



# Short-Term Time Series Modelling Forecasting Using Genetic Algorithm

Haviluddin<sup>1</sup>, Rayner Alfred<sup>2</sup>

<sup>1</sup>Faculty of Computer Science and Information Technology, Universitas Mulawarman, Indonesia

<sup>2</sup>Faculty of Computing and Informatics, Universiti Malaysia Sabah, Jalan UMS, 88400 Kota Kinabalu, Malaysia

The prediction analysis of a network traffic time series dataset in order to obtain a reliable forecast is a very important task to any organizations. A time series data can be defined as an ordered sequence of values of a variable at equally spaced time intervals. By analyzing these time series data, one will be able to obtain an understanding of the underlying forces and structure that are produced by the observed data and apply this knowledge in modelling for forecasting and monitoring. The techniques used to analyze time series data can be categorized into statistical and machine learning techniques. It is easy to apply a statistical technique (e.g., Autoregressive Integrated Moving Average (ARIMA)) in order to analyze time series data. However, applying a genetic algorithm in learning a time series dataset is not an easy and straightforward task. This paper outlines and presents the development of genetic algorithms (GA) that are used for analyzing and predicting short-term network traffic datasets. In this paper, the mean squared error (MSE) is taken and computed as the fitness function of the proposed GA based prediction task. The results obtained are compared with the performance of one of the statistical techniques called ARIMA. This paper is concluded by recommending some future works that can be applied in order to improve the prediction accuracy.

**Keywords:** time series, network traffic, forecasting, genetic algorithm, mean squared error (MSE).

## 1. INTRODUCTION

Time Series Analysis is used for many applications such as Economic Forecasting, Sales Forecasting, Budgetary Analysis, Stock Market Analysis, Yield and Workload Projections, Process and Quality Control, Inventory and Utilities Studies, Census Analysis, Network Monitoring and Analysis and etc. Network monitoring is not an easy task and it is a vital part of a Network Administrators job. Network Administrators are constantly striving to ensure smooth operation of their networks. In any universities, when a network is down, even for a small period of time, the teaching and research productivities within these universities would be affected and the ability to provide essential learning and teaching services would be compromised.

Network Managers need to monitor traffic movement and performance throughout the network in order to maintain smooth operation of their networks. One of the issues that network managers should pay attention to is the bandwidth usage.

Network monitoring and analysis on the bandwidth usage can be performed by using a traffic management system tool. This is important in order to avoid any network congestions in the network due to the density of traffic. The traffic management system has the ability to manage the network by setting variables of network elements, so that it presents the optimum use of real-time bandwidth data during the network data communication processes<sup>6, 11, 24</sup>.

\*Email Address: haviluddin@unmul.ac.id

These network traffic datasets are non-linear time series datasets which can be analyzed and predicted to determine the amount of usage on a daily, weekly, monthly and even yearly. There are many related works conducted to perform the analysis and prediction of these type of time series datasets in order to obtain a good forecast accuracy that includes weather, rainfall, temperature, wind speed forecasting<sup>1,14,20</sup>, financial; stock market, stock price<sup>15,16,21</sup>, tourist demand, tourist quantity<sup>5</sup> and engineering, network traffic, internet traffic<sup>4,6,8-10,12,17,24</sup>.

There are increasing interests in developing more advanced forecasting techniques in learning time series datasets (e.g., network traffic) as it will provide more information to the university's network manager for better decision making results. A Genetic Algorithm (GA) method is one of the machine learning techniques that is capable in solving the problem of forecasting a non-linear time series dataset<sup>18-20</sup>. As a result, the main objective of the paper is to outline and evaluate a genetic algorithm (GA) based prediction algorithm that is developed to model time series datasets. The ICT Universitas Mulawarman statistical data of the daily inbound outbound network traffic recorded for five days will be used as the main datasets. A step-by-step processes involved in the proposed genetic algorithm will be described clearly and the mean squared error (MSE) is taken and computed as the fitness function of the proposed GA based prediction algorithm. The rest of this paper is structured as follows. Section 2 describes the proposed genetic algorithm approach, including both time series models. The dataset is described in Section 3. In Section 4, the results of the forecasting are discussed. Finally, this paper is concluded in Section 5.

