

Forecasting of Air Pollution Index with Artificial Neural Network

Nur Haizum Abd Rahman^a, Muhammad Hisyam Lee^{a*}, Mohd Talib Latif^b, Suhartono^c

^aDepartment of Mathematical Sciences, Universiti Teknologi Malaysia, 81310 UTM Johor Bahru, Johor, Malaysia

^bSchool of Environmental and Natural Resource Sciences, Universiti Kebangsaan Malaysia, Selangor, Malaysia

^cDepartment of Statistics, Institut Teknologi Sepuluh Nopember, Surabaya, Indonesia

*Corresponding author: mhl@utm.my

Article history

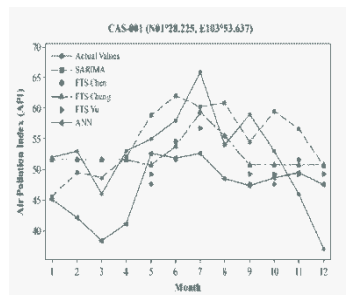
Received :21 January 2013

Received in revised form :

7 May 2013

Accepted :25 June 2013

Graphical abstract



Abstract

In recent years, the arisen of air pollution in urban area address much attention globally. The air pollutants has emerged detrimental effects on health and living conditions. Time series forecasting is the important method nowadays with the ability to predict the future events. In this study, the forecasting is based on 10 years monthly data of Air Pollution Index (API) located in industrial and residential monitoring stations area in Malaysia. The autoregressive integrated moving average (ARIMA), fuzzy time series (FTS) and artificial neural network (ANNs) were used as the methods to forecast the API values. The performance of each method is compare using the root mean square error (RMSE). The result shows that the ANNs give the smallest forecasting error to forecast API compared to FTS and ARIMA. Therefore, the ANNs could be consider a reliable approach in early warning system to general public in understanding the air quality status that might effect their health and also in decision making processes for air quality control and management.

Keywords: Fuzzy time series; artificial neural network; ARIMA; Air Pollution Index (API); time series; forecasting

Abstrak

Sejak kebelakangan ini, pencemaran udara di kawasan bandar mendapat perhatian di seluruh dunia. Bahan pencemar udara boleh menjejaskan kesihatan kepada semua hidupan. Ramalan siri masa adalah kaedah penting pada masa kini kerana mampu untuk meramalkan peristiwa pada masa hadapan. Dalam kajian ini, 10 tahun data bulanan Indeks Pencemaran Udara (IPU) dari stesen pemantauan yang terletak di kawasan industri serta penempatan di Malaysia. Purata bergerak bersepadu autograsi (ARIMA), siri masa kabur (FTS) dan rangkaian neural tiruan (ANN) telah digunakan sebagai kaedah untuk meramal nilai IPU. Prestasi setiap kaedah akan dibandingkan menggunakan punca min ralat persegi (RMSE). Hasilnya, ANN memberi ralat terkecil untuk meramal IPU berbanding FTS dan ARIMA. Oleh itu, ANN boleh dipertimbangkan sebagai pendekatan dalam sistem amaran awal kepada orang awam dalam memahami status kualiti udara yang memberi kesan kepada kesihatan mereka dan juga dalam proses membuat keputusan bagi kawalan kualiti udara dan pengurusan.

Kata kunci: Siri masa kabur; rangkaian neural tiruan; ARIMA; Indeks Pencemaran Udara (IPU); siri masa; ramalan

© 2013 Penerbit UTM Press. All rights reserved.

1.0 INTRODUCTION

Time series is an important area of forecasting with active research in variety of areas. The historical observations of the variable are analyzed in developing a model that describes the relationship between time and variable. Then, the model is used to extrapolate the time series into future. Much effort in development and improvement of time series forecasting models remains as the big issues for over several decades. Generally, the time series methods divided into classical methods and modern methods. Conventionally, the researchers tend to employ the classical methods of time series as the initial way for analysis, modelling and forecasting. The most widely important classical

methods is the autoregressive integrated moving average (ARIMA) model [1]. ARIMA provides reasonable accuracy in forecasting however the major limitation is ARIMA model only can captured the linear form of time series data.

