

24 and 12 hours Exceedances prediction of particulate matter 10 Using extreme value distribution: A case study in Kuala Lumpur City Centre

Mohd Hafiz Zawawi^{1,*}, Law Zhe Yuan¹, Nurul Izma Mohammad², Ahmad Zia Al-Saufie³, Mohd Zakwan Ramli¹, Nor Azalina Rosli⁴

¹Department of Civil Engineering, College of Engineering, University Tenaga Nasional, Kajang, Malaysia

²Department of Civil Engineering, Faculty of Engineering, University Technology Petronas, Perak, Malaysia

³Department of Mathematics, University Technology Mara, Pulau Pinang, Malaysia

⁴Faculty of Engineering, University Malaysia Sarawak, Kuching, Malaysia

Abstract: In air pollution studies, one of the main purposes is to compute the air pollutants composition from sources of emission. One of the most frequently used statistical methods in determining the air pollutants composition is Extreme Value Distributions (EVD). This study determines and differentiates the effectiveness of Particulate Matter 10 (PM10) result analysis based on one 24-hour interval and two 12-hours intervals in a day at Kuala Lumpur City Centre (KLCC), Kuala Lumpur. The EVD analysis consists of Gumbel Distribution (Type I), Frechet Distribution (Type II) and Weibull Distribution (Type III). There are two performance indicators namely Root Mean Square Error (RMSE) from error measures and Coefficient of Determination (R²) from accuracy measures that will be used to assess distribution analysis. The two 12-hours intervals will be just on 13 November 2014 using Frechet Distribution for detailed EVD analysis. From the 24-hours result analysis, the probability of exceedances for 13 November 2014 is 0.0754 with 78 predicted and 52 actual number of unhealthy minutes based on the best distribution for that day itself which is Frechet Distribution (Type II). For 14 November 2014, the probability of exceedances is 0.0845 with 40 predicted and 39 actual number of unhealthy minutes based on Gumbel Distribution (Type I). As for the 12-hours intervals, the probability of exceedances is 0.0372 for the first interval and 0.1103 for the second interval. The predicted and actual number of unhealthy minutes is 14 and 8 respectively. All the observed and predicted PM10 values will be compared with the threshold value 61ug/m³ to determine the exceedances.

Key words: Extreme value distribution; PM10; Exceedances

1. Introduction

In Malaysia, the Department of Environment (DOE) has been observing closely the ambient air quality from the 52 air monitoring stations situated throughout the whole country. Through the observation, the air pollution composition consists mainly of six criteria pollutants which are ozone (O₃), lead (Pb), carbon monoxide (CO), sulfur dioxide (SO₂), nitrogen dioxide (NO₂) and particulate matter (PM) especially PM₁₀ (Ahmat et al., 2015). The originating sources of particulate matter can be the through natural phenomenon processes such as volcanic eruptions, sand storms in deserts, and woodland fires. The cases of human activities involve the combustion of fossil fuels and industrial processes such as power plants (Omidvarborna et al., 2015). The air pollutants pose a potential risk and damage to human health, crops and environment (Jamal et al., 2004). Results from previous studies indicate that the air pollutants in urban areas in Malaysia are having a steady increase from time to time (Afroz et al., 2003; Talib et al., 2002).

In recent years, statistical analysis and probability distributions have been extensively implemented in the investigation of air pollution data. In the past researches, various types of probability distributions have been applied to fit the air pollutant concentrations including Weibull distribution, lognormal distribution, Gamma distribution and Rayleigh distribution (Noor et al., 2011). Disparate area of studies have been published several years ago in the application of extreme value theory for example in the field of risk management operations (Yao et al., 2012), exposure to volatile organic compounds (Su et al., 2012), future market (Kao and Lin, 2010), capital requirements (Tsai and Chen, 2011), wind speed (Torrielli et al., 2013; Reynolds, 2012), wave heights (Petrov et al., 2013) and storm (Reeve et al., 2012). Although the extreme value theory is being extensively applied in hydrological studies, the contribution and importance of the theory in air pollution studies also cannot be neglected (Ahmat et al., 2014). EVD is frequently applied in assessment and estimation of air pollution concentration (Dasgupta and Bhaumik, 1995; Horowitz and Barakat, 1979; Kuchenhoff and Thamerus, 1996; Lu, 2002; Lu and Fang, 2003; Quintela and Fernández, 2011; Reyes et al., 2010;

* Corresponding Author.

