# Forecasting stock price direction using an EMD-KPCA-based SVM

Anass Nahil1, Abdelouahid Lyhyaoui2

*1Laboratory of Innovative Technologies, AbdelmalekEssaadi University, Tangier, 90000, Morocco,*
*2Laboratory of Innovative Technologies, AbdelmalekEssaadi University, Tangier, 90000, Morocco*
*Corresponding Author: Anass Nahil*

**ABSTRACT:**
*Stock trend prediction is regarded as one of the most challenging tasks of financial time series prediction. Due to the non-linear and non-stationary characteristics of stock market time series, it is generally difficult to model and predict such series by a single forecasting model. In this paper, a novel hybrid model, which combine Empirical Mode Decomposition (EMD), feature selection using Kernel Principal Component Analysis (KPCA) and Support Vector Machine (SVM), is proposed for predicting direction of stock's movement. First, the original stock price time series is decomposed into a set of sub-series by EMD. Second, a feature selection process (KPCA) is introduced to constitute the relevant and informative features. Finally, a predictive model (SVM) is established using these selected features. A comprehensive parameters setting is performed to improve its prediction performances. The effectiveness of the proposed model is assessed on the data sets of six Moroccan banks listed in the Casablanca Stock Exchange. Compared with the single SVM, traditional EMD-based SVM and traditional KPCA-based SVM, the experimental results show that the proposed model has the best performance and is suitable for the stock price direction prediction.*
**KEYWORDS:** *Empirical mode decomposition, Kernel principal component analysis, Stock pricedirection, Support vector machine.*

## I.      INTRODUCTION

Stock price prediction is one of the most important and challenging subjects of financial markets. Therefore, many researches have been achieved in this field. However, because of the non-linearity and non-stationarity of the stock market, conventional statistical methods [1,2] used to forecast stock price are not sufficiently effective. Consequently, many machine learning techniques, such as Artificial Neural Networks (ANN) or Support Vector Machine (SVM), are used to improve the prediction accuracy. These machine learning techniques mainly focus on two purposes. One is the prediction of the future price based on the historical prices and the technical indicators [3]-[16], and the other is the forecasting of the direction of stock price [17]-[22].

There is a rich literature on prediction of stock markets with machine learning techniques. Here, we only choose to discuss several related works. One could find more in the references therein. Different techniques have already been explored for stock price prediction. One of the best performing algorithms in the financial world appears to be SVM [3,17]. Other well-known techniques are Neural Networks [28] and Decision Trees [7]. SVM was used by Kim (2003) to predict the direction of daily stock price change in the Korea composite stock price index (KOSPI). Twelve technical indicators were selected to make up the initial attributes. This study compared SVM with back-propagation neural network (BPN) and case-based reasoning (CBR). Huang et al. (2005) [17] investigated the predictability of financial movement direction with SVM by forecasting the weekly movement direction of NIKKEI 225 index. They compared SVM with Linear Discriminant Analysis, Quadratic Discriminant Analysis and Elman Back-propagation Neural Networks. The experiment results showed that SVM outperformed the other classification methods. Ou and Wang (2009) [18] used total ten data mining techniques to predict price movement of Hang Seng index of Hong Kong stock market. The approaches include Linear discriminant analysis (LDA), Quadratic discriminant analysis (QDA), K-nearest neighbor classification, Naive Bayes based on kernel estimation, Logit model, Tree based classification, neural network, Bayesian classification with Gaussian process, Support vector machine (SVM) and Least squares support vector machine (LS-SVM). It was evident from the experimental results that the SVM and LS-SVM generated superior predictive performance among the other models. Mantri et al. (2010) [23] calculated the volatilities of Indian stock markets using GARCH, EGARCH, GJR-GARCH, IGARCH & ANN models. This study used Fourteen years of data of BSE Sensex & NSE Nifty to calculate the volatilities. The experiment results showed that there is no difference in the volatilities of Sensex & Nifty estimated under the GARCH, EGARCH, GJR GARCH, IGARCH & ANN models. Tsai et al. (2011) [5] investigated the prediction performance that utilizes the classifier ensembles method to analyze stock returns. The hybrid methods of majority voting and bagging were considered. Moreover, performance using two types of classifier ensembles were compared with those using single baseline classifiers (neural networks, decision trees, and logistic regression). The results indicated that multiple classifiers outperform single classifiers in terms of prediction

