

# Short Term Stock Price Prediction with Machine Learning

Lucas Breinlinger<sup>1</sup>

## Abstract

Accurately predicting stock prices is a highly sought-after mastery. In recent years Long Short Term Memory (LSTM), a type of recurrent neural network, has proven to be highly promising in achieving better predictions. Using a simple LSTM with only 25 neurons already can predict loosely the future stock prices based on the stock prices of the previous 30 days. The refined model with two LSTM layers with 150 neurons and two dense layers is very successful in predicting the future close price.

<sup>1</sup> Contact: [lucas.nm.breinlinger@bmw.com](mailto:lucas.nm.breinlinger@bmw.com)

## Contents

1	<b>Project Definition</b>
2	<b>Data Analysis</b>
3	<b>Methodology</b>
4	<b>Results</b>
5	<b>Conclusion</b>
	<b>References</b>

## 1. Project Definition

The accurate forecast of stock prices, traded on an exchange, is inarguably one of the most challenging topics in the field of asset pricing. The stock market is described as an unpredictable, dynamic and non-linear construct. Predicting the prices of its stocks is an ambitious undertaking and depends on many variables including but not limited to the global economy, the company's metrics and local and global political situation.

Historically there are two main prediction methods. The fundamental analysis which is a qualitative analysis of the company is interested in finding the true value of a stock and comparing it to the actual traded value. The evaluator utilizes several performance criteria e.g. the P/E ratio to truly assess the underlying stock. Secondly, the technical analysis, which is solely based on the past price of the stock e.g. in form of closing or opening prices as time-series. It is rather a short-term prediction using factors like the Moving Average (MA) or the Exponential Moving Average (EMA). Its basic assumption is that every significant information about the stock is already considered in the stock price [1].

The fast computational development has led to the point that Machine learning techniques have a significant application in financial problems. The use of artificial neural networks has

found more and its way into the field of stock price prediction. Here a recurrent neural network (RNN) has been found very proven, more precisely the Long Short Term Memory (LSTM). Its advantage is being able to process entire sequences of data rather than only one single data point. It has proven to be very practical with time series data such as our historical stock prices.

In this project, I create an application that is going to predict the **closing price** of any given stock on which it is trained. For the sake of convenience, this report only considers the stock price prediction of the **Apple Inc. (\$AAPL)** stock.

A train-test cycle is used to test the accuracy of the prediction. The root mean squared error (1) is used to measure the accuracy of the model:

$$RMSE = \sqrt{\frac{\sum_i (\hat{y}_i - y_i)^2}{N}} \quad (1)$$

where  $\hat{y}_i$  is the prediction of the  $i - th$  stock price  $y_i$  and  $N$  is the total number of predicted stock prices. The RMSE is a good measure to compare prediction with actual values and makes it easier to compare different models

In the next section, I will analyze the data and give some insights about the data. Afterwards I will explain the methodology and discuss the results. In the end, I give an outlook for further developments and improvements.

## 2. Data Analysis

The dataset was collected with the help of `yfinance` python plugin. In my case, it is replacing the not working Yahoo Finance API. The collected data frame contains

- 1 108 837 rows
- 10 columns
- price information for 457 unique stocks
- Data of the last 10 years (including 17/07/2022)

**Table 1.** First lines of *Apple Inc.* stock prices dataset

Ticker	Name	Date	Open	High	Low	Close	Volume	Dividends	Stock Splits
AAPL	Apple Inc.	2012-07-19	18.77	18.90	18.61	18.87	436861600	0	0
AAPL	Apple Inc.	2012-07-20	18.83	18.87	18.54	18.56	397471200	0	0
AAPL	Apple Inc.	2012-07-23	18.25	18.61	18.05	18.54	487975600	0	0
AAPL	Apple Inc.	2012-07-24	18.65	18.72	18.38	18.45	565132400	0	0
AAPL	Apple Inc.	2012-07-25	17.64	17.84	17.51	17.66	877312800	0	0

**Table 2.** Statistics of the Apple stock prices. Which represent our used features

	Open	High	Low	Close
count	2514.0	2514.0	2514.0	2514.0
mean	55.0	55.6	54.4	55.0
std	46.0	46.6	45.4	46.0
min	12.1	12.5	12.0	12.2
25%	22.8	23.1	22.7	22.8
50%	36.5	36.7	36.2	36.5
75%	66.7	68.5	66.0	67.5
max	182.6	182.9	179.1	182.0

and was collected on 18/07/2022. I use a list of 457 randomly selected stock tickers to generate the dataset.

