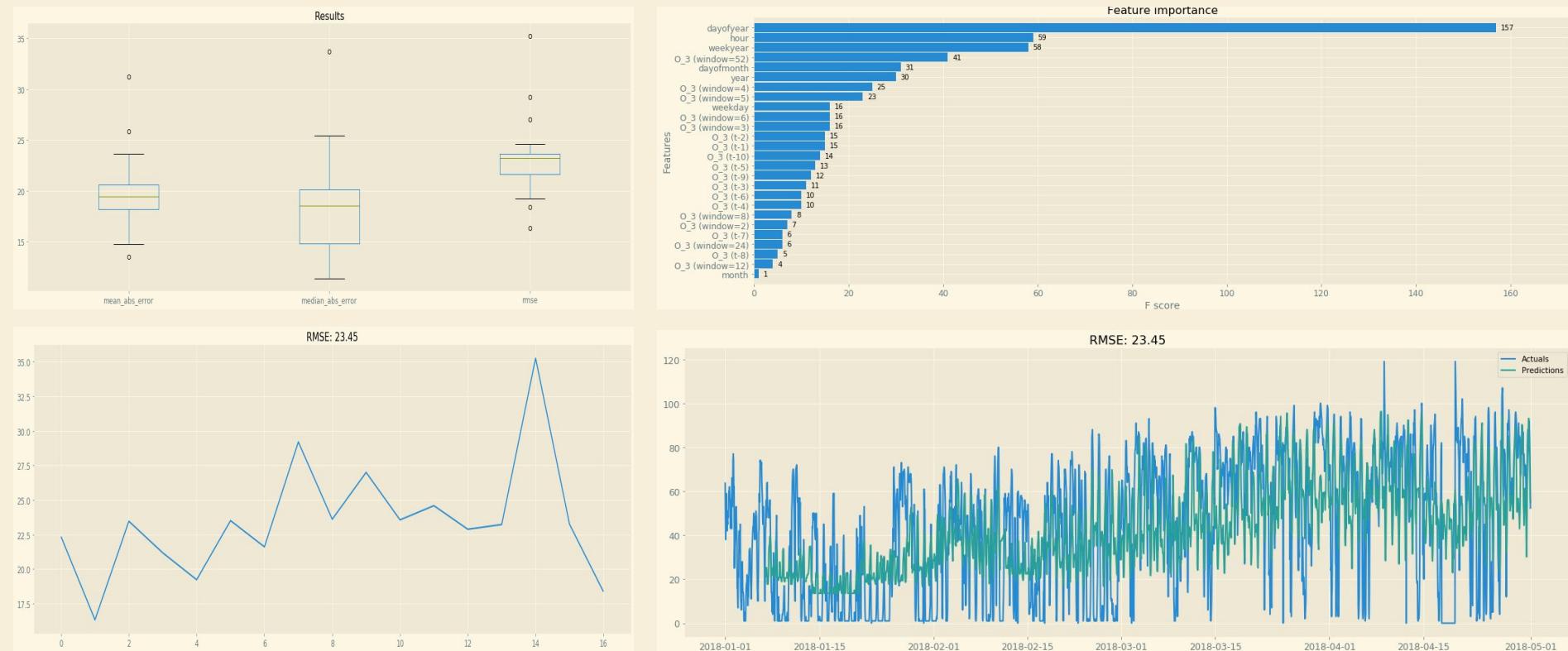


O3

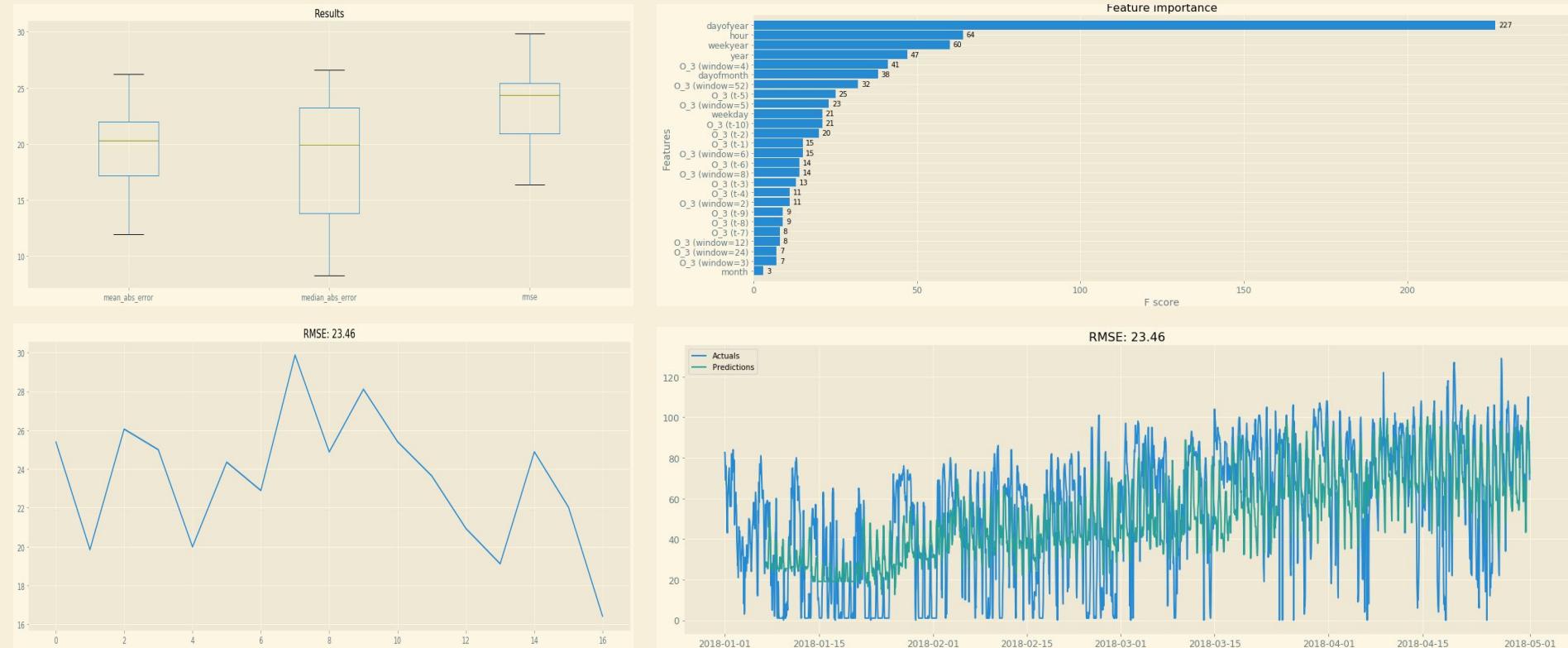
Station nr 28079008 - Parque del Retiro/ Escuelas Aguirre



