# Identification of significant factors for air pollution levels using a neural network based knowledge discovery system

Kit Yan Chan[a] and Le Jian[b*]

[a] Department of Electrical and Computer Engineering, Curtin University, Perth, AUSTRALIA

[b] School of Public Health, Curtin Health Innovation Research Institute, Curtin University, Perth, AUSTRALIA

**\*Corresponding author**:

Dr. Le JIAN

Address:

    School of Public Health

    WHO Collaborating Centre for Environmental Health Impact Assessment

    Curtin Health Innovation Research Institute

    Faculty of Health Sciences

    Curtin University

    Kent Street

    GPO BOX U1987

    Perth, Western Australia 6845

Phone: +61 8 9266 4250

Fax:    +61 8 9266 2958

Email:  l.jian@curtin.edu.au

**Abstract**

Artificial neural network (ANN) is a commonly used approach to estimate or forecast air pollution levels, which are usually assessed by the concentrations of air contaminants such as nitrogen dioxide, sulfur dioxide, carbon monoxide, ozone, and suspended particulate matters (PMs) in the atmosphere of the concerned areas. Even through ANN can accurately estimate air pollution levels they are numerical enigmas and unable to provide explicit knowledge of air pollution levels by air pollution factors (e.g. traffic and meteorological factors). This paper proposed a neural network based knowledge discovery system aimed at overcoming this limitation in ANN. The system consists of two units: a) an ANN unit, which is used to estimate the air pollution levels based on relevant air pollution factors; b) a knowledge discovery unit, which is used to extract explicit knowledge from the ANN unit. To demonstrate the practicability of this neural network based knowledge discovery system, numerical data on mass concentrations of PM2.5 and PM1.0, meteorological and traffic data measured near a busy traffic road in Hangzhou city were applied to investigate the air pollution levels and the potential air pollution factors that may impact on the concentrations of these PMs. Results suggest that the proposed neural network based knowledge discovery system can accurately estimate air pollution levels and identify significant factors that have impact on air pollution levels.

*Keywords*: Artificial neural network; Main effect analysis; Air pollution; Air monitoring; Meteorological factors; Particulate matter

## 1. Introduction

Air pollution is a major environmental risk to health in many developed and developing cities of the world. The air pollution levels are usually determined by the concentrations of air pollutants such as nitrogen dioxide, sulfur dioxide, carbon monoxide, ozone and suspended particulate matters (PMs). PMs are defined by the U.S. Environmental Protection Agency as "very small pieces of solid or liquid matter, such as particles of soot, dust, fumes, mists, or aerosols" [10]. They are usually produced by air pollution factors such as energy production from power plants, burning of fossil fuels in factories, power plants, industrial processes, residential heating, combustion of gasoline, diesel and hydrocarbon fuels in vehicles, etc. [5, 21]. Meanwhile, unfavorable meteorological factors may also affect the formation and growth of new air pollutants and the ability of the atmosphere to disperse air pollutants [3, 25, 38, 40]. Severe air pollution levels can be life threatening, can cause breathing difficulty, headache, dizziness, and result in heart attack [20]. Long term exposure to air pollutants can result in chronic respiratory and cardiovascular diseases including cancers [9, 17, 24]. Therefore, it is essential to monitor criteria air pollutants in the atmosphere by developing accurate models which can indicate the relationship between air pollution levels and air pollution factors. However, it is difficult to develop such models using traditional statistical methods, as they are unable to model complex nonlinear relationships between air pollution factors [1]. More recently, the universal estimator [22, 41, 42], namely artificial neural network (ANN), has been demonstrated their capability to model non-linear relationships between input and output variables to estimate, evaluate and forecast air pollution levels [12]. ANNs are unsupervised learning techniques whereby collected air pollution data are trained in order to create a black-box model, which maps between two domains, namely i) the domain for air pollution factors and ii) the domain for air pollution levels. In the literature, ANNs have been applied to estimate air pollutant levels such as concentrations of sulfur dioxide [6], carbon monoxide [27, 29], PMs [13, 33, 36, 43], and ozone [4, 8].
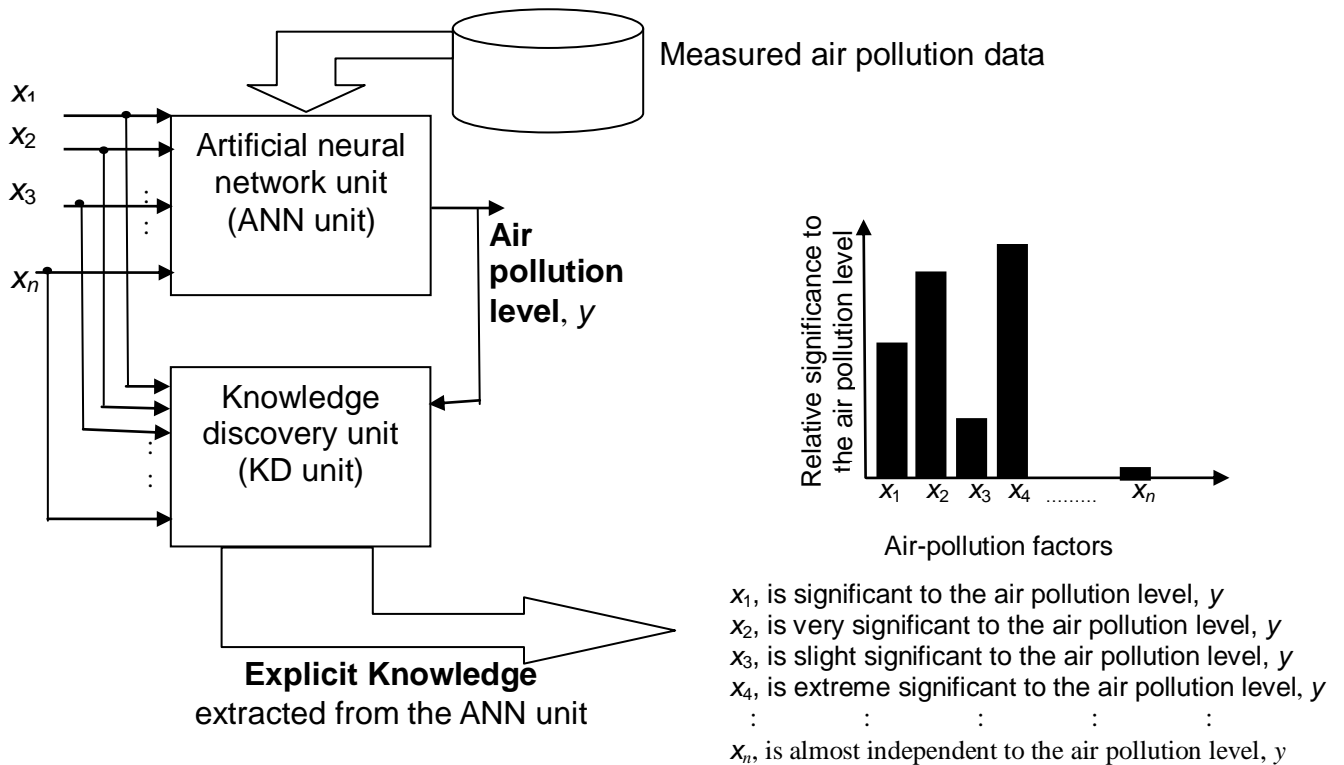
However, ANNs are black-box structures which yield no explicit knowledge [2]. Because of this limitation, these traditional ANNs may not be appropriate to use in estimating air pollution levels, even though they may achieve more accurate estimates than other explicit modeling methods such as statistical methods [30].

