

A PM 2.5 Forecasting Model Based on Air Pollution and Meteorological Conditions in Neighboring Areas

Muhammad Adrezo^{1,2}, Yo-Ping Huang¹ and Frode Eika Sandnes³[0000-0001-7781-748X]

¹Department of Electrical Engineering, National Taipei University of Technology
Taipei, Taiwan 10608

²Department of Computer Science, UPN Veteran Jakarta, Jakarta, 12450 Indonesia

³Faculty of Technology, Art and Design, Oslo Metropolitan University, Oslo, Norway
muhammad.adrezo@gmail.com; yphuang@ntut.edu.tw; frodes@oslomet.no

Abstract. Air pollution has received much attention in recent years, especially in the most densely populated areas. Sources of air pollution include factory emissions, vehicle emissions, building sites, wildfires, wood-burning devices, and coal power plants. Common and dangerous air pollutants include nitrogen dioxide (NO₂), ozone (O₃), carbon dioxide (CO₂), particulate matter 10 (PM 10) and particulate matter 2.5 (PM 2.5). This study focused on PM 2.5 because it has an aerodynamic diameter less than or equal to 2.5 μm . The small size of this pollutant makes it easily inhaled by humans and may end up deep in the lungs or even the bloodstream. Such pollutants can trigger health problems such as asthma, respiratory inflammation, reduced lung function and lung cancer. The purpose of this work was to forecast the next hour of PM 2.5 based on air pollution concentrations and meteorological conditions. The approach also uses station location data to cluster the area and to determine the neighboring areas of each station. Forecasting is based on the Long Short-Term Memory (LSTM). The result shows that the proposed approach can effectively forecast the next hour of PM 2.5 pollution.

Keywords: Air Pollution; PM 2.5; Forecasting; Long Short-Term Memory (LSTM).

1 Introduction

Air pollution has received much attention in recent years, especially in densely populated areas. The American Lung Association [1] estimated that nearly 134 million people in the US, that is, over 40 % of the population, are at risk of disease and premature death because of air pollution. Bad outdoor air quality caused an estimated 4.2 million premature deaths in 2016. According to the World Health Organization [2] about 90 percent of premature deaths due to poor air quality occurred in low GDP per capita countries. Indoor smoke is an ongoing health threat to the 3 billion people who cook and heat their homes by burning biomass, kerosene, and coal. Examples of common

pollutants include soot, smoke, mold, pollen, nitrogen dioxide (NO₂), ozone (O₃), carbon dioxide (CO₂), particulate matter 10 (PM 10), and particulate matter 2.5 (PM 2.5). High concentrations of such substances may cause health problem to people in the affected area. Researchers have unearthed many health effects which are believed to be associated with exposure to air pollution. Effects caused by air pollution include respiratory diseases (including asthma and reduced lung function), cardiovascular diseases, cancers, and adverse pregnancy outcomes (such as preterm birth).

Throughout history, there has been many tragedies caused by air pollution resulting in diseases and deaths. Some of the worst tragedies caused by air pollution during the 19th century includes The Donora Smog of 1948 (Pennsylvania), The Great Smog of 1952 (London), The 1983 Melbourne dust storm and The 1997 Southeast Asian haze. The Donora Smog affected almost half of the population of Donora, killed 20 people and caused respiratory problems for 6000 people [3]. The Great Smog of London caused reduced visibility and even penetrated indoor areas. At least 4000 people were killed, and many more become ill [4]. Causes of air pollution includes factory emissions, vehicle emissions, building construction, wildfires, wood-burning devices, and coal-fired power plants.

One of the most dangerous air pollutants includes particulate matter 2.5 (PM 2.5). PM 2.5 is one of the primary indicators of air pollution because it affects more people than any other pollutant. PM 2.5 has an aerodynamic diameter of 2.5 μm or less. Common components of PM include sulfate, nitrates, sodium chloride, ammonia, mineral dust, black carbon, and water. High concentrations of PM is related to human health as it can easily be inhaled by humans and thereby affect the respiratory system and the cardiovascular system, and even damage the blood and nervous system and ultimately may cause death [7].

