

# Time Series based Air Pollution Forecasting using SARIMA and Prophet Model

K. Krishna Rani Samal  
Department of CSE  
NIT Rourkela, India  
Mob: +91 9874641997  
517cs6019@nitrkl.ac.in

Santosh Kumar Das  
Department of ECE  
NIT Rourkela, India  
Mob: +91 9437940105  
dassk@nitrkl.ac.in

Korra Sathya Babu  
Department of CSE  
NIT Rourkela, India  
Mob: +91 8249892662  
ksathyababu@nitrkl.ac.in

Abhirup Acharaya  
Department of ECE  
NIT Rourkela, India  
Mob: +91 9861169611  
116ec0248@nitrkl.ac.in

## ABSTRACT

Air pollution severely affects many countries around the world causing serious health effects or death. Increasing dependency on fossil fuels through the last century has been responsible for the degradation in our atmospheric condition. Pollution emitting from various vehicles also cause an immense amount of pollution. Pollutants like RSPM, SO<sub>2</sub>, NO<sub>2</sub>, SPM, etc. are the major contributors to air pollution which can lead to acute and chronic effects on human health. The research focus of this paper is to identify the usefulness of analytics models to build a system that is capable of giving a rough estimate of the future levels of pollution within a considerable confidence interval. Rendered linear regression techniques are found to be insufficient for the time-dependent data. In this regard, we have used time series forecasting approach for predicting the future levels of various pollutants within a considerable confidence interval. The experimental analysis of the forecasting for the air pollution levels of Bhubaneswar City indicates the effectiveness of our proposed method using SARIMA and Prophet model.

## CCS Concepts

• Software and its engineering → Software notations and tools  
→ General programming languages

## Keywords

Pollution; Time Series; SARIMA model; Prophet model; SO<sub>2</sub>; NO<sub>2</sub>; RSPM; SPM

## 1. INTRODUCTION

Human population growth is becoming a tensed issue and this overpopulation has brought many undesirable effects on the environment. Many social and economic factors are having a very bad influence on the environment. The results of this harmful influence lead to the release of hazardous pollutants such as SO<sub>2</sub>, NO<sub>2</sub>, RSPM, SPM, etc. into the atmosphere. Prolonged to such particles can cause several acute as well as chronic health issues. Every day around 93 percent of the world's children suffer from

environmental health risk due to air pollution. Odisha health care data is analyzed to evaluate the health effects of pollutants and concluded that 4.80 percentage of death occurred per 100000 populations in Odisha during 2014 is due to respiratory diseases as shown in Figure 1.

% Death-year-wise						
2010	2009	2011	2012	2013	2014 - % Deaths	
0.10%	0.40%	0.40%	0.20%	0.10%	0.10%	Arunchal Pradesh
0.30%	0.40%	0.50%	0.20%	0.40%	0.30%	Sikkim
	1.00%	0.70%	0.90%	0.90%	2.70%	Chhatisgarh
0.50%	0.20%	2.20%	1.20%	1.00%	0.90%	Manipur
	1.00%	0.40%	0.30%	0.30%	1.30%	Punjab
0.60%	0.50%	1.30%	1.00%	0.70%	1.60%	Mizoram
1.10%	2.40%	2.50%	5.00%	2.20%	4.00%	Rajasthan
2.30%	2.50%	1.90%	0.90%	1.00%	0.70%	Haryana
3.30%	2.60%	2.20%	2.30%	3.20%	3.30%	Uttarakhand
4.40%	5.90%	5.10%	1.10%	1.80%	1.00%	Kerala
5.90%	5.90%	7.90%	5.40%	10.70%	22.70%	Uttar Pradesh
6.50%	2.50%	4.10%	5.60%	5.00%	3.90%	Delhi
6.70%	1.90%	5.40%	3.20%	3.70%	1.90%	Tripura
	3.60%	10.80%	7.10%	6.70%	4.80%	Odisha
	5.40%	6.20%	3.30%	4.70%	8.30%	Himanchal Pradesh
7.20%	6.50%	7.30%	7.90%	7.30%	1.40%	Karnataka
8.50%	8.60%	7.30%	7.50%	6.10%	7.10%	Madhya Pradesh
8.70%	0.30%	0.90%	0.50%	0.10%	0.50%	Tamil Nadu
9.80%	17.60%	9.50%	12.10%	10.60%	5.30%	Andhra Pradesh**
16.10%	23.30%	21.20%	18.20%	21.40%	22.90%	West Bengal

