



Stock Market Portfolio Allocation Using Twitter Sentiment Analysis

**Kulik Jan
Pooja Grewal Kaur
Rachel Loo Jia Yin
Wisely Neo Boon Hao**

School of Physical and Mathematical Sciences

2022

Abstract

The market is driven by people and people by emotions. Emotions of individuals are easily swayed by external factors, such as social media and financial news. Recently, there have been numerous events where sentiments have been more influential in driving stock prices than any other factors pertaining to fundamental or technical aspects of the stock. Two notable examples are cryptocurrency mania and Gamestop. We believe that if the sentiment of an asset is quantified it will be possible to trade that asset profitably and minimise risk. In order to do that, unstructured tweets are first cleaned and Word2Vec is used to extract qualitative representation of tweets' features. Data are partitioned using training-test rate of 80:20. Furthermore, three models – Random Forest, Long Short-Term Memory (LSTM) and Transformers are trained and tested using the pre-processed data. Subsequently, we developed two trading strategies to utilise the acquired sentiment scores. Those strategies are then back-tested using walk-forward optimisation method and evaluated using several performance metrics. Finally, we estimate the cost of running the model and minimal investment.

Table of Contents

Abstract	i
1 Introduction	1
2 Methods	2
2.1 Datasets	2
2.1.1 Unsupervised Dataset	2
2.1.2 Supervised Dataset	2
2.1.3 Testing Dataset	3
2.2 Portfolio Allocation	3
2.3 Investment Strategy	4
2.3.1 Rebalance Strategy	4
2.3.2 Utilising Sentiment Score	5
2.3.3 Portfolio Weighting	5
2.4 Backtesting	6
3 Models	7
3.1 Word2Vec	7
3.1.1 Background	7
3.1.2 Implementing the model	7
3.2 Random Forest	7
3.2.1 Implementation and Results	7
3.3 LSTM	8
3.3.1 Background	8
3.3.2 Running the model	9
3.4 Transformers	9
3.4.1 Background	9
3.4.2 Model selection	10
3.4.3 Fine-tuning the model	11
3.5 Models comparison	11
4 Results	12
5 Discussion	15
6 Conclusion	16
References	17
A Appendix	A-1
A.1 Data processing	A-1
A.2 Word2vec	A-2
A.3 Random Forest	A-2
A.4 LSTM	A-4

A.5	roBERTa	A-6
A.6	Results	A-7

Introduction

Stock market prediction is of interest to many people and entities as it provides significant investment opportunities for individuals and companies to grow their assets. Stock traders predict trends in the stock market to determine when to buy or sell a stock in an attempt to maximise profits. If these trends are adequately predicted, they are able to earn a considerable profit margin. The financial markets, however, are influenced by a plethora of factors ranging from economic to national policies, disasters and seasons. The stock market has long been regarded as one of the most volatile markets and consequently, is difficult to predict. Traditionally, it was predicted using quantitative variables, mainly on the stock and the market conditions. In this report we adopt a different strategy and attempt to analyse the Twitter sentiment associated with particular stocks in order to minimise the risk of investment.

Methods

2.1 Datasets

A total of 2 datasets related to financial tweets were sourced from Kaggle, the first contains over 23,000 unlabeled tweets that will be used to train the Word2Vec model for feature generation. The second is a labeled dataset of over 5,000 tweets that will be used to train, validate and test models. Stock symbols were scraped over various sources and used to transform datasets for consistency. We transformed stock tickers to the value of “@stock” and numbers from integers to the word “@number” since numbers are arbitrary in this case and might lead to increased inconsistencies.

2.1.1 Unsupervised Dataset

This dataset consists of financial tweets that have a mix of stocks, currencies, and commodities. To generate features for the model consistent with the target outcome, we filtered the dataset with the scraped stock tickers that match up with those in the Nasdaq stock exchange. Subsequently, we did a simple exploratory data analysis (EDA) to check that there were a total of 23,598 tweets available to be fed into the feature-generating algorithm.

2.1.2 Supervised Dataset

A similar pre-processing algorithm was applied to this supervised dataset and a value of 5,791 tweets were returned for training, validation, and testing. Before splitting the data, we performed EDA on the dataset to obtain more information.

2.1.2.1 Labeled Sentiment Comparison

From the following graph it can be observed that there is a slight imbalance between negative sentiments and positive sentiments.

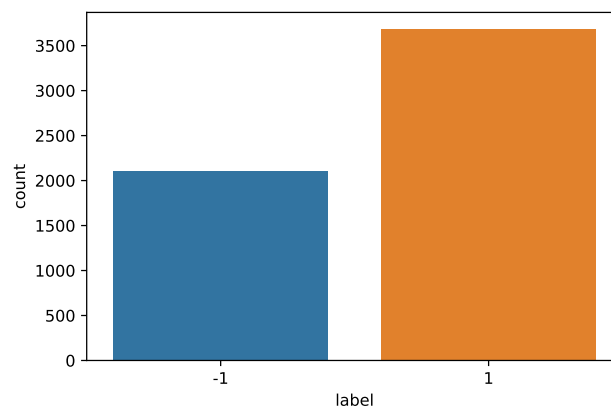


Figure 2-1: Countplot for Supervised Data

to predict the market movement and allocate portfolio accordingly.

Evaluation of collective sentiment towards a particular stock can be efficiently performed on social media platforms such as Twitter. It has enabled regular investors as well as bigger entities to share their thoughts, predictions and analyses about the stock market. Furthermore, many investors seek financial advice on social media thus potentially taking those opinions into consideration when creating their financial strategy. This results in a complex social structure, where opinions and beliefs can eventually become influential enough to have a real impact on the stock market thus impacting the price [4].

2.3 Investment Strategy

First, we need to establish a strategy to evaluate sentiment. For this purpose, we construct several models using Natural Language Processing (NLP) techniques. Those models are subsequently compared and the most accurate one is used to compute tweets' sentiment. The model processes each tweet and assigns a sentiment score, which essentially is the probability of the sentiment conveyed in a particular tweet being positive.

Furthermore, Twitter's key-word search functionality is used in order to retrieve tweets corresponding to a particular stock that is being analysed. Sentiment score of each stock is evaluated by searching for all tweets from a particular period that contained the respective ticker symbol. The models can only evaluate general sentiment of a tweet, so if multiple stocks are mentioned it is impossible for them to distinguish whether the sentiment towards a particular one is positive or negative. Therefore, those tweets have to be filtered out. After analysing all available tweets over the desired period, we averaged individual sentiment scores coming from single tweets in order to obtain the final sentiment score of a stock, which is represented by Equation 2.1

$$SS = \frac{1}{N} \cdot \sum_{i=1}^N IVS_i \quad (2.1)$$

where SS is the sentiment score, N is the total number of tweets corresponding to a particular stock and IVS is the individual sentiment score. Once the method of assigning sentiment scores is developed, it is possible to apply several strategies in order to use that information. These strategies are discussed in the following subsections.

2.3.1 Rebalance Strategy

The first important issue is defining the period over which the sentiment score is aggregated and portfolio rebalanced. Evaluation has to be performed at regular intervals after which a decision whether a particular stock should be bought, held or sold is made. Previously conducted studies have shown that the impact of sentiment on the overall stock market is rather short, around a span of one day [5]. Furthermore, Katayama et al. [6] evaluated the performance of a daily, weekly and monthly rebalance strategy in the context of sentiment analysis and concluded that the former significantly outperforms the latter two. Therefore, we decided to use the daily rebalance strategy.

Due to a limited number of tweets available each day it is crucial to select a subset of stocks to be considered such that there is enough data available to make an informed decision. Therefore,

we decided to only consider stocks from the S&P100 index, which contains 101 leading US stocks. Those are the most popular stocks that are widely discussed on English-speaking social media platforms and thus, there should be enough data available to consistently evaluate the sentiment of each of those stocks.

2.3.2 Utilising Sentiment Score

There are various strategies that can be applied in terms of utilising the obtained sentiment scores. The first one involves using the absolute value of sentiment. After evaluating all sentiment scores for the S&P100 index we select some number of stocks with highest scores to be included in the portfolio. If one of the selected stocks is being held in the portfolio from the previous day, no further action towards that stock is taken. If one of the stocks that is being held in the previous portfolio is not among the selected stocks, it is sold. Finally, if the stock is selected but is not contained in the previous portfolio, it is bought.

