

# Evolutionary Multiobjective Optimization Approach for Evolving Ensemble of Intelligent Paradigms for Stock Market Modeling

Ajith Abraham<sup>1</sup>, Crina Grosan<sup>2</sup>, Sang Yong Han<sup>1</sup> and Alexander Gelbukh<sup>3</sup>

<sup>1</sup>School of Computer Science and Engineering  
Chung-Ang University, Seoul 156-756, Korea

<sup>2</sup>Department of Computer Science  
Babeş-Bolyai University, Cluj-Napoca, 3400, Romania

<sup>3</sup>Centro de Investigacin en Computacin (CIC)  
Instituto Politcnico Nacional (IPN), Mexico

ajith.abraham@ieee.org, cgrosan@cs.ubbcluj.ro, hansy@cau.ac.kr,  
gelbukh@gelbukh.com

**Abstract.** The use of intelligent systems for stock market predictions has been widely established. This paper introduces a genetic programming technique (called Multi-Expression programming) for the prediction of two stock indices. The performance is then compared with an artificial neural network trained using Levenberg-Marquardt algorithm, support vector machine, Takagi-Sugeno neuro-fuzzy model and a difference boosting neural network. As evident from the empirical results, none of the five considered techniques could find an optimal solution for all the four performance measures. Further the results obtained by these five techniques are combined using an ensemble and two well known Evolutionary Multiobjective Optimization (EMO) algorithms namely Non-dominated Sorting Genetic Algorithm II (NSGA II) and Pareto Archive Evolution Strategy (PAES) algorithms in order to obtain an optimal ensemble combination which could also optimize the four different performance measures (objectives). We considered Nasdaq-100 index of Nasdaq Stock Market and the S&P CNX NIFTY stock index as test data. Empirical results reveal that the resulting ensemble obtain the best results.

## 1 Introduction

Prediction of stocks is generally believed to be a very difficult task. The process behaves more like a random walk process and time varying [20],[5]. The obvious complexity of the problem paves way for the importance of intelligent prediction paradigms [21], [6]. During the last decade, stocks and futures traders have come to rely upon various types of intelligent systems to make trading decisions [1], [2],[4],[17],[13]. In this paper, we first perform a comparison between five different intelligent paradigms. Two well-known stock indices namely Nasdaq-100 index of Nasdaq<sup>SM</sup> [11] and the S&P CNX NIFTY stock index [12] are used in experiments. Nasdaq-100 index reflects Nasdaq's largest companies across

major industry groups, including computer hardware and software, telecommunications, retail/wholesale trade and biotechnology. The Nasdaq-100 index is a modified capitalization-weighted index, which is designed to limit domination of the index by a few large stocks while generally retaining the capitalization ranking of companies. Similarly, S&P CNX NIFTY is a well-diversified 50 stock index accounting for 25 sectors of the economy [12]. It is used for a variety of purposes such as benchmarking fund portfolios, index based derivatives and index funds. The CNX Indices are computed using market capitalization weighted method, wherein the level of the Index reflects the total market value of all the stocks in the index relative to a particular base period.

Our research is to investigate the behavior of five different techniques for modeling the Nasdaq-100 and NIFTY stock market indices so as to optimize the performance indices (different error measures and correlation coefficient) and also to find an ensemble combination of these techniques in order to further optimize the performance. The five techniques used in the experiments are: an artificial neural network trained using the Levenberg-Marquardt algorithm, support vector machine [18], difference boosting neural network [16], a Takagi-Sugeno fuzzy inference system learned using a neural network algorithm (neuro-fuzzy model) [7] and Multi-Expression Programming (MEP) [14], [15]. In order to find an optimal combination of these paradigms, the task is to evolve five coefficients (one for each technique) so as to optimize the four performance measures (objectives) namely Root Mean Squared Error (RMSE), Correlation Coefficient (CC), Maximum Absolute Percentage Error (MAP) and Mean Absolute Percentage Error (MAPE). For this purpose, the problem is formulated as a multiobjective optimization problem using NSGA II and PAES. Results obtained by the evolved ensemble are compared with the results obtained by the five techniques.

We analyzed the Nasdaq-100 index value from 11 January 1995 to 11 January 2002 and the NIFTY index from 01 January 1998 to 03 December 2001. For both the indices, we divided the entire data into almost two equal parts. In section 2, we formulate the evolutionary multiobjective approach for the ensemble design followed by experimentation setup and results in Section 3. Some conclusions are also provided towards the end.