As the alternative solution, many researchers tend to develop new methods that can overcome the limitation of classical methods such as artificial neural network (ANN) and fuzzy time series (FTS). ANN is the nonlinear time series forecasting methods with the particular model form is unspecified [2, 3]. In previous study, ANN is one of the most accurate and widely used forecasting models in many applications [4, 5]. Mostly used in comparison study between classical model ARIMA or time series regression[6]. FTS is the branch of artificial intelligence which is

appropriate in decision making process in complex systems when the situation of the problem is often unclear [7]. The FTS has been used in the field of air pollution cases by several authors [7-9]. In the past research, the fuzzy set applied in terms of membership function which means describing the characteristics of air pollution concentrations in terms of good, medium or poor quality [7, 8, 10]. Then, the characteristics of air pollution concentration become the input variables for the fuzzy set. In contrast, the input variables for the application of FTS in this study is obtain from the available data set with more than one fuzzy time series approaches (i.e. Chen's, Yu's and Cheng's).

Currently, major environmental issues can be seen mainly in climate change and pollutions. Air pollution is the fundamental pollution problem that described as multiple spatial and temporal scales, including complex chemical and physical mechanism. It may escalates as a function of human activity and it is highly nonlinear as a problem [11]. It firstly concern the effect towards human health such as asthma, headache and dizziness and secondly towards the environment [8, 12]. For long term consequences, the air pollution tends to increase the earth risk on global warming and greenhouse [13]. Numerous researches have been done in order to monitor the air pollutants [14, 15]. Mostly in monitoring the main attention in air pollution problem namely as particulate matter (PM₁₀), ozone (O₃), carbon monoxide (CO₂), sulphur dioxide (SO₂) and nitrogen dioxide (NO₂) [13, 16-18]. However, some research done by looking at the air quality status based on the most dominant air pollutant with the highest concentration considered as the contaminants that will determine the value of Air Pollution Index (API) [12, 19].

Accurate forecasts of air quality are the essential for efficient planning by the various sectors of the related to the economics performances. Additionally, the forecast accuracy is important in strategic management to maintain the air quality. Time series modeling and forecasting is necessary to prevent the air pollution situation worsen in the long run [15, 20]. This approach is selected to be used in air quality management to help making future planning and helps to shape the better air quality since time series analysis is the major task for researchers used in development [21]. The main objective of this work is to construct and develop the accurate statistical forecasting models to predict the monthly API and to evaluate such models in order to monitor the air quality status.

2.0 MATERIALS AND METHODS

2.1 Description of the Sampling Site

The Southern Coastal Region of the Peninsular Malaysia, Johor Bahru, Johor has been chosen as the study site which is the second largest metropolitan area in Malaysia after the capital city, Kuala Lumpur [22]. The data measured at continuous monitoring stations located at N01°28.225, E103°53.637 (CAS 001) in Johor Bahru, established on October 1995. Johor Bahru is the capital of Johor State with presence of large number of residential, industries and commercial hotspot in the region. The 10 years monthly data was divided into training data set from 2000 – 2008 (108 observations) to identify the API model and testing data set in 2009 with total 12 observations to check the model performance.

2.2 Air Pollution Index (API)

A simple generalized way in describing the air quality status is through the index system, API. Through the API, general public can understand easily the air quality status for their own health

precaution. The API was classify based on the highest index value of five main air pollutants namely as PM₁₀, O₃, CO₂, SO₂ and NO₂. In order to assess the air quality status on human health, the API scale and terms used in describing the air quality status. This information with different ranges reflects as “Good (0-50), Moderate (51-100), Unhealthy (101-200), Very Unhealthy (201-300) and Hazardous (301 and above)” can be as the benchmark of air quality management or data interpretation of decision making processes.

