# A New Stock Price Prediction Method
# Based on Pattern Classification

Zeng Zhanggui, Hong Yan and Alan M.N. Fu
School of Electrical and Information Engineering,
University of Sydney, NSW 2006, Australia
E-mail: zzeng@ee.usyd.edu.au

## Abstract

*In this paper, a new method is proposed to predict the trend of a stock price based on pattern analysis. The probabilistic relaxation algorithm is employed to classify the probability vectors of the patterns related to a considered stock price. A series of experiments have been carried out on the real stock price to verify the effectiveness of the proposed method.*

*Keywords : Stock price prediction, Pattern classification, Probabilistic relaxation*

## 1   Introduction

Advanced financial analysis methods can be useful for improving the investment performance in the financial markets.[1] Due to stochastic human-factors and the nonlinear and multivariable inputs in the financial system, it is difficult to conduct analysis and prediction using conventional techniques such as the linear filters, the martingale strategy and the risk free rate strategy.[2]

Many prediction methods[3] are confined to deriving only a future value and do not yield a probability of successful prediction of the future prices. The probability of a correct prediction makes sense when making an investment decision. It can be defined as the *credit degree*, and is a similar concept to that used in statistic and probability studies.[4] Every prediction can be assessed before knowing the actual value in terms of the *credit degree*.

This paper proposes a new prediction algorithm based on pattern classification which has a relatively high prediction accuracy and can provide the *credit degree* of prediction.

## 2   Dynamic System Model of a Stock Price

The variation of a stock price with respect to time can be described as a dynamic stochastic system with multiple inputs in which the unit of time is a day.

Let $y_0(k)$ denote the raw price of a stock at time $k$ and $y(k)$ denote the price of the stock out a moving average filter, i.e.

$$y(k) = \frac{1}{\tau} \sum_{i=1}^{\tau} y_0(k - i + 1), \; k = 1, 2, ...$$

where $\tau$ is an integer to denote the time constant of the filter. The function $y(k)$ provides a mathematical description of the stock price. To consider the variation of $y(k)$ in a simple manner, we define a characteristic function $x(k)$ associated with $y(k)$ in the following way;

$$x(k) = \begin{cases} 1 & \text{if } y(k) \geq y(k-1) \quad k = 1, 2, ... \\ 0 & \text{otherwise.} \end{cases}$$

where $k$ is the current time.

The function $x(k)$ describes two trends of a stock price, the increase and the decrease, which are the main issues in predicting a stock price.

Since $k$ is discrete, $x(k)$ refers to a characteristic sequence of the stock price. Let

$$Z(k) = [z_1, ..., z_d, ...z_D]^T,$$

where

$$z_d = x(k - d + 1), d = 1, 2, ...D; k > D,$$

denote the $D$-dimensional vector, in which the $d$th component is the price trend $x(k - d + 1)$, $d$=1, 2, ...,$D$, where T means matrix transpose. $Z(k)$ is called a lag vector[8] which represents a basic pattern. Because any element $x(.)$ has two possible values 1 or 0, there are a total of $M = 2^D$ basic patterns or lag vectors available in the $D$-dimensional space.

A lag vector can be characterized by a quantity called weighted sum $m$ of $Z(k)$, which is defined as

$$m = \sum_{d=1}^{D} 2^{(D-d)} z_d, 0 \leq m \leq M - 1 \qquad (1)$$

where $M = 2^D$.

Equation (1) is actually a transformation from a binary number to a decimal one. A typical pattern is illustrated in Figure 1.


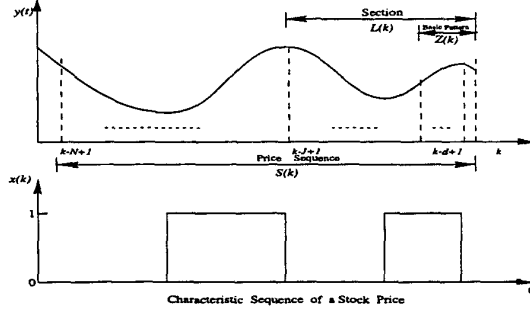
Figure 1: The pattern analysis of dynamic system of stock price.

