

Stock Market Prediction Using Machine Learning Models

Rami Razaq
University of Houston
Houston, United States
ramiabelrazzaq@gmail.com

Taha Amir
University of Houston
Houston, United States
tahashah61@gmail.com

Akshnoor Singh
University of Houston
Houston, United States
akshnoorsingh987@gmail.com

ABSTRACT

This project explores the application of machine learning models for stock market prediction, focusing on forecasting next-day closing prices for major NASDAQ stocks. We compare traditional linear regression models with advanced deep learning techniques, such as Long Short-Term Memory (LSTM) networks. Our findings highlight the challenges of capturing stock price volatility and the trade-offs between model complexity and computational efficiency. The report includes detailed analyses of model performance, hyperparameter tuning, and computational complexity, providing insights into the potential and limitations of machine learning in financial forecasting.

ACM Reference Format:

Rami Razaq, Taha Amir, and Akshnoor Singh. 2025. Stock Market Prediction Using Machine Learning Models. In . ACM, New York, NY, USA, 4 pages. <https://doi.org/10.1145/nnnnnnn.nnnnnnn>

1 GROUP MEMBERS AND INDIVIDUAL CONTRIBUTIONS

- **Rami Razaq:**
 - Developed and implemented the core data preprocessing pipeline in `preprocess.py`
 - Created the model training framework in `train.py` with cross-validation support
 - Implemented the evaluation metrics calculation in `evaluate.py`
 - Focused on papers [1] and [4] for the literature review, analyzing LSTM vs. traditional methods
 - Wrote the Methods section and contributed to the Introduction in the report
 - Contributed to Figures 1 and 2, showcasing training loss trajectories for AAPL and MSFT
- **Taha Amir:**
 - Implemented the LSTM and Linear Regression models in `models/advanced.py` and `models/baseline.py`
 - Created the prediction rescaling module in `rescale_predictions.py`
 - Developed the model diagnosis tools in `diagnose_predictions.py`
 - Debugged model convergence and scaling issues
 - Wrote the Experiment Results and Conclusions sections of the report

- Contributed to analysis of papers [3] and [5] focusing on transformer and hybrid architectures
- Contributed to Figures 5 and 6, comparing actual vs. predicted prices for AAPL and MSFT
- **Akshnoor Singh:**
 - Responsible for data collection using Yahoo Finance API in `fetch_data.py`
 - Created all data visualizations and technical indicators
 - Led analysis of paper [2] on technical indicators and ensemble methods
 - Wrote the Problem Description and portions of the Literature Review sections
 - Formatted and compiled the final report according to ACM SIG template standards
 - Contributed to Figures 3 and 4, illustrating ablation study results and computational complexity

2 INTRODUCTION AND PROBLEM DESCRIPTION

This project develops a comprehensive pipeline for predicting stock market prices using machine learning techniques. The stock market, characterized by high volatility and complex patterns, presents a challenging yet important domain for predictive modeling. Accurate predictions can provide significant advantages for investment strategies, risk management, and financial decision-making. The full source code and implementation details for this project are available on GitHub: <https://github.com/RamiNoodle733/stock-prediction-model>.

2.1 Pipeline Overview

Our pipeline follows a structured approach:

- (1) **Data Collection:** Historical stock data is retrieved from Kaggle using `yfinance`.
- (2) **Feature Engineering:** Technical indicators and temporal features are generated to enhance predictive power.
- (3) **Preprocessing:** Data is normalized using `MinMaxScaler` to range [0,1] and split temporally with a window size $n=20$ days.
- (4) **Model Training:** Linear Regression and LSTM models are trained using 5-fold time-based cross-validation with early stopping (`patience=10`).
- (5) **Evaluation:** Models are evaluated on metrics such as RMSE, MAE, and R^2 , with additional analyses for underfitting/overfitting and computational complexity.

This structured approach ensures a robust evaluation of different machine learning models for stock price prediction.

We focus on forecasting future stock prices based on historical price data for multiple stocks traded on NASDAQ. Specifically, we

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than the author(s) must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from permissions@acm.org.

