

# A comparison between ANNs and traditional stock market forecasting techniques



Alexandros Antoniou

Centre for Computational Finance and Economic Agents

CSEE

University of Essex

Submitted in partial satisfaction of the requirements for the

Degree of Master of Science

in

Computational Finance

*Supervisor* Dr Maria Kyropoulou  
*Second Supervisor* Dr John O'Hara

August 2020

# Acknowledgements

I would like to thank my supervisor, Dr Maria Kyropoulou, for helping me with this thesis.

I would also like to thank my friends and family for keeping me sane during this lock-down and overall worldwide panic. It has not been a fun time.

# Abstract

TODO

---

*Keywords: Deep Learning, Artificial Neural Networks, Stock Market, Time Series prediction, Neural Networks*

# Table of Contents

<b>1</b>	<b>Introduction</b>	<b>1</b>
1.1	Problem statement . . . . .	1
1.2	Outline of Paper . . . . .	2
<b>2</b>	<b>Overview of Neural Networks</b>	<b>3</b>
2.1	Neurons . . . . .	3
2.2	Activation Functions . . . . .	4
2.2.1	Rectified Linear Unit . . . . .	4
2.2.2	Sigmoid . . . . .	5
2.2.3	Hyperbolic tangent . . . . .	5
2.3	Backpropagation . . . . .	6
2.4	LSTM . . . . .	7
<b>3</b>	<b>Autoregressive models</b>	<b>8</b>
<b>4</b>	<b>Methodology</b>	<b>9</b>
4.1	Time-Series Forecasting . . . . .	9
4.2	ARIMA . . . . .	9
4.3	Model Description . . . . .	11
4.4	Model Evaluation . . . . .	13
<b>5</b>	<b>Results</b>	<b>15</b>
5.1	Dataset . . . . .	15
5.2	Evaluation Metrics . . . . .	15
5.3	Results . . . . .	15
<b>6</b>	<b>Conclusions</b>	<b>16</b>

# 1 Introduction

The ability to predict changes in stock prices is extremely important to the financial world as it influences trading strategies and reduces risks in the market. Forecasting has long been a problem for the business and technology communities and has seen little advances until quite recently, with the advent of neural networks and deep learning.

Before artificial neural networks, the finance world used other methods to model the time series that arose from the continuous updating of stock prices. Models like the autoregressive integrated moving average model (ARIMA) and the generalised autoregressive conditional heteroscedasticity model (GARCH) have become key econometric methods for forecasting time series and are still widely used in finance.

The focus of this project is to provide a comparison between the traditional methods for time series forecasting mainly the ARIMA model, and simple implementations of artificial neural networks, in the context of financial time series prediction.

Forecasting stock prices with moving averages belongs to the technical analysis category of financial analysis. Kirkpatrick and Dahlquist (2006) define *technical analysis* as the study of prices in freely traded markets with the intent of making profitable trading or investment decisions,<sup>[16]</sup> hence the models will only use daily prices for the selected stocks, ignoring company financial data for both forecasting methods, ANNs and autoregressive models.

## 1.1 Problem statement

The Efficient Market Hypothesis states that stock prices reflect all available information and respond to any changes in that information through price changes. An implication of this relation between available information and the market is that consistently predicting a change in an asset price is hard. Followers of this hypothesis believe that the market instantly responds to new information made available and so technical analysis of assets is moot. The Efficient Market Hypothesis finds its roots in the works of Bachelier (1900), Mandelbrot (1963) and Samuelson (1965) but more closely associated to Fama whose published review paper, *‘Efficient Capital Markets: A Review of Theory and Empirical Work’*, pro-

posed the market efficiency types and the *joint hypothesis problem* that introduces certain requirements for testing market efficiency.<sup>[2][19][24][8]</sup> Of course, the EMH being near or even completely untestable means that a lot of research has been done to provide insight to the problem of predicting the stock market, both for the hypothesis and against. Grossman and Stiglitz (1980) propose a model that suggests that the market is not perfectly efficient as information is costly and in an efficient market where asset prices fully reflect all available information, there would be no incentive for investors to acquire information, decreasing the efficiency of the market. Both Grossman and Stiglitz and Black agree on the requirement of investors with different levels of information available to them, suggesting that it makes markets inefficient.<sup>[12][3]</sup>

The *random walk hypothesis* is the financial concept that asset price fluctuations can be modelled as random walks, making them unpredictable. A random walk is a stochastic or random process that involves random changes of a value, in the context of finance it's most often described as the change in price, positive or negative, by a set step point.<sup>[6][7][4]</sup> If the *random walk hypothesis* and the closely related *efficient market hypothesis* hold true, it should be extremely hard to consistently model and attempt to forecast asset prices. The following sections however, will go through modern and traditional methods of stock market modelling and forecasting.