As PM 2.5 cause several diseases the Environment Protection Agencies (EPA) of several countries around the world are monitoring and forecasting PM 2.5. Prediction is important to issue early pollution warnings, for decision making and pollution control, thereby improving the life quality of the population. Traditional techniques and artificial intelligence (AI) techniques have both been applied to forecast PM 2.5. Mathematical and statistical techniques are used for traditional PM 2.5 forecasting in which a physical model was designed, and then data was calculated using mathematical differential equations. However, the traditional techniques have several shortcomings such as difficulties of processing large data sizes, long computation time and limited accuracy and inability to predict extreme points. However, with the advancement in technology, many researchers moved from mathematical and statistical techniques to computational techniques and AI techniques. AI techniques can overcome some of the challenges faced by the traditional techniques. Specific approaches include artificial neural network (ANN), machine learning, and deep learning.

This study focused on forecasting PM 2.5 concentrations using AI techniques based on air pollutants (NO₂, CO, and O₃) and meteorological conditions (wind speed, wind direction, and rain) in each area and its neighboring areas. Air pollutants and meteorological conditions such as wind from neighboring areas are needed because PM 2.5 is a tiny particle that is easily carried by the wind from one area to other areas. This study used two datasets from the EPA of Taiwan. The first dataset is a station location

dataset consisting of longitude and latitude of each station. The second dataset is an air pollution dataset consisting of air pollutant measurements and meteorological observations.

2 Related work

Researchers and EPAs around the world have focused on air pollution, especially PM 2.5 concentrations. The research attention is in line with the public concern about the dangers of air pollution, especially PM 2.5 concentrations. PM 2.5 has received more attention than other air pollutants as it can easily be inhaled and cause many health problems. Researchers and the governments of many countries have deployed many systems to forecast PM 2.5 concentrations to be able to issue early PM 2.5 concentration warnings. Several engines and system architectures have been proposed for forecasting air pollution, especially PM 2.5 concentrations.

Ganesh et al. [5] focused on forecasting air quality index, not air pollutant concentrations such as PM 2.5, PM 10, CO, and O₃. They presented different regression models such as Support Vector Regression (SVR) and linear models such as multiple linear regression consisting of stochastic gradient descent, mini-batch gradient descent, and gradient descent to forecast air quality index based on air pollution index data. Shaban and Rezk [6] collected air quality data wirelessly from monitoring nodes that were equipped with an array of gaseous and meteorological sensors. They focused on the monitoring system and the forecasting module. They investigated three machine learning models, namely Support Vector Machine (SVM), M5P model tree, and ANN.

Gu, Qiao and Lin [7] proposed a heuristic recurrent air quality predictor to infer air quality using SVR. The authors forecasted air pollution using air pollutant concentrations and the meteorological conditions in the local area. The authors compared the proposed approach to three popular predictors (Voukantsis, Vlachogianni, and Kaboodvandpour). Meteorological condition and air pollutant data were also used by Tsai, Zheng and Cheng [8]. Recurrent Neural Network (RNN) with Long Short-Term Memory (LSTM) was proposed by the authors as an approach to forecast PM 2.5 concentrations. Oprea, Mihalache and Popescu [9] tried to compare two computational intelligence techniques, namely ANNs and adaptive neuro-fuzzy inference system (ANFIS) to forecast PM 2.5 concentrations based on air pollutant concentrations and meteorological conditions.

Utilizing the data of surface meteorological observation and air pollution PM 2.5 concentrations in Wuhan City was conducted by Chen, Qin and Zhou [10]. They used multiple regression analysis and back-propagation neural network to develop an air pollution PM 2.5 index forecasting model. According to previous research we know that air pollutant concentrations correlates with meteorological conditions. The wind can carry small particles including PM 2.5 that has a size of 2.5 μm from one area to another. It can also affect PM 2.5 concentrations in other areas. Therefore, the conditions in neighboring areas is also important as an indicator for forecasting PM 2.5 concentrations.

3 Proposed approach

The proposed approach involves two primary processes to forecast PM 2.5 concentrations. The first process involves selecting neighboring areas of each station. The second process involves making a PM 2.5 forecasting model. The system framework for the PM 2.5 forecasting is shown in Fig. 1.

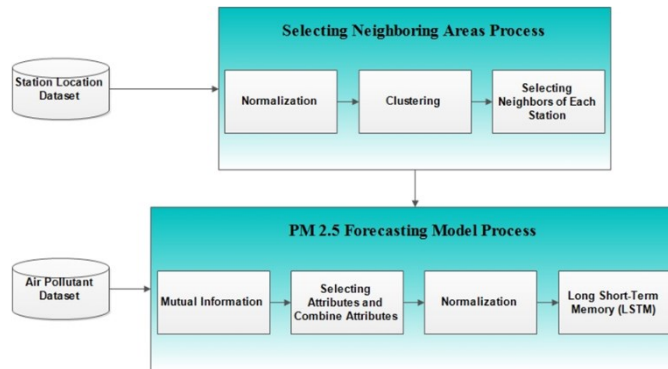


Fig. 1. System framework for PM 2.5 forecasting.