## 2 Evolutionary Multiobjective Optimization Approach for Constructing Ensemble of Intelligent Paradigms

The goal is to optimize several error measures: Root Mean Squared Error (RMSE), Correlation Coefficient (CC), Maximum Absolute Percentage Error (MAP) and Mean Absolute Percentage Error (MAPE):

$$RMSE = \sqrt{\sum_{i=1}^N |P_{actual,i} - P_{predicted,i}|}$$

$$CC = \frac{\sum_{i=1}^N P_{predicted,i}}{\sum_{i=1}^N P_{actual,i}},$$

$$MAP = \max \left( \frac{|P_{actual,i} - P_{predicted,i}|}{P_{predicted,i}} \times 100 \right),$$

$$MAPE = \frac{1}{N} \sum_{i=1}^N \left[ \frac{|P_{actual,i} - P_{predicted,i}|}{P_{actual,i}} \right] \times 100,$$

where  $P_{actual,i}$  is the actual index value on day  $i$ ,  $P_{predicted,i}$  is the forecast value of the index on that day and  $N$  = total number of days. The task is to have minimal values of RMSE, MAP and MAPE and a maximum value for CC. The objective is to carefully ensemble the different intelligent paradigms to achieve the best generalization performance. Test data is then passed through these individual models and the corresponding outputs are recorded. Suppose the daily index value predicted by DBNN, SVM, NF, ANN and MEP are  $a_n, b_n, c_n, d_n$  and  $e_n$  respectively and the corresponding desired value is  $x_n$ . The task is to combine  $a_n, b_n, c_n, d_n$  and  $e_n$  so as to get the best output value that maximizes the CC and minimizes the RMSE, MAP and MAPE values.

## 2.1 Ensemble Approach

Evolve a set of five coefficients (one for each technique) in order to obtain a linear combination between these techniques so as to optimize the values of RMSE, CC, MAP and MAPE. We consider this problem as a multiobjective optimization problem in which we want to find solution of this form:  $(coef_1, coef_2, coef_3, coef_4, coef_5)$ , where  $coef_1, \dots, coef_5$  are real numbers between -1 and 1, so as the resulting combination:

$$coef_1 * a_n + coef_2 * b_n + coef_3 * c_n + coef_4 * d_n + coef_5 * e_n$$

would be close to the desired value  $x_n$ . This means, in fact, to find a solution (an array of five real numbers) so as to simultaneously optimize RMSE, CC, MAP and MAPE. This problem is equivalent to finding the Pareto solutions of a multiobjective optimization problem (objectives being RMSE, CC, MAP and MAPE). We used the two very known Multiobjective Evolutionary Algorithm (MOEA): NSGA II and PAES. For a detailed description of these techniques please refer to [3] for NSGA II and [8], [9] and [10] for PAES.

### 3 Experiment Results

We considered 7 year’s month’s stock data for Nasdaq-100 Index and 4 year’s for NIFTY index. Our target is to develop efficient forecast models that could predict the index value of the following trade day based on the opening, closing and maximum values of the same on a given day. For the Nasdaq-100index the data sets were represented by the ‘opening value’, ‘low value’ and ‘high value’. NIFTY index data sets were represented by ‘opening value’, ‘low value’, ‘high value’ and ‘closing value’. The assessment of the prediction performance of the different paradigms and the ensemble method were done by quantifying the prediction obtained on an independent data set.

#### 3.1 Parameter Settings

We used a feed forward neural network with 4 input nodes and a single hidden layer consisting of 26 neurons. We used tanh-sigmoidal activation function for the hidden neurons. The training using LM algorithm was terminated after 50 epochs and it took about 4 seconds to train each dataset. For the neuro-fuzzy system, we used 3 triangular membership functions for each of the input variable and the 27 *if-then* fuzzy rules were learned for the Nasdaq-100 index and 81 *if-then* fuzzy rules for the NIFTY index. Training was terminated after 12 epochs and it took about 3 seconds to train each dataset. Both SVM (Gaussian kernel with  $\gamma = 3$ ) and DBNN took less than one second to learn the two data sets [2]. Parameters used by MEP are presented in Table 1.

Table 1. MEP parameter settings

Parameter		Value
Population size	Nasdaq	100
	Nifty	50
Number of iterations	Nasdaq	60
	Nifty	100
Chromosome length	Nasdaq	30
	Nifty	40
Crossover Probability		0.9
Functions set		+, -, *, /, sin, cos, sqrt, ln, lg, log <sub>2</sub> , min, max, abs

#### 3.2 Ensemble Design Using MOEA

**MOEAs Parameter Settings** The main parameters used in the experiments by the evolutionary algorithms (ensemble) are presented in Table 2.

