

# Mahalanobis Based k-Nearest Neighbor Forecasting versus Time Series Forecasting Methods

Vindya Kumari Pathirana<sup>1,\*</sup> and Kandethody M. Ramachandran<sup>2</sup>

<sup>1</sup>Department of Mathematics, University of Connecticut, Waterbury, 06702, Connecticut, United States

<sup>2</sup>Department of Mathematics and Statistics, University of South Florida, Tampa, 33620, Florida, United States

\*Corresponding Author: vindya.pathirana@uconn.edu

Copyright ©2016 by authors, all rights reserved. Authors agree that this article remains permanently open access under the terms of the Creative Commons Attribution License 4.0 International License

**Abstract** Foreign exchange (FX) rate forecasting is a challenging area of study. Various linear and nonlinear methods have been used to forecast FX rates. As the currency data are nonlinear and highly correlated, forecasting through nonlinear dynamical systems is becoming more relevant. The k-nearest neighbor (k-NN) algorithm is one of the most commonly used nonlinear pattern recognition and forecasting methods that outperforms the available linear forecasting methods for the high frequency foreign exchange data. As the k neighbors are selected according to a distance function, the choice of distance plays a key role in the k-NN procedure.

The time series forecasting method, Auto Regressive Integrated Moving Average process (ARIMA) is considered as one of the superior forecasting method for time series data. In this work, we compare the performances of Mahalanobis distance based k-NN forecasting procedure with the traditional ARIMA based forecasting method. In addition, the forecasts were transformed into a technical trading strategy to create buy and sell signals. The two methods were evaluated for their forecasting accuracy and trading performances.

**Keywords** Forex Trading, k-Nearest Neighbor, Mahalanobis Distance, ARIMA, Multi-step Ahead Time Series Forecasting

markets averaged 5.3 trillion dollars per day in April 2013. Thus, foreign exchange rates forecasting is one of the challenging and important applications of financial time series prediction.

Foreign exchange rate forecasting is a challenging task due to the non linearity and the highly correlated nature of the data [10, 25]. Nonlinear dynamical systems are becoming more popular and relevant forecasting techniques due to these data structure. Neighbor Algorithms is one of the most popular such non-linear pattern recognition algorithm, which dates back to an unpublished report by Fix and Hodges in 1951, [6]. The basic principle of *k*-nearest neighbor (*k*-NN) rule is to investigate the past behavior of the currency data so that it can fully capture the dependency of the future exchange rates and that of the past. As a pattern recognition algorithm, *k*-NN looks for the repetitions of specific price patterns such as major trends, critical or turning points.

*k*-nearest neighbor forecasting procedure is mainly based on the similarity structure of the past and the present. The recognized “nearest neighbors” are the only data values used in the forecasting algorithm. The term ‘nearest’ is determined by a distance metric. Thus, it is highly important to have a distance function which captures the true nature of the data. Among nearest neighbor methods, Mahalanobis distance proved to be more efficient [19, 20]. In this paper, we will compare Mahalanobis based nearest neighbor method of forecasting to some of the popular time series based methods.

## 1 Introduction

The foreign exchange (FX) market is a non-stop cash market where currencies of nations are traded. Foreign currencies are constantly and simultaneously bought and sold across local and global markets, and traders’ investments increase or decrease in value based upon currency movements. The investors goal in FX trading is to profit from foreign currency movements. A preliminary global study by the Bank for International Settlements from the 2013 Triennial Central Bank Survey of Foreign Exchange and OTC Derivatives Markets Activity show that trading in foreign exchange

In section 2, we will give some background material on *k*-NN & distance measures, times series forecasting methods, and multi-step ahead forecasting of time series.

In section 3 we will briefly discuss some of our previously obtained results of choosing embedding Dimension (*m*), number of nearest neighbors (*k*) and Mahalanobis distance as the distance choice and then present the comparison results of proposed Mahalanobis distance based *k*-NN and ARIMA forecasting models for single step ahead and multi-step ahead forecasting. The discussion and conclusions will be given in section 4 and 5.

## 2 Materials and Methods

### 2.1 $k$ -Nearest Neighbor Algorithm and the Choice of Distance

#### 2.1.1 Background

$k$ -nearest neighbor ( $k$ -NN) algorithm rank the data and chose the  $k$  closest of them based on the distance between the query vector and the historical values. First, we divide the time series data,  $\{x_t\}_{t=1}^n = \{x_1, x_2, \dots, x_n\}$  in to two separate parts; for  $T < n$ , a training (in-sample) set  $\{x_1, x_2, \dots, x_T\}$  and a testing (or out-of-sample) set  $\{x_{T+1}, x_{T+2}, \dots, x_n\}$ . In order to identify behavioral patterns in the data, we transform the scalar time series in to time series vectors. We need to choose an embedding dimension ( $m$ ) and delay time ( $\tau$ ) to create vectors out of the training set. After selecting  $m$  and  $\tau$ , a time series vector at time  $t$  can be written as;

$$x_t^{m,\tau} = (x_t, x_{t-\tau}, \dots, x_{t-(m-1)\tau}) \quad (1)$$

$$\text{where } 1 + (m-1)\tau \leq t \leq T$$

These  $m$ -dimensional vectors are often called as  $m$  - *histories* and the  $m$ -dimensional space  $\mathbb{R}^m$  is referred to be the phase space of the time series [10]. The primary goal of  $k$ -NN method is to use the most relevant vectors out of the training set in the forecasting. The most relevant vectors are the ones having similar dynamic behavior as the delay vector  $x_T^m$ . We compare the distance between the delay vector and all the other  $m$ -history vectors to choose the vectors with similar dynamic behavior [10]. Then we look for the closest  $k$  vectors in the phase space  $\mathbb{R}^m$  such that they minimize the distance function  $d(x_T^m, x_i)$ .

In  $k$ -NN algorithm,  $m$  and  $k$  are predetermined constants. In the literature, the optimal values of  $m$  and  $k$  are quite ambiguous. There have been quite a lot argument and discussions about the optimal choice of  $m$  and  $k$  since the NN rule was first officially introduced by Cover and Hart in 1967 [7, 23]. In section 3 We will discuss the choice of  $m$  and  $k$  for Mahalanobis distance along with other distance choices.