A set of successive lag vectors is known as a section, denoted by $L(k)$, i.e.

$$L(k) = [l_1, ..., l_j, ...l_J]^T$$

$$l_j = Z(k - j + 1), j = 1, 2, ..., J,$$

where $J$ is the time interval of the section (e.g. 30 days). $Z(k)$ is considered as the current basic pattern and $L(k)$ is considered as the current section when the current time is $k$. The section $L(k)$ can be considered being a $D \times J$ matrix in which each column represents a basic pattern of a given stock price.

The probability of a basic pattern is analyzed in one section $L(k)$. Since the lag vectors $Z(k - j + 1)$ $j=1$, $2$ ,..., $J$ in the section $L(k)$ are chosen from $M$ basic patterns, they can be classified into $M$ classes according to their weighted sum $m$ of vector $Z(k - j + 1)$ given in Equation (1).

In the section $L(k)$, those vectors $Z(k - j + 1), j = 1, 2, ..., J$ which have the same weighted sum $m$ are sorted into the $m$th class. Suppose the number of vectors in the $m$th class is $\widehat{N}(m)$. Obviously, $\sum_{m=0}^{M-1} \widehat{N}(m) = J$. Then the probability of the $m$th basic pattern is given as follows;

$$p_k(m) = \frac{\widehat{N}(m)}{J}, \quad m = 0, 1, ..., M - 1, \qquad (2)$$

$$0 \leq \widehat{N}(m) \leq J.$$

Thus the $D \times J$-dimensional section $L(k)$ can be represented by the $M$-dimensional probabilistic vector $P_k$, i.e

$$P_k = [p_k(0), ..., p_k(m), ..., p_k(M - 1)]^T,$$

where $p_k(m)$ is given by Equation (2). The aim of this transformation is that the classification of a larger $D \times J$-dimensional section $L(k)$ can be realized by classifying the smaller $M$-dimensional probabilistic vector $P_k$.

A set of sections denoted by $S(k)$ is then considered. $S(k)$ is given as

$$S(k) = \{s_1, ..., s_n, ..., s_N\},$$

$$s_n = L(k - n + 1), 1 \leq n \leq N,$$

where $N$ is a long period in which the stock will be considered. The set $S(k)$ can be regarded as a $D \times J \times N$-dimensional array. Using the previous results, a set of $N$ probabilistic vectors can be used to represent the set $S(k)$;

$$\Phi_k = \{\phi_1, ..., \phi_n, ..., \phi_N\},$$

$$\phi_n = P_{(k-n+1)}, 1 \leq n \leq N,$$

where $P_k$ is known as the current probabilistic vector.

Now the problem of classification of the sections in the set $S(k)$ is reduced to the problem of classifying the probabilistic vector in $\Phi_k$. The probabilistic relaxation classification algorithm[5] which will be used to solve the problem is presented in the next section.

Suppose that the probabilistic vectors in $\Phi_k$ are classified into $M$ classes, i.e.

$$\Phi_k = \phi_k^{(0)} \bigcup \phi_k^{(1)} \bigcup \cdots \bigcup \phi_k^{(m)} \cdots \bigcup \phi_k^{(M-1)},$$

where the $m$th $(m = 0, 1, ..., M - 1)$ class is denoted as

$$\phi_k^{(m)} = \{P_{k_{m1}}, ..., P_{k_{ml}}, ..., P_{k_{mI_{km}}}\},$$

$$k - N + 1 \leq k_{ml} \leq k,$$

$$l = 1, 2, .., I_{km} \quad 1 \leq I_{km} \leq N,$$

"$\bigcup$" is the set operator "union" and $I_{km}$ is the number of probabilistic vectors in class $\phi_k^{(m)}$.