Roberts, 1979; Smith, 1989; Surman et al., 1987). Throughout the years, the applications of statistical analysis and probability distributions were employed to grasp the current air quality and to anticipate the PM10 concentration in the time ahead. Weibull distribution (Wang et al., 2004; Lu et al., 2000), Gamma Distribution (Lu et al., 2003; Karaca et al., 2005) and Lognormal distribution (Lu et al., 2000; Hitzenberger and Tohno, 2000) are few of the probability distributions to be fit the air pollutant data. The modelling along with generalized extreme value distribution is a popular candidate for applications in risk management, finance, insurance, economics, hydrology, material sciences, telecommunications, and other extreme events (Kotz and Nadarajah, 2000). The distinguished feature of Extreme Value Analysis (EVA) is to evaluate the events' abnormality or scarcity such as the distribution concentrations in terms of maximum or minimum, exceedances or data frequencies. The type of distribution analysis applied in Extreme Value Theory (EVT) is the EVD (Ahmat et al., 2015). The EVT is regarding the probability calculations and statistical inference related to extreme values of random processes (Leong et al., 2002). The Fig. 1 shows the particulate matter emission load by sources from the year 2004 to 2012 in Malaysia (Ahmat et al., 2015) (Fig. 2).

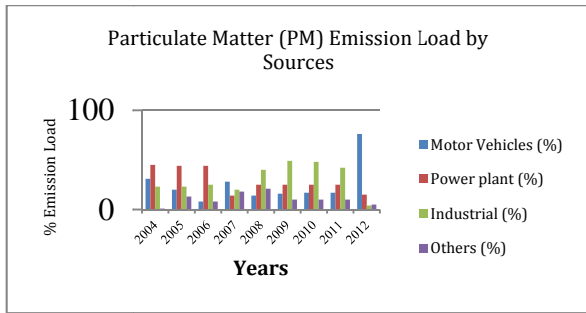


Fig. 1: Particulate matter emission load by sources from year 2004 to 2012 (Malaysia Environmental Quality Report)

2. Study area

The area of study is in KLCC, Kuala Lumpur. The location is considered to be in the heart of Kuala Lumpur city, Malaysia with the coordinates (3.1578° N, 101.7116° E). Geographically, the location is strategically located in the air pollution prone areas due to the vast vehicular emissions on the roads and highways. In addition, the southern part of Peninsular Malaysia is prone to the trans-boundary smoke due to forest fires from the Sumatera regions which contributed to the higher PM10 contributions (Ahmat et al., 2015).



Fig. 2: Location of Kuala Lumpur, Malaysia

3. Probability distribution

The series of data acquired will be computed and analyzed using the EVD method. Under the EVD method, there are three types of distribution which are the two-parameter Gumbel Distribution (Type I), three-parameter Frechet Distribution (Type II) and two-parameter Weibull Distribution (Type III). Table 1 shows the Probability Density Function (PDF) and Cumulative Distribution Function (CDF) equations for each distribution.

