

# Air Quality in India

## Forecasting and Analytics

Anit Mathew and Ritwik Katiyar

### Summary

*Our project aim is to analyze air pollution in India. To predict what the situation will be in the next year with reference to time. It turns out that winter is the most polluted time of the year. Pollution levels decrease during the summer and monsoons. The lowest contamination levels were recorded in August and September. Air pollution in June 2022 was the highest compared to June in the last five years. Two of our forecast models indicate that there will be no change in the air pollution levels and one of our model indicates that the pollution levels would slightly decrease. Due to such discrepancy we can not make any strong claims regarding predictions and would possibly need a different machine learning model to archive substantial results.*

## Introduction

Air Pollution is the contamination of the indoor or outdoor environment by any chemical, physical or biological agent that modifies the natural characteristics of the atmosphere. PM2.5 are tiny particles in the air that reduce visibility and cause the air to appear hazy when levels are elevated. It is important to understand what level of Air pollution is surrounding us. India has been at the top of the index for the last few years. We have gathered aggregated data on Air pollution in India over 5 years from 2017 to 2022. In this project we will try to find in which month of the year air pollution is the highest in the country, is there any relationship between time of the year? Further, we will try to answer the question of whether air pollution will increase or decrease in India through machine learning.

The dataset used was uploaded by 'fedesoriano' as [Kaggle-Dataset: Air Quality Data in India](#). The dataset contains the following variables:

- Timestamp Sort: Date and time when the data was collected.
- Year-sort: Year when the data was collected.
- Month-sort: Month when the data was collected.
- Day-sort: The day when the data was collected.
- Hour-sort: Hour when the data was collected.
- PM2.5: PM2.5 level of the country when the data was collected.

The first five rows of our dataset are as follows:

	Timestamp	Year	Month	Day	Hour	PM2.5
0	2017-11-07 12:00:00	2017	11	7	12	64.51
1	2017-11-07 13:00:00	2017	11	7	13	69.95
2	2017-11-07 14:00:00	2017	11	7	14	92.79
3	2017-11-07 15:00:00	2017	11	7	15	109.66
4	2017-11-07 16:00:00	2017	11	7	16	116.5

The Queries for our project are as follows:

1. Analyzing the air pollution dataset year-wise and visualizing the current scenario. Is the air pollution situation in India every year or is there any difference?
2. Identifying if there is any relationship between Air pollution and months of the year. Is the air pollution in India the same throughout the year or does it fluctuate on a month to month basis.?
3. Predicting Air pollution based on PM 2.5 levels. What will be the air pollution for the year 2023, will it increase or decrease?

## Methodology

The data was collected from Kaggle and was easy to download from the website. Upon looking at the website, it was hard to understand the data. The data consisted of hourly data of Air pollution in India for the past years.

1. Data Cleaning: The data contained more than 5000 records. To start with the research, it was important to clean the data.
2. For analyzing the year-wise scenario, we combined the data in python for every year using a group by function and calculated the average air pollution for that particular year. We pulled statistical data and plugged the data into bar plots to better visualize.
3. For identifying the relationship between hour, month, and year, we segregated the data in python for every month wise, hourly wise, for each year separately using the group by function and calculated the average air pollution for that particular year. We pulled statistical data and plugged the data into bar plots to better visualize the scenario. The different approach helps us to gain insight into the actual situation.
4. For Machine learning, before we even start to run any models we need to make sure that our data is ready for machine learning. To accomplish that, we will be altering our dataset to be organized by day instead of an hour, for easy processing. We will be accomplishing that by taking the mean pollution level for every recorded day. We will also be removing any extra columns within our data set such as “day”, “month”, and “year” variables. Finally, we will be creating testing and a training split for our dataset. This would be useful to test our model against pre-existing values. We will be splitting the data into 70% training and 30% testing.

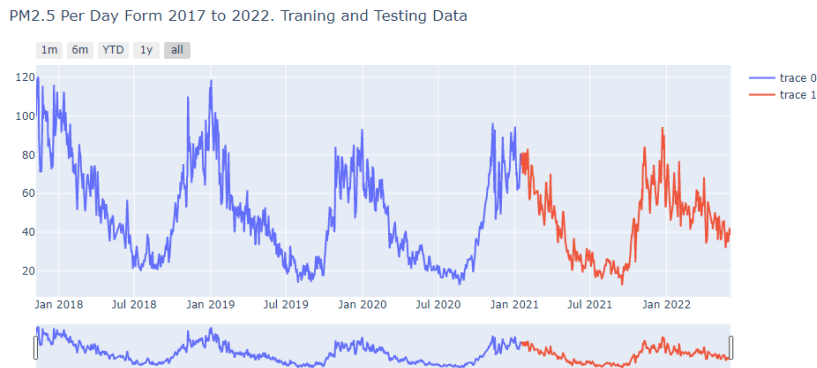
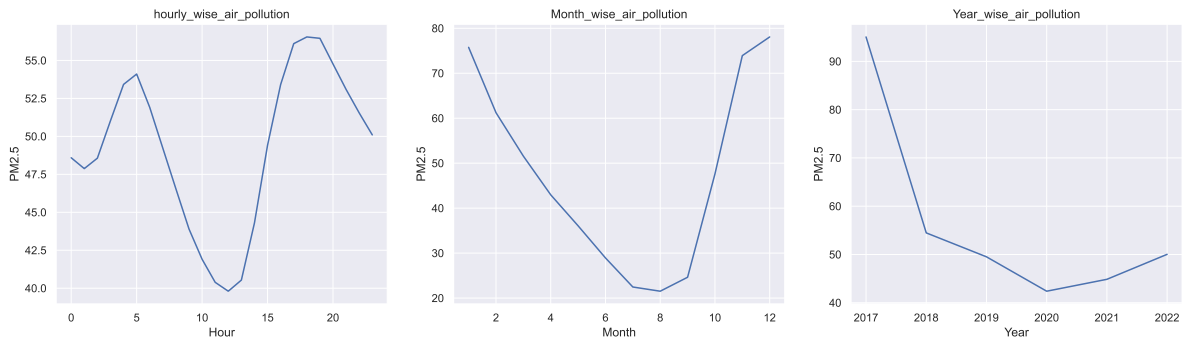