The structure and the first 5 lines of the table filtered for *Apple Inc.* are displayed in table 1. For each day and each stock (denoted by its ticker (*Ticker*) and company name (*Name*)) the dataset contains the price of the stock in dollar when opening the market (*Open*), closing the market (*Close*). Furthermore there the days lowest (*Low*) and highest intra-day price (*High*) is available, as well as the daily traded volume (*Volume*).

Looking at the *Apple Inc.* stock there are 2514 columns regarding the apple stock with a minimum closing price of 12.17 and a maximum closing price of 182.01. There are no missing values (see table 2) in the dataset, and the date and value columns are correctly set as a datetime respectively float columns. In table 2, the statistics of the distinctive variables are shown.

**Figure 1.** Closing price of the *Apple Inc.* in the last 10 years.

In figure 1, the highly volatile closing price of the stock

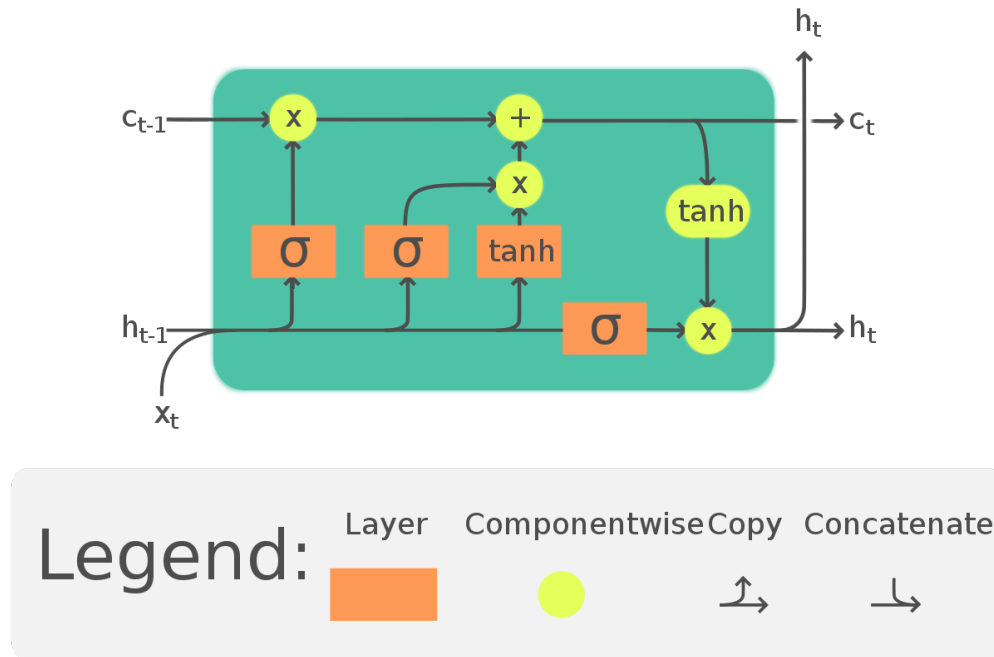
is shown for its last 10 years. In recent years the stock prices have increased roughly sevenfold. In the next section, I am going to explain the methodology, followed by a chapter discussing the results.

### 3. Methodology

Here, I am using a Recurrent Neural Network for the stock price prediction. More specifically a Long Short Term Memory Network is used to build the model. These networks are capable of learning long-term dependencies and are commonly used for handling and forecasting time-series data. In figure 2, a recurrent cell of a LSTM network is shown. Each cell consists of three gates, an input gate, an output gate and a forget gate. [3]

