

BIRKBECK COLLEGE, UNIVERSITY OF LONDON

MSc. DATA SCIENCE MASTER THESIS

---

# The Applicability of Machine Learning Methods For Currency Trading

---

*Author*

Edward GILL

*Supervisor*

Dr. George MAGOULAS

September 12, 2019

This proposal is substantially the result of my own work, expressed in my own words, except where explicitly indicated in the text. I give my permission for it to be submitted to the JISC Plagiarism Detection Service. The proposal may be freely copied and distributed provided the source is explicitly acknowledged.

## **Abstract**

The availability of immense computing power has opened a veritable Pandora's box of highly complex and computationally intensive models which can now be applied to any dataset, whether appropriate or not. While machine learning models have led to rapid and impressive use cases across a host of different fields of research, this project will focus on the efficacy of machine learning models in global financial markets and specifically currency markets.

Applying complex methods within the financial domain requires a different approach from other areas such as medical research or consumer behaviour forecasting given the inherent non-stationarity and noisiness of financial time series. The aim of this project is to build on existing research and try to gauge the efficacy of applying machine learning methods to the investment process. The motivation for this thesis is to try to cut through some of the recent hype generated in financial media and try to uncover how machine learning may best be used by financial practitioners.

# Contents

<b>1</b>	<b>Introduction</b>	<b>3</b>
1.1	Project Goal . . . . .	3
1.2	Project Objectives . . . . .	4
1.2.1	Data Retrieval, Cleaning and Standardisation . . . . .	4
1.2.2	Feature Extraction . . . . .	5
1.2.3	Model Training and Architecture Selection . . . . .	5
1.2.4	Model Validation and Strategy Backtesting . . . . .	5
1.2.5	Testing on Fully Withheld Data . . . . .	5
1.3	Project Organisation . . . . .	5
<b>2</b>	<b>Description of the Problem</b>	<b>6</b>
<b>3</b>	<b>Literature Review</b>	<b>7</b>
3.1	Heteroskedasticity in Financial Markets . . . . .	7
3.2	Machine Learning Model Description . . . . .	8
3.2.1	Artificial Neural Networks . . . . .	8
3.2.2	Support Vector Machines . . . . .	11
3.2.3	Ensemble Methods: Random Forests . . . . .	12
<b>4</b>	<b>Methodology</b>	<b>14</b>
4.1	System Specifications . . . . .	14
4.2	Data Description . . . . .	14
4.3	Data Pipeline . . . . .	15
4.3.1	Data Cleaning . . . . .	15
4.3.2	Standardisation . . . . .	15
4.3.3	Feature Extraction . . . . .	16
4.3.4	Justification of Features . . . . .	17
4.4	Training and Testing Methodology . . . . .	19
4.5	Model Evaluation Methods . . . . .	19
4.5.1	Assessing Model Performance . . . . .	20
<b>5</b>	<b>Project Design and Implementation</b>	<b>22</b>
5.1	Model Architecture . . . . .	22
5.1.1	Trading Model Methodology . . . . .	22
5.2	Machine Learning Model Parameter Specifications . . . . .	26
<b>6</b>	<b>Results, Analysis and Discussion</b>	<b>27</b>
6.1	Performance of Model One: Trend Capture . . . . .	30
6.1.1	Decision Trees: Random Forest . . . . .	30
6.1.2	Support Vector Classifiers . . . . .	31
6.1.3	Long Short Term Memory Trading Model . . . . .	33
6.2	Performance of Model Two: Regime Estimation . . . . .	33
6.3	Performance on Blind Test Data . . . . .	37
<b>7</b>	<b>Conclusions</b>	<b>40</b>

# 1 Introduction

The velocity of technological change that society has encountered since the turn of the millennium is striking, the world has embarked on a computational revolution which has brought the availability of immense processing power to the masses. Computational power at a scale that government sponsored computer science programs of the early nineties could only dream of. In fact, the processing power of most smartphones are larger than the supercomputers of yesteryear [1].

Moore's Law, the well known law which states that the number of transistors in a dense integrated circuit doubles about every two years, has held for over fifty years [2].

The allure of applying machine learning methods to finance is well entrenched, where researchers dream of uncovering repeatable patterns within the data which guarantees consistent alpha generation (returns unexplained by conventional market beta factors [3]). A quick search of Google, another revolution of big data processing and machine learning algorithms, uncovers a vast array of academic papers which profess of having found patterns in the data that could only have been uncovered by complex models with non-linear capabilities.

Using machine learning models in conjunction with big data can and do provide fantastic results in the area of the biological sciences or as recommender systems for movies or advertisements. However applying these same methods to financial markets is a different and more challenging task due to the underlying unpredictability of asset prices.

This paper delves into the world of financial machine learning and explores certain techniques to try to uncover how machine learning methods can be best applied to financial markets.

## 1.1 Project Goal

The main goal of this project is

- To create a fully systematic trading algorithm which uses various technical price based features to capture the trend in an asset price timeseries.
- To overlay a trend capture model using a separate machine learning model to identify regimes which are conducive to the positive performance of trend models.
- Finally to test and validate this particular trading model architecture across a host of different parameter variations, testing timeframes and features to thoroughly test the models generalisation capability.

The methodology behind creating prediction models in financial markets will differ vastly from similar projects in non financial domains, financial markets are inherently unstable and contain a substantial amount of noise which makes applying complex algorithms in finance susceptible to out of sample failure [4].

The task of uncovering repeatable patterns in financial data remains incredibly precarious, financial markets are non-stationary, exhibit leptokurtic (fat tailed) distributions and can shift regimes at highly unpredictable speeds. Much has been written on the complexity of financial markets and how researchers must overcome such dynamic systems [5]. While the explosion of increased computational power aided the application of machine learning in very useful tasks such as speech recognition and image processing, the applicability of these methods in finance are, on the face of it at least, questionable [6], thus this project has placed much greater emphasis

on the robustness testing of the selected features and models, while also trying to keep the underlying features as intuitive and simple as possible.

## 1.2 Project Objectives

The main aim of the project is to understand how machine learning methods can be applied to currency forecasting and the classification of macroeconomic regimes in which a trend trading strategy may perform best.

The project attempts to outline best practices to use when applying machine learning techniques to currency trading. As the project progressed, it was clear that feeding the model masses of data in the hope of finding a successful pattern would not generalise very well in the future.

Much has been written about the lack of verifiable and replicable results in machine learning research which suggests the capacity to overfit is high [7]. Hence the approach taken in this project was to assign one machine learning model with a very narrow goal, not to predict the next days prices based on a host of different technical and fundamental features, but to simply capture the *trend* in the price series.

The difference between a model created to capture trend vs a model created to predict the next day or week ahead price is subtle but important.

Firstly, trend can be captured in very simple terms with few features (thus reducing any dimensionality issues) while it is also not reliant on the stability of correlations between non price data (such as economic indicators) and the price series itself. The only goal of a trend capture model is to identify the direction of the trend in the price series with the expectation that the machine learning model can identify the non linearities inherent in trend reversals which more simple linear trend models may not be able to capture. Thus the problem domain that the model needs to assist in solving becomes not one of outright price prediction but of trend prediction. Thus if there exists auto correlation in the price series, then the model is able to identify this trend and trade in the direction of that trend, while if no trend exists (or the asset price is exhibiting a mean reverting behaviour) then the model is expected to underperform.

This is where the second model is assigned the task of using economic and fundamental data to try to identify when the underlying price series has a higher probability of exhibiting sustained auto-correlation. Thus the key here to profitability, is not the ability to forecast the underlying price series on a day to day basis, but to create a strategy which has an asymmetric risk profile, where in trending markets the combinations of models perform well while in mean reverting regimes the second model reduces or switches off the first trend capture model and by corollary reduces the severity of losses. The appeal of this type of model architecture is that it is price series agnostic, such that it should work as well on currency data as it should on equity price data. This means one could relatively easily create a cross asset trading model which would likely mimic the returns of a Commodity Trading Advisor (CTA), a type of hedge fund which profits from trending price series.

The main objectives in order to achieve the stated aim of the projects are shown below.

### 1.2.1 Data Retrieval, Cleaning and Standardisation

- One of the most important aspects of a trading model is to ensure that the underlying data is of high quality and doesn't contain erroneous values. The proposed model will use hourly spot exchange rates for G10 currencies <sup>1</sup>, various economic data from the Federal

---

<sup>1</sup>Hourly FX Spot Data Provided by Citi Foreign Exchange

Reserve Economic Database (FRED), the International Monetary Fund (IMF) and the Organisation for Economic Cooperation and Development (OECD) will be used to try to uncover exploitable patterns in the data.

- This data will also need to be cleaned and mapped to correct date scales so that the model only uses data available at the time of trade signal generation. The raw data itself will need to be standardised to allow the model to understand the true relationships between the feature vectors and the output.

### **1.2.2 Feature Extraction**

- Feature extraction for financial time series prediction is a very challenging task due to the inherent noise in the data. Including features which have an intuitive justification for predicting future prices is important [8]. Feature selection driven by mining the training set greatly increases the risk of overfitting and should be avoided when working with financial times series.

### **1.2.3 Model Training and Architecture Selection**

- The project itself will test the efficacy of a range of algorithms in trying to predict currency movements, this project will also make use of non price data in order to assist the model in understanding potential regime changes in the trading environment.

### **1.2.4 Model Validation and Strategy Backtesting**

- Once a model architecture has been selected, it is important to understand not only classification rates but also how the model would perform in a live market setting. This can be approximated by creating a backtesting 4.5 framework which depicts the model's trading performance and also includes the impact of transaction costs.

### **1.2.5 Testing on Fully Withheld Data**

- While the test data allows us to understand how the model performs, there will also be an opportunity to analyse how the selected model will work on new unseen data which provides a useful guide as to how the model will generalise in the future.

## **1.3 Project Organisation**

The project is structured as followed,

- Section 2 provides a description of the underlying problem of applying machine learning (ML) to financial timeseries.
- Section ?? provides an overview of the latest research being undertaking in the financial machine learning domain.
- Section 4 dives into the data set needed and the specifications of the system architectures used in this project alongside the model features used and the reason for their selection.
- Section 5 delves into the core of the model structure and how two separate methods are combined to create the full life cycle of the strategy.

- Section 6 looks into the performance metrics of each model as well as the final model architecture used and the structure of that model.
- Section ?? looks at the main takeaways of the results and offers thought on the direction for future work that should be undertaken when applying machine learning models to financial trading.
- Section 7 concludes with final thoughts on the project overall and highlights some interesting points which undertaking the project has uncovered.

## 2 Description of the Problem

The problem to overcome when applying ML to trading is the level of signal to noise in the data. ML practitioners in finance face a very challenging task as the past is rarely a perfect guide to the future. For example in image processing, if a neural network is fed a matrix of pixels which form the image of cows, and then trains based on the historical patterns it observes, the trained neural network will be able to identify a cow with high classification rates. However, if we provide an image of a mouse, the model likely classifies the mouse as a small cow due to the historical patterns that have been learned. Machine learning models, no matter how complex or sophisticated the underlying features are, require regularity for prediction. In finance, you can learn all the historical patterns and probabilities surrounding an asset and the model may even perform well in the testing period, however, in finance the rules can change rapidly. Structural breaks can and do occur in past relationships. Even the most robust artificial intelligence would likely not survive the impact of a 280 character tweet from Donald Trump entirely changing the rules of the game. Thus, the approach taken in this project is to narrow the scope of the is required from the machine learning model. This project asks each model to perform a very clear narrowly defined task which is to recognise the patterns of a trending asset price. This could be considered akin to using ML models in finance like a surgical scalpel rather than a hammer for every nail.

Correlations and patterns in markets can be extremely fickle, complex and dynamic in nature [9]. Why is this so one may ask?

Well part of the issue lies in the efficiency of informational flow across markets, barriers to entry are much smaller than ten or twenty years ago and there has likely been a compression in the informational edge that smart money investors (such as hedge funds) previously have had<sup>2</sup> [10]. The market price itself is a reflection of all these moving parts and thus even if a robust pattern has been found, it is likely to be arbitrated away as more market actors step in to trade away any mis-pricing [7]. This is quite literally like an image recognition model having to recognise cats which have the ability to morph into dogs.

One issue that is particularly pertinent in the financial domain is the issue of data dimensionality and the underlying size of each individual data set where even if there were fifty years of daily financial data available, it amounts to less than 20,000 data points<sup>3</sup>. This means that one must think carefully before adding features to the model.

The use of exhaustive grid search algorithms will simply overfit the model and should not be used. This is due to the fact when we change and search for optimal hyperparameters, given the number of possible variations and only one small set of live price points, the probability of

---

<sup>2</sup>Not to mention the crack down on information sharing between investment banks and hedge funds

<sup>3</sup>Assuming a 252 trading day calendar

finding a strategy by luck remains extremely high. Arnott et al. provide an example of a trading strategy which would have passed even the most stringent evaluation methods in both training, validation and testing phases which actually turned out to have been created using a ranking of stocks based on alphabetical ordering [8].

Financial markets have been shown to exhibit heteroskedasticity in their returns [11] and as such, creating a robust systematic trading strategy relies on more than just a simple data input/signal output type of architecture. In trading systems, the classification of regimes and subsequent buy or sell signals must also be compared vs the expected value of the that trade. A model with an 80% classification success rate will under perform if it mis-classifies the 20% of samples which result in outsized asset moves.