Figure 1: Training and Testing Split

## Data Analysis and Visualizations

What is the situation in India in terms of air pollution? The data is vast and confusing to read. To better analyze the situation, we plotted the data on bar plots. First, we took a broader approach, we cumulated 5-year data into an hour, month and year. The grading of air quality is as follows:

Air Quality Category	PM2.5 $\mu\text{g}/\text{m}^3$ Averaged Over an Hour	PM2.5 $\mu\text{g}/\text{m}^3$ Averaged Over 24 Hours
Good	Less than 25	Less than 12.5
Fair	25-50	12.5-25
Poor	50-100	25-30
Very poor	100-300	50-150
Extremely Poor	More than 300	More than 150

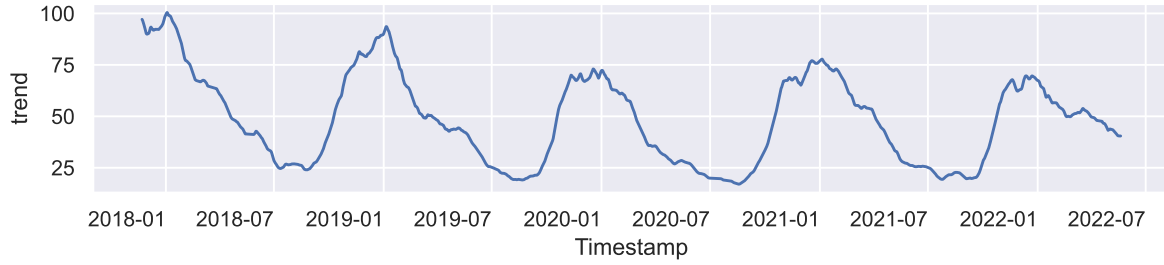


- In the first plot, we cumulated the hour-wise plot. We can see that there is a wave-like formation. Overall the average air pollution is above 40 which brings somewhere between poor and very poor categories. The pollution level peaks two times during the day. One is from 4 to 5 and the other one is from 5 to 7. The lowest it drops during noon.
- In the second plot, we cumulated a month-wise data plot. We can see that the pollution level is at its peak during the start and end of the year. By looking at the plot, we see that July and August have the lowest level of air pollution which is less than 30 which brings them under the fair category.
- In the third plot, we cumulated year-wise data to plot. We can see that except for 2017, each year has average air pollution between 40 to 60. In 2017, air pollution was above 80. To better understand the results, we tried to dig a little deeper.

Bar plot may not provide accurate visualization which in turn hamper our analysis. To accurately summarize the data, we cumulated it and plotted it on a heat map. Hence, The plot below provides a deeper analysis of air pollution in India. We can take note that winters are the most polluted time of the month. The Pollution level has been declining with every coming year. During the summer and monsoon, pollution level declines.

**[Please Refer to the Figure 1 of the Appendix]**

To visualize the general trend of our data and to see if there is an overall increase or decrease in pollution levels. We decided to create a trend plot, which maps out the air pollution trend of our data.



From the overall trend, we can notice that there is a decrease in the peaks or the max values as the days go on. However, we also notice that there are long periods of low pollution levels as well. This trend particularly existed during the COVID years of 2020 and 2021 when the majority of the population was working from home and the overall pollution levels were down. Another way to visualize our data is through box plots to see what the variations of pollution levels were for every day within a month we constructed the following plot.

**[Please Refer to the Figure 2 of the Appendix]**

The plot shows us the distribution of pollution levels per month for every year within our data set. Interestingly we notice that the minimum pollution levels recorded were in the month of September of the year 2021. We can also observe that in January 2019 there was a huge spike in pollution on one of the days. We can also observe that the amount of variation in pollution levels for various months seems to decrease every year.

Now, the question arises, will pollution increase or decrease in 2023? To answer this question, we focused on running models on our data to make predictions based on the various trends and observations we made from our plots.

## Forecasting

### Linear Regression Model

As our first attempt at forecasting air pollution levels, we began by using a simple linear regression model. The linear regression model is used to find a relationship between two variables by fitting a linear equation to the observed data. This is a simple model that should be able to utilize the relationship between time and air quality to generate predictions on what air quality would be like.

The equation for the linear regression line would be as follows for our model:

$$AirQuality = \beta_0 + \beta_{Time} + \epsilon$$

**[Please Refer to the Liner Regression Summary Section of the Appendix]**

From the summary, of our model, we notice that our R-squared value is just 0.063 meaning that our model is only able to explain a mere 6% of the variation in our dataset. We can't use a model that can only explain such a small proportion of variance. We can plot our model's fitted line onto our dataset to better understand why our model failed.

**[Please Refer to the Figure 3 of the Appendix]**

The figure shows us that we can't use a simple linear line to explain such a complex relationship. We noticed from visualizing the data that there is a repeating pattern and if we wish to forecast the repeating pattern based on the past repeating pattern we can't use a simple regression line. To obtain better results, we need to use a different model.

## **ARIMA Modeling**

ARIMA stands for autoregression integration of a moving average. This model is specifically used for forecasting time series data. The equation for the ARIMA model consists of three parts.

The first part is autoregression: Auto regression refers to the changing variable within the model that regresses on its own and has prior values.

The second part is integrated: This refers to the differencing between observations.

The third part is the Moving Average: This takes into account the dependency between an observation and a residual error from the moving average model.

The ARIMA parameters for the model are as follows:

P = liner combination lags.

d = Number of times the raw observations difference.

q = Liner combination of lagged forecast errors. Simply, the size of the moving average window.

Before we even run the model we need to first determine the values for p,q, and d.

To get a better understanding of what ARIMA is mathematically the formula for ARIMA can be given as: Auto Regressive Formula:

$$AR = Y_t = \alpha + \beta_1 Y_{t-1} + \beta_2 Y_{t-2} \dots \beta_p Y_{t-p} + \epsilon_t$$

Moving Average Formula:

$$MA = Y_t = \alpha + \epsilon_t + \phi_1 \epsilon_{t-1} + \phi_2 \epsilon_{t-2} \dots \phi_q \epsilon_{t-q}$$

Hence, the combined ARIMA formula:

$$Y_t = \alpha + AR + MA$$

In non-mathematical terms the formula means Predicted value = a constant + liner combination lags of Y + liner combination of lagged forecast errors.

Before we even run the model we need to first determine the values for p,q, and d.