## 2. METHODOLOGY

### 2.1 The Principle of Genetic Algorithm

The basic concept of GA can be found at the University of Michigan, United States of America by John Holland in 1975 as outlined in a book entitled "Adaptation in Natural and Artificial Systems". Then, it was popularized by one of his students, David Goldberg in the 1980s. GA is an algorithm that seeks to apply an understanding of the natural evolution of problem-solving tasks. The approach taken by this algorithm is to randomly combine a wide selection of the best solutions in a set to get the next generation of the best solution based on a condition that maximizes compatibility called fitness. Then, this generation will represent improvements on the initial population<sup>7,16,19</sup>.

Based on this concept, a GA can be described as a computational abstraction of biological evolution that has worked with a population of possible solutions. A chromosome is normally used to represent the problem-solutions. The initial population that consists of a set of chromosomes is normally generated randomly. Each chromosome will go through an evaluation process using

a measure called the fitness function in which this fitness value of a chromosome will show the quality of the chromosomes in the population. Then, the next population, which is also known as offspring, is generated from the process of evolution of chromosomes through iterations called generations. A new chromosome is formed by combining a pair of chromosomes through the crossover and mutation processes<sup>2,3,18,22</sup>.

### 2.2 The Genetic Algorithm Cycle

In general, the implementation of the GA will go through a simple cycle consisting of four stages that include (1) Constructing a population consisting of several strings of chromosome called initialized population, (2) Evaluation of each string of chromosome value called using predefined fitness function, (3) Performing the selection process to get the best string of chromosome called individual selection, and (4) Genetic manipulation in order to create a new population of chromosomes called reproduction<sup>13,22</sup>. Fig. 1 illustrates the cycle of the GA implementation.

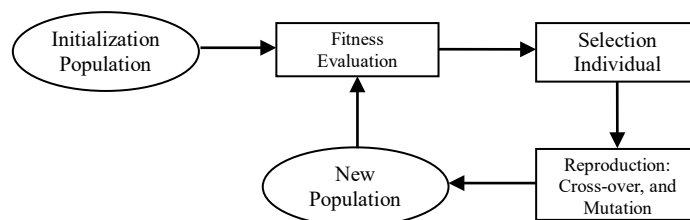


Fig.1. The Genetic Algorithm Cycle

The GA method that will be used to solve the problem of forecasting a non-linear time series dataset is as follows<sup>22</sup>;

- Step 1: Encoding schemes: Coding genes on *chromosome* using Real Number *Encoding* (RNE) and each chromosome represents a possible solution.
- Step 2: Generating Initial Population: Value of genes in each *chromosome* is generated randomly. The size of the population depends on the problem to be solved and the type of genetic operators that will be implemented.
- Step 3: Evaluation function: Individual chromosome is evaluated based on a predefined function because the value of fitness will greatly affect the performance of genetic algorithms.
- Step 4: Selection: using the method *roulette-wheel*, *random* and *tournament*.
- Step 5: Forming a New Generation: A new generation is formed by using two operators; namely crossover and mutation. The crossover is done by using a *one-point crossover*. Then, the mutation process is carried out by using the *uniform multi point mutation* criteria that is choosing a gene that will be modified based on the probability of mutation.

Step 6: Go to Step 3. This continues until the stopping criteria are met.

