

Linear Forecasting: Predictive Stock Market Trends

Rohan Shah and Joshua Kannagala

Abstract This research explores the application of machine learning techniques to predict stock market prices, integrating historical stock data with sentiment analysis. We examine various machine learning models, including Long Short-Term Memory (LSTM) networks and Principal Component Analysis (PCA), to assess their effectiveness in predicting stock trends. After careful evaluation, sentiment analysis using social media and news data was incorporated to capture market sentiment, providing a more comprehensive approach to stock price prediction. The model developed in this study demonstrates a 3.2% improvement in accuracy when sentiment analysis is included, as measured by the reduction in Mean Absolute Error (MAE). Despite the model's enhanced accuracy, challenges such as the limited availability of historical sentiment data and the potential for overfitting are acknowledged. The findings suggest that combining quantitative financial data with qualitative sentiment insights can lead to more accurate stock price predictions. This research contributes to the growing body of work focused on improving predictive models in finance by incorporating both traditional and sentiment-based analysis techniques.

Rohan Shah
University of Illinois, Champaign, IL, e-mail: rohan.m.s1211@gmail.com

Joshua Kannagala
University of Illinois, Champaign, IL e-mail: joshuakann23@gmail.com

1 Introduction

Predicting stock prices is a complex and challenging task, yet it remains such a fundamental objective for investors and financial analysts aiming to make informed decisions. The motivation for this project arises from the constantly changing nature of the stock market, influenced by economic indicators, market sentiment, and global events. In a world where financial markets can turn on a dime, harnessing the power of machine learning presents an exciting opportunity to enhance the accuracy of stock price predictions [1].

This research explores the potential of machine learning techniques to predict stock prices, providing a thorough analysis of various models and their performance. By combining historical stock data with sentiment analysis, the study aims to uncover patterns and insights that can improve prediction models.

2 Related Work and Background

Previous research on stock price prediction has implemented many different machine learning methods, including LSTM models, PCA for reducing data complexity, and sentiment analysis. We combined the best ideas and positive outcomes of these approaches to build a model that uses both past data and sentiment scores to improve prediction accuracy.

2.1 Stock Prediction with High Accuracy using Machine Learning Techniques

In this article, authors Bansal, Goyal, and Choudhary examine and compare the use of five different machine learning algorithms to accurately predict stock market trends in the Indian Stock Market. These five models include Linear Regression, Support Vector Machines (SVM), and Long Short-Term Memory (LSTM) networks, K-Nearest Neighbors, and Decision Tree Regression networks. Through a comparative analysis of the performances of each method, the study concludes that deep learning algorithms proved to be the most accurate [2].

2.2 Stock Price Prediction using LSTM and its Implementation

In this article, Siddharth M explores how accurate Long Short-Term Memory (LSTM) networks are in analyzing complex patterns in historical stock price data. The article provides a detailed tutorial on using LSTM networks for stock price prediction, covering the theoretical background of LSTM, the challenges of traditional methods, the architecture of LSTM, and a practical implementation using Python. By leveraging LSTM's ability to capture long-term dependencies, the study demonstrates how LSTM networks can effectively analyze time-series data and predict stock prices [3].

2.3 Stock Price Prediction using Principal Components (PCA)

In this article, authors Ghorbani and Chong investigate the application of Principal Component Analysis (PCA) in predicting stock market prices. The study uses PCA to reduce dataset dimensionality, focusing on principal components that capture the most variance. The authors apply PCA to historical stock data and use these components in machine learning models to predict future prices. Their findings demonstrated that PCA enhances prediction accuracy by isolating key patterns and

trends, making the models more robust compared to those without PCA integration [4].