## Determining the d-value

D-value refers to differencing. Which in turn refers to making the time series stationary. To test if our dataset is stationary we can utilize what is known as the Augmented Dickey-Fuller test (ADF). ADF is a unit root test, in other words, the test looks for the presence of a 'unit root' (A unit root essentially means a systematic pattern within a data set that is unpredictable.) To determine if the series is non-stationary. we can set up our hypothesis test as follow by taking an  $\alpha = 0.05$ .

$H_o$  = The Data set is not stationary

$H_a$  = The Data set is stationary

ADF Statistic: -2.5720220909338787

p-value: 0.09891869282208893

We reject the Null hypothesis if the p-value < 0.05

From the summary above we can observe that the p-value is greater than 0.05. Which means we do not have significant evidence to reject the null hypothesis. Hence this means our data is non-stationary.

Since our data is non-stationary we can try to differentiate our data set. Differencing simply refers to calculating the difference between each variable within our dataset.

$$y_t = Air_{pollution}$$

$$y'_t = y_t - y_{t-1}$$

After differencing we can then run the ADF test again to see if differencing helped make our data stationary. If our data is still not stationary we would need to differentiate again. Running the test on the differenced data we get.

Setting up a new hypothesis test with an  $\alpha$  of 0.05

$H_o$  = The once differenced data set is not stationary

$H_a$  = The once differenced data set is stationary

We reject the Null hypothesis if the p-value < 0.05

ADF Statistic: -13.175156123883413

p-value: 1.227599209964447e-24

After differencing our p-value has decreased significantly meaning we now have enough evidence to reject the null and conclude that our data is now stationary and we were able to accomplish that via only one differencing once. Hence meaning that  $d = 1$ .

[\*\*To See the Full ADF Output Please Refer to the Appendix: ADF Output\*]

## Determining p

We can determine the value of the Auto regressive term p by taking a look at the Partial Autocorrelation plot (PACF). PACF is a plot that displays the correlation between the series and its lag.

Mathematically speaking:

$$PACF = \alpha_0 + \alpha_1 Y_{t-1} + \alpha_1 Y_{t-1} \dots \alpha_q Y_{t-q}$$

Where each term represents a lag.

The PACF for our 1<sup>st</sup>-order differenced data set looks as follows:

**[To See Figure 4 of the Appendix Section]**

We can observe that the PACF lag 2 is quite significant since it is above the significant line (in blue). There seem to be other significant lags however, we will take the lowest number of lags.

### **Determining q**

Similar to how we determined the value of p we can now look at the Autocorrelation Function plot (ACF) to determine the value of q. The value of ACF is a reference to the number of moving average terms or the moving average error of the lag. The following plot shows us the ACF of our data.

**[To See Figure 5 of the Appendix Section]**

From the plot, we can ascertain that the ACF lag 2 is once again quite significant since it is above the significant line (in blue). Now that we have our p, q, and d values we can finally get ready for modeling. However, before we can model we need to first split our model into

### **Model 1 - ARIMA Model (p = 2, d = 1, q = 2)**

We are now ready to run the ARIMA model. However, before we run any predictions we can take a look at our model to see if there are glaring issues that would need to be corrected.

**[Please Refer to the Model-1 Section of the Appendix for Plots and Summary of this model]**

We can see from the “standardized residual for”P” plot above that the residuals are generally spread out which is good since clustering could mean there is a serious issue with our model. We can also notice from the quantile-quantile plot (q-q plot) that our data doesn’t seem to show that our residuals may be a bit skewed. However, since most of the data points fall onto a linear line aside from the data at the ends we can assume the normality for most of our residuals. We can also see that there are no strong correlations within our residuals from the correlogram. Strong correlations would mean that there are correlations our model is unable to explain. Finally, the KDE curve is a little higher than a normal  $N(0,1)$  curve line. Suggesting once again that our residuals could be non-normally distributed. Despite our residuals not being completely perfect we can go ahead and try to make predictions using our model. Keeping in mind that our predictions could be incorrect and better time series models (perhaps a neural network) could yield a better result.

From the summary if we set up a hypothesis test for our AR and MA lags. We get the following:

$H_0$ : The AR Lag 1 is not significant

$H_a$ : The AR Lag 1 is significant

Reject Null if P-value is less than an  $\alpha$  of 0.05

The P-value is almost 0 so we reject null



$H_0$ : The AR Lag 2 is not significant

$H_a$ : The AR Lag 2 is significant

Reject Null if P-value is less than an  $\alpha$  of 0.05

The P-value is almost 0.616 so we do not have evidence to reject null

$H_0$ : The MA Lag 1 is not significant

$H_a$ : The MA Lag 1 is significant

Reject Null if P-value is less than an  $\alpha$  of 0.05

The P-value is almost 0 so we reject null

$H_0$ : The MA Lag 2 is not significant

$H_a$ : The MA Lag 2 is significant

Reject Null if P-value is less than an  $\alpha$  of 0.05

The P-value is almost 0.162 so we do not have evidence to reject null

We can notice that the lags two for both AR and MA seem to not be significant predictors. It could be possible that we could remove one of the lags. However, we can stick with our current model and try different variations after this model.



Figure 2: Results From ARIMA Model-1

The plot displays the upper and lower limits as predicted by the model. We can also observe the mean which seems to be a straight line. Meaning that based on our model the air pollution level would essentially remain the same for the coming year and there will be no decrease or increase in overall air pollution.

To test the strength of our model we can also take a look at the mean squared error. For the number of errors, our model had predicting the values. The mean squared error for our model is 34.03, which is relatively high. However, due to the nature of our data set, this is the best our model can produce. Based on manual analysis and determination of the p,q, and d values. Although, there is an ARIMA model that uses an automatic stepwise function to determine the best model based on the lowest Akaike Information Criterion (AIC).

