

---

# Stock Prediction Using LSTM, Decision Trees and Linear Regression

---

Varun Ghantasala  
gvarun@vt.edu

## Abstract

In today's world of finance understanding stock prices is of utmost importance. A lot of stockbrokers are dive deep into various analysis for stock price prediction to place their bets. This paper explains an approach to understand the volatility of stock prices using Machine Learning Techniques. This approach is ready to be used for any stock to conduct a detailed analysis and predict stock prices of the future. In this paper Machine Learning strategies called Linear Regression, LSTM and Decision Trees are used to predict stock prices. Their accuracy and reliability shall be explored in the paper. A comparison between the methods shall also be made.

## 1 Introduction

The volatility of stock prices makes it very hard and can result in inconsistent analysis if done by humans. Quantitative stock brokers try to find patterns in stock to invest in them. In today's world, Machine learning techniques are greatly used in understanding stock since a computer can compute prices and take into various variables to predict stocks. Some of these Machine Learning methods have proven to be useful. With methods that are depicted in the paper, a step towards understanding pattern and prediction of data can be made. In this paper, Linear Regression, Decision Trees and LSTM models have been used to predict stock prices.

### 1.1 LSTM

Long short term memory is a type of neural network under the subset of Recurrent Neural Networks. Unlike other neural networks which have forward feed, LSTM uses feedback connections. An LSTM unit has a cell, input gate, output gate and forget gate. LSTM can keep track of arbitrary long term correlations in input sequences. LSTM can also handle the vanishing gradient problem faced by RNN due to forget gate. The forget gates chooses whether to discard previous inputs. The input gate tries to learn new information. Whereas the output gate passes the new information from current timestamp to the next. LSTM can be used for classifying, forecasting and processing. Major uses of LSTM include handwriting and speech recognition, time series analysis and video games.

### 1.2 Linear Regression

Linear Regression is one of the most well-known algorithms in machine learning. It was initially developed in the field of statistics. It understands relationship between input and output variables in the form of a linear relationship by reducing the cost function error. In other words it attempts to find the relationship between two variables by fitting a linear equation to the observed. The simplest form of the regression equation with one dependent and one independent variable is  $y = a + b \cdot x$  where  $y$  is expected value,  $a$  is the constant,  $b$  is called the regression coefficient and  $x$  is the true value. Major uses of Linear Regression include forecasting an event, trend forecasting, understanding the correlation between two values.

### 36 1.3 Decision Trees

37 As the name suggests, Decision tree has a tree like model of decisions. Decision Tree algorithm is  
38 a sub-category of supervised algorithms. Simply put, Decision Trees tries to create a model from  
39 training data by understanding inherent decision rules. This algorithm uses multiple algorithms to  
40 create and split a node into sub-nodes. As the tree grows it decides what feature to choose and what  
41 functions to use for splitting. It can solve both regression and classification problems. It could be  
42 used for both categorical and continuous variables. It could be used in ecommerce, diagnosis of  
43 diseases and to detect frauds.

## 44 2 Methodology

### 45 2.1 LSTM

46 The input values are Adjusted Close stock data prices of AAPL. The data for X label had an array  
47 of 30 values with a label of 31st value. The data fed was first split in to train and test data with 4:1  
48 ratio respectively. The model had 5 LSTM layers including dropout layer. After data preprocessing,  
49 the data was fed into LSTM model. Additionally, the loss parameter was mean squared error and  
50 optimizer used was Adam Optimizer. The layers were initially taken from research papers as standards  
51 and tested while changing the layer parameters to arrive at the best solution. After this, the number of  
52 epochs was also decided based on the epochs vs loss graph. Furthermore, batch size was notably of  
53 the order 2 which was also experimented to arrive at the best batch size model. The optimal epochs  
54 and parameters were used to create and fit the model. After which, the training data was used to  
55 predict values. This was compared using the RMSE performance metrics.