In this paper, a neural network based knowledge discovery system is proposed not only to estimate air pollution levels, but also to generate relative significances of air pollution factors to air pollution levels. This knowledge is important because it provides researchers with a better insight into the extent to which a particular air-pollution factor impacts on air-pollution levels. As transportation is a major source of air pollution worldwide, this knowledge is also useful for the transportation and infrastructure sectors to estimate pollution levels before planning or developing new transportation infrastructures.

## 2. A neural network based knowledge discovery system

A schematic representation of the proposed neural network based knowledge discovery system is depicted in Fig. 1. The system consists of two main units: a) the artificial neural network unit, namely ANN unit, and b) the knowledge discovery unit, namely KD unit. The ANN unit is developed to estimate the specified indicators of air pollution level, $y$ (e.g., nitrogen monoxide or PMs), based on the $n$ air pollution factors, $x_1$, $x_2$, … and $x_n$, which consist of meteorological and traffic flow factors. The indicators of air pollution level can be the levels of PMs and gaseous pollutants such as ozone, nitrogen monoxide and nitrogen dioxide. The meteorological factors can be temperature, relative humidity, barometric atmosphere pressure and wind speed. When data on both air pollution levels and air pollution factors are collected, an ANN unit can be developed. The KD unit is developed to extract informative knowledge from the ANN unit, which traditionally is a black-box or implicit in nature, and no explicit information can be indicated clearly between $x_1$, $x_2$, … $x_n$, and $y$. Based on the new KD unit, the relative

significance of the concerned air pollution factors, that mainly affect the air pollution level, can be indicated. If a slight change in an air pollution factor results in a great change in the air pollution level, this air pollution factor is considered as a critical contributor to the air pollution level. Therefore, it is essential and important to consider those significant factors in analyzing and controlling air pollution level.



**Fig. 1** The framework of neural network based knowledge discovery system

*2.1 The artificial neural network unit*

The ANN unit consists of three types of neural nodes: a) ***the input nodes***, which feed the air pollution factors $x_1$, $x_2$, … and $x_n$ into the ANN unit; b) ***the output node***, which estimates the air pollution level, $y$; and c) ***the hidden nodes***, which link input nodes of the air pollution factors and the

5

output nodes of the air pollution level. The ANN unit is suitable for complex and nonlinear interactions between air pollution factors in order to estimate the air pollution level based on the input-output functional relationship, $f$, which is denoted by the equation (1):

$$y = f\left(x_1, x_2, ..., x_n\right) = \sum_{j=1}^{n_h} w_j \Psi\left(\sum_{i=1}^{n}\left(v_{ij} x_i - b_j\right)\right) - b \tag{1}$$

where $n_h$ denotes the number of hidden nodes of the ANN unit;

$w_j$, denotes the weight of the link between the $j$-th hidden node and the output node for the air pollution level, $y$ with $j$=1, 2, ..., $n_h$;;

$v_{ij}$, denotes the weight between the input node for the $i$-th air pollution factor, $x_i$, and the $j$-th hidden node with $i$=1,2,..., $n$ and $j$=1, 2, ..., $n_h$;

$b_j$ and $b$, denote the biases for the $j$-th hidden nodes and the output node respectively;

$\Psi$ is the transfer function of the hidden set in which the sigmoid function is used.