## Model 2 - Auto ARIMA

Auto ARIMA is a stepwise ARIMA model that picks the best  $p$ ,  $q$ , and  $d$  values based on the lowest AIC score. The AIC score is known as an estimator of the quality of our statistical model for the given data. The AIC estimates the quality of the model, relative to other models for the same data, hence serving as a way to select a model.

```
Best model: ARIMA(1,1,2)(0,0,0)[0]
Total fit time: 6.539 seconds
```

Figure 3: Results From Auto ARIMA

**[Please Refer to the Model-2 Section of the Appendix for Plots and Summary of this model]**

We can see that model suggests that we use  $p = 1$ ,  $d = 1$ , and  $q = 2$  to better fit our model. This is a little off from what we did in our previous model. However, we will go with the model's output and try to validate if the auto ARIMA did provide us with better results. We can also observe that in the previous model we concluded by looking at the summary that we can remove one of the lags. Well, based on the auto ARIMA summary results if we change the AR lags from 2 to 1 we notice that all our lags are now significant and have a p-value close to 0 meaning they are all significant. The residual plots for our data, however, still look the same as they did before. Meaning that the conclusions made in the previous model regarding the data set still hold.



Figure 4: Results From Auto ARIMA Model-2

The model's predicted output looks very similar to what we had in our previous model. Since there is only a minor change ( $p = 1$  instead of 2) it makes sense that we don't notice a substantial increase or decrease in the overall predictions. The mean squared error for our first model was approximately 34.03. However, the mean squared error for our auto ARIMA model is approximately 34.10. There is a slight decrease in the quality of our model, and overall there doesn't seem to be a major difference between the two variations. However, it was interesting to see the results nonetheless.

### Model 3 - Assuming Our data is Stationary

If we go back to the point where we determined the  $d$  value. We noticed that the  $p$ -value given by the ADF test on our data set was:

$p$ -value: 0.09891869282208893

If we set  $\alpha$  of 0.1 We can reject the null and assume that our data is stationary. Also from looking at our PACF and ACF plots one can notice that we don't seem to be able to see strong collinearity among the lags. Which could mean that we are over stationarizing our data. If we run auto-ARIMA again, however, this time we force Auto-ARIMA to assume that our data is stationary. We get the following results:

```
Best model: ARIMA(2,0,2)(0,0,0)[0] intercept
Total fit time: 12.999 seconds
```

Figure 5: Results From Auto ARIMA Assuming Stationary

Based on the auto ARIMA, we need to set our  $p$  and  $q$  values as 2 for both. Setting these values we can run our model and take a look at the model residual diagnostics one again to see if there are any major issues.

**[Please Refer to Model-3 Section Of the Appendix for plots and summary of our Model]**

We can see that our model doesn't seem to have too many glaring issues, the plots for the residuals look almost identical to the plots that were computed for model 1. Hence, we will be making the same assumptions as we did before. We can also take note of the fact that our summary indicates that all our selected lags are significant. Hence, we do not need to make any more changes to our model.

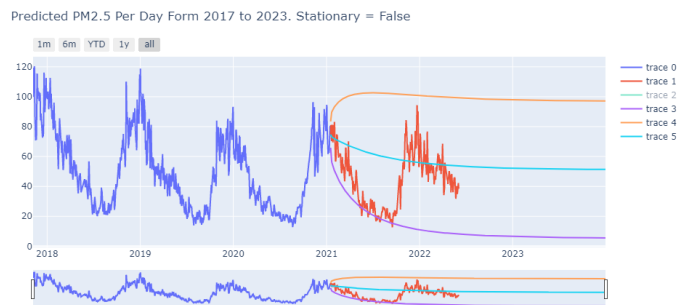


Figure 6: Results From Auto Arima Assuming Stationary

After using our model to make predictions we notice that our mean squared error has gone down substantially to 24.04 from around 34.10 in our previous models. We are not sure why the mean square error has gone down. However, it could be possible that we may be overfitting our model. We can notice from figure 6 that our upper and lower bound predictions have gotten narrower compared to the

previous models. However, despite the shortcomings, it is interesting to note that out of the other two models this is the only model that predicts that the air pollution level would decrease. Although, the decrease in pollution level is rather minor.

## Discussion

Forecasting air quality is a complex task due to environmental dynamics, unpredictability, and changes in pollutant status and time. The serious consequences of air pollution on people, animals, plants, monuments, climate and environment require continuous monitoring and analysis of air quality, especially in developing countries like India. The results that we can take back from this case study are:

- Winters are the most polluted time of the month.
- Winter of 2017 was the most polluted time in the past 5 year.
- Pollution level has been declining with every coming year.
- During the summer and monsoon, pollution level declines.
- Best month is the month of August and September.
- January 2022 had the lowest air pollution level during winters comparatively and gradually declining.
- But as per the plot, the month of June 2022 had the highest air pollution level if compared with June of the last 5 years.
- Two of our prediction model suggests that there will be neither an increase nor a decrease in pollution levels for the next year. However, one of our models suggests a minor decline. Due to varying results, we can not say if we want to make a definitive conclusion. Also, we noticed a trend of residuals from the models potentially being skewed. Which could indicate that perhaps ARIMA isn't the best choice of a prediction model for this particular data set and, in the future we need to employ a better machine-learning model. A Neural network (RNN) or other such models could perhaps be what we are looking for.

Having said that, there are still a lot of factors, which can be a reason for air pollution. There is further scope to this project. We can utilize the results from this data set and combine them with other datasets to try and interpret factors at a much deeper level. Also, we could try and use other factors to predict air pollution rather than the time series itself. However, it was interesting to try and utilize ARIMA and learn about time series data as a whole.

