Bhavyasree Bhuvanagiri (A0262837W), Peiwen Jiang (A0262785R), Liwei Wang (A0262669N), Yumo Yao (A0262687N), Yunkai Zhang (A0262770A)

Github link: https://github.com/YYM-yym/Stock-Price-Prediction

Abstract

- This project aims to develop a stock price prediction system using a combination of sentiment analysis, Long Short-Term Memory (LSTM), and Generative Adversarial Networks (GANs). The system takes into account the historical stock prices, relevant news articles, and social media sentiments, to predict future stock prices.
- The sentiment analysis component uses natural language processing techniques to analyze tweets related to the target stock, extracting sentiment scores. These sentiment scores are then used as inputs to the LSTM model, which is trained on historical stock prices and sentiment data to predict future stock prices.
- To improve the accuracy of the predictions, the system incorporates a GAN model, which is trained on the historical stock prices and sentiment data. The GAN model then generates synthetic data that resembles the real data and then used to augment the training dataset for the LSTM model.
- The performance of the system is evaluated using various metrics, such as mean squared error, mean absolute error, and root mean squared error.
- **Keywords**: Sentiment Analysis, LSTM, GAN, stock price prediction

1. Problem Description

The Efficient Market Hypothesis (EMH) suggests that financial markets are "informationally efficient," and that current stock prices already reflect all available information, making it impossible to consistently predict future price movements. However, even in an efficient market, there may be moments of irrational exuberance or panic that could temporarily distort stock prices.

One potential source of such distortions is social media, where users can express their opinions, emotions, and reactions to real-time market events. While social media data is not a perfect predictor of stock prices, it can provide valuable insights into the collective sentiment of investors, which can sometimes anticipate changes in market trends.

What is the trend of the future stock price? How is this correlated with the sentiment on social media? In this project, we aim to address these questions using deep learning methods including sentiment analysis and time series prediction, and thus predict future stock price fluctuations based on both the historical prices and tweets sentiments with high accuracy.

2. Data Collection and Exploration

2.1 Data Sources

Our data is downloaded from Kaggle. The tweets part contains tweets for top 25 most watched stock tickers on Yahoo Finance from 30-09-2021 to 30-09-2022. And the stock price part contains 6700 stock market price and volume data for corresponding dates and stocks.

2.2 Data Description

Our data contains two parts, the first part is tweets. The components of it are the following:

- Date date and time of tweet
- Tweet full text of the tweet
- Stock Name full stock ticker name for which the tweet was scraped
- Company Name full company name for corresponding tweet and stock ticker

And the second part is stock price, it contains the following components:

- Date corresponding to tweet dates
- Open, High, Low, Close, Adj Close daily stock price
- Volume Volume for the corresponding date
- Stock Name corresponding to tweet stock names

3. Sentiment Analysis

Sentiment analysis is a natural language processing technique that involves the analysis of text to determine the overall sentiment or emotional tone of a given piece of content. One popular application of sentiment analysis is the analysis of tweets. With over 500 million tweets being posted daily, Twitter provides a vast source of data for sentiment analysis.

By analyzing tweets, businesses and organizations can gain insights into how their brand, products, or services are being perceived by the public, and can adjust their marketing strategies accordingly. Additionally, sentiment analysis on Twitter can be used to monitor public opinion on a variety of topics, including politics, current events, social issues, and stocks.

In our project, sentiment analysis using tweets can provide valuable insights into the attitudes and opinions of individuals and the public as a whole, which can better help us to foresee society's attitude toward stock, and help us make investments.

3.1 Lemmatization & textBob

We firstly define a Python function to preprocess a given tweet by performing a series of text cleaning and normalization tasks.

The function firstly removes any URLs, user mentions, and hashtags from the tweet using regular expressions (re module). It then removes any punctuation from the tweet using the translate method of the string class, and converts the tweet to lowercase using the lower method.

Next, the tweet is tokenized into individual words using the word_tokenize function from the nltk (Natural Language Toolkit) library. The function also removes any stop words (common words such as "the" and "and") using a predefined set of stop words from the stopwords module.