The ANN weights are determined based on $N_D$ pieces of collected air pollution data in the form of

$$d\left(k\right) = \left[y\left(k\right), \varphi\left(k\right)\right] \quad \text{with } k\text{=1, 2, ... } N_D; \tag{2}$$

where $y\left(k\right)$ and $\varphi\left(k\right) = \left[x_1\left(k\right), x_2\left(k\right), ..., x_n\left(k\right)\right]$ are the $k$-th air pollution data with respect to the air pollution level and the $n$ air pollution factors, respectively. The ANN unit is evaluated based on the mean absolute relative error, $e_{MARE}$, formulated in equation (3), where both small and large errors have the same weights. $e_{MARE}$ indicates the differences between the actual observations and the estimates of the ANN unit:

$$e_{MARE} = \frac{1}{N_D}\sum_{k=1}^{N_D}\left|\frac{y\left(k\right) - \hat{y}\left(k\right)}{y\left(k\right)}\right|, \tag{3}$$

where $\hat{y}\left(k\right)$ is the estimate based on equation (1) with respect to $\varphi\left(k\right)$ and $y\left(k\right) \neq 0$. The Levenberg-

Marquardt algorithm is then used to train the ANN unit by minimizing $e_{MARE}$ [15]. It starts by randomly generating the first two initial ANN weights, $w(0)$ and $w(1)$, at the 0-th and the 1-st iterations, where

$$w(0) = \left[ w_1(0), w_2(0),..., w_{n_h}(0), v_{11}(0), v_{12}(0),..., v_{1n_h}(0), v_{21}(0), v_{22}(0),..., v_{2n_h}(0), v_{n1}(0), v_{n2}(0),..., v_{nn_h}(0) \right]$$
and

$$w(1) = \left[ w_1(1), w_2(1),..., w_{n_h}(1), v_{11}(1), v_{12}(1),..., v_{1n_h}(1), v_{21}(1), v_{22}(1),..., v_{2n_h}(1), v_{n1}(1), v_{n2}(1),..., v_{nn_h}(1) \right],$$

respectively.

It then updates the ANN weights at the $(l+1)$-th iteration using the following formulation:

$$w(l+1) = w(l) + \left[ J^T(w) J(w) + \mu I \right]^{-1} J^T(w) R \qquad (4)$$

where $R = \left[ (y(1) - \hat{y}(1)) \quad y(2) - \hat{y}(2) \quad ... \quad y(N_D) - \hat{y}(N_D) \right]^T$. The details of the determination of the Jacobian matrix, $J(w)$, please refer to Hagan and Menhaj's early work [15].

## 2.2 The knowledge discovery unit

Although the air pollution level ($y$) can be estimated based on the air pollution factors ($x_n$) by using the ANN unit ($f$) formulated in equation (1), it is difficult to obtain explicit knowledge from $f$ solely based on equation (1), as it is an implicit or a black-box structure. Based on $f$, it is impossible to know which air pollution factor has more impact on the air pollution level and which air pollution factor can be ignored in analysis. In the KD unit illustrated in Figure 1, the main effect analysis [26] is proposed to determine the significant contribution of each air pollution factor to the air pollution level. It involves the study of the relative significance of variables of a system as in the traditional experimental design. It determines the relative significance of an air pollution factor by measuring the difference of the estimated air pollution level, when the value of the air pollution factor is changed from a level to another.

The significance of the air pollution factors, $x_1$, $x_2$, …, and $x_n$, are denoted by $d_1$, $d_2$, …, and $d_n$, respectively, where $d_i$, indicates the significance of the $i$-th air pollution factor, $x_i$, to the air pollution level, $y$. To determine $d_i$, $x_i$ is first quantized into $N_{div}$ levels based on equation (5):

$$x_{ij} = x_i^{\min} + \frac{x_i^{\max} - x_i^{\min}}{N_{div}} \cdot (j-1) \qquad (5)$$

where $j$=1,2,…,$N_{div}$; $x_i$ is within the range of $\left[ x_i^{\min}, x_i^{\max} \right]$; and $x_{ij}$ is the $j$-th quantization of $x_i$.

Then, the main effect of the $i$-th air pollution factor at level $j$, $\delta(x_{ij})$, can be calculated based on equation (6):

$$\delta(x_{ij}) = \frac{\sum_{k=1}^{N_s} f(y_k)}{N_s} \qquad (6)$$

where $y_k = \left[ y_k^1, y_k^2, ..., y_k^n \right]$ is a random sample with $k$=1, 2…,$N_s$; all $y_k^i = x_{ij}$; but all $y_k^m$ with $m \neq i$ are generated randomly within the range of the $m$-th air pollution factor, $x_m$, i.e. $x_m \in \left[ x_m^{\min}, x_m^{\max} \right]$. The numerator represents the total effect of $x_{ij}$, with respect to $f$, when all the other air pollution factors are in random states except the $i$-th air pollution factor which is a constant. Hence, the main effect of $x_{ij}$ with respect to $f$ can be estimated based on this set of random samples.