**2.3 Time Series Data**

A time series data can be described as a period course of action model that illuminates a variable regarding its own past and a spasmodic exacerbation term<sup>5, 23</sup>. In principle, a time series model is used to predict the current value of data,  $X_t$ , based on the data ( $X_{t-n}, \dots, X_{t-2}, X_{t-1}$ ), where  $n$  is the number of past observation and  $t$  is the current time of observation made. Time series models have been widely used for forecasting in the past four decades, with the dominance of Artificial Neural Network models. In this work, the time series data that has been taken by the software CACTI, which is one of the open source software in network management protocol will be fed into the proposed GA based prediction algorithm. Table 1 shows the inbound and outbound of the network traffic real data obtained from the Universitas Mulawarman statistical data.

Table.1. Network Traffic Real Data.

Date	Time	Inbound+ Outbound	Date	Time	Inbound+ Outbound
6/21/2013	1 0:00:00	6293000	6/23/2013	97 0:00:00	10517000
	2 0:30:00	5185000		98 0:30:00	6715000
	3 1:00:00	5404000		99 1:00:00	13109000
...	...	...	...	...	...
	47 23:00:00	12390000		143 23:00:00	7121000
	48 23:30:00	11661000		144 23:30:00	5236000
6/22/2013	49 0:00:00	8390000	6/24/2013	145 0:00:00	4528000
	50 0:30:00	7307000		146 0:30:00	3603000
	51 1:00:00	7972000		147 1:00:00	5926000
...	...	...	...	...	...
	95 23:00:00	10444000		191 23:00:00	6190000
	96 23:30:00	14530000		192 23:30:00	5969000

**2.4 Data and Implement Setting**

In order to demonstrate the process of forecasting the nonlinear time series, a four days daily network traffic data from 21 – 24 June 2013 (192 samples series data) was taken and the GA based prediction algorithm is applied. The training data was 75% (144 samples) and testing data was 25% (48 samples). Before training, the inputs and tests data will be normalized. The aim of the normalization process is to get the data with a smaller size that represents the original data without losing its own characteristics. In this experiment, a MATLAB R2013b was used to perform the process of analyzing and forecasting. The normalization formula form is as follow,

$$\bar{X} = \frac{X - X_{min}}{X_{max} - X_{min}} \tag{1}$$

where,  $\bar{X}$  : actual value of samples,  $X_{max}$  : maximum value,  $X_{min}$  : minimum value. The data after normalization show in Table 2. Based on the data outlined in Table 2, a function can be defined to learn this time series data as shown in Equation 2,

$$X_t = a_{t-n}X_{t-n}(k) + \dots + a_{t-1}X_{t-1}(k), \tag{2}$$

where  $X_t$  is the target output, the sequence of  $a_{t-n}, \dots, a_{t-1}$  is

a positive real number that represents the weights,  $X_{t-n}, \dots, X_{t-1}$  is a sequence of time series data representing the network traffic data.

Table.1. Network Traffic Data after Normalization.

Group	Period = [ $X_{t-5}, X_{t-4}, X_{t-3}, X_{t-2}, X_{t-1}$ ]	Input					Target
		$X_{t-5}$	$X_{t-4}$	$X_{t-3}$	$X_{t-2}$	$X_{t-1}$	$X_t$
Train Group	1	0.262	0.231	0.237	0.201	0.154	0.139
	2	0.231	0.237	0.201	0.154	0.139	0.164
	3	0.237	0.201	0.154	0.139	0.164	0.145
.....	.....	.....	.....	.....	.....	.....	
	144	0.232	0.213	0.187	0.251	0.246	0.211
	145	0.213	0.187	0.251	0.246	0.211	0.162
Test Group	146	0.187	0.251	0.246	0.211	0.162	0.163
	.....	.....	.....	.....	.....	.....	
	192	0.253	0.262	0.231	0.237	0.201	0.154

**2.5 Applying GA in Learning Time-Series Data**

In order to predict the network traffic using the proposed GA based prediction algorithm, the time series data must be arranged in order of time in one period. The purpose of this study is to measure the changes of data by minimizing the value of the difference between the actual and predicted values. The analysis of time series data using the proposed GA has been carried out as follows:

Step 1: Encoding schemes: Each gene in the chromosome is coded using a *real number encoding*. In other words, the chromosome is represented as a sequence of real numbers (describing a sequence of events). Where each chromosome  $x$  corresponds to a predefined fitness function  $f(x)$ .