2.4 Stock trend prediction using sentiment analysis

In this article, Xiao and Ihnaini introduce a novel weighted sentiment index for predicting stock trends using Twitter and news data, analyzed with VADER and FinBERT, respectively. They compare the predictive performance of tweets and news across different time divisions, employing algorithms like K-Nearest Neighbors, Decision Trees, SVM, Random Forests, Naïve Bayes, and Logistic Regression. The study finds that Naïve Bayes consistently delivers the best accuracy, particularly with skewed datasets, demonstrating that combining news and tweet sentiment scores improves prediction accuracy [5].

3 Data and Methodologies

To build a comprehensive stock price prediction model, various tools and methods were utilized to ensure accuracy and relevance. Google Cloud APIs played an integral role in efficiently retrieving data using Python. This allowed large volumes of information from historical stock data sources to be automated onto spreadsheets. Yahoo Finance served as the primary source for historical and real-time stock data, providing extensive and reliable market information. Additionally, sentiment scores and data were gathered using VADER and BeautifulSoup to incorporate market sentiment into the model, enhancing its predictive power by combining both quantitative financial data and qualitative insights [6].

3.1 Data Collection

After gathering historical stock data from many companies over the past ten years, the data was organized in Google Sheets. The spreadsheets focused on technology companies, including major companies like Apple, Microsoft, and Google. Following the collection of basic information from Yahoo Finance, additional metrics were calculated using Python and subsequently stored in the spreadsheets for further analysis with machine learning.