Both NSGA II and PAES use a binary representation of solutions.

**Table 2.** Parameters used by NSGA II and PAES

Parameter	Value
Population size /Archive size	250
Number of function evaluations	125,000
Chromosome lenght	30

**Results Analysis and Discussions** Table 3 summarizes the results achieved for the two stock indices using the five intelligent paradigms (SVM, NF, ANN, DBNN, MEP) and the ensemble approach using NSGA II and PAES. Using the MOEA- ensemble approach, we obtained a population of feasible solutions. In Table 3, we present one of the solutions from the final population obtained by NSGA II and from the archive obtained by PAES respectively.

**Table 3.** Performance comparison of the results obtained by the intelligent paradigms and MOEAs (NSGA II and PAES)

	SVM	NF	ANN	DBNN	MEP	NSGA II	PAES
<b>Test results - NASDAQ</b>							
RMSE	0.0180	0.0183	0.0284	0.0286	0.021	0.01612	0.01614
CC	0.9977	0.9976	0.9955	0.9940	0.999	0.9994	0.998
MAP	481.50	520.84	481.71	116.98	96.39	94.989	94.976
MAPE	7.170	7.615	9.032	9.429	14.33	10.559	10.542
<b>TEST results – NIFTY</b>							
RMSE	0.0149	0.0127	0.0122	0.0225	0.0163	0.01317	0.01319
CC	0.9968	0.9967	0.9968	0.9890	0.997	0.999	0.999
MAP	72.53	40.37	73.94	37.99	31.7	28.50	29.75
MAPE	4.416	3.320	3.353	5.086	3.72	2.933	2.910

The ensemble obtained using NSGA II for Nasdaq is:

$$0.245357 * b_n + 0.77028 * c_n + 0.000978 * d_n + 0.00097 * e_n.$$

The ensemble obtained using PAES for Nasdaq is:

$$0.016756 * a_n + 0.242174 * b_n + 0.749939 * c_n + 0.0016604 * d_n + 0.0005028 * e_n$$

$e_n$

The ensemble obtained using NSGA II for Nifty is:

$$0.276637 * a_n + 0.220919 * b_n + 0.520039 * c_n + 0.642229 * d_n + 0.032258$$

$* e_n$

The ensemble obtained using PAES for Nifty is:

$$0.0700763 * a_n - 0.05659 * b_n + 0.4931 * c_n + 0.1541 * d_n + 0.3338 * e_n$$

The best result for Nasdaq, obtained by ensemble using NSGA II for RMSE is 0.01611. The other results are: CC = 0.999, MAP = 94.99, MAPE = 10.56

The best result for Nasdaq, obtained by ensemble using NSGA II for MAP is 94.32. The other results are: RMSE = 0.0323, CC = 0.931, MAPE = 12.80

The best result for Nasdaq, obtained by ensemble using NSGA II for MAPE is 10.417. The other results are: RMSE = 0.0171, CC = 0.993, MAP = 94.68

The best result for Nasdaq, obtained by ensemble using PAES for RMSE is 0.01611. The other results are: CC = 0.999, MAP = 95.009, MAPE = 10.58

The best result for Nasdaq, obtained by ensemble using PAES for MAP is 94.49. The other results are: RMSE = 0.0538, CC = 0.877, MAPE = 17.45

The best result for Nasdaq, obtained by ensemble using PAES for MAPE is 10.51. The other results are: RMSE = 0.0163, CC = 0.995, MAP = 94.94

The best result for Nifty, obtained by ensemble using NSGA II for RMSE is 0.01245. The other results are: CC = 0.999, MAP = 45.39, MAPE = 2.81

The best result for Nifty, obtained by ensemble using NSGA II for MAP is 24.54. The other results are: RMSE = 0.0283, CC = 0.952, MAPE = 6.49

The best result for Nifty, obtained by ensemble using NSGA II for MAPE is 2.770. The other results are: RMSE = 0.0127, CC = 0.994, MAP = 45.86

The best result for Nifty, obtained by ensemble using PAES for RMSE is 0.01256. The other results are: CC = 0.999, MAP = 34.806, MAPE = 2.824

The best result for Nifty, obtained by ensemble using PAES for MAP is 24.28. The other results are: RMSE = 0.02159, CC = 0.970, MAPE = 4.94

The best result for Nifty, obtained by ensemble using PAES for MAPE is 2.780. The other results are: RMSE = 0.01266, CC = 0.997, MAP = 35.47

The results are further graphically illustrated. In Figure 1, the values for RMSE, CC, MAP and MAPE obtained by NSGA II and PAES for Nasdaq test data are depicted. Figure 2 depicts the values for RMSE, CC, MAP and MAPE obtained by NSGA II and PAES for Nifty test data.