2.3 Box-Jenkins Modeling Approach

The Box-Jenkins or ARIMA is classified as linear models that capable in presenting both stationary and non-stationary time series. Most of researchers use this model to forecast univariate time series data. Three main steps that must be considered in building the model for forecasting are tentative identification, parameter estimation and diagnostic checking [23]. ARIMA remain as the popular models until nowadays because of the flexibility in represents different types of time series, autoregressive (AR), moving average (MA), autoregressive moving average (ARMA) and autoregressive moving average model (ARIMA) can be abbreviated ARIMA (p, d, q) when the dataset is nonstationary, the difference will take part to stationer the data. In the case where seasonal components are included in the model, the ARIMA model is called as SARIMA model can be abbreviated as SARIMA (p, d, q) (P, D, Q)^S. Generally the SARIMA model can be written as:

$$\varphi_p(B)\Phi_p(B^S)(1-B)^d(1-B^S)^D Y_t = \theta_q(B)\Theta_Q(B^S)\varepsilon_t \quad (1)$$

where:

$$\begin{aligned} \varphi_p(B) &= 1 - \varphi_1 B - \varphi_2 B^2 - \dots - \varphi_p B^p \\ \Phi_p(B) &= 1 - \Phi_1 B^S - \Phi_2 B^{2S} - \dots - \Phi_p B^{pS} \\ \theta_q(B) &= 1 - \theta_1 B - \theta_2 B^2 - \dots - \theta_p B^p \\ \Theta_Q(B) &= 1 - \Theta_1 B^S - \Theta_2 B^{2S} - \dots - \Theta_Q B^{QS} \end{aligned}$$

where B denoted as the backward shift operator, d and D denote as the non-seasonal and seasonal order of difference respectively.

2.4 Neural Network

Artificial neural networks are flexible in capturing the nonlinear time series data and have been suggested as an alternative model in time series forecasting [6, 24]. For this study, we used the multi-layer perceptron (MLP) consists of processing element (PEs), called as neurons. An MLP is typically composed of several layers of nodes with the first layer is an input layer where the external information is received. The last or the highest layer is an output layer where the problem solution is obtained. The input layer and output layer are separated by one or more intermediate layers called the hidden layers. MLP characterized for three layers (input, hidden, output) of simple processing units connected by acyclic links shown in Figure 1.

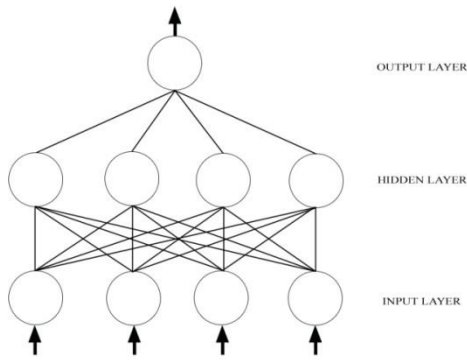


Figure 1 A simple neural network structure

Arrows from input to hidden layer and hidden layer to output layer indicate the strength of each connection and can be measured by a quantity called as weight. In the input layer and hidden layer, each unit processes it with a transfer function or activation function and lastly distributes the result to the output layer [4]. Commonly, the activation function used is the logistic function, $f(x) = 1 / \{1 + \exp(-x)\}$ in the hidden layer meanwhile the linear function, $f(x) = x$ used at the output stage.

The simplest MLP with input layer and hidden layer can be calculates with this function

$$y_t = w_0 + \sum_{i=1}^n w_i y_i \tag{2}$$

Simple network structure that has a small number of hidden nodes often works well in out-sample forecasting. Single hidden layer is the most widely used model form for time series modeling and forecasting. Hidden layer included in neural network system to increase the modeling flexibility. The function can be written as:

$$y_t = w_0 + \sum_{j=1}^q w_j \cdot g(w_{0,j} + \sum_{i=1}^p w_{i,j} \cdot y_{t-i}) + \varepsilon_t \tag{3}$$

where w_j ($j=0, 1, 2, \dots, q$) and w_{ij} ($i=0, 1, 2, \dots, p; j=1, 2, \dots, q$) are the model parameters often called as the connection weights; p is the number of input nodes and q is the number of hidden nodes.

2.5 Fuzzy Time Series Analysis

Fuzzy time series has been used in the field of air pollution by several authors [8]. According to Song and Chissom[25, 26], generally the concepts of fuzzy time series can be defined as let U be the universe of dicourse, where $U = \{u_1, u_2, \dots, u_b\}$ and $U = [D_{\min} - D_1, D_{\max} + D_2] = [\text{begin}, \text{end}]$. A fuzzy set (A_i) of U is defined as $A_i = f_{A_i}(u_1) / u_1 + f_{A_i}(u_2) / u_2 + \dots + f_{A_i}(u_b) / u_b$, where f_{A_i} is the membership function of the fuzzy set A_i , $f_{A_i} : U \rightarrow [0,1]$. u_a is a generic element of fuzzy set A_i , and $f_{A_i}(u_a)$ is the grade of membership of u_a in A_i , where $f_{A_i}(u_a) \in [0,1]$ and $1 \leq a \leq b$.