Figure 1. Percentage of death due to respiratory syndrome

The air pollution has become so severe that the public should be timely informed about pollution level and environmental changes so that they can be cautious to keep them safe. Therefore, there are many forecasting models being in use to predict the pollution level in advance, however there is still a requirement of a more

accurate mathematical model to forecast the pollution level and air quality index, which has the negative impact on human health. Air pollution level turns so severe in India [1] so that this has become a leading factor of cancer, heart diseases, cancer, many respiratory infections. Exposure to NO<sub>2</sub>, SPM can affect our lungs and respiratory system, hence the effects of air pollution are alarming [2]. In that case, everyone has to be alerted regarding this danger sign. Due to government economic budget, it is very difficult to deploy efficient sensors at everywhere to measure the pollution levels to alert the public in advance. Therefore, it is better to develop a model which can forecast pollution level in advance without the direct use of any sensors in real time [3].

## 2. RELATED WORKS

Principle component analysis technique is used to forecast the pollutant value in one-day advance [1]. It also works well, when

© 2019 Association for Computing Machinery. ACM acknowledges that this contribution was authored or co-authored by an employee, contractor or affiliate of a national government. As such, the Government retains a nonexclusive, royalty-free right to publish or reproduce this article, or to allow others to do so, for Government purposes only.

ITCC 2019, August 16–18, 2019, Singapore, Singapore

© 2019 Association for Computing Machinery.

ACM ISBN 978-1-4503-7228-2/19/08...\$15.00

DOI: <https://doi.org/10.1145/3355402.3355417>



first checks for stationarity and seasonality, then identifies AR, MA parameter. It follows the differencing process to convert non stationarity data to stationary, which generates the ARIMA (Auto regressive integrated moving average) model.

ARIMA model can be used for both data generation and forecasting [16] [17]. General notation follows for an ARIMA model is  $ARIMA(p; d; q)$ , where  $p$  is the number of auto regressive terms,  $d$  is the order of differencing,  $q$  is the number of moving average terms. This model can be used for non-stationary data. Presence of non-stationary can be checked by various statistical method. Differencing is one of the methods. First, it applies differencing to make the series of data stationary and then apply the ARIMA model. After differencing steps, it finds out AR, MA parameters then it uses a particular model [18]. ARIMA has two different types of models based on the seasonal effects, such as ARIMA and SARIMA model. Seasonal ARIMA can be used to forecast the values during special holidays [19]. The generalized form of the ARIMA model is given in Equation 1,

$$(1 - \sum_{i=1}^p \phi_i L^i)(1 - L)^d X_t = \delta + (1 + \sum_{i=1}^q \theta_i L^i) \varepsilon_t \quad (1)$$

where  $L$  is the lag operator.

$\phi_i$  is the moving average part parameter

$\varepsilon_t$  is the error term.

SARIMA (Seasonal auto regressive integrated moving average) model is similar to the ARIMA model but this model is preferable when the time series exhibits seasonality. Mathematically it can be expressed in terms of a composite model which can be denoted as

$ARIMA(p, d, q)(P, D, Q)S$ . Here, the model parameters  $p, d$  and  $q$  represent the non-seasonal AR order, no seasonal differencing, non-seasonal MA order respectively. Further, the model parameters  $P, D, Q$  and  $S$  are corresponding to the seasonal AR order, seasonal differencing, seasonal MA order, and time span of repeating seasonal pattern respectively. However, the model can be further expressed in simple form without differencing part as mentioned in Equation 2 [16], [20].

$$\phi_p(B)\Phi_p(B^S)(1-B)^d(1-B^S)^D Y_t = \theta_q(B)\Theta_q(B)^S \varepsilon_t \quad (2)$$

where,

$$\begin{aligned} \phi_p(B) &= 1 - \phi_1 B - \phi_2 B^2 - \phi_3 B^3 - \dots - \phi_p B^p \\ \Phi_p(B) &= 1 - \phi_1 B^S - \phi_2 B^{2S} - \dots - \phi_p B^{pS} \\ \theta_q(B) &= 1 - \theta_1 B - \theta_2 B^2 - \dots - \theta_q B^q \\ \Theta_q(B) &= 1 - \theta_1 B^S - \theta_2 B^{2S} - \dots - \theta_q B^{qS} \end{aligned}$$