accuracy and returns on investment. Nair et al. (2011) [6] predicted the next day's closing value of five international stock indices using an adaptive Artificial Neural Network based system. The system adapted itself to the changing market dynamics with the help of genetic algorithm which tunes the parameters of the neural network at the end of each trading session. Mishra et al. (2011) [24] tested for the presence of nonlinear dependence and deterministic chaos in the rate of returns series for six Indian stock market indices. The result of analysis suggested that the returns series did not follow a random walk process. Rather it appeared that the daily increments in stock returns were serially correlated and the estimated Hurst exponents were indicative of marginal persistence in equity returns. Sun and Li (2012) [25] proposed new Financial distress prediction (FDP) method based on SVM ensemble. The algorithm for selecting SVM ensemble's base classifiers from candidate ones was designed by considering both individual performance and diversity analysis. Experimental results indicated that SVM ensemble was significantly superior to individual SVM classifier. Liu and Wang (2012) [26] investigated and forecasted the price fluctuation by an improved Legendre neural network by assuming that the investors decided their investing positions by analyzing the historical data on the stock market. Garg et al. (2013) [27] worked to analyze the effect of three model selection criteria across two data transformations on the performance of GP while modeling the stock indexed in the New York Stock Exchange (NYSE). It was found that FPE criteria have shown a better fit for the GP model on both data transformations as compared to other model selection criteria. Ballings et al. (2015) [22] benchmarked ensemble methods against single classifier models in predicting stock price direction.

In this paper, we will combine signal decomposition algorithms, feature extraction and machine learning techniques to predict the direction of the movement of six Moroccan banks listed in the Casablanca Stock Exchange. Ten indicators constructed from historical data are selected and used as input in the experiment. After decomposing the original time series into a set of subseries, the initial input-output pairs are constructed from all the sub-series and the original series (EMD). Then a feature selection process is introduced to constitute relevant and informative features (KPCA). Finally, a predictive model is established using these selected features (SVM).

The reminder of this paper is organized as follows. Sections 2, 3 and 4 present the models used in the experiments (EMD, KPCA and SVM). Section 5 describes the research data. Section 6 exposes the proposed hybrid model for stock price direction prediction. The experimental procedure and results are presented and discussed in Section 7. Finally, conclusions are drawn in Section 8.

## II. EMPIRICAL MODE DECOMPOSITION

Empirical Mode Decomposition (EMD) is a data-driven algorithm which can deal with non-linear and non-stationary signals. The main idea of EMD is to decompose a complicated signal into a finite and small number of oscillatory modes based on the local characteristic time scale by itself. Each oscillatory mode is expressed by an intrinsic mode function (IMF), which has to satisfy the following two conditions [29]: (1) In the whole dataset, the number of extrema and the number of zerocrossings must either be equal or differ at most by one; (2) At any point, the mean between the upper and lower envelopes, which are defined by the local maxima and minima, must be zero. Let $y(t)$ be a given time series, the computational steps of EMD is described as follows [29]:

- **Step 1**: Identify all the local extrema of $y(t)$ and then connect all the local maxima and local minima with an interpolation method (e.g. cubic spline) to generate an upper envelope $y_{up}(t)$ and a lower envelope $y_{low}(t)$, respectively;

- **Step 2**: Calculate the mean envelop $m(t)$ from the upper and lower envelopes $m(t) = 1/2\left[y_{up}(t) - y_{low}(t)\right]$ and then subtract it from the original time series to obtain a detailed component $d(t) = y(t) - m(t)$;

- **Step 3**: Check whether $d(t)$ is an IMF. If $d(t)$ is an IMF then set $c(t) = d(t)$ and meantime replace $y(t)$ with the residual $r(t) = y(t) - c(t)$. Otherwise, replace $y(t)$ with $d(t)$ and repeat Steps 1-2 until the following termination criterion is satisfied:

$$\sum_{t=1}^{l} \frac{\left[d_{j-1}(t) - d_j(t)\right]^2}{\left[d_{j-1}(t)\right]^2} \leq \delta (j = 1,2,\dots; t = 1,2,\dots,l) \qquad (1)$$

where $l$ is the length of the signal and $j$ denotes the number of iterative calculation. A typical value for $\delta$ is usually set between 0.2 and 0.3;

- **Step 4**: Repeat Steps 1-3 until all the IMFs and the residual are obtained. Finally, the original time series $y(t)$ can be decomposed as follows:

$$y(t) = \sum_{i=1}^{n} c_i(t) + r_n(t) \qquad (2)$$

where $c_i(t)(i = 1,2,\dots,n)$ represents different IMFs and $r_n(t)$ is the final residual.

## III. KERNEL PRINCIPAL COMPONENT ANALYSIS

Principal Component Analysis (PCA) is a linear dimension reduction method, which can only handle linearly correlated data. KPCA is developed to over-come such weakness by conducting nonlinear process monitoring. KPCA projects original data space into a high-dimensional feature space before implementing PCA operation.