There are two hyperparameters that have to be tuned for this strategy: the number of stocks that are bought every day and the minimal number of tweets for the sentiment score to be considered as valid. The second threshold stems from the fact that some stocks are much more popular than others and thus, the number of available tweets can differ significantly. If the number of tweets is very low, the aggregated sentiment score might fail to represent the real sentiment towards a particular stock. Therefore, it is necessary to filter out those instances. Every time the number of available tweets about a certain stock is below the pre-defined threshold, we set its sentiment score to 0 in order to disregard it from the subsequent analysis. In case there are not enough non-zero sentiment scores available to meet the number of stocks defined by the first hyperparameter, only the stocks with non-zero scores are included.

The second strategy involves considering a daily change in the value of sentiment rather than the absolute value. This approach stems from the fact that some stocks consistently have lower sentiment scores than others, but they still might be worth investing in. Therefore, we apply a very similar approach to the previously described one, but the percentage change of sentiment with respect to the previous day is used instead of the absolute value of sentiment.

2.3.3 Portfolio Weighting

After the set of stocks to be included in the portfolio is selected, it is necessary to establish the weighting between them. First, it is possible to assign an equal weight to each of the stocks, as in Equation 2.2 [6]

$$w_{eq,i,T} = \frac{1}{N} \quad (2.2)$$

where N is the number of selected stocks, i is the index of the stock and T is the point in time. Moreover, it is possible to perform portfolio weighting based on the market capitalisation, as in Equation 2.3 [6]

$$w_{cap,i,T} = \frac{NAV_{i,T}}{\sum_{i=1}^N NAV_{i,T}} \quad (2.3)$$

where $NAV_{i,T}$ is the market capitalisation of stock i at time T , N is the number of selected stocks, i is the index of the stock and T is the point in time. Katayama et al. [6] investigated

both weighting methods and concluded that the first one performs better both in terms of the expected revenue and risk. Therefore, we choose it as the final portfolio weighting method.

2.4 Backtesting

After the model is ready, a backtest has to be performed in order to tune the hyperparameters and evaluate the actual performance. Due to the relatively small amount of data, we decided to perform a walk-forward optimisation, which is shown in Figure 2-3.

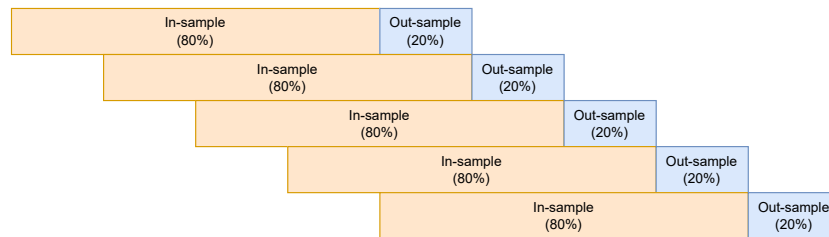


Figure 2-3: Walk-forward optimisation

The algorithm involves splitting the data into some number of folds and subsequently finding the optimal values of hyperparameters on the in-sample data of each fold. Then, the optimised model is tested on the out-sample data in order to obtain an estimation of the actual performance. Then, these steps are repeated for each fold. This procedure allows to use a much bigger portion of the available data while still providing meaningful results. The final model is backtested using a walk-forward optimisation with 10 folds.

Models

3.1 Word2Vec

This section introduces Word2Vec, an unsupervised method used to produce word embeddings. Word embedding is a set of language modeling and feature generation techniques under NLP. Word embedding captures the context of a word in a document, semantic and syntactic similarity, and relative association with other words in the document. In NLP, texts from the documents are mapped to vectors.

3.1.1 Background

Developed in 2013, it was discovered that the language model can be trained in two steps: using a simple model to train continuous word vectors followed by an N -gram feedforward neural network model [7]. Two methods were proposed then, the first being Continuous Bag of Words (CBOW) and the other being Continuous Skip-Gram Model. This project will be focusing on the use of the latter. The skip-gram model uses the current word to predict the context of surrounding words within a set window. An increased range of the window creates better word vectors but more complex computation is required [7]. In other words, the current word is the input layer with contextual words at the output. Numerous studies have used a number between 50 to 300 as the number of dimensions for their word embeddings, where it is also noted that performance for intrinsic tasks maximises at approximately 300 dimensions [8]. Since more dimensions require more computation, this project will use 100 dimensions for each word embedding due to time and computational constraints.

3.1.2 Implementing the model

The model used is from the notable “gensim” package. As mentioned in the previous chapter, a large amount of financial tweets was used to train the Word2Vec model to generate word embeddings. The generated word vector matrix then output dimensions of 6062 by 100, where each word is represented by a vector of 100 dimensions. In addition to the dimensions specified, a window size of 5 is used.

3.2 Random Forest

3.2.1 Implementation and Results

The Random Forest model was trained using the aforementioned Word2Vec model. The dataset was first split into training and testing sets, before mapping each tweet to a feature vector by averaging the word embeddings of all words in the tweet. We took the average of the word embeddings of the word vectors in a given tweet as the tweets dataset was of variable-lengths. By taking the individual word vectors and transforming them into a feature set that is the same length for every tweet, the variable-length tweet issue was resolved.

Since each word is a vector in a 100 dimensional space, vector operations can be used to combine the words in each tweet. We calculated the average feature vectors for training and

testing sets and used them to train a Random Forest.

To tune the hyperparameters and prevent overfitting we used GridSearch cross-validation, as GridSearch alone could produce a model that performs poorly on test data and highly on training data [9]. It was found that the optimal values of the hyperparameters are as follows:

- n_estimators: 100
- max_depth: 20

3.3 LSTM

An innovation of Recurrent Neural Networks (RNNs), Long Short-Term Memory (LSTM) network is capable of doing more than a typical RNN. Where RNNs face the problem of the inability to retain information over long sequences, the design of LSTM overcomes that issue as it uses a memory cell to store information from previous states [10]. As the nature of this project requires context from prior sequences when using neural networks, LSTM is a preferred choice over RNN.

3.3.1 Background

A vanilla LSTM model has 3 layers, the input layer, hidden layer and output layer. The LSTM model consists of 3 gates: forget gate, input gate and output gate. As information moves from one cell to another, the current cell (C_t) will receive information from the previous cell (C_{t-1}). The gates come together to form a chain of sigmoid and hyperbolic functions as opposed to just a sigmoid function in Vanilla-RNN models.

Forget Gate: The deciding factor if the information from C_{t-1} should be kept or forgotten. Controlled by a sigmoid (σ) function, it outputs a value $[0, 1]$ where if the value is 0 it tells the function to forget (ignore) everything, while 1 to retain everything.

$$f_t = \sigma(W_f \cdot [H_{t-1}, x_t] + b_f)$$

Input Gate: Takes in the current input vector (x_t) and processes the information to decide which vector should be updated.

$$\begin{aligned} i_t &= \sigma(W_i \cdot [H_{t-1}, x_t] + b_i) \\ \tilde{C}_t &= \tanh(W_C \cdot [H_{t-1}, x_t] + b_C) \end{aligned}$$

Similar to the forget gate, the sigmoid function present in i_t serves to restrict values to an interval of $[0, 1]$. The hyperbolic function outputs a value between $[-1, 1]$ to determine what values should be subtracted or added towards the cell state respectively. The values of the forget gate and input gate come together to form the current cell state as seen in the following equation:

$$C_t = i_t \cdot \tilde{C}_t + f_t \cdot C_{t-1}$$

Output Gate: Output using filtered values of the current cell state C_t . These output values are the values that the following cell C_{t+1} uses.

$$\begin{aligned} o_t &= \sigma(W_o \cdot [H_{t-1}, x_t] + b_o) \\ H_t &= o_t \cdot \tanh(C_t) \end{aligned}$$

3.3.2 Running the model

Transformed vectors obtained from the trained Word2Vec model are loaded into an embedding matrix layer of the LSTM model. The embedding layer as the name suggests provides the model with details of the input, namely the dimensions and vectorized values for each word trained, where these vectors will represent the words. Dropout and recurrent dropout layers have also been specified at 0.5 to combat overfitting as well as the lack of sufficient labelled data.

