

Stock Market Prediction

Kushwah Pratiksha, Aditya Kumar Sah, Vishal Kumar, Smilepreet Kaur

Abstract

Stock market prediction has become a prominent application domain for artificial intelligence and big data technologies. This study presents a deep learning-based approach utilizing Long Short-Term Memory (LSTM) networks for forecasting stock prices, using historical Google stock data. The primary goal is to analyze patterns in large-scale time-series data and enhance predictive performance using deep sequential models. The proposed method integrates data preprocessing, feature engineering, and sequential LSTM modeling with performance metrics including MSE, RMSE, R^2 , and MAE. Experimental results demonstrate a strong predictive capability with an R^2 score of 0.9921, highlighting the potential of LSTM in financial forecasting. Future work aims to incorporate sentiment analysis and real-time data streaming for improved volatility handling and decision-making.

Keywords: Stock prediction · Deep learning · LSTM · Financial forecasting · Big data analytics

1 Introduction

In an era where financial decisions are increasingly reliant on rapid and accurate data interpretation, stock market prediction has emerged as a critical area of research. The inherent volatility and complexity of market behavior, influenced by diverse economic, political, and psychological factors, pose a significant challenge for traditional forecasting models. With the proliferation of data generated through financial transactions, news feeds, and social media, the need for scalable and intelligent predictive systems has become more pronounced.

Artificial intelligence, and specifically deep learning, has shown substantial promise in decoding nonlinear patterns embedded in time-series data. Among the architectures available, Long Short-Term Memory (LSTM) networks have proven particularly effective due to their ability to capture long-term dependencies. This research leverages LSTM networks to model and predict the future price movements of Google stock using historical data. By building a robust data processing pipeline and a deep learning framework, the study aims to deliver insights into the feasibility and efficacy of such models in real-world financial forecasting scenarios.

Motivation In today's fast-paced financial landscape, even milliseconds can make the difference between profit and loss. The idea of harnessing deep learning to navigate this complexity is both timely and essential. Our motivation stems from a desire to bridge the gap between raw data and informed investment decisions. By focusing on Google stock data as a case study, this project aims to provide a scalable blueprint for applying AI in financial markets. Moreover, we were driven by the challenge of creating a system that could learn from the past and assist in making forward-looking decisions—something every modern investor craves.

2 Related Work

Traditional approaches to stock forecasting, such as linear regression and ARIMA, fail to capture the nonlinear temporal dependencies in stock data. Machine learning models like Support Vector Machines (SVM) and Random Forest (RF) improved accuracy but lacked memory mechanisms. LSTM networks, capable of retaining long-term dependencies, have outperformed conventional models in financial time-series tasks.

Fischer & Krauss (2018) applied LSTM models on the S&P 500 index and demonstrated significantly improved predictions over shallow networks. Zhang et al. (2022) incorporated sentiment data into LSTM and showed performance gains across multiple stock datasets. Hybrid models like CNN-LSTM and Bi-LSTM have also been tested for noise filtering and temporal smoothing. These studies indicate the evolving trend towards deep, hybrid, and ensemble learning for market prediction.

3 Problem Statement and Objectives

Problem Statement: Stock price movements are driven by a multitude of interdependent variables, making financial time-series forecasting a notoriously complex task. Classical statistical models often struggle with high volatility and nonlinearities inherent in market data. There is a compelling need for advanced models that can learn temporal correlations, adapt to varying market conditions, and deliver consistent performance across diverse scenarios. This research addresses this gap by designing and implementing an LSTM-based model tailored for sequential forecasting of Google stock prices.

Objectives:

1. To acquire and preprocess historical stock data specific to Google Inc. for time-series analysis.
2. To develop a deep learning model using Long Short-Term Memory (LSTM) networks for price prediction.
3. To conduct hyperparameter optimization to fine-tune model performance.
4. To evaluate the model using industry-standard metrics such as MSE, RMSE, MAE, and R-squared.
5. To visualize and interpret the predictive outcomes, highlighting trends, limitations, and areas for future enhancement.

4 Dataset Description

The dataset used consists of daily Google stock price data from 2010 to 2023, obtained via Yahoo Finance API. It includes attributes such as Open, High, Low, Close, Volume, and Adjusted Close. The 'Close' price was selected as the target variable due to its relevance in financial decision-making. The dataset contains 3,300+ data points and covers multiple market cycles, including periods of high volatility such as the COVID-19 crash and subsequent recovery.