Given the inherent challenges faced when developing systematic trading strategies, one may wonder if creating a strategy based on machine learning methods is a fruitless endeavor. Due to the unique issues related to applying machine learning in the financial domain, this project took a different view on how best to use machine learning by first creating a model to only predict trends using price data and then allowing a second model to try to capture the likelihood of auto correlation in the price series increasing. This diversifies the trading model as now the problem of alternative/changing drivers is reduced (as the model only trades previous price trend which could be driven by any factor) while the second model will increase the allocation of risk capital to the model when it deems the probability of trend working is high, thus it is trying to approximate the probability distribution of future trend regimes.

### 3 Literature Review

#### 3.1 Heteroskedasticity in Financial Markets

Figure 1: Growth of \$100 Invested Since 1975



Source: Market Watch [12]

Complex systems that continuously evolve make applying ML techniques in a way that results in a robust model very difficult [8]. Financial markets have been shown to exhibit heteroskedasticity in their returns [11] and as such, creating a robust systematic trading strategy



relies on more than just a simple data input/signal output type of architecture. In trading systems, classification of regimes and subsequent buy or sell signals must also be compared vs the expected value of the that trade. A model with an 80% classification success rate will under perform if it mis-classifies the the 20% of samples which result in large moves in the asset.

The corollary of that is by only targeting classification of outliers you reduce your sample set and thereby reduce the statistical validity of the results. There is rarely a free lunch in financial markets and as such there will be no free lunch when using machine learning methods to predict financial times series.

For example figure 1, shows that missing out on the best 25 days of the Standard and Poors 500 companies equity index (S&P 500) over the last fifty years would have been very costly [12]. Given that 25 trading days represents only 0.005% of the total trading days , returns are significantly impacted. This shows the importance of being correctly positioned for outlier events and also the fact that one model can be severely impacted by one off or black swan type events [13]. Trend based models tend to exhibit positive skew as they can pick up and ride the wave on this outlier events XXXXciteGreenwood2017XXXX. Previous research has centered on price based data possibly as it is easily available and also because it is one dataset that is likely to have higher frequency data available ( [14], [15], [16]).

## 3.2 Machine Learning Model Description

### 3.2.1 Artificial Neural Networks

Artificial Neural Networks (ANN) were conceptually first written about in the early 1940s [17] when neurophysiologist Warren McCulloch and mathematician Walter Pitts wrote a paper on the theoretical functioning of neurons. Donal Webb then expanded on this in his paper "The Organization of Behavior" while the first computerised neural networks were developed by IBM in models called "ADALINE" and "MADALINE" [17]. At their core, ANNs are a series of large matrix multiplications where weights and biases, which represent the strength of connections between nodes, are optimized to provide the best approximation to the desired output.

The popularity of ANN methods have grown as enhanced computational power led to the lowering of barriers to entry in the field. One of the more popular networks used in practice is the mutlilayer perceptron, initially introduced by Frank Rosenblatt in 1958 [18]. A perceptron is the simplest form of a layered network and consists of a single neuron with an adjustable weight. The simplest form of the perceptron can divide the hyperplane into a 2-dimensional space which is then used to classify the given inputs. The hyper plane is defined by the linearly separable function below.

$$\sum_{i=1}^n x_i w_i - \theta = 0 \quad (1)$$

Multi Layered Perceptrons (MLP) can be defined as

$$NeuralNet_j^l = \sum_{i=1}^{N_l} w_{ij}^{l-1} y_i^{l-1} \quad (2)$$

where  $NeuralNet_j^l$  is the  $l$ th layer of the  $j$ th neuron which gives us the weighted sum of its inputs. The weights of the connections to the next layer  $l = 1$  are used to calculate the weighted

sum of inputs for the subsequent layer in the network.

$$Y_j^l = \text{ActivationFunction}(\text{NeuralNet}_j^l) \quad (3)$$

This weighted sum of inputs is then passed through to an activation function which maps the output of the hidden layer to a space defined by the activation function, in this step we can introduce non linearity into the model itself by using a non linear activation function such as the softmax function described below.

$$y = \frac{\exp(x_i)}{\sum_{j=0}^k \exp(x_j)} i = 0, 1, 2, \dots k \quad (4)$$

This maps any series of values onto a plane of [0,1] and the output depends on the values produced by the node which feeds the activation function.

While there are a vast array of model training techniques, one method to train an ANN is to use backpropagation, where updates to the network's weights are fed backwards through each layer. Each model has an error function which is used to grade performance, which could be the squared difference between the models output and the desired output (mean squared error).

The universal approximation theorem [19] of the MLP model in neural networks states that ANNs can approximate any function and thus in theory, should be able to identify any patterns within the data that another method would also identify.

The temptation to introduce multiple layers with thousands of connections in order to find a solution is ever present, however (and especially in the realm of finance) this increases the risk of over-fitting or selecting a model based on sheer luck, which fails when tested on live unseen data. This project uses a recurrent neural network architecture as standard neural networks have limitations. They rely on the assumption of independence among the training and test examples. After processing each data point, the entire state of the network is lost. If the data points are generated independently (such as an individuals medical diagnosis), this presents no problem. But if data points are related in time or space, this approach is flawed.

In financial assets, the market reacts and learns from the history of the price series, thus data points are indeed related. Recurrent neural networks look promising as they learn the objective function by analysing sequences of points (hence they have been used extensively in text sentiment analysis where context of words within the sentence is important).

This memory element can be very useful in financial timeseries as past behaviour tends to dictate the future prices path dependancy. The specific model architecture this paper uses is the Long Short Term Memory model (insert reference). LSTM has traditionally been used in the prediction of sequences of text in natural language processing. The use of LSTM model in finance is gaining popularity with XXXXX insert references all using LSTM for time series prediction. LSTM was introduced by [20] and makes use of various "gates" to control the flow of information across various nodes.

$$F_t^i = \sigma(W_F u_t^i + V_F h_{t-1}^i + bias_f) \quad (5)$$

$$I_t^i = \sigma(W_I u_t^I + V_I h_{t-1}^I + bias_i) \quad (6)$$

$$O_t^i = \sigma(W_O u_t^i + V_O h_{t-1}^i + bias_o) \quad (7)$$

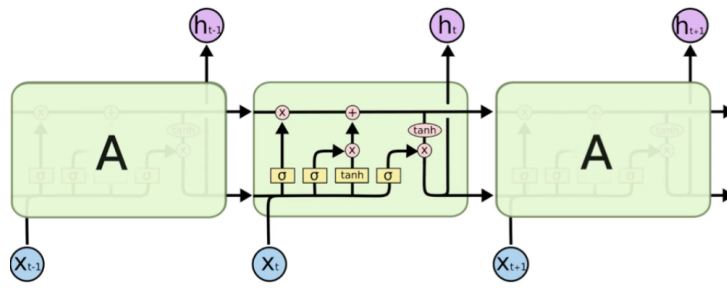
$$C_t^i = F_t^i * C_{t-1}^i + I_t^i * \tanh(W_C u_t^i + V_C h_{t-1}^i + bias_C) \quad (8)$$

$$H_t^i = O_t^i * \tanh(C_t^i) \quad (9)$$

$$Z_t^i = g(W_z H_t^i + bias_z) \quad (10)$$

where  $*$  denotes element wise multiplication,  $\sigma$  is the sigmoid activation function,  $W$  and  $V$  are weights matrices of different layers, and  $F_t^i$ ,  $I_t^i$ ,  $O_t^i$  are the forget, input, output gates and  $C_t^i$ ,  $H_t^i$ ,  $Z_t^i$  are the cell, hidden and outout states of the LSTM. The LSTM uses the cell states summarise past information and controls the flow of historical memory uses the forget gate while the pass through of new information is set by the input gate. Thus the LSTM model should be better equipped than an ANN in learning the long term relationships relevant to assisting in the prediction of the target variable. Figure 2 depicts the structure of a Long Short Term Memory Model.

Figure 2: Long Short Term Memory Model



Source: colah.github.io, Google Images

**3.2.1.1 Background Research** ANNs have been used extensively in academic financial journals, Czekalski et al. made use of feedforward multilayer perceptrons to predict the direction of the EURUSD currency rate [21], while Wang et al. applied probabilistic neural networks and Least Square Support Vector Machine models to assist in the prediction of Chinese equities [16]. Literature on using ANNs to forecast currencies has a twenty year history as even at turn of the millennium Joarder et al. tested three separate ANN based forecasting models using standard backpropagation, scaled conjugate gradient and backpropagation with Bayesian regularisation to predict movements in various Australian dollar currency crosses [22]. Gunduz et al. tested a variety of ANN models and Support Vector Machines (SVM) to try to predict the price of equities traded on the Istanbul Stock Exchange, they found model classification

accuracy rates consistently around 70% on test data [23] .

Most of the research cited makes use of various technical indicators which are derived directly from the price series, for example, Gunduz et al. make use of ten price based technical indicators to uncover patterns within Turkish equities [23]. Prior research suggested above 70% classification rates were attainable on average, however in contrast to this, there is evidence that ANNs often show inconsistent performance on noisy data ( [24], [25], [26]), which highlights the challenges in applying complex models as the risk of overfitting the data remains an ever present threat.

### 3.2.2 Support Vector Machines

One model which can be used to find patterns in financial data are Support Vector Machines (SVM) and Support Vector Regressions (SVR). SVMs were initially introduced by Cortes and Vapnik [27] which led to the development of non-linear models which provide binary classifications of target vectors while SVRs output real-valued predictions. In similar fashion to an ANN, the SVM tries to classify an n-dimensional feature vector according to where that data lies on an n-dimensional space.

The model makes use of the kernel trick [28] to transform the original data into a higher dimensional space such that we can then uncover a separable hyperplane. The function updates the parameters of the kernel function such that it can more accurately separate and classify the data points. The error correction model in the case of a SVM or SVC (Support Vector Classifier) is a maximisation problem where the model tries to maximise the distance between all the separating hyper planes (maximal margin classifier) it has constructed. The solution to the optimisation problem can take the form of the following

$$Maximise_{\beta_0, \beta_1, \dots, \epsilon_1 \dots M} M_{margin} \quad (11)$$

Subject to

$$\sum_{j=1}^p \beta_j^2 = 1 \quad (12)$$

$$y_i(\beta_0 + \beta_1 x_1 + \beta_2 x_2 + \dots + \beta_j x_{ij}) \geq M(1 - \epsilon_i) \quad (13)$$

such that  $\epsilon_i \geq 0$  and  $\sum_{i=1}^N \epsilon_i \leq C$ . Where  $M$  is the width of the margin to be maximised,  $\beta_i$  denotes the weight of the coefficients and  $x_i$  is the  $i^{th}$  observation ( $i = 1, 2, \dots, N$ ) of the dataset,  $\epsilon_i$  is a variable whose value sets the degree to which individual observations are allowed on the wrong side of the margin or hyperplane.  $C$  is a non negative tuning parameter which sets out the total error allowed for observations to be on the wrong side of the margin [29].

**3.2.2.1 Background Research** Kara et al. utilised SVMs to discover patterns in Turkish equities achieving classification success rates of around 60% when tested on hold out data [30]. Both polynomial and radial basis SVMs were trained on a host of price based technical data with the radial SVM outperforming the polynomial SVM. The polynomial SVM also had a much higher variability in classification success rates and thus the radial basis may be best

suitable to predicting financial time series. The results of the radial SVM were significant with a test statistic of 3 and p-value of 1%. The radial basis kernel is used in this project.

Similarly, Patel et al. tested a selection of four different machine learning methods (ANN, SVM, Random Forest and Naive Bayes) using ten technical indicators [31]. The predictor space was also transformed into trend deterministic states, which map the original continuous values into discrete values based on the trend of the continuous valued indicator (downward trend = -1, upward trend = 1). Results suggested that model performance is improved when trend deterministic data is used instead of continuous data. They postulate that the models perform better on the trend deterministic data as this transformation inputs the underlying trend of the indicator which then allows the model to learn the relationship between the underlying predictor trends and the actual output trend.

SVM classification rates using this approach (for both radial and polynomial kernels) approached nearly 90% on average. Huang et al. also tested the performance of an SVM model (radial kernel) to predict Nikkei futures prices<sup>4</sup> and achieved classification rates of 73% on test data. The authors also found that combining model forecasts increased out-of-sample test results [14]. Overall, prior research has found interesting and impressive results for using SVMs as a method to predict future time series such as equities ([32], [15], [24]) and currency movements [33].