Since tweets have a character count limitation, an approximation of padding to length 50 was deemed sufficient as the max number of words in the dataset is observed to be less than 40. After the LSTM layer, the network will be activated using a sigmoid function as the model is attempting to do a binary classification. The model is then evaluated while training using a validation set, and after training, the model is evaluated over a test set by observing the accuracy through a classification report.

3.4 Transformers

Transformers are one of the most recent developments in the area of NLP [11]. Similarly to RNNs, they are designed to process sequential data, which renders them widely applicable in the area of text processing. However, contrary to RNNs, Transformers process the entire input at once which allows to identify global dependencies between data more efficiently. Furthermore, Transformers allow for significantly more parallelisation thus greatly reducing training time [11].

3.4.1 Background

Architecture of Transformers consists of two stacks of layers for the encoder and decoder, which are shown in the left and right sections of Figure 3-1, respectively. They both consist of 6 fully connected, identical layers, where each layer uses the mechanism of self-attention in its sub-layer. Encoder is used to process the input and generate encodings which represent the relevancy of different segments of the input to each other. On the other hand, decoder is used to decompose the encodings and produce the final output. Self-attention mechanism is used to scan each part of the input and evaluate the corresponding relevance of other parts. When producing the output sequence, the model is able to use the information from any of the previous states and weights them based on the previously evaluated relevance. Moreover, self-attention ensures that all positions in a sequence are connected with a constant number of operations, which allows for a better parallelisation compared to commonly used recurrent layers. Self-attention layers also have lower time complexity than recurrent layers if the sequence length is smaller than the representation dimensionality, which is usually the case. Finally, self-attention leads to models that can be interpreted more easily [11].

Apart from the self-attention sub-layers, encoder and decoder layers also use fully connected feed-forward networks, which apply linear transformations with a ReLU function, as in Equation 3.1.

$$\text{FFN}(x) = \max(0, x \cdot W_1 + b_1) \cdot W_2 + b_2 \quad (3.1)$$

Furthermore, the model uses a mechanism called multi-headed attention in order to improve the performance of the self-attention layers. It allows the model to comprehend information from different representation subspaces at different positions in a sequence [11].

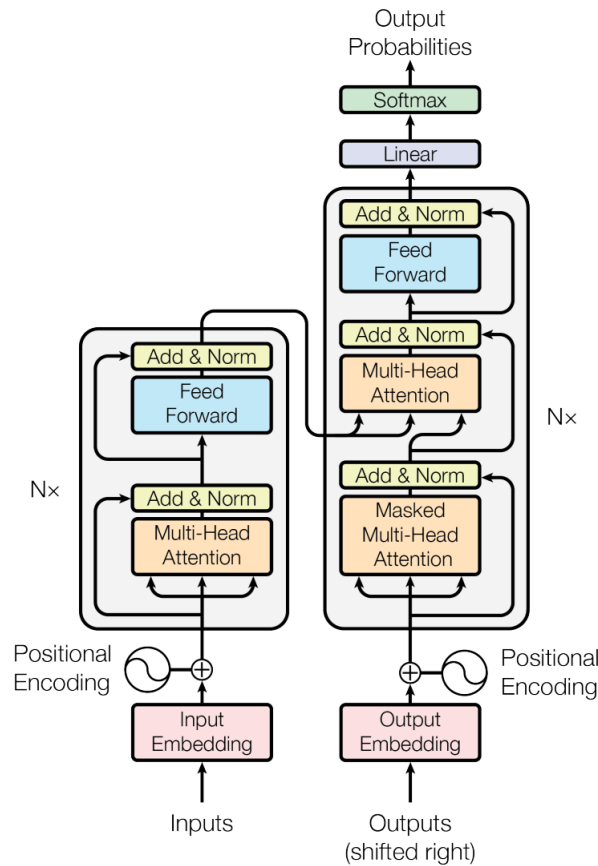


Figure 3-1: Architecture of Transformers [11]

The most efficient method of Transformers implementation is using a pre-trained model and fine-tuning it to perform a specific task. One of the most notable pre-trained models is BERT [12]. It was developed by Google in 2018 and commonly used in NLP applications since then. Since being widely adopted in the scientific community, BERT has been extensively analysed and modified, which resulted in many improved models being derived from it. One of such models is RoBERTa [13], which introduced several improvements to BERT. They primarily included training advancements, such as increased training time and batch size, longer training sequences and dynamically changing the masking pattern applied to the training data. As a result, RoBERTa is capable of performing better than other post-BERT methods [13]. For that reason, we chose to train a RoBERTa based model in order to achieve a state-of-the-art performance in sentiment analysis.

3.4.2 Model selection

The selected model was trained on a dataset of 124 million tweets in order to perform sentiment analysis. The training dataset has been based on tweets from January 2018 to December 2021 and thus the model is also capable of correctly understanding and classifying texts that relate to COVID-19. This ability may be crucial in order to correctly evaluate the market sentiment during potential future pandemic outbreaks.

3.4.3 Fine-tuning the model

The model was trained on various types of tweets. Therefore, it is necessary to fine-tune it on a dataset containing tweets about stock market in order to develop an understanding of specific financial jargon and the implications of particular words. For example, the baseline model may classify the word “bear” as neutral, but the fine-tuned model should classify it as negative since it might be associated with a bear market.

Furthermore, the baseline model classifies each piece of text as either positive, negative, or neutral. The training and test datasets contain tweets that are classified as either positive or negative, without the neutral class. Since the object of interest is binary classification, the output of the final model is converted by disregarding the probability of the neutral class and normalising the probabilities of positive and negative classes such that the entries add up to 1.

3.5 Models comparison

The following table presents test accuracy achieved by each model. Accuracy was measured on the same test set in order to obtain a meaningful comparison.

Table 3-1

Model	Test accuracy
Random forest	0.686
LSTM	0.745
roBERTa	0.827

Results

In order to gain more insight into the characteristics and dynamics of the sentiment scores, we computed and plotted them over a period of two months with the daily reevaluation strategy. This analysis was performed for some of the major stocks due to a better data availability. A resulting graph for Meta and Tesla can be seen in Figure 4-1.

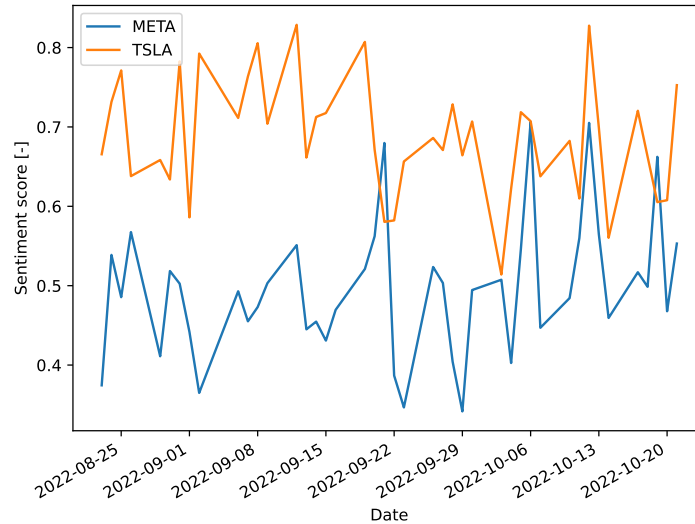


Figure 4-1: Sentiment scores of Meta and Tesla stocks over a span of two months

As it can be seen, the average value for Tesla is significantly higher than for Meta. Furthermore, similar dynamics between the two curves can be observed around the sudden spikes and falls, where they follow the same trajectory. This indicates that those changes are probably a market response to the general price trends, not necessarily to some event related to a particular company. Nevertheless, there are regions where the curves behave differently, which implies that the sentiment score can also serve as an indicator of the relative sentiment of different stocks. Furthermore, we plotted the sentiment score against the actual stock price. A resulting graph for Apple can be seen in Figure 4-2.