Table 1: EVD distributions with respective PDF and CDF equations

EVD	Probability Density Function (PDF)	Cumulative Distribution Function (CDF)
Gumbel	$f(x; \mu, \sigma) = \frac{1}{\sigma} \exp \left[-\frac{x-\mu}{\sigma} - \exp \left(\frac{x-\mu}{\sigma} \right) \right]$	$f(x; \mu, \sigma) = \exp \left[-\exp \left(\frac{x-\mu}{\sigma} \right) \right]$
Frechet	$f(x; \mu, \sigma, \lambda) = \frac{1}{\sigma} \left[1 + \lambda \left(\frac{x-\mu}{\sigma} \right)^{-\lambda} \right]^{-1} \exp \left\{ - \left[1 + \lambda \left(\frac{x-\mu}{\sigma} \right)^{-\lambda} \right] \right\}$	$f(x; \mu, \sigma, \lambda) = \exp \left[- \left(\frac{x-\mu}{\sigma} \right)^{-1-\lambda} \right]$
Weibull	$f(x; \sigma, \lambda) = \frac{\lambda}{\sigma} \left(\frac{x}{\sigma} \right)^{\lambda-1} \exp \left[- \left(\frac{x}{\sigma} \right)^\lambda \right]$	$f(x; \sigma, \lambda) = 1 - \exp \left[- \left(\frac{x}{\sigma} \right)^\lambda \right]$

4. Performance indicators (PI)

There are six PI which are divided into two types of measures which are the accuracy measures and error measures. Under accuracy measures are the PA, R2 and IA. On the other hand, under error measures are RMSE, NAE and MAE (Table 2 and 3).

The best distribution is selected based on the highest accuracy measures and lowest error measures. The distribution model is considered to be favorable when the error measures have values closer to zero while the accuracy measures are

closer to one. 4 Time plot analysis were shown in Figs 3 - 5.

5. 24 hours EVD analysis

Table 4 and 5 illustrates the EVD distributions with their respective parameters and each performance indicator with the ranking matrix. The best distribution is chosen based on the highest accuracy measure which is closest to the value of one and the lowest error measure which is closest to the value of zero. Based on Table 5, Weibull is chosen as the best distribution for each day due to the lowest

sum of total ranks for performance indicators among the three distributions. Fig. 6 shows the Cumulative Distribution Function (CDF) graph curves of the best distribution for each day while Table 6 shows the best distribution for each day with the probability of exceedances and percentage. On 13 November 2014, the probability of predicted exceedances for Frechet

Distribution (Type II) is 0.0754 ($F(x > 61) = 0.0754$) and the predicted number of minutes that exceed 61ug/m3 is 54. The number of the predicted PM10 exceedances in Frechet Distribution is 96% in agreement with the actual number of two minute intervals of 52 intervals among the three types of distribution.

Table 2: Performance indicators

Indicators	Equations
Accuracy Measures	
Prediction Accuracy (PA)	$PA = \sum_{t=1}^n \frac{(P_t - \hat{P})(O_t - \hat{O})}{(n-1)\sigma_p\sigma_o}$
Coefficient of Determination (R ²)	$R^2 = 1 - \frac{\sum_{t=1}^n (O_t - P_t)^2}{\sum_{t=1}^n (O_t - \bar{O})^2}$
Index of Accuracy (IA)	$IA = 1 - \frac{\sum_{t=1}^n (P_t - O_t)^2}{\sum_{t=1}^n (P_t - \bar{P} + O_t - \bar{O})^2}$
Error Measures	
Root Mean Square Error (RMSE)	$RMSE = \sqrt{\frac{\sum_{t=1}^n (O_t - P_t)^2}{n}}$
Normalized Absolute Error (NAE)	$NAE = \frac{\sum_{t=1}^n (P_t - O_t)}{\sum_{t=1}^n O_t}$
Mean Absolute Error (MAE)	$MAE = \frac{\sum_{t=1}^n (O_t - P_t) }{n}$

Table 3: Descriptive analysis

Descriptive Analysis	Date		
	13 November 2014	14 November 2014	15 November 2014
N	716	473	719
Mean	32.4772	32.6399	26.4392
Median	27.805	28.31	22.64
Variance	303.6712	286.0024	212.7682
Std. Deviation	17.4262	16.9116	14.5866
Minimum	10.34	0.02	8.35
Maximum	127.43	103.08	119.81
Range	117.09	103.06	111.46