### 3.2.3 Ensemble Methods: Random Forests

Ensemble methods are a technique in machine learning which uses the outputs of multiple models to form a consensus around the prediction, while there are many different techniques one can use, this paper proposes to explore whether ensemble methods can be used to reinforce the accuracy of the model on unseen data. One approach is to use tree based methods [34], which involve segmenting the predictor space into a number of regions. Figure 3 looks at the prior 12 hour (x-axis) and 24 hour (y-axis) return of EURUSD with the subsequent 24 hour ahead move in the currency colour coded from red (negative return) to green (positive return). Even in such a simple split of the predictor space, we see some evidence that more subsequent positive days occur in the bottom left of the chart (when prior 12 and 24 hour moves in the currency have been negative). Tree based methods use this type of inference to make predictions on future data points. The goal of the decision tree is to find multiple regions which then minimise the squared error of predictions

$$\sum_{j=1}^J \sum_{i=1}^{R_j} (y_i - y_{R_j}^*)^2 \quad (14)$$

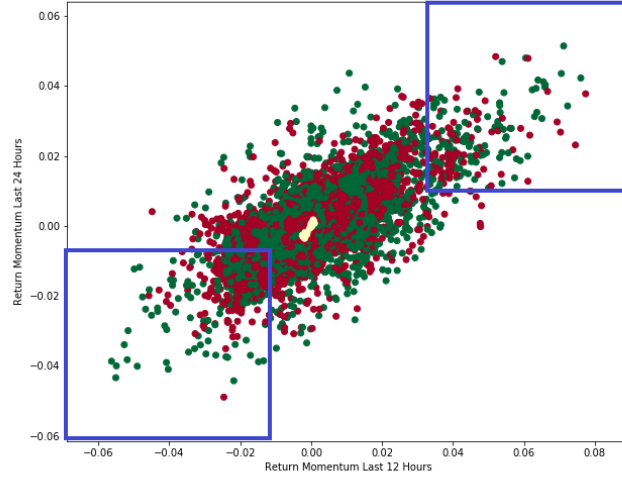
where  $y_{R_j}^*$  is the mean response for the training observation within the  $j^{th}$  region. The methods tested in the project will revolve around Random Forests (introduced by Leo Breiman [35]) which averages the prediction across hundreds of single decision trees. The trend estimation model uses a random forest architecture to identify trend conducive regimes,

**3.2.3.1 Background Research** Research has shown that Random Forests can generalise better than single decision trees as they maintain high variance while reducing the bias of the model [36]. Random forests use bootstrapping to take multiple random samples from the training data and a random subset of feature vectors to create a diverse array of decision trees which aid

---

<sup>4</sup>Nikkei 225 is one of the main Japanese Stock Market Indices

Figure 3: Tree Based Methods Segment the Predictor Space into Numerous Regions



Source: EURUSD FX Rate, Kaggle Database

model generalisation. This method helps to create decision trees which are uncorrelated to each other and thus have a lower probability of overfitting the data. Patel et al. showed that it was possible to achieve classification rates of over 80% for various stock price predictions when using the random forest algorithm and input features which consisted of technical price based indicators [31]. The research used a random selection of three feature vectors (out of a total of ten) to create each single decision tree in the random forest model and tested the results from 20 to 200 trees using a majority vote for the final classification.

Kumar et al. performed extensive tests to compare the performance of support vector machines to random forests and found that on average support vector machines outperformed random forests when predicting the direction of the Indian stock exchange [25]. Three predictors were considered for each node while the number of trees for each model ranged from 200 to 1500. There was no significant reduction in the error rate for the out-of-sample data based on the number of trees. The RF model achieved classification rates of over 67% on out of sample data. Chatzis et al. found that RF methods outperformed SVMs and logistic regression techniques when predicting stock market crisis events, they also note that RF models are relatively robust to overfitting due to each forest only being exposed to a subset of the available feature vectors [37].

## 4 Methodology

This section provides a deep dive into the underlying methodology of the full life cycle of the model, including how the data pipeline was structured, the training and testing methodology and various model evaluation techniques. The model architecture selected was to combine two different models, the first model with the explicit goal of capturing the trend in the time series while the second model would act as a regime identifier and allocate capital to the trend model when the current regime was conducive to the positive performance of trend. The model design is covered in greater detail in section 5.

### 4.1 System Specifications

The Python programming language (Python 3.6, [38]) will be used as the main programming application to achieve the objectives, with the python packages of Numpy [39], Scipy [40], Sci-Kit Learn [41], Matplotlib [42] and Pandas [43] all providing very useful data manipulation functions to help store, clean and process data. For more complex ANN learning, the machine learning package Keras [44] will be used in the training and testing of Neural network models. In order to train the more intensive models such as the Long Short Term Memory model, a cloud server was used, the specifications of the machine were as follows, 1 CPU with 1.7GB RAM and 250GB SSD. While not computationally very powerful, the machine was selected due to its low running costs which allowed the RNN to train on the dataset overnight. Total time spent running the server was over 250 hours.

### 4.2 Data Description

One of the most important aspects of a trading model is to ensure that the underlying data is of high quality and doesn't contain erroneous values. The trading model uses hourly spot exchange rates for G10 currencies<sup>5</sup> alongside various economic data and currency valuation data from the CitiFX and the Organisation for Economic Cooperation and Development (OECD) and national statistics offices.

The data used for this project consisted of hourly foreign exchange data for the EURUSD currency pair<sup>6</sup> and related economic and currency valuation data. Figure 4 shows the performance of EURUSD over the last sixteen years formed of hourly snapshots and its 55 day and 55 week moving averages. In total the hourly data was sourced from 2001 to the start of 2019, which consisted of just over 120,000 rows of data. One issue encountered was that training the LSTM (even using powerful GPU machines on the cloud server) took an extended amount of time to train. This made cycling through various model iterations challenging, with the process often left to run overnight. There is evidence of the existence of momentum in currencies [45] and this project will explore ways of capturing the trending nature of currencies more accurately than a simple linear model. Economic data can be obtained from a variety of sources, with the Federal Reserve Economic Database (FRED)<sup>7</sup>, IMF and OECD all providing vast treasures of low frequency economic data to create feature vectors. Relative economic growth<sup>8</sup> (as measured by momentum across indicators such as consumer spending, GDP growth, industrial production and business confidence) has been shown to assist in the prediction of currency movements [46].

---

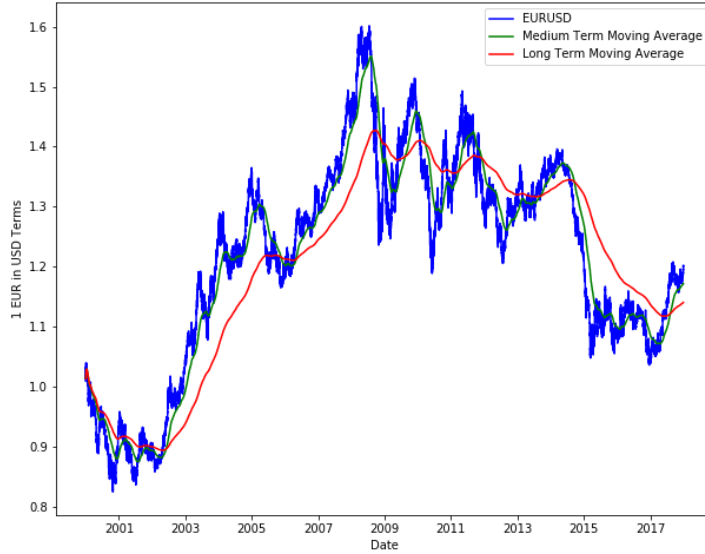
<sup>5</sup>Hourly FX Spot Data Provided by Citi Foreign Exchange

<sup>6</sup>Value of 1 European EUR in terms of the United States Dollar

<sup>7</sup>Successfully used in [37] to assist with forecasting stock market crises

<sup>8</sup>Citi 2019

Figure 4: EURUSD Rate Since 2001



Source: EURUSD FX Rate, Kaggle Database

### 4.3 Data Pipeline

#### 4.3.1 Data Cleaning

The data was also screened for errors and missing data points. Cleaning of the data was broken down into outlier detection (to capture any erroneous inputs) and cross checking the data vs known events in time (such as Brexit, election dates) to see if the observed outliers match the expected time frames. This data was also screened and mapped to correct date scales so that the model only used data available at the time of trade signal generation. Once the data had been cleaned appropriately, the next step was to standardise the data and then choose the model features.

#### 4.3.2 Standardisation

The data was stored and manipulated using Pandas dataframes while Sci-Kit Learn provided the functionality to run both the random forest and support vector regressions. While the package also has some very useful functions such as StandardScaler to standardise the data, in this project a bespoke scaling method was used where the data is standardised by the standard deviation of the feature over a rolling three year lookback window. This is a more realistic approach to use when creating systematic trading strategies as this method uses all available data up to the point of trading and expresses all point in relation to its most recent history, allowing the denominator (the standard deviation ) to vary over time as the market changes regime.

Standardising the data remains an important initial step that has been used across a wide



variety of prior research. Subtracting the mean and dividing by the standard deviation (z-scoring) [47], min-max standardisation<sup>9</sup> ([23], [25]) and other techniques such as mapping functions to scale feature vectors onto a pre-defined scale ([30], [16], [31]) all help to achieve the overall goal of scaling, which is to ensure that larger or more volatile features do not overwhelm features of smaller magnitude and lower volatility.

---

```
def standardise_data(dataset, full_cols, standardised_cols, window):
    '''
    This function computes the standardised returns on a rolling basis looking
        backwards over the window specified.
    This is the most realistic in terms of a trading strategy and also means the test
        data is standardised on the correct basis using the
    latest data available at each timestep.
    :param dataframe:
    :param cols:
    :return: standardised dataframe of values
    '''
    train_standardised =
        dataset[standardised_cols].subtract(dataset[standardised_cols].rolling(window).mean())
    train_standardised =
        train_standardised.divide(dataset[standardised_cols].rolling(window).std())
    # we will only return the data which is outside the initial window standardisation
    period
    # add non standardised features
    for feature in full_cols:
        if feature not in standardised_cols:
            train_standardised[feature] = dataset[feature]
    # This function now returns the necessary file with both standard and non
    standardised columns.
    return train_standardised.loc[window:,:]
```

---

### 4.3.3 Feature Extraction

The model features were chosen with the goal of capturing the underlying trend of the time series, it is very important in finance that the features are not chosen using a grid search of the hyperparameters. Feature selection driven by mining the training set greatly increases the risk of overfitting and should be avoided when working with financial time series [8]. Thus the features which had an intuitive justification for defining the trend are shown in the following table. Financial time series prediction is a very challenging task given the inherent noise in the data, thus the number of available features were kept to a minimum and in effect used only the deviation of the currency price from the short, medium and long term price averages.

Dimensionality reduction (via principle component analysis) was also performed on the features, however this did not have any positive impact on the models performance. The second model then allocate risk capital to the trend strategy based on longer term fundamental economic data. The feature selection process for model two (regime prediction) was based on economic intuition as this can lessen the probability of over-fitting [8].

---

<sup>9</sup>Map observations onto a scale of [0,1] with 0 being the minimum value of the observation and 1 denoting maximum

Table 1: Model One: Trend Capture

Indicator	Description	Type
Short Term Momentum	5,10, 21, 55 Period Moving Average	Technical
Medium Term Momentum	100, 155, 200, 255 Period Moving Average	Technical
Long Term Momentum	300, 350, 400, 455 Period Moving Average	Technical

As patterns and correlations in prices can change quite quickly, it is important not to overfit the model on the training data. The technical price based indicators are captured in the first model, while the second model measures non price based information such as economic momentum and valuation. Table 1 outlines the features which will be used in model two.

Table 2: Model Two: Trend Strength Estimation

Indicator	Description	Type
Relative Economic Momentum	1 Year Economic Momentum	Economic
Global Economic Momentum	1 Year Global Economic Momentum	Economic
Currency Volatility	Annualised Volatility	Technical Stress Factor

The approach of combining two separate models each with specific goals is a novel way to use machine learning in finance, in essence one model controls the price prediction while the second model is used to gauge how likely the first model is to generate positive returns.

#### 4.3.4 Justification of Features

The aim of the features used in this project was to try to juxtapose higher frequency technical price data alongside lower frequency economic data to assist in a model’s prediction accuracy. While price based metrics offer the easiest route to higher frequency data, they also lack the ability to capture the fundamental economic factors that can drive currency direction [48]. Thus combining two separate models to handle each aspect of the trade signal process can create synergies between each model. Below we showcase factors that can impact currency direction.

**4.3.4.1 Economic Fundamentals** One of the most important factors that sets the broad direction of a currency pair is the relative economic performance between two countries. Dahlquist et al. show that measuring relative economic momentum between different countries has the ability to forecast currency direction over medium term horizons [46]. That paper showed statistically significant results that economic fundamentals have an important role in relative capital flows and subsequent currency direction.

**4.3.4.2 Price Momentum** Price momentum here refers to the historical trend in the currency, research does suggest there exists exploitable momentum within the currency space [49]. The reasoning for the existence of momentum in currencies is likely due to decision delays amongst different investor segments who build up positions in currencies over time. Auto-correlation in currency returns may also exist as information is priced slowly over time (such as economic data) and thus as more economic data is published, this adds further confirmation to a trading signal and thus conviction and trade size increases.

**4.3.4.3 Investor Positioning** Currency positioning relates to the magnitude of investor holdings of a currency. Positioning and currency flows can have informational value that impacts a model's performance [50]. For instance, if the finalised model triggered a buy recommendation on a currency, but if the wider market had already built up a sizable long position in that currency, then the likelihood of further buyers coming in to drive the price higher is diminished.

