

urban hydrology involve flood and pollution prevention. However, our study only considers flood prevention based on the river flow prediction. If an undeveloped area is then transformed into a developed area, the conditions of the soil structure will be disturbed (Hall, 1984). These factors can change the magnitude of the river flow. The volume of runoff will increase significantly with the increase in the magnitude of the river flow due to the impervious areas and the lack of drainage. Hence, downstream flooding problems exist in urban areas. There are several methods that can be used to estimate the river flow in a watershed that is located in an urban area, such as the empirical and physical process methods. Referring to the empirical method for urban hydrology research, the behaviour of river flow in the downstream area is important to provide accurate information for the whole river flow (Viesmann and Lewis, 1996). This information can help in planning, development and flood prevention of the downstream area.

The Langat River, which is one of the longest rivers in the state of Selangor, Malaysia, is used as a case study. This research focuses on the downstream area at Kajang, which is well-known for experiencing flood hazards. Figure 1 shows the four gauging stations along the Langat River. The Langat River flows from east to southeast, which is from Lui River to Kajang. The total length of the upstream and downstream is about 34.4 km and the downstream area has been identified as a flood risk area (Mohammed et al., 2011). Checkpoint 1 (station number 3118445) is located at the Lui River gauging station (upstream) and Checkpoint 2 (station number 2917401) is located at the Kajang gauging station (downstream). The Langat River at the Kajang gauging station has been used for the river flow analysis and prediction using the nonlinear prediction method. This area had a population of 229 655 people in 2000, which increased to 342 657 people in 2010 (Department of Irrigation and Drainage Malaysia, 2005). The increase in population in this area reflects the development in the Kajang area. Furthermore, the study area is adjacent to an industrial area and pig farms. Flooding in this area can cause damage to the industrial area and pollution in the Langat River basin. Thus, studies of the downstream area (Kajang) are important to provide information

14333

about the flow downstream. This study was conducted at this point so that the release of water from Checkpoint 2 could be estimated for a certain length of time. The results of this study could help to identify the preventive measures that could be undertaken in this downstream area.

The analysis and prediction of river flow could provide the information about the dynamics of the river flow system. However, the flow of the river is not dependent on rainfall alone. The characteristics of an area, such as shape, slope, land, soil structure and climate change, can also affect the flow of the river in an area (Viesmann and Lewis, 1996). Thus, the application of stochastic methods is often used to analyse complex natural conditions, such as the river flow. Developments in the study of nonlinear time series analysis is growing with some revolutionary methods. One particular method that provides important findings is known as chaos theory, which explains that a complex system can be analysed by deterministic methods that use a minimum number of the system's variables (Islam and Sivakumar, 2002). Several decades ago, a number of studies were performed to obtain information on characterizing, modelling and predicting hydrological phenomena as a deterministic system (e.g. Jayawardena and Lai, 1993; Sivakumar, 2000; Ghorbani et al., 2010). The results showed that the river flow prediction and other hydrological processes are in good agreement with the actual data values (Sivakumar, 2003; Regonda et al., 2005; She and Yang, 2010; Khatibi et al., 2012). In addition, prediction using chaos theory can reveal the number of variables that affect the dynamics of the river flow.

Studies on river flow analysis and prediction in Malaysia have been done and improved for a variety of purposes, such as providing information for flood prevention. Several methods, such as support vector machine method (Shabri and Suhartono, 2012), neural network model (Ahmad and Juahir, 2006) and hydrodynamic modelling (Ghani et al., 2010), have been used for river flow prediction. However, several methods have yet to be explored for the purpose of river flow prediction in Malaysia, such as chaos theory, Bayesian methods and wavelet methods. River flow prediction using chaos theory involves a single variable (river flow data) albeit there are other dominant