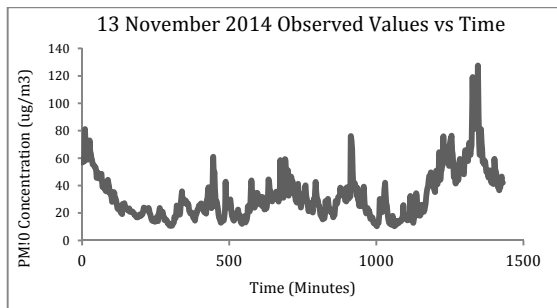


Fig. 3: Time series plot for PM10 concentration at KLCC on 13 November 2014

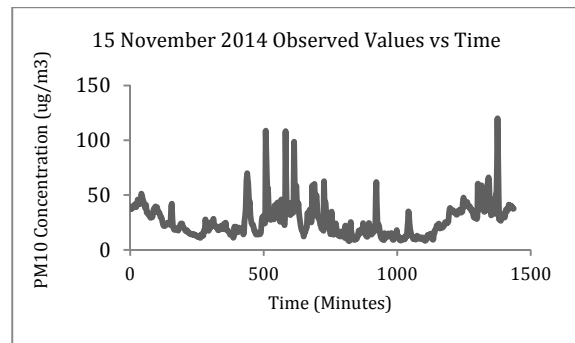


Fig. 5: Time series plot for PM10 concentration at KLCC on 15 November 2014

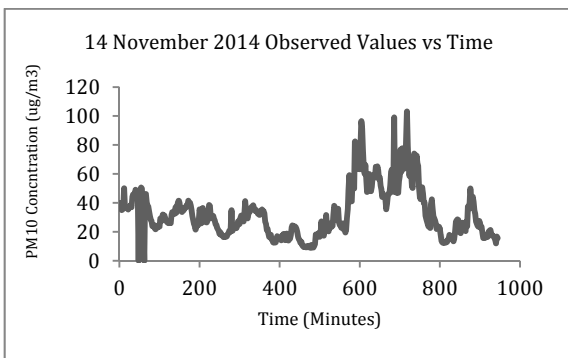


Fig. 4: Time series plot for PM10 concentration at KLCC on 14 November 2014

Therefore, this shows that the Frechet Distribution has the closest predicted number of unhealthy minutes to the actual number of unhealthy minutes which is 52. On 14 November 2014, the probability of having predicted PM10 distribution concentration more than 61 ug/m3 in a day for Gumbel Distribution is equal to 0.0845 ($F(x > 61) = 0.0845$). Therefore the predicted exceedances or the predicted number of minutes that exceed PM10 concentration of 61 ug/m3 is 40. With three types of distribution in comparison, the number of the predicted PM10 values in Gumbel Distribution is 98% in agreement with the actual number of two minute intervals higher than 61ug/m3 of 39 intervals. Therefore, the Gumbel Distribution has the closest predicted number of

unhealthy minutes to the actual number of unhealthy minutes which is 39. 15 November 2014 results show that the probability of predicted exceedances for Weibull Distribution is 0.0188 ($F(x > 61) = 0.0188$) in which the predicted number of minutes that exceed 61 ug/m³ is 18. The percentage

of accuracy is 78% between the number of predicted PM10 exceedances and the number of actual exceedances. Therefore, Weibull Distribution has the closest predicted number of unhealthy minutes to the actual number of unhealthy minutes which is 14 (Table 7).