For the forecasting we can incorporate variety of Statistical and time series predicting methods with NN algorithm. In the literature of  $k$ -NN forecasting, the most commonly used forecasting method is locally weighted simple linear regression [3, 10]. Thus the future forecasts were obtained using the following locally adjusted linear regression model [7]:

$$\hat{x}_{T+1} = \sum_{n=0}^{m-1} \hat{a}_n x_{T-m\tau} + \hat{a}_m \quad (2)$$

The coefficients were fitted by the linear regression of  $x_{t_j+1}^m$  on  $x_{t_j}^m = (x_{t_j}, x_{t_j-\tau}, \dots, x_{t_j-(m-1)\tau})$  for  $j = 1, 2, \dots, k$ . Thus the estimated coefficients  $\hat{a}_i$  are the values of  $a_i$  that minimize

$$\sum_{j=1}^k (x_{t_j+1} - a_0 x_{t_j} - a_1 x_{t_j-\tau} - \dots - a_{m-1} x_{t_j-(m-1)\tau} - a_m)^2 \quad (3)$$

The data used in equation (3) are the only  $k(m+1)$  data values obtained from the  $k$ -neighbor vectors of

size  $m$  and the corresponding next values,  $x_{t_j+1}^m$  for  $j = 1, 2, \dots, k$  chosen neighboring vectors, not the entire data.

As the forecasting is completely based on the selected  $k$  nearest neighbors, it is highly important to use a distance function which captures the relevant behavior of the data accurately. Many researchers have pointed out the difficulty of choosing a distance measure for the NN algorithm that works well for different types of data. Over the past decades, the most common choice of distance was Euclidean distance [7, 10]. The way it is defined, the Euclidean distance is unable to capture the trend of the highly volatile (hence random) and highly correlated foreign exchange data when choosing the neighbors for the NN algorithm. Apart from Euclidean distance, several other distance measures such as Manhattan, Minkowski, and Hamming distances have been used in the algorithm for various types of classification problems [13, 23].

Even though the asymptotic probability of error of the NN is independent of the choice of metric, classification performance of finite sample nearest neighbor algorithm is not independent of the distance function [13, 17]. As Nearest neighbor rule is highly sensitive to outliers, selecting irrelevant neighbors can cause increase in forecasting error. In their work, Fukunaga & Hostetler showed that using a proper distance measure, the variance of the finite sample estimate can be minimized [13]. Short & Fukunaga investigate the relation between the distance function in  $k$ -NN and the error measure [13]. They concluded that the error can be minimized by using an appropriate distance metric without increasing the number of sample vectors.

In time series pattern recognition, an appropriate distance function can categorize data in to clusters by capturing the similarity or dissimilarity between the data. The *Euclidean* and *Manhattan (Absolute)* distances are the commonly used distances measures in nearest neighbor classification and forecasting algorithm.

*Euclidean distance* calculates the real straight line distance between two points and it's the most common distance of choice in NN algorithms. It works well for low dimensional data, it performs poorly when the data are high dimensional. Also, Euclidean is not the best distance choice when the data are highly correlated as it does not account the for correlation among the vectors.

*Manhattan distance* gets its name from the rectangular grid patterns of the streets in Manhattan [18]. It looks at the absolute difference between the coordinates. It is also recognized as a computationally simplified version of Euclidean distance. Manhattan distance is preferred to Euclidean distance in practice sometime, because the distance along each axis is not squared, a large difference in one of the dimensions will not affect the total outcome.

#### 2.1.2 Mahalanobis distance

*Mahalanobis distance* was introduced by P. C. Mahalanobis in 1936 by considering the possible correlation

among the data [9].

Consider  $n$ -dimensional vectors  $x = (x_1, x_2, \dots, x_n)$  and  $y = (y_1, y_2, \dots, y_n)$  in  $\mathbb{R}^m$ . The Mahalanobis distance between two vectors  $x$  and  $y$  is defined as:

$$d(x, y) = \sqrt{(x - y)' \Sigma^{-1} (x - y)} \quad (4)$$

Here,  $\Sigma^{-1}$  is the inverse of variance-covariance matrix  $\Sigma$  between  $x$  and  $y$  and  $'$  denotes the matrix transpose. The major difference in Mahalanobis to any other distance measure is that it takes the covariance into account. Due to this reason it is also called Statistical distance as well.

Mahalanobis distance belongs to the class of generalized ellipsoid distance defined by

$$d(x, y) = \sqrt{(x - y)' M (x - y)} \quad (5)$$

Here  $M$  is a positive definite, symmetric matrix. In the case the Mahalanobis distance, the matrix  $M$  becomes the inverse of variance-covariance matrix. Obviously, this includes Euclidean distances as a special case when  $M$  is the identity matrix.

When using Euclidean distance, the set of points equidistant from a given location is a sphere. The Mahalanobis distance stretches this sphere to correct for the respective scales of the different variables, and to account for correlation among variables [24]. As the axes of ellipsoidal sphere can assume any direction depending upon the data, this is more applicable in the area of time series pattern recognition. Thus, unlike dimensional Euclidean distance, it is possible to express the correlation and weight between dimensions using Mahalanobis distance. Due to these advantages, Mahalanobis distance captures the correlation and the trend of the time series, better compared to other distances [7, 17].

In our earlier work, we proposed to use Mahalanobis distance in  $k$ -NN algorithm for FX data. We compared the performance of the Mahalanobis distance based  $k$ -NN algorithm with popular Euclidean and Manhattan distance based algorithm. Some of the earlier results are given in the next section.

The performance of the Mahalanobis distance based  $k$ -nearest neighbor algorithm was compared with the time series forecasting technique, ARIMA in two ways:

- (i) Forecast accuracy
- (ii) Transforming their forecasts into a technical trading rule

In the former case, our goal is to capture the deviation of the fitted values against the actual observations. In the latter case, we are interested in looking at the forecasts in financial point of view.

### 2.1.3 Measures of Forecasting Accuracy

To capture the deviation of fit, we used commonly used accuracy measures, Mean square error (MSE),

Means absolute percentage error (MAPE), and Normalized Root Mean Square Error (NRMSE). Apart from these traditional accuracy measures, the following version of Theil's  $U$ -statistic ( $U$ ) to compare the forecasting accuracy of our model.

$$U = \frac{\sqrt{\sum_{t=1}^n (\hat{x}_t - x_t)^2}}{\sqrt{\sum_{t=1}^n (\hat{x}_t)^2} + \sqrt{\sum_{t=1}^n (x_t)^2}} \quad (6)$$

Here  $x_t$  is the actual value and  $\hat{x}_t$  is the fitted value.