3.1 Data

Our study relied on two datasets (station location dataset and air pollution dataset) from 67 stations in Taiwan that were collected by the EPA in Taiwan. The station location dataset was used for selecting neighboring areas. The air pollution dataset was used for building the PM 2.5 forecasting model. The air pollution dataset contains hourly data in 6 years from 2012 to 2017. The data were divided into training data, validation data and testing data with ratios of 4: 1: 1, respectively. Training and validation data were used to build the PM 2.5 forecasting model. Testing data were used to measure the quality of the PM 2.5 forecasting model. To forecast the next hour of PM 2.5 concentrations, this study used three-hour window of observations.

3.2 Selecting Neighboring Areas

This phase involves determining the neighboring areas of each station. The flowchart of this process is shown in Fig. 2. First, the data is normalized using min-max normalization. Next, the data is clustered using the X-means method. Finally, the neighboring areas of each station is selected using the radius of each station using the Euclidean distance. However, using clustering results from only the locations which have a closer distance are not selected as neighboring areas if such stations lie on the edge of a cluster. To overcome this issue, we determined neighboring areas based on the radius of each station as shows in Fig. 3. The radii were determined as follows:

1. Based on the clustering result, the distance is calculated using Euclidean distance (1) between the station and the center cluster of each cluster:

$$EU(a,b) = \sqrt{\sum_{i=1}^n (x_i - y_i)^2} \quad (1)$$

$EU(a,b)$ is the distance between vector a and vector b , where vector a is the location station coordinate and vector b is the cluster center coordinate.

2. The mean distance of all clusters is calculated based on the result from step 1.

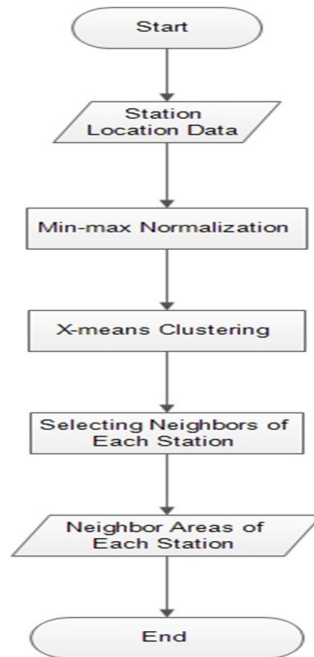


Fig. 2. Selecting neighboring areas.

3. To determine neighboring areas of the station the cluster of the station is checked.
4. The mean distance of the cluster is set as a radius of the station.
5. The Euclidean distance between the station and other stations is calculated. Then,

$$\begin{cases} dist \leq radius ; True \\ dist > radius ; False \end{cases} \quad (2)$$

If True, then the station is selected as the neighboring area. If False, the station is not selected as the neighboring area.

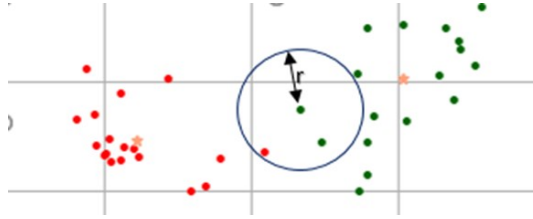


Fig. 3. Neighboring areas of the station based on radius.

6. Steps 3 and 4 are repeated until neighboring areas of all stations are determined.

3.3 PM 2.5 Forecasting Model

A PM 2.5 forecasting model is applied to each station. Fig. 4 illustrates the process. First, mutual information of air pollutant dataset of each station is calculated. After that, the attributes of each station are selected based on the mutual information result by setting a threshold. The mutual information result is used for determining the attributes for the next step by setting a threshold. The threshold is based on the mean mutual information score of the attributes.