The input matrix can be obtained as $X = [x_1, x_2, \dots, x_n]^T \in R^{n \times m}$ where $x_i$ is the observation vector at time i. Mapping from dataspace to feature space is implemented with the following nonlinear mapping function $\phi(.)$:

$$R^m \xrightarrow{\phi(.)} F^h \qquad (3)$$

The observation vector $x_i$ becomes $\phi(x_i)$ in feature space. The covariance matrix is then represented as:

$$S^F = \frac{1}{n}\sum_{i=1}^{n} \phi(x_i)\phi(x_i)^T \qquad (4)$$

where $\phi(x_i)$ is scaled as zero mean. However, $\phi(x_i)$ cannot be acquired directly. The eigenvalue decomposition [14,15] in kernel space can be obtained as:

$$\lambda v = S^F v = \left(\frac{1}{n}\sum_{i=1}^n \phi(x_i)\phi(x_i)^T\right) v = \frac{1}{n}\sum_{i=1}^n \langle \phi(x_i), v\rangle \phi(x_i) \qquad (5)$$

where $\lambda$ and $v$ are eigenvalue and eigenvector of $S^F$, respectively, and $\langle .,.\rangle$ denotes inner product. Each $v$ with $\lambda \neq 0$ lies in the span of training data in kernel space. The coefficients $\alpha_{i\in\{1,2,...,n\}}$ exist [14,15], such that:

$$v = \sum_{i=1}^n \alpha_i \phi(x_i) \qquad (6)$$

By combining equations (3) and (4) through multiplication from both sides with $\phi(x_k)$, the equation can be written as:

$$\lambda \sum_{j=1}^n \alpha_j \langle \phi(x_j), \phi(x_k)\rangle = \frac{1}{n}\sum_{i=1}^n \sum_{j=1}^n \alpha_j \langle \phi(x_j), \phi(x_i)\rangle \langle \phi(x_k), \phi(x_i)\rangle \quad (7)$$

Thus, the solution to decompose the covariance matrix is to obtain the inner product in nonlinear space. Kernel matrix K, which is denoted by $n \times n$, is defined as follows:

$$K_{ij} = \langle \phi(x_i), \phi(x_j)\rangle \qquad (8)$$

The inner product can be acquired by introducing a kernel function. Kernel functions have three kinds, namely, polynomial, sigmoid, and Gaussian kernels. In this paper, Gaussian kernel is adopted as follows:

$$K(x,y) = \exp\left(\frac{-|x-y|^2}{\sigma}\right) \qquad (9)$$

where $\sigma$ is a constant. Kernel matrix should be centralized with:

$$K - I_n K - K I_n + I_n K I_n \rightarrow K, I_n = \frac{1}{n}I_1 \in R^{n\times n} \qquad (10)$$

Equation (5) is rewritten as:

$$\lambda\alpha = \frac{1}{n}K\alpha, \alpha = [\alpha_1, \alpha_2, ..., \alpha_n]^T \qquad (11)$$

The selection of kernel principal component (PC) is based on $\lambda$ value. PCs that correspond to large $\lambda$ should be kept in PC space, and PCs with small $\lambda$ are placed into residual space. The $j^{th}$ extracted PC is calculated by mapping training data $\phi(x)$ in feature space onto eigenvector $v_j$ as follows:

$$t_j = \langle v_j, \phi(x)\rangle = \sum_{i=1}^n \alpha_i^j \langle \phi(x), \phi(x_i)\rangle, j = 1,2,...,k \qquad (12)$$

where k is the number of principal components extracted in principal component space.

## IV.  SUPPORT VECTOR MACHINE

SVM is popular non-linear model developed by Vapnik and his co-workers based on statistical learning theory [30]. SVM has been used extensively for time series forecasting [31].

Let's consider a set of data $(x_1, y_1), (x_2, y_2), ..., (x_n, y_n)$, where $x_i \in R_n$ is feature vector and $y_i \in -1, +1$ class vector. The two classes are separated by a hyperplane.