$U$ -statistic is a measure of the degree to which the forecasted values differ from the actual values and is independent of the scale of the variable. The way it is constructed,  $U$ -statistic necessarily lies between zero and one, with zero indicating a perfect fit. However, it does not provide information on forecasting bias, which is better captured by mean square error.

### 2.1.4 Trading Decisions

As in any other financial market, in FX market also a trader's main goal is to make more money out of foreign currency fluctuations. The primary goal of foreign exchange rate forecasting has to be making proper trading signals: *buy* and *sell* at each time step so that the trader makes more money. To satisfy this main aspect, first we need to transform forecasts into trading signals.

The forecasts were transformed into a simple technical trading strategy using the trading rule used by Fernandez-Rodriguez, Sosvilla-Rivero, and Andrada-Felix in their work [6, 7]. Let  $\hat{r}_t$  given by

$$\hat{r}_t = \ln(\hat{x}_{t+1}) - \ln(1 + i'_t) - \ln(1 + i_t) \quad (7)$$

be the estimated return from a foreign currency position over the period  $(t, t + 1)$  based on the forecasted FX rate at time  $t$ . Here  $x_t$  represents the spot exchange rate at time  $t$ ,  $\hat{x}_{t+1}$ , is the forecasted value for  $x_{t+1}$  is the domestic (US) daily interest rate and  $i'_t$  is the foreign country daily interest rate. The trading signals at time  $t$  are made based on the estimated return  $\hat{r}_t$ . The positive returns are executed as long positions (buy) and the negative returns are executed as short position (sell) [6, 7]. So the trading decision can be given as

$$\hat{z}_t = \begin{cases} 1 & ; \text{if } \hat{r}_t > 0 \\ -1 & ; \text{if } \hat{r}_t < 0 \end{cases} \quad (8)$$

Based on estimated return, we calculate *estimated total (logaccess) return* of the trading strategy over the time period  $(1, n)$  as

$$\hat{R}_n = \sum_{t=1}^n \hat{z}_t r_t \quad (9)$$

Here  $r_t$  is the actual return at time given by

$$r_t = \ln(x_{t+1}) - \ln(x_t) - \ln(1 + i'_t) - \ln(1 + i_t)$$

We also consider the popular performance measure: *Sharpe ratio* to compare the results along with the estimated total return. The Sharpe ratio,  $S_R$  used here

is the mean daily total return of the trading strategy over its standard deviation,

$$S_R = \frac{\mu_{\hat{R}_n}}{\sigma_{\hat{R}_n}} \quad (10)$$

Higher values of Sharpe ratio indicates that the model is performing better.

## 2.2 $k$ -Nearest Neighbor Algorithm and Mahalanobis Distance

### 2.2.1 Data

The data used here are exchange rates of Euro (EUR), British pound sterling (GBP), Swiss franc (CHF), Japanese Yen (JPY), and Canadian dollar (CAD) vis-à-vis American dollar (USD) obtained from the ProQuest Statistical Datasets. These are the daily spot rates of the currencies from *January 2006 to December 2010*. In order to make the comparison more effective, we have considered 1250 data from each currency, and taken 1000 data values as our training sample. The remaining 250 values were considered as the test sample. The coefficients of the model were updated every time when new information arrived.

### 2.2.2 Embedding Dimension ( $m$ ) and Number of Nearest Neighbors ( $k$ ).

The choice of embedding dimension,  $m$ , and the number of nearest neighbors,  $k$  in the  $k$ -NN forecasting procedure is a key issue need to be addressed prior to making trading signals. Therefore, first we conducted an empirical investigation to find the optimal values of  $m$  and  $k$ . We wanted to figure out whether the choices for  $m$  and  $k$  are data dependent, and also distance dependent. The forecasting accuracy was compared using all the error measures mentioned in section 2.1.3, by varying the value of  $m$  and  $k$  along with different distance functions. 80% of the data was considered as the training set, and the remaining 20% was taken as the testing set. After analyzing the results, the key parameters  $m$  and  $k$  of the algorithm were chosen as 3 and 20, respectively. The complete results of choosing the embedding dimension ( $m$ ) and neighborhood size ( $k$ ) can be found in [20].

## 2.3 Time Series Forecasting Methods

The autoregressive process (AR) and the moving average process (MA) were very popular representations among the time series community over the past. Both of these models are only applicable to stationary time series data. Each method has its own pros and cons. The ARMA model combines the AR and MA processes to have a better forecasting in time series by taking advantages of both AR and MA methods.

### 2.3.1 The General mixed Autoregressive Moving Average (ARMA) Process

The General ARMA( $p,q$ ) process is a combination of an autoregressive process of order,  $p$ , and a moving average process of order,  $q$ . Herman Wold was the person who first put together AR and MA models to create ARMA process in 1938. Since then, this method has been used in many areas of time series. ARMA( $p,q$ ) process is defined as;

$$x_t = \alpha + \phi_1 x_{t-1} + \dots + \phi_p x_{t-p} + \epsilon_t + \theta_1 \epsilon_{t-1} + \dots + \theta_q \epsilon_{t-q} \quad (11)$$

or

$$\Phi_p(L)x_t = \alpha + \Theta_q(L)\epsilon_t$$

where

$$\Phi_p(L) = 1 - \phi_1 L - \phi_2 L^2 - \phi_3 L^3 - \dots - \phi_p L^p$$

and

$$\Theta_q(L) = 1 + \theta_1 L + \theta_2 L^2 + \theta_3 L^3 + \dots + \theta_q L^q$$

The ARMA process is invertible if the roots of  $\Theta_q(L) = 0$  lie outside the unit circle and stationary if the roots of  $\Phi_p(L) = 0$  lie outside the unit circle [?, 16]. Note that we need to make the assumption of  $\Theta_q(L) = 0$  and  $\Phi_p(L) = 0$  sharing no common roots [?, 16].