With $\delta(x_{ij})$, the significance, $d_i$, of the $i$-th air pollution factor, $x_i$, with respect to $f$ is defined by

$$d_i = \frac{\sum_{j=1}^{N_{div}} \left( \delta(x_{ij}) - \bar{\delta}_i \right)^2}{N_{div}}, \qquad (7)$$

where $\bar{\delta}_i = \frac{\sum_{j=1}^{N_{div}} \delta(x_{ij})}{N_{div}}$. $\qquad (8)$

8

The numerator in equation (7) represents the total change with respect to $f$, when $x_{ij}$ changes from a level to another level. When the total change is small, $x_{ij}$ is not significant to $f$. Hence, the outcome of $f$ changes slightly, even though $x_{ij}$ changes largely. In another extreme, if the total change is large, $x_{ij}$ is significant to $f$. Hence, the outcome of $f$ changes largely, even $x_{ij}$ only change slightly. We can consider a simple case with only two levels, 'high' and 'low' levels, where $N_s=2$. The significance of $x_{ij}$ can be evaluated by the difference between the main effect in 'low' level and the main effect in 'high' level. When the difference is small, the outcome of $f$ changes slightly whenever $x_{ij}$ is in 'low' or 'high' level. Hence, $x_{ij}$ is not too significant. Otherwise, when the difference is large, the significance of $x_{ij}$ is large.

The relative significance of $x_i$, is given by:

$$d'_i = \frac{d_i}{\sum_{j=1}^{n} d_i} \cdot 100\% . \tag{9}$$

where $d'_i$ indicates the relative significance of the $i$-th air pollution factor, $x_i$, with respect to the air pollution level, $y$, and $d'_i$ is relative to the total significance of all air pollution factors. The analysis of relative significances can be organized more efficiently, as $d'_i$ represents the rate of the $i$-th air pollution factor to the total significance of all air pollution factors.

## 3. A case study of ambient particulate matter concentrations

In order to illustrate the operation of the neural network based knowledge discovery system, data from a case study of ambient air monitoring was undertaken to estimate the air pollution level with respect to the ambient PM concentrations, which are important indicators of ambient air quality because of their detrimental effects on health and visibility impairment [33, 40]. Here, the concentrations of PM2.5 (with

aerodynamic diameter of PM ≤ 2.5 μm) and PM1.0 (with aerodynamic diameter of PM ≤ 1.0 μm) were studied because they are the indicators of fine and ultrafine particles that can enter the thorax and lower respiratory tract. Although a number of studies in the last decade have quantified and characterized PM2.5 in China [7, 11, 14, 16, 43, 44], there is relatively scarce research on traffic related PM1.0 in China [23, 35].

This case study was conducted near a busy road (Zhong He Viaduct) in the city centre of Hangzhou in 2010. The total number of newly registered on-road cars, buses and trucks in 2008 was 3.6 times greater than the numbers in 2000; the number of gasoline fueled vehicles increased 4.4 times and diesel fueled vehicles increased 1.2 times [18]. Zhong He Viaduct is a two-way vehicle only viaduct, with two lanes in each direction. The length of the viaduct is about 20 km from north to south with 10 exits on each side. In order to develop the models for estimating the concentrations of PMs at the roadside, the data on PM2.5 and PM1.0, as well as meteorological variables and the traffic flow were measured. The concentrations of PM2.5 and PM1.0 were measured by TSI DustTrak DRX Aerosol Monitor 8533, which was calibrated by the manufacturer to the respirable fraction of the standard ISO 12103, A1 Arizona road dust. The DustTrak DRX can simultaneously measure multiple size segregated mass fractions of the sampled aerosol, including PM2.5 and PM1.0 [39]. Two models namely $f_1$ and $f_2$, which estimate PM2.5 and PM1.0 at time $t$ namely $\hat{y}_1(t)$ and $\hat{y}_2(t)$, were developed based on equation (10) and (11) respectively:

$$\hat{y}_1(t) = f_1\left(x_1(t), x_2(t), x_3(t), x_4(t), x_5(t), y_1(t-T_s)\right) \tag{10}$$

$$\text{and } \hat{y}_2(t) = f_2\left(x_1(t), x_2(t), x_3(t), x_4(t), x_5(t), y_2(t-T_s)\right) \tag{11}$$

where $T_s$ is the sampling time; $x_1(t)$, $x_2(t)$, $x_3(t)$ and $x_4(t)$ are denoted as the meteorological variables temperature (ºC), relative humidity (RH, %), wind speed (m$^{-s}$) and barometric pressure (hPa), respectively, at time $t$. The meteorological data were collected by using TSI 9555A Advanced

Anemometer; $x_5(t)$ is denoted as the traffic flow at time $t$ and was measured by a real-time traffic surveillance system automatically counting the number of vehicles passing the surveillance point. This traffic flow data was provided by the Hangzhou City Traffic Control and Administration Center; $y_1(t-T_s)$ and $y_2(t-T_s)$ are PM2.5 and PM1.0 measured at time $t$, respectively. Previous studies [18, 19, 37] indicated that traffic flow was a significant predictor of particles from vehicle emissions, and those meteorological factors were also significant estimators in forecasting roadside atmospheric concentrations of submicron particles [45]. Hence, those air-pollution factors were selected in this study. All those data were collected from 14[th] May to 16[th] May from 7:30 am to 15:30 pm with a sampling interval time of 1 minute (i.e. $T_s = 1$). Hence, 24 hours of data or a total of 1440 pieces of data were collected.

Apart from using ANN unit, linear regression was used to develop models for $f_1$ and $f_2$, which are formulated in equations (10) and (11) respectively. To evaluate the performance of all these models, cross validation, namely repeated random sub-sampling validation, was carried out using the same data mentioned above. All those data were collected from 14[th] May to 16[th] May from 7:30 am to 15:30 pm with a sampling interval time of 1 minute (i.e. $T_s = 1$). Hence, 24 hours of data or a total of 1440 pieces of data were collected.