### 56 2.2 Linear Regression

57 The input values are the same as considered in LSTM. Whereas in this case the input X and Y values  
58 for fitting the model was training data and its respective label after a fixed number of days. In my  
59 case the fixed number of days was set to 30. In other words, the first data points label shall be 30th  
60 data point. In this was the model was resourced from scikit learn to fit X and Y values. After which  
61 training data was used to predict the expected values which was then compared with test values.

### 62 2.3 Decision Trees

63 The input values are the same as considered in LSTM and Linear Regression i.e. Adjusted Close  
64 stock prices. Whereas in this case the input X and Y values for fitting the model was training data  
65 and its respective label after a fixed number of days. In my case the fixed number of days was set  
66 to 30. In other words, the first data points label shall be 30th data point. In this was the model was  
67 resourced from scikit learn to fit X and Y values. After which training data was used to predict the  
68 expected values which was then compared with test values.

## 69 3 Data Collection and Preprocessing

70 Date Collection is fundamental module needed for any type of analysis. The data of stock prices  
71 under analysis in this paper is that of Apple. Inc with stock ticker abbreviation 'AAPL'. The data is  
72 extracted from Yahoo Finance under the historical prices section. The data collected 1st Jan 2019  
73 to 31st December 2019. This data is limited to 2019 since COVID related events has caused an  
74 inconsistency in data patterns from 2020 and has been shown in figure 1.

### 75 3.1 LSTM

76 LSTM For data preprocessing for LSTM requires an input of three-dimensional arrays. The first step  
77 was to normalize or scale the data to range (0,1) since LSTM is very sensitive to the scale of data.  
78 This was followed by reshaping into 3-dimensional arrays as (samples, time steps, features).