Those cells are linked together to a chain, where each cell looks at the input gate at the previous hidden state  $h_{t-1}$ , the current input  $x_t$ , and the previous cell state  $c_{t-1}$ . The forget cell is the first part where  $h_{t-1}$  and  $x_t$  are passed through the sigmoid activation function and pointwise multiplied with the cell state  $c_{t-1}$ . The input gates combine a sigmoid and tanh layer, here it is decided which parts are going to be updated. Multiplying both layers determines which information is used to update the cell state by pointwise adding it to the output of the forget gate, resulting in the final  $c_t$ . In the output gate, the same input ( $h_{t-1}$  and  $x_t$ ) is passed through a sigmoid function and finally multiplied with the current cell state  $c_t$ , yielding the final output the hidden state  $h_t$ . [4]

For the LSTM to work properly, the stock closing prices have to be normalized. This will help our model to converge faster and be more stable. Furthermore, when using the scaler on more than one input feature with different scales all features will contribute equally to the fitting and there will be no bias due to a different scaling [5]. I am using the MinMaxS-



**Figure 2.** The recurrent cell of a LSTM network [2].

caler (2) of sci-kit learn.

$$x_{i,scaled} = \frac{x_i - \min(x)}{\max(x) - \min(x)} \quad (2)$$

Here, I am trying to predict future stock prices  $y_t$  (target) with the help of the previous  $n$  stock prices. For example, I am predicting the stock price from day 11 using day 1 to 10 as my features, for day 12, using day 2-11 and so on. In order to do this, I defined a splitting function that generates a data sample consisting of  $n$  samples (the previous stock prices) and the target value  $n + 1$ .  $n$  is denoting the size of my window which I am using to predict the prices. This function is applied to my previously scaled stock prices.

Finally, I split this data set into a training and a test set, since we are dealing with a time-series we cannot randomly select the samples, rather I split the sample at some point. In my experiment, I am using 80% as a training set and the remaining 20% as my test set. Using the predicted stock price  $\hat{y}_i$  and the real stock price  $y_i$  of the test set the RMSE can be easily calculated using eq. 1.

**Table 3.** Overview of the model parameters.

Model	LSTM 1	LSTM 2	Dense 1	Dense 2	Epoch	$n$
simple	10	-	1	-	1	30
refined 1	10	10	25	1	1	30
refined 2	50	50	25	1	1	30
refined 3	150	150	25	1	5	30

The first **simple model** I build, consisted of a single LSTM layer with 10 neurons and a dense Layer with 1 neuron. The

dense layer is needed to get the dimensionality of the output of the LSTM  $h_t$  back to a 1D output. I am using the previous thirty close prices ( $n = 30$ ) as my input  $x_t$ .

My refined model is a deeper neural network since the relations between stock prices are quite complicated. Thus, I added an LSTM layer and a second dense layer with 25 neurons. The refined model is run with different parameters, the first run is 10 neurons for both LSTM layers, the next run is 50 neurons for both layers and in my last model, I am using 150 neurons and 5 epochs. All the time  $n$  is kept at 30 days.

In table 3 the model parameters are shown in an overview. In the next section, I will discuss the results of these models.

## 4. Results

**Table 4.** The resulting *RMSE* of the four models.

Model	<i>RMSE</i>
simple	42.49
refined 1	28.34
refined 2	9.59
refined 3	3.74

As expected the *RMSE* for the refined model is smaller than for the simple model. Furthermore, refined model 3 with 150 neurons per LSTM and 5 epochs also performed better than refined models 1 and 2 with fewer neurons and only one epoch.

For comparison see the values for all four models in table 4. In figure 1 the predicted stock prices for all models are plotted

versus the actual stock price. Furthermore figure 1 is a detailed view of the stock price prediction for the last 60 days.

As one can see in the detailed graphics for each model, the more complicated the model, the better the prediction. The poor performance of the simple model can be seen in figure 3a and 3b. The prediction follows loosely the real price, however, it is much higher than the real price. The refined model 1 already performs slightly better than the simple model (compare RMSE) and figure 3c and 3d, but also fails to make an accurate prediction.

A huge leap in prediction accuracy is made when the number of neurons is increased from refined model 1 to refined model 2. This accuracy can even be improved by tripling the number of neurons and going through the dataset 5 times (epochs). The predicted price follows closely the actual price, see 3f.

The refined model 3 is the most accurate of my models. It has the lowest RMSE and in figure 3h it is visible that it follows closely the actual price.