**4.3.4.4 Currency Valuation** Uncovering the true fundamental fair value of a currency has been the focus of international academic research for many years [51], and it remains an elusive and rather ambiguous task as currency drivers can undergo structural shifts. Purchasing Power Parity currency levels, which rely on the Bassala Samuelson effect and the law of one price [52] are one method of trying to anchor the value of a currency to one fundamental price level. While currencies can deviate from fair value, there is some evidence that developed market<sup>10</sup> currencies tend to mean revert in the long term [53].

While the above drivers of currency markets are not exhaustive, it can provide a frame work for constructing a systematic trading strategy. Currency drivers can shift over the short and medium term and while this can all be partially captured by price action, it is also prudent to try to understand the dynamics of other potential drivers as discussed in the previous section. Therein lies the main problem when combining machine learning and finance, that other areas such as the biological or physical sciences do not encounter, is that at the core of each price, are human or human created trading algorithms which have collectively decided on a price level at which to transact.

Every price of a trade in financial markets consist of the dreams, fears and greed of human emotion. Not only this, but market investors are influenced by past patterns and incorporate new information (including newly published research) into future decisions, thus altering the environment in which the historical pattern has been learned. Bartram at al. provide evidence that anomaly (excess) profits published in research papers tend to decay substantially in the period after publication, which suggests any new findings are arbitrated away<sup>11</sup> quickly as the new information is incorporated into the market [7].

While early academic research espoused the efficient market hypothesis, there are certain times where the market is not fully rational [54]. Thus opportunities do exist, but finding repeatable opportunities requires more than just analysing price patterns in data. Given the randomness of markets, even poor strategies can perform well for long periods of times due to sheer luck. This creates new challenges for practitioners of machine learning in finance, in fact, new research has also started to focus on protocols and best practices around data mining in finance [8].

Being able to identify trends would allow the model to outperform when there is strong presence of momentum while not overly churning the portfolio if there is a lack of momentum in any direction. Building the model to this asymmetric specification can assist in the generalisation of the model.

Thus it was decided to first train a model to identify trends in each series and then overlay that model with another capable of using the non price data to identify periods when currencies exhibit trending behaviour. For example, periods of elevated economic dispersion could be conducive to trend strategies outperforming.

---

<sup>10</sup>Developed Market Currencies consist of USD, EUR, JPY, CHF, GBP, AUD, NZD, NOK, SEK, CAD

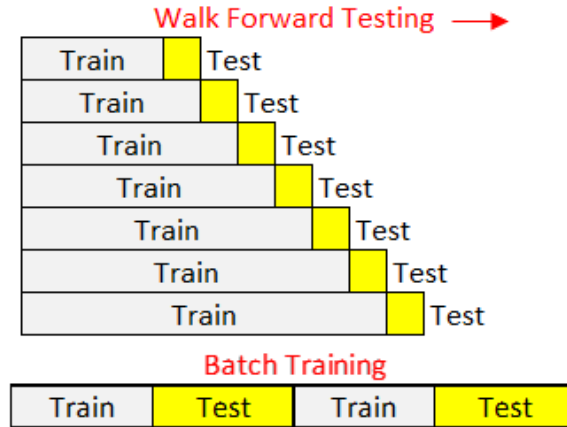
<sup>11</sup>Of course a more sinister reason could be model overfitting on historical data

#### 4.4 Training and Testing Methodology

The model training stage allows the learning algorithm to understand and capture the relationship between the features and the classification target vector (the target vector will be the future direction of the asset price across various time periods). Prior research has generally used binary classification of the future asset price move when developing models( [55], [23], [37]).

Model training allows us to understand which combination of models (RNN, SVM, RF) and underlying architecture (hidden layers, number of decision trees, kernel type) are best suited to capturing the nuances in the data. This stage of the process also provided an opportunity to experiment with various models to see the differences in performance. The training and testing phases were performed over multiple different windows such as walk forward testing, K-fold testing 5 and simple train and test splits over the full dataset. The structure of the different training/testing methods are shown below.

Figure 5: Walk Forward and K-Fold Train/Test Split



The results of training and testing on the price data will be further explored in section 6.

#### 4.5 Model Evaluation Methods

In order to validate the model, the standard procedures of analysing the mean-squared error, classification hit rates and precision and recall will be analysed. However, given this is a trading system, the selected model and its structure will need to be tested in a similar setting to real world trading. This requires the construction of a backtesting engine which can simulate each models performance taking into the signal lag <sup>12</sup> and the currency performance per each trading signal. This will provide a much more accurate analysis of the selected trading strategy performance<sup>13</sup>. While classification success rates will be important in gauging model accuracy, high classification rates may not lead to positive strategy performance if large currency moves are mis-classified. The code for the trade backtesting engine is shown below .

<sup>12</sup>The time between signal generation and trading time

<sup>13</sup>The vast majority of research papers referenced in this thesis did not perform this type of testing

---

```

def backtester(results, test, trade_horizon):
    """
    Calculate the returns of the trading strategy and store the results in the test
    file.
    :param results: the trading signals generated by the model.
    :param test: this contains the test data
    :param trade_horizon: this is the length of time of the trading signal prediction
    :return: A dataframe which has the test data and also the results of the trading
            strategy in pct terms of log returns.
    """
    # This needs to change to handle the change in the target
    predictions = pd.DataFrame({"Date": test['Date'], "Predictions": results})
    test_results = pd.merge(test, predictions, how="left", on="Date").fillna(0)
    # calculate the returns of the signal
    test_results["erf_signal"] = test_results['Predictions'].apply(erf)
    test_results["scaled_signal"] =
        test_results['Predictions'].shift(2).rolling(trade_horizon).sum() /
        trade_horizon
    test_results["scaled_signal_erf"] =
        test_results['erf_signal'].shift(2).rolling(trade_horizon).sum() /
        trade_horizon
    # no shift needed as we have already done that in previous step
    test_results['strat_returns'] = test_results['logret'] *
        test_results['scaled_signal']
    test_results['strat_returns_sum'] = test_results['strat_returns'].cumsum()
    # calculate returns for the error function
    test_results['strat_returns_erf'] = test_results['logret'] *
        test_results['scaled_signal_erf']
    test_results['strat_returns_sum_erf'] = test_results['strat_returns_erf'].cumsum()
    strat_return = test_results['strat_returns'].sum()
    # provide the information ratio on an annualised basis.
    information_ratio = (test_results['strat_returns'].mean() * 260) /
        (test_results['strat_returns'].std() * np.sqrt(260))
    return (test_results, strat_return, information_ratio)

```

---

#### 4.5.1 Assessing Model Performance

**4.5.1.1 Mean Squared Error** The mean squared error (MSE) is used to measure the prediction accuracy of the model. It looks at the average distance of the predicted value from the actual value.

$$MeanSquaredError = \frac{1}{N} \sum_{i=1}^n (y_i - F(x_i))^2 \quad (15)$$

where  $y_i$  is the actual value and  $F(x_i)$  is the predicted value for the  $i_{th}$  observation. Normally it is best practice to choose the model which has the lowest test MSE.

**4.5.1.2 Precision and Recall** Another method of model validation is to look at the precision and recall of the model [31], whose values can be calculated from a confusion matrix, which

shows the number of correct and incorrect classifications per class. Table 3 shows the confusion matrix for a two class classification problem.

Table 3: Confusion Matrix

Actual (row) Predicted (col)	Buy	Sell
Buy	True Positive	False Negative
Sell	False Positive	True Negative

The confusion matrix allows us to calculate the precision and recall of the model.

$$Precision = \frac{TruePositive}{TruePositive + FalsePositive} \quad (16)$$

$$Recall = \frac{TruePositive}{TruePositive + FalseNegatives} \quad (17)$$

The precision captures the correct positive signals ( e.g. in this case the correct number of buy signals as a percentage of all buy signals) that were actually triggered. Recall measures model accuracy as a function of the total number of positives in the data set itself. These statisitcs are important to analyse when assessing trading models as certain periods can be heavily biased towards one class, for example in the case of a financial asset which has a strong trend will by design have a higher occurrence of one class vs another.

**4.5.1.3 Deflated Sharpe Ratio** Given the inherent issues of overfitting in finance ([8], [4]), it is also important to analyse the overall Sharpe Ratio [56] which measures the return of a trading model per unit of risk (or volatility of returns).

High Sharpe ratios can be obtained by a genuinely good prediction model or through pure overfitting/selection bias, hence it is useful to deflate the Sharpe Ratio of every strategy tested to take into account the number of model tests as well as the skew of each strategy's distribution [57]. The Deflated Sharpe Ratio is a test statistic which can outline the probability of having found a true positive. It is defined below as

$$\widehat{DSR} = \left[ \frac{\widehat{SR} - \widehat{SR}_0 \sqrt{T-1}}{\sqrt{1 - \gamma_3 \widehat{SR} + \frac{\gamma_4 - 1}{4} * \widehat{SR}^2}} \right] \quad (18)$$

where  $\widehat{SR}_0 = \sqrt{V[\widehat{SR}_n] \left( (1 - \gamma) Z^{-1} \left[ 1 - \frac{1}{N} \right] + \gamma Z - 1 \left[ 1 - \frac{1}{N} \exp(-1) \right] \right)}$  and  $\widehat{SR}$  is the estimated Sharpe Ratio to be tested,  $T$  is the sample length,  $\gamma_3$  is the skewness of the returns distribution and  $\gamma_4$  is the kurtosis of the selected strategy,  $V[\widehat{SR}_n]$  is the Sharpe Ratio variance across the number of strategies tested,  $N$  is the number of independent trials of strategies run (i.e. the number of models tested) and  $Z$  is the cumulative function of the Gaussian distribution<sup>14</sup>.

<sup>14</sup>Further information and theoretical proofs of the deflated sharpe are available from [57]

## 5 Project Design and Implementation

This section outlines the specific steps taken to create each model and provides some justification for the selection of various model parameters and specifications used. Various Python packages such as Sci-Kit Learn and Keras were used to create and test the classification success rates of each type of model. Fischer et al. made use of Keras (on top of Google's Tensorflow) to implement Recurrent Neural Networks (RNN) [47] while using Sci-Kit Learn for logistic regression, SVM and Random Forest methods.

One issue encountered was that the length of time needed to train the LSTM meant that training multiple variations across parameters was challenging from a time perspective.

### 5.1 Model Architecture

Three model architectures were used in this project.

- Recurrent Neural Networks: Long Short Term Memory Models
- Decision Trees: Random Forests
- Support Vector Machines: both classification and regression methods.

One novel approach in selecting the target vector (i.e. the future return of the model) was to use a risk adjusted return and predict the risk adjusted return as opposed to the binary return [?], this ensures that the trend model can also estimate the positioning size and negates the need for a separate trade risk adjustment to be made.

During the course of the project, it became clear that using historical price data to train the model produced erratic performance in the testing period across all three models and various training and testing phases. Perhaps this was due to the nature of the underlying security being EURUSD (one of the most highly traded instruments in finance and thus was highly efficient) or that the model tended to over fit the training data (even when using techniques such as drop out rates and tree pruning was performed).

#### 5.1.1 Trading Model Methodology

Thus given the failure of the model to learn robust profitable trading patterns its own price data, it was necessary to go back to basics and understand what the goal of the model and its features ultimately was. We are trying to predict the trend in the price series. If the series we are trying to trend has only exhibited trend in patches across the full time frame of available data, then the model is likely unable to pick this trend up over the full sample, as it ends up training on conflicting patterns. This was one fact that did emerge when training using K fold cross validation. In some instances the training data set happened to be during regimes of strong upward or downward trends, and subsequently the model would perform well in the test phase<sup>15</sup>. Thus for the model to be able to capture trends, it needed to be trained to recognise a trend, therefore, I decided to switch the training data away from the historical price series data and onto randomly generated trending data series which could replicate trends in asset price series. Of course when training on randomly generated data one must also include some noise term such

---

<sup>15</sup>The results of various model performance is covered in section 6

that we create series which mimic financial asset prices. Hence the trending data was injected with a noise term. Training on randomly generated data has two very appealing traits, one it reduces the chances of overfitting the data. The ability to overfit is reduced but not removed as there is a chance the randomly generated data series happens by chance to mimic closely the testing phase of the data (hence training of various iterations of the randomly generated data is important). Secondly, it opens up the full live price series data to exist as the test dataset. Hence by training the model to recognise trends, we can now test on the full 120,000 rows of EURUSD price data.