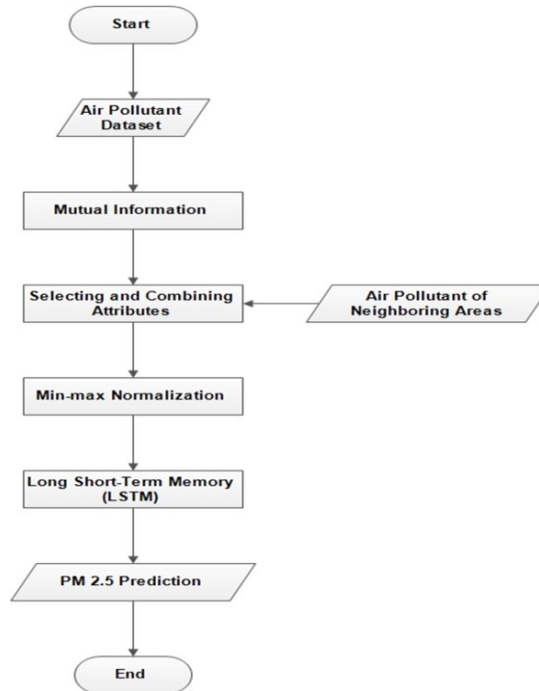


Fig. 4. The PM 2.5 forecasting model.

The threshold is determined using equation (3)

$$T = \frac{1}{2} avg_MI \quad (3)$$

where avg_MI is the mutual information score mean of the attributes. The attributes that have a mutual information score equal to or above the threshold are selected. After that, neighboring areas determined in the previous steps are used. The selected attributes of the station combined with PM 2.5, PM 10, wind speed and the wind direction in neighboring areas. After combining the data, the data is normalized. Finally, the normalized data is used as input to the Long Short-Term Memory (LSTM) for making PM 2.5 forecasts. Our approach uses a one-to-one LSTM architecture. This architecture consists of one input, one LSTM layer, and one output. Inside the LSTM layer we used 64 LSTM units and an Adam optimizer. Root mean square error was used as a loss function during training. To measure the quality of the PM 2.5 forecasting model, RMSE and MAE were used. The proposed LSTM structure is shown in Fig. 5.

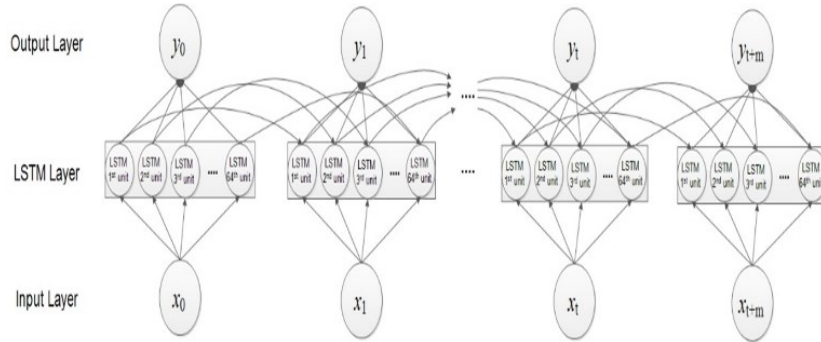


Fig. 5. Proposed LSTM structure.

4 Experimental evaluation and results

Fig. 6 shows the results of experiment with different learning rates. Figs. 6 (a-f) show learning rates of 0.005, 0.001, 0.0005, 0.0001, 0.00005 and 0.00001, respectively. These experiments used 150 iterations and 32 LSTM units. The air pollution dataset from the Annan measurement station was used for these experiments. These experiments aimed to determine the effect when the learning rate decreases.

Based on the results of the experiments with different learning rates, learning rates of 0.005, 0.001, and 0.0005 showed fluctuations in training loss and validation loss. On the other hand, learning rates of 0.0001, 0.00005, and 0.00001 exhibited stable training loss and validation loss even though the error is higher with a learning rate of 0.00001. This happens because, with a small learning rate, more iterations are needed to reach the optimal model.