## 1.2 Outline of Paper

Sections 2 and 3 will provide a comprehensive review of relevant literature for both neural networks and autoregressive integrated moving average models. Section 4 begins with a brief description of the preprocessing performed to the data, followed by descriptions of the ARIMA models and the LSTM model, including optimisation techniques. Section 4 ends with a description of the evaluation metrics chosen for the comparison. Section 5 will present the results of the project and a comparison between the two models, showing the LSTM models outperform the traditional ARIMA models. Section 6 will offer discussion on the advantages and limitations of neural networks compared to traditional numerical methods for time-series forecasting. The appendix contains the totality of graphs and values extracted from the training and testing of the models.

## 2 Overview of Neural Networks

In this section, we provide a brief history of neural networks in the context of finance and time series prediction, as well as a more detailed description of the architecture used in the paper.

Neural networks are systems largely belonging to the study of Machine Learning, typically associated with solving computational problems that other models and algorithms struggle to. Loosely based on biological neural networks, like the brain, artificial neural networks are collections of neuron-like nodes and links that connect them. Their resemblance to the brain in terms of architecture is part of what allowed neural networks to grow in computing power and popularity over the past decades, from their emergence in the 1940s.<sup>[17][20]</sup>

### 2.1 Neurons

Neurons are the artificial equivalent of a biological neuron and are the fundamental components of a neural network. Neurons perform three tasks, receive input from other neurons, apply an activation function to the input and output a value to other neurons. Neurons in the network are connected and each of these connections carries a signal and has certain weight attached to it, which affects the signal carried.<sup>[26][23]</sup> Graphically, a neuron could be represented in figure 2.1.

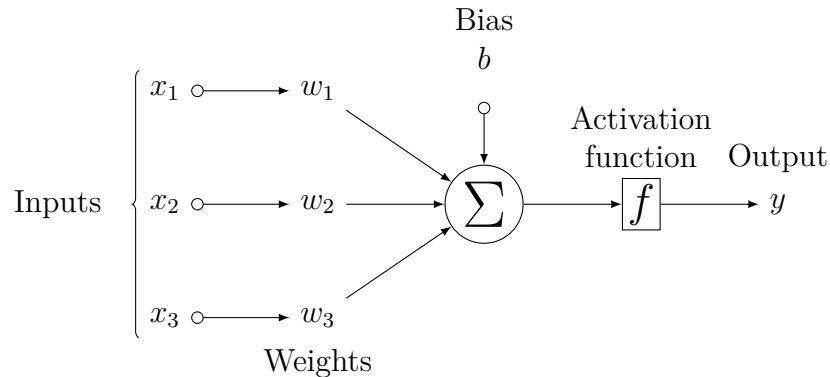


Figure 2.1: Graphical representation of a Neuron

## 2.2 Activation Functions

Wilson (2009) defines the activation function of a neural network as the function that governs the output behaviour of a neuron, given a set of input values. There are several types of activation functions, the simplest being the step function.

In mathematical terms, the step function, or the unit step function, is defined as

$$H(x) := \begin{cases} 0, & \text{for } x < 0 \\ 1, & \text{for } x \geq 0 \end{cases} \quad (2.1)$$

where  $H(x)$  is the Heaviside step function<sup>[1]</sup>. At value 0, an output of  $H(0) = 1$  is selected and passed down to the next layer of neurons.

### 2.2.1 Rectified Linear Unit

The Rectified Linear Unit (ReLU) is an activation function mathematically defined as

$$f(x) = [x]^+ = \max(0, x) \quad (2.2)$$

where  $x$  is the neuron's input value. Plotted on cartesian axes, the graph shows a linear relationship for  $f(x)$  and  $x$  where  $x > 0$ . Generating a plot of ReLU further explains its activation range.