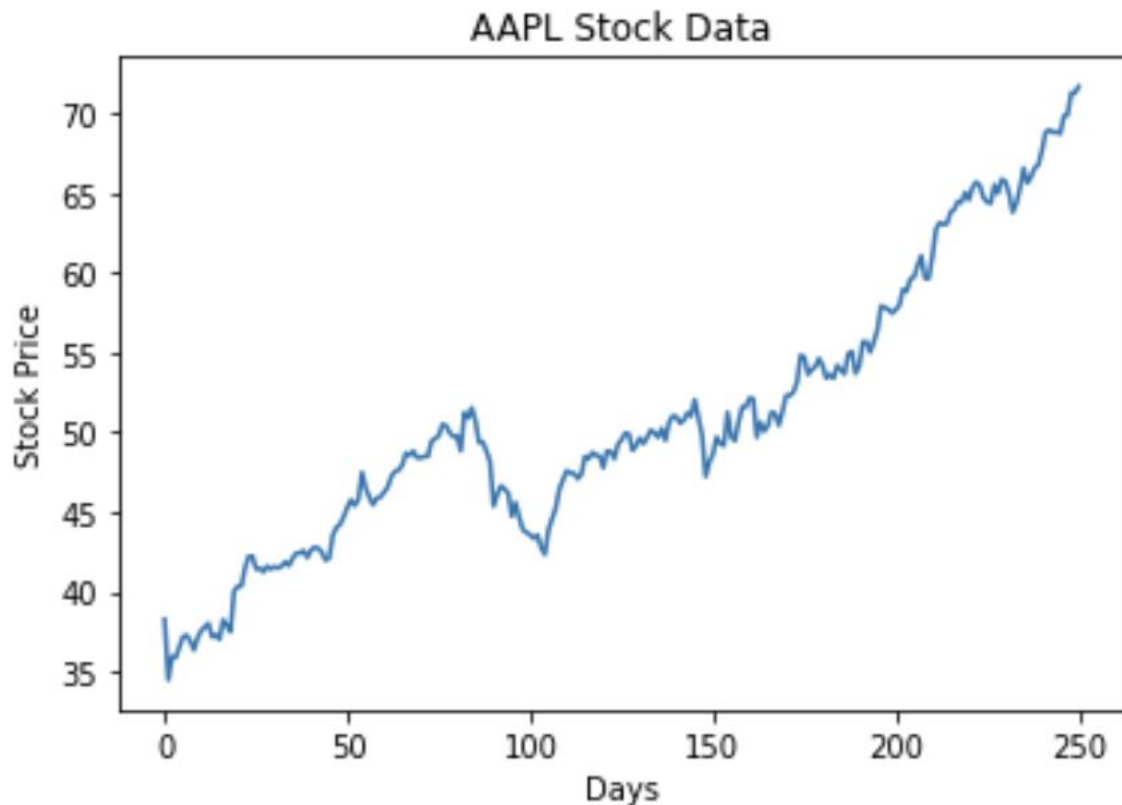


Figure 1: Original APPL Stock Price Data

### 79 3.1.1 Linear Regression and Decision Trees

80 Not much preprocessing of data was done since the data was already available in data frames. The  
81 target data was the 30th data point of the first data point which is the X variable in this case. The data  
82 was split in training and test in 4:1 ratio. The X and Y data points were finally fed into the Linear  
83 Regression model and Decision Trees models respectively.

## 84 4 Results

### 85 4.1 LSTM

86 The LSTM predictions with respect to original stock data is shown below in Fig 2. The variations do  
87 not seem to be much from the visual aspect which can also be proven through RMSE values. The  
88 RMSE value for LSTM test predictions is 1.23 and train predictions is 0.983. The training loss vs  
89 epochs for the LSTM model is shown below the epoch chosen for my case is 500.

90 The epoch values were selected based on Fig.3 to arrive at an optimum epoch value

### 91 4.2 Linear Regression

92 The results on APPL stock prices predictions are shown below. The variations from the graph seem  
93 to be higher than the other two models. The RMSE value is also high and comes out to be 7.035

### 94 4.3 Decision Trees

95 The results of prediction are given below in the graph. The variations in the graph are not as much  
96 compared to Linear Regression. The prediction error is 4.037