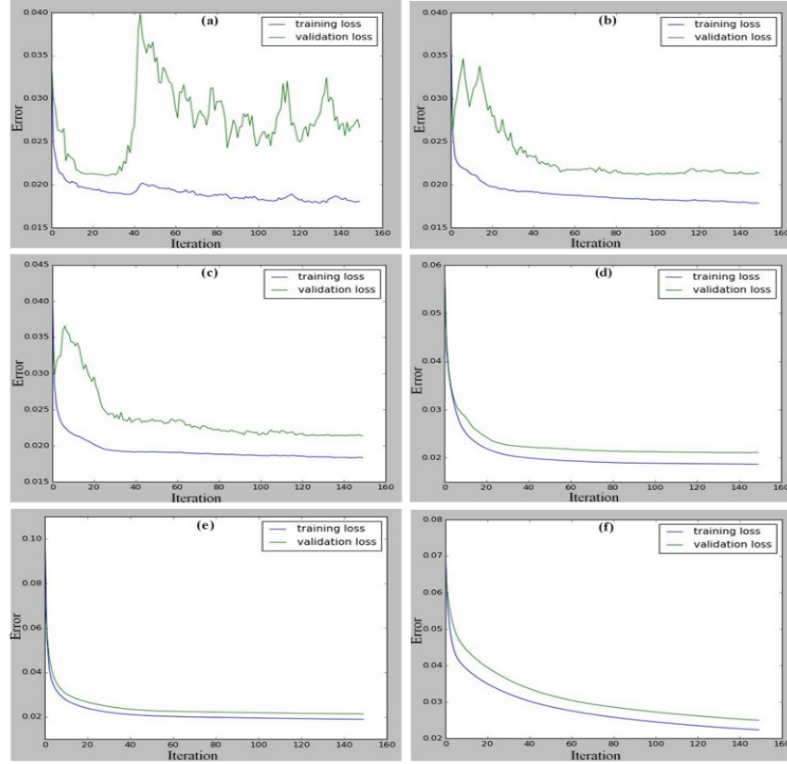


Fig. 6. Training loss and validation loss with different learning rates

Table 1. Mean RMSE and mean MAE of different LSTM units

	Mean RMSE of 67 Stations	Mean MAE of 67 Stations
32 LSTM Units	4.2022	3.1421
64 LSTM Units	4.1936	3.1309
128 LSTM Units	4.2147	3.1588

Table 1 shows the experiment with three different numbers of LSTM units (32 units, 64 units and 128 units). A PM 2.5 forecasting model was built for each station and the quality of the model was based on the mean error of the results from the 67 stations. A learning rate of 0.00005 with 350 iterations was used. This experiment aimed to determine the effect of number of LSTM units on the model. The experiments showed that the LSTM model with 64 units yielded the best result based on the mean error of the 67 stations with a mean RMSE of 4.1936 and MAE of 3.1309. But, several units in the LSTM model did not have any significant effects. This assessment was based on the limited reduction in errors.

Also 350 iterations was used to forecast the next hour of PM 2.5 concentrations with a learning rate of 0.00005 and 64 LSTM units. The result is shown in Table 2.

Table 2. Mean RMSE and mean MAE of different number of data

Input Data window	Mean RMSE of 67 Stations	Mean MAE of 67 Stations
1-hour	4.4541	3.3472
3-hours	4.1936	3.1309
6-hours	4.2731	3.2019
12-hours	4.1698	3.1117
24-hours	4.1831	3.1202
48-hours	4.1750	3.1348
72-hours	4.1474	3.0973

The 72-hour data window gave the best result with a mean RMSE of 4.1474 and mean MAE of 3.0973. However, more time was needed to train the model. Using a 3-hour window of data to predict the next hour of PM 2.5 concentrations yielded a mean RMSE of 4.1936 and mean MAE of 3.1309. The difference of errors when using 3-hour windows and 72-hour windows was 0.0462. Even though the difference was small, the processing time was very different.

We also explored the use of data from neighboring areas. This experiment was configured with 3-hour data windows, 350 iterations, 64 LSTM units and a learning rate of 0.00005. The results are shown in Table 3.

Table 3. Result of the model with neighbors and without neighbors

	Mean RMSE of 67 Stations	Mean MAE of 67 Stations
With neighbors	4.1936	3.1309
Without neighbors	4.3926	3.2466

The results show that including data about the neighboring areas gave more beneficial results than without data about neighboring areas. Clearly, air pollution in one area can affect neighboring areas.

Finally, an experiment was conducted using only PM 2.5 data. This experiment was configured with a 3-hour data windows, 350 iterations, 64 LSTM units and a learning rate of 0.00005. The results are shown in Table 4.

Table 4. Result of the model that only used PM 2.5 and other important attributes

	Mean RMSE of 67 Stations	Mean MAE of 67 Stations
Use important attributes	4.1936	3.1309
Only use PM 2.5 attributes	4.3709	3.2443

The mean results of the 67 stations showed that the model which included the important attributes gave a better performance than the model that only used PM 2.5 attributes. This means that other attributes such as meteorological conditions and other air pollutants played an important role in the PM 2.5 forecasts.