Step 2: Generating Initial Population: Initial population process is to determine the value of each gene in the chromosome to generate random numbers. The solution (or the structure of the chromosome) for the problem is defined based on the formula,  $X_t = a_{t-n}X_{t-n}(k) + \dots + a_{t-1}X_{t-1}(k)$ , and the structure of the chromosome used to model the data shown in Table 2 will be  $[a_{t-5}, a_{t-4}, a_{t-3}, a_{t-2}, a_{t-1}]$ . The initial population size is 200.

Step 3: Evaluation function: Individual chromosome is evaluated based on a predefined function:  $X_t = a_{t-n}X_{t-n}(k) + \dots + a_{t-1}X_{t-1}(k)$ , where the values for  $X_t, X_{t-1}, X_{t-2}, X_{t-3}, X_{t-4}$  and  $X_{t-5}$  are taken from Table 2. In other words, the GA is defined to minimize the Mean Squared Error (MSE) between the  $X_t$  and  $a_{t-n}X_{t-n}(k) + \dots + a_{t-1}X_{t-1}(k)$ .

Step 4: Selection: The selection process is to establish a set of mating pool in accordance with the number of chromosomes to produce new offspring. In this experiment, three models of selection, namely the *roulette wheel*, *random* and *tournament*. In the Roulette wheel process, individual with the best fitness is not necessarily

ected at the next generation but have a better chance of being elected. This process is done by generating random numbers ( $r$ ), and then be checked against the values of  $a_1, a_2, a_3, a_4, a_5$  to the number of population so that  $r \leq p_c$ . In the Random selection process, the individual with the best fitness randomly selected from the population. In the Tournament process, the individual with the best fitness randomly selected and chosen as a parent with a size parameter value between 2 to  $N$ .

Step 5: Forming a New Generation: A new generation is formed by using two operators; namely crossover and mutation. A *one-point* method of crossover  $p_c$  with crossover rate of 0.2 and a *uniform multi point mutation* method with mutation rate of 0.005, and number of *iteration* of 100 times, and finally three selection processes will be used that includes the *roulette wheel*, *random* and *tournament* selections.

### 3. EXPERIMENTAL RESULT

This section presents the results obtained as shown in Table 3. The iteration process shows that the *roulette wheel* and *random* selections produced MSE values of 0.00497. But, the *random* selection has longer time estimation iteration than the *roulette wheel* selection which is 337.815s. Table 3 also shows that the *tournament* selection has MSE value of 0.005 and 339.632s for longest time estimation iteration. The GA based prediction algorithm has a relative long time estimation iteration process but this process depends on the set of input values. However, the MSE performance of the proposed GA has obtained good results. Fig 2, shows the graphs training and testing of three selections methods and the final MSE performance values which is 75% of the samples. In comparison, the MSE value obtained using the ARIMA (1,0,1)<sup>12</sup> is 0.00411 which is comparable with the result obtained using the GA based algorithm.

Table.3. Setting and Performance of GA.

GA Setting	Selection Method		
	Roulette Wheel	Random	Tournament
MSE	<b>0.004</b>	0.004	0.005
Time Estimation	337.744s	337.815s	339.632s
Population	200	200	200
$p_c$	0.2	0.2	0.2
$p_m$	0.005	0.005	0.005
Iteration	100	100	100

Therefore, in the first training, a *population* size of 200 real number chromosomes,  $p_c$  with *one-point* method of 0.2 and  $p_m$  with *uniform multi point mutation* of 0.005, and *iteration* value of 100 are used to produce the optimal output with low MSE values.