5 Methodology

This section delineates the comprehensive methodological framework adopted for implementing the LSTM-based stock prediction model. The approach encompasses data acquisition, preprocessing, model architecture design, training strategy, and performance evaluation.

5.1 Data Acquisition and Preprocessing

Daily historical stock prices of Google (GOOG) were sourced from Yahoo Finance for the period spanning 2010 to 2023. The dataset includes key features such as Open, High, Low, Close, Volume, and Adjusted Close. The target variable selected for prediction was the 'Close' price, owing to its strategic relevance in trading decisions.

To ensure model readiness, the data underwent the following preprocessing steps:

- **Missing Value Handling:** Any missing or anomalous values were addressed through interpolation or elimination.
- **Normalization:** A MinMaxScaler was applied to scale numerical features within the range [0, 1], enhancing model convergence.
- **Sequence Generation:** Time-series windows of 60 days were used as input features to predict the subsequent day's closing price, enabling the LSTM model to learn from historical trends.
- **Data Splitting:** The dataset was split in an 80:20 ratio into training and testing subsets to validate model generalization.

5.2 Model Architecture and Training

A sequential LSTM architecture was employed using the TensorFlow-Keras framework. The model comprises:

- Four stacked LSTM layers, each with 100 memory units.
 - Dropout layers (rate = 0.2) after each LSTM layer to mitigate overfitting.
 - A final Dense layer with a single neuron for regression output.
- The model was compiled using the Mean Squared Error (MSE) loss function and optimized using the Adam optimizer. Training was conducted over 20 epochs with a batch size of 32, determined through grid-based hyperparameter tuning.

5.3 Evaluation Metrics

To assess predictive accuracy and robustness, the model was evaluated using the following performance metrics:

- **Mean Squared Error (MSE):** Measures the average squared difference between actual and predicted values.
- **Root Mean Squared Error (RMSE):** Provides an interpretable error measure in the same units as the stock prices.
- **Mean Absolute Error (MAE):** Reflects the average magnitude of prediction errors.
- **R-squared (R^2):** Indicates the proportion of variance explained by the model, with values closer to 1 denoting better performance.

6 Results and Discussion

The model achieved an MSE of 0.0043, RMSE of 0.0655, MAE of 0.0517, and an R^2 score of 0.9921. These results indicate a high level of predictive accuracy, with the predicted values closely tracking actual stock trends. The training loss curve stabilized within 15 epochs, demonstrating quick convergence.

Limitations:

- The model does not incorporate macroeconomic indicators, news sentiment, or geopolitical events.
- Volatility spikes are harder to capture without external features.
- Risk of overfitting exists despite the use of dropout layers.

Comparison with Other Models:

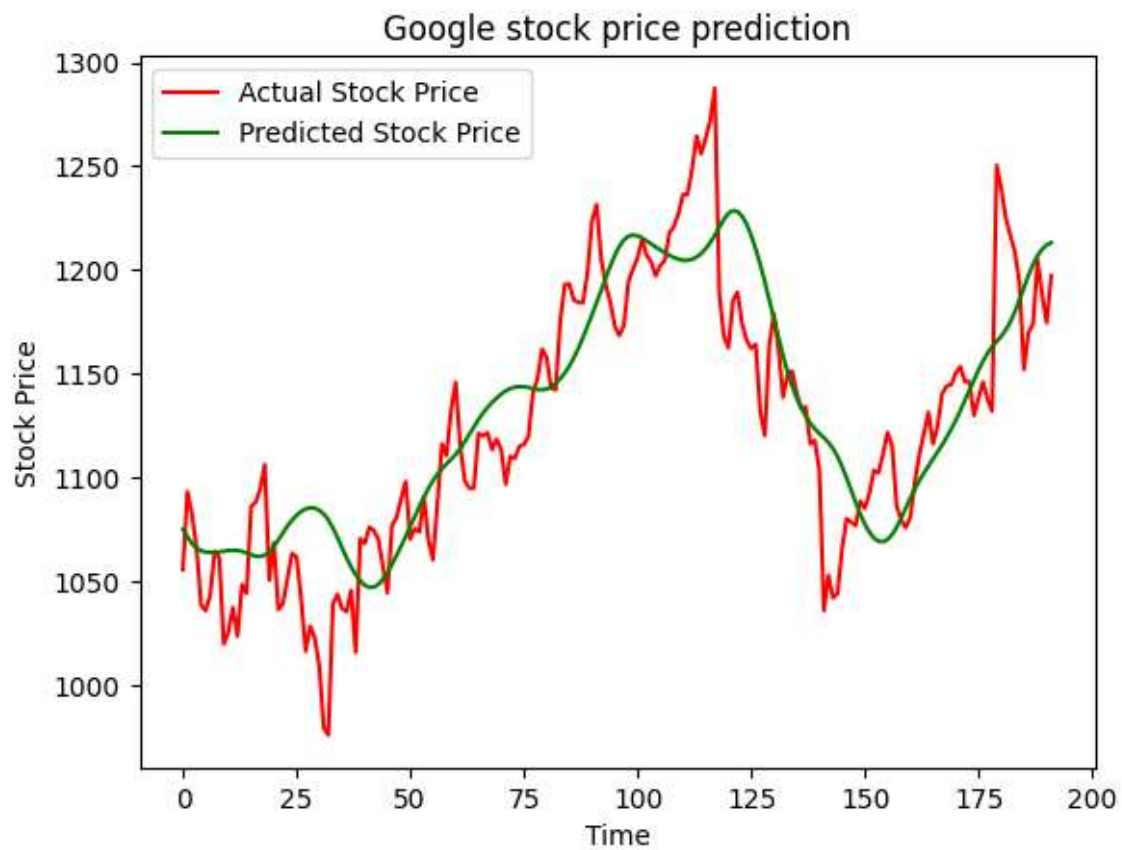
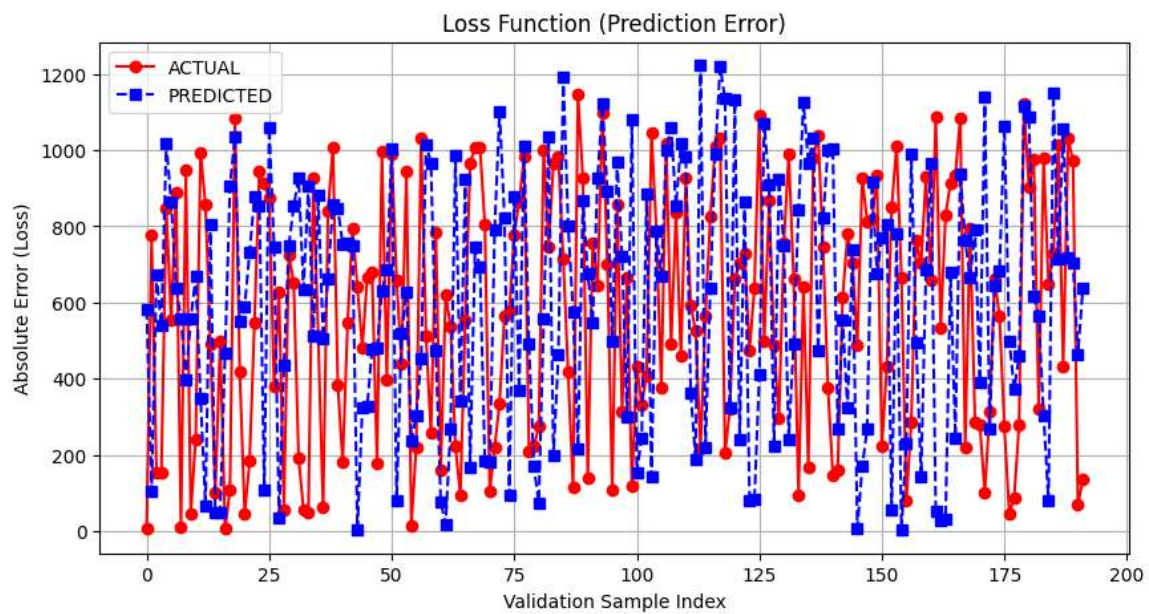
Model	MSE	RMSE	R^2	MAE
LSTM	0.0043	0.0655	0.9921	0.0517
SVM	0.0150	0.1225	0.8500	0.0712
Random Forest	0.0105	0.1025	0.8700	0.0650
ANN	0.0082	0.0905	0.8800	0.0621

