# A Prediction Comparison for PM$_{2.5}$ Between Neural Network and Multiple Regression Models in Rabigh, Saudi Arabia

## Issam Mohammed Aquil Alghanmi[1], Ibrahim Abdelaziz Al-Darrab[2], Osman Imam Taylan[3], Omar Seraj Aburizaiza[4]

[1]Department of Industrial Engineering, King Abdulaziz University, Jeddah, Saudi Arabia

[2, 3]Professor, Department of Industrial Engineering, King Abdulaziz University, Jeddah, Saudi Arabia

[4]Professor, Unit for Ain Zubaida Rehabilitation & Ground Water Research, King Abdulaziz University, Jeddah, Saudi Arabia

**Abstract:** *A high concentration of fine particulate in the atmosphere has negative consequences for human health and wellbeing. Therefore, the prediction of the concentration of particles in atmospheric air is imperative so that the public is well aware of the atmospheric condition, and the standard of air quality can be properly managed. The study explores the feasibility of using neural network methods in replacement of the generally-used statistical models for prediction of the daily average concentration of PM$_{2.5}$ (particulate matter having diameter $\leq$ 2.5 um). 24-h PM$_{2.5}$ observation from May 6$^{th}$to June 17$^{th}$, 2013, at a specific spot in Rabigh city revealed high chronological changes with an average of (36.97 ± 16.22 ug/m$^3$). The results showed that the concentration surpassed the limit specified by (WHO) guideline (25 ug/m$^3$). Nine toxic Trace Elements (TEs) that are dangerous for human health were considered in this study, including (V, S, Lu, Ni, Cl, Zn, Cu, Pb, and Cr). These trace elements were found in abundance in PM$_{2.5}$ (ug/m$^3$). These trace elements were used as input and served as a basis for the formulation of NN models and (MLR) models. The study drew a contrast between the two models was found to be (2.017)-(10.596). The result showed that properly formulated and trained ANNs are effective in resolving the issues associated with for prediction cast of particulate pollution.*

**Keywords:** PM$_{2.5}$, Trace Elements, ANNs, MLR, Rabigh

## 1. Introduction

From the 1990s has been adequate proof that particulate matter (PM) is a severe threat to human health even in comparatively tiny concentrations in ambient air [1]. The World Health Organization (WHO) announced that the standards for particulate matter that $\leq$ 2.5 um (PM$_{2.5}$) In 24-hour must not be exceeded (25 ug/m$^3$) per day [2]. The correct predictions of PM$_{2.5}$ concentration levels at the right time will prevent the population from PM and would support the actions taken to find a better solution [1]. The concentration of particle matter can be defined by utilizing a direct simulation of all its related processes, physical or chemical requires an examination of a large number of parameters that can describe aerosol generation, formation, transport, and removal in the atmosphere; also that will lead to complex and higher attempts. Besides, statistical analysis is the tool utilized over and over the modeling purpose [1].

## 2. Inclusion of ANNs

Being deeply rooted in the atmospheric pollution predicting sector, ANNs are considered multiple regression models' alternative competitors. This paper's objective is to examine the performance of ANNs compared to the performance of multiple regression based on the concentration prediction of PM$_{2.5}$ in Rabigh. In addition, the paper will investigate the capability of the models in predicting events exceedance, acknowledging them as essential methods at operational levels for authorities. At least for the research area studied, ANNs can handel a quantifiable improvement over predictions derived from regression [1].

There have been various researches conducted on the application of Artificial Neural Networks (ANNs) for the prediction of the concentration of pollutants in atmospheric air. This was first explored by Boznar, et al. [1]. Later contributions were made in the matter by Yi and Prybutok, Comrie, and Spellman [1]. They investigated the application of ANN for daily prediction of ozone (O$_3$) concentration. The experts also contrasted the ANN with the conventional multiple regression models. ANN application in the prediction of concentration of nitrogen oxides (NOx)-nitrogen dioxide (NO$_2$) and O$_3$ in the atmosphere was investigated by Gardner and Dorling, who came to conclude that ANNs gave more accurate predictions as compared to linear models [1]. Hence they recommended that the properly formulated and trained neural networks give equivalent or better results than the linear model provided the same data is used by both models. This conclusion was agreed upon by Chaloulakou et al. [1]. Who supported the results for a prediction made for ozone in Athens, Greece.

Moreover, Chelani et al. [1] also supported the results with reference to SO$_2$ in New Delhi, India. Limited use of ANNs was seen with respect to particulate pollutants. ANNs were applied for a 1-hour prediction of PM$_{2.5}$ concentration in Santiago, Chile, by Perez et al. [1]; ANNs were applied by Kohlemainen et al. [1] for prediction the daily average and maximum concentration of PM$_{10}$ in Kuopio, Finland; and lastly, they were applied by Chelani et al. [1]. For the prediction of the daily average concentration of PM$_{10}$ in Jaipur, India. Models developed on the basis of neural networks were also employed by Lu et al. [1] for the prediction of respirable suspended particulate (RSP) in Hong

Kong. The neural networks were found to deliver equal or slightly better results as compared to the regression models for prediction of daily average and daily maximum concentration of $PM_{10}$ and $PM_{2.5}$.