**Prophet forecasting Model:** This model is developed by Facebook, available in python and R [21] [22]. Due to its three main features, ie. trend, seasonality, holidays [23] and

demand for the high quality of forecasting are the main reason for building this model. It can be represented as in Equation 3,

$$y(t) = g(t) + s(t) + h(t) + \varepsilon_t \quad (3)$$

where the model parameters  $g(t), s(t), h(t), \varepsilon_t$  are piecewise linear curve for modeling non-periodic changes in time series, periodic changes, the effects of holidays with irregular schedules, error term accounts for any unusual changes not accommodated by the model respectively. To fit the proposed model with seasonality effects and forecast based on it, it uses a Fourier series which provides a flexible model. Seasonal effects  $s(t)$  can be represented as in Equation 4 [24],

$$s(t) = \sum_{n=1}^N a_n \cos\left(\frac{2\pi nt}{p}\right) + b_n \sin\left(\frac{2\pi nt}{p}\right) \quad (4)$$

where,  $p$  represents regular period

## 5. MEASURE OF ACCURACY

RMSE and MSE are the two criteria chosen to measure the performance of time series forecasting model as shown in Equation 5 and Equation 6 where error is  $e_i, i = 0, 1, 2, 3, \dots, n$ . The model which is having the least value of RMSE and MSE is selected as the best pollution forecasting model.

$$RMSE : \sqrt{\frac{1}{n} \sum_{i=1}^n e_i^2} \quad (5)$$

$$MSE : \text{Mean square error} : \frac{1}{n} \sum_{i=1}^n e_i^2 \quad (6)$$

## 6. RESULTS AND DISCUSSIONS

In this part, we present the time series pollution forecasting using historical pollution dataset of Bhubaneswar city, India. The analyzed data included year wise data from 2005 to 2015. The dataset contains pollutant values i.e. SO<sub>2</sub>, NO<sub>2</sub>, SPM, RSPM values with GPS coordinate for various pollution monitoring stations of Odisha, India. The data is processed according to the requirement of the forecasting model and missing values are tackled using mode void fill method, using backward and forward fill wherever deemed fit and necessary.

The time series plots illustrate that there is roughly a constant level of certain pollutants. In addition to that, there is also a constant level of seasonal fluctuation and random fluctuations over time. Differencing process is done to handle this type of situation before forecasting model development. The weekly, monthly and yearly seasonality checked. The required parameters also fed to the model for a more precise forecast. The Indian holiday list feed is expected to significantly boost the model performance as it would add an extra parameter for better correlation amongst the dates and the pollutant levels correspondingly. There are three main steps i.e. stationary test, model identification, and forecasting in building the forecasting model.

### 6.1 Stationary Test

The present paper implemented Dickey-Fuller test to check the stationarity of the data before the SARIMA and prophet model implementation. Results of dickey fuller test that conducted year wise for each pollutant are shown in Figure 3 - Figure 6. In summary, we concluded from this test that the test statistics is less than the p-value for each time series pollutant value which implies that series are not stationary. Log transformation is also used to stabilize the non-constant variance of time series.

