

Predicting the Stock Market: Modeling 30 Years of S&P 500 Trading

Aria Alaghemand, Alec Anderson, Reed Oken

Applied Artificial Intelligence, University of San Diego

AAI 501: Introduction to Artificial Intelligence

Professor Ying Lin

April 17, 2023

Introduction

In today's world, the stock market is a key indicator of economic stability and a source of investment opportunities for many individuals and businesses. Accurately predicting future trends in the stock market is a challenging and complex task that requires extensive data analysis and the use of advanced machine learning algorithms (Vij et al., pg 1). This project analyzes 30 years of pricing data for the S&P 500 and compares several models with the goal of identifying a model which can make the most accurate predictions. A variety of machine learning algorithms including regression, time series analysis, and deep learning will be used to train models on historical data and predict future market trends. Multiple course topics will be applied, of note, regression, supervised and unsupervised machine learning, deep learning, and neural networks. Our system should be able to handle many years of data, deal with noisy data, and adapt to changing market conditions. Additionally, we aim to test the extent that past performance in the market can influence future prices and provide insights into which models can accurately predict these trends.

Moreover, the project aims to create a robust system that can accurately predict market trends while being adaptable to new data and market conditions. We will employ various data preprocessing techniques such as normalization, outlier detection, and imputation to prepare the data for analysis. Additionally, we will perform feature engineering to extract relevant features that may influence the stock market trends. The evaluation of model performance will be based on metrics such as mean absolute error (MAE), mean squared error (MSE), and root mean squared error (RMSE). Ultimately, the goal of this project is to provide a reliable system that can assist individuals and businesses in making informed investment decisions based on accurate predictions of future market trends.

Data Cleaning/Preparation

Data cleaning and preparation are crucial aspects of any data analysis project. For the S&P 500 stock data analysis project, there are several steps that were taken to prepare the data for analysis. Data was initially collected from Kaggle, where 30 years of daily S&P 500 trading data dating back to 1993 is freely available for download as a .csv. This data can also be easily obtained via Yahoo finance. Data was loaded into a Pandas dataframe in Python to prepare for cleaning and an exploratory data analysis. Minimal cleaning was required for this dataset, as the historical data collection has been well maintained across many platforms. Of note, the dataset used for this project included the trading date as an object, which was transformed into a date string for easier use in the project. There were no null entries in the dataset at this point, so there was no need to take any additional steps.

Exploratory Data Analysis

The exploratory data analysis for the project primarily focused on expanding the available information that was available for each day of trading through feature extraction, in addition to identifying an appropriate target for modeling. Each day of trading data initially contained the date, daily trading prices (open, close, high, low), daily trading volume, and the date split into individual columns; year, month, day of month, week of the year, and day of the week. To build upon this, there were several key features identified for extraction.

In order to provide more features for models to train and predict on, common financial technical indicators were added to the dataset through the use of the Python library `pandas_ta`. Technical indicators are statistical calculations or models which can be applied to financial markets to provide insight on price trends, market patterns, and trading activity, and are largely derived from stock price and trading volume. The first technical indicator added to the dataset

was the relative strength index (RSI). RSI is a momentum oscillator which measures the magnitude of recent price changes to identify overbought or oversold conditions on a financial asset. As RSI is a bounded oscillator, values range from 0 to 100 and is calculated from the average gain and loss of an asset over a time period, for this project a time period of 15 days was used. In addition to RSI, three forms of exponential moving average (EMAx) were introduced to the dataset, these being EMA with fixed timeframe (EMAF), EMA with adaptive multiplier (EMAM), and EMA with support and resistance (EMAS). Moving averages are used to smooth out price fluctuations and to provide a better view of underlying price trends for an asset. EMAF is calculated using average asset price over a specific period, with greater weight given to more recent prices. EMAM adjusts weighting based on the volatility of the asset, with asset price during periods of high volatility having a greater weight than during low volatility. EMAS accounts for price levels where sellers tend to enter (support) or exit (resist) the market, helping to identify areas where price trends may reverse. As all of these technical indicators are calculated over a period of time, their inclusion in the dataset resulted in several null values for these indicators at the start of the dataset, where enough time had not accrued to calculate the indicator. There were 150 days worth of null values which were dropped from the dataset so that no day of trading contained null values.