Table 4: Extreme value distribution parameters

EVD		Time		
Distribution	Parameters	13 November 2014	14 November 2014	15 November 2014
Gumbel	μ	42.1146	41.8767	34.6319
	σ	23.7268	21.1417	22.3243
Frechet	μ	0.2549	0.0644	0.2211
	σ	10.4838	0.0644	9.1037
	λ	23.4246	24.7055	18.9934
Weibull	σ	36.8581	36.8250	29.9661
	λ	2.0002	2.0155	1.9415

Where: μ – location parameter, σ – scale parameter, λ – shape parameters

Table 5: Performance indicators for EVD distributions

Time	EVD	Performance Indicators				Sum
		RMSE	Rank	R ²	Rank	
13 November 2014	Gumbel	8.6104	3	0.7559	3	6
	Frechet	7.3763	2	0.8232	2	4
	Weibull	6.4057	1	0.8649	1	2
14 November 2014	Gumbel	7.5234	3	0.8021	3	6
	Frechet	6.0896	2	0.8703	2	4
	Weibull	5.4703	1	0.8954	1	2
15 November 2014	Gumbel	12.6212	3	0.2513	3	6
	Frechet	4.8317	2	0.8903	2	4
	Weibull	4.4083	1	0.9087	1	2

Table 6: Distribution type analysis

Date	EVD Distribution	Probability	Percentage (%)
13 November 2014	Frechet	0.0754	96
14 November 2014	Gumbel	0.0845	98
15 November 2014	Weibull	0.0188	78

Table 7: Predicted and Actual Number of unhealthy days

Date	Predicted no. of unhealthy days	Actual no. of unhealthy days
13 November 2014	78	52
14 November 2014	40	39
15 November 2014	18	14

6. 12 hours extreme value distribution

two 12-hour intervals using only Frechet Distribution (Table 8).

The detailed EVD analysis is carried out on only one day which is 13 November 2014 but divided into

Table 8: PM₁₀ descriptive analysis on 13 November 2014

Descriptive Analysis	Date (13 November 2014)	
	00:00 - 12:00	12:00 - 00:00
N	358	359
Mean	28.358	36.609
Median	24.73	32.29
Variance	176.88	395.49
Std. Deviation	13.3	19.887
Minimum	10.7	10.34
Maximum	81.08	127.43
Range	70.38	117.09

Based on Table 8, the maximum concentration of PM10 recorded in 13 November 2014 from 12am to 12pm is 81.08ug/m³ while the minimum concentration of PM10 recorded is 10.7ug/m³. The

mean PM10 concentration is 28.358ug/m³. Meanwhile, the highest PM10 concentration in 13 November 2014 from 12pm to 12am is 127.43ug/m³

while the lowest PM10 concentration is 10.34ug/m³ (Table 9 and 10).

Table 9: Performance indicators for Frechet Distribution on 13 November 2014

Time	EVD Distribution	Performance Indicators			
		RMSE	Rank	R ²	Rank
00:00 - 12:00	Frechet (Type II)	3.455	1	0.9325	1
12:00 - 00:00		9.625	2	0.7658	2

Table 10: Performance indicators for Frechet Distribution on 13 November 2014

Time	Probability	Predicted No. of Unhealthy Minutes	Actual No. of Unhealthy Minutes
13-Nov-14			
00:00 – 12:00	0.0372	14	8
12:00 – 00:00	0.1103	40	44
Sum		54	52
Percentage (%)		96	

From 12am to 12pm, the probability of having predicted PM10 distribution concentration more than 61 ug/m³ in a day is equal to 0.0372 ($F(x > 61) = 0.0372$). Therefore the predicted exceedances or the predicted number of minutes that exceed PM10 concentration of 61 ug/m³ is 14. The probability of exceedances from 12pm to 12am is 0.1103 ($F(x > 61) = 0.1103$) and the predicted number of minutes that exceed 61 ug/m³ is 40. The total sum of predicted number of unhealthy minutes in two intervals is 54 while the total sum of actual number of unhealthy minutes is 54 which are as the same as Table From 12am to 12pm, the number of actual observed PM10 distribution concentration more than 61 ug/m³ in a day is 8. The number of actual observed PM10 values from 12pm to 12am is 44. Therefore, the total sum of actual number of unhealthy minutes is 52 similar to table and that leads to the same percentage of accuracy between the predicted and actual number of minutes which is 96%.