Again, it can be seen that both curves follow similar dynamics, with overlapping spikes and falls. This indicated that the sentiment score can be used as a price indicator, and potentially also to minimise the risk associated with price decline. Finally, we applied the two different approaches for portfolio allocation as discussed in subsection 2.3.2 and evaluated them through the method of back-testing. The result of this analysis is presented in Figure 4-3.

The third strategy called “Benchmark” serves as a baseline model for comparison with the strategies involving sentiment analysis. It involves creating an equally-weighted portfolio consisting of all S&P100 stocks and just waiting for a given amount of time. These three strategies will be evaluated and compared based on several metrics:

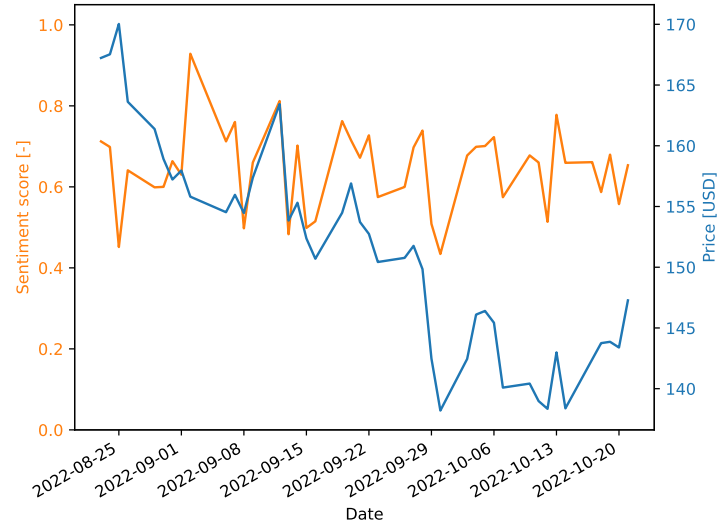


Figure 4-2: Sentiment score and price of Apple stock over a span of two months

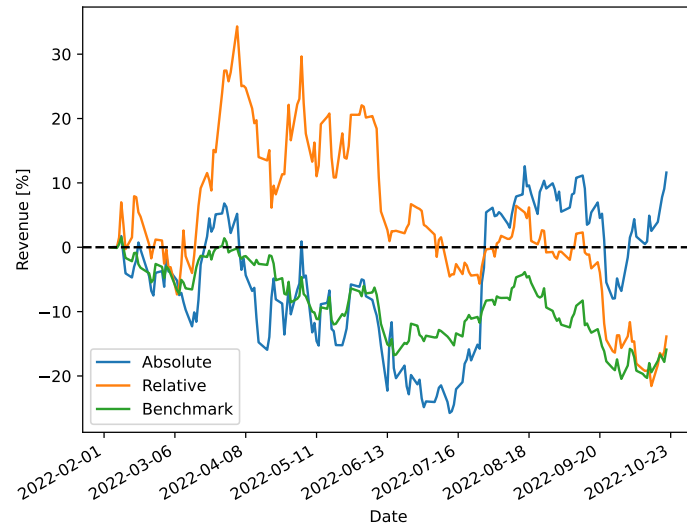


Figure 4-3: Cumulative returns of different trading strategies

- Cumulative return: total value that the strategy has earned over a given period of time. It is calculated using Equation 4.1.

$$\text{Cumulative return} = \frac{\text{Final value of investment} - \text{Initial value of investment}}{\text{Initial value of investment}} \quad (4.1)$$

- Annualised return: total value that the strategy would earn over a period of one year. It is calculated using Equation 4.2.

$$\text{Annualised return} = \left(1 + \text{Cumulative return}^{\frac{365}{\text{no. of days}}} \right) - 1 \quad (4.2)$$

- Annualised volatility: a measure of risk that is obtained by calculating standard deviation of the portfolio returns.
- Sharpe ratio: an indicator that compares the return of an investment to its risk. It is based on “risk-free” returns, which come from assets that are not subject to risk, such

as government bonds. For the sake of this comparison the U.S. 10 Year Treasury Note will be used with the value of 4.21% as of the 21st of October, 2022. Sharpe ratio is calculated using Equation 4.3.

$$\text{Sharpe ratio} = \frac{\text{Portfolio returns} - \text{Risk free returns}}{\text{Standard deviation of the portfolio returns}} \quad (4.3)$$

- Maximum drawdown: a measure of the maximal loss from peak to trough of a portfolio. It is calculated using Equation 4.4.

$$\text{Maximum drawdown} = \frac{\text{Price at trough} - \text{Price at peak}}{\text{Price at peak}} \quad (4.4)$$

Performance metrics of the three strategies are presented in Table 4-1. As can be seen, the absolute sentiment strategy achieves significantly better returns than the other two strategies. However, it also suffers from higher volatility and maximum drawdown, especially compared to the benchmark strategy. Nevertheless, the Sharpe ratio, which takes into account both returns and volatility, indicates that the absolute sentiment strategy actually performs better than the other two strategies.

Table 4-1: Performance metrics of different trading strategies

	Absolute	Relative	Benchmark
Cumulative return	11.61%	-13.87%	-15.90%
Annualised return	16.62%	-18.86%	-21.53%
Annualised volatility	49.87%	43.12%	22.96%
Sharpe ratio	0.55	-0.28	-0.96
Maximum drawdown	-34.06%	-41.61%	-21.66%

Discussion

Next, the cost of implementing the developed Transformers model to perform portfolio allocation is compared against the cost of just hiring a team of data scientists and analysts. As seen in diagram Table 5-1 below, on average, the total investment for a five-year total cost of ownership (TCO) ranges from approximately S\$60,750 to S\$94,500 [14]. This range is largely due to the two structures - Bare-Bones and MLOps Framework. The bare-bones structure lacks essential elements, resulting in performance degradation over time. The MLOps on the other hand, consist of not just developing a machine learning model, but also implementing it as a production system. This would increase automation and improve the quality of production models. The cost breakdown for both approaches is seen in Table 5-1.

On the other hand, the cost of hiring a full team of data scientists is significantly greater, with a rough estimation in the neighbourhood of S\$1.05 million to S\$2.31 million for the same five-year period. The three main roles in a data science team consist of a data scientist, data/business analyst, and subject matter expert [15]. The pay range for those roles can be seen in Table 5-2.

Cumulatively, the total cost of hiring a team of data scientists far exceeds the total investment of implementing the developed Transformers model. Ergo, it would be optimal and more economical to develop and implement the developed model with the MLOps framework approach. Furthermore, given that the average annualised return of the developed model is 16.62% (see Table 4-1), it would be necessary to invest at least S\$43,809 for a period of 5 years in order to start making profit. It has to be noted that the potential profits are not guaranteed and the model exhibits relatively high volatility.

Table 5-1: Operational Cost of ML solution

	Bare-Bones	MLOps Framework
Model Infrastructure	\$9,000/yr	\$8,000/yr
Data Support	\$6,750/yr (labour)	\$10,000 (labour) + \$3,200/yr
Engineering/Deployment	\$9,000 (labour)	\$24,500 (labour) + \$516/yr
Total Investment 5 year TCO	\$60,750	\$94,500

Table 5-2: Data Science Team without Transformer Model

	Pay Range	
Data Scientist	\$9,000/month	\$30,000/month
Subject Matter Expert	\$3,500/month	\$3,500/month
Data/Business Analyst	\$5,000/month	\$5,000/month
Total Investment 5 year TCO	S\$1,050,000	S\$2,310,000

Conclusion

With emotions as a factor driving the ever-unpredictable stock market, the objective of this project was to gauge the sentiments of the community and determine a viable trading strategy. To gauge those strategies, we used notable methods known in the machine learning industry – RF, LSTM, and Transformers. Word embeddings were generated from unsupervised tweets before applying them to the RF and LSTM algorithms, while Transformers have a more unique independent model. The model with the highest test accuracy was chosen to conduct back-testing on tweets and stock prices scraped in real-time. As the most recently developed of the three methods, the Transformers model was used for further analysis. Back-testing included 3 different investment strategies and analysed over 5 performance metrics to determine what would be the most optimal strategy on top of the sentiment analysis by the model. The sentiment score was applied to the respective strategical equations, weighted and evaluated together with a walk-forward optimisation of 10-folds. In conclusion, the highest returns for the model was noted to be the Absolute sentiment strategy, outperforming the Relative sentiment strategy and Benchmark strategy in terms of returns and Sharpe ratio. The annualised return of the final strategy was estimated to be 16.62% on average. Furthermore, the cost of running the model for 5 years in an MLOps framework was estimated to be S\$94,500, which implies that it would be necessary to invest at least S\$43,809 for a period of 5 years in order to start making profit.