Following the inclusion of technical indicators within the dataset, the next step in the exploratory data analysis was selecting a target feature for prediction. Given the historical trading data for a stock up until a certain day, the most obvious choice for a target is the average price for the next day of trading. With this in mind, an average price was calculated for each trading day, by averaging the intra-day high and low prices. This average was then shifted one day in the dataset, such that for each day, the next day's average was available as a target.

However, there are several limitations with this approach. Of note is the possible significant changes which can happen in a single day of trading. Additionally, as the stock market trends up over time, a price gain of 10 dollars early in the dataset when the average price is less than 100 dollars is far more significant than a gain of 10 dollars later in the dataset when the average price is well over 300 dollars. Considering this, the daily percent change was identified as another feature for extraction. Using the current day's average price and the next day's average price, the price change between the two days was calculated as a percentage of the first day's price. This percentage will be impacted significantly less by inflation and thus could prove to be a better target feature over the course of the dataset. Model selection then proceeded with both features identified above as prospective targets.

Model Selection

When attempting to forecast for a time series such as the stock market, there are many different models available to choose from. Models chosen for this project include recurrent neural networks (RNN) in the form of long short-term memory (LSTM), linear regression, support vector machine (SVM), decision tree, gradient boosting, and random forest. In this project, models were provided only with historical stock market trading data on one index stock, S&P 500. With this in mind, none of the models selected for the project will be able to account for any external factors such as economic indicators, geopolitical events, and overall market sentiment, which affect market movements and trends.

Long Short-Term Memory Model

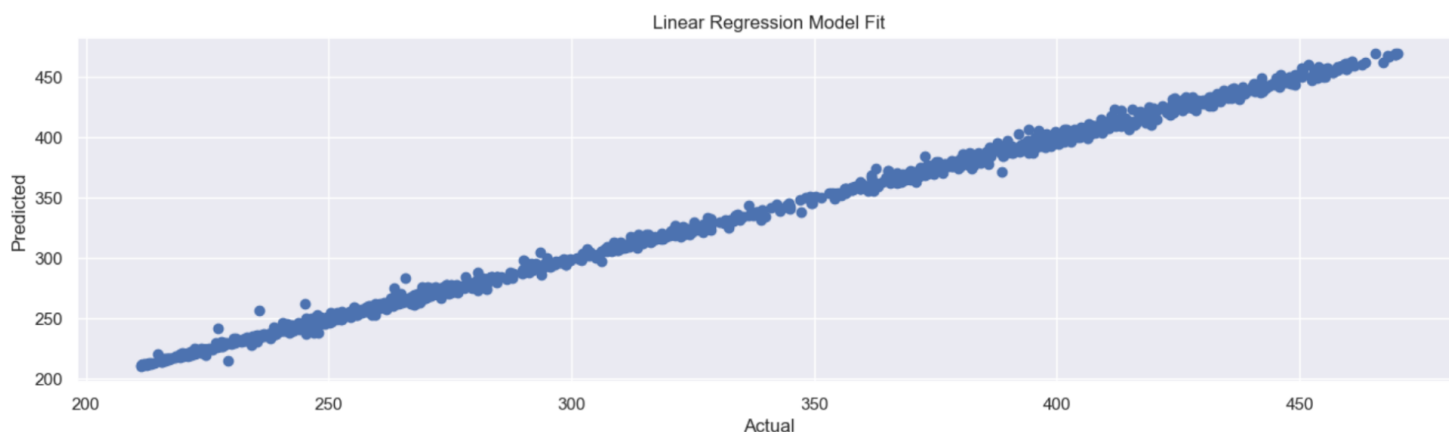
LSTM was selected for investigation in this project as the model is able to capture temporal information of past market movements and patterns. While LSTMs are generally not able to account for external factors such as those discussed above, none of the models in this