### 2.3.2 The General mixed Autoregressive Integrated Moving Average (ARIMA) Process

In reality, most of the time series are non-stationary. For non-stationary time series, roots of the AR polynomial do not lie outside the unit circle. Therefore, we are not able to use the general mixed ARMA( $p,q$ ) model for forecasting. In such cases, the time series can be converted to a stationary process by differencing. This is also known as the *integrated* part of the algorithm., which transforms the general stationary ARMA process in to non stationary ARIMA( $p,d,q$ ) process. Here  $d$  is the degree of differencing. The difference filter is normally given by

$$(1 - L)^d \quad \text{where} \quad L^j x_t = x_{t-j} \quad (12)$$

Generally,  $d$  will be a positive integer and represents the number of times  $x_t$  must be differenced to achieve a stationary transformation. Typically,  $d \in \{0, 1, 2, \dots, d\}$ . When  $d = 0$ , the ARIMA process becomes stationary ARMA process. Thus the autoregressive integrated moving average, ARIMA( $p,d,q$ ) can be written as

$$\Phi_p(L)(1 - L)^d x_t = \alpha + \Theta_q(L)\epsilon_t \quad (13)$$

where

$$\Phi_p(L) = 1 - \phi_1 L - \phi_2 L^2 - \phi_3 L^3 - \dots - \phi_p L^p$$

and

$$\Theta_q(L) = 1 + \theta_1 L + \theta_2 L^2 + \theta_3 L^3 + \dots + \theta_q L^q$$

Sometimes, selecting the best order of the ARIMA( $p, d, q$ ) is a challenging task. As the forecasting is strongly depending on the order of the model, it is highly important to pick the correct order. The procedure needs to be completed in two steps. First we need to figure

out the differencing order, of the process. To determine the correct order of differencing, we continue the differencing procedure until the time series becomes stationary. The Kwiatkowski-Phillips-Schmidt-Shin (KPSS) test and Augmented Dickey-Fuller unit root test are normally used to determine the stationarity of a time series [21].

Once the correct differencing order,  $d$  is determined, the order of AR polynomial,  $p$ , and MA polynomial,  $q$  are determined using either Akaike information criterion (AIC). The AIC normally measures the quality of each model, relative to each of the other models. It is defined as

$$\ln(L) = 2M - 2\ln(L) \quad (14)$$

where,  $M$  is the number of parameters in the model, and  $\ln(L)$  is the unconditional log-likelihood function given by

$$\ln(L) = -\frac{n}{2}\ln(2\pi\sigma^2) - \frac{1}{2}\sigma^2 \sum_{i=1}^n (x_i - \mu)^2 \quad (15)$$

Here,  $\mu$  and  $\sigma$  are the mean and the standard deviation of time series respectively. The AIC is calculated by changing the values of  $p$  and  $q$  in the ARIMA model, and the model with the smallest AIC is usually selected for forecasting.

## 3 Results

### 3.1 $k$ -Nearest Neighbor Forecasting with Mahalanobis Distance

#### 3.1.1 Embedding Dimension ( $m$ ) and Number of Nearest Neighbors ( $k$ )

The Table 1 gives the  $U$ -statistic values for different choices of  $m$  with Mahalanobis distance. We were able to obtain similar results with other error measures as well. We also investigate the choice of neighborhood size,  $k$  with difference distance measures. The Figure 1 gives the  $U$ -statistic values for different numbers of nearest neighbors, ( $k$ ) with Mahalanobis inductance.

The complete results of choosing the embedding dimension ( $m$ ) and neighborhood size ( $k$ ) can be found in [20]. After analyzing the results, the key parameters  $m$  and  $k$  of the algorithm were chosen as 3 and 20, respectively.

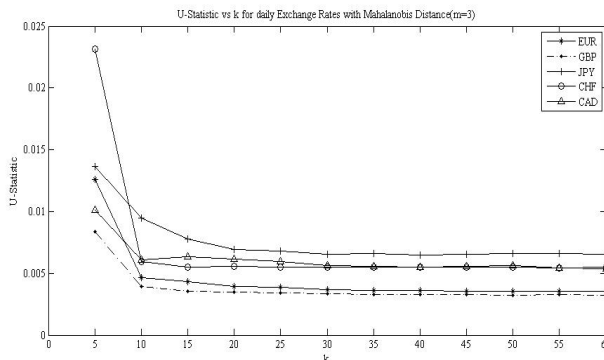


Figure 1. U-Statistic vs.  $k$  with Mahalanobis Distance.

### 3.1.2 Forecasting Accuracy and Trading Performances of Mahalanobis Distance vs. Other Distance Measures.

As discussed in section 2.1, the choice of distance plays a key role in the nearest neighbor algorithm. In our earlier published work, we have provided sufficient evidence supporting Mahalanobis distance as the choice of distance in  $k$ -NN procedure [19, 20]. We came to this conclusion by comparing the forecasting accuracy as well as the trading performances of the 5 currency data sets with Mahalanobis distance and other traditional distances such as Euclidean and absolute distances. The Mahalanobis distance outperforms the traditional distance functions for all the data sets with respect to the forecasting accuracy and trading performances. The details results can be found in [19] and [20].

## 3.2 Foreign Exchange Rates Forecasting with general ARIMA Process

In this section, we will first introduce the data preparation procedure for the same five currency data sets we used in section 2.2. Then, we will determine the appropriate ARIMA model for each data set, and finally compare the ARIMA approach with the Mahalanobis distance based  $k$ -NN forecasting method. In this paper also, the comparison will be performed according to two main aspects of forecasting. As the primary step, we will consider the different error measures discussed in section 2.1.3 to compare the forecasting accuracy. As the secondary step, the ARIMA forecasts will be transformed into trading signals using the same technical trading strategy discussed in section 2.1.4 and compare with the trading performances of  $k$ -NN procedure.

As discussed in Section 2.3.2, the order of differencing will be determined using the Kwiatkowski-Phillips-Schmidt-Shin (KPSS) test. For that, we will keep on differencing the series until the data becomes stationary. To figure out the order  $p$  of AR polynomial and  $q$  of MA polynomial, we are considering a positive constant  $m = 5$  with  $p+q = m$ . Then, we will vary the values of  $p$  and  $q$  such that  $p+q \leq m$  and estimate the parameters;  $\phi_1, \phi_2, \dots, \phi_p, \theta_1, \theta_2, \dots, \theta_q$  of each ARIMA( $p,d,q$ ) model. The Akaike information criterion (AIC) was computed for each model to choose the model with the minimum AIC.