## Appendix

Cumulative air pollution level of India 2017-2022

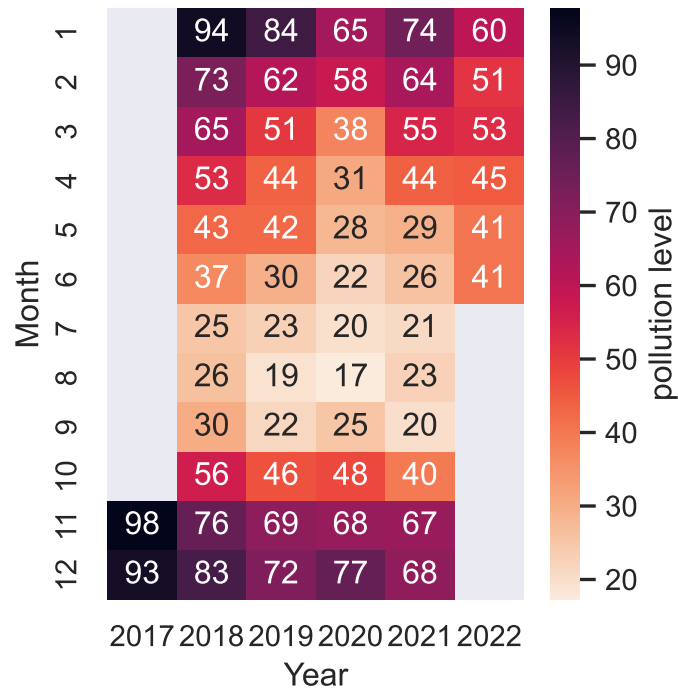


Figure 1: Heat Map

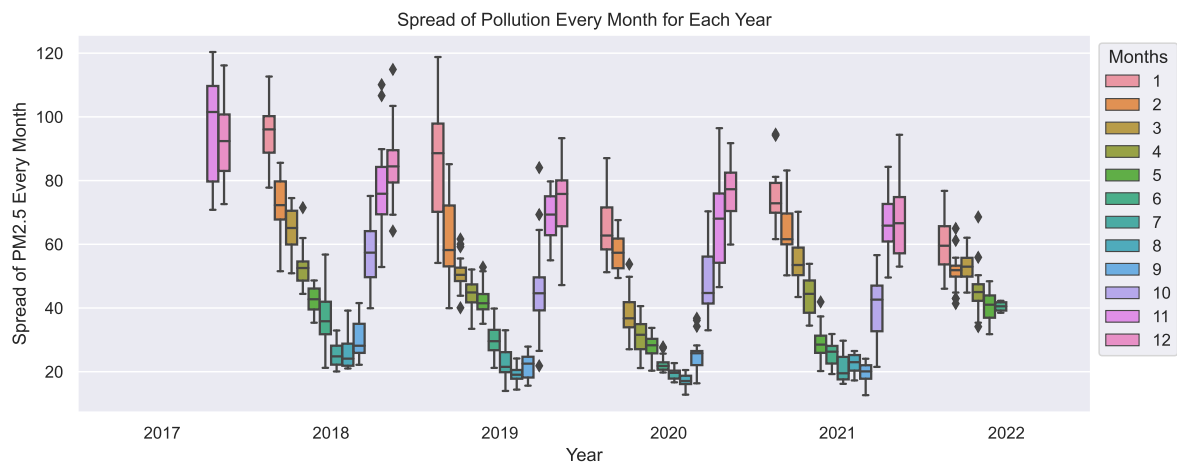


Figure 2: Box Plot

## Liner Regression Model Summary

Table 1: OLS Regression Results

Dep. Variable:	PM2.5	R-squared:	0.063
Model:	OLS	Adj. R-squared:	0.063
Method:	Least Squares	F-statistic:	109.3
Date:	Wed, 21 Dec 2022	Prob (F-statistic):	8.63e-25
Time:	16:35:19	Log-Likelihood:	-7328.5
No. Observations:	1616	AIC:	1.466e+04
Df Residuals:	1614	BIC:	1.467e+04
Df Model:	1		
Covariance Type:	nonrobust		

	coef	std err	t	P> t	[0.025	0.975]
const	59.0993	1.133	52.182	0.000	56.878	61.321
x1	-0.0123	0.001	-10.453	0.000	-0.015	-0.010

Omnibus:	193.489	Durbin-Watson:	0.055
Prob(Omnibus):	0.000	Jarque-Bera (JB):	78.784
Skew:	0.339	Prob(JB):	7.80e-18
Kurtosis:	2.156	Cond. No.	1.94e+03

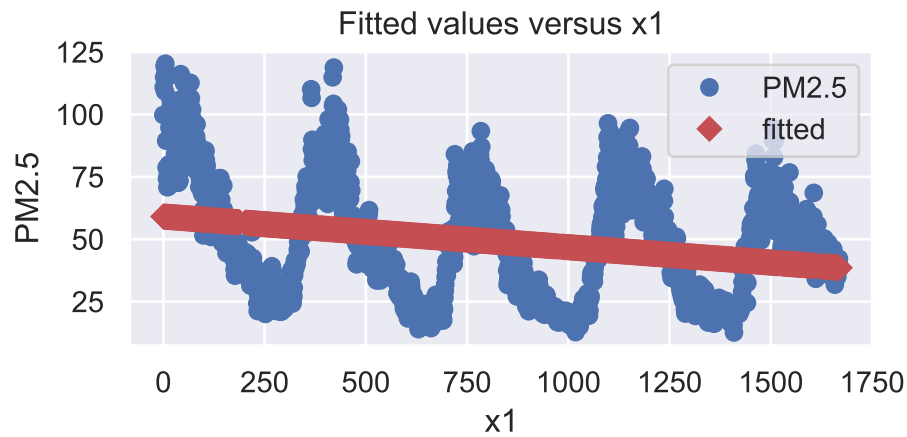


Figure 3: Liner Regression Fitted Line Plot

## ADF Output

### First ADF Test

ADF Statistic: -2.5720220909338787  
n\_lags: 0.09891869282208893  
p-value: 0.09891869282208893  
Critical Values:  
1%, -3.4344410778805936  
Critical Values:  
5%, -2.86334697077965  
Critical Values:  
10%, -2.567731995333179

### Second ADF Test

ADF Statistic: -13.175156123883413  
n\_lags: 1.227599209964447e-24  
p-value: 1.227599209964447e-24  
Critical Values:  
1%, -3.4344436389240722  
Critical Values:  
5%, -2.8633481011816406  
Critical Values:  
10%, -2.567732597265625