The randomly generated price series was created using the following formula.

$$X_t = X_{t-1} * \alpha + \epsilon_{t-1} \quad (19)$$

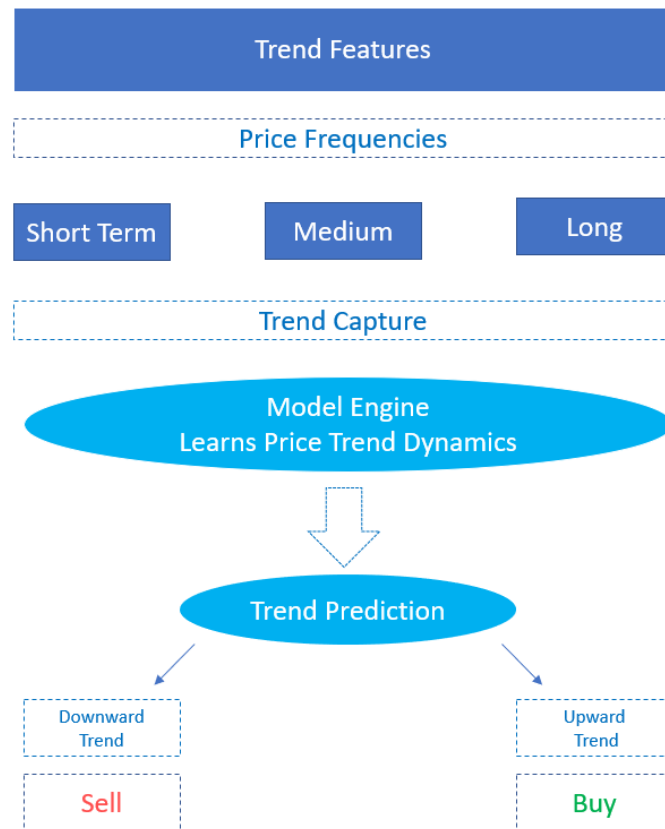
where  $X_t$  is the simulated price at time  $t$ ,  $\alpha$  is a parameter to control the strength of the trend and  $\epsilon$  is a randomly distributed i.i.d error term used to simulate the noise of financial time series.

The trend capture model (model one) is described in the below diagram 6. The figure outlines how the underlying features feed into the trend model and also how each selection of features are tailored to capture short, medium and long term trend frequencies. Model one handles the signal generation process, by measuring the underlying relationship between the features and the target return variable. As mentioned in 1, the features were created from price based technical signals, such as deviations from exponential moving averages across both short, medium and long frequencies. This purely price based factor model is used to capture trend only, hence we feed in only factors which are applicable to signaling a trend in either direction.

The overall goal of model one is to use machine learning models to identify the pattern and non-linearities inherent in trending prices which more simple linear based models fail to capture. The type of model also varied from using regression methods to estimate a continuous valued 24 hour period ahead and also a more simple classification on whether the price was predicted to move higher or lower. This one day ahead trend prediction was then used as the input into model two (trend estimation, 7).

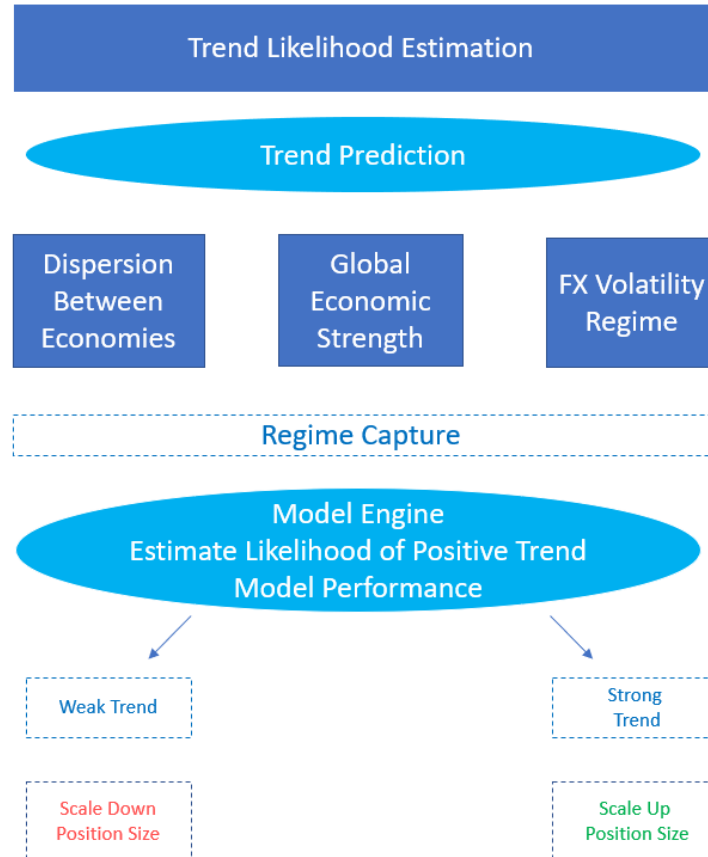


Figure 6: Model Architecture :Trend Capture



The second model (figure 7), then tries to capture the prevailing market regime, and uses much slower moving fundamental data as features (as described in 2) in order to assist in position sizing of the trend model, by learning when historical fundamental features have tended to coincide with the positive performance of the trend strategy. This work on the intuitive reasoning behind why a currency exhibits trending behaviour, one process of thought is that economic dispersion between countries [?] can lead to large trends between currencies as the capital is pulled towards the currency of the economically stronger country. Decision delays by different segments of investors can also influence the slow grind of a higher price, while central bank monetary policy can also lead to large deviations between different countries currencies. Central bank policy can often be driven by economic data and thus looking at economic data we maybe able to identify regimes which are conducive to trending prices series. We also include a volatility feature as well as global and G10 economic performance, which can act as a potential harbinger of further increases in volatility when trend is driven by factors such as risk appetite (where investors flock back to safe haven assets). The target for this model is the performance of the trend model one historically over the testing period. Thus we combine two different models in order to incase the chances of out of sample model generalisation.

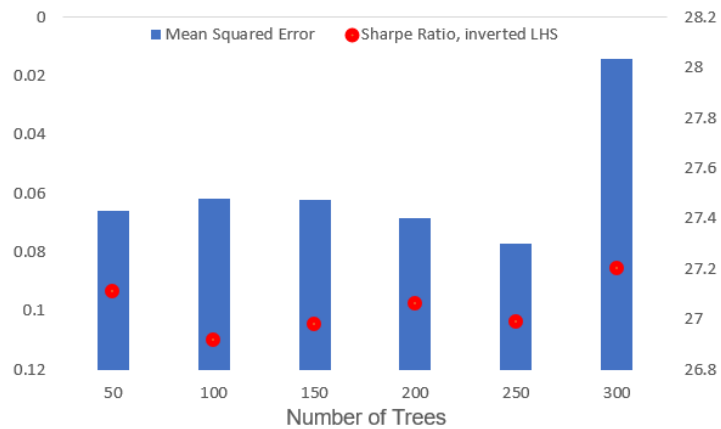
Figure 7: Model 2: Trend Estimation



## 5.2 Machine Learning Model Parameter Specifications

As mentioned previously, three different machine learning models were used in the project. The testing and training datasets remained standardised across each model so that comparison of results would be as accurate as possible. The python package Sci-Kit was used to build both the random forest and SVC (SVR) models. Parameter variation was performed across a different number of decision trees and from that analysis, 200 trees were used in the final model phase as it was relatively quick to train and also provided good results (??). When looking at the results of batch training and testing, we can see that adding more trees to the random forest model did not necessarily add more prediction power. Figure ?? shows the mean squared error scores and the sharpe ratio <sup>16</sup> of each model trained across various number of trees and using a 50/50 train and test split. The maximum number of features to be selected at each node (i.e. the max features parameter) was kept low to reduce the correlation between trees and increase model generalisation.

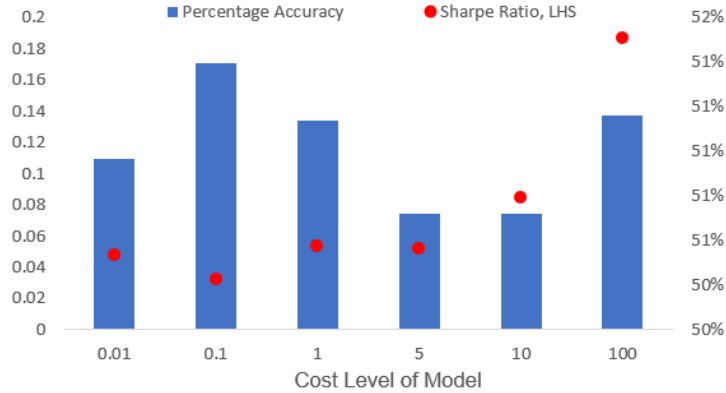
Figure 8: Random Forest Classifier Performance Across Various Number of Trees



Similarly for the Support Vector regression/classifier models used, a review of the appropriate cost to use was performed. The kernel used was the radial basis kernel to try to capture the non linear patterns in the price series. Interestingly, across various different cost levels, it was noted that higher costs levels tended to have slightly higher accuracy scores which also tended to translate into higher Sharpe ratios in the test period 9. The test period of the accuracy scores and Sharpe ratios shown here was from October 2014 until December 2017 ( the period of data post 2017 was truncated at the start of the project to be used for the final blind model testing). Overall, the fact that the higher cost allowances generally performed better in the test period shows the value of not over fitting the training model when using financial time series. The recurrent neural network is one model which is very powerful in terms of the complexity and the nature of problems it can in theory solve, however the increased complexity also comes with a loss of generalisation given the propensity to over fit the data. While this in mind, the following model architecture was chosen to train and test the model, an initial first layer of 32 nodes followed by a layer of 16 nodes where the soft max activation function was. The main parameter which was varied in order to assess model performance was the look back sequence

<sup>16</sup>assuming a risk free rate of 0%

Figure 9: Support Vector Classifier Performance Across Various Cost Levels



and the number of epochs used to train the model. Generally as the project progressed , the higher the number of training epochs and look back sequences the better the model performance, increasing the epoch and look back sequence size noticeably increased the computational time spent training the model.

---

```
# Initialise the model the Recurrent Neural Network model
model = Sequential()
# add a layer of 32 nodes, where the sequences across each observation are returned
model.add(LSTM(first_layer,batch_input_shape = (None,look_back,no_features),
    return_sequences = True))
# second layer has half the nodes of the previous layer and uses a softmax
    activation function.
model.add(LSTM(second_layer, return_sequences = False, activation="softmax"))
# finally the model is compiled using the MSE and the adam optimiser
model.compile(loss = "mean_absolute_error", optimizer="adam", metrics =
    ['accuracy'])
```

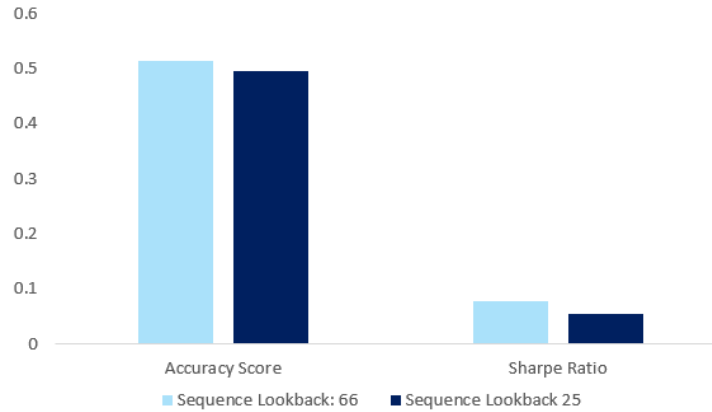
---

Figure 10 , shows the accuracy score and the sharpe ratio (in the testing period) for look back of 25 periods and 66 periods, the 66 period look back outperformed the shorter horizon, possibly as this longer sequence allowed the model to place current price action in the context of a longer window, similar to a human day traders using recent history to base his trading decision on when identifying a trend in the price.

## 6 Results, Analysis and Discussion

The initial model training and testing on the actual price series data did not perform very well, perhaps due to the model overfitting the problem or perhaps from the fact that if the model trained on periods when there was a weak trend, then the model was unable to identify the correct parameters and weights needed to capture any future trends in the price series. In figure 11, we show a scatter plot of the annualised return of trading EURUSD across various machine learning models, training and testing periods and model parameters. For a strategy to be robust in finance, a promising trading strategy would show clustering of the points (i.e.

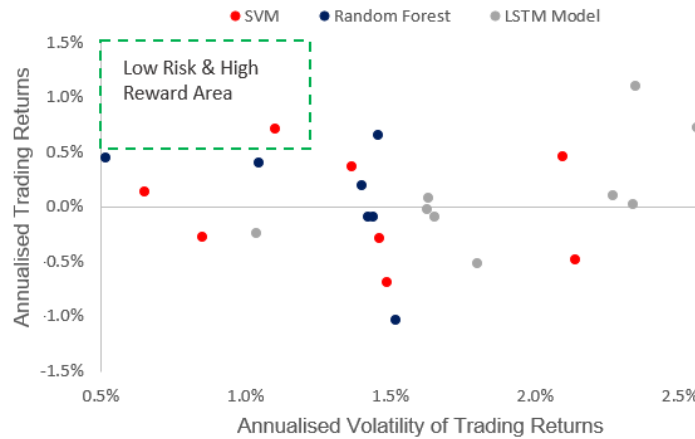
Figure 10: LSTM Model Performance Across Lookback Sequences



strategy variations )at the top left of the chart (which would indicate an area of desirable return per unit of risk taken) , however in this chart we can see quite scattered points, suggesting there is a severe lack of robustness in the underlying dynamics of the trading model.

This may be due to the features used or it may be due to the specification of each underlying model. However, given the justification for the features used in the model, one would be wary of using tweaked priced based features on the sole fact that they perform better, as that is blatant overfitting. One reason for erratic performance of each model is that each model training period may have occurred in periods where the currency was in a regime of mean reversion or weak trend, this could inhibit the performance of the model when in being able to recognise trend patterns. While when the model was trained on a trending price series, the model could learn the patterns associated with auto correlation in the price, and trade this profitably in the testing period.