However, while the numbers seem to determine this as a promising strategy, further research is required due to the constraints of data and time. The data obtained is rather biased due to the short time frame which does not consider the possibility of huge changes in society such as wars, disasters, and other inexhaustible reasons. Other methods could also be further optimised such as increasing the number of tweets to train Word2Vec, as published works have mentioned 300 dimensions as a peak parameter for the method.

References

-
- [1] B. Jacobs and K. Lev, “The complexity of the stock market,” 1989.
 - [2] E. F. Fama, “The behavior of stock-market prices,” 1965.
 - [3] P. Koratamaddi, K. Wadhwani, M. Gupta, and S. G. Sanjeev, “Market sentiment-aware deep reinforcement learning approach for stock portfolio allocation,” 2021.
 - [4] J. Smailović, M. Grčar, N. Lavrač, and M. Žnidaršič, “Stream-based active learning for sentiment analysis in the financial domain,” 2014.
 - [5] D. Katayama and K. Tsuda, “A method of measurement of the impact of japanese news on stock market,” 2018.
 - [6] D. Katayama and K. Tsuda, “A method of using news sentiment for stock investment strategy,” 2020.
 - [7] M. Tomas, C. Kai, C. Greg, and D. Jeffrey, “Efficient estimation of word representations in vector space,” 2013.
 - [8] P. Kevin and B. Pushpak, “Towards lower bounds on number of dimensions for word embeddings,” 2017.
 - [9] S. Kale, “Grid search tuning of hyper parameters in random forest classifier,” 2022.
 - [10] Y. Wang, Y. Liu, M. Wang, and R. Liu, “LSTM model optimization on stock price forecasting,” pp. 173–177, 2018. DOI: 10.1109/DCABES.2018.00052.
 - [11] A. Vaswani, N. Shazeer, N. Parmar, *et al.*, “Attention is all you need,” 2017.
 - [12] J. Devlin, M.-W. Chang, K. Lee, and K. Toutanova, “BERT: Pre-training of deep bidirectional transformers for language understanding,” 2018.
 - [13] L. Yinhan, O. Myle, G. Naman, *et al.*, “Roberta: A robustly optimized BERT pretraining approach,” 2019.
 - [14] R. Coop, “What is the cost to deploy and maintain a machine learning model?,” 2021.
 - [15] J. Saltz, “8 key data science team roles,” 2022.

Appendix

A.1 Data processing

```
1 import pandas as pd
2 import re
3 from nltk.corpus import words
4
5
6 def cleanText(text):
7     text = re.sub(r'https?\S+', 'http', text)
8     text = text.replace('user', '@user')
9     text = re.sub(r'\S+', '@user', text)
10    text = re.sub(r'\d+', '@number', text)
11
12    upperCase = re.findall(r'\b[A-Z](?:\s+[A-Z])*\b', text)
13    for i in range(len(upperCase)):
14        upperCaseSplit = upperCase[i].split()
15        for j in range(len(upperCaseSplit)):
16            if upperCaseSplit[j] in nasdaqSymbols or upperCaseSplit[j].lower() not in
               englishWords:
17                text = text.replace(upperCaseSplit[j], '@stock')
18
19    text = re.sub('[^A-Za-z0-9@]', ' ', text)
20    text = re.sub(' +', ' ', text)
21    text = text.lower()
22    text = text.strip()
23    return text
24
25
26 df = pd.read_csv("data/stock_tweets.csv")
27 nasdaq = pd.read_csv("data/nasdaq_stocks.csv")
28 nasdaqSymbols = nasdaq.iloc[:, 0].values
29 englishWords = set(words.words())
30
31 df['Text'] = df['Text'].apply(cleanText)
32
33 df = df.rename(columns={'Text': 'text', 'Sentiment': 'label'})
34 df.to_csv("data/cleaned_tweets.csv", index=False)
35
```

Listing A.1: Data cleaning

```
1 from datasets import load_dataset
2 import pandas as pd
3 import matplotlib.pyplot as plt
4 import numpy as np
5
6 dataset = load_dataset('csv', data_files="data/cleaned_tweets.csv", split='train')
7 dataset = dataset.train_test_split(test_size=0.2, shuffle=True, seed=1)
8
9 train_data = pd.DataFrame(dataset['train'])
10 test_data = pd.DataFrame(dataset['test'])
11
12 train_data.to_csv("data/train_data.csv", index=False)
13 test_data.to_csv("data/test_data.csv", index=False)
```

Listing A.2: Data splitting

A.2 Word2vec

```

1 import pandas as pd
2 from gensim.models import Word2Vec, KeyedVectors
3
4 df = pd.read_csv("data/cleaned_tweets.csv")
5 data = [row.split(' ') for row in df['text']]
6
7 model = Word2Vec(data, min_count=3, vector_size=100, workers=3, window=5, sg=1,
8                 epochs=100)
9 word_vectors = model.wv
10 word_vectors.save("data/vectors.kv")
11
12 reloaded_word_vectors = KeyedVectors.load("data/vectors.kv")
13 print(reloaded_word_vectors)

```

Listing A.3: Word2vec

A.3 Random Forest

```

1 from datasets import load_dataset, Dataset
2 #Relabel the sentiment value to [0,1] where 0 = neg and 1 = pos
3 train_data['label'] = train_data['label'].replace(-1,0)
4 test_data['label'] = test_data['label'].replace(-1,0)
5 #Transform dataset from pandas to dataset for consistency
6 train_df = Dataset.from_pandas(train_data)
7 test_df = Dataset.from_pandas(test_data)
8 trainval_df = train_df.train_test_split(test_size=0.2)
9 trainval_df, test_df
10 x_train = trainval_df['train']["text"]
11 x_val = trainval_df['test']["text"]
12 x_test = test_df["text"]
13 y_train = trainval_df['train']["label"]
14 y_val = trainval_df['test']["label"]
15 y_test = test_df["label"]
16
17 def make_feature_vec(words, model, num_features):
18     """
19     Average the word vectors for a set of words
20     """
21     feature_vec = np.zeros((num_features,), dtype="float32") # pre-initialize (for
22     speed)
23     nwords = 0.
24     index2word_set = set(model.wv.index_to_key) # words known to the model
25
26     for word in words:
27         if word in index2word_set:
28             nwords = nwords + 1.
29             feature_vec = np.add(feature_vec, model.wv[word])
30
31     feature_vec = np.divide(feature_vec, nwords)
32     return feature_vec
33
34 def get_avg_feature_vecs(reviews, model, num_features):
35     """
36     Calculate average feature vectors for all tweets
37     """
38     counter = 0