## 5. Conclusion

In this paper, I looked at the problem of forecasting accurately stock prices. In the beginning, I give an overview of why this topic matters so much for the financial world. Additionally, I present some real-life data on the stock prices of Apple stock and some data insights before explaining the methodology. Stock prices can be treated as time-series data, thus LSTM neural networks are predestinated as a suitable algorithm for predicting stock prices. The simple model is falling quite short in predicting the stock prices, however, my most complicated model is roughly 11 times better than the simple one.

Even though the model is quite good at predicting stock prices, in a real-life scenario it probably will perform poorly due to unforeseen and disruptive events (e.g. COVID19 or Ukraine War).

Generating test and training sets were more complex than usual since I am dealing with time series data. Also, I needed to put together a timeframe of past stock prices as my feature variables with its corresponding stock price prediction. This was quite tricky to generate.

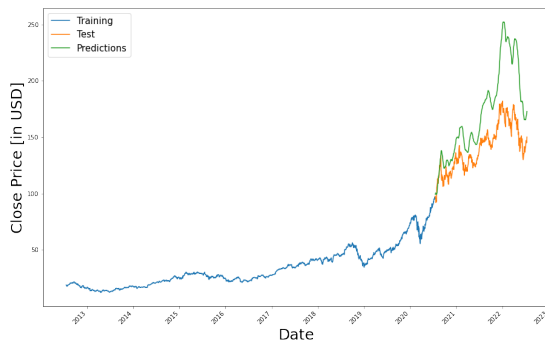
For future research, I will try and compare different machine learning algorithms and compare their performance with the LSMT. Furthermore one could try to play with the parameters, and layers in my current network to make even better predictions of the future stock price. Also, I will use different stocks and see if the performance of the models is going to change.

## References

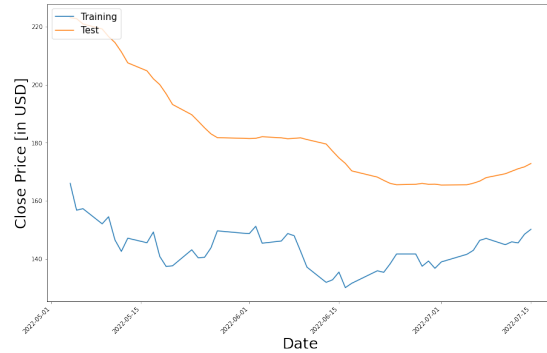
- [1] Stock market prediction. [https://en.wikipedia.org/wiki/Stock\\_market\\_prediction](https://en.wikipedia.org/wiki/Stock_market_prediction). Accessed: 2022-07-16.
- [2] Long short-term memory. [https://en.wikipedia.org/wiki/Long\\_short-term\\_memory#](https://en.wikipedia.org/wiki/Long_short-term_memory#)

/media/File:LSTM\_Cell.svg. Accessed: 2022-07-19.

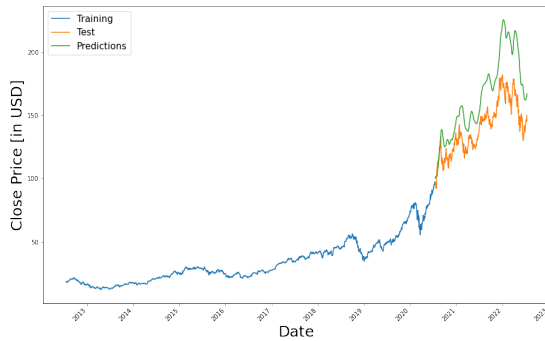
- [3] Gentle introduction long short-term memory networks. [https://machinelearningmastery.com/gentle-introduction-long-short-term-memory-network#:~:text=Long%20Short%20Term%20Memory%20\(LSTM,complex%20area%20of%20deep%20learning](https://machinelearningmastery.com/gentle-introduction-long-short-term-memory-network#:~:text=Long%20Short%20Term%20Memory%20(LSTM,complex%20area%20of%20deep%20learning.). Accessed: 2022-07-19.
- [4] Introduction to long short-term memory. <https://www.analyticsvidhya.com/blog/2021/03/introduction-to-long-short-term-memory-lstm/>. Accessed: 2022-07-19.
- [5] Why data should be normalized before training a neural network. [https://towardsdatascience.com/why-data-should-be-normalized-before-training-a~:text=Among%20the%20best%20practices%20for,and%20leads%20to%20faster%20convergence](https://towardsdatascience.com/why-data-should-be-normalized-before-training-a~:text=Among%20the%20best%20practices%20for,and%20leads%20to%20faster%20convergence.). Accessed: 2022-07-22.