Figure 11: Annualised Returns (Y-axis) vs Volatility (X-axis) Across Various Strategies

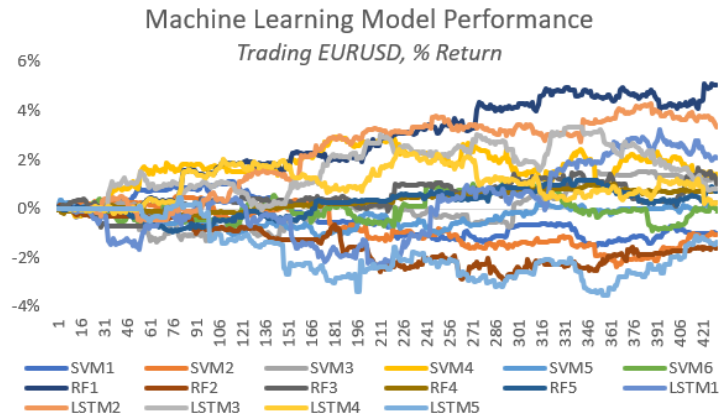


The poor performance of the various models on the live price data in capturing trend was disappointing, however, we refrained from running an exhaustive grid search type optimisation on the data as this will only lead to overfitting of the problem, when one runs so many simulations across different parameters (especially in finance) then it is inevitable that we uncover an acceptable strategy by luck (similar to a coherent novel being written by a billion monkeys typing on a typewriter). Given the inherent non stationarity of financial timeseries, it is vital to ensure the model performs in a stable way across a variety of parameters. This rules out running complex transformations or optimisations on the features (outside of kernel methods inherent in the SVM model) as the probability of having selected the "best" parameters by pure chance is high. In financial trading, model robustness is key.

As pointed out by Mclean et al. [58], financial model performance can differ vastly prior to publication than after it<sup>17</sup>, whether due to over fitting or post publication anomaly decay [7], this means we must be extra vigilante when assessing model performance.

In the sections below, we show the return profile for a host of different trading models performance across support vector machines, recurrent neural networks and random forest. The returns shown in 12 display the model performance for various different training and testing periods. The returns are calculated from the backtester outlined in section 4.5. The cumulative returns shown are the log returns of each specific model when trading EURUSD, overly randomly sampled 400 day periods. The performance is erratic and scattered across winning and losing strategies . It was due to this under performance and lack of robustness which led to the need to look at alternative avenues when training the model.

Figure 12: Erratic Returns of Machine Learning Based Trading Models



Due to this lack of stability in the each models performance, a different methodology was needed and instead each model was now trained on randomly generated data as outlined in ?? Using randomly generated price series allows us to calibrate the strength of trend (and the significance of the noise) when training each model on how to recognise trending price series. The goal here is to allow the model to recognise the dynamics and patterns of a trending price series, and thus removes the risk that the training period of the live price series lacks any sustainable trend.

<sup>17</sup>Research suggested a 26% fall in performance out-of-sample and 58% decline post publication

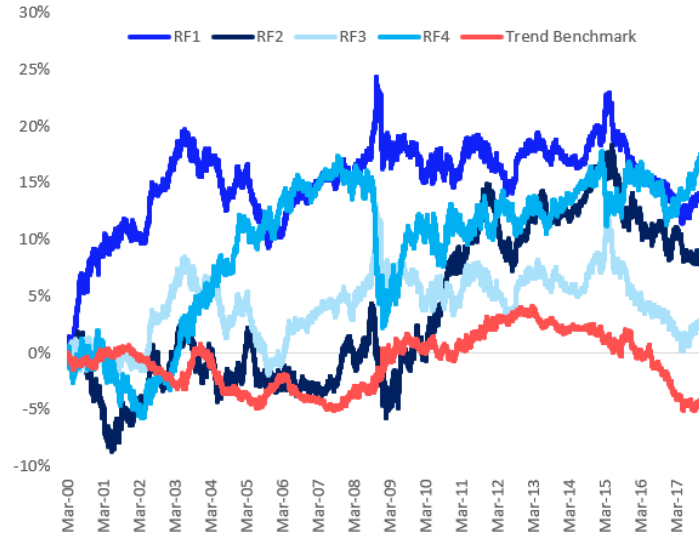
## 6.1 Performance of Model One: Trend Capture

This section outlines how various machine learning models performed in trying to predict EU-RUSD direction having been trained only on randomly generated trending price series. Each of the models here were trained as classifiers such that precision and recall could be calculated but also due to the fact that the returns tended to be more robust. This is possibly due to the fact that the regression models tended to have been skewed by large values that can occur in financial time series, which meant that when a classifier of the move (up or down) was used, this helped to dampen the volatility of the signal and removed the impact of outliers on the models ability to learn the patterns of the model.

### 6.1.1 Decision Trees: Random Forest

The first model tested after training on the randomly generated trend was the random forest machine learning model. Figure 13 shows the cumulative log return of trading EURUSD across various iterations of the models trained on different randomly generated datasets. The models below used 200 trees and a maximum of four features when using selecting node splits. The testing period here is the full dataset which spans from the start of 2002 to the end of 2017. This testing period covers various different market regimes, such as the Great Financial Crisis, unorthodox central bank policies from 2010 onwards and the end of the commodity super cycle. In terms of results, the returns are still somewhat scattered, however the performance

Figure 13: Cumulative Returns: Random Forest Based Trading Model



on average is positive. The red line in the chart denotes the performance of a simple linear benchmark strategy of trading in the direction of the price trend over the last 24 hours, we use this 24 hour period as it matches the trade horizon that the model used as its target variable. In effect, saying that a simple linear strategy predicts what will happen in the next 24 hours by looking at what occurred in the previous 24 hours. In addition to the cumulative returns of the strategy, we also shown a break down of each models performance statistics in figure 14,

Figure 14: Performance Statistics: Random Forest Based Trading Model

	Benchmark	RF1	RF2	RF3	RF4
Annual Return	-0.21%	0.74%	0.44%	0.11%	0.91%
Annual Volatility	1.75%	3.21%	3.83%	3.15%	3.77%
Sharpe	-0.12	0.23	0.12	0.03	0.24
Max	0.56%	0.74%	0.81%	0.73%	0.75%
Min	-0.74%	-0.80%	-0.94%	-0.80%	-0.84%
Sample Size (N)	111654	111654	111654	111654	111558
True Positive	22577	23460	28088	30004	29200
True Negatives	23172	22910	27789	30023	29027
False Positives	21327	22722	20541	18464	19390
False Negatives	22375	22938	20532	18257	19330
Accuracy (Hit Ratio)	39%	41%	44%	43%	44%
Precision	49%	51%	50%	50%	50%
Recall	50%	51%	58%	62%	60%

which are broken down to provide the annualised return and volatility, the sharpe ( risk taken to reward received), the ratio of largest gain to largest loss and finally the precision and recall of each model.

### 6.1.2 Support Vector Classifiers

The performance of the support vector classifier was mixed, with some models tending to under perform the simple linear benchmark. One particular model (labelled SV2) in figure 15, performed very well up until the fina few years of the sample, perhaps as this coincided with a period of mean reverting behaviour in EURUSD. Overall , the SVM model performed somewhat disappointingly over the full sample.

The performance statistics shown in 16 show the under performance of the models in general, as they show low or negative Sharpe ratios alongside underwhelming max gain to max loss ratios. Interestingly, the precision and recall scores are the same if not partly better than same scores for the random forest model, and this showcases the fact that classification accuracies do not tell us the full story when analysing financial timeseries. While the SVM models teneded to have relatively impressive precision and recall scores, this did not translate into strong trading signal performance as the model likely misclassified the higher volatility time periods, leading to larger average losses than equivalent average gains.



Figure 15: Cumulative Returns: Support Vector Classifier Based Trading Model

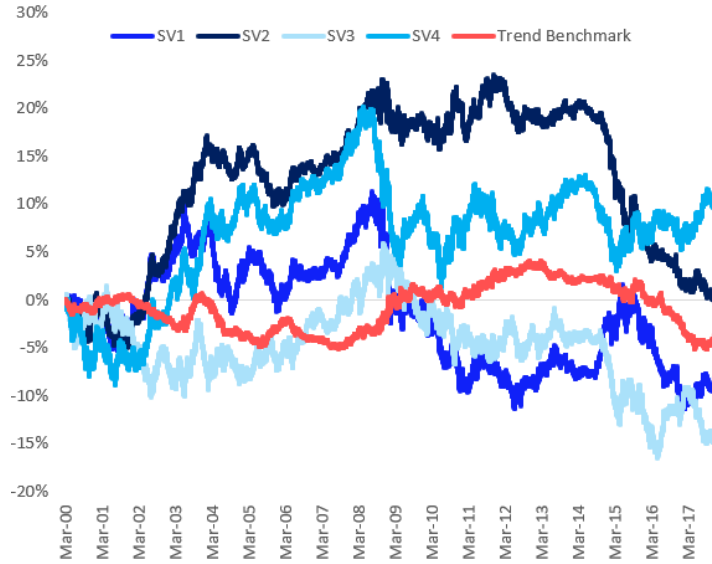


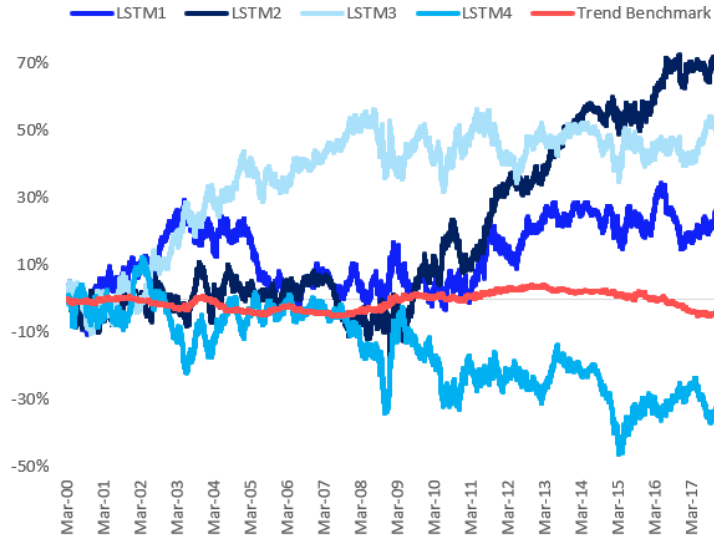
Figure 16: Performance Statistics: Support Vector Classifier Based Trading Model

	Benchmark	SVM1	SVM2	SVM3	SVM4
Annual Return	-0.21%	-0.45%	0.05%	-0.85%	0.63%
Annual Volatility	1.75%	4.12%	3.53%	4.15%	4.48%
Sharpe	-0.12	-0.11	0.01	-0.20	0.14
Max	0.56%	0.73%	0.64%	0.86%	0.81%
Min	-0.74%	-0.84%	-0.84%	-1.04%	-0.86%
Sample Size (N)	111654	111654	111654	111558	111654
True Positive	22577	33009	33725	11257	45036
True Negatives	23172	32901	33135	11139	44578
False Positives	21327	15528	14784	20850	3655
False Negatives	22375	15517	15265	21088	3686
Accuracy (Hit Ratio)	39%	43%	43%	29%	44%
Precision	49%	50%	50%	50%	50%
Recall	50%	68%	69%	35%	92%

### 6.1.3 Long Short Term Memory Trading Model

The recurrent neural network such as the LSTM provided probably the most impressive results, with the majority of models outperforming the linear benchmark 18. All of the LSTM models here used an initial layer of 32 nodes , with an additional layer containing 16 nodes and using the softmax activation. The length of the lookback sequence here was 90 periods.

Figure 17: Cumulative Returns: Long Short Term Memory Based Trading Model



In terms of the statistics, we do see on average positive returns from the LSTM models and they do outperform the linear benchmark in terms of sharpe ratio as well. In fact model 3 and 4 in the sample have positive max gain to max loss ratios, suggesting they may exhibit positive skewness to their returns., Overall the LSTM models looked most promising in terms of capturing the trend in the price series.

The outperformance of the LSTM model, does suggest that it is important for the model to understand the historical sequence of the data points such that the patterns of trend can sustain long enough for the model to recognise the sequence of points which form a part of a trending price series. This result is also seen in citeLim2019 where their LSTM model outperformed other machine learning models. In citeLIM2019, they used 100 epochs for training , while in this paper 500 epochs were used.