There is no particular rule for separating the data into training and validation data. Typically, a large portion of data is used for training and a small portion of data is used for validation, where the model was fitted to the training data and estimation accuracy was evaluated using the validation data. For each validation, the data was selected randomly and split into training and validation data in order to ensure that the characteristics of the training and validation data are unbiased thus to minimize the effects of data discrepancies. A total of 1260 pieces of data (21 hours of data) were selected randomly for training the models. The remaining 180 pieces of data (or 3 hours of data) were used for testing the
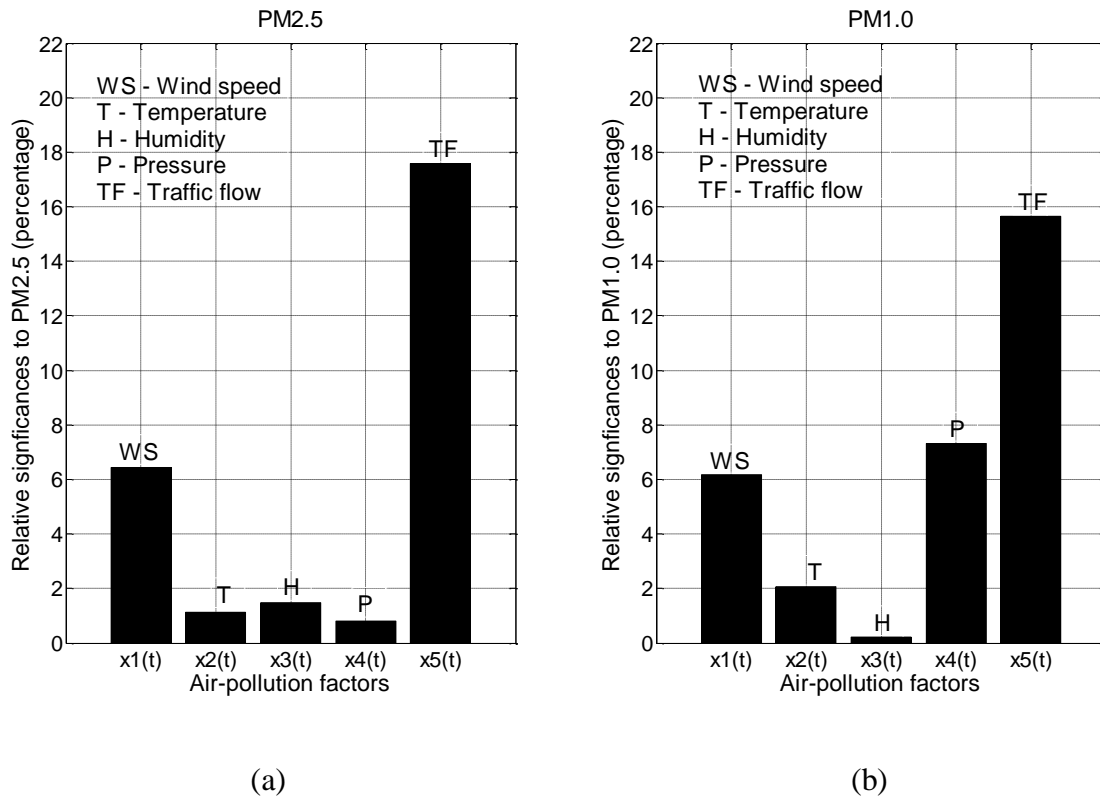
generalization capability of the models.

The cross validation results for both ANN unit and linear regression with respect to both PM2.5 and PM1.0 are shown in Table. 1. In the first validation test for PM2.5, linear regression resulted in 8.64 % of test error, while ANN unit only obtained 7.74 % of test error. Also the first validation test of PM1.0 showed that linear regression resulted in 7.88% of test error compared with 5.92% of that in ANN unit. Hence, the test errors obtained by ANN unit are smaller than those obtained by linear regression in both PM concentrations. Here the number of cross validations was not initially pre-defined, as there is no particular rule for setting the number of cross validations. In general, if more cross validations are conducted, more convincing conclusion can be reached. Hence, we kept performing the cross validations until a convincing conclusion can be made. We observed that the validation results obtained by ANN were generally better than those obtained by the linear regression model, after performing 15 cross-validations. Therefore, we stop performing cross validation at the 15th cross-validation. Table 1 further shows that the models developed by the ANN unit can yield smaller validation errors than those by the linear regression models for both PM2.5 (9.69% vs. 7.93%, $t$ value = 8.83) and PM1.0 (10.24% vs. 7.85%, $t$ value = 8.18). Therefore, the generalization capability of the models developed by the ANN unit in predicting both PM2.5 and PM1.0 levels is significantly better than that of the statistical regression method.