Hyperplane, g(x) which accurately separates the data into its corresponding classes is given by:

$$W^t x_i + b = 0 \qquad (13)$$

where W is a vector with real values and b is a constant. Their values should be derived in such that the unknown data are classified accurately. This is possible by maximizing the separation margin between the classes. This can be achieved by maximizing $W^t x_i + b = 0$. For maximizing m, W should be minimized. For a given set of linearly separable data, this can be formulated as a quadratic optimization problem:

$$\begin{cases} \min \frac{1}{2}\|W\|^2 \\ \text{subject to} \quad y_i(W^T x_i + b) \end{cases} \qquad (14)$$

It can also be solved in terms of Lagrange multipliers, $a_i$ :

$$\begin{cases} \max L(\alpha) = \sum_{i=1}^N \alpha_i + 2^{-1}\sum_{i=1}^N \sum_{j=1}^N \alpha_i \alpha_j y_i y_j \langle x_i, x_j\rangle \\ \quad \text{subject to } \alpha_i x_i = 0 \end{cases} \qquad (15)$$

where $\langle x_i, y_j\rangle$ is an inner product.

In this case, support vectors are designated by $\alpha_i^*$, where $\alpha_i^* > 0$ :

$$W^* = \sum i = 1N\alpha_i^* x_i y_i \qquad (16)$$
$$b^* = y_{sv} - \sum i = 1N\alpha_i^* y_i\langle x_i, x_s v\rangle \qquad (17)$$

The optimal function is given by:

$$f(x) = \text{sign}(\sum_{i\in SV} \alpha_i^* y_i\langle x_i, x_{sv}\rangle + b^*) \qquad (18)$$

In case if data cannot be separated linearly then (W*, b*) does not exist. Then the input data should be mapped first from n-dimensional space ($R^n$) to higher dimensional feature space ($F^m$):

$$\Phi: R^n \rightarrow F^m, x_i \rightarrow \phi(x_i) \qquad (19)$$

Then another function f is used to map the data from feature space to decision space ($Y^2$):

$$f: F^m \rightarrow Y^2, \Phi(x_i) \rightarrow f(\phi(x_i)) \qquad (20)$$

Kernel functions are used for mapping nonlinearly separable data into higher dimensional feature space. Thus:

$$f(x) = sign(\sum_{i\in SV} \alpha_i^* y_i . k(x_i, x_{sv}) + b^*) \qquad (21)$$

where $k(x_i, x_{sv})$ is kernel function.

## V.  RESEARCH DATA

This section describes the research data and the selection of predictor attributes. The research data used in this study is the daily stock price of 6 Moroccan banks listed in the Casablanca Stock Exchange. The entire data set covers the period from January 2, 2012 to June 30, 2016 (Fig. 1) The total number of cases is 1118 trading days.
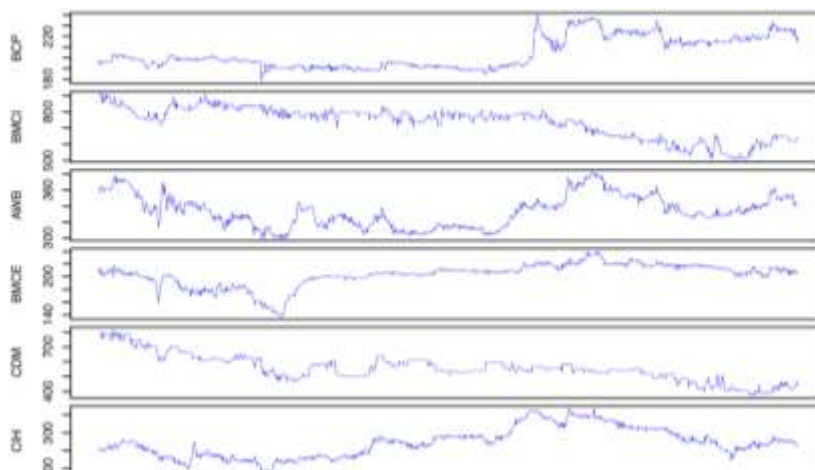
**Figure 1:** The closing price of the major Moroccan banks from January 2, 2012 to June 30, 2016

The total number of cases with increasing direction is 4493 (67%)while the number of cases with decreasing direction is 2209(33%). The number of cases for each bank is given in Table 1.

**Table 1:** The number of increase and decrease cases percentage for each bank

| Banks | Increases count | Increases percentage | Decreases count | Decrease percentage |
|-------|-----------------|----------------------|-----------------|---------------------|
| BCP | 687 | 62 | 430 | 38 |
| BMCI | 851 | 76 | 266 | 24 |
| AWB | 638 | 57 | 479 | 43 |
| BMCE | 695 | 62 | 422 | 38 |
| CDM | 957 | 86 | 160 | 14 |
| CIH | 665 | 60 | 452 | 40 |
| Total | 4,493 | 67 | 2,209 | 33 |

Ten technical indicators for each case were used as input variables. Many fund managers and investors in the stock market generally accept and use certain criteria for technical indicators as the signal of future market trends (Kim, 2003) [3]. A variety of technical indicators are available. Some technical indicators are effective under trending markets and others perform better under no trending or cyclical markets. In the light of previous studies, it is hypothesized that various technical indicators may be used as input variables in the construction of prediction models to forecast the direction of movement of the stock price index. We selected ten technical indicators as feature subsets by the review of domain experts and prior researches (Kim, 2003 [3]; Kim & Han, 2000 [4]). Table 2 summarizes the selected technical indicators and their formulas.