Station nr 28079018 - San Isidro/Calle Farolillo

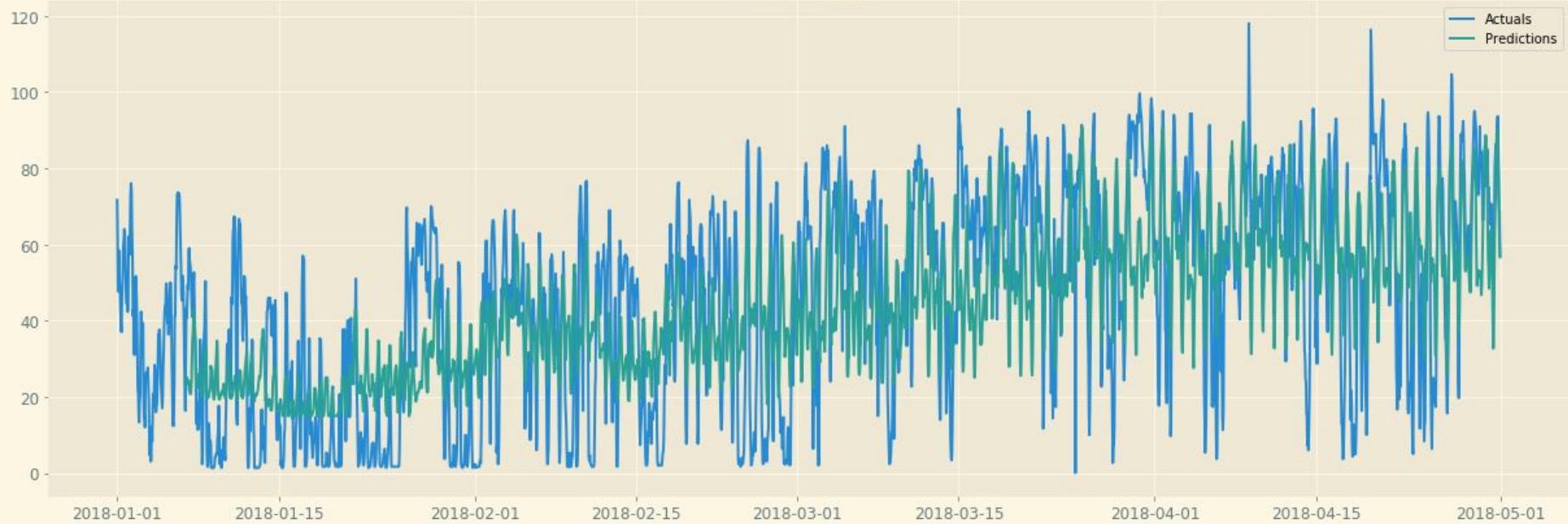


Station nr 28079024 - Casa de Campo



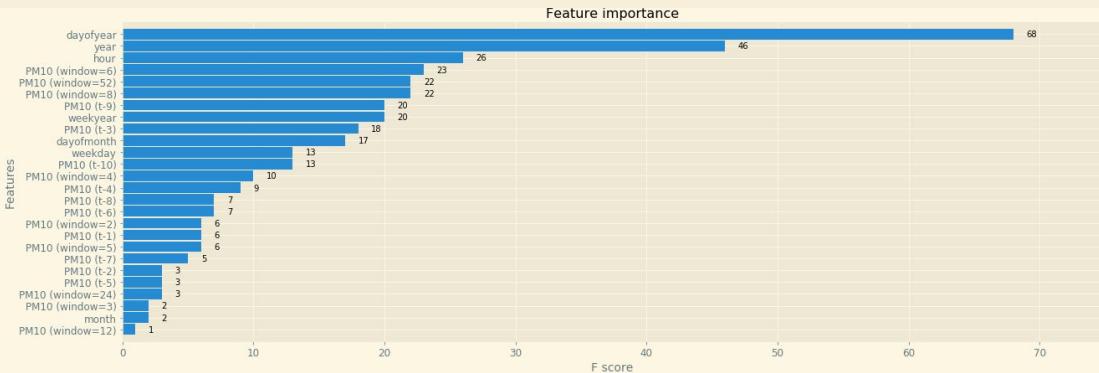
Average of the 3 stations

RMSE: 21.29

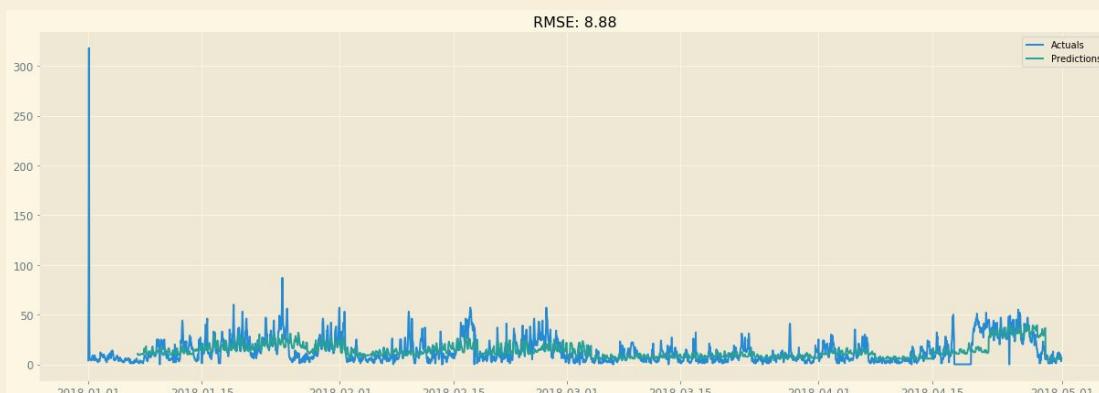