**Table 1** Cross validations for PM2.5 and PM1.0

| Validation number | PM2.5 | | PM1.0 | |
|---|---|---|---|---|
| | Linear regression | ANN | Linear regression | ANN |
| 1 | 8.64 | 7.74 | 7.88 | 5.92 |
| 2 | 6.02 | 5.16 | 14.96 | 13.45 |
| 3 | 10.85 | 9.91 | 9.74 | 6.71 |
| 4 | 8.71 | 6.47 | 16.00 | 13.39 |
| 5 | 7.13 | 4.26 | 10.50 | 8.27 |
| 6 | 11.41 | 7.81 | 11.29 | 8.91 |
| 7 | 10.62 | 9.04 | 6.87 | 4.86 |
| 8 | 7.89 | 6.19 | 7.86 | 6.37 |
| 9 | 11.34 | 10.34 | 7.71 | 4.02 |
| 10 | 7.67 | 5.78 | 13.9 | 11.34 |
| 11 | 10.02 | 7.78 | 7.93 | 6.78 |
| 12 | 8.34 | 6.41 | 8.72 | 4.18 |
| 13 | 9.70 | 8.02 | 16.25 | 11.76 |
| 14 | 19.34 | 18.20 | 7.95 | 6.70 |
| 15 | 7.69 | 5.84 | 6.02 | 5.04 |
| Mean ± SD | 9.69 ± 3.12 | 7.93 ± 3.32 | 10.24 ± 3.44 | 7.85 ± 3.22 |
| $t$ value | 8.83 | | 8.18 | |

Based on the KD unit, the main effect of each variable in $f_1$ (i.e. PM2.5) and $f_2$ (i.e. PM1.0) can be determined and are shown in Fig. 2(a) and 2(b), respectively. The results indicate the relative significance of each air pollution factor with respect to the concentrations of PM2.5 and PM1.0. Compared with meteorological factors (wind speed, $x_1(t)$, temperature, $x_2(t)$, relative humidity, $x_3(t)$ and barometric pressure, $x_4(t)$), traffic flow, $x_5$ $(t)$, provides more contribution to the estimated concentrations of PM2.5 and PM1.0. This is unsurprising, as the increase of the total number of vehicles passing through the area is associated with increased emission of PM2.5 and PM1.0. In addition, the four measured meteorological factors also have influence on the concentrations of PM2.5 and PM1.0,

13

but at a less significant extent compared with that from the traffic flow. Results from Fig. 2 also indicate that wind speed plays more important role than other meteorological factors on PM2.5 concentration, but barometric pressure, exceed wind speed, impacts more on PM1.0 concentration. In summary, this KD unit overcomes the limitation of current neural network approaches whereby no explicit information can be indicated within the neural networks.



(a)                                                          (b)

**Fig. 2** Relative significance of air pollution factor contributing to PM2.5 (a) and PM1.0 (b)

Then, a hypothesis test was used to evaluate the significance of the air pollution factors based on the result of linear regression. The $t$-values are used to indicate whether the air pollution factors are significant or not. When the $t$-value of the corresponding air pollution factor is less than 2.09, this air pollution factor is insignificant to the air pollution level with 98% of confidence level. Otherwise, this air pollution factor is significant. The $t$-values of the air pollution factors with respect to emissions of

PM2.5 and PM1.0 are shown in Table 2. The *t*-values are bolded, when they are large than 2.09. For PM2.5, the *t*-values for wind speed and traffic flow are large than 2.09. Hence, it indicates that wind speed and traffic flow are significant factors to PM2.5. For PM1.0, it also indicates that wind speed and traffic flow are the two significant factors. The *t*-value of barometric pressure is 2.00 which is near to the significant level. Hence, barometric pressure is somewhat important to PM1.0 compared with the other insignificant weather factors. These analysis of hypothesis tests is similar to those obtain by KD, where both wind speed and traffic flow are significant to PM2.5 and PM1.0, and barometric pressure is relatively significance to PM1.0 compared with the other insignificant pollution factors.

**Table 2** T-values of pollution factors obtained by linear regression models

| | Wind speed $x_1(t)$ | Temperature $x_2(t)$ | Relative humidity $x_3(t)$ | Barometric pressure $x_4(t)$ | Traffic flow $x_5(t)$ |
|---|---|---|---|---|---|
| **PM2.5** | **2.91** | 0.74 | 0.63 | 0.99 | **4.11** |
| **PM1.0** | **2.66** | 1.14 | 1.55 | <u>2.00</u> | **5.42** |

## 4. Conclusions

In this paper, a neural network based knowledge discovery system has been developed to estimate air pollution levels based on a set of measured air monitoring data. Cross validations and the case study results have demonstrated that this new system is able to overcome the limitation of traditional ANNs and generate accurate explicit knowledge of the significant contribution of each pollution factor on air pollution levels (PM2.5 and PM1.0) within the existing ANNs for estimating air pollution levels. In another word, the system is able to estimate PM2.5 and PM1.0 concentrations based on traffic flow, meteorological conditions at a busy traffic roadside, and past measured PM concentrations. Based on this explicit knowledge, researchers can gain a better insight into the significance and influence of a

particular air pollution factor on the mass concentrations of PMs. Thus, the system has potential application value in planning future air monitoring programs to achieve cost-effective outcomes. For the future work, we will develop a rule discovery system [31, 32] which extracts symbolic rules from neural networks, in order to illustrate relations between air-pollution factors and air-pollution levels. As we only used the simple neural network with one-hidden layer in this research, we will enhance the effectiveness of the neural network by integrating the mechanism of hybrid approaches in our next project [28].