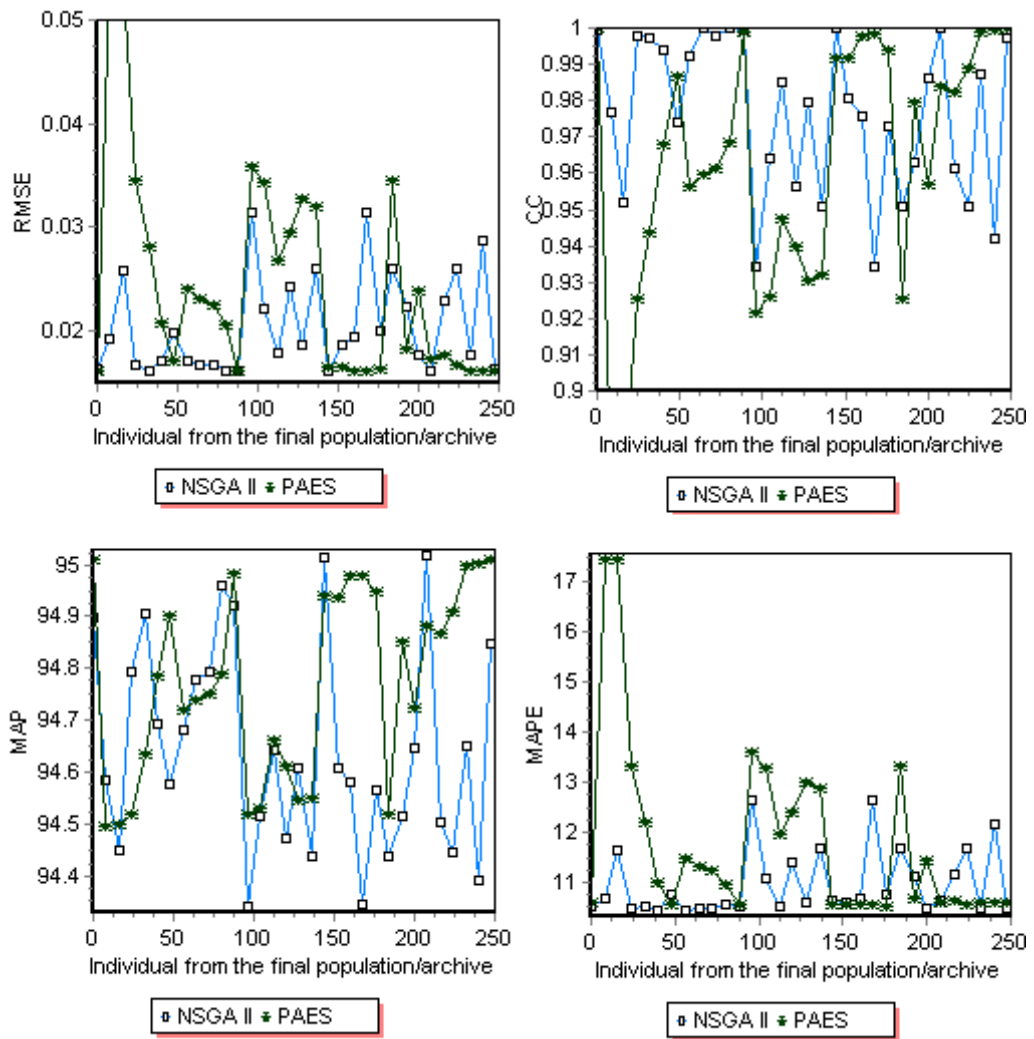
As evident from Figures 1 and 2, it is difficult to say one of the MOEAs could successfully obtain the best results for all indices. As an example, for Nifty, quality of solutions in the final population for RMSE obtained by NSGA II is better than the solutions obtained by PAES in the final archive. At the same time, for Nifty index, the quality of solutions in the final population for MAP obtained by NSGA II is comparatively poorer than the solutions obtained by PAES in the final archive.