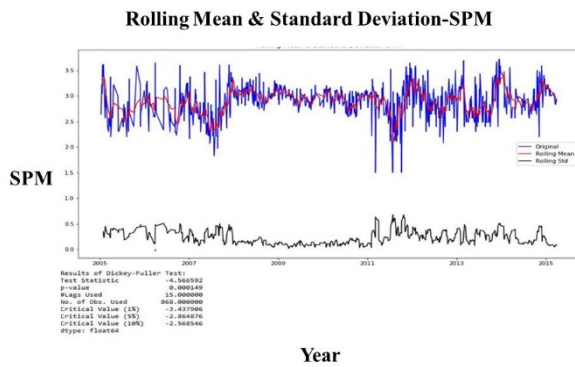


Figure 3. Dickey-Fuller test for SPM

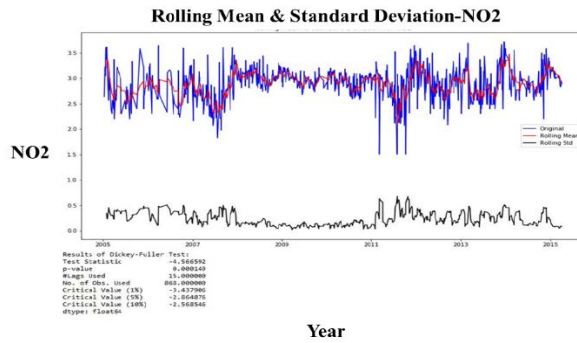


Figure 4. Dickey-Fuller test for NO2

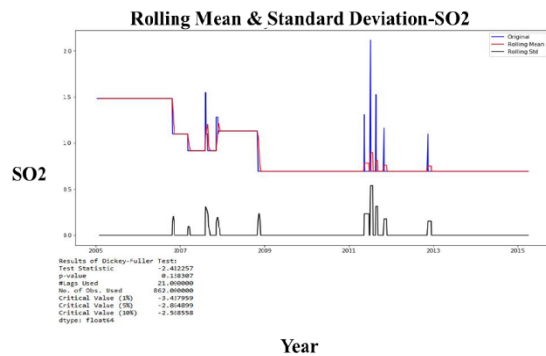


Figure 5. Dickey-Fuller test for SO2

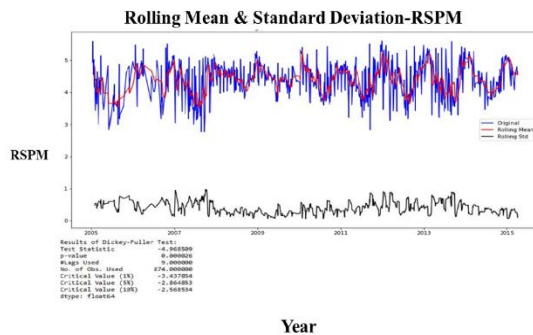


Figure 6. Dickey-Fuller test for RSPM

## 6.2 Model Identification

Akaike Information Criteria (AIC) and Bayesian information criterion(BIC) are two statistical measure of SARIMA model [25]. The model with lower AIC and BIC value is better while comparing two models. Hence, the combination of these measure is used to identify the best order of the SARIMA model for pollution

forecasting. Table 1 shows the lower value of these measure which considered to select the best order of the SARIMA model for each pollutant.

Table 1: AIC and BIC values to find the best order of SARIMA model

Pollutant	SPM	NO2	SO2	RSPM
Order	(0,1,2)	(0, 1, 2)	(0, 1, 2)	(1, 1, 1)
Seasonal order	(1, 1, 1,12)	(1, 1, 1,12)	(1, 1, 1,12)	(1, 1, 1, 12)
AIC	20.157	20.157	-1244.481	552.886
BIC	45.691	45.691	-1233.532	578.421

## 6.3 Forecasting using SARIMA and Prophet Model

Log transformation is used in this paper while developing a forecasting model to convert nonstationary time series into a stationary time series to achieve better performance. The actual results and forecasting results of SARIMA model are shown in Figure 7 - Figure 10. The actual results and predicted results of time series Prophet on log model are shown in Figure 11 - Figure 14.