Finally, the function lemmatizes the remaining words (reduces them to their base or dictionary form) using the WordNetLemmatizer class from the nltk library. The resulting tokens are joined back together into a single string separated by spaces, and the preprocessed tweet is returned.

3.2 Flair

We also used the Flair library to perform sentiment analysis on a preprocessed list of tweets.

Firstly, import the Flair library and load a pre-trained sentiment analysis model for English (en-sentiment). The code then loops over each preprocessed tweet in the preprocessed_tweet column of a dataframe. For each tweet, a Flair Sentence object is created, and the sentiment model is used to make a prediction on the sentence.

The numerical score of the prediction is appended to the probs list, and the sentiment value (either "POSITIVE" or "NEGATIVE") is appended to the sentiments list.

Finally, the code adds the probs and sentiments lists as new columns in the original dataFrame, with the names probability and sentiment, respectively.

3.3 Sentiment Intensity Analyzer

We firstly import the SentimentIntensityAnalyzer class from the nltk.sentiment.vader module. Then, a new instance of the SentimentIntensityAnalyzer class is created and assigned to the variable sentiment_analyzer.

Next, the code loops over each row in the dataFrame. For each row, the Tweet column value is retrieved and normalized using the unicodedata.normalize function. The VADER polarity_scores method is then called on the normalized tweet text to compute a sentiment score, which is a dictionary containing four scores: negative, neutral, positive, and compound.

The code then uses the at method of the dataFrame to update the sentiment_score, Negative, Neutral, and Positive columns with the corresponding values from the sentiment score dictionary.

3.4 Method Comparison

TextBlob is a simple and easy-to-use tool that can be installed as a Python package. It offers good performance on short and simple text and also provides polarity scores for each sentence in addition to the entire text. However, it may not perform well on complex or nuanced text, and its pattern-based algorithm may not work well in all situations.

Flair stands out for its state-of-the-art performance on a wide range of NLP tasks, including sentiment analysis. It combines traditional statistical NLP techniques with deep learning methods to achieve this high level of performance. However, it may not be as simple and straightforward to use as TextBlob or require more computational resources than some other options.

Sentiment Intensity Analyzer is based on the VADER lexicon, which has been shown to perform well on a variety of different types of text. It calculates polarity scores for both positive and negative sentiment and provides separate scores for positivity, negativity, and neutrality. However, it may not perform well on highly subjective or nuanced text, and it relies exclusively on the VADER lexicon, which may not capture all nuances of sentiment.

Comparing these three methods, Flair's state-of-art performance on NLP tasks offers a high accuracy. Also, other features it provided, such as probability, are very useful for future analysis and prediction. Also, given the complexity of the Tweets and the non-standard languages used on social media, Flair is shown to be satisfactory when handling these issues.

4. Exploratory Data Analysis

4.1 Text Visualization - Word Cloud

Word cloud could be generated when visualizing text-based data. Text frequency and correlation between texts can be shown in a word cloud, by the size and position of texts in the plot. It is an effective way to summarize key themes and ideas within a given text.

4.1.1 All Tweets Visualization



From the word cloud chart, we can see that many stocks are highly topical in their own right, such as Tesla and Amazon. In addition to this, there are many time-related words, such as today, tomorrow and will, indicating that people on social media are also actively predicting stocks. Also, we may notice that Tesla has been a heated topic.

4.1.2 Individual Stock Tweets Visualization



Tweets about Tesla stock are quite self-focused. This is to say, compared to other tweets, those comments seldom

mention other stocks. Also, those tweets include strong opinions about the market and the company operations, so that words like "vehicle" and "elon musk" are frequently mentioned.



Tweets about Amazon stock, on the other hand, have been frequently mentioning other stocks, including other companies from FANNG. Interestingly, the comments are more focused on finance, with words "stock", "investor", "growth", "buying", "holding", etc., and less focused on company products and operations.

Tesla and Amazon are two examples that present different concerns and focuses of the public. Generally, the tweets might include opinions or comments on recent news or developments related to the company, updates on the stock price, discussions on the company's products or services, and predictions or forecasts about the future performance of the stock. Additionally, there may be tweets from financial analysts or investors providing analysis or recommendations on whether to buy or sell the stock. These tweets could provide useful insights on predicting the company's stock price and the whole market.