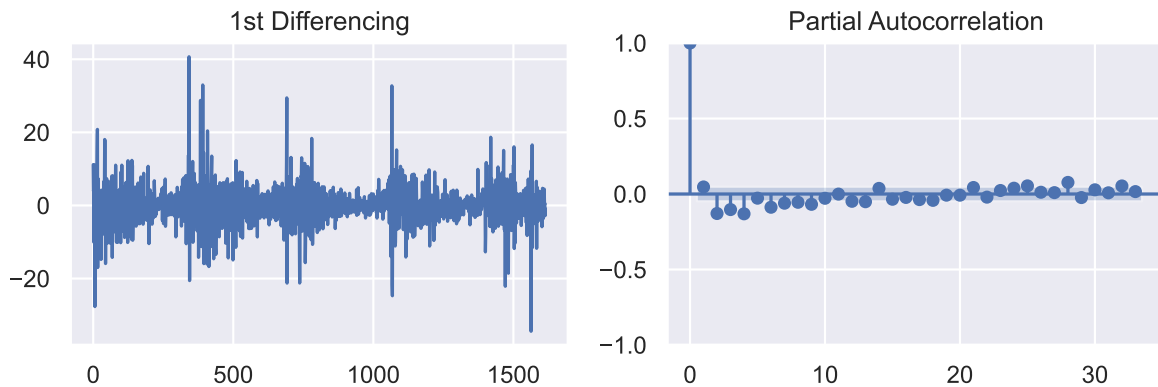


Figure 4: PACF Plot

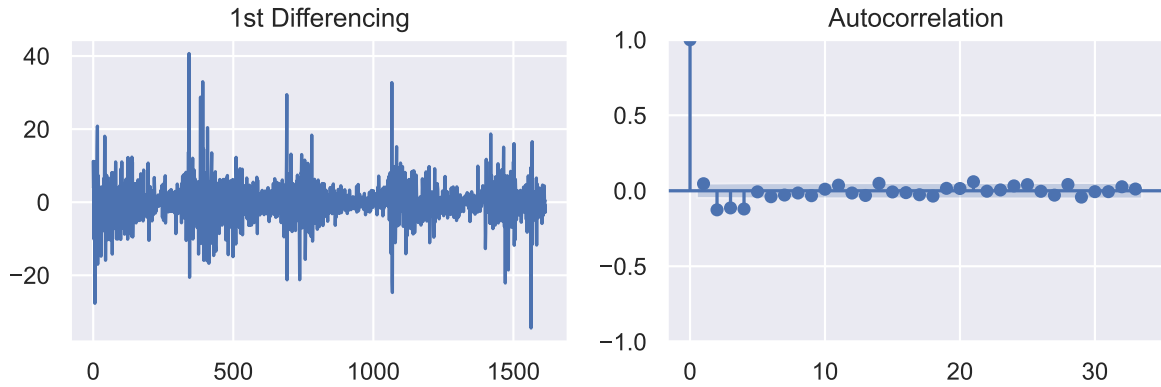


Figure 5: ACF Plot

## Model-1

### SARIMAX Results

```
=====
Dep. Variable:          PM2.5    No. Observations:          1131
Model:                  ARIMA(2, 1, 2)    Log Likelihood          -3499.838
Date:                   Wed, 21 Dec 2022    AIC              7009.676
Time:                   16:35:20    BIC              7034.826
Sample:                 0    HQIC              7019.178
                        - 1131
```

Covariance Type: opg

```
=====
              coef    std err          z      P>|z|      [0.025      0.975]
-----
ar.L1          0.6870     0.116     5.921     0.000     0.460     0.914
ar.L2         -0.0507     0.101    -0.501     0.616    -0.249     0.147
ma.L1         -0.6783     0.118    -5.758     0.000    -0.909    -0.447
ma.L2         -0.1532     0.110    -1.397     0.162    -0.368     0.062
sigma2         28.6825     0.541    53.032     0.000    27.622    29.743
=====
```

```
=====
Ljung-Box (L1) (Q):          0.00    Jarque-Bera (JB):          3201.83
Prob(Q):                    0.98    Prob(JB):              0.00
Heteroskedasticity (H):      0.61    Skew:                  0.82
Prob(H) (two-sided):         0.00    Kurtosis:              11.08
=====
```

### Warnings:

[1] Covariance matrix calculated using the outer product of gradients (complex-step).