## 3. Data and Procedures

### Utilized Data and Research Area

The Saudi Arabian city of Rabigh is located in the province of Makkah. The city touches the Red Sea at its eastern coast and lies on the Tropic of Cancer amide $22^{nd}$ and $23^{rd}$ latitudes North of Equator. As per Shedrack et al. [3], the population of Rabigh is nearly 180, 352 people in 2014. The city has exceptionally hot weather in summers while a warm one in the winter season. The humidity is on the rise in summers with sudden short rains. There is not much precipitation during the rest of the year. The city sees an elevation of temperatures that starts from the month of April and continuously rises until it goes beyond 45°C during the months from July to September. Rabigh is suffering from unhealthy air as it contains large numbers of factories in the South and Southwest while the city and the adjacent areas generally in the North, as shown in (Fig.1) [3].

The daily $PM_{2.5}$ samples were collected in a previous study of Rabigh [3], on 2 µm Whatman filters pore-size PTFE 46.2 mm were supported by weighed in advance polypropylene ring, which was sequentially numbered, when a sampling pump using low volume air was used. Installed in Rabigh which is a fixed site and for a period between $6^{th}$ May to $17^{th}$ June 2013, the $PM_{2.5}$ sampler was equipped with the following; a housing unit, a gooseneck, a power supply, rubber stopper with an inner diameter of 5.27 cm, filter holder with a 47 mm diameter, data logger, mass flow meter, an air volume totalizer, a flow controlled and elapsed time indicating pump, and a 16.67 L min-1 flow rate, which was optimal for sampling $PM_{2.5}$, operating aluminum cyclone separator which had a 2.5 µm cut size. To ensure that ambient $PM_{2.5}$ representation is well obtained, a height of 3-5 meters above the ground level was ensured when fixing the sampler inlets. The height also ensured that the sampler inlets avoided dust for effectiveness. Furthermore, the Trace Elements in $PM_{2.5}$ samples were analyzed by a Thermo Scientific ARL QUANT'X energy dispersive X-ray fluorescence spectrometer (ED–XRF) (model AN41903–E 06/07C, Ecublens Switzerland) using six secondary fluorescers (Si, Ti, Fe, Cd, Se, and Pb). In this study, nine trace elements are chosen (Sulfur, Chlorine, Vanadium, Chlorine, Nickel, Copper, Zink, Lutetium, lead).

## 4. Variables for Prediction

The Nine predictor variables were chosen from all trace elements in $PM_{2.5}$ (ug/m$^3$) as they havea higher concentration in $PM_{2.5}$ (ug/m$^3$), and they also, are toxic to human health.



**Figure 1:** Location of the $PM_{2.5}$ sampling site and the industrial area in Rabigh, Saudi Arabia
Source: Adapted from [3].

## 5. Results and Discussion

**Multiple linear regressions**

In multiple linear regression (MLR), a linear combination of two or more predictor variables is used to explain the variation in response. In order to calculate the relationship between n input variables (x) and the target variable (y) we could use the linear equation:

$$y = a_0 + a_1x_1 + a_2x_2 + \cdots + a_nx_n$$

or

$$y = a_0 + \sum_{i=1}^{n} a_i x_i + \varepsilon$$

in our research, we used the multiple linear regression to identify the most effective elements to $PM_{2.5}$, where the variables defined as follow:

The dependent variable: $y = PM_{2.5}$
The independent variables:
($x_1 = S$, $x_2 = Cl$, $x_3 = Cr$, $X_4 = Ni$, $x_5 = Zn$, $X_6 = Cu$, $x_7 = Zn$, $x_8 = Lu$, $x_9 = Pb$)

First step: we calculate the correlation coefficient between $PM_{2.5}$ and the ninetrace elements.

**Table 1:** Correlation between $PM_{2.5}$ (ug/m$^3$) and 9 independent variables

| $PM_{2.5}$ (ug/m$^3$) | S-(ug/m$^3$) | Cl-(ug/m$^3$) | V-(ug/m$^3$) | Cr-(ug/m$^3$) | Ni-(ug/m$^3$) | Cu-(ug/m$^3$) | Zn-(ug/m$^3$) | Lu-(ug/m$^3$) | Pb-(ug/m$^3$) |
|---|---|---|---|---|---|---|---|---|---|
| r-value | 0.062 | 0.137 | -0.264 | 0.867 | 0.627 | 0.75 | 0.417 | 0.755 | 0.23 |
| p-value | 0.705 | 0.399 | 0.1 | 0.00 | 0.00 | 0.00 | 0.01 | 0.00 | 0.89 |
| Conclusion | NoRelation | NoRelation | NoRelation | Relation | Relation | Relation | Relationbutnotlinear | Relation | NoRelation |

As shown in Table 1, four variables have a linear relation to $PM_{2.5}$. In contrast, four variables have no relation to $PM_{2.5}$, and one variable has relation but not linear.