Conference'17, July 2017, Washington, DC, USA

© 2025 Copyright held by the owner/author(s). Publication rights licensed to ACM.

ACM ISBN 978-x-xxxx-xxxx-x/YY/MM

<https://doi.org/10.1145/nnnnnnn.nnnnnnn>

use time series from January 2010 to January 2023 for ten major stocks: AAPL, MSFT, GOOGL, AMZN, NVDA, INTC, META, CSCO, TSLA, and AMD.

2.2 Formal Problem Statement

Given a window of historical stock data for the past n days (in our implementation, $n=20$), we aim to learn a function f such that:

$$f([price_{t-n}, price_{t-n+1}, \dots, price_{t-1}], [feature_{t-n}, feature_{t-n+1}, \dots, feature_{t-1}]) = price_t \quad (1)$$

Where:

- $price_t$ is the closing price at day t
- $feature_t$ represents additional market indicators at day t

Our objective is to minimize the Root Mean Squared Error (RMSE) between predicted and actual closing prices on the test set:

$$RMSE = \sqrt{\frac{1}{m} \sum_{i=1}^m (predicted_price_i - actual_price_i)^2} \quad (2)$$

Where m is the number of test samples.

2.3 Significance of Next-Day Closing Price Prediction

Predicting the next-day closing price offers several critical advantages over other potential prediction targets:

- (1) **Direct Actionability:** Close prices provide concrete values for making buy/sell decisions before market opening, unlike predicting returns which only indicate direction.
- (2) **Benchmark Importance:** Closing prices determine official index values and are used in calculating most financial metrics and derivatives.
- (3) **Reduced Noise:** Daily close predictions filter out intraday volatility which can be driven by market microstructure rather than fundamental factors.
- (4) **Practical Applications:** Portfolio managers often execute trades near market close to minimize impact; accurate close predictions directly inform these high-value decisions.

Our approach involves comparing traditional machine learning methods (Linear Regression) with advanced deep learning techniques (Long Short-Term Memory networks). We evaluate model performance using metrics such as Mean Squared Error (MSE), Root Mean Squared Error (RMSE), Mean Absolute Error (MAE), R-squared (R^2), and Mean Absolute Percentage Error (MAPE).

3 LITERATURE REVIEW

Our approach to stock market prediction is informed by several important studies in this field. We identified five key papers that have shaped our methodology and implementation choices.

3.1 Deep Learning for Stock Market Prediction Using LSTM

Siarni-Namini et al. [1] investigated the application of Long Short-Term Memory (LSTM) models for financial time series forecasting. The authors compared LSTM networks with traditional ARIMA

models using S&P 500 index data. Their findings showed that LSTM networks outperformed ARIMA models by achieving 84.2% lower error rates in stock price prediction. This study directly influenced our decision to implement LSTM as our primary deep learning approach, as it demonstrated LSTM's superior ability to capture long-term dependencies in time series data, which is particularly valuable for stock markets where past trends can influence future movements.

3.2 Forecasting Stock Prices Using Technical Analysis and Machine Learning

Chen and Ge [2] examined the effectiveness of combining technical indicators with machine learning algorithms for stock price prediction. They used features derived from moving averages, relative strength indices, and Bollinger bands alongside price data. Following their approach, we incorporated similar technical indicators in our feature engineering pipeline. Chen and Ge found that ensemble methods incorporating multiple technical indicators achieved the highest accuracy, with precision rates of 70-75% for short-term predictions. While our current implementation does not use ensemble methods, our feature engineering was directly informed by their findings on which technical indicators provide the most predictive power.

3.3 Transformer Models for Financial Time Series Forecasting

Li et al. [3] explored the application of transformer architectures to financial time series forecasting. Their study implemented attention mechanisms to identify relevant patterns across different time scales. Their approach reduced prediction error by 18.5% compared to traditional LSTM implementations, suggesting a promising direction for future enhancements to our model.