14334

The correlation dimension is based on the correlation integral introduced by Grassberger (1986):

$$C(r) = \lim_{N \rightarrow \infty} \frac{2}{N(N-1)} \sum_{i,j=1}^N H(r - \|Y_i - Y_j\|) \quad (4)$$

where H is the Heavyside function, which has the value 0 or 1 and can be defined as:

$$H(r - \|Y_i - Y_j\|) = \begin{cases} 1, & 0 \leq (r - \|Y_i - Y_j\|) \\ 0, & < (r - \|Y_i - Y_j\|) \end{cases} \quad (5)$$

and acts as a barrier to the Euclidean distance between two points on the attractor Y_i and Y_j . The correlation function $C(r)$ is calculated for the pair of points (Y_i, Y_j) with a distance less than the radius r . In the limit to infinite amount of data ($N \rightarrow \infty$) and sufficiently small $r (r \rightarrow 0)$, the relation $C(r) \cong \alpha r^{D_2}$ is expected (Men et al., 2004). The correlation dimension D_2 and correlation exponent ν can be defined as:

$$D_2 = \lim_{r \rightarrow 0} \nu \quad (6)$$

$$\nu = \frac{\delta[\log C(r)]}{\delta[\log r]} \quad (7)$$

Several steps are required to identify the value of the correlation dimension. The first step is to draw a graph $\ln C(r)$ vs. $\ln(r)$ with a given m . Then the gradient (correlation exponent ν) of the m -dimensional curve values has to be determined. The gradient of the graph can be measured by the least squares method for determining the scaling. For finite data and where the value of r exceeds the diameter, there exists a saturated area of the graph. The saturated area is the scaling region. A better way to estimate the gradient is to use $\delta[\log C(r)]/\delta[\log r]$. To examine if there is a chaotic nature, the correlation exponent (slope ν) vs. m -dimensional has to be plotted. If the value of the

14337

correlation exponent is finite, low and non-integer, the system is considered to be of low dimensional chaotic nature (Men et al., 2004). If the correlation value increases without limit as m increases, the system should be studied as a stochastic system.

The false nearest neighbour method (FNN) is an effective method for finding the embedding dimension m for the reconstruction phase space. This method has been used to analyse river flow time series (Wu and Chau, 2010; Ghorbani et al., 2012). This paragraph describes how FNN is implemented. Suppose the dimension increases then the distance between the point and the nearest neighbour should not change if it is indeed the nearest neighbouring point. Computation of the distance between the point and the nearest neighbour is by the Euclidean distance.

FNN can be calculated using the following algorithm. Assume that $Y_i = \{X_i, X_{i+\tau}, X_{i+2\tau}, \dots, X_{i+(m-1)\tau}\}$ has nearest neighbour $Y_i^{NN} = \{X_i^{NN}, X_{i+\tau}^{NN}, X_{i+2\tau}^{NN}, \dots, X_{i+(m-1)\tau}^{NN}\}$. Then, calculate the Euclidean distance $\|Y_i - Y_i^{NN}\|$.

For all points i in vector space, equation $\frac{|X_{i+m\tau} - X_{i+m\tau}^{NN}|}{\|Y_i - Y_i^{NN}\|} > R_T$ is used and the value of false nearest neighbour can be calculated. R_T is a value between 10 and 30. In this study, the value of 15 is used (Wu and Chau, 2010). Repeat the algorithm with different embedding dimensions and the value of the false nearest neighbour that is close to zero is used as the embedding dimension.

2.2 Prediction

In this study, the prediction of river flow has been performed by using the local linear approximation method. This method was proposed by Lorenz (1969). Application of the local linear approximation method is to (1) examine whether the river flow at the downstream areas can be predicted, (2) to compare the prediction results for Models I and II. The local linear approximation method is used to predict river flow in downstream areas as follows. The first step is to reconstruct the phase space. The combination of the preliminary parameter pair (τ, m) is important for reconstruction of the phase space