## 6.2 Performance of Model Two: Regime Estimation

The section looks at the performance of the trend estimation model, where we feed in the performance of each trend model ( ) and use fundamental features, such as economic growth differentials, historical volatility and the state of global growth to identify periods which can be conducive to the positive performance of currencies. The approach taken in this stage of the trading process was split the dataset into one batch training which consisted of the period from 2002 until 2010 and a test period form 2010 until the end of 2017. The model chosen for trend estimation was a simple random forest decision tree as this type of model does not

Figure 18: Performance Statistics: Long Short Term Memory Based Trading Model

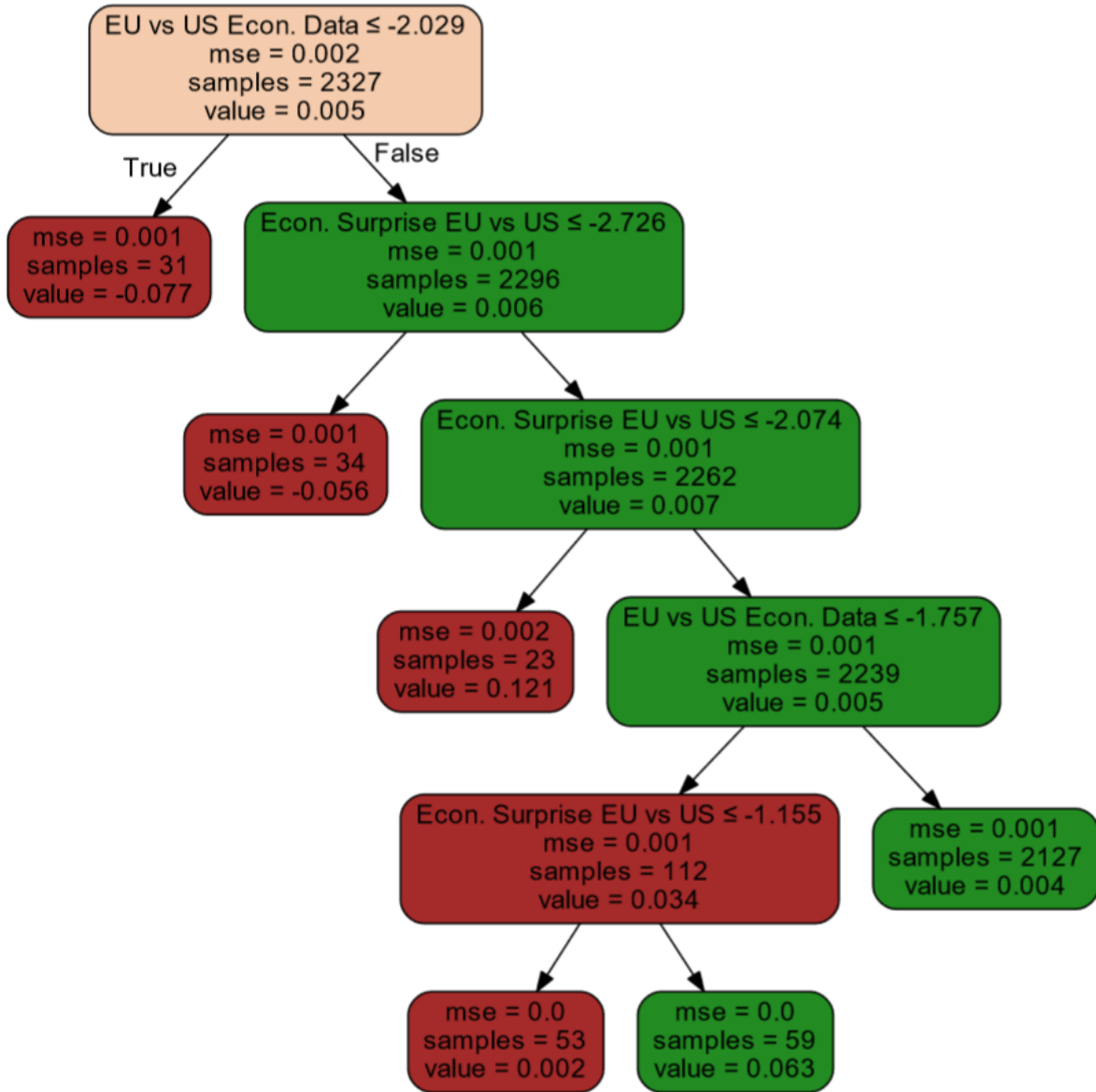
	Benchmark	LSTM1	LSTM2	LSTM3	LSTM4
<b>Annual Return</b>	-0.21%	1.56%	3.78%	3.00%	-1.99%
<b>Annual Volatility</b>	1.75%	10.15%	10.15%	10.15%	10.14%
<b>Sharpe</b>	-0.12	0.15	0.37	0.30	-0.20
<b>Max</b>	0.56%	1.81%	2.03%	2.25%	2.25%
<b>Min</b>	-0.74%	-2.25%	-2.25%	-2.03%	-1.47%
<b>Sample Size (N)</b>	111654	111654	111654	111654	111654
<b>True Positive</b>	22577	30262	20741	45837	17351
<b>True Negatives</b>	23172	29702	20147	45169	17387
<b>False Positives</b>	21327	18497	28052	3020	30744
<b>False Negatives</b>	22375	18490	28011	2896	31332
<b>Accuracy (Hit Ratio)</b>	39%	44%	44%	44%	43%
<b>Precision</b>	49%	50%	51%	50%	50%
<b>Recall</b>	50%	62%	43%	94%	36%

require a large training dataset and also as the model itself can allow us to build simple charts showing the decision making process of the model itself. A simplified tree structure is shown for one of the trend models in 19. In this instance, only three features were used, global economic data momentum, the economic performance differential between the United States and Euro Area and the Economic surprise differential between the US and Euro Area <sup>18</sup>. The graphic of the decision tree structure shows that when the economic performance differential between Euro Area and the US has been below -2 standard deviations, then the trend model in this case under performed. Hence we can see that the goal of the trend estimation model is to align the fundamental factors with periods then the trend capture model has performed well, allowing the second trend estimation model to try to switch on and off the trend model when the regime is not conducive to the performance of trend (based on those historical patterns).

---

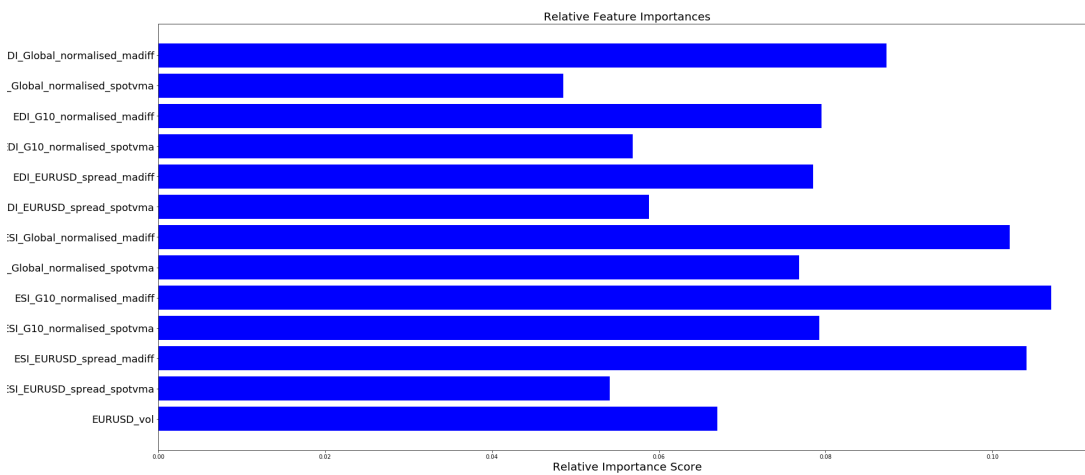
<sup>18</sup>An economic surprise measures the difference between the actual economic release value and the economists consensus expectation

Figure 19: Decision Tree Structure Example of the Trend Estimation Model



When running the trend estimation model, given the number of possible features across various different lookback horizons and the number of feature combinations, a variable importance plot was created to try to understand which features may have the most explanatory power, a variable importance plot is created by shuffling the features and running the model each time without one feature to try to uncover how much value that feature adds to the overall prediction ability of the model, this feature selection study occurred in the training period only. In this case, one trains on daily data only and hence we look to reduce the models dimensionality by selecting a few sub components of the most interesting features. The variable importance chart is shown below in figure 20.

Figure 20: Performance Model 2 Trend Estimation Across Each ML Model Trading Model



We can see that certain features do trend to outperform other features in the level of their predictability. This this allows us to reduce the dimensionality of the model as some of the features are likely redundant. The trend estimation model worked on a daily basis and predicted the performance the trend model one month ahead. Figure 21, shows the performance of a selection of trend models from each of the three models which have been fed into the trend estimation model and trained on how each individual model has performed given the economic, macro risk and volatility features shown above.

Figure 21, shows the performance statistics of the models once having been filtered by the trend. The testing period ran from December 2010 until December 2017. The red lines denote the average cumulative performance of all the different trend models created using a particular machine learning model. The blue lines denote the average performance across the models tested of the same strategies when they have been filtered by the trend estimation model. The trend estimation model would reduce the exposure of the trend model to zero if it had foretasted a loss in the strategy over the next one month.

The performance statistics on average across all three models once having been filtered by the trend estimation model are shown in 22, on average across all the models, applying

Figure 21: Performance Model 2 Trend Estimation Across Each ML Model Trading Model



trend estimation either increased the average return or reduced a trading strategies' maximum drawdowns (a trading models high to low cumulative losses).

Overall, the two models do tend to work well together and there does seem to be some value in diversifying the decision making process across different models with different mandates.

### 6.3 Performance on Blind Test Data

In this section, we look at the performance of the finalised models and model architectures on a fully unseen set of EURUSD data, this dataset was only used in the last stage of the project to provide an indication of how the model may have performed if it had been created at the end of 2017. We show the performance of the trend model and the overlay of the trend estimation model. The linear benchmark is shown in red (once again, this is the performance of trading in the direction of the previous 24 hour trend of the currency). One clear takeaway is that the machine learning models were much less volatile than the benchmark, and two of the machine learning models outperformed the linear benchmark over the full period of the test. While there is not enough data to draw statistical conclusions from the blind test results, there is evidence that the models and proposed model architecture proposed in this paper can indeed be used in the domain of currency trading. The fact that the trend estimation model was able to identity the fact of a lack of sustained trend allowed the support vector machine trade trending price behaviour only when the probability of trend working was high.

The model statistics are shown in ??, the accuracy hit rates of the trend models after applying the model two trend estimation is positive for all models. Each machine learning

Figure 22: Performance Statistics Model 2 Trend Estimation Average Performance Per Trading Model

	<b>LSTM</b>	<b>RF</b>	<b>SVM</b>
<b>Annual Return</b>	1.25%	0.56%	-0.72%
<b>Annual Volatility</b>	6.80%	2.31%	1.90%
<b>Average Sharpe</b>	0.19	0.19	-0.37
<b>Max</b>	2.8%	0.8%	1.0%
<b>Min</b>	-3.4%	-1.0%	-1.0%
<b>Accuracy (Hit Ratio)</b>	24.1%	28.5%	19.1%
<b>Max Gain/Loss Ratio</b>	0.12	0.04	0.52
<b>Drawdown</b>	-19.3%	-5.9%	-6.8%
<b>Calmar Ratio</b>	-7.4%	-10.4%	10.5%
<b>DSR Test Statistic</b>	47%	30%	3%

model also outperforms its linear benchmark and has lower standard deviation of returns and higher average annualised returns.

Figure 23: Trading System Performance Fully Unseen Live Data

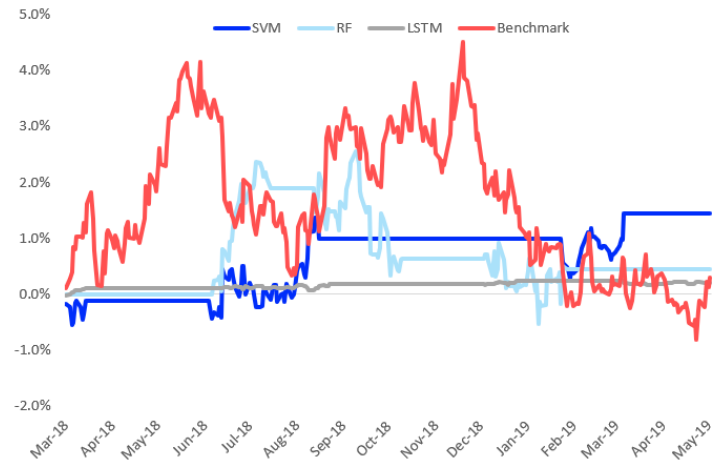


Figure 24: Trading System Performance Fully Unseen Live Data

	Benchmark	SVM	RF	LSTM
Average Return	0.16%	0.97%	0.38%	0.15%
Average Volatility	4.79%	1.82%	2.70%	0.16%
Sharpe	0.03	0.53	0.14	0.89
Max	0.90%	0.86%	0.76%	0.04%
Min	-1.11%	-0.49%	-0.83%	-0.04%
Accuracy (Hit Ratio)	47%	52%	53%	58%
Max Gain/Loss Ratio	0.80	1.75	0.92	1.04
Drawdown	-5.32%	-1.10%	-3.09%	-0.09%
Calmar Ratio	0.0	-0.9	-0.1	-1.6



## 7 Conclusions

Allow the thesis to bite itself in the tail, blind testing results, the possibility that this model actually works and how can you train/test better? Look at your self in the mirror, what have you learned and what are the takeaways for future work.

In this section , we provide some discussion on the results and the possible implications of the work shown here. How did it perform? what worked best, what are some of the interesting takeaways and further work that needs to be explored here?

additional commentary - to be removed/ or included. Objective 3: Model Training/Fine Tuning and Architecture Selection The model training stage allows the learning algorithm to understand and capture the relationship between the features and the classification target vector (the target vector will be the future direction of the asset price across various time periods, both regression and classification models will be tested). Prior research has generally used binary classification of the future asset price move when developing models( [55], [23], [37]).