RMSE: 34.03192530105746

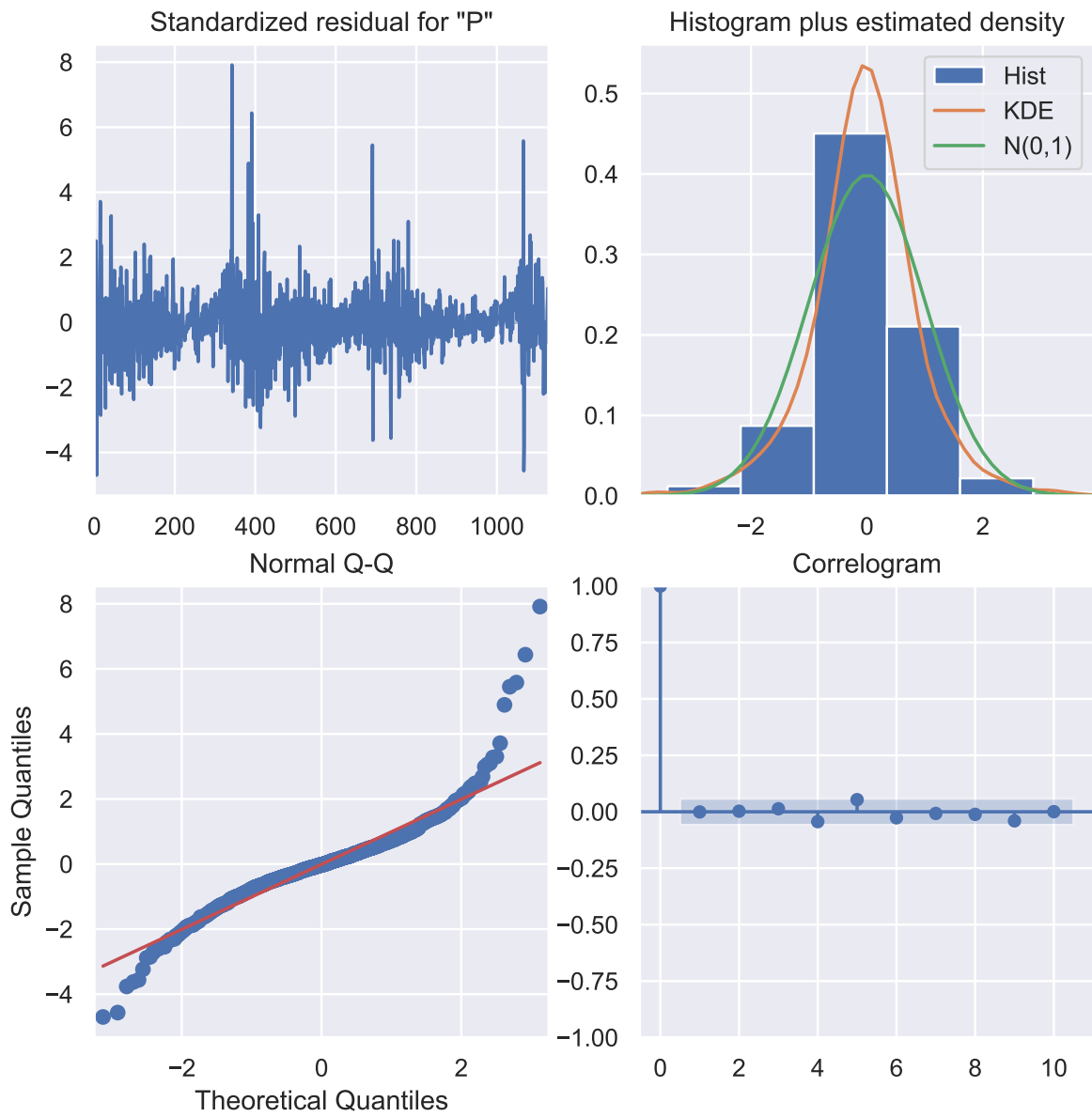


Figure 6: Model Diagnostic Plots - Model 1 (p=2,d=1,q=2)

## Model-2

### SARIMAX Results

```
=====
Dep. Variable:          PM2.5    No. Observations:          1131
Model:                  ARIMA(1, 1, 2)    Log Likelihood          -3499.871
Date:                   Wed, 21 Dec 2022    AIC              7007.741
Time:                   16:35:20    BIC              7027.861
Sample:                 0    HQIC              7015.343
                        - 1131
Covariance Type:        opg
=====
```

	coef	std err	z	P> z	[0.025	0.975]
ar.L1	0.6269	0.042	14.906	0.000	0.544	0.709
ma.L1	-0.6188	0.043	-14.489	0.000	-0.702	-0.535
ma.L2	-0.2087	0.020	-10.488	0.000	-0.248	-0.170
sigma2	28.6841	0.540	53.127	0.000	27.626	29.742

```
=====
Ljung-Box (L1) (Q):          0.00    Jarque-Bera (JB):          3208.26
Prob(Q):                     1.00    Prob(JB):              0.00
Heteroskedasticity (H):      0.61    Skew:                  0.83
Prob(H) (two-sided):         0.00    Kurtosis:              11.09
=====
```

#### Warnings:

[1] Covariance matrix calculated using the outer product of gradients (complex-step).  
RMSE for Auto ARIMA: 34.03192530105746

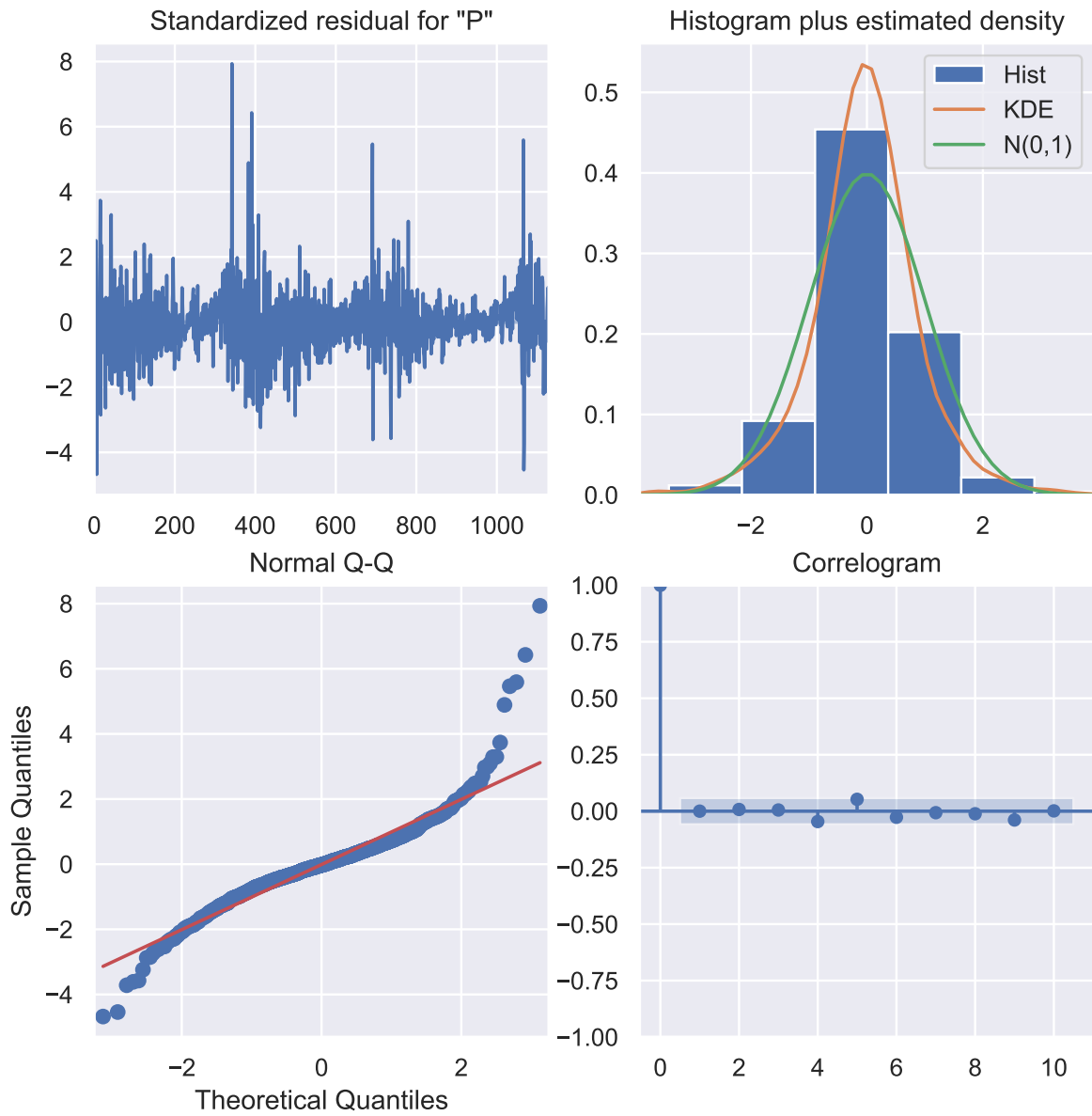


Figure 7: Model Diagnostic Plots - Model 2 ( $p=1, d=1, q=2$ )