```

```

39     review_feature_vecs = np.zeros((len(reviews), num_features), dtype='float32') #
        pre-initialize (for speed)
40
41     for review in reviews:
42         review_feature_vecs[counter] = make_feature_vec(review, model, num_features)
43         counter = counter + 1
44     return review_feature_vecs
45
46 # calculate average feature vectors for training and test sets
47 num_features = 100
48 clean_train_reviews = x_train
49 trainDataVecs = get_avg_feature_vecs(clean_train_reviews, model, num_features)
50
51 clean_test_reviews = x_test
52 testDataVecs = get_avg_feature_vecs(clean_test_reviews, model, num_features)
53 trainDataVecs.shape
54
55 # remove instances in test set that could not be represented as feature vectors
56 nan_indices = list({x for x,y in np.argwhere(np.isnan(testDataVecs))})
57 if len(nan_indices) > 0:
58     print('Removing {:d} instances from test set.'.format(len(nan_indices)))
59     testDataVecs = np.delete(testDataVecs, nan_indices, axis=0)
60     x_test.drop(x_test.iloc[nan_indices, :].index, axis=0, inplace=True)
61     assert testDataVecs.shape[0] == len(x_test)
62
63 # Fit a random forest to the training data before hyperparameter tuning
64 forest = RandomForestClassifier()
65 print("Fitting a random forest to labeled training data...")
66 forest = forest.fit(trainDataVecs, y_train)
67
68 #print classification report before hyperparameter tuning
69 print("Predicting labels for test data before hyperparameter tuning")
70 result = forest.predict(testDataVecs)
71 print(classification_report(y_test, result))
72 accuracy = round(accuracy_score(y_test, result),3)
73 print(accuracy)
74
75 #Now we'll do some hyperparameter tuning using GridSearchCV.
76 def print_results(results):
77     print('BEST PARAMS: {}\n'.format(results.best_params_))
78     means = results.cv_results_['mean_test_score']
79     stds = results.cv_results_['std_test_score']
80     for mean, std, params in zip(means, stds, results.cv_results_['params']):
81         print('{} (+/-{}) for {}'.format(round(mean, 3), round(std * 2, 3), params))
82
83 from sklearn.model_selection import GridSearchCV
84 rf = RandomForestClassifier()
85 parameters = {
86     'n_estimators': [5,50,100],
87     'max_depth': [2,10,20,None]
88 }
89
90 cv = GridSearchCV(rf,parameters)
91 cv.fit(trainDataVecs,y_train)
92 print_results(cv)
93
94 #using the best parameter, we train a new random forest rf2
95 rf2 = RandomForestClassifier(n_estimators = 100, max_depth = 20)
96 rf2.fit(trainDataVecs, y_train)
97
98 #print classification report after hyperparameter tuning
99 print("Predicting labels for test data after hyperparameter tuning")
100 predictions = rf2.predict(testDataVecs)

```



```

101 print(classification_report(y_test, predictions))
102 accuracy1 = round(accuracy_score(y_test, predictions), 3)
103 print(accuracy1)

```

Listing A.4: Random Forest

A.4 LSTM

```

1 from gensim.models import Word2Vec, KeyedVectors
2 from datasets import load_dataset, Dataset
3 import pandas as pd
4 reloaded_word_vectors = KeyedVectors.load('vectors.kv')
5
6 dataset = load_dataset('csv', data_files="data/cleaned_tweets.csv", split='train')
7 dataset = dataset.train_test_split(test_size=0.2, shuffle=True, seed=1)
8 #Load training and test dataset from cleaned and split data
9 train_data = pd.read_csv('data/train_data.csv')
10 test_data = pd.read_csv('data/test_data.csv')
11 #Relabel the sentiment value to [0,1] where 0 = neg and 1 = pos
12 train_data['label'] = train_data['label'].replace(-1,0)
13 test_data['label'] = test_data['label'].replace(-1,0)
14 #Transform dataset from pandas to dataset for consistency
15 train_df = Dataset.from_pandas(train_data)
16 test_df = Dataset.from_pandas(test_data)
17 #Split into training and validation dataset
18 trainval_df = train_df.train_test_split(test_size=0.2)
19
20 #Specify hyperparameters
21 vocab_size = embedding_matrix[0]
22 mxlen = 50
23 nb_classes = 2
24
25 tokenizer = Tokenizer(num_words=vocab_size)
26 tokenizer.fit_on_texts(trainval_df['train']['text'])
27 sequences_train = tokenizer.texts_to_sequences(trainval_df['train']['text'])
28 sequences_test = tokenizer.texts_to_sequences(test_df['text'])
29 sequences_val = tokenizer.texts_to_sequences(trainval_df['test']['text'])
30
31 word_index = reloaded_word_vectors.key_to_index
32 for i in word_index: #to add 1 to the indexes given by Word2Vec as index 0 is
    reserved for padding
33     if word_index['@stock'] < 2:
34         word_index[i] += 1
35     else:
36         word_index['@stock'] = 1
37         break
38
39 print('Found %s unique tokens.' % len(reloaded_word_vectors.key_to_index))
40 res = [(key,val) for key, val in word_index.items()]
41 print('Top few parts of token dict',res[:5])
42
43 #Pad vectors to same length
44 x_train = sequence.pad_sequences(sequences_train, maxlen=mxlen)
45 x_test = sequence.pad_sequences(sequences_test, maxlen=mxlen)
46 x_val = sequence.pad_sequences(sequences_val, maxlen=mxlen)
47
48 #Transform sentiment value into categorical classes
49 y_train = np_utils.to_categorical(trainval_df['train']['label'], nb_classes)
50 y_test = np_utils.to_categorical(test_df['label'], nb_classes)
51 y_val = np_utils.to_categorical(trainval_df['test']['label'], nb_classes)
52

```

```

53 batch_size = 32
54 nb_epoch = 20
55
56 embedding_layer = Embedding(reloaded_word_vectors.vectors.shape[0],
57                             reloaded_word_vectors.vectors.shape[1],
58                             weights=[reloaded_word_vectors.vectors],
59                             trainable=False)
60
61 model = Sequential()
62 model.add(embedding_layer)
63 model.add(LSTM(128, recurrent_dropout=0.5, dropout=0.5))
64 model.add(Dense(nb_classes))
65 model.add(Activation('sigmoid'))
66 model.summary()
67
68 model.compile(optimizer='adam', loss='binary_crossentropy', metrics=['accuracy'])
69 model.save("LSTM Model")
70
71 rnn = model.fit(x_train, y_train, epochs = nb_epoch, batch_size=batch_size,
72                shuffle=True, validation_data=(x_val, y_val))
73 score = model.evaluate(x_val, y_val)
74 print("Val Loss: %.2f%%" % (score[0]*100))
75 print("Val Accuracy: %.2f%%" % (score[1]*100))
76 score_2 = model.evaluate(x_test, y_test)
77 print("Test Loss: %.2f%%" % (score_2[0]*100))
78 print("Test Accuracy: %.2f%%" % (score_2[1]*100))
79
80 plt.figure(0)
81 plt.plot(rnn.history['accuracy'], 'r')
82 plt.plot(rnn.history['val_accuracy'], 'g')
83 plt.xticks(np.arange(0, nb_epoch+1, nb_epoch/5))
84 plt.rcParams['figure.figsize'] = (8, 6)
85 plt.xlabel("Num of Epochs")
86 plt.ylabel("Accuracy")
87 plt.title("Training vs Validation Accuracy LSTM l=50, epochs=20") # for max length =
88     10 and 20 epochs
89 plt.legend(['train', 'validation'])
90
91 plt.figure(1)
92 plt.plot(rnn.history['loss'], 'r')
93 plt.plot(rnn.history['val_loss'], 'g')
94 plt.xticks(np.arange(0, nb_epoch+1, nb_epoch/5))
95 plt.rcParams['figure.figsize'] = (8, 6)
96 plt.xlabel("Num of Epochs")
97 plt.ylabel("Training vs Validation Loss LSTM l=50, epochs=20") # for max length = 10
98     and 20 epochs
99 plt.legend(['train', 'validation'])
100 plt.show()
101
102 #Check data
103 y_pred = model.predict(x_test)
104 # Convert Y_Test into 1D array
105 yy_true = [np.argmax(i) for i in y_test]
106 yy_scores = [np.argmax(i) for i in y_pred]
107
108 test_df = test_data
109 test_df['pred'] = yy_scores
110 test_df
111
112 #Reporting Scores
113 from sklearn.metrics import classification_report, confusion_matrix
114 target_names = ['Neg', 'Positive']

```

```

112 print(classification_report(test_df['label'], test_df['pred'],
    target_names=target_names))