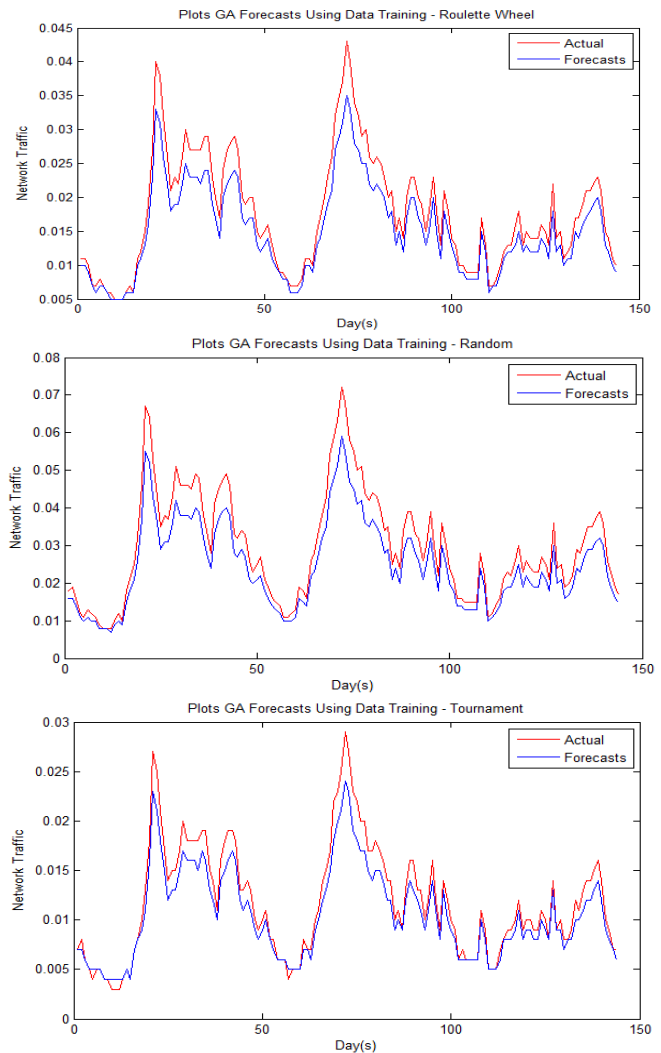


Fig.2. Plots of results for Training Data using the GA Modelling; *roulette wheel*, *random*, *tournament* selections

### 4 CONCLUSIONS

This paper examined a time series forecasting with genetic algorithms. The results shown that the proposed genetic algorithm has a pretty good value between training and testing data observed and predicted. Then, this algorithm can be used as an alternative modeling methodology in analyzing and forecasting time series data. Based on the experimental results obtained, it can be concluded that the GA setting with the *population* size of 200, *real number chromosomes*, a *one-point* method of crossover  $p_c$  with crossover rate of 0.2 and a *uniform multi point mutation* method with mutation rate of 0.005, with *roulette wheel* selection and number of *iteration* of 100 times, the time that is required to obtain an optimal output is approximately 337.815s and the obtained MSE is quite encouraging. It means that the GA setting has been able to achieve the performance goals, and comparable to the result obtained using the ARIMA method. Therefore, combining the GA with neural network method can be done in order to optimize the weights and biases or the structure for generate a higher accuracy of MSE.