#### 3.2.1 ARIMA Forecasting Model for EUR/USD Daily Rates

Following the step-by-step procedure we introduced above, the forecasting model with minimum AIC for the EUR/USD exchange rates data set was ARIMA(1,1,1), that is a combination of first order autoregressive (AR), and a first order moving average (MA) with the first difference filter ( $d = 1$ ). The model can be explicitly written with the estimated parameters as below:

$$(1 - 0.1329L)(1 - L)x_t = 0.00019 + (1 - 0.1323L)\epsilon_t \quad (16)$$

**Table 1.**  $U$ -statistic values for different  $m$ 

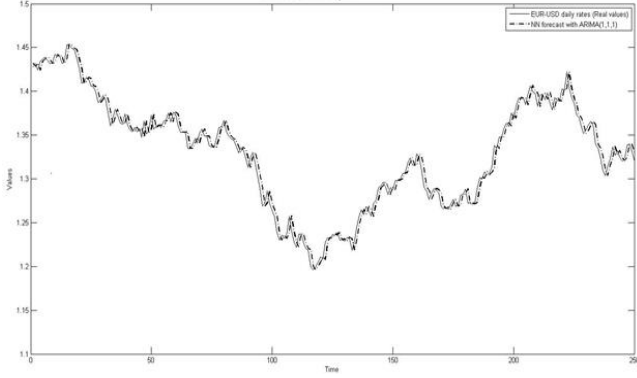
Currency	$m = 3$	$m = 4$	$m = 5$	$m = 6$
EUR	0.00389801	0.00392875	0.00430221	0.00456411
GBP	0.00345430	0.00348754	0.00359710	0.00362095
JPY	0.00690517	0.00717526	0.00810838	0.00771370
CHF	0.00555915	0.00572837	0.00586630	0.00590578
CAD	0.00597986	0.00642774	0.00612342	0.00643873

**Table 2.**  $U$ -statistic values for different distance measures ( $m = 3$  &  $k = 20$ )

Currency	Mahalanobis distance	Euclidean distance	Absolute distance
EUR	0.00389801	0.00405640	0.00401522
GBP	0.00345430	0.00377776	0.00374190
JPY	0.00690517	0.00971223	0.01331857
CHF	0.00555915	0.00755257	0.00771336
CAD	0.00597986	0.00648498	0.00619930

By letting  $\epsilon_t = 0$ , we have the one day ahead forecasting time series for EUR/USD currency data as

$$x_t = 0.000192 + 1.132915x_{t-1} - 0.132915x_{t-2} - 0.132342\epsilon_{t-1} \quad (17)$$

**Figure 2.** ARIIMA(1,1,1) Forecasts and real values for EUR/USD daily exchange rates.

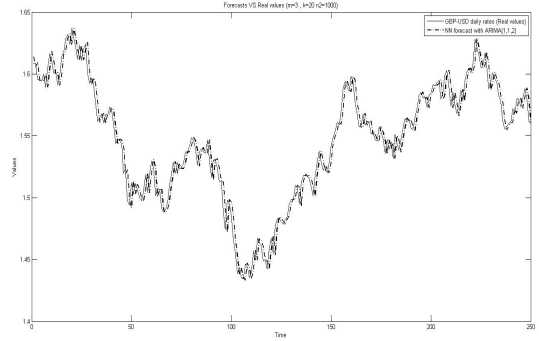
### 3.2.2 ARIMA Forecasting Model for GBP/USD Daily Rates.

The forecasting model with minimum AIC for the GBP/USD exchange rates data set was ARIIMA(1,1,2) model. This process is a combination of first order autoregressive (AR), and a second order moving average (MA), with the first difference filter. The model can be explicitly written with the estimated parameters as below:

$$x_t = -0.000218556 + 0.363691x_{t-1} + 0.636309x_{t-2} + 0.652109\epsilon_{t-1} + 0.063216\epsilon_{t-2} \quad (18)$$

### 3.2.3 ARIMA Forecasting Model for JPY/USD Daily Rates.

After comparing AIC for JPY/USD rates data set, we came up with the following ARIMA(1,1,2) model, that is a combination of second order autoregressive (AR),

**Figure 3.** ARIMA(1,1,2) Forecasts and real values for GBP/USD daily exchange rates.

and a first order moving average (MA), with the first difference filter.

$$(1 - 0.664950L)(1 - L)x_t = 0.00000098 + (1 - 0.400612L - 0.234753L^2)\epsilon_t \quad (19)$$

Expanding the autoregressive operator and the difference filter and then letting  $\epsilon_t = 0$ , we obtained one day ahead forecasting time series for JPY/USD currency data as

$$x_t = 0.00000098 + 1.0.664950x_{t-1} - 0.664950x_{t-2} - 0.400612\epsilon_{t-1} - 0.234753\epsilon_{t-2} \quad (20)$$

### 3.2.4 ARIMA Forecasting Model for CHF/USD Daily Rates.

For the CHF/USD daily rates, we found that ARIIMA(1,1,1) as the model with smallest AIC by varying the values  $p$  and  $q$ , after deciding the deference degree as one. The ARIMA process can be explicitly written with the estimated parameters as:

$$x_t = 0.000231 + 0.918785x_{t-1} + 0.0881215x_{t-2} - 0.0881215\epsilon_{t-1} \quad (21)$$

### 3.2.5 ARIMA Forecasting Model for CAD/USD Daily Rates.

The forecasting model with the minimum AIC value for the CAD/USD daily rates was a combination of first order autoregressive (AR), and a third order moving average (MA) with the first difference filter, namely, ARIMA(1,1,3).

$$x_t = 0.000055 + 1.481798x_{t-1} - 0.481798x_{t-2} - 0.182836\epsilon_{t-1} - 0.130882\epsilon_{t-2} + 0.067018\epsilon_{t-2} \quad (22)$$

### 3.3 $k$ -Nearest Neighbor Forecasting vs. ARIMA Forecasting - Forecasting Accuracy

In section 3.2, we discussed in details the general autoregressive integrated moving average forecasting models for the five daily exchange rates data sets. The given figures (Figure 2 & 3 ) indicate that the ARIMA forecasts follow the actual values pretty well as similar to the case of  $k$ -NN forecasting. In this section, our goal is to compare the forecasting accuracy of ARIMA models with our proposed Mahalanobis distance based  $k$ -nearest neighbor forecasting procedure.