3.1.1 Company Historical Stock Price Data

| Date | Open | High | Low | Close | Volume | Dividends | Stock Splits | Symbol | SMA_20 | EMA_20 | Volatility_20 | BB_upper | BB_lower |
|---------------------|-------------|---------------|---------------|---------------|-----------|-----------|--------------|--------|---------------|---------------|---------------|---------------|---------------|
| 2024-08-13 00:02:19 | 0.099994509 | 221.889990308 | 219.009994506 | 221.270004272 | 44155300 | 0.0 | 0.0 | AAPL | nan | 221.270004272 | nan | nan | nan |
| 2024-08-12 00:02:16 | 0.070007324 | 219.509994506 | 215.600006103 | 217.529998776 | 38028100 | 0.25 | 0.0 | AAPL | nan | 220.913813273 | nan | nan | nan |
| 2024-08-09 00:02:11 | 0.854702451 | 216.520374480 | 211.724037866 | 215.900005493 | 422016000 | 0.0 | 0.0 | AAPL | nan | 220.444879198 | nan | nan | nan |
| 2024-08-08 00:02:12 | 0.863619282 | 213.952355450 | 208.588588707 | 213.063385009 | 47161100 | 0.0 | 0.0 | AAPL | nan | 219.741879752 | nan | nan | nan |
| 2024-08-07 00:02:08 | 0.860766653 | 213.362998602 | 206.151380776 | 209.374230951 | 62519400 | 0.0 | 0.0 | AAPL | nan | 218.773036261 | nan | nan | nan |
| 2024-08-06 00:02:05 | 0.826409492 | 209.747220941 | 200.837535824 | 206.904002216 | 69669500 | 0.0 | 0.0 | AAPL | nan | 217.651604447 | nan | nan | nan |
| 2024-08-05 00:02:18 | 0.858622302 | 213.253166215 | 195.773384902 | 209.028000913 | 119548600 | 0.0 | 0.0 | AAPL | nan | 216.830314587 | nan | nan | nan |
| 2024-08-02 00:02:18 | 0.896619619 | 225.339173497 | 217.458296314 | 219.605804433 | 105568600 | 0.0 | 0.0 | AAPL | nan | 217.094648954 | nan | nan | nan |
| 2024-08-01 00:02:24 | 1.105900033 | 224.220463437 | 216.769088842 | 218.107543945 | 62501000 | 0.0 | 0.0 | AAPL | nan | 217.191113334 | nan | nan | nan |
| 2024-07-31 00:02:21 | 1.038607739 | 223.561235668 | 220.374921605 | 221.933242187 | 50030300 | 0.0 | 0.0 | AAPL | nan | 217.632268463 | nan | nan | nan |
| 2024-07-30 00:02:18 | 0.906591348 | 220.075272755 | 215.870133336 | 218.547042846 | 41643800 | 0.0 | 0.0 | AAPL | nan | 217.716380633 | nan | nan | nan |
| 2024-07-29 00:02:16 | 0.709167255 | 219.040458194 | 215.500559497 | 217.967686157 | 36311800 | 0.0 | 0.0 | AAPL | nan | 217.744941864 | nan | nan | nan |
| 2024-07-26 00:02:18 | 0.447142492 | 219.236237652 | 215.760250151 | 217.708007812 | 41601300 | 0.0 | 0.0 | AAPL | nan | 217.741424335 | nan | nan | nan |
| 2024-07-25 00:02:18 | 0.676878259 | 220.594871876 | 214.371863876 | 217.238555905 | 51391200 | 0.0 | 0.0 | AAPL | nan | 217.60532104 | nan | nan | nan |
| 2024-07-24 00:02:23 | 1.01016993 | 224.540305102 | 216.87864786 | 218.287322598 | 61777600 | 0.0 | 0.0 | AAPL | nan | 217.750083617 | nan | nan | nan |
| 2024-07-23 00:02:24 | 1.10587962 | 226.677623652 | 222.422539431 | 224.749847412 | 39960300 | 0.0 | 0.0 | AAPL | nan | 218.416727788 | nan | nan | nan |
| 2024-07-22 00:02:28 | 0.74541939 | 227.519655988 | 222.832075789 | 223.701080322 | 48201800 | 0.0 | 0.0 | AAPL | nan | 218.919994581 | nan | nan | nan |
| 2024-07-19 00:02:24 | 0.60076292 | 226.537785626 | 223.021851244 | 224.050659178 | 49151500 | 0.0 | 0.0 | AAPL | nan | 219.408633717 | nan | nan | nan |
| 2024-07-18 00:02:20 | 0.137604653 | 236.173579137 | 222.013028715 | 223.820868841 | 186034600 | 0.0 | 0.0 | AAPL | nan | 219.838384495 | nan | nan | nan |
| 2024-07-17 00:02:29 | 0.84700036 | 231.192405643 | 226.377971235 | 228.815389602 | 57345900 | 0.0 | 0.0 | AAPL | 218.389613422 | 220.674271397 | 5.57062238010 | 232.530858102 | 207.248385829 |
| 2024-07-16 00:02:34 | 0.728307528 | 235.996843506 | 232.061396246 | 234.548522949 | 43234300 | 0.0 | 0.0 | AAPL | 219.053538276 | 221.995628687 | 6.62373163705 | 232.301002550 | 205.806078002 |
| 2024-07-15 00:02:36 | 0.206594889 | 236.955727793 | 232.820514767 | 234.128997802 | 62631300 | 0.0 | 0.0 | AAPL | 219.883486272 | 223.151187651 | 7.41539838462 | 234.714265996 | 205.052892457 |
| 2024-07-12 00:02:28 | 0.653309112 | 232.371030924 | 228.415602801 | 230.273452780 | 63046500 | 0.0 | 0.0 | AAPL | 220.587961588 | 223.629498813 | 7.7029433031 | 236.003102012 | 205.191812979 |
| 2024-07-11 00:02:21 | 1.22475568 | 232.12319408 | 225.508978093 | 227.308000244 | 64710900 | 0.0 | 0.0 | AAPL | 221.309637341 | 224.160678706 | 7.62775570587 | 236.565348752 | 206.054225930 |