14338

because this phase space result will be used in making a prediction. The difference between Models I and II is in the reconstruction phase space. Models I and II involve $\tau = 1$ but involve different methods in determining the value of m . Model I uses the correlation dimension while FNN is employed for Model II. Assume that the reconstruction of phase space is like $Y_i = \{x_i, x_{i+\tau}, x_{i+2\tau}, \dots, x_{i+(m-1)\tau}\}$. The nearest neighbour for Y_t is required to predict Y_{t+1} . Assume that the vector of the minimum distance to the nearest neighbour is Y_M . Next, for the local linear approximation method, the values of Y_M and Y_{M+1} are used to satisfy the linear equations $Y_{M+1} = AY_M + B$. The constants A and B are calculated using the least squares method. Thus, the predictive value Y_{t+1} can be calculated using $Y_{t+1} = AY_t + B$.

2.3 Performance evaluation

The assessment of the prediction accuracy of the models for predicting the daily river flow is evaluated by using the mean absolute error (MAE), root mean square error (RMSE) and correlation coefficient (CC). The MAE, RMSE and CC are as follows:

$$MAE = \frac{1}{n} \sum_{t=1}^n |y_t^o - y_t^f| \tag{8}$$

$$RMSE = \sqrt{\frac{1}{n} \sum_{t=1}^n (y_t^o - y_t^f)^2} \tag{9}$$

$$CC = \frac{\frac{1}{n} \sum_{t=1}^n (y_t^o - \bar{y}_t^o) (y_t^f - \bar{y}_t^f)}{\sqrt{\frac{1}{n} \sum_{t=1}^n (y_t^o - \bar{y}_t^o)^2} \sqrt{\frac{1}{n} \sum_{t=1}^n (y_t^f - \bar{y}_t^f)^2}} \tag{10}$$

where y_t^o is the observed and y_t^f is the forecast value at time t , and n is the number of data points. MAE and RMSE can provide information on the predictive ability of the

14339

models involved. Meanwhile, the correlation coefficient CC can measure the correlation between the prediction and the observed data.

3 Description of data

Langat River is one of the longest rivers in Selangor and its river basin is transboundary, inasmuch as it crosses three states – Selangor, Negeri Sembilan, and the Federal Territory of Kuala Lumpur and Putrajaya (Department of Irrigation and Drainage Malaysia, 2011). The Langat River flows from Mount Nuang in Hulu Langat district to the Straits of Malacca in Kuala Langat. The Langat River catchment area covers a total of 1815 km² and is located between latitude 2°40'152'' N and 3°16'15'' N, and longitude 101°19'20'' E and 102°1'10'' E (Juahir et al., 2011). There are two water reservoirs located in this area – Langat Dam and Semenyih Dam. Langat Dam was built with an area of 54 km² and Semenyih Dam has an area of 41 km². Both of these dams were built to deliver water for domestic and industrial use. In addition, the Langat Dam is also used to generate electricity for the use of residents in the vicinity of the Langat Valley. There are several towns and villages built along the Langat River – Cheras, Semenyih, Dengkil and Kajang. Since 1976, Langat River has also been acknowledged to be an area that regularly suffers flooding.

The variations of daily river flow data for Checkpoint 2 are shown in Fig. 2. The irregular patterns in data for Kajang River show that the river in this area is a complex system. The overall data were taken from the Department of Irrigation and Drainage Malaysia. Missing data constitute about 0.018% and were filled using the results of linear interpolation calculation. The statistical parameters of the data cover a period of four years (January 2002 to December 2005) and are shown in Table 1.

4 Results and Discussion

River flow prediction using NLP involves the reconstruction of the phase space and prediction. Thus, the discussion of the findings is divided into two parts. The first part is to determine the parameters for the reconstruction of the phase space for Models I and II. Meanwhile, the description of the prediction results are discussed in the second part.