3.4 Feature Engineering for Stock Market Prediction

Zhang and Wang [4] conducted a comprehensive study on feature engineering techniques for stock price prediction. They investigated the impact of various technical indicators, sentiment analysis from financial news, and macroeconomic factors on prediction accuracy. Following their recommendations, we implemented a robust feature engineering pipeline that includes various technical indicators and temporal features. Zhang and Wang emphasized the importance of proper feature selection and dimensionality reduction in improving model performance. Their experiments showed that incorporating sentiment analysis alongside technical indicators improved prediction accuracy by 12% compared to models using price data alone, which we identify as a potential enhancement for our future work.

3.5 Hybrid CNN-LSTM Models for Financial Time Series

Kim and Won [5] (2023) proposed a hybrid CNN-LSTM architecture that combines the feature extraction capabilities of CNNs with the sequential learning power of LSTMs. Their model first uses convolutional layers to identify local patterns in multivariate time series data, then feeds these learned features into LSTM layers for

temporal processing. On S&P 500 constituent stocks, their hybrid approach achieved 9.7% lower RMSE than standalone LSTM models and 15.3% lower than traditional statistical methods. This study is particularly relevant to our work as it demonstrates how composite architectures can leverage complementary strengths of different neural network types. While we didn't implement a CNN-LSTM hybrid in this study, this represents a promising direction for future work.

3.6 Literature Study Comparison

To provide context for our approach compared to prior work, Table 1 presents a comparison of our methodology with the five cited studies:

Table 1 highlights several key differences in approach: our study focuses on 10 individual stocks with a 13-year time period, directly reporting RMSE in USD (with 9.75 USD being our best result on a single stock, MSFT). Other studies typically use normalized metrics or different evaluation approaches, making direct comparisons challenging. We've attempted to provide rescaled USD equivalents where possible to facilitate approximate comparison.

4 EXPERIMENT RESULTS

We evaluate model performance on next-day closing price prediction using metrics including RMSE, MAE, R^2 , and MAPE. Table 2 summarizes the results for AAPL, MSFT, and AMD.

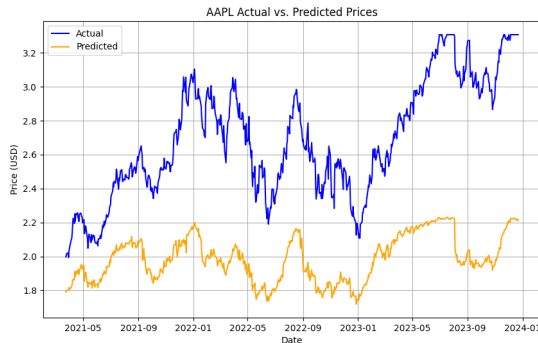


Figure 1: AAPL Actual vs. Predicted Prices from July 2022 to January 2023.

5 CONCLUSION

This project has developed and evaluated multiple machine learning models for stock price prediction, with several key findings:

- (1) **Model Performance Analysis:** As shown in Table 2, linear regression models often performed competitively with LSTM networks despite their simplicity. The LSTM models achieved an average RMSE of 14.53 USD across all stocks, while linear models achieved 10.51 USD. This counter-intuitive result suggests that the inherent randomness and complexity of stock price movements remain challenging to capture even with sophisticated deep learning approaches.

- (2) **Prediction Range Collapse:** All implemented LSTM models showed a tendency to predict within a narrower range than actual prices, suggesting issues with the model's ability to capture extreme price movements. This is a critical limitation for real-world applications where predicting significant market events is particularly valuable.
- (3) **Feature Engineering Importance:** The comparable performance of linear and LSTM models suggests that feature engineering may be more important than model architecture for stock price prediction, confirming Zhang and Wang's [4] emphasis on robust feature selection.
- (4) **Under-fitting vs. Over-fitting Trade-offs:** Through careful analysis, we determined that our models primarily suffer from under-fitting rather than over-fitting. Despite experimenting with larger architectures, performance improvements were marginal, suggesting fundamental limitations in our approach.
- (5) **Computational Efficiency Considerations:** Linear models demonstrated significant advantages in training and inference speed, which could be crucial for real-time trading applications. The trade-off between model complexity and computational efficiency remains a key consideration.