Fig. 1 AAPL Stock Data Metrics

3.1.2 Sentimental Scores

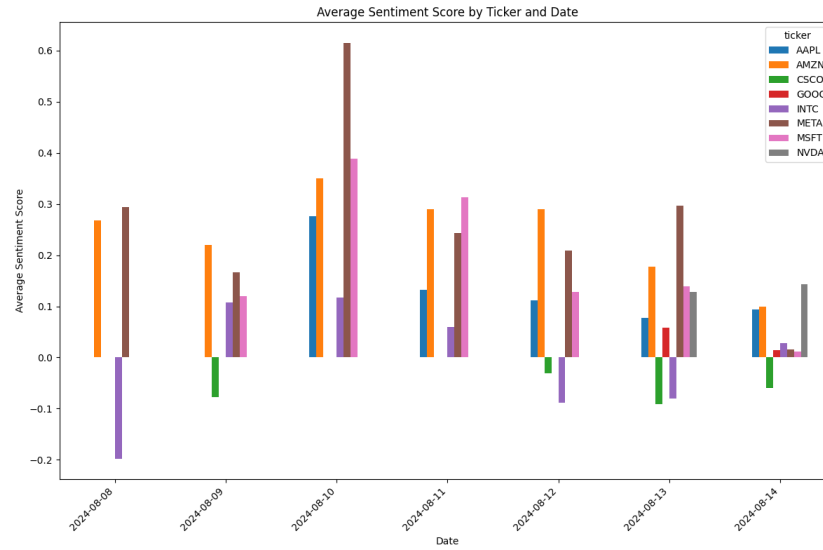


Fig. 2 Average Sentiment Scores of Several Technology Companies in a week

3.2 Methods

3.2.1 Visual Studio Code

Visual Studio Code served as the main coding platform due to its flexibility and extensive support for Python libraries. As a result, Python was also chosen as the main programming language, as many of these libraries were needed for machine learning. Key libraries included:

- **Pandas:** Used for data manipulation and analysis, particularly for organizing stock price data into data frames for easy computation.
- **Yfinance:** A Python library providing easy access to historical market data from Yahoo Finance, allowing the retrieval of stock prices, volume, and other important financial metrics.
- **Scikit-learn (Sklearn):** Employed for implementing machine learning models, particularly for preprocessing data, training models, and evaluating their performance.

3.2.2 Linear Regression

Linear regression is a statistical method used to model the relationship between a dependent variable (in this case, stock prices) and one or more independent variables (such as historical prices and sentiment scores). This approach was employed to predict future stock prices by identifying the linear relationship between these variables. The model calculates a best-fit line that minimizes the difference between predicted and actual stock prices [7].

3.2.3 Mean Absolute Error

Mean Absolute Error (MAE) is a metric used to evaluate the accuracy of regression models. It measures the average magnitude of errors between predicted and actual values, providing an indication of how close the predictions are to the real stock prices. In this study, MAE was utilized to assess the performance of the linear regression model, serving as a key metric to determine the model's accuracy and effectiveness [8].

$$\text{MAE} = \frac{1}{n} \sum_{i=1}^n |x_i - \hat{x}|$$

Fig. 3 Mean Absolute Error Formula

3.2.4 Sentiment Methods

To incorporate sentiment analysis into the stock price predictions, several tools and methodologies were utilized:

- **VADER Sentiment Analysis:** Used to analyze the sentiment of text data from social media posts and news articles. VADER is particularly well-suited for analyzing short, informal text from sources like Twitter, making it ideal for this project.
- **BeautifulSoup:** A Python library used for web scraping, which allows the collection of news articles from various financial news websites. This helped create a comprehensive dataset of news sentiments to be analyzed alongside stock prices.
- **Finviz:** An online stock screener that was used as the main source of additional sentiment data, news articles, and financial statements.

4 Results

To evaluate whether sentiment analysis had an affect on the model's accuracy, the mean absolute error was calculated on both the training and testing data sets for the eight technology stocks. Through the tables and column charts below, a comparison can be drawn between the MAE of the training and testing sets with and without sentiment analysis.