The phase diagram can provide information about the dynamics of a system through the trajectories in the phase space. The trajectories that are of interest focus on a subspace called the attractor. In addition, the observation of attractor trajectories in the phase space can provide information about the chaotic behaviour of the system. Hence, the phase diagram and the observation data involved are plotted. Figure 3 shows the phase diagram in two and three dimensions with $\tau = 1$. The trajectories in the phase space can indicate the presence of chaotic behaviour of the data (Sivakumar, 2002). Referring to Fig. 3, the trajectories of the attractor are clearly shown in the two phase diagrams. Thus, the data involved in this analysis are chaotic. Therefore, the dynamics of the system can be studied using chaos theory without involving stochastic methods.

This study involved data from January 2002 to December 2005 (1433 days). Three years of data are used in the reconstruction of the phase space to predict the behaviour one year ahead. Reconstruction of the phase space is based on the embedding dimension. In Model I, the embedding dimension is based on the calculation of the correlation dimension. Graph $\ln C(r)$ vs. $\ln(r)$ in Fig. 4a shows the behaviour of the correlation function ν vs. radius r for the increasing m -dimensional. In general, the increasing value of the m -dimensional gradient occurs at the beginning of the curve from left $m = 1$ to right $m = 10$. Meanwhile, the graph of the correlation dimension estimation, the relationship between the correlation exponent ν for different values of m is shown in Fig. 4b. The relationship between the value of correlation dimension d_2 and m -dimension can be seen in Fig. 4c, which is a graph of d_2 vs. m . The value of the correlation dimension

14341

increased as the value of the m -dimensional increased. The increase in m -dimension can be seen up to a scaling region where the correlation dimension is saturated. The situation in which the value for the correlation dimension is saturated might indicate the existence of deterministic dynamics in the system. The saturated conditions for the d_2 value is in the interval (2.5, 3). The saturation value for d_2 is known as the correlation dimension attractor (Sivakumar, 2000). In general, the sufficient condition for the value of the smallest integer m is m greater than $2D_2$ (Wu and Chan, 2010). Thus, the value of $m = 6$ is related to the Langat River flow time series in Kajang. The correlation dimension d_2 is finite and shows low levels of correlation dimension. Hence, Sungai Langat is a chaotic and deterministic system. Model I involves a combination of preliminary parameters (1, 6) in the phase space reconstruction. Model II involves the calculation of m using FNN to find a combination of the preliminary parameters for RPS. Figure 5 shows the percentage of false nearest neighbours vs. m . Thus, the optimal value for the embedding dimension identified is $m = 14$. Model II involves a combination of parameters (1, 14) for the reconstruction phase space.

4.1 River flow prediction for Model I and Model II

The combination of preliminary parameters for Model I is (1, 6) while for Model II it is (1, 14). Thus, for both models, the combination of the preliminary parameters (τ, m) has been applied to construct the phase space. Figure 6 and Table 2 provide a summary of the river flow prediction results in terms of MAE, RMSE and CC. Overall, the results show good performance prediction for chaos theory in predicting the future value of the river flow for the downstream area. Referring to Table 2, a comparison of prediction performance shows that the prediction results for Model II are better than Model I. The correlation coefficient for Model II (0.6360) is slightly higher compared to Model I (0.6103). Thus, analysis and prediction of the Langat River can provide information in which the selection of a combination of preliminary parameters in the reconstruction phase space is essential for better prediction results. In this study, Model II uses FNN

14342

to calculate the embedding dimension m and is more appropriate than the correlation dimension method.