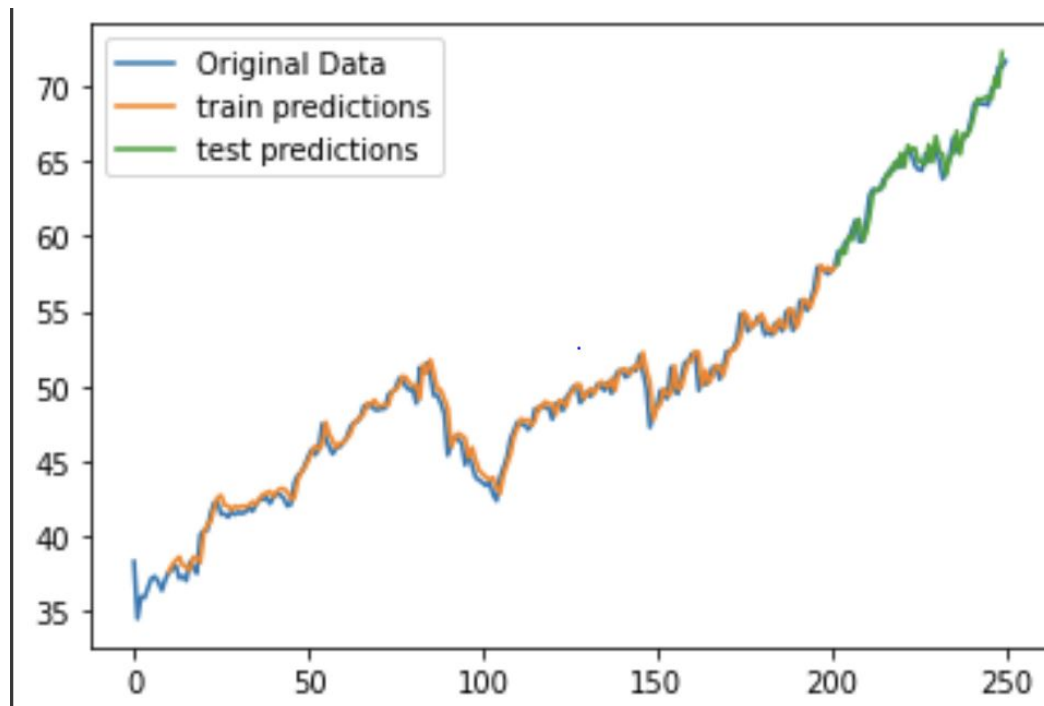


Figure 2: LSTM Prediction Results

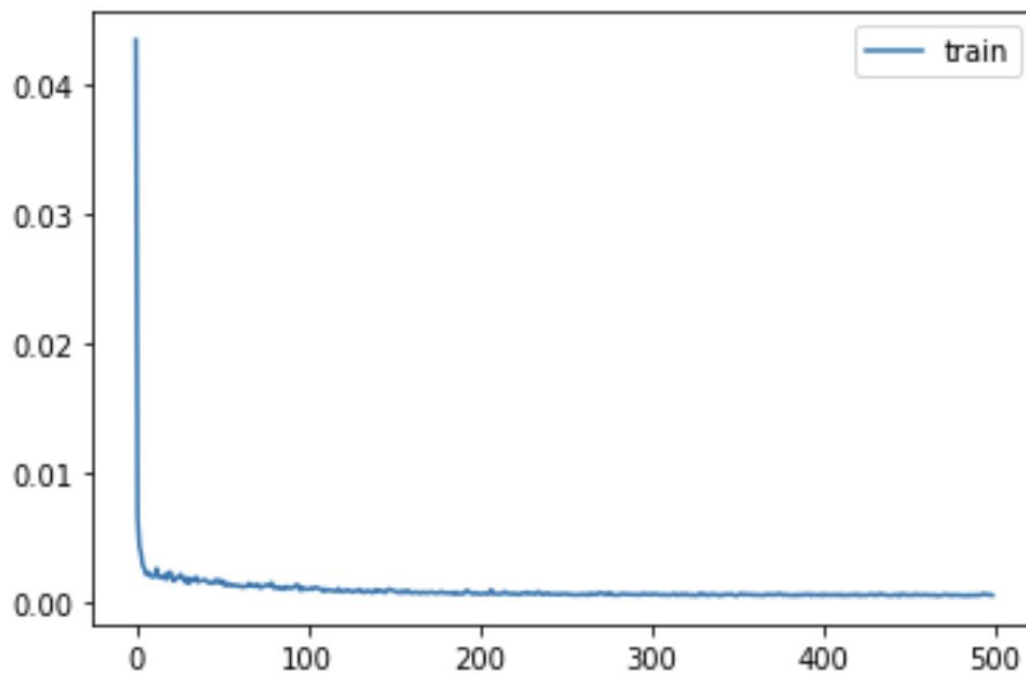


Figure 3: training loss

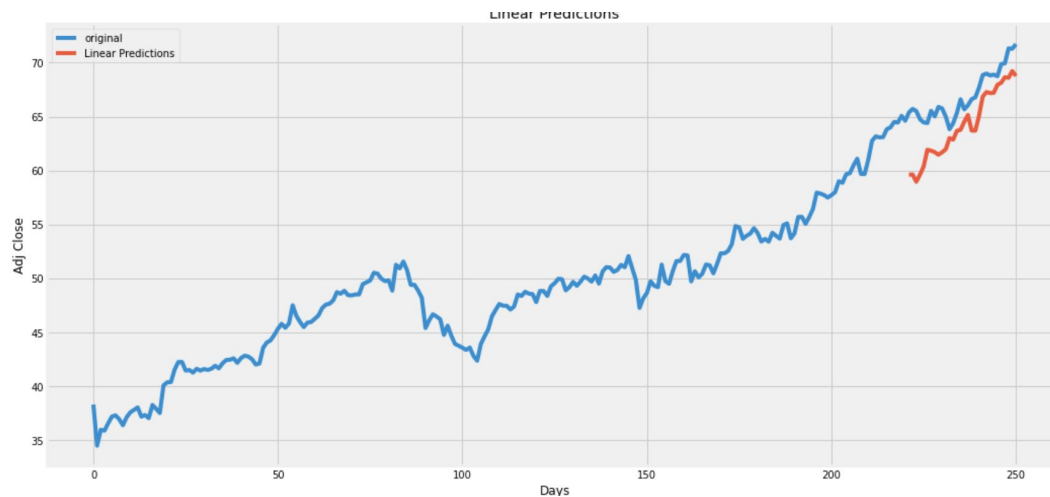


Figure 4: Linear Regression Prediction Results

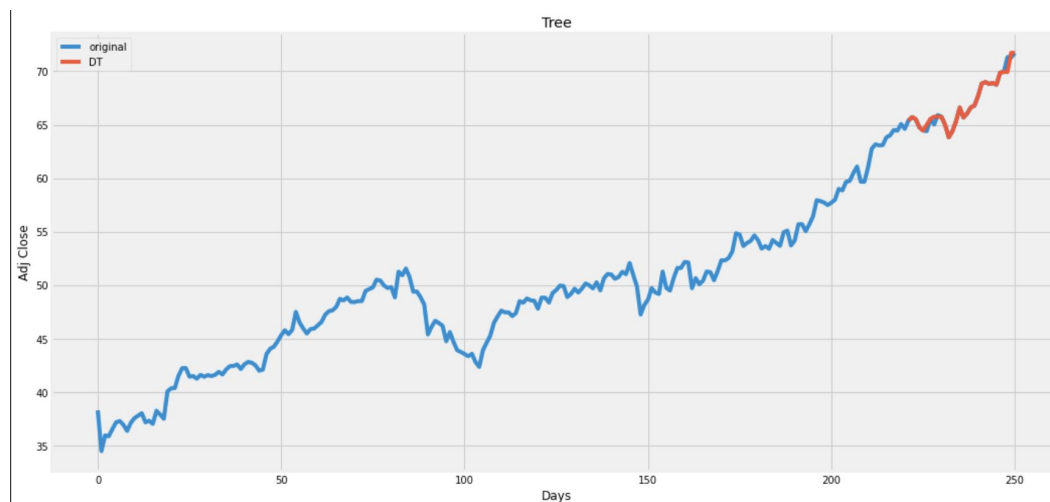


Figure 5: Decision Tree Prediction Results

## 97 5 Future Scope

98 Future scope of this is to try on more data and also different stock price data to understand the  
 99 efficiency of the models in use such as LSTM, Decision Trees and Linear Regression. Additionally,  
 100 further fine tuning of the models also needs to be experimented. Another effective approach is to  
 101 include more data features such as volume, open price, and close stock price to get more efficient  
 102 analysis.

## 103 6 Conclusions

104 The graph plots explain that LSTM performed better than Decision Trees followed by Linear  
 105 Regression. The RMSE value also explains the same comparison which is that LSTM is more  
 106 efficient in predicting values than Decision Trees followed followed by Linear Regression. A more  
 107 detailed analysis with more features and optimum parameters will give more efficient results which  