Based on Song and Chissom[25, 26], Chen [27] improved the establishment step of fuzzy relationships which it used a

simple operation and instead of complex matrix operations. In Chen, the repeated or the recent identical FLRs were simply ignored since the same FLR may not reflect the real world situation. For an example below, the Chen method ignore the relationship for $(t = 3)$ and $(t = 4)$. Therefore, the fuzzy logical relationship group (FLRG) left as $A_1 \rightarrow A_1, A_2$.

- $(t = 1) A_1 \rightarrow A_1$,
- $(t = 2) A_1 \rightarrow A_2$,
- $(t = 3) A_1 \rightarrow A_1$,
- $(t = 4) A_1 \rightarrow A_1$

In contrast, Yu [28] proposed that the same FLR must be considered in forecasting since the recent FLR has the greater weight. To illustrate, the Yu FTS can be shown as below. The most recent FLR ($t = 4$) is assigned the highest weight of 4 indicate the high probability occurrence in the future. Conversely, the most aged with $(t = 1)$, is assigned with the lowest weight of 1 indicate the lowest probability.

- $(t = 1) A_1 \rightarrow A_1$ with weight 1,
- $(t = 2) A_1 \rightarrow A_2$ with weight 2,
- $(t = 3) A_1 \rightarrow A_1$ with weight 3,
- $(t = 4) A_1 \rightarrow A_1$ with weight 4,

The probability of weight appearance and the importance of chronological FLR for the same recent identical FLRS are the main focused FTS proposed by Cheng [29]. The weights illustrate as follows:

- $(t = 1) A_1 \rightarrow A_1$ with weight 1,
- $(t = 2) A_1 \rightarrow A_2$ with weight 1,
- $(t = 3) A_1 \rightarrow A_1$ with weight 2,
- $(t = 4) A_1 \rightarrow A_1$ with weight 3,

The method for forecasting the air pollution index (API) using FTS Chen, Yu and Cheng algorithm is simply can be presented as follows:

Step 1: Define and partition the universe of discourse $U = [D_{\min} - D_1, D_{\max} + D_2]$ into several equal intervals denoted as

$$u_1, u_2, \dots, u_m.$$

Step 2: Based on SARIMA model, the fuzzy logical relationship (FLR) are determined. This procedure also has been proposed by Faraway [30] in order to select the neural network input variable.

Step 3: In order to select the best input for FLR, different combination input will be attempt. Initially with single input followed by two inputs, three inputs and lastly four inputs for station located in Muar.

Step 4: The optimum length of intervals was calculated following the average-based length [31]

Step 5: Calculate the forecasted outputs. Three different fuzzy time series, Chen's [32], Yu's [28] and Cheng's [29] are used in this study.

Step 6: The forecasted result were compared based on the smallest error.

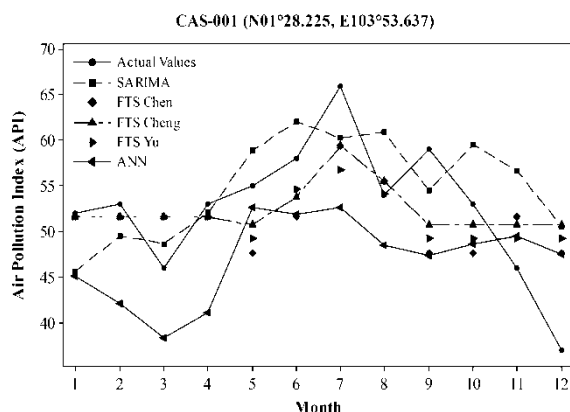


Figure 3 Time series plot of API testing

In general, the result based on residual error evaluation performance give the similar result of the best model used in API forecasting. However, the RMSE is the most common evaluation performance based on the difference between the observed and forecast values with high sensitivity in extreme values due to the power term. Therefore, in short, the model performance testing data set used in forecasting of each station will be based on the smallest RMSE values as shown in Table 2. The actual values and predicted values time series plot for stations CAS 001 shown in Figure 3. The MLP give the best result in forecasting the API with the neural model $p=3$ and $q=2$.