5 Conclusions

Analysis and prediction for testing the presence of chaotic behaviour in daily river flow data recorded at Langat River involving the station at Kajang, Selangor, Malaysia, has been performed. The station is located in the downstream area, which is a flood prone area. The analysis was carried out on the river flow data for a period of 4 yr (2002–2005). The focus of this study was to identify the chaotic behaviour of the river flow data in the downstream area and determine whether the river flow can be predicted when chaotic behaviour of the river exists downstream. Chaos theory, together with NLP, were used in the analysis. The reconstruction phase space clearly shows the existence of a chaotic attractor. Hence, the data involved in this analysis are chaotic. Next, was the attempt to make a prediction for one year ahead with the observed data using the results of the reconstruction of the phase space for three years. Two combinations of preliminary parameters were used. Model I used $\tau = 1$ for which m is the result of the calculation of the correlation dimension, while Model II used $\tau = 1$ for which m is the result of FNN calculation. Using these methods, the optimal combination for Model I was (1,6) and for Model II it was (1,14). The overall prediction results showed that both models could give a good prediction for the river flow downstream. However, the combination of the preliminary parameters for Model II using the FNN algorithm provided a better prediction result than Model I, which used the correlation dimension. The results showed that Langat River in Kajang, which is in the downstream area, is chaotic and predictable using NLP. Therefore, the results of the analysis and prediction of river flow in the downstream area could provide information on river flow for the authorities to take appropriate control of the downstream flooding.

14343

Acknowledgements. The authors are grateful to Universiti Kebangsaan Malaysia for providing financial support via the grant UKM-DIP-2012-31.

References

- Abarbanel, H. D. I.: Analysis of observed chaotic data, Springer-Verlag, Inc, New York, 1996.
- Adenan, N. H. and Noorani, M. S. M.: Behaviour of daily river flow: Chaotic?, Proceedings of the 20th National Symposium on Mathematical Sciences: Research in Mathematical Sciences: A Catalyst for Creativity and Innovation, Malaysia, 221–228, 2013.
- Ahmad, Z. and Juahir, H.: Neural Network Model for Prediction of Discharged from the Catchments of Langat River, Malaysia, IIUM Engineering Journal, 7, 25–34, 2006.
- Department of Irrigation and Drainage Malaysia: Langat River Integrated River Basin Management Study, Final Report, Technical Studies Part 1 of 4, 2005.
- Department of Irrigation and Drainage Malaysia: Review of the National Water Resources (2000–2050) and Formulation of National Water Resources Policy, Selangor, Federal Territory of Kuala Lumpur and Putrajaya, 2011.
- Ghani, A. A., Ali, R., Zakaria, N. A., Hasan, Z. A., Chang, C. K., and Ahmad, M. S. S.: A Temporal Change Study of the Muda River System Over 22 Years, International Journal of River Basin Management, 8, 25–37, 2010.
- Ghorbani, M. A., Kisi, O., and Aalinezhad, M.: A probe into the chaotic nature of daily streamflow time series by correlation dimension and largest Lyapunov methods, Appl. Math. Modell., 34, 4050–4057, 2010.
- Ghorbani, M. A., Daneshfaraz, R., Arvanagi, H., and Pourzangbar, A.: Local Prediction in River Discharge Time Series, Online Journal of Civil Engineering and Urbanism, 2, 51–55, 2012.
- Grassberger, P.: Do climatic attractors exist?, Nature, 323, 609–612, 1986.
- Hall, M. J.: Urban Hydrology, Elsevier Applied Science Publishers, New York, 1984.
- Islam, M. N. and Sivakumar, B.: Characterization and prediction of runoff dynamics: a nonlinear dynamical view, Adv. Water Res., 25, 179–190, 2002.
- Jayawardena, A. W. and Lai, F.: Chaos in hydrological time series, Proceeding of the Yokohama Symposium – Extreme Hydrological Events: Precipitation, Floods and Droughts, Yokohama, 59–66, 1993.

14344

Table 1. Statistics for river flow series at Kajang station (Checkpoint 2).

Number of data	Average	Max	Min	Standard Deviation	Skew	Kurtosis
1433	8.6492	212	0.3	13.779	6.733	68.438

14347

Table 2. Prediction performance at Kajang station (Checkpoint 2).