4.1.3 Tweets Visualization by Sentiments



From the positive tweets, positive words including "well", "support", "thank", "good" pop out with high frequency. Also, presence of the verb "buy" indicates the profitability.



Negative tweets, unexpectedly, show quite a lot words that indicate hope and expectation. This probably reflects the mindset of stockholders: happy when the stock is up, and hopeful when the stock is down.

4.2 Text Visualization - Other Plots

Other than the word cloud, other methods are used to visualize text-based data. A bar chart is used to visualize the most common words in the tweets, and most of them are stock names. From the pie chart, we can see that more than a half of the tweets are negative ones, with 43,498 in total. The number of positive tweets is 37,295.





The plot shows the number of tweets changed and the sentiment over time. There are several peaks during this period, 2022/1/20 - 2022/2/10 and 2022/4/15 - 2022/4/30, with high tweet volumes and extremely high volume of negative tweets.



4.3 Stock Price Visualization

We have also visualized the stock price using line charts. Despite the fluctuation, there are indeed significant drops during these two periods mentioned above. This is to say that social media comments do indicate the trend of the stock market.



With all these exploratory data analysis, we can have some idea about the content of the tweets and the stock trend. Also, we can find that tweets do have some correlation with stock prices.

5. Prediction Models

In the following section, we used two deep learning models LSTM and GAN to compare their performances, and since the price fluctuates a lot during the original period, we also create another structure using time step to make the plots of the prices smoother. We also completed the hyperparameter tuning by random search to further improve our prediction.

5.1 Data Preparation

In order to better fit the deep learning model with the preprocessed data, we merge the datasets and create filters for both Tesla and Amazon stocks so that we can use two stock prices to evaluate our models.

To be more specific, firstly we have two datasets:

- Dataset d1 contains information about the sentiment-flair model. It has 80793 data concerns Date, Tweet, Stock Name, Company Name, time, preprocessed_tweet, probability and sentiment.
- Dataset d2 contains information of daily change of stock prices. For instance, Date, Stock Name, Open, High, Low, Close, Adj Close, Volume.

We merge them on the columns Date and Stock Name, then filter the merged dataset to extract data for two specific stock names, "TSLA" and "AMZN" for further processing. We also convert the Date column of d1 and d2 to datetime format using the pd.to_datetime() function to help further reading and unifying. Then we build two functions to further prepare the data by firstly selecting the relevant features to rescale the data with MinMaxScaler and setting the Date column as the index.

Secondly, the function prepare_data splits the data into input and output sequences, where the input sequence is a window of the past 'lookback' number of days, and the output sequence is the close price for the next day. It also stores the date for each output sequence. The output of the function is three arrays:

- X, which contains the input sequences of length lookback;
- y, which contains the corresponding output values (in this case, the Close price of the stock);
- dates, which contains the dates corresponding to each output value.

In the last step of data preparation, we split the dataset into training and testing sets using the train_test_split function with 20% data in the testing set and 80% in the training set. The data is not shuffled since we need to preserve the chronological order of the time series data.

5.2 LSTM Model Prediction

5.2.1 Introduction on LSTM Stock price prediction

LSTM (Long Short-Term Memory) is a type of recurrent neural network that is commonly used for time-series forecasting, including stock price prediction. LSTM networks are well-suited for modeling time series data because they can learn long-term dependencies in the data and handle input sequences of variable length. In the context of stock price prediction, LSTM models can be trained on historical price and volume data to predict future stock prices.

Stock price prediction is a challenging task due to the complex and non-linear relationships between market variables, as well as the impact of external events and news on the stock market. However, LSTM models have shown promising results in this domain and are widely used by traders and investors to inform their decision-making processes.

In our project, we used varied historical stock prices to train LSTM and then optimize the model using various techniques such as regularization and hyperparameter tuning. Finally, the trained model can be used to generate predictions for future stock prices, which can be used to inform trading strategies and investment decisions.

5.2.2 LSTM without the time step

The LSTM model is built using the build_lstm_model function. It shows the architecture that includes three LSTM layers and a Dense layer. The LSTM layers have 50 units, and the dropout rate is set to 0.2 to prevent overfitting. The model is then compiled with the Adam optimizer and uses the mean squared error as the loss function.