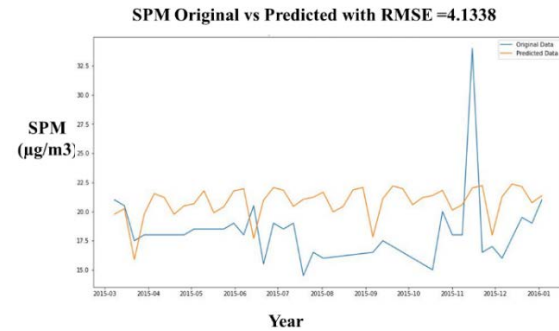


Figure 7. SARIMA model-SPM

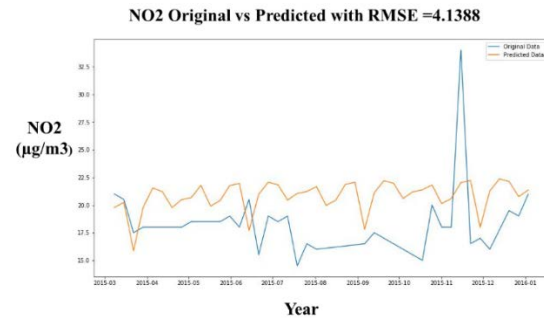


Figure 8. SARIMA model-NO2

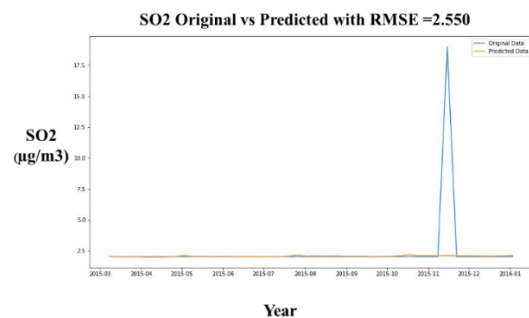


Figure 9. SARIMA model-SO2



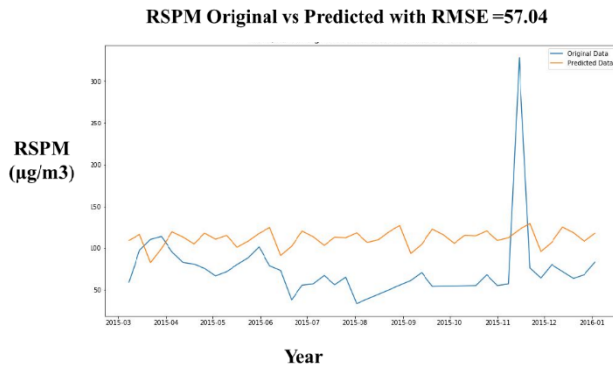


Figure 10. SARIMA model-RSPM

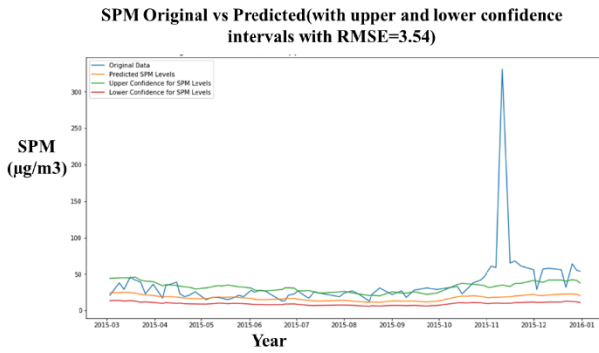


Figure 11. Prophet on log model-SPM

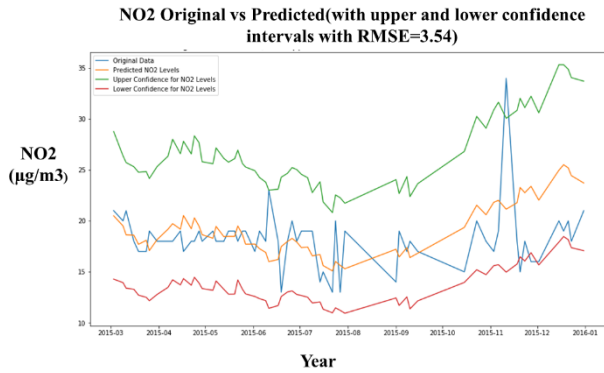


Figure 12. Prophet on log model-NO2

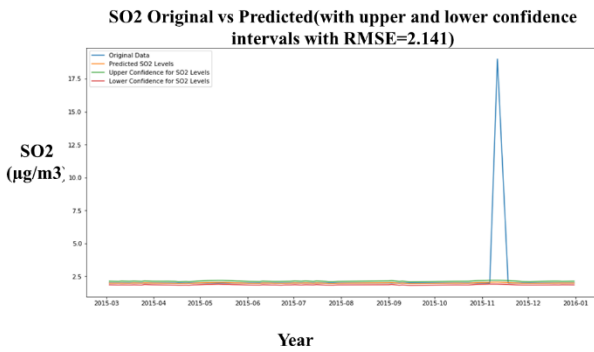


Figure 13: Prophet on log model-SO2

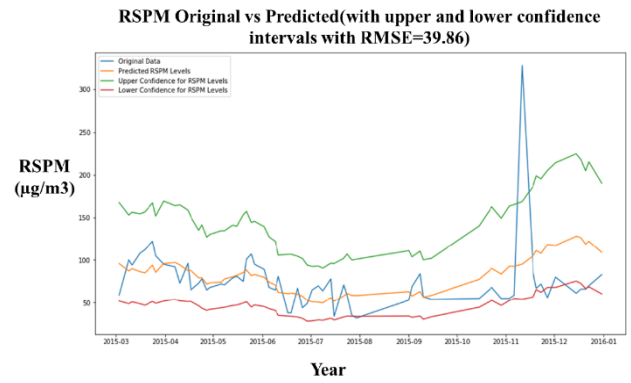


Figure 14. Prophet on log model-RSPM

## 7. COMPARATIVE ANALYSIS

The predictions of air pollution level are location dependent. Further, there is no previous reported work of prediction of air pollution levels corresponding to our research location, i.e. Bhubaneswar. However, similar types of techniques are applied to one of the cities of Bulgaria [11]. Hence, comparison is made between the general and logarithmic models to identify the successful model for forecasting. The experimental results show that logarithmic model works better at adapting to the long term trends and anomalies in the pollutant levels. The results show that the log-model gives comparatively slightly better values as compared to the general model on a larger scale. Comparison of performance metrics for SARIMA Model and Prophet Model are shown in Table 2-Table 4.

Table 2- Accuracy Metrics for SARIMA general model

Metric/Pollutant	SPM	NO2	SO2	RSPM
RMSE	4.13	4.13	2.55	57.04
MSE	17.12	17.12	6.50	3254.158

Table 3. Accuracy Metrics for Prophet general model