We considered all the accuracy measures mentioned in section 2.1.3 and compared the performances of each procedure based on how accurate their forecasts were. The following tables give the  $U$ -statistic, mean square error, and normalized root mean square error for the currencies EUR, GBP, JPY, CHF, and CAD with Mahalanobis distance based  $k$ -NN algorithm and ARIMA forecasting models.

According to the obtained results given by tables 3, 4, & 5, we can see that the majority of time (3 out of 5) Mahalanobis distance based  $k$ -NN forecasting model out performs the ARIMA method. In the cases of EUR and GBP, the general ARIMA process seems to forecast relatively better compare to the nearest neighbor forecasting. Even though our primary goal in this paper is to compare the trading performances of both methods, it is necessary to further analyze the data, and come up with an explanation behind this situation. For this purpose, we calculated the following statistical measures for all the data sets:

- Total Variation -  
The total variation or the total sum of squares (SST) is a measure of the observed values around the mean. It is comprised the sum of the squares of the differences of each data value with the mean.

$$Total\ variation = \sum_{t=1}^n (x_t - \bar{x})^2 \quad (23)$$

- Standard Deviation -  
In statistics, the standard deviation is a measure of the spread of scores within a set of data. It is a measure that is used to quantify the amount of variation or dispersion of a set of data values. Smaller the standard deviation, closer the data points to its

mean.

$$Standard\ deviation, \sigma = \sqrt{\frac{\sum_{t=1}^n (x_t - \bar{x})^2}{n}} \quad (24)$$

Considering the calculated values for total variation and standard deviation (given in table 6), we observed that the EUR and GBP daily rates have relatively higher total variation and standard deviation compare to the remaining data sets.

### 3.4 $k$ -Nearest Neighbor Forecasting vs. ARIMA Forecasting - Comparing Trading Decisions

As it is obvious that currency trader's main goal is to make more money, in this section we evaluated these two prediction models ( $k$ -NN and ARIMA) considering their trading performances. We transformed the ARIMA forecasts in to trading signals, *buy* and *sell* using technical trading strategy discussed in section 2.1.4. Then, the performance measures, *total (log access) return* and *Sharpe ratio* were calculated and compared with those of  $k$ -nearest neighbor forecasting technique. Higher values of these measures indicate that the model is performing better.

The estimated total return and Sharpe ratio for the technical trading strategy under  $k$ -NN algorithm and ARIMA process are given in tables 7 and 8. The final conclusion of forecasting model is pretty much same as that of error measures. Proposed Mahalanobis distance based  $k$ -NN method outperforms the ARIMA process majority of the time. According to the forecasting accuracy, both EUR and GBP daily exchange rates data sets support ARIMA model. However, when comparing total return and Sharpe ratio, GBP/USD daily rates pretty much gave the same numerical values for both the models. Therefore, the results for trading decisions also indicated that the  $k$ -nearest neighbor forecasting model producing more accurate and profitable trading signals compared to the general ARIMA process.

The results from section 3.3 and section 3.4 motivates to investigate more on the behavior of time series data and the most appropriate forecasting technique. The primary goal of the next section is to study the forecasting accuracy of simulated time series data with both Mahalanobis distance based  $k$ -NN method and the general ARIMA forecasting models.

### 3.5 Simulation Data Analysis

Time series data simulation plays an important role in many areas of time series data analysis such as economics & finance, environmental studies , and engineering. It is a whole different area of research, where the researchers have paid much more attention in the recent history. Generating financial time series such as exchange rates data is a challenging task compared to most of the other time series data simulation. A huge amount of empirical contributions been made towards

**Table 3.**  $U$ -statistics with  $k$ -NN and ARIMA models

Currency	$k$ -NN forecasting	ARIMA forecasting
EUR	0.003898012	0.003456137
GBP	0.003454303	0.00318494
JPY	0.00690517	0.008302838
CHF	0.005559151	0.007286362
CAD	0.005979865	0.00731294

**Table 4.** Mean square error with  $k$ -NN and ARIMA models

Currency	$k$ -NN forecasting	ARIMA forecasting
EUR	0.00010905	0.00008479
GBP	0.00011439	0.00009725
JPY	0.00024810	0.00035878
CHF	0.00011441	0.000196566
CAD	0.00013689	0.000199935

this topic, and variety of economical, financial and time series models been proposed and experimented by many academic and industrial researchers during the last two decades. As most of the traditional financial theory based methods failed to match the features displayed by the actual data, many alternative models were proposed to overcome the issues of these traditional theory based models [2].

The purpose of any foreign currency generating algorithm is to replicate a certain exchange rate by considering all the financial and economical factors related to those two countries, which is a complicated task. In their work Bianchi, Pantanella, and Pianese claimed that using their proposed multifractional process with random exponent, they were successfully able to replicate EUR/JPY and EUR/USD [2] exchange rates. Also, Oyediran & Afieroho have worked on developing an algorithm to simulate many different FX rates such as European euro, British pound sterling and the US dollar against the Nigeria naira [18]. Their simulation models were also developed after analyzing the historical data of the corresponding currency rates.

All these simulation algorithms have one main goal in common. Their goal was to develop a procedure well capture the behavior of a given currency rate, which was not our intention of simulation study in this work. The goal here is to capture the behavior(s) of a time series to decide which forecasting algorithm ( $k$ -NN or ARIMA) would be more beneficial. Even though our primary interest is forecasting and decision making in foreign exchange market, for the simulation study we considered time series data in general.

Auto regressive (AR), moving average (MA), and general and mix ARIMA models are the most popular time series data simulation techniques among the time series research community. These time series processes have been used by many researchers over the recent history to replicate time series data using different computer software such as MATLAB and R [15]. For the simulation data analysis, several time series data sets were simulated in MATLAB environment with the use of the built-in MATLAB functions “*arima*” and “*simulate*”. Since the data were simulated using ARIMA process, there is always a possibility of having an advantage of using an ARIMA forecasting model.