We set the number of epochs to 50, and the batch size to 32 during the model training. The training data of each stock is fed to the model, and the validation data is set directly to the testing set so that the training session can print the loss to monitor the performance of the LSTM model from each epoch.

Since our goal is to predict the close prices for Tesla and Amazon stocks in the testing set, we use the mean squared error to evaluate the differences between the predicted and actual close prices. Thus, we build the evaluate_model function to do the final evaluation for the model performance. And the LSTM without time step model reaches the MSE = 4.782 for the Tesla stock price, MSE = 11.915 for the Amazon stock price.

As for the visualization, we plot the actual versus predicted stock close price for each stock in the same graph over Oct 1, 2021 to Nov 22, 2021.





As discussed above, we also build another LSTM model with the time step. The structure and hyperparameters stay the same in this section, but we use the time index instead of the exact date. In the predictions plot, "Time" represents the index or order of the test data points used for the prediction. The time steps don't necessarily correspond to specific dates, but rather show the order in which the test data is used for evaluating the model's performance.

The test data is not shuffled before splitting, so the time steps in the plot are in chronological order. The plot is used to visualize how well the LSTM model's predictions match the actual stock prices over the given test data sequence.



5.2.4 LSTM hyperparameter tuning

Since we are dealing with deep learning models, RandomizedSearchCV is a computationally efficient method for hyperparameter tuning. It randomly tests hyperparameters from a predefined search space and selects the best ones to optimize the model's performance on the testing dataset.

We define a hyperparameter search space with two hyperparameters:

- lstm_units: The number of LSTM units in the model.
- dropout_rate: The rate at which the model randomly drops out units during training.

We select the following values for each hyperparameter:

- lstm_units: [30, 50, 70, 100]
- dropout_rate: [0.1, 0.2, 0.3, 0.4, 0.5]

The best hyperparameters are 'lstm_units' = 30 and 'dropout_rate' = 0.1 for the Tesla stock price dataset, and the model reaches a mean squared error of 4.799, which does not make improvements on the basic LSTM model. However, for the Amazon stock price data, random search returns a mean squared error of 2.437, which is much lower than the original model, with the same best hyperparameters 'lstm_units' = 30 and 'dropout_rate' = 0.1.

Also, we plot the corresponding graphs for price variation during the time index. In this section, we prefer to visualize the model with the time step so that the curves overlap more often than original ones, showing more direct matching actual and predicted data for stock price changes.



5.2.5 Conclusion

Three LSTM models are built, trained, and evaluated for stock price prediction using TSLA and AMZN data. The best LSTM models are found to have a mean squared error of 4.782 for TSLA and 2.437 for AMZN, indicating good predictive performance.

5.3 GAN Model Prediction

5.3.1 Introduction on GAN Stock price prediction

A Generative Adversarial Network (GAN) is a type of neural network that involves two neural networks working together in a game-theoretic framework. One network, called the generator, is responsible for creating new data, while the other network, called the discriminator, is responsible for distinguishing between real and fake data. The two networks are trained together in a process called adversarial training, where the generator tries to create realistic data to fool the discriminator, while the discriminator tries to correctly identify real and fake data.

The most well-known applications of GANs would be image generation, style transfer, and data augmentation. Moreover, in a previous study, GANs can be used to identify patterns in the data that may not be easily visible through traditional statistical analysis since GANs can help to uncover underlying relationships and correlations between the Twitter data and the stock price.

5.3.2 GAN without the time step

Like we have two versions of LSTM models for different time periods in the previous section, during the implementation of GAN models, we also choose to create with and without the time step versions to help visualize by smoothing the results after presenting the prediction by date. In this section, we will first discuss the GAN model without the time step.

Firstly, we define a basic GAN model structure by constructing three functions as discussed above: build_generator, build_discriminator, and build_gan, which create the generator, discriminator, and GAN model itself respectively.

The generator function takes a latent dimension and output shape as input and returns a model that generates synthetic data. It consists of two dense layers with leaky ReLU activation and batch normalization, followed by a dense layer with a tanh activation function that outputs data of the desired output shape.