**Table 2 :** Selected technical indicators and their formulas

| Name of indicators | Formulas |
|--------------------|----------|
| Simple 10-day moving average | $\dfrac{C_t + C_{t-1} + \cdots + C_{t-9}}{n}$ |
| Weighted 10-day moving average | $\dfrac{(10)C_t + (9)C_{t-1} + \cdots + C_{t-9}}{n + (n-1) + \cdots + 1}$ |
| Momentum | $C_t - C_{t-9}$ |
| Stochastic K% | $\dfrac{C_t - LL_{t-(n-1)}}{HH_{t-(n-1)} - LL_{t-(n-1)}} \times 100$ |
| Stochastic D% | $\dfrac{\sum_{i=0}^{n-1} K_{t-i}}{10}\%$ |
| Relative Strength Index (RSI) | $100 - \dfrac{100}{1 + \left(\sum_{i=0}^{n-1} UP_{t-i}/n\right) / \left(\sum_{i=0}^{n-1} DW_{t-i}/n\right)}$ |
| Moving Average Convergence Divergence (MACD) | $MACD(n)_{t-1} + \dfrac{2}{n+1} \times (DIFF_t - MACD(n)_{t-1})$ |
| Larry William's R% | $\dfrac{H_n - C_t}{H_n - L_n} \times 100$ |
| A/D (Accumulation/Distribution) Oscillator | $\dfrac{H_t - C_{t-1}}{H_t - L_t}$ |
| CCI (Commodity Channel Index) | $\dfrac{M_t - SM_t}{0.015 D_t}$ |

$C_t$ is the closing price, $L_t$ is the low price and $H_t$ the high price at time t ; $DIFF_t = EMA(12)_t - EMA(26)_t$ ; EMA is exponential moving average, $EMA(k)_t = EMA(k)_{t-1} + \alpha \times (C_t - EMA(k)_{t-1})$, $\alpha$ is a smoothing factor which is equal to $\frac{2}{k+1}$ ; k is the time period of k-day exponential moving average, $LL_t$ and $HH_t$ implies lowest low and highest high in the last t days, respectively. $M_t = (H_t + L_t + C_t)/3$, $SM_t = (\sum_{i=1}^{n} M_{t-1+1})/n$, $D_t = (\sum_{i=1}^{n}|M_{t-i+1} - SM_t|)/n$, $UP_t$ means upward price change while $DW_t$ is the downward price change at time t.

The original data were scaled into the range of [-1,1]. The goal of linear scaling is to independently normalize each feature component to the specified range. It ensures that the larger value input attributes do not overwhelm smaller value inputs, and helps to reduce prediction errors (Kim, 2003) [3].

## VI.     EXPERIMENTAL DESIGN

In this paper, an EMD-KPCA-based SVM is proposed for stock price direction prediction. The framework of the model is shown in figure 3. and the building procedure is described as follows:

- Step 1: Ten technical indicators are calculated from the original stock prices series.
- Step 2: Each technical indicator is decomposed by EMD into a number of IMFs and a residual. An example of the decomposition for the simple 10-day moving average for the BMCI is given in figure 2.
- Step 3: The original features are constructed from all the IMFs and the residual, and they constitute the potential input variables for the predictive models.
- Step 4: The optimal feature subset is selected using KPCA.
- Step 5: The forecasting model (SVM) is built using the selected feature subset to perform the stock price direction prediction.
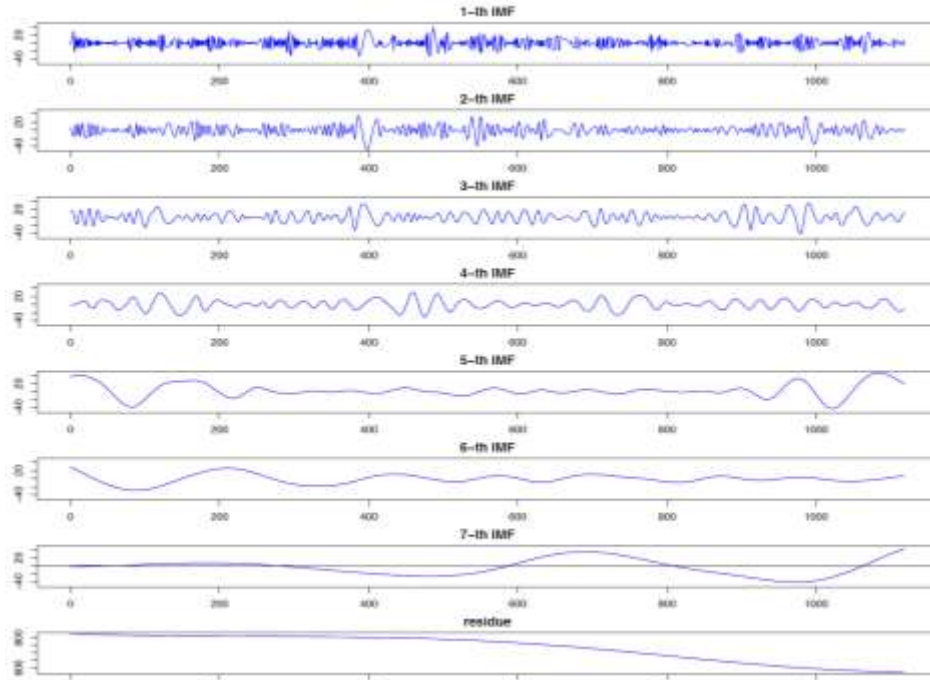