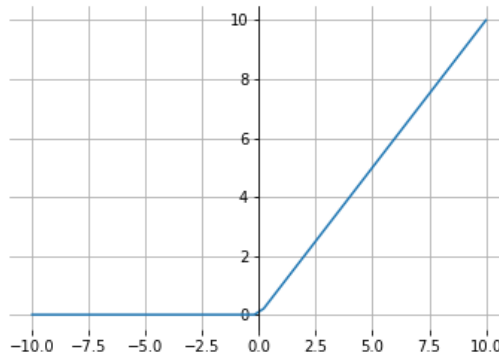


Figure 2.2: Plot of the Rectified Linear Unit function



ReLU was first introduced in the 2000 and 2001 Hahnloser papers.<sup>[13][14]</sup> ReLU is currently the most widely used activation function<sup>[21]</sup> and has found great success in training deep neural networks primarily in the fields of computer vision<sup>[10]</sup> and speech recognition.<sup>[18]</sup>

### 2.2.2 Sigmoid

The Sigmoid function is a logistic activation function, mathematically defined as

$$f(x) = \frac{L}{1 + e^{-k(x-x_0)}} \quad (2.3)$$

where

$L$  is the maximum point of the curve,

$k$  is the steepness of the curve,

$x_0$  is the value of the midpoint of the curve.

Sigmoid functions are non-linear, differentiable, defined for all real input values and have an output greater than zero for all inputs, in contrast to the ReLU activation function whose output is greater than zero only with positive inputs. A well known example of a sigmoidal curve is the cumulative distribution function of the normal distribution.

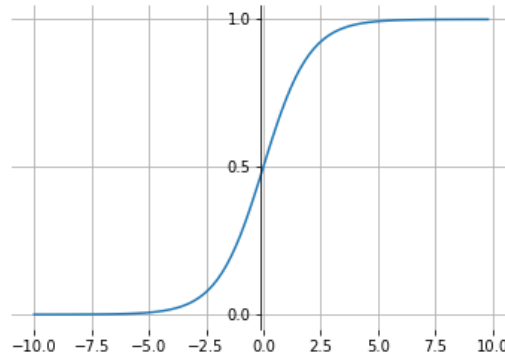


Figure 2.3: Plot of the sigmoid function

### 2.2.3 Hyperbolic tangent

The hyperbolic tangent is a specialised case of the sigmoid function, mathematically defined as

$$f(x) = \tanh(x) = \frac{e^{2x} - 1}{e^{2x} + 1} \quad (2.4)$$

Unlike the standard sigmoid function,  $\tanh$  can output negative values for negative inputs giving it double the range that the standard sigmoid has.

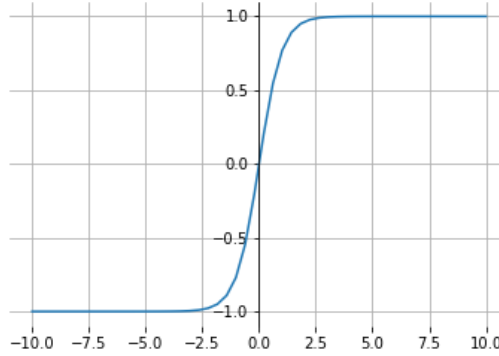


Figure 2.4: Plot of the hyperbolic tangent function

## 2.3 Backpropagation

The backpropagation algorithm is a key feature of neural networks as it allowed for fast and efficient training of multi-layered networks by distributing error values back through the layers of the network, adjusting the weights in the connections between neurons respectively.<sup>[25]</sup> In more detail, backpropagation functions compute the gradient  $\nabla_x f(\mathbf{x}, \mathbf{y})$  numerically in a simple and inexpensive procedure.<sup>[22]</sup> What computing the gradient really means is calculating the derivative of the loss function in the model, with the loss function being a selected error function.

The chain rule is a way to calculate the derivate of functions by decomposing a function to other functions whose derivative is known. Let  $x$  be a real number and  $f$  and  $g$  be functions that map from  $\mathbb{R}$  to  $\mathbb{R}$ . By generalising the chain rule of calculus

$$\frac{dz}{dx} = \frac{dz}{dy} \frac{dy}{dx} \quad (2.5)$$

beyond the scalar case into vector notation, such that

$$\nabla_x z = \left(\frac{\delta y}{\delta x}\right)^T \nabla_y z \quad (2.6)$$

where  $\frac{\delta y}{\delta x}$  is a Jacobian matrix of  $g$ . The backpropagation algorithm is essentially a fast computation of this Jacobian gradient for each connection and node in the network.<sup>[11]</sup>

## 2.4 LSTM

Long Short-Term Memory (LSTM) is a type of recurrent neural network first introduced in 1997 by Hochreiter and Schmidhuber. Like all recurrent networks, LSTM differs from feed-forward networks by incorporating cyclical connections between neurons.