```

Listing A.5: LSTM

A.5 roBERTa

```

1 from datasets import load_dataset
2 from transformers import AutoTokenizer
3 from transformers import AutoModelForSequenceClassification
4 import numpy as np
5 import evaluate
6 from transformers import TrainingArguments, Trainer
7 import pandas as pd
8
9
10 def adjustLabels(label):
11     return label + 1
12
13
14 def findMaxLen(dataset):
15     maxLen = 0
16     for sequence in dataset:
17         if len(sequence.split()) > maxLen:
18             maxLen = len(sequence.split())
19     return maxLen
20
21
22 def tokenizeFunction(data):
23     return tokenizer(data['text'], padding="max_length", truncation=True,
24                     max_length=maxLen)
25
26
27 def computeMetrics(eval_pred):
28     logits, labels = eval_pred
29     # setting up all logits corresponding to neutral to nan in order to filter them
30     # out
31     logits[:, 1] = np.nan
32     predictions = np.nanargmax(logits, axis=-1)
33     return metric.compute(predictions=predictions, references=labels)
34
35
36 # adjusting labels to 0 -> negative and 2 -> positive
37 df = pd.read_csv("data/cleaned_tweets.csv")
38 df['label'] = df['label'].apply(adjustLabels)
39 df.to_csv("data/tweets_roberta.csv", index=False)
40
41
42 dataset = load_dataset('csv', data_files="data/tweets_roberta.csv", split='train')
43 maxLen = findMaxLen(dataset['text'])
44
45
46 dataset = dataset.train_test_split(test_size=0.2, shuffle=True, seed=1)
47 tokenizer =
48     AutoTokenizer.from_pretrained("cardiffnlp/twitter-roberta-base-sentiment-latest")
49 tokenized_datasets = dataset.map(tokenizeFunction, batched=True)
50
51
52 train_dataset = tokenized_datasets['train']
53 eval_dataset = tokenized_datasets['test']
54
55
56 model =
57     AutoModelForSequenceClassification.from_pretrained("cardiffnlp/twitter-roberta-
58                                                         base-sentiment-latest")

```

```

51 training_args = TrainingArguments(output_dir="training_checkpoints",
    evaluation_strategy="epoch")
52
53 metric = evaluate.load('accuracy')
54
55 trainer = Trainer(
56     model=model,
57     args=training_args,
58     train_dataset=train_dataset,
59     eval_dataset=eval_dataset,
60     compute_metrics=computeMetrics,
61 )
62
63 trainer.train()
64 trainer.save_model('tuned_model')
65 tokenizer.save_pretrained('tuned_model')

```

Listing A.6: roBERTa training

```

1 from transformers import AutoModelForSequenceClassification
2 from transformers import AutoTokenizer, AutoConfig
3 import numpy as np
4 from scipy.special import softmax
5 import re
6
7
8 def preprocess(text):
9     text = re.sub(r'https?\S+', 'http', text)
10    text = re.sub(r'@\S+', '@user', text)
11    text = re.sub(r'\d+', '@number', text)
12    text = re.sub(r'[$] [A-Za-z]\S*', '@stock', text)
13    text = re.sub(r'[^A-Za-z0-9@]', ' ', text)
14    text = re.sub(r' +', ' ', text)
15    text = text.lower()
16    text = text.strip()
17    return text
18
19
20 def getScore(text):
21     text = preprocess(text)
22     encoded_input = tokenizer(text, return_tensors='pt')
23     output = model(**encoded_input)
24     scores = output[0][0].detach().numpy()
25     scores = softmax(scores)
26     scores[1] = 0
27     scores = scores / np.sum(scores)
28     return scores[2]
29
30
31 MODEL = f"tuned_model"
32 tokenizer = AutoTokenizer.from_pretrained(MODEL)
33 config = AutoConfig.from_pretrained(MODEL)
34 model = AutoModelForSequenceClassification.from_pretrained(MODEL)

```

Listing A.7: roBERTa sentiment score

A.6 Results

```

1 import numpy as np
2 import pandas as pd
3 import datetime as dt

```

```

4 import re
5 import snsrape.modules.twitter as sntwitter
6 from joblib import Parallel, delayed
7 from roberta import getScore
8 import time
9
10
11 def checkSymbols(symbols):
12     check = True
13     for symbol in symbols:
14         if symbols[0] != symbol:
15             check = False
16             break
17     return check
18
19
20 def historicalSentiment(i):
21     stopDate = dt.datetime.strptime(dates[i], '%Y-%m-%d') + dt.timedelta(days=int(1))
22     stopDate = stopDate.strftime('%Y-%m-%d')
23     startDate = dates[i]
24
25     scores = np.zeros(2 * len(symbols))
26     for j in range(len(symbols)):
27         query = '$' + symbols[
28             j] + ' lang:en since:' + startDate + ' until:' + stopDate + '
29             -filter:links -filter:replies'
30         symbolScores = np.array([])
31         for k, tweet in enumerate(sntwitter.TwitterSearchScraper(query).get_items()):
32             if len(symbolScores) >= 100:
33                 break
34             if checkSymbols(re.findall(r'[$][A-Za-z]\S*', tweet.content)):
35                 symbolScores = np.append(symbolScores, getScore(tweet.content))
36
37         scores[2 * j] = np.average(symbolScores) if len(symbolScores) > 0 else 0
38         scores[2 * j + 1] = len(symbolScores)
39     return scores
40
41
42 symbols = pd.read_csv("data/S&P100_symbols.csv")['Symbol']
43
44 columns = []
45 for i in range(len(symbols)):
46     columns.append(symbols[i])
47     columns.append(symbols[i] + ' TweetCount')
48
49 date = dt.datetime.strptime('2022-02-02', '%Y-%m-%d')
50 dates = []
51 for i in range(365):
52     if date.weekday() < 5:
53         dates.append(date.strftime('%Y-%m-%d'))
54     date = date - dt.timedelta(days=int(1))
55
56 start = time.time()
57
58 output = pd.DataFrame(columns=columns, index=dates)
59 results = np.array(Parallel(n_jobs=2)(delayed(historicalSentiment)(i) for i in
60     range(2)))
61
62 output = pd.DataFrame(results, columns=columns, index=dates[:results.shape[0]])
63 output.index.name = 'Date'
64 output.to_csv("data/historical_sentiment.csv")

```

```

65 end = time.time()
66 print('Time elapsed:', round((end - start) / 60, 2))

```

Listing A.8: Historical sentiment

```

1  import pandas as pd
2  import numpy as np
3  import matplotlib.pyplot as plt
4  import matplotlib.dates as mdates
5  import datetime as dt
6
7
8  def benchmark(indexDateStart, indexDateStop):
9      stockWallet = np.zeros(prices.shape)
10     wealth = np.array([1])
11
12     for i in range(indexDateStop - indexDateStart):
13         date = dates[indexDateStart + i]
14         if date not in sentiment.index:
15             break
16
17         if i > 0:
18             wealth = np.append(wealth, np.sum(stockWallet[i - 1] *
19                                                 prices.loc[date].to_numpy()))
19
20         for j in range(stockWallet.shape[1]):
21             stockWallet[i, j] = wealth[i] / stockWallet.shape[1] / prices.loc[date,
22                                     prices.columns[j]]
23
24     return wealth
25
26 def backtestAbsolute(numStocks, minTweetCount, indexDateStart, indexDateStop,
27     startWealth):
28     stockWallet = np.zeros(prices.shape)
29     wealth = np.array([startWealth])
30
31     for i in range(indexDateStop - indexDateStart):
32         date = dates[indexDateStart + i]
33         if date not in sentiment.index:
34             break
35
36         if i > 0:
37             wealth = np.append(wealth, np.sum(stockWallet[i - 1] *
38                                                 prices.loc[date].to_numpy()))
39
40         tweetCount = sentiment.loc[date].to_numpy()[1::2]
41         sentimentScores = sentiment.loc[date].to_numpy()[::2]
42         sentimentScores = np.where(tweetCount > minTweetCount, sentimentScores, 0)
43
44         if np.count_nonzero(sentimentScores) < numStocks:
45             numStocks = np.count_nonzero(sentimentScores)
46
47         sortedTickers = np.argsort(-sentimentScores)
48
49         for j in range(stockWallet.shape[1]):
50             if np.where(sortedTickers == j)[0][0] < numStocks:
51                 stockWallet[i, j] = wealth[i] / numStocks / prices.loc[date,
52                                         prices.columns[j]]
53
54     if indexDateStop < len(dates):
55         finalWealth = np.sum(stockWallet[indexDateStop - indexDateStart - 1] *
56                               prices.loc[dates[indexDateStop]].to_numpy())