project were provided insight to external factors, and the project instead focused on predicting patterns and trends in the stock market from raw trading data. As such, LSTM is an appropriate model for this project, with its downsides being minimized due to the scope of the project. The metrics we used to evaluate our model are mean squared error (MSE), mean average error (MAE), and root mean squared error (RMSE). MSE is the average squared difference between the actual and predicted values, representing how close our regression line is to the dataset. The lower our MSE values, the better our model is at forecasting and making predictions (Orsel & Yamada, pg. 3). MAE is the average absolute error between the actual and predicted values. Similar to MSE, the closer our MAE values are to zero, the more accurate our model is. Lastly, RMSE tells us the square root of the average squared difference between the actual and predicted values, with a lower score representing how well the model fits the dataset (Orsel & Yamada, pg. 3). The primary difference between RMSE and MSE metrics is that MSE will penalize larger errors more severely. Looking at our metrics for the LSTM model, we observed a high MSE score of 58.72, with MAE and RMSE scores of 6.09 and 7.66 respectively. Overall, this was our second best performing model after linear regression, which will be discussed in the next section. Even with the high MSE score, this LSTM model performed relatively better with the MSE and RMSE scores than the majority of models in this report.



Linear Regression Model

Linear regression was also chosen for investigation as a potential model for the stock market. The choice of linear regression was based on the ability of the model to identify a linear relationship between the set of input variables and output variable, the predicted next day average price. The primary limitation of linear regression is that the model is not able to account for nonlinear relationships between the input variables and the output variable, potentially limiting its accuracy in forecasting complex stock market trends and patterns. However, by including several commonly used stock market indicators (RSI, EMAF, EMAM, and EMAS) within the input variables, this otherwise significant limitation of linear regression was minimized for this forecasting model. Linear regression is also limited in that it assumes that the relationship between input and output variables is constant over time, which of course is not the case for time series forecasting. In testing our linear regression model, we found this model fit our data the best and performed very well with a MSE score of 7.87, MAE score of 1.85 and a RMSE score of 2.80. This model outperformed the previous LSTM model in every metric and showed significant gains primarily with MSE, which was reduced from 58.72 in the previous LSTM model.



Support Vector Regression Model

In this section of our report, we tested several other models such as support vector regression, decision tree regression, gradient boosting regression and random forest regression to evaluate whether any alternative models could outperform LSTM or linear regression. Starting with support vector regression, this model was chosen as support vector machines (SVM) provide a supervised learning model that can analyze data for regression analysis (Vazirani et al., pg 2). Although support vector regression (SVR) is optimal for nonlinear relationships, we chose this model for its ability to handle outliers, noise, and overfitting. We used the linear SVR kernel and it performed better than the decision tree, gradient boosting and random forest alternative models. The R-squared score of .61 represents a moderate effective size, however the MSE value of 2193.54, MAE value of 42.99 and RMSE value of 46.84, suggest that the model performed slightly worse than LSTM, and significantly worse than linear regression.



Decision Tree, Gradient Boosting, and Random Forest Regression Models

We tested decision tree, gradient boosting, and random forest regression algorithms to verify whether they could outperform the previous LSTM, LR and SVR models. These models

share similarities with SVR as they are supervised machine learning algorithms that can solve both regression and classification tasks (Vazirani et al., pg 2). Starting with decision tree regression, the advantages are that it is non-parametric and makes no assumptions about distributions and it is not influenced by outliers.. However, we did not expect this model to perform well as stock prices are continuous variables which are typically not utilized well in decision tree regression models. Next, we aimed to include gradient boosting regression as unlike our decision tree model, gradient boosting regression is better at handling continuous variables. Lastly, random forest regression was chosen for its ability to handle noisy data and its similar framework to the decision tree regression model. Overall, these three models performed poorly with high MSE, MAE and RMSE values, with MSE values exceeding 16,000 and MAE/RMSE values exceeding 100. These metrics suggest that linear regression still performs the best, followed by LSTM and support vector regression. These algorithms were worth including for comparison purposes but we were unable to achieve results that outperformed our previous models.