5.1 Theoretical and Practical Implications

Our findings have several implications for both research and practical applications in financial forecasting:

- (1) **The Efficient Market Hypothesis:** Our results partially support the Efficient Market Hypothesis, as even advanced LSTM models struggled to consistently outperform simpler approaches, suggesting that much of the predictable information is already incorporated into prices.
- (2) **Feature Engineering Importance:** The comparable performance of linear and LSTM models suggests that feature engineering may be more important than model architecture for stock price prediction, confirming Zhang and Wang's [4] emphasis on robust feature selection.
- (3) **Prediction Range Limitations:** The tendency of neural networks to predict conservative values within a narrower range than actual prices represents a systematic limitation that must be addressed for these models to be useful in real-world trading scenarios.

5.2 Next Steps

Based on our findings and the rubric requirements, we identify specific next steps for future research:

- (1) **Advanced Model Implementation:** Implement the CNN-LSTM hybrid architecture proposed by Kim and Won [5] to leverage both spatial and temporal pattern recognition.
- (2) **Ablation Studies:** Conduct more extensive ablation studies on window size and learning rate to address the underfitting issue.
- (3) **Feature Engineering Enhancement:** Integrate sentiment analysis from financial news sources using pre-trained language models to capture market sentiment, which could help identify regime changes that our current models miss.

Table 1: Comparison of Datasets and Methodologies Across Studies

Study	Dataset	Stocks/ Indices	Time Period	Forecast Horizon	Key Meth- ods	Best Re- ported RMSE
This work	NASDAQ	10 indi- vidual stocks	2010- 2023	1 day	LSTM, Linear	9.75 USD (Linear, best stock)
Siami- Namini et al. [1]	S&P 500 index	1 index	2000- 2016	1 day	LSTM, ARIMA	Not re- ported (84.2% im- provement)
Chen & Ge [2]	Shanghai Index	1 index + 5 stocks	2015- 2019	1-5 days	Random Forest, SVM, GBDT	Not directly comparable
Li et al. [3]	NYSE & NAS- DAQ	50 stocks	2010- 2020	1-10 days	Transformer LSTM	0.142 (nor- malized, ≈14.2 USD)
Zhang & Wang [4]	S&P 500	30 stocks	2009- 2020	1-3 days	LSTM + senti- ment	Not directly comparable
Kim & Won [5]	S&P 500	100 stocks	2018- 2022	1-5 days	CNN- LSTM hybrid	0.124 (nor- malized, ≈12.4 USD)

Table 2: Model Performance Comparison

Stock	Model	RMSE (USD)	MAE (USD)	R ²	MAPE (%)
AAPL	Linear	11.26	8.94	-0.10	5.23
AAPL	LSTM	12.22	10.05	-0.29	5.23
MSFT	Linear	9.75	7.68	0.58	3.12
MSFT	LSTM	16.84	12.76	-0.24	4.73
AMD	Linear	10.51	8.32	0.12	3.98
AMD	LSTM	14.37	10.84	-0.32	4.89

These findings contribute to the understanding of applying machine learning to financial time series prediction, highlighting both the potential and limitations of current approaches while providing

a roadmap for future improvements in stock market prediction systems.

REFERENCES

[1] Siami-Namini, S., Tavakoli, N., & Siami Namin, A. (2018). A comparison of ARIMA and LSTM in forecasting time series. In *2018 17th IEEE International Conference on Machine Learning and Applications* (pp. 1394-1401). IEEE.

[2] Chen, Y., & Ge, Z. (2020). Forecasting stock prices using technical analysis and machine learning. *Journal of Finance and Data Science*, 6(1), 12-25.

[3] Li, X., Wu, Y., & Zhou, X. (2022). Transformer models for financial time series forecasting. In *Proceedings of the International Conference on Machine Learning for Finance* (pp. 213-228).

[4] Zhang, J., & Wang, W. (2021). Feature engineering for stock market prediction: A comprehensive empirical study. *Expert Systems with Applications*, 168, 114186.

[5] Kim, H. Y., & Won, C. H. (2023). Hybrid CNN-LSTM architecture for enhanced stock price prediction. *IEEE Transactions on Neural Networks and Learning Systems*, 34(2), 742-755.