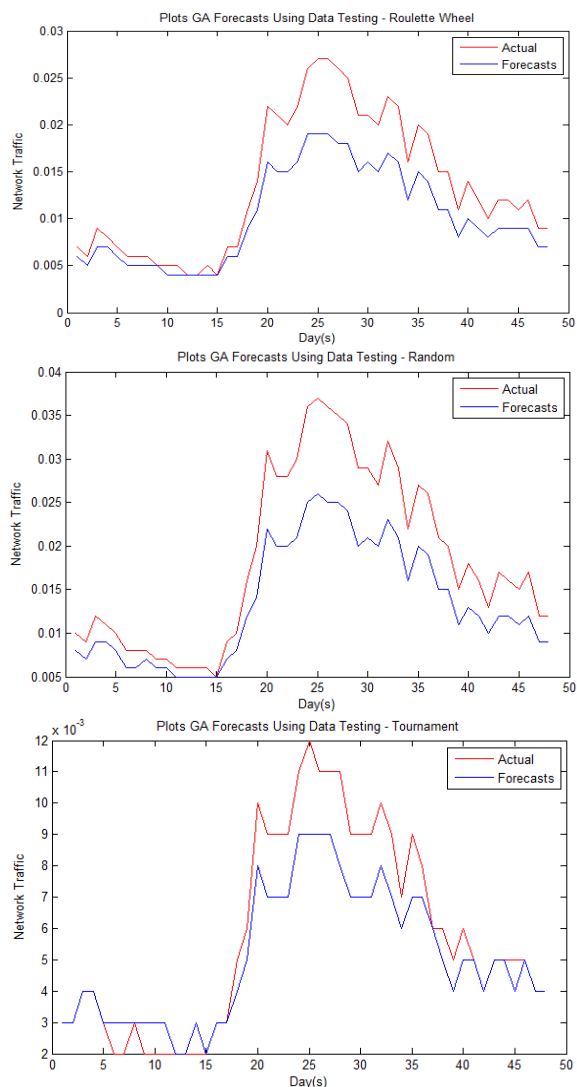


Fig.2. Plots of results for Testing Data using the GA Modelling; roulette wheel, random, tournament selections