The LSTM was introduced specifically to address the two problems of the *vanishing gradient problem* and the *exploding gradient problem*. The two problems are closely linked problems regarding the propagation of errors in deep neural networks. In short, the two gradient problems appear when attempting to train networks with gradient-based learning, iterative methods for computing the gradient at each sample point, and either cause the gradient to be "vanishingly" small (*vanishing gradient*) or to "explosively" increase (*exploding gradient*).<sup>[27][28]</sup>

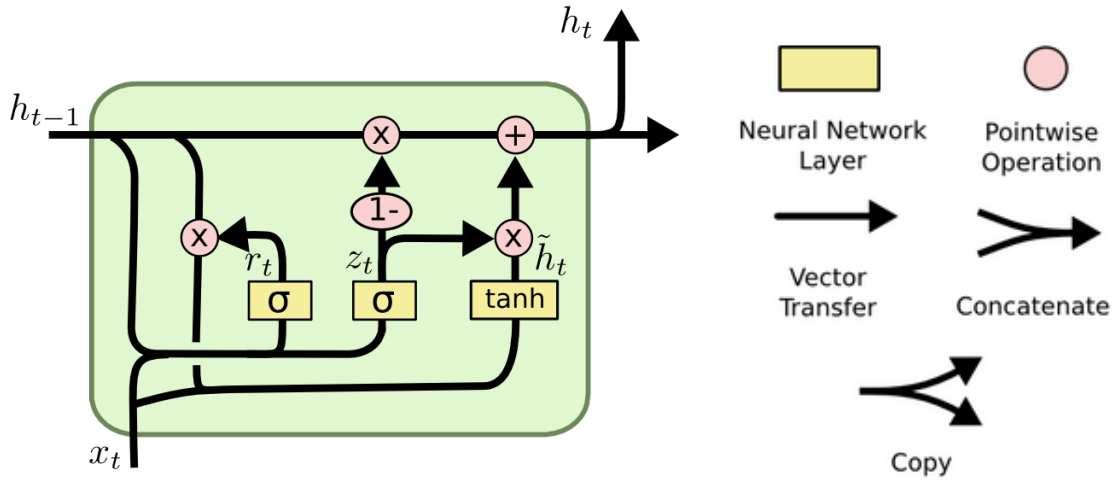


Figure 2.5: Diagram of the LSTM cell

### 3 Autoregressive models

The prices of stocks can be modelled as non-linear time series, which have been at the centre of attention in the finance world since the 1970s with George Box and Gwilym Jenkins popularised their Box-Jenkins method for finding the best-fit of a time series model<sup>[9]</sup>.

## 4 Methodology

The following section provides details in the construction of the model for predicting stock prices, as well as a description of the ARIMA model to which we compare the network.

### 4.1 Time-Series Forecasting

Financial data are discrete in time and as such can be modelled as time-series with calculated means and standard deviations.

### 4.2 ARIMA

To fit an ARIMA model we need to first perform some basic data exploration to identify key features in the time series. Our aim is to attempt to identify the order of the model through observations of the plots instead of relying on inefficient but exhaustive grid search.

In total, there are three model parameters we need to identify:

- $p$ , autoregressive terms
- $d$ , integrated terms
- $q$ , moving average terms

We begin by plotting autocorrelation functions for the entire dataset and each stock individually. Slow degradation of the ACF plot tends to point to autoregressive terms while the opposite is true for moving average terms. Figure 4.1 shows us the autocorrelation plot for all the stocks. We observe a slow fall in the ACF plot, which is consistent with autoregressive terms in the model.