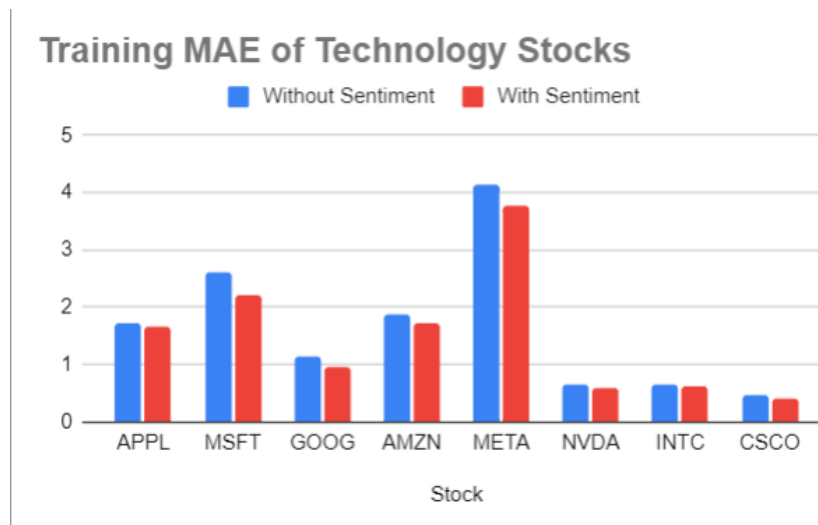
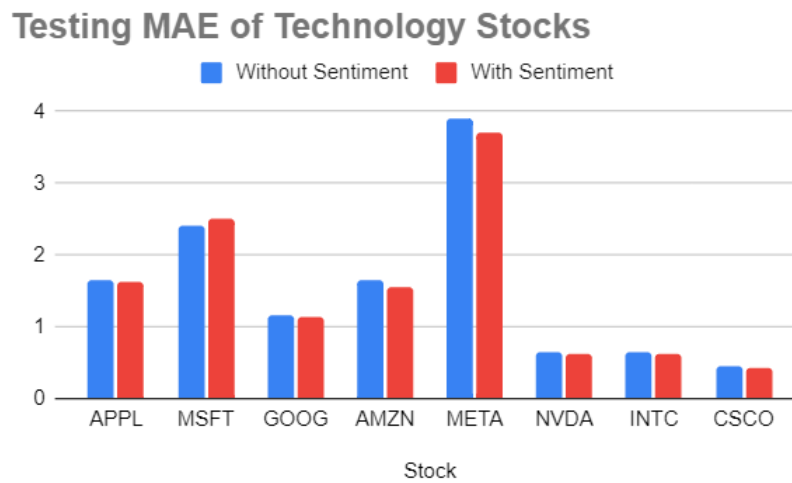
4.1 Model Evaluation

| Stock | Training Dataset (MAE) | Testing Dataset (MAE) |
|-------|------------------------|-----------------------|
| APPL | 1.72 | 1.64 |
| MSFT | 2.6 | 2.4 |
| GOOG | 1.15 | 1.15 |
| AMZN | 1.88 | 1.65 |
| META | 4.15 | 3.9 |
| NVDA | 0.66 | 0.64 |
| INTC | 0.63 | 0.63 |
| CSCO | 0.46 | 0.45 |

Fig. 4 MAE without Sentiment Analysis

| Stock | Training Dataset (MAE) | Testing Dataset (MAE) |
|-------|------------------------|-----------------------|
| APPL | 1.65 | 1.63 |
| MSFT | 2.2 | 2.5 |
| GOOG | 0.95 | 1.12 |
| AMZN | 1.73 | 1.54 |
| META | 3.78 | 3.71 |
| NVDA | 0.59 | 0.62 |
| INTC | 0.6 | 0.62 |
| CSCO | 0.41 | 0.43 |

Fig. 5 MAE with Sentiment Analysis

**Fig. 6** MAE in Training Set**Fig. 7** MAE in Testing Set

Through observing the results, the integration of sentiment analysis led to a 3.2% reduction in Mean Absolute Error, enhancing the model's accuracy. Despite the challenges posed by the limited historical sentiment data, the inclusion of sentiment metrics provided valuable insights, particularly for stocks with high public engagement.