The discriminator function takes an input shape as input and returns a model that distinguishes between real and fake data. It consists of two convolutional layers with leaky ReLU activation and dropout, followed by batch normalization, flattening, and a dense layer with sigmoid activation that outputs a binary value representing the validity of the input data. Then we use the built basic GAN model to generate predictions for stock prices of TSLA and AMZN. The first step is to set the hyperparameters for the defined models, including the size of the latent space, number of epochs, batch size, and learning rate. Then the 'train_gan' function is called twice to train the GAN models for both TSLA and AMZN training sets. As stated in the LSTM section, we feed the model with the training data, and evaluate the results using the testing data. Finally, the function returns the MSE values of 0.688 and 0.743 for Tesla and Amazon stock pricing prediction. We also visualize our results like we do in the LSTM section, but with different scales.



5.3.3 GAN with the time step

Though in the last section, we do not improve much in Mean Square Errors of the stock Tesla by adding a time step to our LSTM model, it seems the time step efficiently helps with the GAN model to produce more accurate data for Tesla predicted price. The reason for that may be due to the larger fluctuations in the GAN model for TSLA, and by using the time step, we can decrease the influence. After applying the time step, the Mean Square Error reduces to 0.599.

5.3.4 GAN hyperparameter tuning

As we discuss in the LSTM section, given that we are working with deep learning models, RandomizedSearchCV is an efficient way to tune hyperparameters. It randomly tests hyperparameters from a predefined search space and identifies the optimal ones to enhance the model's performance on the testing dataset. Thus, we choose to follow the similar steps to complete hyperparameter tuning for GAN models.

The best hyperparameters are

- 'learning_rate': 0.05;
- 'batch_size': 64;
- 'num_layers': 3;
- 'num units': 128

And through applying the best results to our model, we reach the minimal Mean Square Error 0.214 for TSLA and 0.040 for AMZN stock price prediction.

In this section, we opt to display the model's performance by creating graphs that illustrate the price variation over time. We also choose to use the time step model for visualization purposes, as it enables the actual and predicted data for stock price changes to align more closely, resulting in a better match between the two curves.



5.3.5 Conclusion

In conclusion, our study has demonstrated the potential of using Generative Adversarial Networks (GANs) for predicting stock prices. By incorporating NLP features and time steps, our GAN models achieved competitive results with a mean squared error (MSE) of 0.214 for TSLA and 0.040 for AMZN. Our results suggest that the proposed model can be an effective tool for predicting stock prices.

6. Conclusion

After conducting a thorough analysis of the stock market, utilizing sentiment analysis, LSTM, and GAN, we have successfully achieved an MSE of 0.214 in our stock price prediction project. This result indicates that our predictive model was able to accurately forecast future stock prices with a high degree of precision, making it a valuable tool for investors seeking to make informed investment decisions.

The integration of sentiment analysis, LSTM, and GAN enabled us to extract valuable insights from large volumes of data, while also accounting for fluctuations in the market and investor sentiment. Our model's ability to learn from historical data and adjust its predictions based on changing market conditions makes it a highly effective tool for long-term investment strategies.

Overall, our research demonstrates the potential for advanced machine learning techniques to provide highly accurate stock price predictions, which can aid investors in making informed decisions. With further refinement and development, we believe that our approach can be used to help investors achieve greater returns on their investments while minimizing risks.

7. Discussion

First, the finding that the proposed LSTM and GAN models with NLP features and time steps can effectively predict stock prices has significant implications for investors. Investors can use these models to inform their investment decisions, potentially leading to better returns and reduced risk. As the stock market is notoriously difficult to predict, any tool that can provide accurate predictions is highly valuable.

Moreover, this study has important implications for researchers in the field of financial forecasting. The results suggest that incorporating NLP features and time steps into models can improve prediction accuracy. This finding is particularly important given the vast amounts of data that are now available through social media platforms like Twitter. Researchers can continue to explore the use of social media and alternative data sources to further improve stock price prediction accuracy.