**Second step:** we made Regression between the nine elements and $PM_{2.5}$ (ug/m$^3$)

**Table 2:** Regression between the 9 elements & $PM_{2.5}$ (ug/m$^3$)

| Source | DF | Adj SS | Adj MS | F-Value | P-Value |
|---|---|---|---|---|---|
| Regression | 9 | 9354.5 | 1039.39 | 34.48 | 0.000 |
| S-(ug/m$^3$) | 1 | 89.2 | 89.25 | 2.96 | 0.096 |
| Cl-(ug/m$^3$) | 1 | 0.0 | 0.00 | 0.00 | 0.996 |
| V-(ug/m$^3$) | 1 | 9.6 | 9.59 | 0.32 | 0.577 |
| Cr-(ug/m$^3$) | 1 | 311.2 | 311.23 | 10.32 | 0.003 |
| Ni-(ug/m$^3$) | 1 | 63.2 | 63.22 | 2.10 | 0.158 |
| Cu-(ug/m$^3$) | 1 | 107.0 | 106.95 | 3.55 | 0.069 |
| Zn-(ug/m$^3$) | 1 | 4.5 | 4.47 | 0.15 | 0.703 |
| Lu-(ug/m$^3$) | 1 | 31.2 | 31.20 | 1.03 | 0.317 |
| Pb-(ug/m$^3$) | 1 | 129.8 | 129.76 | 4.30 | 0.047 |
| Error | 30 | 904.4 | 30.15 | | |
| Total | 39 | 10258.9 | | | |

As shown in the table, 2the independent variables together have a significant effect on $PM_{2.5}$; in contrast, some variables have no significant impact.

**Table 3:** Showed the Model Summary of the nine independent variables

| S | R-sq | R-sq (adj) | R-sq (pred) |
|---|---|---|---|
| 5.49056 | 91.18% | 88.54% | 84.17% |

**Table 4:** Showed the Coefficients of the 9 independent variables and VIF

| Term | Coef | SE Coef | T-Value | P-Value | VIF |
|---|---|---|---|---|---|
| Constant | -0.27 | 6.81 | -0.04 | 0.969 | |
| S-(ug/m$^3$) | 2.10 | 1.22 | 1.72 | 0.096 | 4.01 |
| Cl-(ug/m$^3$) | -0.01 | 1.12 | -0.01 | 0.996 | 1.83 |
| V-(ug/m$^3$) | -213 | 377 | -0.56 | 0.577 | 16.41 |
| Cr (ug/m$^3$) | 4050 | 1260 | 3.21 | 0.003 | 18.90 |
| Ni-(ug/m$^3$) | 2037 | 1407 | 1.45 | 0.158 | 16.39 |
| Cu-(ug/m$^3$) | -1554 | 825 | -1.88 | 0.069 | 5.77 |
| Zn-(ug/m$^3$) | 200 | 519 | 0.39 | 0.703 | 11.49 |
| Lu-(ug/m$^3$) | 1380 | 1356 | 1.02 | 0.317 | 16.98 |

| | | | | | |
|---|---|---|---|---|---|
| Pb-(ug/m$^3$) | -854 | 412 | -2.07 | 0.047 | 3.56 |

As shown in Table 4, some values of VIF that greater than ten meaning there is multicollinearity in the model.

**Third step:**

To diagnosis and solve the problem of multicollinearity in the model, we will begin to remove the variables causing the multicollinearity gradually and re-estimate the model again without these variables where the steps as follows:

1) The elements of Vanadium and Cooper were removed because they have the highest value of the Variance Inflation Factor, and they have no significant impact on $PM_{2.5}$.
2) After the two elements removed in the first step and re-estimate the model, we found that there is a multicollinearity problem by another two elements they are zinc and luteum. By considering the P-Value, the two elements have no impact on the $PM_{2.5}$, also, by comparison, the correlation coefficient of the two variables with $PM_{2.5}$ found that Zinc has less correlation coefficient to $PM_{2.5}$, for that the Zinc element was removed.
3) After removing the Zinc element and re-estimate the model, we found that there is one element causing multicollinearity, that is Nickel element, for that Nickel element was removed.
4) After the Nickel element removed, the model containing five elements was estimated. We found that this model has no multicollinearity.

**Fourth step:**

By using five remaining elements, we will determine the best model by considering the coefficient of determination

**Table 5:** Showing the Response to $PM_{2.5}$ (μg/m$^3$):

| Vars | R-sq | R-sq (ADJ) | R-sq (Pred) | Mallows CP | Std Dev | S | Cl | Cr | Lu | Pb |
|---|---|---|---|---|---|---|---|---|---|---|
| 1 | 75.1 | 74.5 | 72.7 | 45.9 | 8.1917 | X | | | | |
| 1 | 57 | 55.9 | 50.4 | 105.7 | 10.773 | | | X | | |

| 2 | 86.7 | 86 | 84 | 9.7 | 6.0643 | X | | X | | |
| 2 | 82.4 | 81.4 | 79.2 | 24.1 | 6.9925 | | | | | X |
| 3 | 87.3 | 86.2 | 84.1 | 9.9 | 6.0154 | X | X | X | | |
| 3 | 87.2 | 86.2 | 84 | 10.1 | 6.0321 | X | X | X | | |
| 4 | 89.6 | 88.4 | 86.4 | 4.4 | 5.5297 | X | | X | X | X |
| 4 | 87.3 | 85.9 | 83.1 | 11.8 | 6.0998 | X | X | X | X | |
| 5 | 89.7 | 88.2 | 85.9 | 6 | 5.579 | X | X | X | X | X |

As shown in table 5 the best model is the model that containing [S-(ug/m$^3$), Cr (ug/m$^3$), Lu-(ug/m$^3$), Pb-(ug/m$^3$) ]. Because it has the highest value of the adjusted coefficient of determination is (88.4).