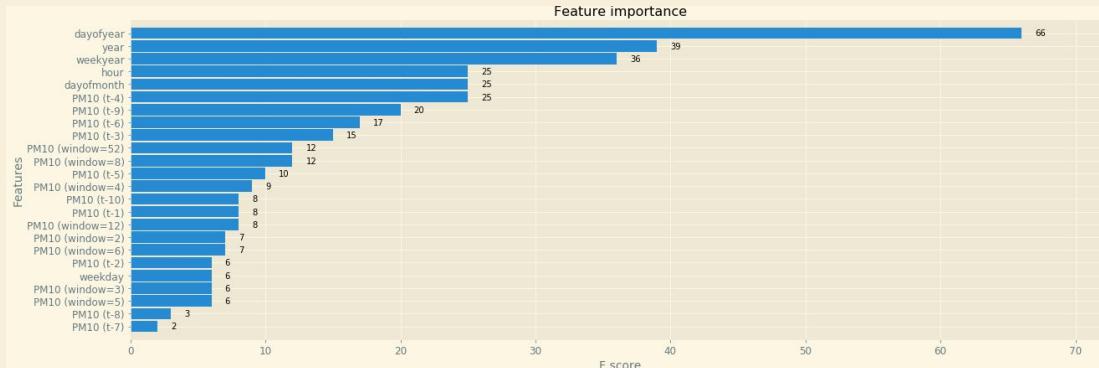
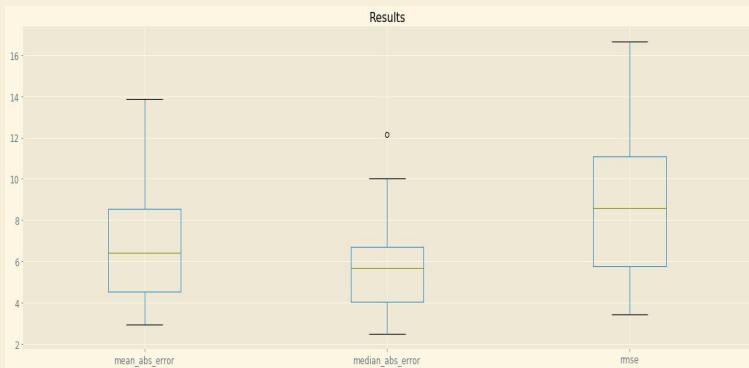


PM10

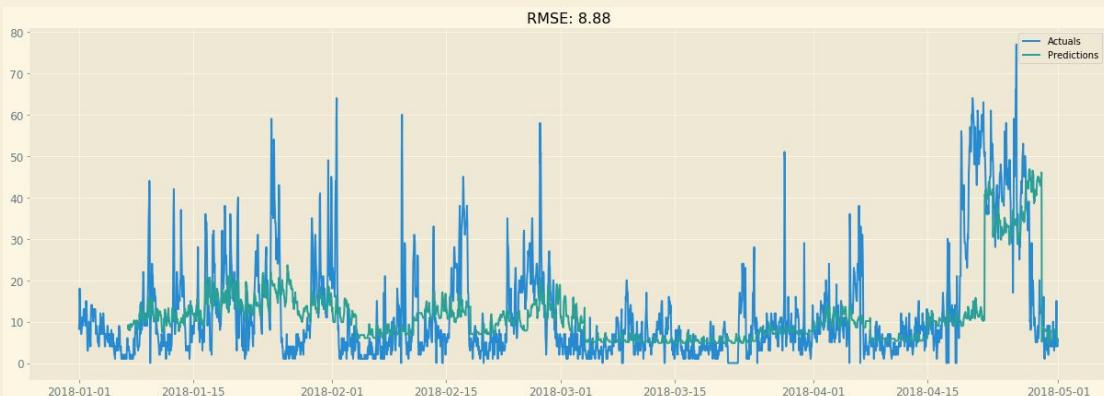
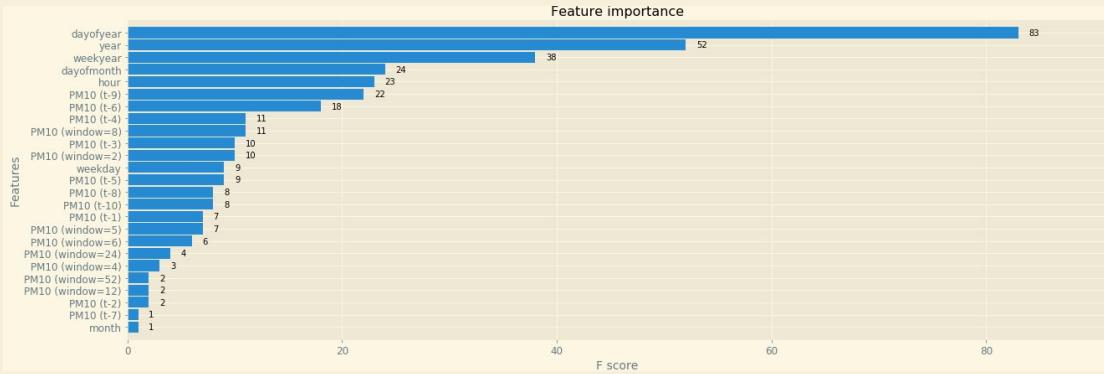
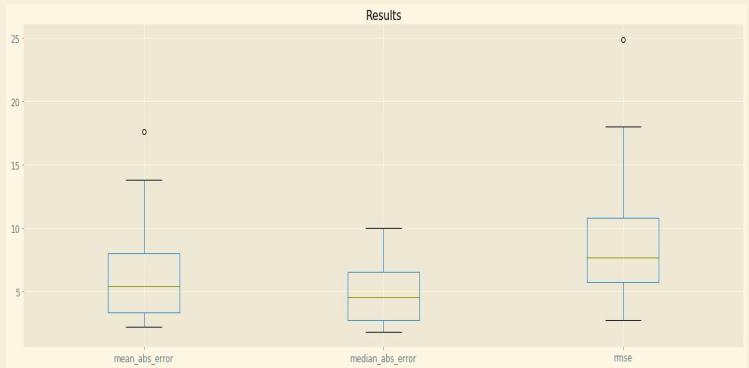
Station nr 28079008 - Parque del Retiro/ Escuelas Aguirre