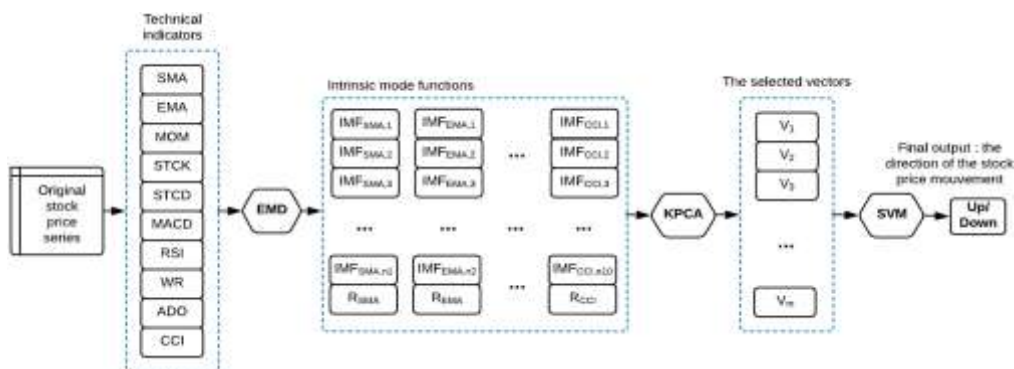**Figure 2 :** The decomposition for the simple 10-day moving average for the BMCI



**Figure 3 :** Framework of the proposed EMD-KPCA-based SVM

## VII.     RESULTS AND DISCUSSION

All the six data sets are partitioned into the training set, the validation set, and the testing set according to the time sequence. There is a total of 784 data patterns in the training set, and 167 data patterns in both the validation set and the test set in each of the data sets.

The Gaussian function $\exp\left(-(x_i - x_j)^2/\delta^2\right)$ is used as the kernel function of KPCA because the Gaussian kernel tends to give good performance under general smoothness assumptions (Smola, 1998) [32]. The best value of $\delta^2$ is chosen based on the cross-validation method.

In the implementation of KPCA, one major problem is to search for the optimal number of principal components n. The following procedure is used in the experiment. Firstly, the principal component calculated using the eigenvector corresponding to the largest eigenvalue is used as the inputs of SVM to perform the prediction. The principal components are then increased one by one (corresponding to the eigenvectors sorted in a descending order of eigenvalues) at each step. The

number of principal components from the minimum value 1 to the maximum value (305) are all investigated. The value which produces the best performance (the biggest accuracy) is used.

Accuracy and f-measure are used to evaluate the performance of all models. Computation of these evaluation measures requires estimating Precision and Recall which are evaluated from True Positive (TP), False Positive (FP), True Negative (TN) and False Negative (FN). These parameters are defined in Eqs (22):

$$\text{Precision}_{\text{positive}} = \frac{TP}{TP+FP}, \quad \text{Recall}_{\text{positive}} = \frac{TP}{TP+FN}$$
$$\text{Precision}_{\text{negative}} = \frac{TN}{TN+FN}, \quad \text{Recall}_{\text{negative}} = \frac{TN}{TN+FP} \tag{22}$$

Precision is the weighted average of precision positive and negative while Recall is the weighted average of recall positive and negative.

$$\text{Accuracy} = \frac{TP+TN}{TP+FP+TN+FN} \tag{23}$$

$$\text{F} - \text{measure} = \frac{2 \times \text{Precision} \times \text{Recall}}{\text{Precision} + \text{Recall}} \tag{24}$$

For training SVM, The Gaussian function is also used as the kernel function of SVM. The optimal values of the Gaussian parameter $\delta^2$, C and $\epsilon$ are chosen based on the cross-validation method.