### Model-3

#### SARIMAX Results

```

=====
Dep. Variable:          PM2.5    No. Observations:          1131
Model:                ARIMA(2, 0, 2)    Log Likelihood          -3503.621
Date:                Wed, 21 Dec 2022    AIC              7019.242
Time:                16:35:21    BIC              7049.427
Sample:                0    HQIC              7030.646
                        - 1131
Covariance Type:          opg
=====

```

	coef	std err	z	P> z	[0.025	0.975]
const	51.1881	17.755	2.883	0.004	16.390	85.987
ar.L1	1.6191	0.045	36.217	0.000	1.531	1.707
ar.L2	-0.6208	0.044	-14.003	0.000	-0.708	-0.534
ma.L1	-0.6126	0.046	-13.427	0.000	-0.702	-0.523
ma.L2	-0.2076	0.020	-10.364	0.000	-0.247	-0.168
sigma2	28.6346	0.551	51.945	0.000	27.554	29.715

```

=====
Ljung-Box (L1) (Q):          0.00    Jarque-Bera (JB):          3310.22
Prob(Q):                    0.96    Prob(JB):              0.00
Heteroskedasticity (H):      0.61    Skew:                  0.92
Prob(H) (two-sided):         0.00    Kurtosis:              11.18
=====

```

#### Warnings:

[1] Covariance matrix calculated using the outer product of gradients (complex-step).  
 RMSE for Auto ARIMA: 24.040078724006772

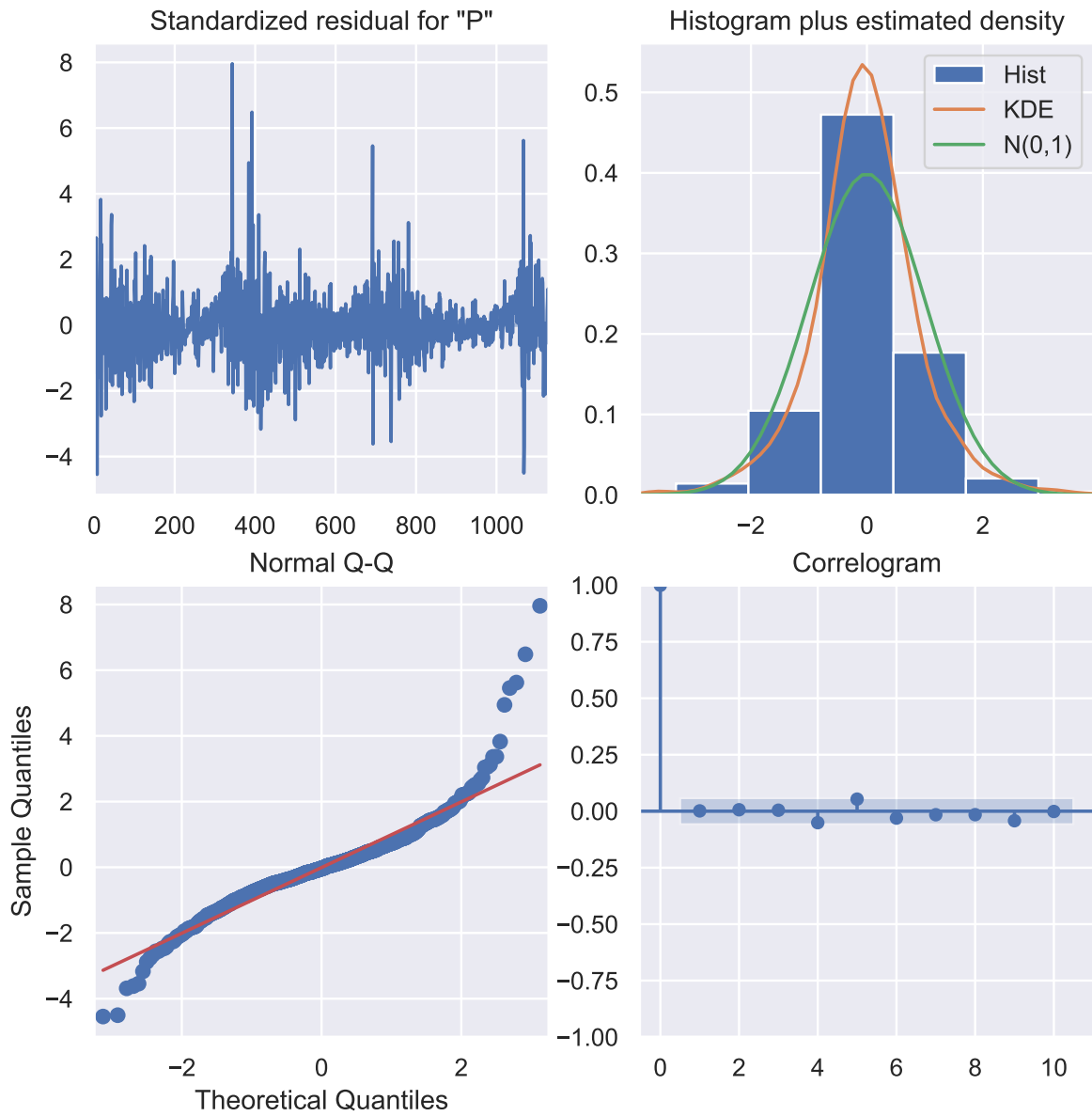


Figure 8: Model Diagnostic Plots - Model 3 ( $p=2, d=0, q=2$ )

## References

- [1] fedesoriano. *Air Quality Data in India (2017 - 2022)*. <https://www.kaggle.com/datasets/fedesoriano/air-quality-data-in-india>. 2022.
- [2] Rob J Hyndman and George Athanasopoulos. *Forecasting: Principles and Practice (3rd ed)*. Accessed on: December 19th 2022. OTexts: Melbourne, Australia. OTexts.com/fpp3., 2010.
- [3] Selva Prabhakaran. “ARIMA Model Complete Guide to Time Series Forecasting in Python”. In: *Machine Learning Plus. com* (2021). Accessed: December 19th 2022.