**Table 6:** Showing the Analysis of Variance of the Four Independent Variables

| Source | DF | Adj SS | Adj MS | F-Value | P-Value |
|---|---|---|---|---|---|
| Regression | 4 | 9188.7 | 2297.18 | 75.13 | 0.000 |
| S-ug/m$^3$ | 1 | 624.7 | 624.68 | 20.43 | 0.000 |
| Cr (ug/m$^3$) | 1 | 1501.9 | 1501.89 | 49.12 | 0.000 |
| Lu-(ug/m$^3$) | 1 | 232.5 | 232.48 | 7.60 | 0.009 |
| Pb-(ug/m$^3$) | 1 | 239.7 | 239.70 | 7.84 | 0.008 |
| Error | 35 | 1070.2 | 30.58 | | |
| Total | 39 | 10258.9 | | | |

As showing in table 6, the independent variables together have a significant effect on $PM_{2.5}$ and each individual variable also has significant impact on $PM_{2.5}$.

**Table 7:** Showing the Model Summary of the Four Independent Variables

| Ssq (pred) | R-sq | R-sq (adj) | R- |
|---|---|---|---|
| 5.52965 | 89.57% | 88.38% | 86.38% |

**Table 8:** Showing Coefficients of the Four Independent Variables

| Term | Coef | SE Coef | T-Value | P-Value | VIF |
|---|---|---|---|---|---|
| Constant | 0.50 | 2.70 | 0.18 | 0.855 | |
| S-(ug/m$^3$) | 3.828 | 0.847 | 4.52 | 0.000 | 1.90 |
| Cr-(ug/m$^3$) | 3942 | 563 | 7.01 | 0.000 | 3.71 |
| Lu-(ug/m$^3$) | 1941 | 704 | 2.76 | 0.009 | 4.51 |
| Pb-(ug/m$^3$) | -976 | 349 | -2.80 | 0.008 | 2.52 |

**Finally calculating the best Regression Equation with most effective variables**

$$PM_{2.5} \ (\mu g/m^3) = [0.50 + 3.828 \ S\text{-}(ug/m^3) + 3942 \ Cr \ (ug/m^3) + 1941 \ Lu\text{-}(ug/m^3) - 976 \ Pb\text{-}(ug/m^3) \ ]$$
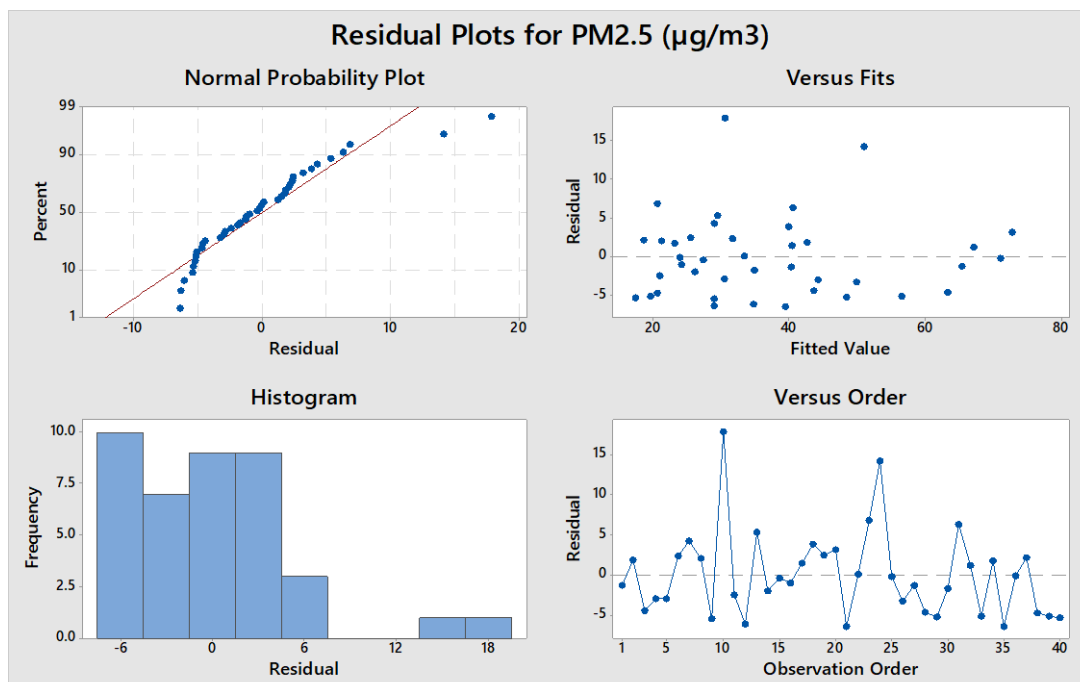


**Figure 2:** Residual Plots for $PM_{2.5}$ (ug/m$^3$)

Figure 2 showing the Residual plots that validate the assumption of regression that residuals are normally plotted and independent.