The results obtained in all the data sets are given in Table 3. The experiments show that our result obtained in the EMD-KPCA-based SVM model can always converge to a better accuracy on the test set than single SVM, EMD-based SVM or KPCA-based SVM. This demonstrates the fact that decomposition coupled with dimension reduction can improve the generalization performance of SVM.

**Table 3** : The accuracy and F-measure of the different models

| Banks | SVM | | EMD-based SVM | | KPCA-based SVM | | EMD-KPCA-based SVM | |
|---|---|---|---|---|---|---|---|---|
| | Accuracy | F-measure | Accuracy | F-measure | Accuracy | F-measure | Accuracy | F-measure |
| BCP | 0.7189 | 0.7314 | 0.7460 | 0.7599 | 0.7673 | 0.7837 | 0.8271 | 0.8301 |
| BMCI | 0.7411 | 0.7499 | 0.7763 | 0.7904 | 0.7544 | 0.7723 | 0.8358 | 0.8401 |
| AWB | 0.7622 | 0.7767 | 0.7439 | 0.7682 | 0.8333 | 0.8469 | 0.8463 | 0.8529 |
| BMCE | 0.7822 | 0.8027 | 0.8338 | 0.8348 | 0.8486 | 0.8688 | 0.9333 | 0.9353 |
| CDM | 0.6566 | 0.6683 | 0.6712 | 0.6860 | 0.7050 | 0.7191 | 0.7487 | 0.7523 |
| CIH | 0.7367 | 0.7471 | 0.7324 | 0.7615 | 0.8164 | 0.8295 | 0.8328 | 0.8357 |
| BCP | 0.7189 | 0.7314 | 0.7460 | 0.7599 | 0.7673 | 0.7837 | 0.8271 | 0.8301 |

## VIII. CONCLUSION

The purpose of this paper is to predict direction of movement for stocks and stock price indices. Prediction performance of four models namely SVM, EMD-based SVM, KPCA-based SVM and EMD-KPCA-based SVM is compared based on 1118 trading days. Ten technical parameters reflecting the condition of stock and stock price index are used to learn each of these models. Experiments show that single SVM model exhibits least performance and EMD-KPCA-based SVM shows the highest performance. The improvement in the prediction accuracy of the proposed system can be deployed in real time for stocks' trend prediction, making investments more profitable and secure. Improvement of accuracy with the help of this approach that is based on common investor's methods for stock investing, also promotes the idea of pre-processing the data based on the domain in which machine learning algorithms are used. Ten technical indicators are used in this paper to construct the knowledge base, however, other macro-economic variables like currency exchange rates, inflation, government policies, interest rates etc. that affect stock market can also be used as inputs to the models or in construction of the knowledge base of an expert system.

## REFERENCES

[1]. A.Kazem, E.Sharifi, F.K.Hussain, M.Saberi and O.K.Hussain Support vector regression with chaos-based firefly algorithm for stock market price forecasting. Applied Soft Computing 13: 947-958, 2013.

[2]. Y. Zuo and E. Kita Stock price forecast using Bayesian network.Expert Systems with Applications39: 6729-6737, 2012.

[3]. K.-j. KimFinancial time series forecasting using support vector machines. Neurocomputing 55: 307-319, 2003.

[4]. K.-J.Kim and I.Han Genetic algorithms approach to feature discretization in artificial neural networks for the prediction of stock price index. Expert Systems with Application19: 125-132, 2000.

[5]. C.-F.Tsai, Y.-C.Lin, D. C.Yen and Y.-M.ChenPredicting stock returns by classifier ensembles. Applied Soft Computing 11: 2452-2459, 2011.

[6]. B. B.Nair, S. G.Sai, A.Naveen, A.Lakshmi, G.Venkateshand V.Mohandas A GA-artificial neural network hybrid system for financial time series forecasting. Information Technology and Mobile Communication147: 499-506, 2011.

[7]. M.-C.Wu, S.-Y.Lin and C.-H.Lin An effective application of decision tree to 903 stock trading. Expert Systems with Applications31(2): 270-274, 2006.

[8]. M.Qi and G.P.Zhang Trend time-series modeling and forecasting with neural networks. IEEE Transactions on Neural Networks 19: 808-816, 2008.