7. Conclusion

In conclusion, the overall best distribution is assumed to be Gumbel Distribution (Type I) as the percentage is the highest among the three distributions which is 98% similarity between the predicted and actual number of unhealthy minutes. Another approach is the division of the secondary data into smaller time intervals for statistical analysis for example performing the EVD distributions analysis for four 12 hour intervals for each day instead of just the whole day itself. With that, the CDF graph curves can be illustrated in 12 figures with three distributions for each interval in one day rather than only three CDF graph curves for one day. Besides that, the probability as well as the number of predicted and actual unhealthy minutes provided can be more detailed and gives space for further discussion. Furthermore, the location of study is also important as the occurrence of PM10 observation is important to determine the outcomes of the analysis.

The average PM10 concentration is 36.609ug/m³. On 13 November 2014 from 12am to 12pm, the RMSE is 3.455 while from 12pm to 12am the RMSE is 9.625. Therefore the best error measurement for RMSE is from 12am to 12pm as the error is the nearest to zero to be considered as the best model. The R² from 12am to 12pm is 0.9325 while the R² from 12pm to 12am is 0.7658.

References

- Afroz, R., Hassan, M., N., & Ibrahim, N., A. [2003]. Review air pollution and health impact in Malaysia. *Environ.Res.* 92: 71–77
- Ahmat, H., Yahaya, A., S., & Ramli, N., A. [2015]. PM10 Analysis for Three Industrialized Areas using Extreme Value. *Sains Malaysiana* 44(2): 175–185
- Dasgupta, R. & Bhaumik, D.K. [1995]. Upper and lower tolerance limits of atmospheric ozone level and extreme value distribution. *Sankhya: The Indian Journal of Statistics* 57(B2): 182-199
- Hitzenberger, R., & Tohno, S. [2000]. Black carbon (BC) concentrations and size distributions at two urban sites (Uji, Japan and Vienna, Austria). *Journal of Aerosol Science* 31: 112–113
- Horowitz, J. & Barakat, S. [1979]. Statistical analysis of the maximum concentration of an air pollutant: Effects of autocorrelation and non-stationarity. *Atmospheric Environment* (1967) 13(6): 811-818
- Jamal, H.,H., Pillay, M.,S., Zailina, H., Shamsul, B.,S., Sinha, K., Zaman Huri, Z., Khew, S.,L., Mazrura, S., Ambu, S., Rahimah, A. & Ruzita, M.S. [2004]. A Study of Health Impact & Risk Assessment of Urban Air Pollution in Klang Valley, Malaysia. Kuala Lumpur: UKM Pakarunding Sdn Bhd.
- Kao, T., C., & Lin, C., H. [2010]. Setting Margin Levels in Future Markets: An Extreme Value Method. *Non-linear Analysis: Real World Applications* 11(6): 1704–1713
- Karaca, F., Alagha, O., & Erturk, F. [2005]. Statistical characterization of atmospheric PM10 and PM2.5 concentrations at a non-impacted suburban site of Istanbul, Turkey. *Chemosphere* 59: 1183–1190
- Kotz, S., & Nadarajah, S. [2000]. *Extreme Value Distributions: Theory and Applications*. London: Imperial College Press.
- Kuchenhoff, H. & Thamerus, M. [1996]. Extreme value analysis of Munich air pollution data.