	MAE	RMSE	CC
MODEL I	2.6857	4.4025	0.6103
MODEL II	3.5160	5.6929	0.6360

14348

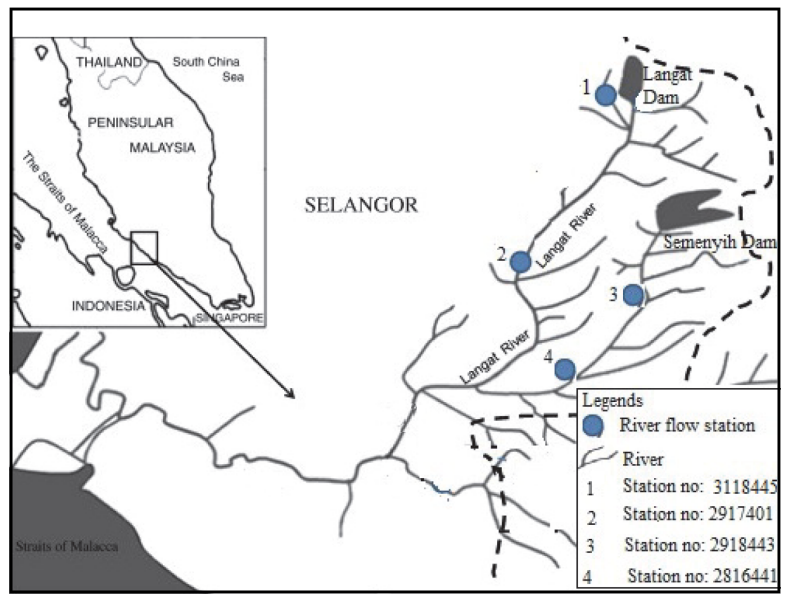


Fig. 1. Location of stations.

14349

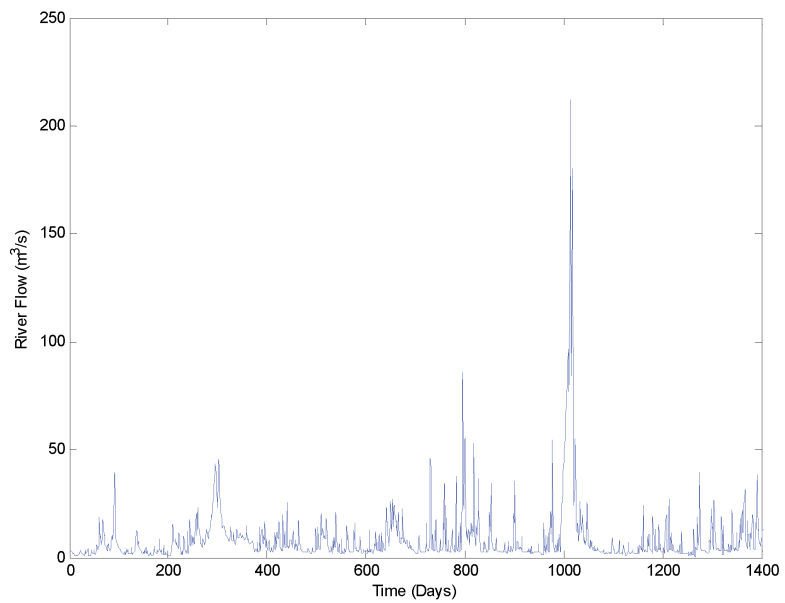


Fig. 2. Variations in data for Kajang station (Checkpoint 2).