Model training will also allow us to tune parameters to help understand which combination of models (ANN, RNN, PNN ,SVM, RF) and underlying architecture (hidden layers, number of decision trees) are best suited to capturing the nuances in the data. This stage of the process will also allow experimentation with various models to see the differences in performance. Various Python packages such as Sci-Kit Learn and Keras were used in order to create and test the classification success rates of each type of model. Fischer et al. made use of Keras (on top of Google’s Tensorflow) to implement Recurrent Neural Networks (RNN) [47] while using Sci-Kit Learn for logistic regression, SVM and Random Forest methods.

One issue encountered was that the length of time needed to train the LSTM meant that training multiple variations across parameters was challenging from a time perspective. Three model architectures were used, LSTM, Random Forests and Support Vector Machines. One novel approach to use in selecting the target vector was to use a risk adjusted return and predict the risk adjusted return as opposed to the binary return, this means that the trend model can also estimate the positioning size and negates the need for a separate trade risk adjustment to be made. This was shown to provide interesting results in a paper on time series momentum from B Lim ADD!

The model was run on a cloud server however it was seriously slow!!!! total LSTM training would have taken over 50 hours!

Objective 4: Model Validation In order to validate the model, the selected model and its structure will need to be tested in a similar setting to real world trading. This requires the construction of a backtesting engine which can simulate each models performance taking into account transaction costs and the funding costs associated with trading currencies. This will provide a much more accurate analysis of the selected trading strategy performance<sup>19</sup>. While classification success rates will be important in gauging model accuracy, high classification rates may not lead to positive strategy performance if large currency moves are mis-classified. Go over the full cycle of how you trained and tested the model please...

Performance Evaluation Analysing and outlining model performance will be one of the key objectives to be satisfied in the final stages of the project. In light of the findings from ( [8], [6]), which highlight in detail the pitfalls inherent in systematic trading, the below outlines an approach to reviewing and evaluating model performance to standards generally set in financial institutions.

Results Across All Model Iterations Keeping track of the number of models being tested is very important as it can help define the statistical likelihood of finding a strategy that is a true positive. Even given 20 randomly selected strategies, there is a non-negligible probability that one strategy will beat the threshold needed for a statistical significance [8]. Thus, the number of trials run must be considered so we can deflate the statistical significance of the strategy performance as a function of the number of trials [59]. Cross Validation Splitting the dataset into various training and test sets will be very important in being able to assess model performance. The most common cross validation techniques are K-fold cross validation

---

<sup>19</sup>The vast majority of research papers referenced in this thesis did not perform this type of testing

and Leave-One-Out cross validation. This project will use a training set, validation set and finally a test dataset which will be held out until the final model has been selected. This will provide a useful gauge of out of sample performance<sup>20</sup>. Problem of Dimensionality In finance, the lack of sufficient data (price) vs the number of possible predictors does lead to the issue of dimensionality, where we can add further features but not necessarily add further data points to be tested. Thus the model proposed here will remain cognizant of the number of features being used and will try to reduce the probability of finding false positive trading strategies. The dimensionality of the current data set will involve one target vector (i.e. future price move) and 17-20 predictor columns as outlined in 1. Out-of-Sample Testing on Unseen Data As stated by Arnott et al [8], selection of model inputs is inherently biased (especially when done by domain experts) as the experts use heuristics to choose what they find the most economically intuitive as predictors, thus even out of sample test results are not fully out of sample. To override some of these concerns, this paper proposes to omit one years worth of data from any analysis until the model and related architecture has been fully completed and the selected model or model architecture has been finalised. Then the chosen model will be blind tested on the final year of data that has been held back to give the best approximation of how the chosen model might have performed.

## References

- [1] S. Bookman, “15 Supercomputers Less Powerful Than Your Phone.” <https://www.theclever.com/15-huge-supercomputers-that-were-less-powerful-than-your-smartphone/>, 2015. [Online; accessed 02-April-2019].
- [2] C. A. MacK, “Fifty years of Moore’s law,” in *IEEE Transactions on Semiconductor Manufacturing*, 2011.
- [3] R. Rebonato, “Asset Management: A Systematic Approach to Factor Investing,” *Quantitative Finance*, 2017.
- [4] M. Lopez de Prado, “The 10 Reasons Most Machine Learning Funds Fail,” *SSRN*, 2018.
- [5] B. Lebaron, “Chaos and Nonlinear Forecastability in Economics and Finance,” *Philosophical Transactions of the Royal Society A: Mathematical, Physical and Engineering Sciences*, 1994.
- [6] D. H. Bailey, J. Borwein, M. Lopez de Prado, and Q. J. Zhu, “Pseudo-Mathematics and Financial Charlatanism: The Effects of Backtest Overfitting on Out-of-Sample Performance,” *SSRN*, 2013.
- [7] S. M. Bartram, L. Djuranovik, and A. Garratt, “Currency Anomalies,” *SSRN Electronic Journal*, 2019.
- [8] R. D. Arnott, C. R. Harvey, and H. Markowitz, “A Backtesting Protocol in the Era of Machine Learning,” *SSRN Electronic Journal*, 2018.
- [9] S. Camargo, S. M. Queirós, and C. Anteneodo, “Bridging stylized facts in finance and data non-stationarities,” *European Physical Journal B*, 2013.

---

<sup>20</sup>In depth study of cross validation techniques outlined in [60]

- [10] N. Kaissar, "Hedge Funds have a Performance Problem." <https://www.bloomberg.com/opinion/articles/2016-03-24/hedge-funds-have-a-performance-problem>, 2016. [Online; accessed 04-April-2019].
- [11] A. Corhay and A. T. Rad, "Conditional heteroskedasticity adjusted market model and an event study," *Quarterly Review of Economics and Finance*, 1996.
- [12] M. Batwick, "How missing out on 25 days in the stock market over 45 years costs you dearly." <https://www.marketwatch.com/story/how-missing-out-on-25-days-in-the-stock-market-over-45-years-costs-you-dearly-2016-01-25>, 2016. [Online; accessed 31-March-2019].
- [13] N. N. Taleb, "Black Swans and the domains of statistics," *American Statistician*, 2007.
- [14] W. Huang, Y. Nakamori, and S. Y. Wang, "Forecasting stock market movement direction with support vector machine," *Computers and Operations Research*, 2005.
- [15] S. Shen, H. Jiang, and T. Zhang, "Stock market forecasting using machine learning algorithms," 2012.
- [16] S. Wang and W. Shang, "Forecasting direction of China security index 300 movement with least squares support vector machine," in *Procedia Computer Science*, 2014.
- [17] B. Widrow and M. A. Lehr, "30 Years of Adaptive Neural Networks: Perceptron, Madaline, and Backpropagation," *Proceedings of the IEEE*, 1990.
- [18] F. Rosenblatt, "The perceptron: A probabilistic model for information storage and organization in the brain," *Psychological Review*, 1958.
- [19] V. Kůrková, "Kolmogorov's theorem and multilayer neural networks," *Neural Networks*, 1992.
- [20] S. Hochreiter and J. Schmidhuber, "Long Short-Term Memory," *Neural Computation*, 1997.
- [21] P. Czekalski, M. Niezabitowski, and R. Styblinski, "ANN for FOREX forecasting and trading," in *Proceedings - 2015 20th International Conference on Control Systems and Computer Science, CSCS 2015*, 2015.
- [22] K. Joarder and R. A. Sarker, "Forecasting of currency exchange rates using ANN: A case study," in *Proceedings of 2003 International Conference on Neural Networks and Signal Processing, ICNNSP'03*, 2003.
- [23] H. Gunduz, Y. Yaslan, and Z. Cataltepe, "Intraday prediction of Borsa Istanbul using convolutional neural networks and feature correlations," *Knowledge-Based Systems*, vol. 137, pp. 138–148, 2017.
- [24] K.-j. Kim, "Predicting Stock Price Direction using Support Vector Machines," *Neurocomputing*, 2003.
- [25] M. Kumar and M. Thenmozhi, "Forecasting Stock Index Movement: A Comparison of Support Vector Machines and Random Forest." 2006.

- [26] 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 Applications*, 2000.
- [27] C. Cortes and V. Vapnik, "Support-Vector Networks," *Machine Learning*, 1995.
- [28] D. Wilimitis, "The Kernel Trick in Support Vector Classification." <https://towardsdatascience.com/the-kernel-trick-c98cdbcaeb3f>, 2018. [Online; accessed 07-April-2019].
- [29] G. James, D. Witten, T. Hastie, and R. Tibshirani, *An Introduction to Statistical Learning*. 2013.
- [30] Y. Kara, M. Acar Boyacioglu, and Ö. K. Baykan, "Predicting direction of stock price index movement using artificial neural networks and support vector machines: The sample of the Istanbul Stock Exchange," *Expert Systems with Applications*, 2011.
- [31] 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*, 2015.
- [32] B. Gui, X. Wei, Q. Shen, J. Qi, and L. Guo, "Financial time series forecasting using support vector machine," in *Proceedings - 2014 10th International Conference on Computational Intelligence and Security, CIS 2014*, 2015.
- [33] J. Kamruzzaman, R. Sarker, and I. Ahmad, "SVM based models for predicting foreign currency exchange rates," 2004.
- [34] V. Podgorelec and M. Zorman, "Decision Tree Learning," in *Encyclopedia of Complexity and Systems Science*, 2015.
- [35] L. Breiman, "Random Forrests," *Machine learning*, 2001.
- [36] R. Genuer, "Variance reduction in purely random forests," *Journal of Nonparametric Statistics*, 2012.
- [37] S. P. Chatzis, V. Siakoulis, A. Petropoulos, E. Stavroulakis, and N. Vlachogiannakis, "Forecasting stock market crisis events using deep and statistical machine learning techniques," *Expert Systems with Applications*, 2018.
- [38] W. McKinney, *Python for Data Analysis: Data Wrangling with Pandas, NumPy, and IPython*. 1976.
- [39] S. Van Der Walt, S. C. Colbert, and G. Varoquaux, "The NumPy array: A structure for efficient numerical computation," *Computing in Science and Engineering*, 2011.
- [40] D. R. Tobergte and S. Curtis, "Learning SciPy for Numerical and Scientific Computing Second Edition," *Journal of Chemical Information and Modeling*, 2013.
- [41] A. Géron, *Hands-On Machine Learning with Scikit-Learn and TensorFlow*. 2017.
- [42] M. A. Wood, *Python and Matplotlib Essentials for Scientists and Engineers*. 2015.

- [43] M. I. Reiff, “Pandas,” *Journal of Developmental and Behavioral Pediatrics*, 2002.
- [44] F. Chollet, “Keras: The Python Deep Learning library,” 2015.
- [45] I. Filippou and M. P. Taylor, “Common Macro Factors and Currency Premia,” *Journal of Financial and Quantitative Analysis*, 2017.
- [46] M. Dahlquist and H. Hasseltoft, “Economic Momentum and Currency Returns.” 2015.
- [47] T. Fischer and C. Krauss, “Deep learning with long short-term memory networks for financial market predictions,” *European Journal of Operational Research*, 2018.
- [48] B. Yoann, “The Truth Nobody Wants to Tell You About AI for Trading.” <https://hackernoon.com/https-medium-com-supernova-the-truth-nobody-wants-to-tell-you-about-ai-for-trading-5d29a297ee93>, 2019. [Online; accessed 31-March-2019].
- [49] T. J. Moskowitz, Y. H. Ooi, and L. H. Pedersen, “Time series momentum,” *Journal of Financial Economics*, 2012.
- [50] L. Menkhoff, L. Sarno, M. Schmeling, and A. Schrimpf, “Currency momentum strategies,” *Journal of Financial Economics*, 2012.
- [51] K. Rogoff, “The Purchasing Power Parity Puzzle,” *Journal of Economic Literature*, 1996.
- [52] F. Hassan, “The price of development: The Penn–Balassa–Samuelson effect revisited,” *Journal of International Economics*, 2016.
- [53] M. Ca’ Zorzi, J. Muck, and M. Rubaszek, “Real Exchange Rate Forecasting and PPP: This Time the Random Walk Loses,” *Open Economies Review*, 2016.
- [54] B. Döme, “Efficiency and Return Predictability.” 2008.
- [55] G. Abreu, R. Neves, and N. Horta, “Currency exchange prediction using machine learning, genetic algorithms and technical analysis,” 2018.
- [56] W. F. Sharpe, “The Sharpe Ratio,” *The Journal of Portfolio Management*, 2009.
- [57] D. H. Bailey and M. Lopez de Prado, “The Deflated Sharpe Ratio: Correcting for Selection Bias, Backtest Overfitting and Non-Normality,” *SSRN*, 2014.
- [58] R. D. Mclean and J. Pontiff, “Does Academic Research Destroy Stock Return Predictability?,” *Journal of Finance*, 2016.
- [59] D. H. Bailey and M. Lopez de Prado, “The Sharpe Ratio Efficient Frontier,” *SSRN*, 2011.
- [60] C. Bergmeir and J. M. Benítez, “On the use of cross-validation for time series predictor evaluation,” *Information Sciences*, 2012.