The observations from section 3.3 and section 3.4

lead to the conclusion that for a time series data with a higher volatility, ARIMA forecasting procedure works better compared to the  $k$ -nearest neighbor method. As can be seen from the table 6, both EUR and GBP data sets have higher volatility measures compared to the rest. Due to this reason, the time series were simulated by varying the standard deviation. We have chosen a range from 0.00126 to 0.896 to capture the range of our data sets’ standard deviations. The simulated 9 data sets and their standard deviations listed in Table 9.

The model comparison was performed using the accuracy measures discussed in section 2.1.3. We only focused on deviation in fit for this comparison. To compare the trading decisions, it is necessary to simulate the interest rates, and also the time series data replicating real FX data of a certain country, which is not our interest here. Also the obtained results in section 2.2 and 3.4 suggest that having more accurate forecasts always lead to a higher trading performances.

We followed the same data preparation procedure discussed in section 2.3.2 to build the best model for each data set when using ARIMA process for forecasting. Even though the data was simulated with the specified orders and parameters, we again tested them for the appropriate differencing order and AR order,  $p$ , and MA order,  $q$ . For Mahalanobis distance based  $k$  nearest neighbor algorithm the parameter  $m$  was set to be 3 and  $k$  was set to be 20 as in section 3.3. One step ahead out of sample forecasts were created for 250 test set and the size of the training window was 1000.

The comparison results of  $U$ -statistic for the simulated data are presented in Table 10. It can be clearly seen that for the data sets 1, 2, and 3, ARIMA based forecasting models had lower  $U$ -statistic values compared to those of Mahalanobis distance based  $k$ -NN forecasting. Those are the data sets with higher standard deviations. When the standard deviation is getting smaller and smaller,  $k$ -NN forecasting algorithm started to perform comparatively better than general ARIMA process. For the data sets 6 through 9, the difference between the  $U$ -statistic values are significant. This supports the claim that for a time series data with a lower standard deviation, the  $k$ -NN method has a higher forecasting accuracy compared to ARIMA.



**Table 5.** Normalized root mean square error with  $k$ -NN and ARIMA models

Currency	$k$ -NN forecasting	ARIMA forecasting
EUR	0.16748184	0.14691536
GBP	0.22251281	0.20423380
JPY	0.30304086	0.35985560
CHF	0.21693213	0.28301529
CAD	0.58311581	0.65600438

**Table 6.** Total Variation and Standard Deviation

Currency	Total variation	Standard deviation
EUR	10.99275618	0.104846346
GBP	36.55015706	0.191180954
JPY	0.000877177	0.000936577
CHF	4.313622939	0.065678177
CAD	4.193662329	0.064758492

The other error measures also support this argument. Even these data were simulated using general ARIMA process, for the time series data with a lower volatility,  $k$ -nearest neighbor forecasting method (with Mahalanobis distance) outperforms the ARIMA forecasting procedure.

We went further and tried to figure out exactly around what value of standard deviation  $k$ -NN procedure starting to work better. Table 10 clearly indicate that somewhere between the values of 0.127 & 0.283,  $k$ -NN forecasting procedure has started performing better. To investigate this further, couple of more data sets were simulated with standard deviation between 0.009 and 0.15. Then, we followed the exact same procedure and predicted 250 future values. From the given results of  $U$ -statistic values in table 11, we can observe that for the standard deviation values below 0.13, the  $k$ -NN has a better forecasting accuracy.

## 4 Discussion

*Still working on it ...*

## 5 Conclusions

In this paper, our main goal was to compare the proposed Mahalanobis distance based  $k$ -NN forecasting with general autoregressive integrated moving average (ARIMA) process, which is assumed to be one of the best time series forecasting technique. As all these forecasting methods are data driven models, giving an optimal forecasting model works with all types of data is practically a difficult task.

From our results, we can conclude that  $k$ -nearest neighbor forecasting algorithm with Mahalanobis distance function outperforms the popular time series forecasting technique, general ARIMA process, majority of the time.

For the data sets with a relatively higher total variation (or highly volatile), ARIMA methods seems to work better compared to the  $k$ -NN forecasting. Our simulation data study supported this claim as well. Considering the accuracy measures ( $U$ -statistic and MSE), we can conclude that for time series data with a smaller standard deviation,  $k$ -NN forecasting procedure more appropriate than the ARIMA process.

The nearest neighbor algorithm is a nonparametric, on-line learning algorithm. Thus, it does not require any distributional assumptions, and data preparation ahead of time. Unlike nearest neighbor, ARIMA process requires model building procedure to select proper differencing order ( $d$ ), autoregressive order ( $p$ ), and moving average order ( $q$ ). The obtained results proved that even with all these model building procedure, still the ARIMA process worked better only for one currency data set according to the trading decisions. We discussed in the previous section (section 2.2) that choosing an appropriate distance in NN algorithm can improve the forecasting significantly. The results obtained in this paper further support our earlier conclusion. Also, we noticed that the  $k$ -NN forecasting method can be further improve by adjusting the algorithm according to the previous forecasting errors, which will be part of our future work.

**Table 7.** Total Return  $k$ -NN and ARIMA models

Currency	$k$ -NN forecasting	ARIMA forecasting
EUR	0.52991777	0.89316418
GBP	4.16807227	4.16807227
JPY	0.67755404	0.47532224
CHF	5.42108879	5.16742874
CAD	4.38589604	4.03711714

**Table 8.** Sharpe Ratio with  $k$ -NN and ARIMA models

Currency	$k$ -NN forecasting	ARIMA forecasting
EUR	0.27890809	0.50803011
GBP	2.41593434	2.41593434
JPY	0.18429771	0.12818451
CHF	1.67419376	1.42328537
CAD	1.26400087	1.04713284