**Artificial neural networks**

In this study, a parallel evaluation investigates the prediction performance of multiple linear regression models and neural networks, implying similar processes, inputs, and data when developing and comparing both methods.

An Artificial Neuron Network (ANN), often known as a Neural Network, is a computer model that is based on the structure and functionality of biological neural networks [1].

In terms of Computer Science, it is similar to an artificial human nervous system for accepting, processing, and transferring data. In a neural network, there are three layers in total:

1) The input layer (This layer receives all of the inputs and feeds them into the model).

2) Output Layer (They maybe several hidden layers that process the information from the input layers).

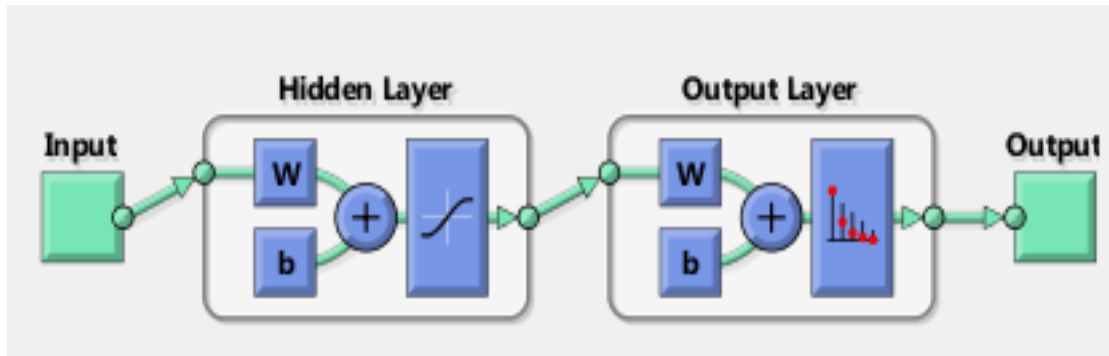3) The Output Layer (The data is made accessible at the output layer once it has been processed).



**Figure 3:** A multi-layered artificial neural network

In this study Backpropagation method used in artificial neural networks to calculate a gradient that is needed in the calculation of the weights to be used in the network asfollows:

**Backpropagation algorithm**:

The error of the network is defined by $C = \frac{1}{2}(y_n - t)^2$

The error gradient of the input vector at a layer (n) is defined as

$$\delta_n = \frac{\partial c}{\partial x_n}$$

The error gradient of the input vector at the last layer N is:

$$\delta_N = \frac{\partial c}{\partial x_N}$$
$$= \frac{\partial}{\partial x_N} * \frac{1}{2} * (y_N - t)^2$$
$$= \left(\frac{\partial}{\partial y_N} * \frac{1}{2} * (y_N - t)^2\right) * \frac{\partial y_N}{\partial x_N}$$
$$= (y_N - t) * \frac{\partial f(x_N)}{\partial x_N}$$
$$= (y_N - t) * f'(x_N)$$

The error gradient of the input vector at an inner layer n is:

$$\delta_n = \frac{\partial c}{\partial x_n}$$
$$= \frac{\partial c}{\partial x_{n+1}} * \frac{\partial x_{n+1}}{\partial x_n}$$
$$= \delta_{n+1} * \frac{\partial x_{n+1}}{\partial x_n}$$
$$= \delta_{n+1} * \frac{\partial w_n * y_n}{\partial x_n}$$
$$= \delta_{n+1} * \frac{\partial w_n * y_n}{\partial y_n} * \frac{\partial y_n}{\partial x_n}$$
$$= \delta_{n+1} * \frac{\partial w_n * Y_n}{\partial y_n} * \frac{\partial f(X_n)}{\partial x_n}$$
$$= (\delta_{n+1}) * w_n * f'(x_n)$$

Therefore, the error gradient of the input vector at a layer n is:

$$\delta n = f'(x_n) * \{(y_n - t)\} \text{ if } n = N$$

$$\delta n = f'(x_n) * \{\delta_{n+1} * w_n\} \text{ if } n < N$$

Hence, the error gradient of the weight matrix $w_n$ is:

$$\frac{\partial c}{\partial w_n} = \frac{\partial c}{\partial x_{n+1}} * \frac{\partial x_{n+1}}{\partial w_n}$$
$$= (\delta_{n+1}) * \frac{\partial w_n * y_n}{w_n}$$
$$= (\delta_{n+1} y_n)$$

Therefore, the change in weight should be:

$$\Delta w_n = -\alpha * \left(\frac{\partial c}{\partial w_n}\right)$$
$$= -\alpha * (\delta_{n+1} y_n)$$

Where α is the learning rate (or rate of gradient descent). Thus, we have shown the significant weight change, from which the implementation of a training algorithm follows trivially [4].