Further research can also investigate the use of other social media platforms besides Twitter, such as Reddit or Facebook, and how incorporating data from these platforms could improve prediction accuracy. Additionally, the LSTM and GAN models could be extended to include additional financial features or indicators to improve prediction accuracy. Finally, further research can also explore the limitations of the model and address any potential ethical concerns related to the use of social media data in financial forecasting.

References

[1]Authors, A. Suppressed for anonymity, 2010.

[2]Duda, R. O., Hart, P. E., and Stork, D. G. *Pattern Classification*. John Wiley and Sons, 2nd edition, 2000.

- [3]Kearns, M. J. *Computational Complexity of Machine Learning*. PhD thesis, Department of Computer Science, Harvard University, 1989.
- [4]Langley, P. Crafting papers on machine learning. In Langley, P. (ed.), Proceedings of the 17th International Conference on Machine Learning (ICML 2000), pp. 1207–1216, Stanford, CA, 2000. Morgan Kaufmann.
- [5]Michalski, R. S., Carbonell, J. G., and Mitchell, T. M. (eds.). *Machine Learning: An Artificial Intelligence Approach, Vol. I.* Tioga, Palo Alto, CA, 1983.
- [6]Mitchell, T. M. The need for biases in learning generalizations. Technical report, Computer Science Department, Rutgers University, New Brunswick, MA, 1980.
- [7]Newell, A. and Rosenbloom, P. S. Mechanisms of skill acquisition and the law of practice. In Anderson, J. R. (ed.), *Cognitive Skills and Their Acquisition*, chapter 1, pp. 1–51. Lawrence Erlbaum Associates, Inc., Hillsdale, NJ, 1981.
- [8]Samuel, A. L. Some studies in machine learning using the game of checkers. *IBM Journal of Research and Development*, 3(3):211–229, 1959.
- [9]Yu, Yong, et al. "A review of recurrent neural networks: LSTM cells and network architectures." Neural computatioStaudemeyer, Ralf C., and Eric Rothstein Morris. "Understanding LSTM--a tutorial into long short-term memory recurrent neural networks." arXiv preprint arXiv:1909.09586 (2019).n 31.7 (2019): 1235-1270.
- [10]Huang, Zhiheng, Wei Xu, and Kai Yu. "Bidirectional LSTM-CRF models for sequence tagging." arXiv preprint arXiv:1508.01991 (2015).
- [11]Greff, Klaus, et al. "LSTM: A search space odyssey." IEEE transactions on neural networks and learning systems 28.10 (2016): 2222-2232.
- [12]Gers, Felix A., Jürgen Schmidhuber, and Fred Cummins. "Learning to forget: Continual prediction with LSTM." Neural computation 12.10 (2000): 2451-2471.
- [13]Venter, J. Craig, et al. "The sequence of the human genome." science 291.5507 (2001): 1304-1351.
- [14]Yeh, J-W., et al. "Nanostructured high-entropy alloys with multiple principal elements: novel alloy design concepts and outcomes." Advanced engineering materials 6.5 (2004): 299-303.
- [15]Ariyo, Adebiyi A., Adewumi O. Adewumi, and Charles K. Ayo. "Stock price prediction using the ARIMA model." 2014 UKSim-AMSS 16th

international conference on computer modelling and simulation. IEEE, 2014.

- [16]Schöneburg, Eberhard. "Stock price prediction using neural networks: A project report." Neurocomputing 2.1 (1990): 17-27.
- [17]Yu, Pengfei, and Xuesong Yan. "Stock price prediction based on deep neural networks." Neural Computing and Applications 32 (2020): 1609-1628.
- [18]Kohara, Kazuhiro, et al. "Stock price prediction using prior knowledge and neural networks." Intelligent Systems in Accounting, Finance & Management 6.1 (1997): 11-22.
- [19]Selvin, Sreelekshmy, et al. "Stock price prediction using LSTM, RNN and CNN-sliding window model." 2017 international conference on advances in computing, communications and informatics (icacci). IEEE, 2017.
- [20]Adebiyi, Ayodele Ariyo, Aderemi Oluyinka Adewumi, and Charles Korede Ayo. "Comparison of ARIMA and artificial neural networks models for stock price prediction." Journal of Applied Mathematics 2014 (2014).