## Acknowledgements

## References

1. B.R. Bakshi, U. Utojo, A common framework for the unification of neural, chemometric and statistical modeling methods, Analytica Chimica Acta 384(3) (1999) 227-247.
2. K.Y. Chan, S.H. Ling, T.S. Dillon, H.T. Nguyen, Diagnosis of hypoglycemic episodes using a neural network based rule discovery system, Expert Systems with Applications 38(8) (2011) 9799-9808.
3. H. Choi, D.S. Choi, S.M. Choi, Meteorological Condition and atmospheric Boundary Layer Influenced upon Temporal Concentrations of PM1, PM2.5 at a Coastal City, Korea for Yellow Sand Event from Gobi Desert, Disaster Advances 3(4) (2010) 309-315.
4. W.G. Cobourn, L. Dolcine, M. French M.C. Hubbard, A comparison of nonlinear regression and neural network models for ground-level ozone forecasting, Journal of the Air & Waste Management Association 50(11) (2000) 1999-2009.
5. B. Croxford, A. Penn, B. Hillier, Spatial distribution of urban pollution: Civilizing urban traffic. Science of the Total Environment, 190 (1996) 3-9.
6. B.M.F. de Castro, J.M.P. Sanchez, W.G. Manteiga, M.F. Bande, J.L.B. Cela, J.J.H. Fernandez, Prediction of SO2 levels using neural networks, Journal of the Air & Waste Management Association 53(5) (2003) 532-539.

7. F.K. Duan, K.B. He, Y.L. Ma, F.M. Yang, X.C. Yu, S.H. Cadle, T. Chan, P.A. Mulawa, Concentration and chemical characteristics of PM2.5 in Beijing, China: 2001-2002. Science of the Total Environment 355(1-3) (2006) 264-275.

8. A.L. Dutot, J. Rynkiewicz, F.E. Steiner, J. Rude, A 24-h forecast of ozone peaks and exceedance levels using neural classifiers and weather predictions. Environmental Modelling & Software 22(9) (2007) 1261-1269.

9. S.T. Ebelt, W.E. Wilson, M. Brauer, Exposure to ambient and nonambient components of particulate matter: a comparison of health effects, Epidemiology 16(3) (2005) 396-405.

10. EPA, Glossary of Climate Change Terms, 2009.

11. M. Fang, C.K. Chan, X.H. Yao, Managing air quality in a rapidly developing nation: China, Atmospheric Environment 43(1) (2009) 79-86.

12. M.W. Gardner, S.R. Dorling, Artificial neural networks (the multilayer perceptron) - A review of applications in the atmospheric sciences, Atmospheric Environment 32(14-15) (1998) 2627-2636.

13. G. Grivas, A. Chaloulakou, Artificial neural network models for prediction of PM10 hourly concentrations, in the Greater Area of Athens, Greece, Atmospheric Environment 40(7) (2006) 1216-1229.

14. Z.P. Gu, J.L. Feng, W.L. Han, L. Li, M.H. Wu, J.M. Fu, G.Y. Sheng, Diurnal variations of polycyclic aromatic hydrocarbons associated with PM(2.5) in Shanghai, China, Journal of Environmental Sciences-China 22(3) (2010) 389-396.

15. M.T. Hagan, M.B. Menhaj, Training Feedforward Networks with the Marquardt Algorithm, IEEE Transactions on Neural Networks 5(6) (1998) 989-993.

16. L.Y. He, M. Hu, Y.H. Zhang, X.F. Huang, T.T. Yao, Fine particle emissions from on-road vehicles in the Zhujiang Tunnel, China. Environ Sci Technol 42(12) (2008) 4461-4466.

17. T.W. Hesterberg, C.M. Long, W.B. Bunn, S.N. Sax, C.A. Lapin, P.A. Valberg, Non-cancer health effects of diesel exhaust: a critical assessment of recent human and animal toxicological literature. Crit Rev Toxicol 39(3) (2009) 195-227.

18. L. Jian, Y.P Zhu, Y. Zhao, Monitoring fine and ultrafine particles in the atmosphere of a Southeast Chinese city, Journal of Environmental Monitoring 13(9) (2011) 2623-2629.

19. S. Kaur, M.J. Nieuwenhuijsen, Determinants of personal exposure to PM2.5, ultrafine particle counts, and CO in a transport microenvironment, Environmental Science and Technology 43 (13) (2009) 4737-4743.

20. N. Kunzli, R. Kaiser, S. Medina, M. Studnicka, O. Chanel, P. Filliger, M. Herry, F. Horak, V. Jr. Puybonnieux-Texier, P. Quenel, J. Schneider, R. Seethaler, J.C. Vergnaud, H. Sommer, Public-health impact of outdoor and traffic-related air pollution: a European assessment. Lancet 356(9232) (2000) 795-801.

21. S. Larssen, M.L. Adams, K.J. Barrett, M. Bolscher, F. de Leeuw, T. Pulles, Air pollution in Europe 1990–2000, European Environment Agency Copenhagen, Topic report 4/2003.

22. H.K. Lam, F.H.F. Leung, Design and training for combinational neural-logic systems, IEEE Transactions on Industrial Electronics 54(1) (2007) 612-619.

23. S.C. Lee, Y. Cheng, K.F. Ho, J.J. Cao, P.K.K. Louie, J.C. Chow, J.G. Watson, PM1.0 and PM2.5 characteristics in the roadside environment of Hong Kong, Aerosol Science and Technology 40(3) (2006) 157-165.

24. K.L. Maier, F. Alessandrini, I. Beck-Speier, T.P.J. Hofer, S. Diabate, E. Bitterle, T. Stoger, T. Jakob, H. Behrendt, M. Horsch, J. Beckers, A. Ziesenis, L. Hultner, M. Frankenberger, Health effects of ambient particulate matter - Biological mechanisms and inflammatory responses to in vitro and in vivo particle exposures, Inhalation Toxicology 20(3) (2008) 319-337.