[9]. L.Yu, S.Wangand K.K.LaiA neural-network-based nonlinear metamodeling approach to financial time series forecasting. Applied Soft Computing 9: 563-574, 2009.

[10]. M.Pulido, P.Melin and O.CastilloParticle swarm optimization of ensemble neural networks with fuzzy aggregation for time series prediction of the Mexican Stock Exchange. Information Sciences 280: 188-204, 2014.

[11]. H.Yu, R. Chen and G. Zhang A SVM stock selection model within PCA. Procedia Computer Science 31: 406-412, 2014.

[12]. S. Choudhury, S. Ghosh, A. Bhattacharya, K.J. Fernandes and M.K. TiwariA real time clustering and SVM based price-volatility prediction for optimal trading strategy. Neuro-computing 131: 419-426, 2014.

[13]. E. Hajizadeh, A. Seifi, M.F. Zarandi and I.Turksen A hybrid modeling approach for forecasting the volatility of S&P 500 index return. Expert Systems with Applications 39: 431-436, 2012.

[14]. L. Cao and F.E. Tay Financial forecasting using support vector machines. Neural Computing & Applications 10: 184-192, 2001.

[15]. Q. Wen, Z. Yang, Y. Song and P. Jia Automatic stock decision support system based on box theory and SVM algorithm. Systems with Applications 37: 1015-1022, 2010.

[16]. A. Nahil and A. LyhyaouiShort-Term Stock Price Forecasting Using Kernel Principal Component Analysis and Support Vector Machines: The Case of Casablanca Stock Exchange. Procedia Computer Science127: 161-169, 2018.

[17]. W. Huang, Y. Nakamori and S.-Y.Wang Forecasting stock market movement direction with support vector machine. Computers & Operations Research 32: 2513-2522, 2005.

[18]. P. Ou and H. Wang Prediction of stock market index movement by ten data mining techniques. Modern Applied Science3: 28-42, 2009.

[19]. L. Luo and X. ChenIntegrating piecewise linear representation and weighted support vector machine for stock trading signal prediction. Applied Soft Computing 13: 806-816, 2013.

[20]. Y.S. Abu-Mostafa and A.F. Atiya Introduction to financial forecasting. Applied Intelligence 6: 205-213, 1996.

[21]. J. Patel, S. Shah, P. Thakkar and K. Kotecha Predicting stock and stock price index movement using Trend Deterministic Data Preparation and machine learning techniques. Expert Systems with Applications 42: 259-268, 2015.

[22]. M. Ballings, D. Van den Poel, N. Hespeels and R. Gryp Evaluating multiple classifiers for stock price direction prediction. Expert Systems with Applications 42(20): 7046-7056, 2015.

[23]. J.K. Mantri, P. Gahan and B. NayakArtificial neural networks-an application to stock market volatility. International Journal of Engineering Science and Technology2: 1451-1460, 2010.

[24]. R.KMishra, S. Sehgal and N. Bhanumurthy A search for long-range dependence and chaotic structure in Indian stock market. Review of Financial Economics 20: 96-104, 2011.

[25]. J. Sun and H. Li Financial distress prediction using support vector machines: Ensemble vs. individual. Applied Soft Computing 12: 2254-2265, 2012.

[26]. F. Liu and J. Wang Fluctuation prediction of stock market index by Legendre neural network with random time strength function. Neurocomputing 83: 12-21, 2012.

[27]. A. Garg, S. Sriram and K. Tai Empirical analysis of model selection criteria for genetic programming in modeling of time series system. Conference on Computational Intelligence for Financial Engineering & Economics: 90-94, 2013.

[28]. S.H.Kim and S. H.Chun Graded forecasting using an array of bipolar predictions: Application of probabilistic neural networks to a stock market index. International Journal of Forecasting 14(3) : 323-337, 1998.

[29]. N.E.Huang, Z.Shen, S.R.Long, M.C.Wu, H.H.Shih, Q.Zheng, N.-C.Yen, C.C.Tungand H. LiuThe empirical mode decomposition and the hilbert spectrum for nonlinear and non-stationary time series analysis. Proceedings of the Royal Society of London A: Mathematical, Physical and Engineering Sciences 454: 903-995, 1998.

[30]. V.Vapnik Statistical Learning Theory. Wiley, New York, NY, 1998.

[31]. N.I.Sapankevych, R.Sankar(2009Time series prediction using support vectormachines: a survey.Comput. Intell. Mag. IEEE 4 (2) : 24-38, 2009.

[32]. A.J.Smola Learning with Kernels. Ph.D. Thesis, GMD, Birlinghoven, Germany, 1998.