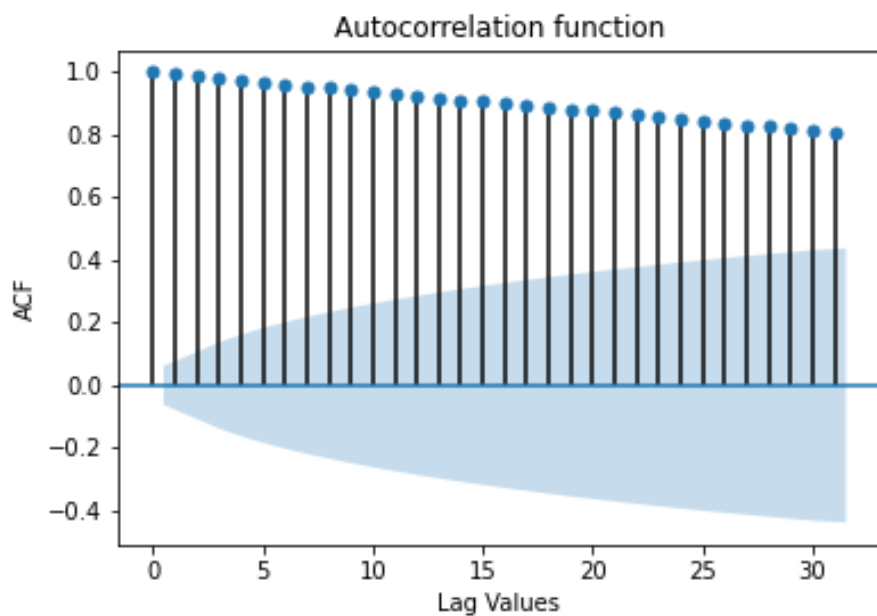


Figure 4.1: Autocorrelation function before any differencing

This graph also tells us that there is significant correlation in the series for over 30 lags. No differencing has been done at this point, the strong correlation could be caused by auto-correlation at lags 1 or 2, which is confirmed by the Partial Autocorrelation function plot below:

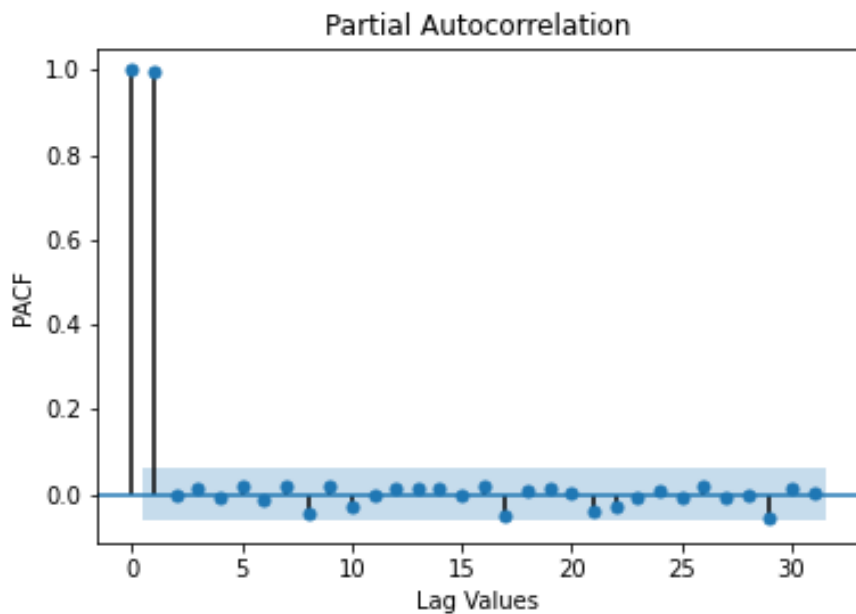


Figure 4.2: Partial Autocorrelation Function

We observe two significant spikes at lags 1 and 2, indicating that the series needs to be differenced twice,  $d = 2$ .

We conclude that a good estimate on the hyperparameters of the ARIMA model we are going to use would be  $(2, 2, 0)$ , for 2 autoregressive terms, 2 differencing terms and 0 moving average terms.

Using the `statsmodel` Python module, we create and validate our ARIMA model using walk-forward validation:

---

```
1     train, test = X[:train_size], X[train_size:]
2     predictions = []
3     for t, _ in enumerate(test):
4         model = ARIMA(train, order=(2,2,0))
5         model_fit = model.fit(dispatch=False)
6         yhat = model_fit.forecast()[0]
7         predictions.append(yhat)
8         train.append(test[t])
9
10    mse = mean_squared_error(test, predictions)
11    mae = mean_absolute_error(test, predictions)
```

---

Figure 4.3: Walk Forward Validation of ARIMA model using the `statsmodels` library

A training and testing split of 80-20 was used to evaluate ARIMA, the same split we use to evaluate the neural network approach.

## 4.3 Model Description

The network is a relatively simple network by most accounts, comprised of a single hidden layer. The simplicity of the model only goes to show the power of neural networks in fitting and forecasting time-series. The model consists of three layers, an input layer, a Long Short-Term Memory (LSTM) layer and the output layer, as shown in 4.4.