25. S. Mikkonen, H. Korhonen, S. Romakkaniemi, J.N. Smith, J. Joutsensaari, K.E.J. Lehtinen, A. Hamed, T.J. Breider, W. Birmili, G. Spindler, C. Plass-Duelmer, M.C. Facchini, A. Laaksonen, Meteorological and trace gas factors affecting the number concentration of atmospheric Aitken (D-p=50 nm) particles in the continental boundary layer: parameterization using a multivariate mixed effects model, Geoscientific Model Development 4(1) (2011) 1-13.
26. D.C. Montgomery, Design and Analysis of Experiments John Wiley and Sons, Inc., New York, 1997.
27. S.M.S. Nagendra, M. Khare, Artificial neural network based line source models for vehicular exhaust emission predictions of an urban roadway, Transportation Research Part D-Transport and Environment 9(3) (2004) 199-208.
28. M. Negoita, D. Neagu, V. Palade, Computational Intelligence: Engineering of Hybrid Systems, Springer Verlag. 2005.
29. H. Niska, T. Hiltunen, A. Karppinen, J. Ruuskanen, M. Kolehmainen, Evolving the neural network model for forecasting air pollution time series, Engineering Applications of Artificial Intelligence 17(2) (2004) 159-167.
30. J.D. Olden, D.A. Jackson, Illuminating the "black box": a randomization approach for understanding variable contributions in artificial neural networks, Ecological Modelling 154(1-2) (2002) 135-150.
31. V. Palade, D. Neagu, R.J. Patton, Interpretation of trained neural networks by rule extraction, Computational Intelligence: Theory and Applications 2206 (2002) 152–161.
32. C.D. Neagu, N. Avouris, E. Kalapanidas, V. Palade, Neural and neuro-fuzzy integration in a knowledge-based system for air quality prediction, Applied Intelligence 17 (2002) 141-169.
33. C.A. Pope, D.W. Dockery, Health effects of fine particulate air pollution: Lines that connect. Journal of the Air & Waste Management Association 56(6) (2006) 709-742.
34. G.A.F Seber, Linear Regression Analysis. , 2nd ed. Hoboken, N.J. : Wiley-Interscience, 2003.
35. Z.X. Shen, Y.M. Han, J.J. Cao, J. Tian, C.S. Zhu, S.X. Liu, P.P. Liu, Y.Q. Wang, Characteristics of Traffic-related Emissions: A Case Study in Roadside Ambient Air over Xi'an, China. Aerosol and Air Quality Research 10(3) (2010) 292-300.
36. T. Slini, A. Kaprara, K. Karatzas, N. Moussiopoulos, PM(10) forecasting for Thessaloniki, Greece, Environmental Modelling & Software 21(4) (2006) 559-565.
37. U.W. Tang, Z. Wang, Determining gaseous emission factors and driver's particle exposures during traffic congestion by vehicle-following measurement techniques. Journal of Air and Waste Management Association 56 (11) (2006) 1532-9.
38. L.T. Wang, C. Jang, Y. Zhang, K. Wang, Q.A. Zhang, D. Streets, J. Fu, Y. Lei, J. Schreifels, K.B. He, J.M. Hao, Y.F. Lam, J. Lin, N. Meskhidze, S. Voorhees, D. Evarts, S. Phillips, Assessment of air quality benefits from national air pollution control policies in China. Part I: Background, emission scenarios and evaluation of meteorological predictions, Atmospheric Environment 44(28) (2011) 3442-3448.
39. X.L. Wang, G. Chancellor, J. Evenstad, J.E. Farnsworth, A. Hase, G.M. Olson, A. Sreenath, J.K. Agarwal, A Novel Optical Instrument for Estimating Size Segregated Aerosol Mass Concentration in Real Time, Aerosol Science and Technology 43(9) (2009) 939-950.
40. J.G. Watson, Visibility: Science and regulation, Journal of the Air & Waste Management Association 52(6) (2002) 628-713.
41. T.C. Wong, A.H.S. Chan, A study of the impact of different direction-of-motion stereotypes on response time and response accuracy using neural network, IEEE Transactions on Systems, Man, and Cybernetics-Part A: Systems and Humans (2012) (In Press)
42. T.C. Wong, K.M.Y. Law, H.K. Yau, S.C. Ngan, Analyzing supply chain operation models with the

PC-algorithm and the neural network, Expert Systems with Applications 38(6) (2011) 7526-7534.

43. S.J. Wu, Q. Feng, Y. Du, X.D. Li, Artificial Neural Network Models for Daily PM(10) Air Pollution Index Prediction in the Urban Area of Wuhan, China. Environmental Engineering Science 28(5) (2011) 357-363.

44. F.M. Yang, B.M. Ye, K.B. He, Y.L. Ma, S.H. Cadle, T. Chan, P.A. Mulawa, Characterization of atmospheric mineral components of PM2.5 in Beijing and Shanghai, China, Science of the Total Environment 343(1-3) (2005) 221-230.

45. L. Jian, Y. Zhao, Y.P. Zhu, M.B. Zhang, D. Bertolatti, An application of ARIMA model to predict submicron particle concentrations from meteorological factors at a busy roadside in Hangzhou, China, The Science of the total environment, 426 (2012) 336-345.