For the analysis by the neural network, the data set was split into three unequal subsets. Namely, the training set included 70% of the available cases, and the remaining cases were equally split into the validation and the test set. The training set should be substantially big and representative so that the network is efficiently developed and has the potential to generalize on the recently presented data of the investigation set. Because it's believed that ANNsfrequently fail to extrapolate successfully on fresh data, it is claimed that the training set contains data values outside the range of those used for testing and validation. The validation set was used in the training procedure for the implementation of the early-stopping art to shun the network's overfitting to the training data. The test set is by determination independent of the training procedure and is only utilized for the network's generalization ability evaluation.-By using the same data used in (MLR) those consist ofnine elements as input variables and the $PM_{2.5}$ as a target variable, we run the training state as a first step.
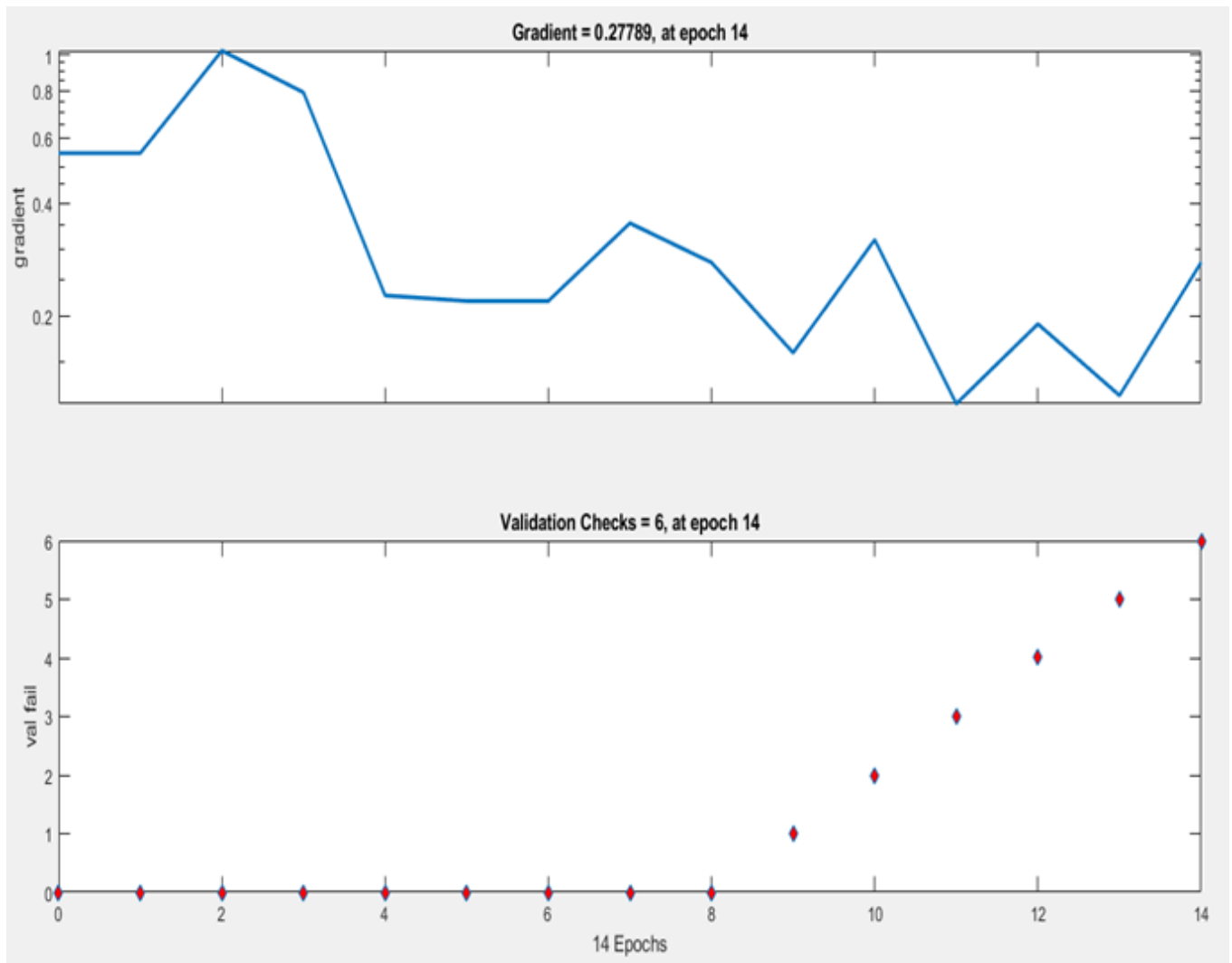
**Figure 4:** Training state

As shown in (Figure.4), the training state of the ANNs Model is presented in this Figure. The final epoch before the repetitions of the error, epoch (14), has the best validation performance, and therefore, the final model weights are picked based on this epoch. This Figure illustrate, the number of error repetitions is equal to 4, which would lead to having the validation check equal to 6.

- In the first diagram, the gradient is equal to (0.27789) at epoch (14), and it's used in the calculation of weights used in the network.
- The second diagram illustrates the number of validation checks at which the network reached to best performance and validation. Validation checks are equal to 6 at epoch (14).

The neural network trained various times to gain the best outcome; after training, the best five neural networks were found and compared in terms of fewer errors to determine the best neural network between them, the results of the best networks in the below:

As shown in (Figure.5), the correlation between the target and the output values is shown in this figure for both the training and validation. Both training and validation show desirable correlation coefficients (R values), and the correlation coefficient shows how strong the association between two variables is. The line passes through most points, and this indicates R's high value.