 108 can be used to predict future data making the job of stockbrokers a lot easier. Though these models  
 109 need to be only used as a second perspective since the market events have lot of factors that could  
 110 affect the stock prices of the future days. I truly believe that the true value of stock prices cannot be  
 111 predicted. Since that would only mean predicting future events which is not possible. In other words,

112 if you cant predict future events then you cant predict stock prices since the stock prices are proven to  
113 be sensitive to the events around us.

## 114 **7 Appendix**

115 Fig1: Original APPL Stock Price Data - Pg. 3

116 Fig2: LSTM Prediction Results - Pg. 4

117 Fig3: Training loss - Pg. 4

118 Fig4: Linear Regression Results - Pg. 5

119 Fig5: Decision Tree Results - Pg. 6

120

## 121 **8 References**

122 Singh, S.,Rehan S., Kumar v. (2022). Stock Price Prediction Using Linear Regression, LSTM and  
123 Decision Tree. *Easy Print*. Preprint no.7805. <https://easychair.org/publications/preprint/rD6d>

124 <https://www.tensorflow.org/apidocs/python/tf/keras/layers/LSTM>

125 Karim, Rezaul Alam, Md Hossain, Md. (2021). Stock Market Analysis Using Linear Regression  
126 and Decision Tree Regression. 1-6. 10.1109/eSmarTA52612.2021.9515762.