The input layer receives a 2-dimensional array of historical stock values, including the pre-

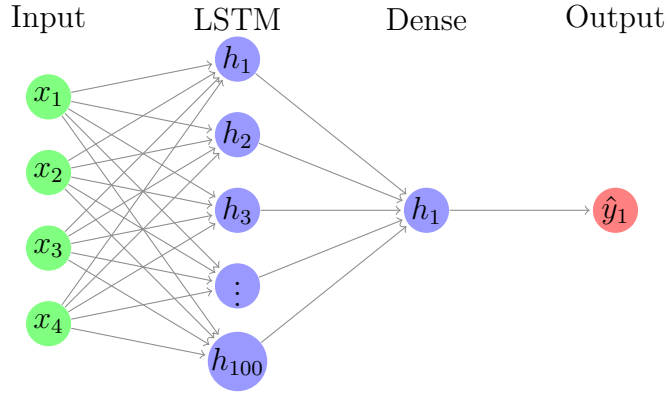


Figure 4.4: Model Architecture

vious day's opening, highest, lowest and closing stock prices. The network then allows the neurons to compete amongst each other and determining an appropriate output. Output is in the shape of a 1-dimensional array containing four values, the future day's predicted stock values for open, high, low and close. The LSTM layer contains 100 nodes and is trained for 100 epochs. The Rectifier Linear Unit was used as its activation function and the Mean Absolute Error was used as its loss function.

Nodes	Epochs	Optimizer	Learning Rate	Activation	Loss
100	100	Adam	0.001	ReLU	MAE

Table 4.1: Description of the hidden LSTM layer

In this project we made use of the Keras API to implement the network. Keras is written in Python and operates on top of TensorFlow, which is an extremely popular and widely used Deep Learning library<sup>[5]</sup>. A simple sequential model using LSTM cells built using Keras would look like this:

---

```

1  model = Sequential()
2  model.add(LSTM(100, activation="relu", input_shape=(n_steps, n_features)))
3  model.add(Dense(n_features))
4  opt = Adam(learning_rate=0.001)
5  model.compile(optimizer=opt, loss="mae", metrics=["mse"])

```

---

With as little as 5 lines of code, we have a working Long Short-Term Memory model ready for training. The LSTM cell easily remembers the long term dependencies in the data and



outputs a 1-dimensional array containing future values. Figure 4.5 contains a more detailed breakdown of the layers used in the network.

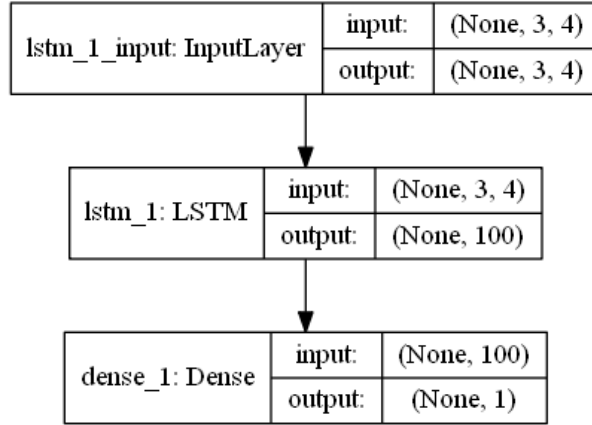


Figure 4.5: Description of layers

A Dense layer is a fully-connected feedforward layer. Feedforward layers are the simplest type of architecture employed in machine learning, which is essentially a collection of neurons that pass signals in one direction only, the layer in front of them. These layers do not form cycles in their connections which is fundamentally different from Recurrent Neural Networks like the LSTM. The purpose of the Dense layer in the model is to reduce the dimensionality of the data from the 100-dimensional space created in the LSTM down to a 1-dimensional array of values, which is then lead to the output layer. The Dense layer performs no operations on the data, it's activation function is linear.

## 4.4 Model Evaluation

To perform an evaluation of the model, we feed into it the testing set gathered from the dataset. The model uses the Mean Absolute Error metric to guide the minimisation of its loss function during training. Mathematically, the Mean Absolute Error is defined as

$$\text{MAE} = \frac{1}{n} \sum_{i=1}^n |\hat{y}_i - y_i| \quad (4.1)$$

and is the sum of the absolute differences of predicted data points to actual data points

divided by the number of data points in the set. MAE is also used in the evaluation of the compiled and trained model.

Evaluation of the model also includes a comparison of Mean Squared Error values. MSE is mathematically defined as

$$\text{MSE} = \frac{1}{n} \sum_{i=1}^n (\hat{y}_i - y_i)^2 \quad (4.2)$$

and is the mean of the squared differences between the predicted values and the actual data points. MAE was chosen for the training of the model because of its linear relationship between penalties and errors, in that it treats a predicted to actual difference of 1 in a way proportionally to how it treats a predicted to actual difference of 5. MSE treats larger differences between predictions and actual data non-linearly.