Let $\phi_k^{(c)}$ denote the class that the current probabilistic vector $P_k$ belongs to, i.e. $P_k \in \phi_k^{(c)}$ (this is known as the current class). Then, the current class $\phi_k^{(c)}$ at time $k$ can be denoted as

$$\phi_k^{(c)} = \{P_k, P_{k_{c1}}, P_{k_{c2}}, ..., P_{k_{cl}}, ..., P_{k_{cI_{kc}}}\}$$

where $I_{kc}$ denotes the total number of probabilistic vectors in the current class and $P_{k_{cl}}$ is the $l$th probabilistic vector at time $k_{cl}$. It is obvious that the probabilistic

vectors $P_{k_{cl}}$ in $\phi_k^{(c)}$ have similar characteristics to $P_k$. The corresponding sections $L(k_{cl})$ in the sequence $S(k)$ have similar dynamic characteristics to the current section $L(k)$.

The aim of this process is to actually classify the sections $L(k-n+1), n = 1, 2, ..., N$ in the sequence $S(k)$ into $M$ classes by classifying the probabilistic vectors $P_{(k-n+1)}, n = 1, 2, .., N$ in the set of probabilistic vectors $\Phi_k$. We retain the current class $\phi_k^{(c)}$ and remove other $M - 1$ classes.

The relative future value $x(k_{cl} + 1)$ of the sections $L(k_{cl}), l = 1, 2, ..., I_{kc}$ is analysed next. In terms of the current class $\phi_k^{(c)}$, these relative future trends $x(k_{cl} + 1)$ $l = 1, 2, ..., I_{kc}$ which are related to the $P_{k_{cl}}$ of the class sections $L(k_{cl})$ can be found according to the historic characteristic sequence. Thus a future value vector $X(k + 1)$ can be developed as

$$X(k + 1) = [X_{k_{c1}}, ..., X_{k_{cl}}, ..., X_{k_{cl_{kc}}}]^T,$$

$$X_{k_{cl}} = x(k_{cl} + 1), = 1, 2, ..., I_{kc}.$$

Note that the element $x(k_{cl} + 1)$ is 1 or 0. Assume $N_1(k)$ denotes the total number of the elements $x(k_{cl} + 1)$ whose value is 1 in the vector $X(k + 1)$, then the probability of the future value taking 1 is defined as

$$p_1(k) = \frac{N_1(k)}{I_{kc}} \tag{3}$$

while

$$p_0(k) = 1 - p_1(k)$$

is the probability of the future value taking 0.

Obviously, the trend which, based on the historic data of similar characteristics, has larger probability should be taken as the prediction value $\tilde{x}(k + 1)$ at $k$ tick .

Hence, the dynamic system prediction model of a stock price can be proposed as follows;

$$\tilde{x}(k + 1) = \begin{cases} 1 & \text{if } p_1(k) \geq p_0(k) \\ 0 & \text{otherwise.} \end{cases}$$

$Crd = p_1(k)$ if $\tilde{x}(k + 1) = 1$,
$Crd = p_0(k)$ if $\tilde{x}(k + 1) = 0$,
where Crd is the reliability of the predictive value.

## 3  Prediction Method

To predict future trends of a stock price, the set of probabilistic vectors $\Phi_k$ are first classified using a probabilistic relaxation classification algorithm.[5] Let $p_i^{(w)}(m)$ denote the $m$th component of the probabilistic vector $P_i^{(w)}$ ($i = k, k - 1, ..., k - N + 1$) at the $w$th iteration and the initial $p_i^{(0)}(m)$ is given by Equation (2). Then,

using the probabilistic relaxation sheme, the updated equation is as follows;

$$p_i^{(w+1)}(m) = \frac{p_i^{(w)}(m) + p_i^{(w)}(m)\nu_i^{(w)}(m)}{R_i^{(w)}} \tag{4}$$