4.0 CONCLUSION

The ARIMA model has become the main forecasting methods in many research and practice. However, the accuracy of time series forecasting is important in decision making process. More recently, the fuzzy time series (FTS) and artificial neural network (ANN) were proposed because of the capability in exploiting the strengths of traditional time series approaches, ARIMA. In this paper, an application of time series forecasting towards air quality problem had been done. The results showed that the neural network is a comparative with FTS (i.e. Chen's, Yu's, Cheng's) and ARIMA with approximately 88% and 89% more accurate respectively. This result shows that the ANNs are flexible intelligence forecasting methods which can provide a useful and effective tool for modeling the complex and poorly understandable processes. Thus, the neural network can be used as the benchmark modern forecasting methods in developing new methods in order to improve the forecasting accuracy. As the conclusion, the methodology developed in this study is acceptable in the sampling site however with this model building approach it may apply in difference topographical features. Moreover, this methodology appears to be very useful for local administration health and environmental institutions in monitor and control the future pollutant trend precisely.

Acknowledgement

The researchers would like to thank Department of Environment, Malaysia for providing pollutants data and also the Ministry of Higher Education, Malaysia for support to accomplish this study.

References

- [1] G. E. P. Box and G. M. Jenkins. 1976. *Time Series Analysis: Forecasting and Control*. 1st ed. San Francisco: Holden-Day.
- [2] G. P. Zhang. 2003. Time Series Forecasting using a Hybrid ARIMA and Neural Network Model. *Neurocomputing*. 50: 159–175.
- [3] F. Li. 2010. Air Quality Prediction in Yinchuan by Using Neural Networks. In *Advances in Swarm Intelligence*. Vol. 6146, Y. Tan, Y. Shi, and K. Tan, Eds. Ed: Springer Berlin/Heidelberg. 548–557.
- [4] K. Moustiris, I. Ziomas, and A. Paliatos. 2010. 3-Day-Ahead Forecasting of Regional Pollution Index for the Pollutants NO₂, CO, SO₂, and O₃ Using Artificial Neural Networks in Athens, Greece. *Water, Air, & Soil Pollution*. 209: 29–43.
- [5] M. Khashei and M. Bijari. 2010. An Artificial Neural Network (p, d, q) Model for Time Series Forecasting. *Expert Systems with Applications*. 37: 479–489.
- [6] V. R. Prybutok, J. Yi, and D. Mitchell. 2000. Comparison of neural Network Models with ARIMA and Regression Models for Prediction of Houston's Daily Maximum Ozone Concentrations. *European Journal of Operational Research*. 122: 31–40.
- [7] F. Bernard. 2003. Fuzzy environmental Decision-making: Applications to Air Pollution. *Atmospheric Environment*. 37: 1865–1877.
- [8] J.-S. Heo and D.-S. Kim. 2004. A New Method of Ozone Forecasting Using Fuzzy Expert and Neural Network Systems. *Science of the Total Environment*. 325: 221–237.
- [9] F. C. Morabito and M. Versaci. 2003. Fuzzy Neural Identification and Forecasting Techniques to Process Experimental Urban Air Pollution Data. *Neural Networks*. 16: 493–506.
- [10] F. Bernard, E. A. 2006. Fuzzy Approaches to Environmental Decisions: Application to Air Quality. *Environmental Science & Policy*. 9: 22–31.
- [11] K. D. Karatzas, G. Papadourakis, and I. Kyriakidis, 2008. Understanding and Forecasting Atmospheric Quality Parameters with the Aid of ANNs," presented at the Neural Networks, 2008. IJCNN 2008. (IEEE World Congress on Computational Intelligence). IEEE International Joint Conference on Hong Kong.