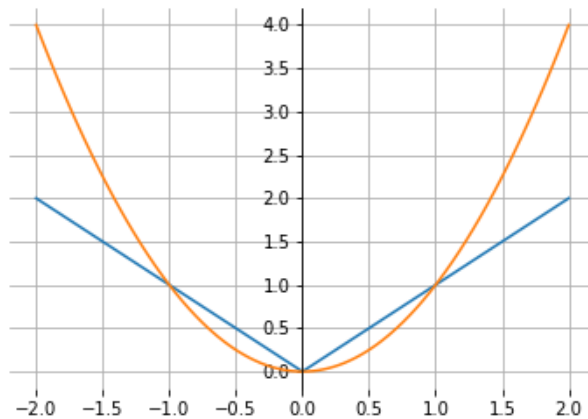


Figure 4.6: Comparison of MAE (blue) and MSE (orange). Note how MAE scales linearly with higher error values while MSE scales quadratically, treating higher errors more harshly.

## 5 Results

### 5.1 Dataset

We retrieve data on tickers available at <https://www.tiingo.com/>. The stock exchanges targeted in this project are the New York Stock Exchange (NYSE) and the National Association of Securities Dealers Automated Quotations (NASDAQ). The data comprise of companies from the technology and the financial services industries.

The data consist of daily values for the period starting January 1, 2016 and ending December 31, 2019. The dataset includes daily information for the `close`, `open`, `high` and `low` prices of each trading day. Below is a table listing all stocks considered for the project:

Stock Name	Symbol	Stock Exchange
Apple Inc.	AAPL	NASDAQ
Amazon Inc.	AMZN	NASDAQ
American Express Company	AXP	NYSE
Boeing Co	BA	NYSE
Bank of America Corp	BAC	NYSE
Citigroup Inc	C	NYSE
Ford Motor Company	F	NYSE
Facebook Inc.	FB	NASDAQ
General Electric Company	GE	NYSE
Alphabet Inc. Class C	GOOG	NASDAQ
Goldman Sachs Group Inc.	GS	NYSE
JPMorgan Chase & Co.	JPM	NYSE
Morgan Stanley	MS	NYSE
Microsoft Corporation	MSFT	NASDAQ
Wells Fargo & Co.	WFC	NYSE

Table 5.1: Stocks included in the study

### 5.2 Evaluation Metrics

### 5.3 Results

## 6 Conclusions

...

# References

- [1] M. Abramowitz and I. A. Stegun, *Handbook of Mathematican Functions with Formulas, Graphs and Mathematical Tables*. National Bureau of Standards, 1972 (cit. on p. 4).
- [2] L. Bachelier, ‘The theory of speculation,’ *Annales scientifiques de l’Ecole Normale Supérieure*, vol. 3, no. 17, pp. 21–86, 1900 (cit. on pp. 1, 2).
- [3] F. Black, ‘Noise,’ *Journal of Finance*, vol. 41, no. 3, 1986. DOI: 10.1111/j.1540-6261.1986.tb04513.x (cit. on p. 2).
- [4] M. Burton, *A Random Walk Down Wall Street*. WW Norton & Company Inc, 1973 (cit. on p. 2).
- [5] F. Chollet *et al.* (2015). ‘Keras,’ [Online]. Available: <https://github.com/fchollet/keras> (cit. on p. 12).
- [6] P. Cootner, *The Random Character of Stock Market Prices*. MIT Press, 1964 (cit. on p. 2).
- [7] E. Fama, ‘Random walks in stock market prices,’ *Financial Analysts Journal*, vol. 21, no. 5, pp. 55–59, 1965. DOI: 10.2469/faj.v21.n5.55 (cit. on p. 2).
- [8] —, ‘Efficient capital markets: A review of theory and empirical work,’ *Journal of Finance*, vol. 25, no. 2, pp. 383–417, 1970 (cit. on pp. 1, 2).
- [9] G. J. George Box, ‘Time series analysis, Forecasting and control,’ 1970 (cit. on p. 8).
- [10] X. Glorot, A. Bordes and Y. Bengio, ‘Deep sparse rectifier neural networks,’ *Artificial Intelligence and Statistics (AISTATS)*, vol. 2011, 2011 (cit. on p. 5).
- [11] I. Goodfellow, Y. Bengio and A. Courville, *Deep Learning*. MIT Press, 2016, <http://www.deeplearningbook.org> (cit. on p. 7).
- [12] S. J. Grossman and J. E. Stiglitz, ‘On the impossibility of informationally efficient markets,’ *American Economic Review*, vol. 70, no. 3, 1980 (cit. on p. 2).
- [13] R. H. R. Hahnloser, R. Sarpeshkar and M. M. et al, ‘Digital selection and analogue amplification coexist in a cortex-inspired silicon circuit,’ *Nature*, vol. 405, pp. 947–951, 2000. DOI: <https://doi.org/10.1038/35016072> (cit. on p. 5).
- [14] R. H. R. Hahnloser and H. S. Seung, ‘Permitted and forbidden sets in symmetric threshold-linear networks,’ *NIPS*, 2001 (cit. on p. 5).
- [15] S. Hochreiter and J. Schmidhuber, ‘Long short-term memory,’ *Neural Computation*, vol. 9, no. 8, 1997 (cit. on p. 7).
- [16] C. D. Kirkpatrick and J. R. Dahlquist, *Technical Analysis: The Complete Resource for Financial Market Technicians*. 2006, ISBN: 978-0-13-153113-0 (cit. on p. 1).