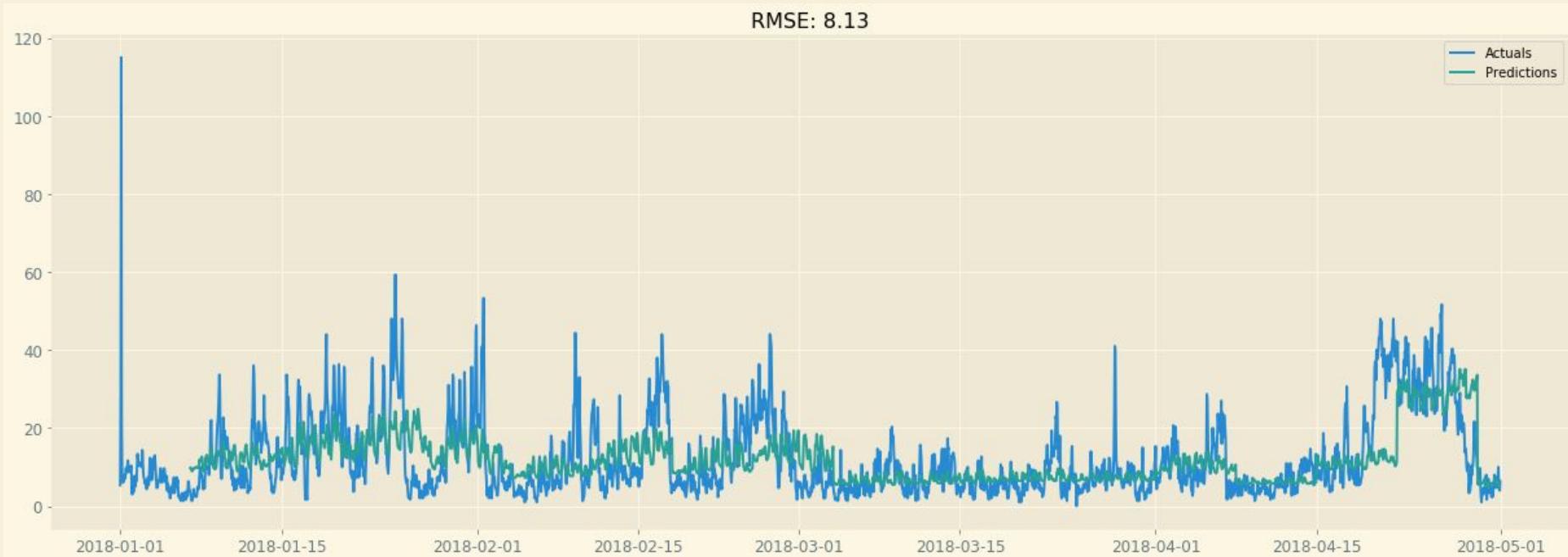
Station nr 28079018 - San Isidro/Calle Farolillo