- Environmental and Ecological Statistics 3: 127-141
- Leong, Y., P., Sleight, J., W., & Torrance, J., M. [2002]. Extreme value theory applied to postoperative breathing pattern. *British Journal of Anaesthesia* 88: 61-4
- Lu, H., C. [2003]. Estimating the Emission Source Reduction of PM10 in Central Taiwan. *Chemosphere* 54: 805-814
- Lu, H.C. & Fang, G.C. [2003]. Predicting the exceedances of a critical PM10 concentration - A case study in Taiwan. *Atmospheric Environment* 37(8): 3491-3499
- Lu, H.C. [2002]. The statistical characters of PM10 concentration in Taiwan area. *Atmospheric Environment* 36(3): 491-502
- Noor, N., M., Tan, C.,Y., Abdullah, M., M., A., B., Ramli, N., A., & Yahaya, A., S. [2011]. Modelling of PM10 Concentration in Industrialized Area in Malaysia: A Case Study In Nilai. *IPCBE 12*
- Omidvarborna, H., Kumar, A., & Kim, D., S. [2015]. Recent studies on soot modeling for diesel combustion. *Renewable and Sustainable Energy Reviews* 48: 635-647
- Petrov, V., Guedes, S., C., & Gotovac, H. [2013]. Prediction of extreme significant wave heights using maximum entropy. *Coastal Engineering* 74(4):1-10
- Quintela-del-Río, A. & Francisco-Fernández, M. [2011]. Analysis of high level ozone concentrations using nonparametric methods. *The Science of the Total Environment* 409(2): 1123-1133
- Reeve, D., T., Randell, D., Ewans, K., C. & Jonathan, P. [2012]. Uncertainty due to choice of measurement scale in extreme value modelling of North Sea storm severity. *Ocean Engineering* 53(10): 164-176
- Reyes, H.J., Vaquera, H. & Villasenor, J.A. [2010]. Estimation of trends in high urban ozone levels using the quantiles of (GEV). *Environmetrics* 21: 470-481
- Reynolds, A., M. [2012]. Gusts within plant canopies are extreme value processes. *Physica A: Statistical Mechanics and Its Applications* 391(11): 5059-5063
- Roberts, E.M. [1979]. Review of statistics of extreme values with applications to air quality data Part I. *Journal of the Air Pollution Control Association* 29(6): 632-637
- Smith, R.L. [1989]. Extreme value analysis of environmental time series: An application to trend detection in ground-level ozone. *Statistical Sciences* 4(4): 367-393
- Su, F., C., Jia, C., & Batterman, S. [2012]. Extreme Value analyses of VOC exposures and risks: A comparison of RIOPA and NHANES datasets. *Atmospheric Environment Atmospheric Environment* 62: 97-106
- Surman, P.G., Boderó, J. & Simpson, R.W. [1987]. The prediction of the numbers of violations of standards and the frequency of air pollution episodes using extreme value theory. *Atmospheric Environment* 21(8): 1843-1848
- Talib, M.L., Rozali, M.O., Norela, S., Ahmad Daud, M.N. & Permata, N.J. [2002]. Air quality in several industrial areas in Malaysia. In *Proceedings of the Regional Symposium on Environment and Natural Resources*, edited by Omar, R., Ali Rahman, Z., Latif, M.T., Lihan, T. & Adam, J.H. April 10-11. Renaissance Hotel, Kuala Lumpur. pp. 703-710
- The Map of Malaysia Internet: <http://www.malaysiasite.nl/map.htm>, January 28, 2015.
- Torrielli, A., Repetto, M., P., & Solari, G. [2013]. Extreme wind speeds from long-term synthetic records. *Journal of Wind Engineering and Industrial Aerodynamics* 115(4): 22-38
- Tsai, M., S., & Chen, L., C. [2011]. The calculation of capital requirement using extreme value theory. *Economic Modelling* 28(1): 390-395
- Wang, X., & Mauzerall, D., L. [2004] Characterizing Distributions of Surface Ozone and Its Impact on Grain Production in China, Japan and South Korea: 1990 and 2020. *Atmospheric Environment* 38: 4383-4402
- Yao, F., Wen, H., & Luan, J. [2012]. CVaR measurement and operational risk management in commercial banks according to the peak value method of extreme value theory. *Mathematical and Computer Modelling* 58(1-2): 15-27