$$m = 0, 1, ..., M - 1, \ i = k, k - 1, ..., k - N + 1,$$

$$w = 0, 1, 2, ...$$

where

$$R_i^{(w)} = 1 + \sum_{m=0}^{M-1} p_i^{(w)}(m)\nu_i^{(w)}(m), \tag{5}$$

$$\nu_i^{(w)}(z) = p_i^{(w)}(z) - s_i^{(w)}(z), \tag{6}$$

$$s_i^{(w)}(z) = \sum_{j=1}^{N} d_{ij} \sum_{l=0}^{M-1} c_{ij}(z, l)p_i^{(w)}(i) \tag{7}$$

where $c_{ij}(z, l)$ is the compatibility coefficient when the $i$th basic pattern belong to the $z$th class and the $j$th basic pattern belong to the $l$th class; $d_{ij}$ represents the influence coefficient of the $i$th basic pattern from the $j$th basic pattern. $c_{ij}(z, l)$ and $d_{ij}$ satisfy

$$\sum_{z=0}^{M-1} c_{ij}(z, l) = 1,$$

$$0 \leq c_{ij}(z, l) \leq 1, z, l = 0, 1, ..., M - 1,$$

$$\sum_{j=1}^{N} d_{ij} = 1,$$

$$0 \leq d_{ij} \leq 1, \ i, j = 1, 2, ..., N.$$

For simplicity, we set $d_{ij} = \frac{1}{N}$ and $c_{ij}(z, l) = \frac{1}{M}$ for $i, j = 1, 2, ..., N$ and $z, l = 0, 1, ..., M - 1$.

Using Equations (4), (5), (6) and (7), the set of probabilistic vectors $\Phi_k$ are classified by the probabilistic relaxation scheme.

After approximately 25 iterations, the dynamic system given by Equations (4), (5), (6) and (7) arrives at a stable state, i.e. each $P_i$ ($i = k, k - 1, ..., k - N + 1$) converges to a basic unit vector $\Lambda^m$ of $M$-dimension space, where $\Lambda^m$ is the basic unit vector with the $m$th component being 1. Therefore, $P_i^{(0)}$ ($i = k, k - 1, ..., k - N + 1$)

in $\Phi_k^{(0)}$ are classified into the same class if they converge with the same basic unit vector.

Next, we pick up the current class $\phi_k^{(c)}$ with the probabilistic vector $P_k$ in the vector group $\Phi_k$ from the classification results. Other classes are not regarded correlative to the current probabilistic vector $P_k$ or section $L(k)$. Hence, they are not useful in predicting the future trends $\tilde{x}(k+1)$.

The future trend vector $X(k+1)$ can now be obtained. Based on the future trend vector $X(k+1)$ and using Equation (3), we can compute the probabilities $p_1(k)$, $p_0(k)$ of future trends and predict the future trend $\tilde{x}(k+1)$ as proposed previously.

In addition, $p_1(k)$ or $p_0(k)$ indicate the probability of a correct prediction. Thus they are respectively defined as the reliability of prediction values.

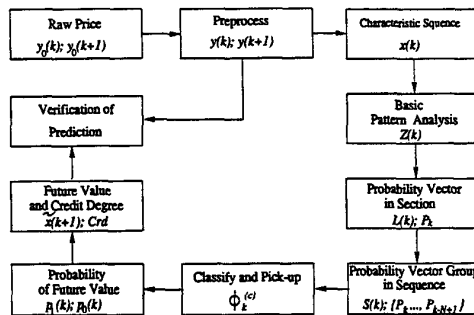Figure 2 depicts this prediction process.



Figure 2: The prediction process based on pattern classification.