- [17] S. Kleene, ‘Representation of events in nerve nets and finite automata,’ *Annals of Mathematics Studies*, pp. 3–41, 1956 (cit. on p. 3).
- [18] A. L. Maas, A. Y. Hannun and A. Y. Ng, ‘Rectifier nonlinearities improve neural network acoustic models,’ *Stanford Press*, 2014 (cit. on p. 5).
- [19] B. Mandelbrot, ‘The variation of certain speculative prices,’ *The Journal of Business*, vol. 36, no. 4, pp. 394–419, 1963. [Online]. Available: <http://www.jstor.org/stable/2350970> (cit. on pp. 1, 2).
- [20] W. McCulloch and W. Pitts, ‘A logical calculus of ideas immanent in nervous activity,’ *A Logical Calculus of Ideas Immanent in Nervous Activity*, vol. 5, no. 4, 1943. DOI: 10.1007/BF02478259 (cit. on p. 3).
- [21] R. Prajit and B. Zoph, ‘Searching for activation functions,’ 2017 (cit. on p. 5).
- [22] D. Rumelhart, G. Hinton and R. Williams, ‘Learning representations by back-propagating errors,’ *Nature*, vol. 323, pp. 533–536, 1986. DOI: 10.1038/323533a0 (cit. on p. 6).
- [23] S. Russel and P. Norvig, *Artificial Intelligence, A Modern Approach*, Third. Pearson, 2010 (cit. on p. 3).
- [24] P. A. Samuelson, *Proof that Properly Anticipated Prices Fluctuate Randomly*. 1965, ch. Chapter 2, pp. 25–38. DOI: 10.1142/9789814566926\_0002 (cit. on pp. 1, 2).
- [25] P. J. Werbos, *Beyond Regression: New Tools for Predictions and Analysis in the Behavioural Sciences*. Harvard University Press, 1975 (cit. on p. 6).
- [26] W. H. Wilson, *The Machine Learning Dictionary*. Saint Elizabeth University, 2009 (cit. on pp. 3, 4).
- [27] B. Y, S. P and F. P, ‘Learning long-term dependencies with gradient descent is difficult,’ *IEEE Transactions on Neural Networks*, vol. 5, no. 2, pp. 157–166, 1994 (cit. on p. 7).
- [28] B. Y, M. T and P. R, ‘On the difficulty of training recurrent neural networks,’ 2012. DOI: 1211.5063 (cit. on p. 7).

# List of Tables

4.1	Description of the hidden LSTM layer . . . . .	12
5.1	Stocks included in the study . . . . .	15

# List of Figures

2.1	Graphical representation of a Neuron . . . . .	3
2.2	Plot of the Rectified Linear Unit function . . . . .	4
2.3	Plot of the sigmoid function . . . . .	5
2.4	Plot of the hyperbolic tangent function . . . . .	6
2.5	Diagram of the LSTM cell . . . . .	7
4.1	Autocorrelation function before any differencing . . . . .	10
4.2	Partial Autocorrelation Function . . . . .	10
4.3	Walk Forward Validation of ARIMA model using the <b>statsmodels</b> library .	11
4.4	Model Architecture . . . . .	12
4.5	Description of layers . . . . .	13
4.6	Comparison of MAE (blue) and MSE (orange). Note how MAE scales linearly with higher error values while MSE scales quadratically, treating higher errors more harshly. . . . .	14