Station nr 28079024 - Casa de Campo

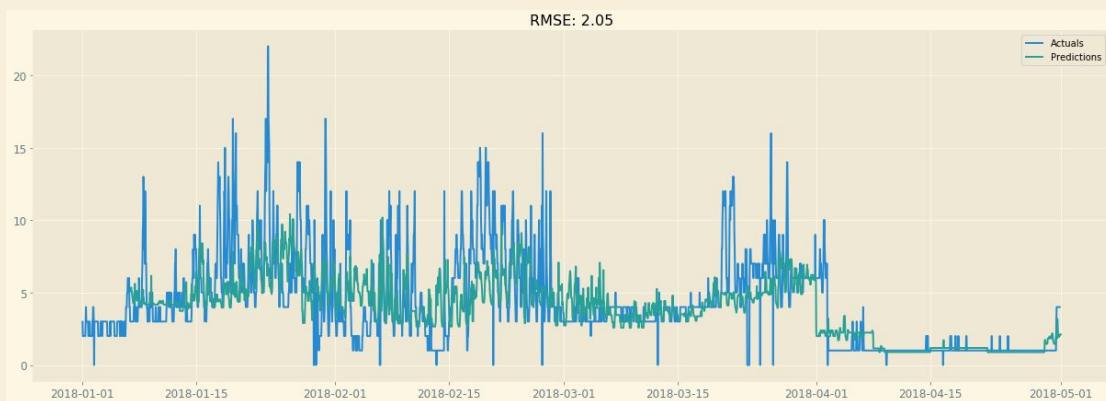
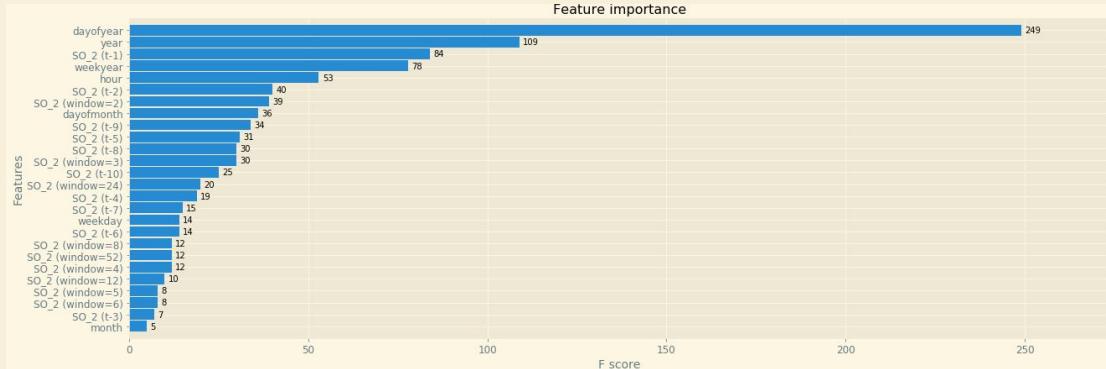
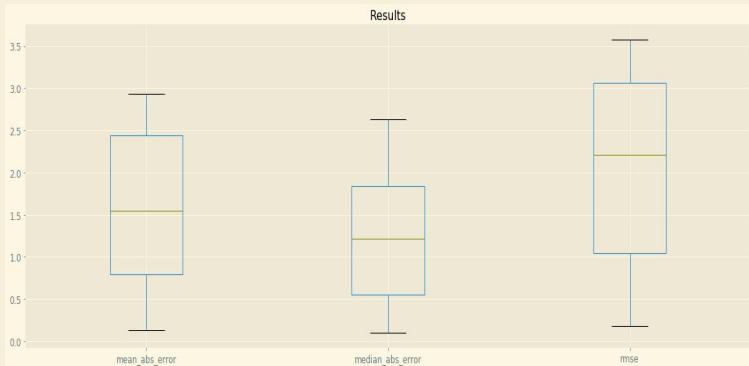