5 Conclusion

The research presented demonstrates that the integration of machine learning techniques with sentiment analysis can enhance the predictive accuracy of stock prices. The model developed showed a Mean Absolute Error reduction of 3.2% when sentiment analysis was utilized, indicating a concrete improvement in prediction accuracy. However, it was observed that the model performed better on the training set than on the testing set, suggesting potential overfitting. This discrepancy points to the need for further refinement in the model to ensure that it generalizes well to unseen data.

5.1 Applications

The methodology and findings from this research have several potential applications in finance. For example, investment firms can utilize the enhanced model to make more informed trading decisions by integrating real-time sentiment data with historical financial metrics [9]. Additionally, the approach could be adapted for use in predicting price movements in other financial markets, such as commodities or foreign exchange. The integration of sentiment analysis can also be extended to other domains where understanding the influence of public sentiment is critical, such as in consumer behavior analysis or political forecasting.

5.2 Concerns and Next Steps

Several concerns were identified during the research process. The limited availability of historical sentiment data restricted the analysis primarily to recent trends, which may have affected the model's overall performance. This is due to the fact that the model relied on historical stock price data from the past ten years, and the sentiment data was not available for every single one of those dates. Additionally, the model was only trained on stocks from the technology sector, which are inherently volatile and may not represent broader market dynamics [10]. Future work should focus on expanding the dataset to include a wider range of industries and a more extensive historical sentiment dataset. This would help limit overfitting and improve the model's generalization to unseen data. Additionally, exploring more sophisticated machine learning algorithms and refining sentiment analysis techniques could be a next step in enhancing the model's robustness and accuracy.

Acknowledgements We would like to thank Professor Christof Teuscher and the altREU faculty mentors for their helpful guidance and feedback throughout this research project.

References

- [1] Y. Chen, S. H. Sim, S. Yan, and M. Dong, "Stock market prediction using artificial neural networks and sentiment analysis," *IEEE International Conference on Systems, Man, and Cybernetics*, pp. 1061-1066, 2013, doi: 10.1109/SMC.2013.187.
- [2] T. M. Dang et al., "A hybrid deep learning model for stock market prediction," *Procedia Computer Science*, vol. 207, pp. 1230-1242, 2022, doi: 10.1016/j.procs.2022.09.121.
- [3] M. S. (2024, July 29). Stock price prediction using LSTM and its implementation. *Analytics Vidhya*.
- [4] Ghorbani, M., & Chong, E. K. P. (2020, March 20). Stock price prediction using principal components. *National Institutes of Health*.
- [5] Xiao, Q., & Ihnaini, B. (2023, March 20). Stock trend prediction using sentiment analysis. *National Institutes of Health*.
- [6] Cristescu, M. P., Nerisanu, R. A., Mara, D. A., & Oprea, S.-V. (2022, November 14). Using market news sentiment analysis for stock market prediction. *MDPI*.
- [7] Mo, H. (2023, December). Comparative analysis of linear regression, polynomial regression, and ARIMA model for short-term stock price forecasting. *ResearchGate*.
- [8] Hodson, T. O. (2022, July 19). Root-mean-square error (RMSE) or mean absolute error (MAE): When to use them or not. *Geoscientific Model Development*.
- [9] Chen, Y., Zhao, H., Li, Z., & Lu, J. (2020, December 4). A dynamic analysis of the relationship between investor sentiment and stock market realized volatility: Evidence from China. *National Institutes of Health*.
- [10] Rasiova, B., & Arendas, P. (2022, December 1). Copula approach to market volatility and technology stocks dependence. *Finance Research Letters*.