14350

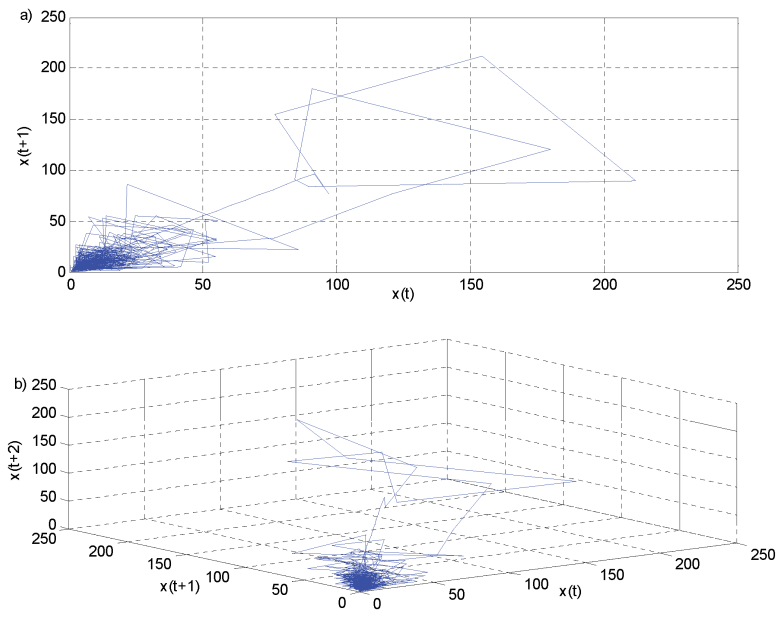


Fig. 3. Phase space diagram.

14351

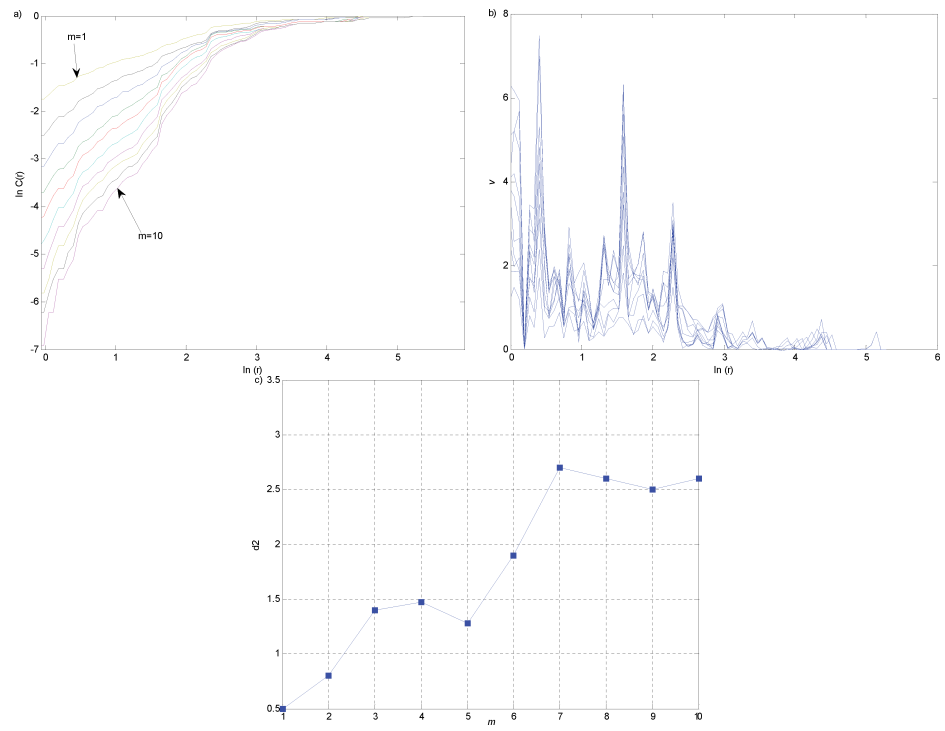


Fig. 4. $\ln C(r)$ vs. $\ln r$ for (a); the estimation of correlation dimension (d_2) for $\tau = 1$ and m is at the interval of [1, 10] (increasing from bottom to top in each pane) (b); relationship between d_2 and m for daily river flow of Langat River (c).