Average of the 3 stations

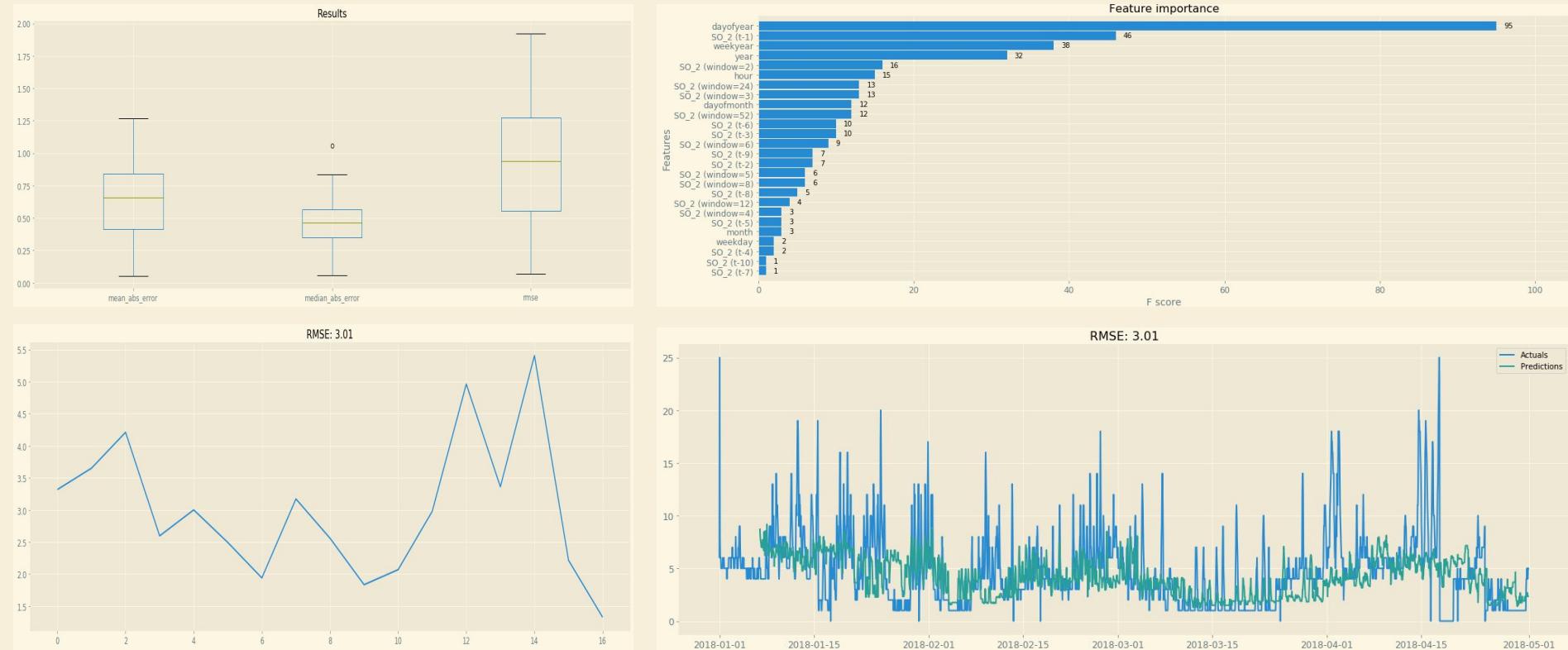


SO₂

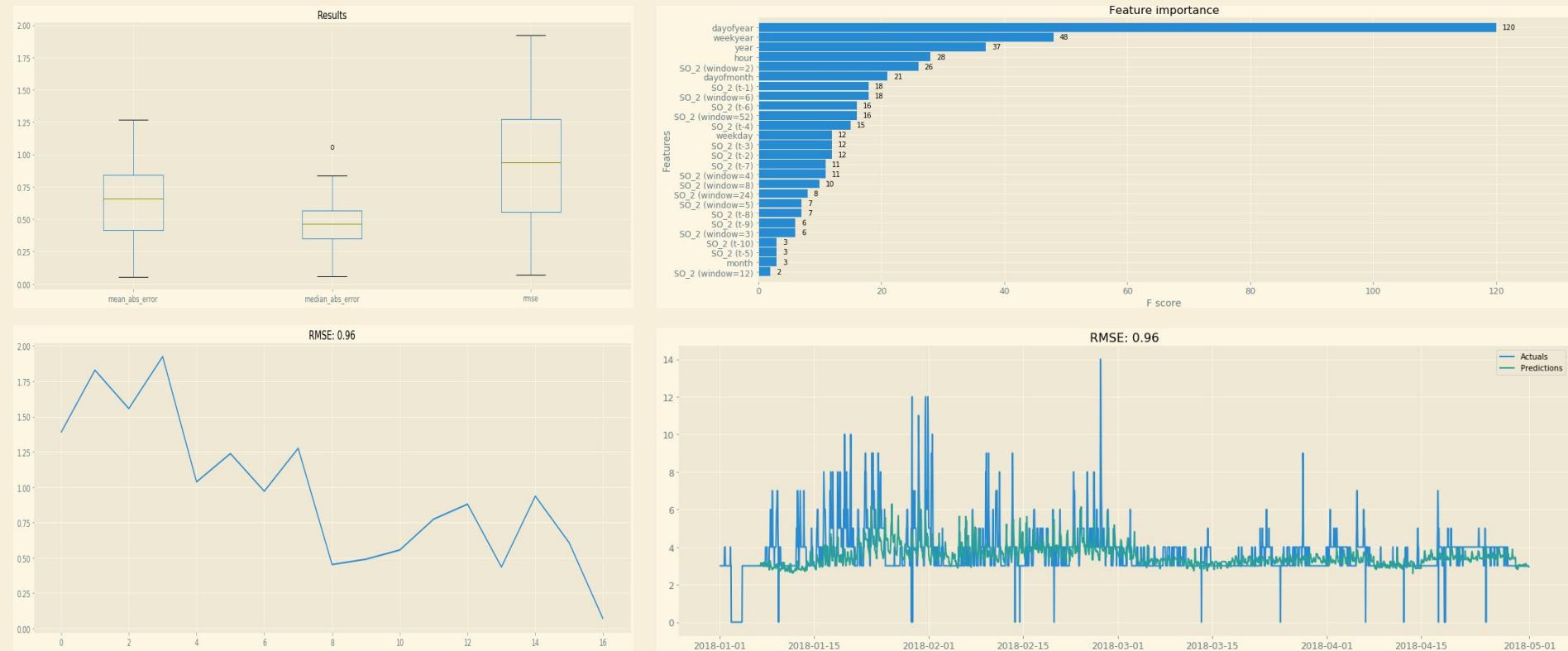
Station nr 28079008 - Parque del Retiro/ Escuelas Aguirre