The results of the 67 stations are shown in Table 5. These results show that Guanshan station achieved the best result with an RMSE of 1.9509 and a MAE of 1.3213. This was followed by Taitung and Cailiao stations with an RMSE of 2.3378 and 2.4441, respectively. Some of the stations exhibited higher error rates such as Linyuan station, Qianzhen station and Xiaogang station. Linyuan station yielded an RMSE of 9.1869

and a MAE of 6.4037. Qianzhen yielded an RMSE of 7.8106 and a MAE of 5.6956. These were followed by Xiaogang with an RMSE of 6.6636 and a MAE of 4.9729. Overall, the proposed method demonstrated good performance with a mean RMSE of 4.1936 and a mean MAE of 3.1309.

5 Conclusion

PM 2.5 particles in the air can affect human health. The PM 2.5 concentration correlates with other pollutants and meteorological conditions in neighboring areas as PM 2.5 is easily carried by the wind from one area to another.

This study used air pollutant concentrations and meteorological conditions to make one-hour forecasts of PM 2.5 concentrations. Neighboring areas are determined based on station location clustering using X-means clustering. Then, we calculated the radius of each station based on the mean distance of each cluster. LSTM was applied as a forecasting engine to make one-hour PM 2.5 concentration forecasts.

Table 5. Station results

Station	RMSE	MAE	Station	RMSE	MAE
Annan	4.2382	3.2652	Puli	3.2002	2.3662
Banqiao	3.0229	2.3447	Puzi	4.3750	3.3637
Cailiao	2.4441	1.8586	Qianjin	5.4130	4.1601
Changhua	5.4659	4.2174	Qianzhen	7.8106	5.6956
Chaozhou	6.5018	4.7824	Qiaotou	3.4655	2.5891
Chiayi	5.9012	4.5605	Renwu	5.0395	3.8873
Dali	4.5240	3.4289	Shalu	4.0696	2.9086
Daliao	4.1197	3.1172	Shanhua	5.0095	3.8823
Dayuan	4.1146	3.0850	Shilin	3.4324	2.6421
Dongshan	4.2207	3.0046	Sinyin	4.7020	3.5364
Douliu	5.6606	4.3191	Songshan	4.0920	3.0924
Erlin	4.5393	3.3183	Tainan	4.5315	3.4037
Fengshan	5.1525	3.8981	Taitung	2.3378	1.7630
Fengyuan	4.6981	3.6365	Taixi	4.4863	3.0750
Fuxing	4.1323	3.0565	Taoyuan	3.9041	3.0017
Guanshan	1.9509	1.3213	Toufen	3.0393	2.1690
Guanyin	4.2012	3.0137	Tucheng	2.5527	1.9375
Guting	3.3802	2.5774	Wanhua	3.8678	3.0053
Hengchun	2.9969	2.1382	Wanli	2.9183	2.0894
Hsinchu	2.9843	2.1693	Xianxi	3.4886	2.5272
Hualian	2.7852	2.1597	Xiaogang	6.6636	4.9729
Hukou	3.6707	2.7523	Xindian	3.1038	2.2755
Jilong	4.5794	3.0979	Xingang	4.5044	3.3591
Linkou	3.3532	2.4243	Xinzhuang	2.5083	1.9533
Linyuan	9.1869	6.4073	Xitun	3.8740	3.0217

Longtan	3.7045	2.9059	Xizhi	2.5024	1.8857
Lunbei	4.7674	3.5696	Yonghe	2.8041	2.1042
Mailiao	4.1385	2.9405	Zhongli	4.1864	3.1976
Meinong	4.7848	3.5040	Zhongming	4.1473	3.2268
Miaoli	4.5506	2.9245	Zhongshan	3.8098	2.9349
Nantou	3.4589	2.4960	Zhudong	3.4594	2.6179
Nanzi	5.4725	4.1397	Zhushan	4.2792	3.3270
Pingtung	5.0589	3.8742	Zuoying	5.6663	4.3459
Pingzhen	3.9752	3.1446	Mean	4.1936	3.1309

Experimental results demonstrated that the proposed approach could effectively make one-hour PM 2.5 concentration forecast for 67 stations. The model achieved a RMSE of 1.9509 and a MAE of 1.3213 for Guanshan station. The overall result also showed a relatively low mean RMSE and MAE for the 67 stations all around Taiwan. The mean RMSE was 4.1936 and the mean MAE was 3.1309.