## 4 Acknowledgements

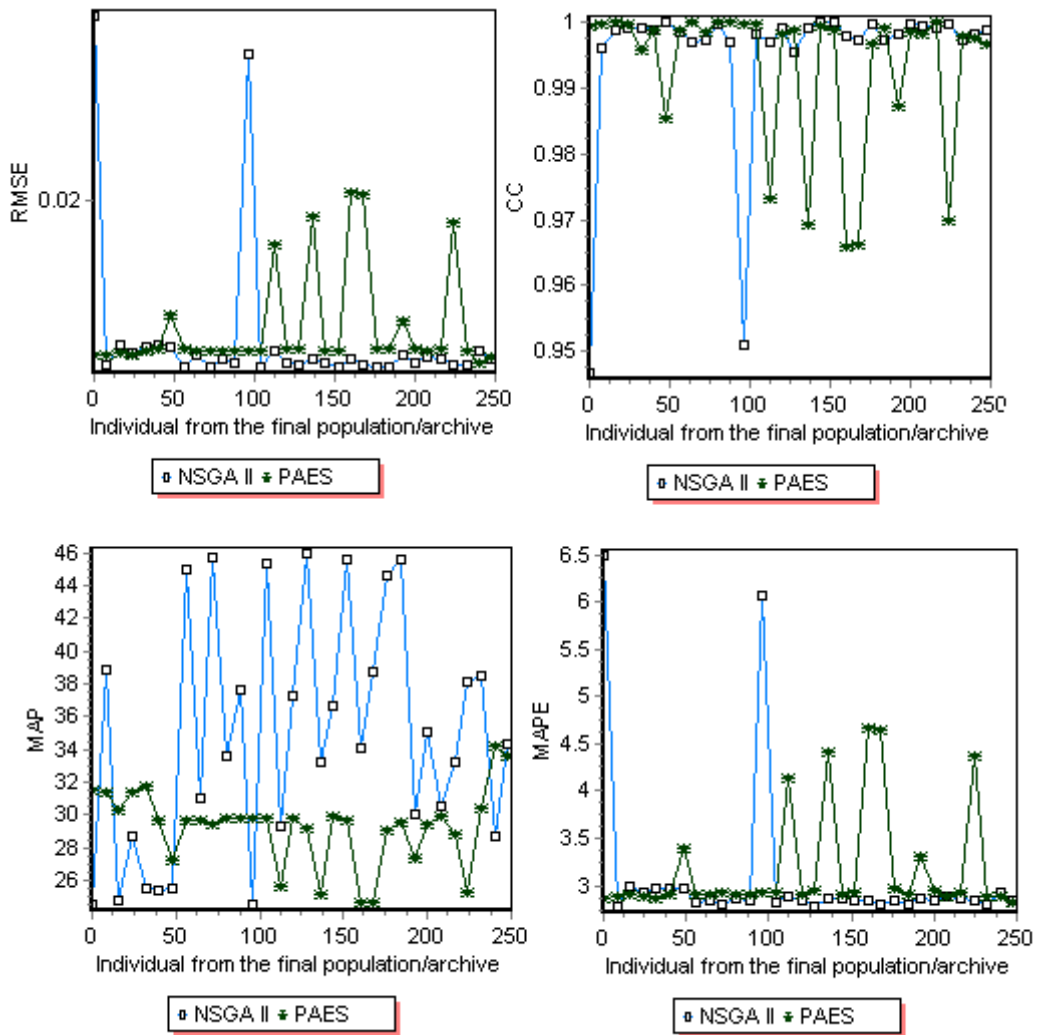
This research was supported by the MIC (Ministry of Information and Communication), Korea, under the Chung-Ang University HNRC-ITRC (Home Network Research Center) support program supervised by the IITA (Institute of Information Technology Assessment).

## 5 Conclusions

The fluctuations in the stock market are chaotic in the sense that they heavily depend on the values of their immediate forerunning fluctuations. This paper presented five techniques for modeling stock indices. Taking into account of the



**Fig. 1.** Values obtained by NSGA II and PAES for RMSE, CC, MAP and MAPE for Nasdaq test data



**Fig. 2.** Values obtained by NSGA II and PAES for RMSE, CC, MAP and MAPE for Nifty test data



No Free Lunch Theorem (NFL) [19], our research using real world stock data also reveals that it is difficult for one of the intelligent paradigms to perform well for different stock indices. Further the different intelligent paradigms were combined using an ensemble approach by two different evolutionary multiobjective algorithms (NSGA II and PAES) so as to optimize several performance measures namely RMSE, CC, MAP and MAPE. We evolved a set of coefficients in order to obtain an ensemble combination of the five techniques by applying NSGA II and PAES. Empirical results also illustrate that a combination of these techniques is very useful. The results obtained by an ensemble of these paradigms clearly outperform results obtained by the techniques individually.