# 4  Experiments

In order to verify the effectiveness of the proposed method in predicting the stock price, a share of Australia's Quantas Airline is used in othe experiments and the share price data are collected from 1 Jan. 1995 to 8 Sept. 1998. A typical sample of the stock price is illustrated in Figure 3.
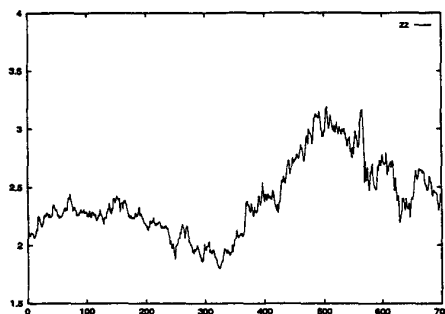


Figure 3: The raw share price of Qantas Airlines.

The experiment is performed with the parameters; N=50, J=15, D=2. To improve the prediction results, a moving average filter of two days is employed to preprocess the price information. The results are given in Table 1 and shown in Figure 4.

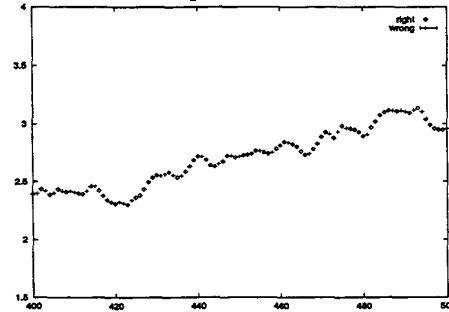| Time(x100 days) | 1 | 2 | 3 | 4 | 5 | 6 | 7 | Average |
|---|---|---|---|---|---|---|---|---|
| Correct prediction | 56% | 70% | 75% | 73% | 68% | 69% | 68% | 68% |

Table 1: Prediction performance after filtering.



Figure 4: The prediction results of the share price of Qantas Airlines.

# 5  Conclusion

In this paper, a new prediction method based on a pattern classification scheme is proposed. The experiments are carried out on Qantas Airline share prices to verify the method. The prediction correction is much better with filtered data than with raw data. However, attention should be paid to the parameters of the filter in case valuable information is lost.

With different parameters, different results are achieved. There is a best set of parameters which produce the best results. This is not surprising as these parameters match the dynamic system best. Coincident to one deduction of the Random Walk Hypothesis[6] is that the most recent value should be used to predict the future trends of a stock price. The best results are achieved with the following parameters; $N = 50, J = 15, D = 2$. Poor results are achieved if N is too large.

| Algorithm | ANN | Martingale | Linear filter | Risk free | Optimum |
|---|---|---|---|---|---|
| Correct prediction | 56.5% | 56.5% | 54.2% | - | 100% |
| Profit(from $100) | 216 | 190 | 187 | 51 | 927 |
| Annual return | 29.5% | 26.0% | 25.4% | 7.0% | 4126.5% |

Table 2: Performance of different investment strategies.

Compared with other prediction methods as shown in Table 2,[2] this new method achieves a higher correct

prediction level and requires less historic data. In addition, the new method is able to produce a prediction credit degree which would be very useful when making selling or buying decisions.

# References

[1] Wunnicke, D. R. Wilson, B. Wunnicke, *Practical Techniques of Financial Engineering*, New York: Wiley 1992

[2] D. T. Pham and L. Xing, *Neural Networks for Identification, Prediction and Control*, London, Berlin, Heidelberg, 1995

[3] D. Kil, *Stock Price Prediction Using An Integrated Pattern Recognition Paradigm, Neural Networks in Financial Engineering*, World Scientific, 270-283, 11-13 October, 1995

[4] X. Zou, H. J. Qing, X. S. Qian and P. C. Yi, *Probability and Statistics*, China Advanced Education, 1985

[5] M. N. Fu and H. Yan, "A New Probabilistic Relaxation Method Based on Probability Space Partition", Pattern Recognition, 1905-1917, 1997

[6] R. J. Tewels, E. S. Bradley and T. M. Tewels, *The Stock Market*, John Wiley and Sons Inc., 1992