Station nr 28079018 - San Isidro/Calle Farolillo

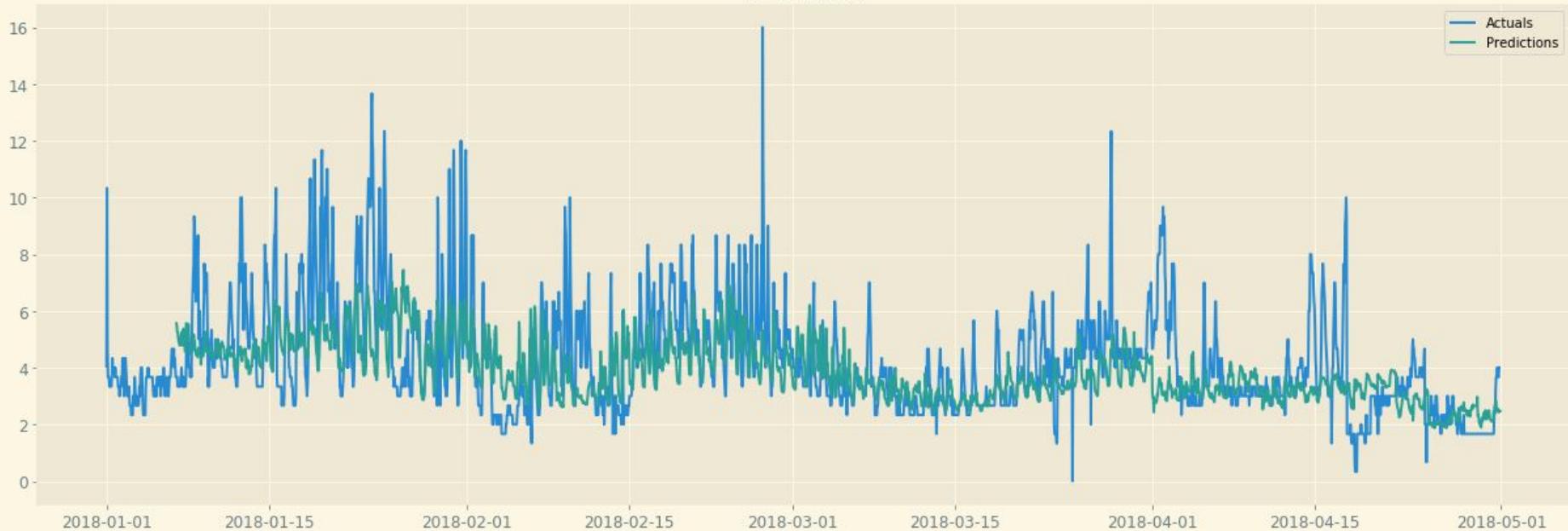


Station nr 28079024 - Casa de Campo



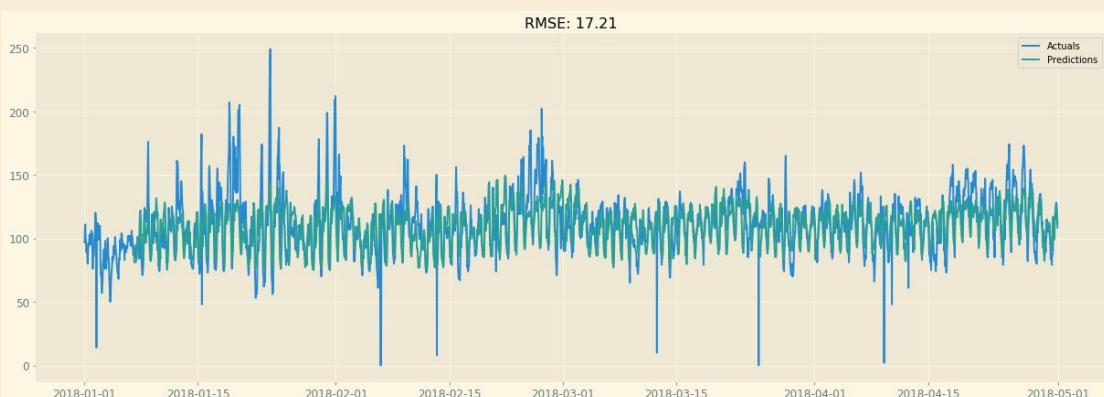
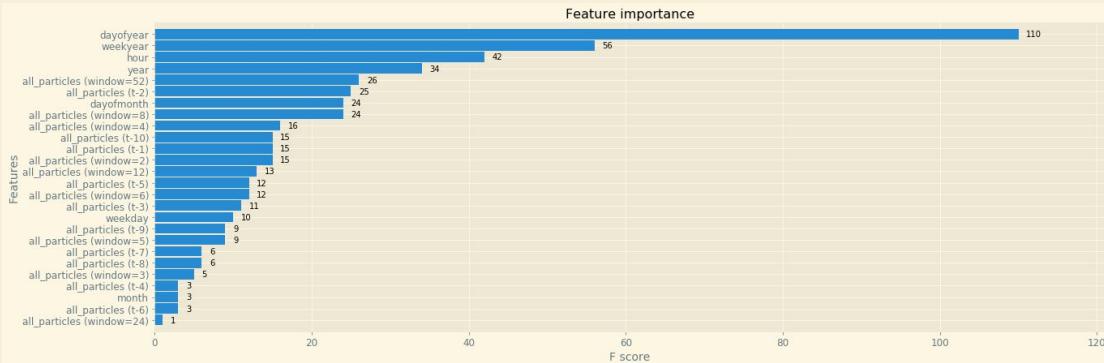
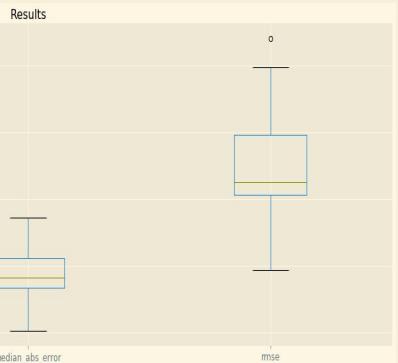
Average of the 3 stations

RMSE: 2.01

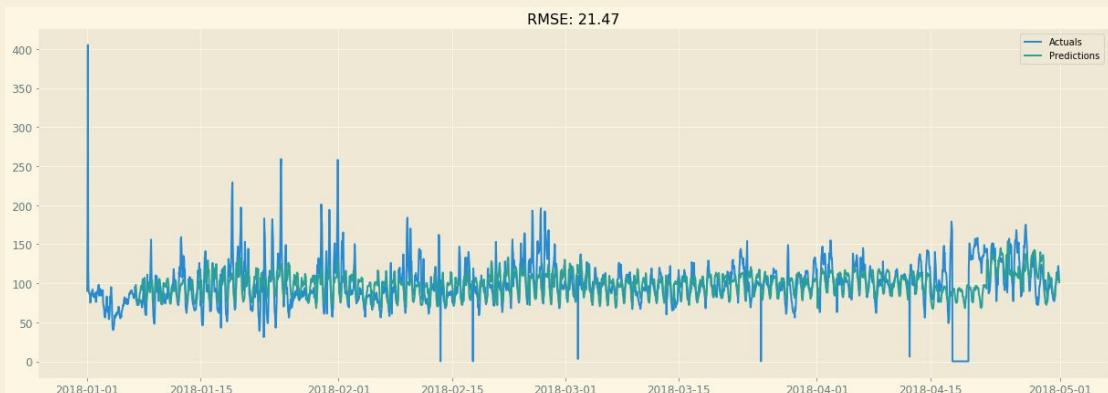
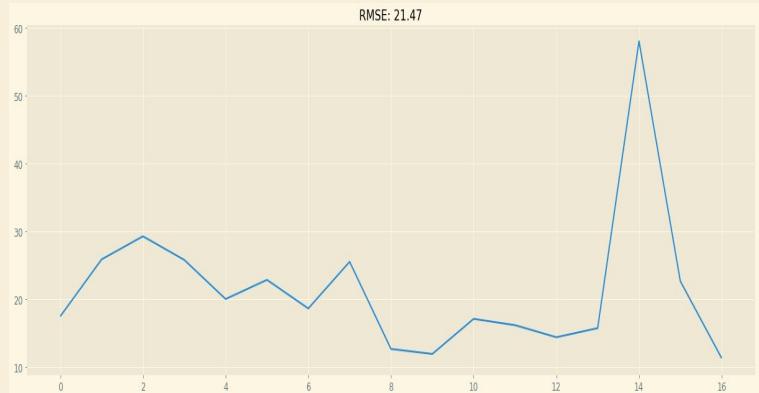
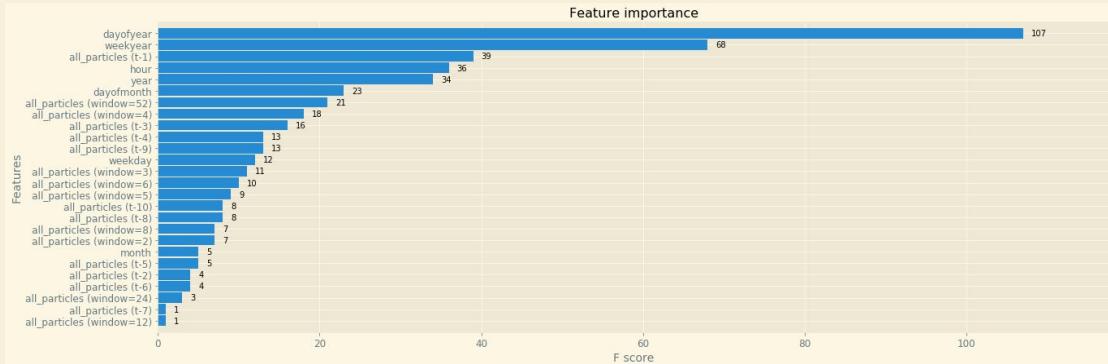
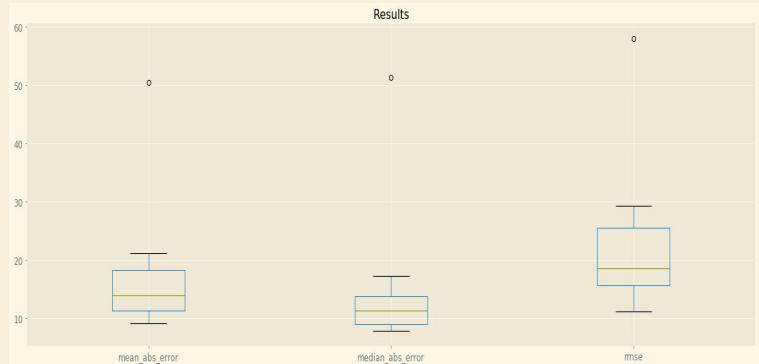