14352

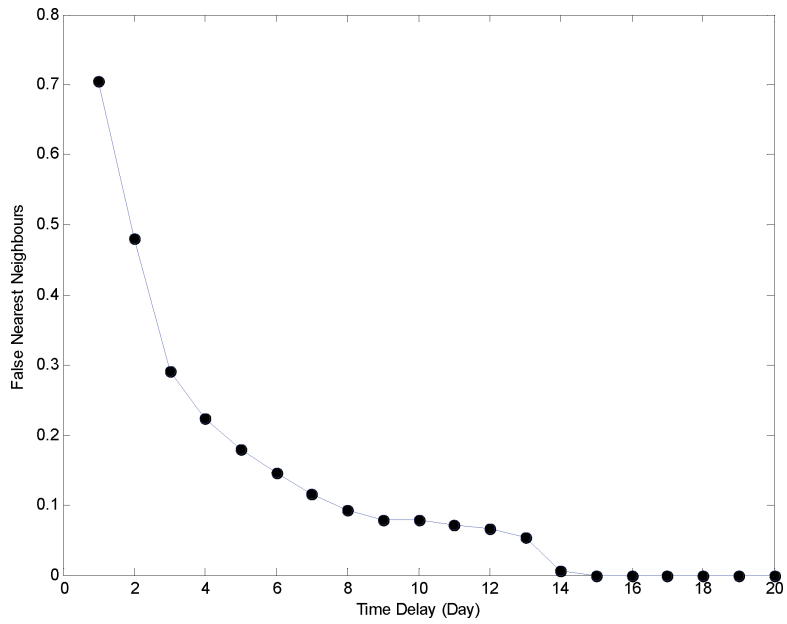


Fig. 5. False nearest neighbour at Kajang station when $\tau = 1$ and $R_T = 15$.

14353

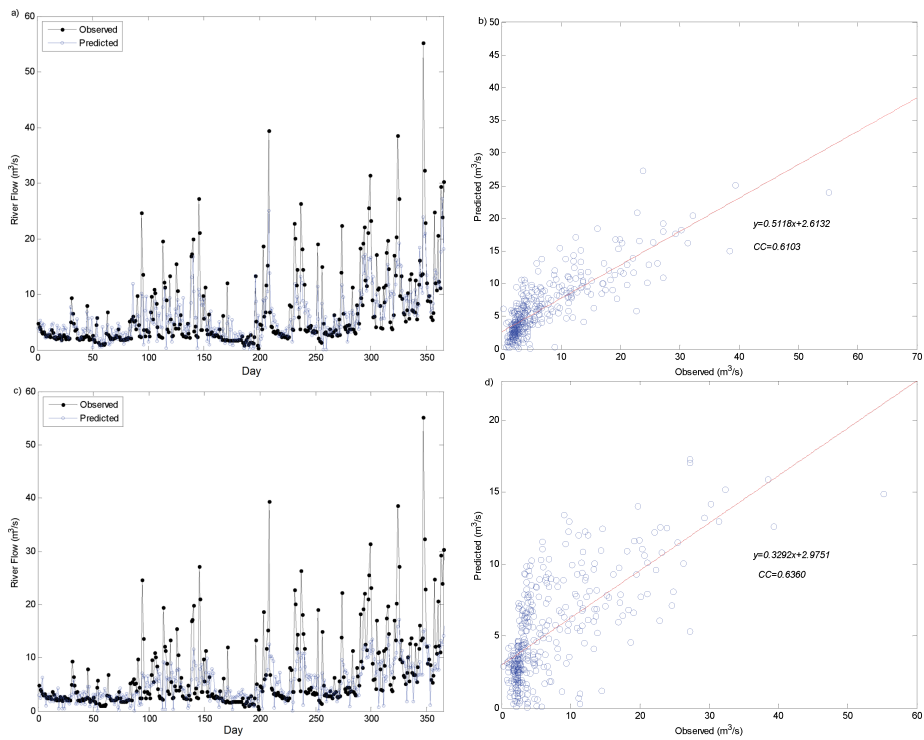


Fig. 6. Comparison of time series and scatter diagrams of prediction results and observed data for (a–b) Model I (1, 6) (c–d) and Model II (1, 14).

14354