- In the training diagram, the blue line passes all points and this is helpful to make powerful learning to the machine.
- In the validation diagram, the green line passes all points and aligned with the dotted line. This proved that accuracy and precision are very high, so R-value is powerful it is equal to (0.998).
- In the test diagram, the red line passes through all test points, and it proved that it succeeded.
- In the final diagram, it illustrates the total correlation coefficient that is equal to (0.9986).
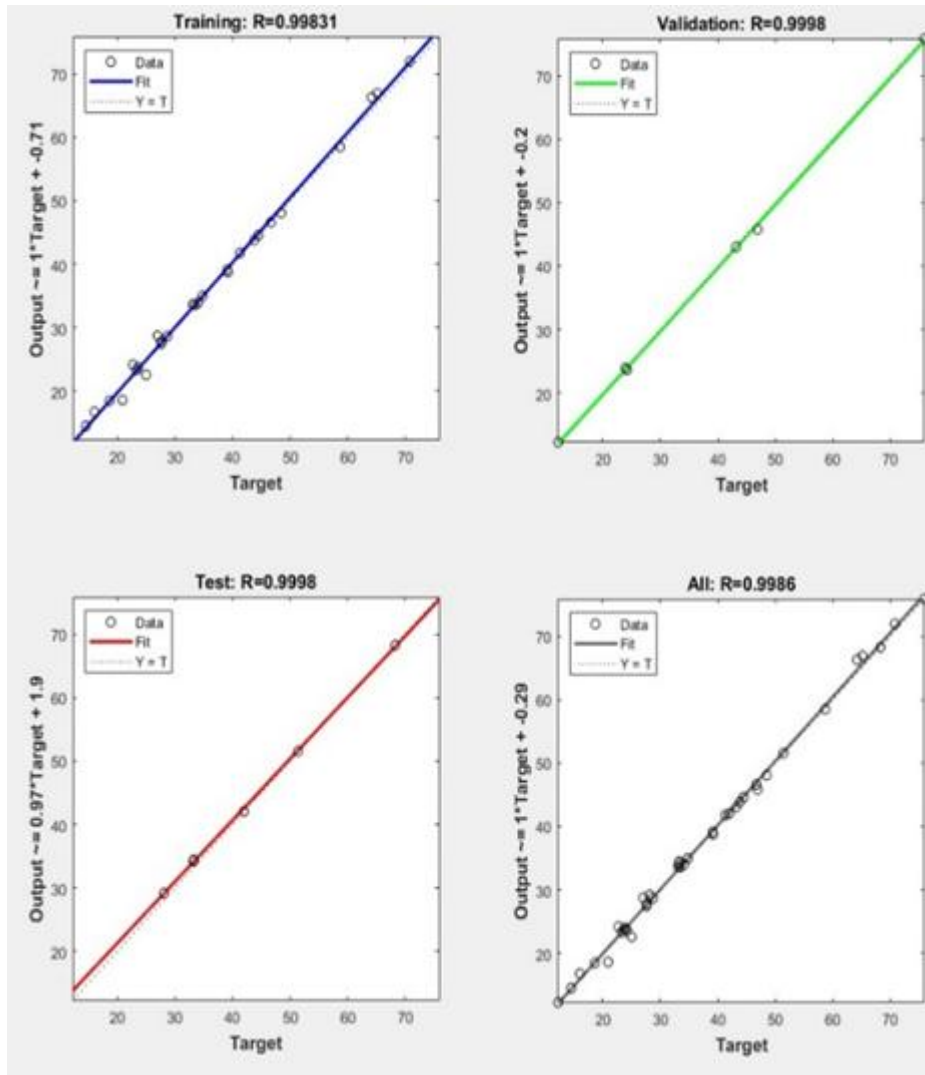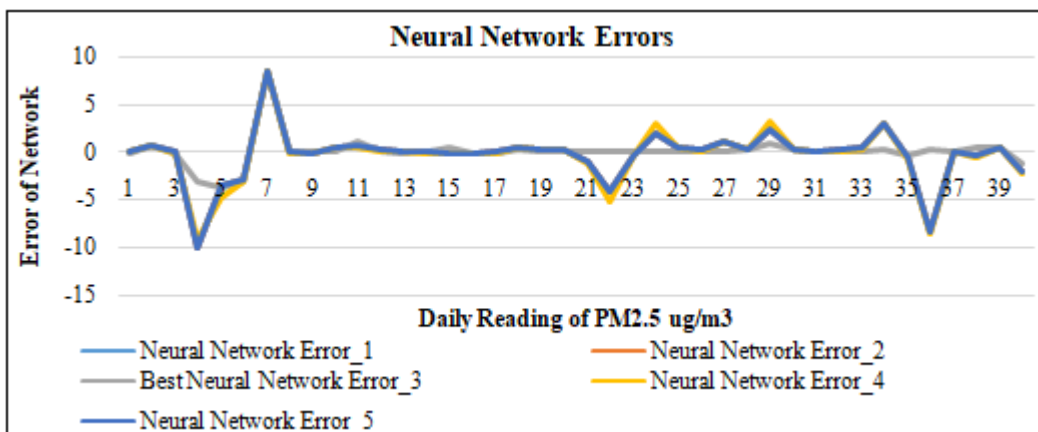
**Figure 5:** Regression