- [12] A. Kumar and P. Goyal, 2011. Forecasting of Daily Air Quality Index in Delhi. *Science of The Total Environment*. 409: 5517–5523.
- [13] A. Kurt and A. B. Oktay, 2010. Forecasting Air Pollutant Indicator Levels with Geographic Models 3 Days in Advance Using Neural Networks. *Expert Systems with Applications*. 37: 7986–7992.
- [14] A. Chaloulakou, M. Saisana, and N. Spyrellis, 2003. Comparative Assessment of Neural Networks and Regression Models for Forecasting Summertime Ozone in Athens. *Science of the Total Environment*. 313: 1–13.
- [15] A. Vlachogianni. 2011. P. Kassomenos, A. Karppinen, S. Karakitsios, and J. Kukkonen. Evaluation of a multiple regression model for the forecasting of the concentrations of NO_x and PM₁₀ in Athens and Helsinki. *Science of The Total Environment*. 409: 1559–1571.
- [16] U. Brunelli, V. Piazza, L. Pignato, F. Sorbello, and S. Vitabile. 2007. Two-days Ahead Prediction of Daily Maximum Concentrations of SO₂, O₃, PM₁₀, NO₂, CO in the Urban Area of Palermo, Italy. *Atmospheric Environment*. 41: 2967–2995.
- [17] Department of Environment. 2004. Malaysia Environment Quality Report 2004. Department of Environment, Putrajaya.
- [18] N. Ghazali, N. Ramli, A. Yahaya, N. Yusof, N. Sansuddin, and W. Al Madhoun. 2010. Transformation of Nitrogen Dioxide into Ozone and Prediction of Ozone Concentrations Using Multiple Linear Regression Techniques. *Environmental Monitoring and Assessment*. 165: 475–489.
- [19] Y. S. Lim, Y. C. Lim, and M. J. W. Pauline. 2008. ARIMA and Integrated ARFIMA Models for Forecasting Air Pollution Index in Shah Alam, Selangor. *The Malaysian Journal of Analytical Sciences*. 12: 257–263.
- [20] T. Slini, K. Karatzas, and N. Moussiopoulos. 2002. Statistical Analysis of Environmental Data as the Basis of Forecasting: An Air Quality Application. *The Science of The Total Environment*. 288: 227–237.
- [21] J. D. Cryer. 1986. *Time Series Analysis*. 1st ed. United States of America: Duxbury Press.
- [22] A. Rizzo and J. Glasson. 2011. Iskandar Malaysia. *Cities*.
- [23] J. E. Hanke and D. W. Wichern. 2005. *Business Forecasting*. 8th ed. Upper Saddle River, NJ: Pearson/Prentice Hall.
- [24] G. Zhang, B. Eddy Patuwo, and M. Y. Hu. 1998. Forecasting with Artificial Neural Networks: The State of the Art. *International Journal of Forecasting*. 14: 35–62.
- [25] Q. Song and B. S. Chissom. 1993. Fuzzy Time Series and its Models. *Fuzzy Sets and Systems*. 54: 269–277.
- [26] Q. Song and B. S. Chissom. 1993. Forecasting Enrollments with Fuzzy Time Series—Part I. *Fuzzy Sets and Systems*. 54: 1–9.

- [27] C. Shyi-Ming. 1996. Forecasting Enrollments based on Fuzzy Time Series. *Fuzzy Sets and Systems*. 81: 311–319.
- [28] Y. Hui-Kuang. 2005. Weighted Fuzzy Time Series Models for TAIEX Forecasting. *Physica A: Statistical Mechanics and its Applications*. 349: 609–624.
- [29] C.-H. Cheng, T.-L. Chen, H. J. Teoh, and C.-H. Chiang. 2008. Fuzzy Time-Series based on Adaptive Expectation Model for TAIEX Forecasting. *Expert Systems with Applications*. 34: 1126–1132.
- [30] N. Sansuddin, N. Ramli, A. Yahaya, N. Yusof, N. Ghazali, and W. Madhoun. 2011. Statistical analysis of PM10 concentrations at different locations in Malaysia. *Environmental Monitoring and Assessment*. 180: 573–588.
- [31] S. Hassanzadeh, F. Hosseinibalam, and R. Alizadeh. 2009. Statistical Models and Time Series Forecasting of Sulfur Dioxide: A Case Study Tehran. *Environmental Monitoring and Assessment*. 155: 149–155.
- [32] S. M. Chen. 2002. Forecasting Enrollments based on High-order Fuzzy Time Series. *Cybernetics and Systems*. 33: 1–16.