Metric/Pollutant	SPM	NO2	SO2	RSPM
RMSE	3.78	3.57	2.15	45.80
MSE	14.34	12.75	4.65	2097.77

Table 4. Accuracy Metrics for Prophet log model

Metric/Pollutant	SPM	NO2	SO2	RSPM
RMSE	3.54	3.54	2.141	39.86
MSE	12.55	12.55	4.58	1589.259

We have also compared the performance of Prophet general model and Prophet log model using their performance metric and concluded that prophet log model provides more accurate forecasting results than Prophet general model.

## 8. CONCLUSION

In this work, we have proposed two approaches for pollution forecasting based on the historical data which contains information from 2005 to 2015. The proposed model predicted pollutant value for 2016. We made a comparison between the model's performance metrics. By looking into the accuracy metric values in Table 2-Table 4, we conclude that both the SARIMA and prophet model provides a good quality of accuracy. However, the best approach is the prophet model on log transformation which has the least minimum RMSE, MSE value. The results show the feasibility of using time series forecasting model, i.e. Prophet model to forecast the future level of pollution and build an early warning system for public safety. This work can be extended by analyzing health care

data to establish health correlation with the pollution level in the future. However, due to the lack of recent data availability, we have restricted our research to the year 2015. Further, the proposed method can be enhanced using a deep learning algorithm to achieve a much higher degree of freedom, versatility, adaptability, and accuracy.