**Figure 6:** Neural Network Errors

As shown in (Figure.6), the NN errors diagram which explained the relationship between daily reading of PM2.5 and the Concentration of $PM_{2.5}$ (ug/m$^3$) for our five neural networks data and after taking the average of each NN error data, we find the best error of data. According to the figure, the best error is the third one, which colored by green. The average is equal to (0.10223).
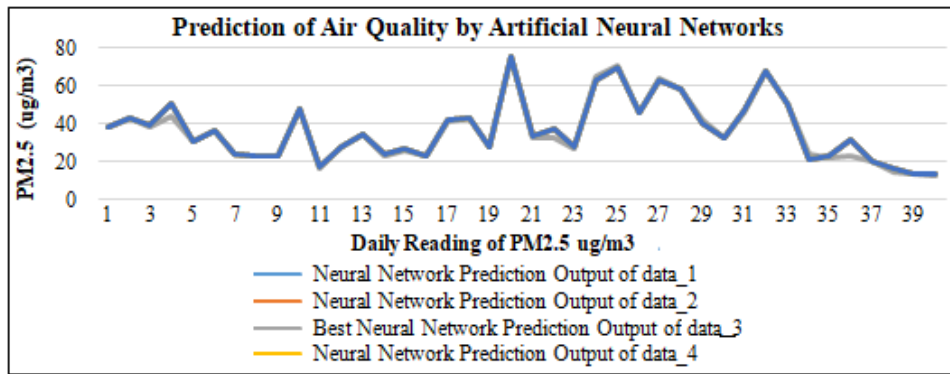
**Figure 7:** Neural Networks Prediction Output of Data

As shown in (Figure.7), the NN output data diagram for 5 neural networks data, and according to it, we find that the best NN output data that colored with green is very close to other data and this proved that it's the best line for output data. The average is equal to (36.8708).

In order to use the neural network to predict the values of the $PM_{2.5}$, 10-day readings were added for the input variables. After training, validation, and testing, the values of the $PM_{2.5}$ were obtained for this period of ten days as the following figure:
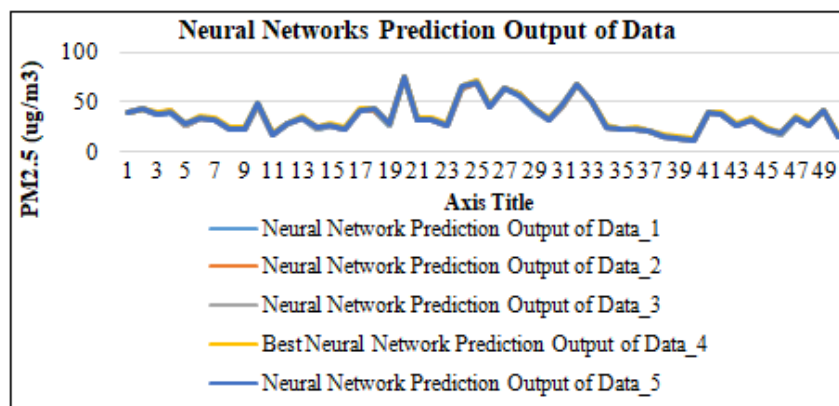


**Figure 8:** Neural Networks Prediction Output of Data

As shown in (Figure.8), the NN output data diagram for five neural networks data, and according to it, we find that the best NN output data that colored with green is very close to other data and this proved that it's the best line for output data. The average is equal to (35.56).

## 6. Conclusion

In this study, ANNs and multiple regression models were used to predict the daily average of $PM_{2.5}$ concentrations. The predictors were the nine trace elements with a high concentration in $PM_{2.5}$ ($ug/m^3$). From the result, neural network models have exceeded the regression models, indicating a nonlinear correlation between $PM_{2.5}$ and the nine predictors. Furthermore, as the contrast between the artificial neural network and multiple regression models are significant and steady, that is not the absolute dominance of the artificial neural networks in the predictability of $PM_{2.5}$ ($ug/m^3$). Based on the best neural network obtained through this study, predictive values of $PM_{2.5}$ ($ug/m^3$) were predicted for ten additional values. These predictions showed a high accuracy of the neural network in predicting the daily concentrations of suspended objects to prevent adverse effects on the population's health.

## References

[1] Chaloulakou, G. Grivas and N. Spyrellis, "Neural Network and Multiple Regression Models for PM10 Prediction in Athens: A Comparative Assessment", Journal of the Air & Waste Management Association, vol.53, no.10, pp.1183-1190, 2003. Available: 10.1080/10473289.2003.10466276.
[2] "WHO Air quality guidelines for particulate matter, ozone…" [Online]. Available: https://apps.who.int/iris/bitstream/handle/10665/69477/WHO_SDE_PHE_OEH_06.02_eng.pdf; sequence=1. [Accessed: 19-Dec-2017].
[3] S. Nayebare et al., "Chemical Characterization and Source Apportionment of PM2.5 in Rabigh, Saudi Arabia", Aerosol and Air Quality Research, vol.16, no.12, pp.3114-3129, 2017. Available: 10.4209/aaqr.2015.11.0658.
[4] Big Theta Θ, " The Math Behind Backpropagation | Big Theta. [Online]. Available: https://bigtheta.io/2016/02/27/the-math-behind-backpropagation.html. [Accessed: 24-Apr-2018].