## References

1. A. Abraham and A. AuYeung. Integrating Ensemble of Intelligent Systems for Modeling Stock Indices, *In Proceedings of 7th International Work Conference on Artificial and Natural Neural Networks*, Lecture Notes in Computer Science- Volume 2687, Jose Mira and Jose R. Alvarez (Eds.), Springer Verlag, Germany, pp. 774-781, 2003.
2. A. Abraham, N. S. Philip and P. Saratchandran. Modeling Chaotic Behavior of Stock Indices Using Intelligent Paradigms. *International Journal of Neural, Parallel & Scientific Computations*, USA, Volume 11, Issue (1&2) pp. 143-160, 2003.
3. K. Deb, S. Agrawal, A. Pratab and T. Meyarivan, A fast elitist non-dominated sorting genetic algorithms for multiobjective optimization: NSGA II. KanGAL report 200001, Indian Institute of Technology, Kanpur, India, 2000.
4. E.B. Del Brio, A. Miguel and J. Perote, An investigation of insider trading profits in the Spanish stock market, *The Quarterly Review of Economics and Finance*, Volume 42, Issue 1, pp. 73-94, 2002.
5. F.E.H. Tay and L.J. Cao. Modified Support Vector Machines in Financial Time Series Forecasting, *Neurocomputing* 48(1-4): pp. 847-861, 2002.
6. W.Huang , S.Goto and M. Nakamura, Decision-making for stock trading based on trading probability by considering whole market movement, *European Journal of Operational Research*, Volume 157, Issue 1, (16), pp. 227-241, 2004.
7. J.S.R. Jang, C.T. Sun and E. Mizutani. *Neuro-Fuzzy and Soft Computing: A Computational Approach to Learning and Machine Intelligence*, Prentice Hall Inc, USA, 1997.
8. J.D. Knowles and D.W. Corne, Approximating the nondominated front using the Pareto archived evolution strategies, *Evolutionary Computation*, 8(2), 149-172, 2000.
9. J.D. Knowles and D.W. Corne, The Pareto archived evolution strategy: A new baseline algorithm for Pareto multiobjective optimization. In *Congress on Evolutionary Computation (CEC 99)*, Volume 1, Piscataway, NJ, 98-105, 1999.
10. J.D. Knowles and D.W. Corne, M-PAES:A memetic algorithm for multiobjective optimization. In *Proceedings of Congress on Evolutionary Computation*, 325-332, 2000.
11. Nasdaq Stock Market<sup>SM</sup>: <http://www.nasdaq.com>.
12. National Stock Exchange of India Limited: <http://www.nse-india.com>.
13. K.J. Oh and K.J. Kim, Analyzing stock market tick data using piecewise nonlinear model, *Expert Systems with Applications*, Volume 22, Issue 3, pp. 249-255, 2002.

14. M. Oltean and C. Grosan. A Comparison of Several Linear GP Techniques. *Complex Systems*, Vol. 14, Nr. 4, pp. 285-313, 2004
15. M. Oltean and C. Grosan. Evolving Evolutionary Algorithms using Multi Expression Programming. *Proceedings of The 7<sup>th</sup> European Conference on Artificial Life*, Dortmund, Germany, pp. 651-658, 2003.
16. N.S. Philip and K.B. Joseph. Boosting the Differences: A Fast Bayesian classifier neural network, *Intelligent Data Analysis*, Vol. 4, pp. 463-473, IOS Press, 2000.
17. R. Rodriguez, F. Restoy and J.I. Pea, Can output explain the predictability and volatility of stock returns? *Journal of International Money and Finance*, Volume 21, Issue 2, pp.163-182, 2002.
18. V. Vapnik. *The Nature of Statistical Learning Theory*. Springer-Verlag, New York, 1995.
19. D.H. Wolpert and W.G. Macready. No free lunch theorem for search. Technical Report SFI-TR-95-02-010. Santa Fe Institute, USA, 1995.
20. W.X.Zhou and D.Sornette, Testing the stability of the 2000 US stock market antbubble, *Physica A: Statistical and Theoretical Physics*, Volume 348, (15), pp. 428-452 , 2005
21. J.D. Wichard, C. Merkwirth and M. Ogorzalek, Detecting correlation in stock market, *Physica A: Statistical Mechanics and its Applications*, Volume 344, Issues 1-2, pp. 308-311, 2004