## 9. REFERENCES

- [1] A. Kumar and P. Goyal, "Forecasting of air quality in delhi using principal component regression technique," *Atmospheric Pollution Research*, vol. 2, no. 4, pp. 436–444, 2011.
- [2] J. S. Pandey, R. Kumar, and S. Devotta, "Health risks of no<sub>2</sub>, spm and so<sub>2</sub> in delhi (india)," *Atmospheric Environment*, vol. 39, no. 36, pp. 6868–6874, 2005.
- [3] L. Chen, J. Xu, L. Zhang, and Y. Xue, "Big data analytic based personalized air quality health advisory model," in *Proc. 13th IEEE Conf. on Automation Science and Engineering (CASE)*, 2017, pp. 88–93.
- [4] C. Zhang, J. Yan, Y. Li, F. Sun, J. Yan, D. Zhang, X. Rui, and R. Bie, "Early air pollution forecasting as a service: An ensemble learning approach," in *Proc. IEEE Int. Conf. on Web Services (ICWS)*, 2017, pp. 636–643.
- [5] S. Taneja, N. Sharma, K. Oberoi, and Y. Navoria, "Predicting trends in air pollution in delhi using data mining," in *Proc. IEEE 1st India Int. Conference on Information Processing (IICIP)*, 2016, pp. 1–6.
- [6] Z. Wang and Z. Long, "Pm<sub>2.5</sub> prediction based on neural network," in *Proc. IEEE 11th Int. Conf. on Intelligent Computation Technology and Automation (ICICTA)*, 2018, pp. 44–47.
- [7] K. B. Shaban, A. Kadri, and E. Rezk, "Urban air pollution monitoring system with forecasting models," *IEEE Sensors Journal*, vol. 16, no. 8, pp. 2598–2606, 2016.
- [8] Y. Zhou, S. De, G. Ewa, C. Perera, and K. Moessner, "Data-driven air quality characterization for urban environments: A case study," *IEEE Access*, vol. 6, pp. 77 996–78 006, 2018.
- [9] B. Yeganeh, M. S. P. Motlagh, Y. Rashidi, and H. Kamalan, "Prediction of co concentrations based on a hybrid partial least square and support vector machine model," *Atmospheric Environment*, vol. 55, pp. 357–365, 2012.
- [10] N.-U. Lee, J.-S. Shim, Y.-W. Ju, and S.-C. Park, "Design and implementation of the sarima-svm time series analysis algorithm for the improvement of atmospheric environment forecast accuracy," *Soft Computing*, vol. 22, no. 13, pp. 4275–4281, 2018.
- [11] D. Voynikova, S. Gocheva-Ilieva, A. Ivanov, and I. Iliev, "Studying the effect of meteorological factors on the so<sub>2</sub> and pm<sub>10</sub> pollution levels with refined versions of the sarima model," in *AIP Conference Proceedings*, vol. 1684, no. 1. AIP Publishing, 2015, p. 100005.
- [12] M. Oprea, S. F. Mihalache, and M. Popescu, "A comparative study of computational intelligence techniques applied to pm<sub>2.5</sub> air pollution forecasting," in *2016 6th International Conference on Computers Communications and Control (ICCCC)*. IEEE, 2016, pp. 103–108.
- [13] N. H. A. Rahman, M. H. Lee, and M. T. L. Suhartono, "Evaluation performance of time series approach for forecasting air pollution index in johor, malaysia," *Sains Malaysiana*, vol. 45, no. 11, pp. 1625–1633, 2016.
- [14] S. Jain and V. Mandowara, "Study on particulate matter pollution in jaipur city," *International Journal of Applied Engineering Research*, vol. 14, no. 3, pp. 637–645, 2019.
- [15] OpenGovernmentDataPlatformIndia. (2017, Oct 16) Ambient air quality data of odisha. [Online]. Available: <https://data.gov.in/catalog/ambientair-quality-data-odisha>
- [16] W. Wang and Y. Guo, "Air pollution pm<sub>2.5</sub> data analysis in los angeles long beach with seasonal arima model," in *Proc. IEEE Int. Conf. on Energy and Environment Technology*, vol. 3, 2009, pp. 7–10.
- [17] G. E. Kulkarni, A. A. Muley, N. K. Deshmukh, and P. U. Bhalchandra, "Autoregressive integrated moving average time series model for forecasting air pollution in nanded city, maharashtra, india," *Modeling Earth Systems and Environment*, vol. 4, no. 4, pp. 1435–1444, 2018.
- [18] wikipedia. (2019, Apr 17) Autoregressive integrated moving average. [Online]. Available: [https://en.wikipedia.org/wiki/Autoregressive\\_integrated\\_moving\\_average](https://en.wikipedia.org/wiki/Autoregressive_integrated_moving_average)
- [19] I. Yenidogan, A. C. Ayir, O. Kozan, T. Dag, and C. Arslan, "Bitcoin forecasting using arima and prophet," in *Proc. IEEE 3rd Int. Conf. on Computer Science and Engineering (UBMK)*, 2018, pp. 621–624.
- [20] M. H. Lee, N. H. A. Rahman, M. T. Latif, M. E. Nor, N. A. B. Kamisan et al., "Seasonal arima for forecasting air pollution index: A case study," *American Journal of Applied Sciences*, vol. 9, no. 4, pp. 570–578, 2012.
- [21] Facebook. (2019, May 15) Automatic forecasting procedure. [Online]. Available: <https://pypi.org/project/fbprophet/>
- [22] S. Taylor. (2019, May 14) prophet: Automatic forecasting procedure. [Online]. Available: <https://cran.r-project.org/web/packages/prophet/>
- [23] FacebookResearch. (2017, Feb 23) Prophet: forecasting at scale. [Online]. Available: <https://research.fb.com/prophet-forecasting-at-scale/>
- [24] G. Borowik, Z. M. Wawrzyniak, and P. Cichosz, "Time series analysis for crime forecasting," in *Proc. IEEE 26th International Conference on Systems Engineering (ICSEng)*, 2018, pp. 1–10.
- [25] J. R. Reddy, T. Ganesh, M. Venkateswaran, and P. Reddy, "Forecasting of monthly mean rainfall in coastal andhra," *International Journal of Statistics and Applications*, vol. 7, no. 4, pp. 197–204, 2017.