Acknowledgments

This study was funded by the Ministry of Science and Technology, Taiwan, under Grants MOST107-2221-E-027-113-, MOST108-2321-B-027-001- and MOST108-2221-E-027-111-MY3.

References

1. MacMunn, A.: More than 4 in 10 Americans live with unhealthy air according to 2018 ‘state of the air’ report. American Lung Association,]. Accessed May 2, 2019 from: <https://www.lung.org/about-us/media/press-releases/2018-state-of-the-air.html>, (2018).
2. WHO.: Ambient (outdoor) air quality and health. Accessed May 2, 2019 from [https://www.who.int/en/news-room/fact-sheets/detail/ambient-\(outdoor\)-air-quality-and-health](https://www.who.int/en/news-room/fact-sheets/detail/ambient-(outdoor)-air-quality-and-health), (2018).
3. Jacobs, E. T., Burgess, J. L., Abbott, M. B.: The Donora smog revisited: 70 years after the event that inspired the clean air act. *American Journal of Public Health* **108**, 85-88 (2018).
4. Wilkins, E. T.: Air pollution aspects of the London fog of December 1952. *Quarterly Journal of the Royal Meteorological Society* **80**(344), 267-271 (1954).
5. Ganesh, S. S., Modali, S. H., Palreddy, S. R., Arulmozhivarman, P.: Forecasting air quality index using regression model: a case study on Delhi and Houston. In: Proc. of Int. Conf. on Trends in Electronics and Informatics, pp.248–254. IEEE (2017).
6. Shaban, K. B., Kadri, A., Rezk, E.: Urban air pollution monitoring system with forecasting models. *IEEE Sensor Journal* **16**(8), 2598–2606 (2016).
7. Gu, K., Qiao, J., Lin, W.: Recurrent air quality predictor based on meteorology- and pollution-related factor. *IEEE Trans. on Industrial Informatics* **14**(9), 3946-3955 (2018).
8. Tsai, Y., Zheng, Y., Cheng, Y.: Air pollution forecasting using RNN with LSTM. In: Proc. of IEEE 16th Int. Conf. on Big Data Intelligence and Computing and Cyber Science and Technology Congress, pp.1074–1079. IEEE (2018).

9. Oprea, M., Mihalache, S. F., Popescu, M.: A comparative study of computational intelligence techniques applied to PM 2.5 air pollution forecasting. In: Proc. of 6th Int. Conf. on Computers Communications and Control, pp.103–108. IEEE (2016).
10. Chen, Y., Qin, H., Zhou, Z.: A comparative study on multi-regression analysis and BP neural network of PM2.5 index. In: Proc. of 10th Int. Conf. on Natural Computation, pp.155–159. IEEE (2014).
11. Han, J., Kamber, M., Pei, J.: Data Mining: Concepts and Techniques. 3rd ed., Morgan Kaufmann (2011).
12. Imamura, K., Kubo, N., Hashimoto, H.: Automatic moving object extraction using x-means clustering. In: Proc. of the 28th Picture Coding Symposium, pp.245–249. IEEE (2010).
13. Nathanson, J. A.: Air pollution. *Encyclopaedia Britannica*. Accessed May 9, 2019 from <https://www.britannica.com/science/air-pollution> (2018).
14. Liu, L., He, G., Shi, X., Song, H.: Metadata extraction based on mutual information in digital libraries. In: Proc. of the First IEEE Int. Symp. on Information Technologies and Applications in Education, pp.209-212. IEEE (2007).
15. Goodfellow, I., Bengio, Y., Courville, A.: Deep Learning. MIT Press Cambridge, MA (2017).
16. Willmott, C. J., Matsuura, K.: Advantages of the mean absolute error (MAE) over the root mean square error (RMSE) in assessing mean model performance. *Climate Research* **30**(79), 79-82 (2005).
17. Lv, B., Cai, J., Xu, B., Bai, Y.: Understanding the rising phase of the PM 2.5 concentration evolution in large china cities. *Scientific Report* **7**(46456), (2017).
18. Niharika, M. V., Padma, S. R.: A survey on air quality forecasting techniques, *International Journal of Computer Science and Information Technologies* **5**(1), 103-107 (2014).