## REFERENCES

- [1]. Kumar Abhishek, M.P. Singh, Saswata Ghosh, and Abhishek Anand, 'Weather Forecasting Model Using Artificial Neural Network', *Procedia Technology*, 4, (2012) (2012), 311 – 18.
- [2]. Rayner Alfred, 'Summarizing Relational Data Using Semi-Supervised Genetic Algorithm-Based Clustering Techniques', *Journal of Computer Science*, 6 (2010), 775-84.
- [3]. Rayner Alfred, and Dimitar Kazakov, 'Aggregating Multiple Instances in Relational Database Using Semi-Supervised Genetic Algorithm-Based Clustering Technique', in *ADBIS 2007*, ed. by Y. Ioannidis, B. Novikov and B. Rachev (Technical University of Varna: Technical University of Varna, 2007), pp. 136-47.
- [4]. Samira Chabaa, Abdelouhab Zeroual, and Jilali Antari, 'Identification and Prediction of Internet Traffic Using Artificial Neural Networks', *J. Intelligent Learning Systems & Applications*, 2, (2010), 147-55.
- [5]. Oscar Claveria, and Salvador Torra, 'Forecasting Tourism Demand to Catalonia: Neural Networks Vs. Time Series Models', *Economic Modelling*, 36 (2014), 220–28.
- [6]. Adriana C. Ferrari-Santos, José Demisio Simões da-Silva, Lília de-Sá Silva, and Milena Prado da-Costa Sene, 'Network Traffic Characterization Based on Time Series Analysis and Computational Intelligence', *Journal of Computational Interdisciplinary Sciences*, 2 (2011), pp. 197-205.
- [7]. Er. Jasmeen Gill, Er. Baljeet Singh, and Er. Shaminder Singh, 'Training Back Propagation Neural Networks with GA for Weather Forecasting', in *8<sup>th</sup> International Symposium on Intelligent Systems and Informatics*, 2010, pp. 465-69.
- [8]. Haviluddin, and Rayner Alfred, 'Comparison of Ann Back Propagation Techniques in Modelling Network Traffic Activities', in *International Conference on Science and Technology For Sustainability (ICoSTechS)* (Batam, 22 October 2014: IEEE, 2014), pp. 224-31.
- [9]. Haviluddin, and Rayner Alfred, 'Daily Network Traffic Prediction Based on Backpropagation Neural Network', *Australian Journal of Basic and Applied Sciences (AJBAS)*, 8 (24), Special 2014 (2014), 164-69.
- [10]. Haviluddin, and Rayner Alfred, 'Forecasting Network Activities Using Arima Method', *Journal of Advances in Computer Networks (JACN)*, 2, (3) September 2014 (2014), 173-79.
- [11]. Haviluddin, and Rayner Alfred, 'A Genetic-Based Backpropagation Neural Network for Forecasting in Time-Series Data', *ICSITech*, 2015.
- [12]. Haviluddin, Rayner Alfred, Joe Henry Obit, Mohd Hanafi Ahmad Hijazi, and Ag Asri Ag Ibrahim, 'A Performance Comparison of Statistical and Machine Learning Techniques in Learning Time Series Data', *Adv. Sci. Lett.* (2015), 3037-41.
- [13]. Mitchell Melanie, 'An Introduction to Genetic Algorithms', © 1996 Massachusetts Institute of Technology, 1996), pp. 1-143.
- [14]. Xue Mei Meng, 'Weather Forecast Based on Improved Genetic Algorithm and Neural Network', in *International Conference on Information Engineering and Applications (IEA) 2012* Lecture Notes in Electrical Engineering 219, © Springer-Verlag London 2013, 2013), pp. 833-38.
- [15]. Fagner A. de Oliveira, Cristiane N. Nobre, and Luis E. Zárata, 'Applying Artificial Neural Networks to Prediction of Stock Price and Improvement of the Directional Prediction Index – Case Study of Petr4, Petrobras, Brazil', *Expert Systems with Applications*, 40 (2013), 7596–606.
- [16]. Yusuf Perwej, and Asif Perwej, 'Prediction of the Bombay Stock Exchange (Bse) Market Returns Using Artificial Neural Network and Genetic Algorithm', *Journal of Intelligent Learning Systems and Applications*, 4, (2012), 108-19.
- [17]. Purnawansyah, and Haviluddin, 'Comparing Performance of Backpropagation and Rbf Neural Network Models for Predicting Daily Network Traffic', in *The 4<sup>th</sup> MICEEI*, pp. 166-69.
- [18]. A. Sedki, D. Ouazar, and E. El Mazoudi, 'Evolving Neural Network Using Real Coded Genetic Algorithm for Daily Rainfall–Runoff Forecasting', *Expert Systems with Applications*, 36, (2009) (2009), 4523–27.
- [19]. Feng Song, and Hongchun Wang, 'Hybrid Algorithm Based on Levenberg-Marquardt Bayesian Regularization Algorithm and Genetic Algorithm', in *The 2013 International Conference on Advanced Mechatronic Systems*, pp. 51-56.
- [20]. K. G Upadhyay, A. K Choudhary, and M. M Tripathi, 'Short-Term Wind Speed Forecasting Using Feed-Forward Back-Propagation Neural Network', *IJEST*, 3, No. 5 (2011), 107-12.
- [21]. Kunwar Singh Vaisla, and Ashutosh Kumar Bhatt, 'An Analysis of the Performance of Artificial Neural Network Technique for Stock Market Forecasting', *(IJCSE) International Journal on Computer Science and Engineering*, 02, (06) (2010), 2104-09.
- [22]. Cheng-Xiang Yang, and Yi-Fei Zhu, 'Using Genetic Algorithms for Time Series Prediction', in *2010 Sixth International Conference on Natural Computation (ICNC 2010)*, pp. 4405-09.
- [23]. Yang, C.-X., & Zhu, Y.-F, 'Using Genetic Algorithms for Time Series Prediction', in *2010 Sixth International Conference on Natural Computation (ICNC 2010)* © IEEE, 2010).
- [24]. Yanhua Yu, Jun Wang, Meina Song, and Junde Song, 'Network Traffic Prediction and Result Analysis Based on Seasonal Arima and Correlation Coefficient', in *2010 International Conference on Intelligent System Design and Engineering Application*, 2010), pp. 980-83.

Received: 22 September 2010. Accepted: 18 October 2010