Conclusion

The goal of this project was to analyze 30 years of S&P 500 data and apply multiple machine learning algorithms to find the best model in predicting prices. In total, we employed six different models including LSTM, linear regression, support vector regression, decision tree regression, gradient boosting regression and random forest regression. This process utilized time-series analysis, regression, and neural networks to obtain metrics that assisted us in evaluating each model's performance. These metrics were mean squared error, mean average error, and root mean squared error. Out of the six models we tested, linear regression performed the best, followed by LSTM and support vector regression. We believe the low MSE, MAE and

RMSE values for linear regression were observed due to our dataset being relatively simple and did not include any nonlinear or complex relationships between variables. Had our dataset included other variables such as market sentiment, we believe the five other models could have performed better under a more complex model. Overall, the six models we tested gave a comprehensive overview of predicting S&P 500 prices and provided a robust starting point for future, more complex analyses.

Recommendations

Several recommendations can be implemented to improve the accuracy and performance of the models covered in this project. First, alternative data sources such as news and social media can give us an insight into market sentiment, which could uncover patterns that explain why the price of the S&P 500 increased or decreased during a given time period. Additionally, we could include hybrid models where two algorithms are appended one after the other (Vazirani et al., pg 3). As shown in the research article, *Analysis of various machine learning algorithm and hybrid model for stock market prediction using python*, their linear regression model performed the best and the results were similar to our study, however they took an additional step by including linear regression in each hybrid model and tested various algorithms such as LR + support vector, LR + decision tree and LR + random forest (Vazirani et al., pg 4). We believe that this approach would improve the accuracy of our model. Lastly, we can implement cross-validation techniques such as K-fold cross-validation. This could be used to test the model on different sets of data, which helps identify whether the model is generalizing well and assists in selecting the best hyperparameters. By implementing these

recommendations, the project can achieve better accuracy and performance and provide more reliable predictions of future trends in the stock market.

References:

Orsel, O., Yamada, S. (2022) *Comparative Study of Machine Learning Models for Stock Price Prediction*. Department of Electrical & Computer Engineering, University of Illinois Urbana-Champaign. 1-6.

<https://doi.org/10.48550/arXiv.2202.03156>

Vazirani, S., Sharma, A., Sharma, P. (2020) *Analysis of various machine learning algorithm and hybrid model for stock market prediction using python*. Dept. of ECE, Amity School of Engineering & Technology, Amity University. 203-207.

<https://doi.org/10.1109/icstcee49637.2020.9276859>

Vij, A., Saxena K., Rana, A. (2021) *Prediction in Stock Price Using of Python and Machine Learning*. AIIT, Amity University Uttar Pradesh. 1-4.

<https://doi.org/10.1109/icrito51393.2021.9596513>

Appendix:

Aria Alaghemand -

Written report: Project proposal, introduction, recommendations, Presentation Slides

Alec Anderson -

Written report: Introduction, Model selection, SVM, Decision tree, Gradient boosting, random forest, recommendations, conclusion, Presentation slides

Jupyter notebook: Linear regression, SVM, Decision tree, random forest, gradient boosting, data visualization

Reed Oken -

Written report: Introduction, Data cleaning and processing, Exploratory data analysis, Model selection, LSTM, Linear regression. Slides for presentation

Jupyter notebook: Data cleaning, processing, exploratory data analysis, model tuning, target feature identification LSTM, data visualization