All Particles

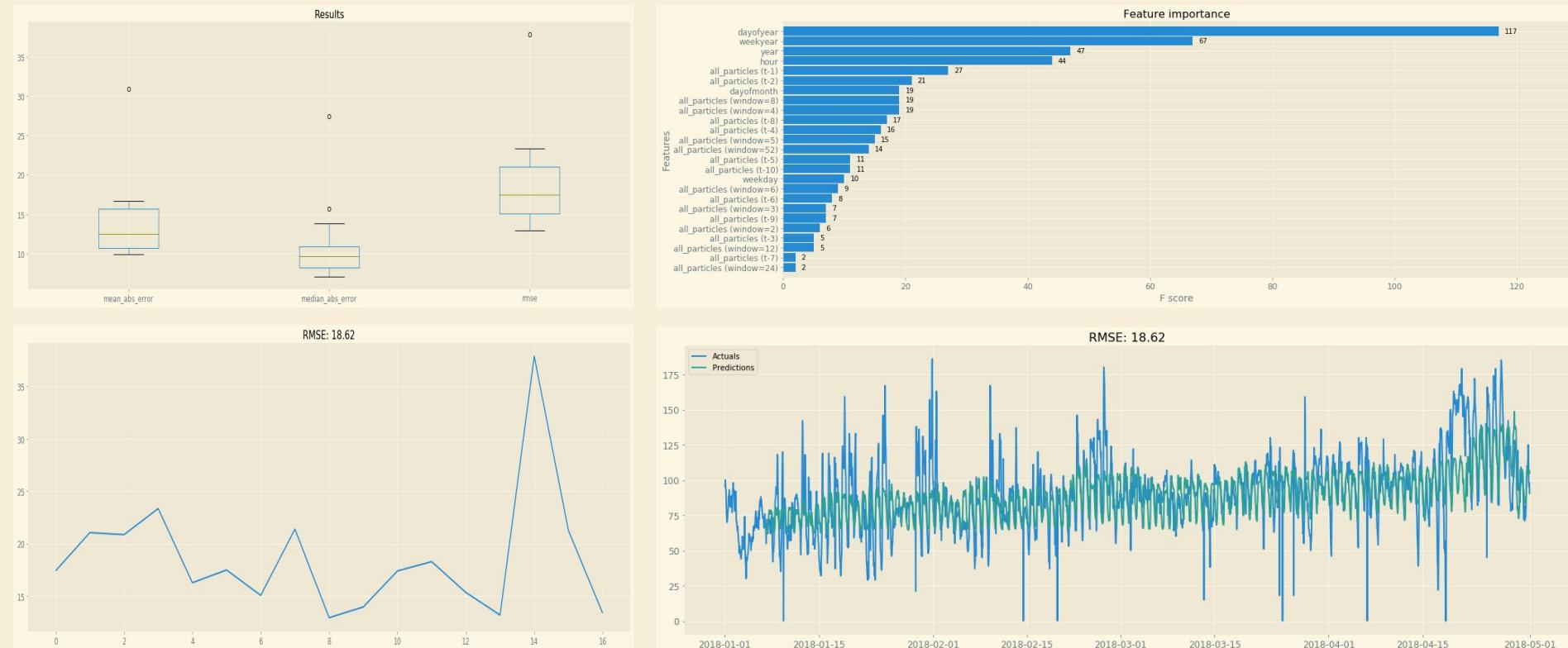
Station nr 28079008 - Parque del Retiro/ Escuelas Aguirre



Station nr 28079018 - San Isidro/Calle Farolillo



Station nr 28079024 - Casa de Campo



Average of the 3 stations

RMSE: 19.10