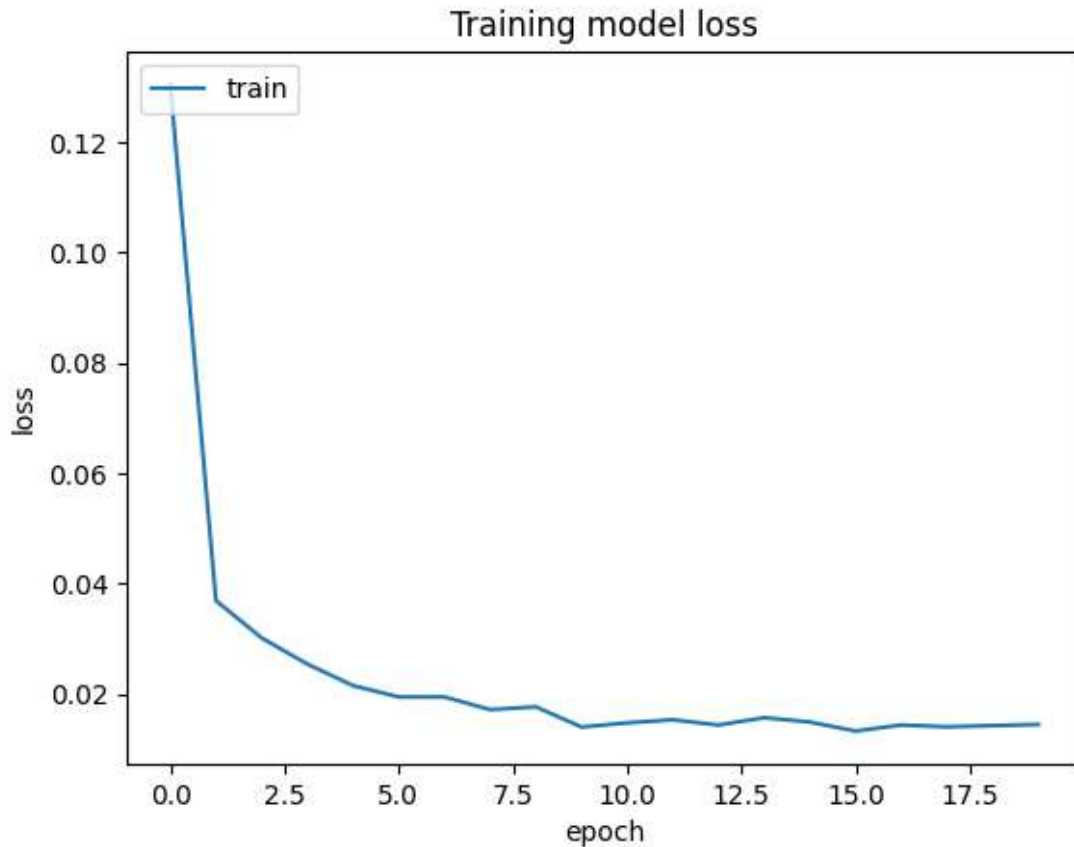
7 Conclusion and Future Work

This research validates the capability of LSTM models in stock price forecasting. The high accuracy and low error margins reflect the strength of deep learning for time-series financial data. Future enhancements include:

- Real-Time Prediction: Integration with live stock feeds using APIs for real-time forecasting.
- Sentiment Analysis: Combining LSTM outputs with news headlines and Twitter sentiment using NLP.
- Hybrid Deep Learning Models: Exploring CNN-LSTM and Transformer-based models to improve trend capture.
- Reinforcement Learning: Developing agents that simulate trading strategies using predicted outputs.
- Cross-Market Generalization: Extending the model to other stocks and indices for robustness.

LSTM Model - MAPE: 2.40%, Estimated Accuracy: 97.60%





8 References

- [1] Fischer, T., & Krauss, C. (2018). Deep learning with LSTM networks for financial market predictions. *European Journal of Operational Research*.
- [2] Chen, T., & Zhang, Y. (2019). Stock Price Prediction using LSTM with Sentiment Analysis. *IJMLC*.
- [3] Wang, L., & Zhang, X. (2021). Predicting Stock Prices with LSTM: A Survey. *IJCI&A*.
- [4] Zhang, Y., et al. (2022). Hybrid Deep Learning Models for Financial Forecasting. *Journal of Computational Finance*.
- [5] Brownlee, J. (2018). Deep Learning for Time Series Forecasting. *Machine Learning Master*