## REFERENCES

- [1] Tao Ban, Ruibin Zhang, Shaoning Pang, Abdolhossein Sarrafzadeh, Daisuke Inoue, Referential  $k$ -NN Regression for Financial Time Series Forecasting in *Neural Information Processing: Lecture Notes in Computer Science* (2013), 8226, 601 - 608.
- [2] Sergio Bianchi, Alexandre Pantanella, and Augusto Pianeese, Modeling and Simulation of Currency Exchange Rates Using Multifractional Process with Random Exponent., *International Journal of Modeling and Optimization*, *ESTSP08* (June 2012) 2, No 3
- [3] M. Casdagli, Nonlinear forecasting, Chaos and statistics, Modeling Complex Phenomena, Part of the series Woodward Conference pp 131-152, Springer, 1992.
- [4] M. Casdagli, Nonlinear prediction of chaotic time series. *Physica D* (1989), 35: 335-356
- [5] M. Casdagli, Chaos and Deterministic versus Stochastic Nonlinear Modeling *Journal of Royal Statistical Society, Series B Vol. 54, No 2* (1992), 303-328 .
- [6] T. M. Cover, P. E. Hart, Nearest Neighbor Pattern Classification, *IEEE Trans. on Information Theory*, *IT-13 (1)* (1967), 21-27.
- [7] L. Devroye, Györfi and G. Lugosi, *A Probabilistic Theory of Pattern Recognition*, Springer-Verlag, New York, 1996.
- [8] D. Farmer and J. Sidorowich, Predicting chaotic time series., *Physical review Letters* (1987), 59, 845-848.
- [9] Fernando Fernández-Rodríguez, S. Sosvilla-Rivero, J. Andrada-Félix, Nearest-Neighbour Predictions in Foreign Exchange Markets, *FEDEA Int. Economics and Finance DEFI* (2002), Working Paper
- [10] Fernando Fernández-Rodríguez, S. Sosvilla-Rivero, J. Andrada-Félix, Technical Analysis in Foreign Exchange Markets: Linear versus Nonlinear Trading Rules, *FEDEA Int. Economics and Finance DEFI* (2000), Working Paper No. 00-02.
- [11] Fernando Fernández-Rodríguez, S. Sosvilla-Rivero, Testing nonlinear forecastability in time series: Theory and evidence from the EMS, *Economics Letters* (1998), 59, 49 - 63.
- [12] Fernando Fernández-Rodríguez, J. Andrada-Félix, S. Sosvilla-Rivero, Combining information in exchange rate forecasting: evidence from the EMS, *Applied Economics Letters* (1997), 4: 7, 441 - 444
- [13] K. Fukunaga and L. D. Hostetler, Optimization of  $k$ -Nearest-Neighbor Density Estimates , *IEEE Transactions on information theory*, Vol. *IT-19* (May 1973), 320-326.
- [14] N. A. Gershenfel and A. S. Weigend, The Future of Time Series: Learning and Understanding , in: *Time Series Predictions: Forecasting Future and Understanding the Past* (1993), 1-70.
- [15] J. D. Hamilton, *Time Series Analysis*. Princeton, NJ: Princeton University Press, 1994
- [16] Gebhard Kirchgassner and Jurgen Wolters, *Introduction to Modern Time Series Analysis* Springer, 2007.
- [17] P. C. Mahalanobis, On the Generalized Distance in Statistics, in: *Proceedings of the national Institute of Science of India 12* (1936), 49-55.
- [18] Oyelami Benjamin Oyediran and Edooghogho Afieroho, Simulation of Exchange Rates of Nigerian Naira Against US Dollar, British Pound and the Euro Currency, in: *Studies in Mathematical Sciences* 6, No. 2 (2013), 58-70.
- [19] Vindya I. K. Pathirana and K. M. Ramachandran, The Distance Choice and Optimal parameter selection in  $k$ -NN algorithm for FX data , in: *The special issue in Communications in Applied Analysis*, 18 (2014), 591-612.
- [20] Vindya I. K. Pathirana and K. M. Ramachandran,  $k$ -Nearest Neighbor Algorithm with Mahalanobis Distance for FX Trading, in: *proceedings of dynamic Systems and Applications*, 6 (2012), 324-328.
- [21] Mohsen Pourahmadi, *Foundations of Time Series Analysis and Prediction Theory*, Wiley Series in Probability and Statistics , New York, 2001.
- [22] Tim Sauer, James A. Yorke, and Martin Casdagli, Embedology. *Journal of Statistical Physics* (1991), 65
- [23] D. Robert Short and K. Fukunaga, The optimal Distance Measure for Nearest neighbor Classification, *IEEE transactions on information theory* (1981), Vol. *IT-27*, No. 5, 622-627.

**Table 9.** Standard deviations of the simulated data

Data Set	Standard deviation
Simulated data set 1	0.896010158
Simulated data set 2	0.400925413
Simulated data set 3	0.283497079
Simulated data set 4	0.126783748
Simulated data set 5	0.040092541
Simulated data set 6	0.012678375
Simulated data set 7	0.004009254
Simulated data set 8	0.001267838
Simulated data set 9	0.001267838

**Table 10.**  $U$ -Statistics for  $k$ -NN forecasts and ARIMA forecasts for Simulated data

Data Set $D_i$	Standard deviation	$k$ -NN $U$ -statistic	ARIMA $U$ -statistic
$D_1$	0.896010158	0.18665837	0.14858764
$D_2$	0.400925413	0.08280900	0.07039754
$D_3$	0.283497079	0.05881467	0.05116939
$D_4$	0.126783748	0.02643441	0.02684713
$D_5$	0.040092541	0.00837700	0.01733388
$D_6$	0.012678375	0.00265053	0.01621586
$D_7$	0.004009254	0.00083831	0.01614590
$D_8$	0.001267838	0.00026511	0.01615349
$D_9$	0.001267838	0.00185742	0.01614677

[24] T. Uemiya, Y. Matsumoto, D. Koizumi, M. Shishibori, K. Kita, Fast Multidimensional Nearest Neighbor Search Algorithm Based on Ellipsoid Distanc, *International Journal of Advanced Intelligence* (2009).

[25] Walters-Williams, and Y. Li, Comparative Study of Distance Functions for Nearest Neighbor, in *Advanced Techniques in Computing and Software Engineering*, (2010), (Ed: Elleithy) 123–456.

**Table 11.** *U*-statistic for *k*-NN forecasts and ARIMA forecasts: Simulated data (standard deviation 0.009 - 0.15)

Simulated data <i>sd</i>	<i>k</i> -NN <i>U</i> -statistic	ARIMA <i>U</i> -statistic
0.15010000	0.031257268	0.030175564
0.13000000	0.02708489	0.02728483
0.12750000	0.026565797	0.026934802
0.12030000	0.0250823020	0.025948689
0.10540000	0.021970607	0.023958385
0.10030000	0.020912901	0.02330971
0.09100000	0.018988001	0.022171638