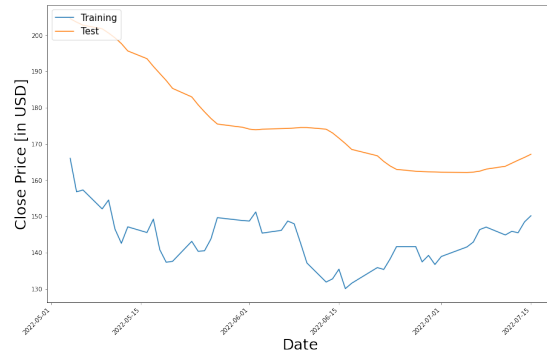
(a) The prediction close price of the simple model versus the real close price.



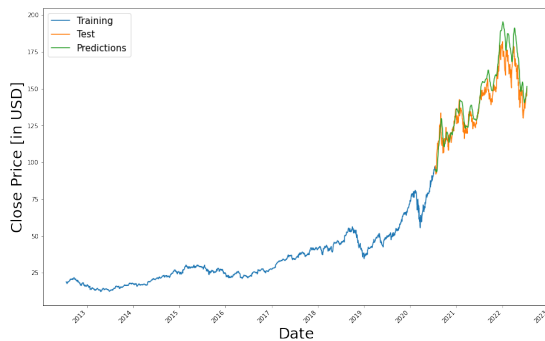
(b) The prediction close price for the last 60 days of the simple model versus the real close price.



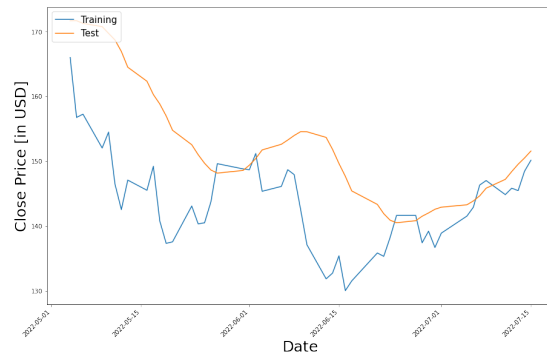
(c) The prediction close price of the refined model 1 versus the real close price.



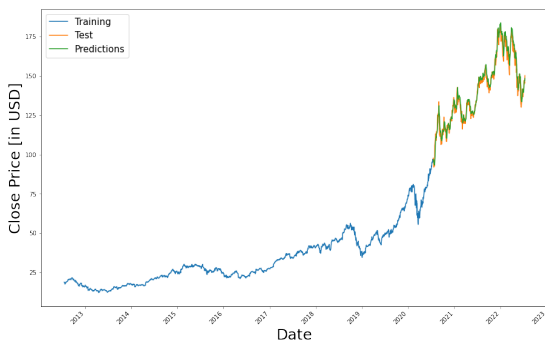
(d) The prediction close price for the last 60 days of the refined model 1 versus the real close price.



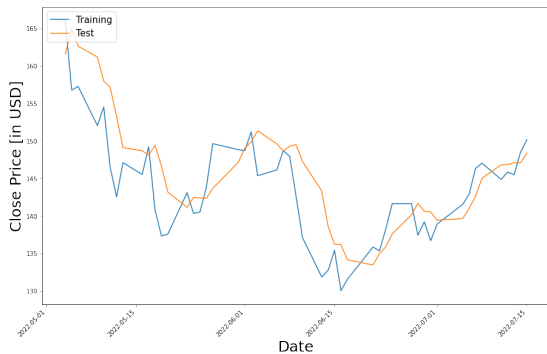
(e) The prediction close price of the refined model 2 versus the real close price.



(f) The prediction close price for the last 60 days of the refined model 2 versus the real close price.



(g) The prediction close price of the refined model 3 versus the real close price.



(h) The prediction close price for the last 60 days of the refined model 3 versus the real close price.

**Figure 3.** The resulting stock price prediction for all models vs. the actual stock price.