```

```

53     else:
54         finalWealth = 0
55     return wealth, finalWealth
56
57
58 def backtestRelative(numStocks, minTweetCount, indexDateStart, indexDateStop,
59                     startWealth):
60     stockWallet = np.zeros(prices.shape)
61     wealth = np.array([startWealth, startWealth])
62
63     for i in range(1, indexDateStop - indexDateStart):
64         datePrevious = dates[indexDateStart + i - 1]
65         date = dates[indexDateStart + i]
66         if date not in sentiment.index:
67             break
68
69         if i > 1:
70             wealth = np.append(wealth, np.sum(stockWallet[i - 1] *
71                                               prices.loc[date].to_numpy()))
72
73         if i > 0:
74             tweetCountPrevious = sentiment.loc[datePrevious].to_numpy()[1::2]
75             tweetCount = sentiment.loc[date].to_numpy()[1::2]
76             sentimentScoresPrevious = sentiment.loc[datePrevious].to_numpy()[::2]
77             sentimentScores = sentiment.loc[date].to_numpy()[::2]
78
79             sentimentScores = np.divide(sentimentScores - sentimentScoresPrevious,
80                                       sentimentScoresPrevious,
81                                       out=np.zeros_like(sentimentScores -
82                                                         sentimentScoresPrevious),
83                                       where=sentimentScoresPrevious != 0)
84             sentimentScores = np.where(np.logical_and(tweetCountPrevious >
85                                                         minTweetCount, tweetCount > minTweetCount),
86                                       sentimentScores, 0)
87
88             if np.count_nonzero(sentimentScores) < numStocks:
89                 numStocks = np.count_nonzero(sentimentScores)
90
91             sortedTickers = np.argsort(-sentimentScores)
92
93             for j in range(stockWallet.shape[1]):
94                 if np.where(sortedTickers == j)[0][0] < numStocks:
95                     stockWallet[i, j] = wealth[i] / numStocks / prices.loc[date,
96                                       prices.columns[j]]
97
98         if indexDateStop < len(dates):
99             finalWealth = np.sum(stockWallet[indexDateStop - indexDateStart - 1] *
100                                prices.loc[dates[indexDateStop]].to_numpy())
101         else:
102             finalWealth = 0
103     return wealth, finalWealth
104
105 def optimise(numStockRange, minTweetCountRange, indexDateStart, indexDateStop,
106             method):
107     xdata = np.zeros(numStockRange * minTweetCountRange)
108     ydata = np.zeros(numStockRange * minTweetCountRange)
109     zdata = np.zeros(numStockRange * minTweetCountRange)
110
111     bestNumStock = 0
112     bestMinTweetCount = 0
113     bestWealth = 0

```

```

109     for i in range(1, numStockRange + 1):
110         for j in range(1, minTweetCountRange + 1):
111             wealth = np.average(method(i, j, indexDateStart, indexDateStop, 1)[0])
112
113             xdata[(i - 1) * numStockRange + j - 1] = j
114             ydata[(i - 1) * numStockRange + j - 1] = i
115             zdata[(i - 1) * numStockRange + j - 1] = wealth
116
117             if wealth > bestWealth:
118                 bestNumStock = i
119                 bestMinTweetCount = j
120                 bestWealth = wealth
121
122             # fig = plt.figure()
123             # ax = fig.add_subplot(projection='3d')
124             # ax.scatter(xdata, ydata, zdata - 1)
125             # ax.set_xlabel('Minimal number of tweets')
126             # ax.set_ylabel('Number of traded stocks')
127             # ax.set_zlabel('Revenue [%]')
128             # fig.tight_layout()
129             # plt.show()
130
131         return bestNumStock, bestMinTweetCount, bestWealth
132
133
134 sentiment = pd.read_csv("data/historical_sentiment.csv", index_col='Date')
135 prices = pd.read_csv("data/historical_prices.csv", index_col='Date')
136 sentiment = sentiment.iloc[:, :-1]
137 prices = prices.iloc[:, :-1]
138
139 dates = list(prices.index)
140
141 split = 0.8
142 folds = 10
143 lenInSample = round(len(dates) * split / folds / (1 - split + split / folds))
144 lenOutSample = round(lenInSample * (1 - split) / split)
145
146 wealthAbsolute = np.array([])
147 wealthRelative = np.array([])
148 finalWealthAbsolute = 1
149 finalWealthRelative = 1
150 for i in range(folds):
151     print("Progress: " + str(round(i / folds * 100, 2)) + "%")
152
153     numStockAbsolute, minTweetCountAbsolute = optimise(20,
154                                                         60,
155                                                         lenOutSample * i,
156                                                         lenOutSample * i + lenInSample,
157                                                         backtestAbsolute)[:2]
158     numStockRelative, minTweetCountRelative = optimise(20,
159                                                         60,
160                                                         lenOutSample * i,
161                                                         lenOutSample * i + lenInSample,
162                                                         backtestRelative)[:2]
163
164     appendedWealthAbsolute, finalWealthAbsolute = backtestAbsolute(numStockAbsolute,
165                                                                     minTweetCountAbsolute,
166                                                                     lenOutSample * i +
167                                                                     lenInSample,
168                                                                     lenOutSample * i +
169                                                                     lenInSample +
170                                                                     lenOutSample,
171                                                                     finalWealthAbsolute)

```



```

169
170     appendedWealthRelative, finalWealthRelative = backtestRelative(numStockRelative,
171                                                                    minTweetCountRelative,
172                                                                    lenOutSample * i +
                                                                    lenInSample,
173                                                                    lenOutSample * i +
                                                                    lenInSample +
                                                                    lenOutSample,
174                                                                    finalWealthRelative)
175
176     wealthAbsolute = np.hstack((wealthAbsolute, appendedWealthAbsolute))
177     wealthRelative = np.hstack((wealthRelative, appendedWealthRelative))
178
179     plotDates = list(prices.index)[lenInSample:]
180     x = [dt.datetime.strptime(d, '%Y-%m-%d').date() for d in plotDates]
181     wealthBenchmark = benchmark(lenInSample, len(dates))
182
183     dailyReturnsAbsolute = np.append([0], np.diff(wealthAbsolute) / wealthAbsolute[:-1])
184     dailyReturnsRelative = np.append([0], np.diff(wealthRelative) / wealthRelative[:-1])
185     dailyReturnsBenchmark = np.append([0], np.diff(wealthBenchmark) /
186                                         wealthBenchmark[:-1])
187
188     print("Cumulative return absolute [%]:", (wealthAbsolute[-1] - 1) * 100)
189     print("Cumulative return relative [%]:", (wealthRelative[-1] - 1) * 100)
190     print("Cumulative return benchmark [%]:", (wealthBenchmark[-1] - 1) * 100, '\n')
191
192     print("Annualised return absolute [%]:", (wealthAbsolute[-1] ** (252 / len(x)) - 1) *
193           100)
194     print("Annualised return relative [%]:", (wealthRelative[-1] ** (252 / len(x)) - 1) *
195           100)
196     print("Annualised return benchmark [%]:", (wealthBenchmark[-1] ** (252 / len(x)) - 1)
197           * 100, '\n')
198
199     print("Annualised volatility absolute [%]:", np.std(dailyReturnsAbsolute) *
200           np.sqrt(252) * 100)
201     print("Annualised volatility relative [%]:", np.std(dailyReturnsRelative) *
202           np.sqrt(252) * 100)
203     print("Annualised volatility benchmark [%]:", np.std(dailyReturnsBenchmark) *
204           np.sqrt(252) * 100, '\n')
205
206     riskFreeRate = 0.0421 / 10 / 252
207     print("Sharpe ratio absolute:",
208           np.sqrt(252) * (np.average(dailyReturnsAbsolute) - riskFreeRate) /
209           np.std(dailyReturnsAbsolute))
210     print("Sharpe ratio relative:",
211           np.sqrt(252) * (np.average(dailyReturnsRelative) - riskFreeRate) /
212           np.std(dailyReturnsRelative))
213     print("Sharpe ratio benchmark:",
214           np.sqrt(252) * (np.average(dailyReturnsBenchmark) - riskFreeRate) /
215           np.std(dailyReturnsBenchmark), '\n')
216
217     print("Maximum drawdown absolute [%]:",
218           (np.min(wealthAbsolute) - np.max(wealthAbsolute)) / np.max(wealthAbsolute) *
219           100)
220     print("Maximum drawdown relative [%]:",
221           (np.min(wealthRelative) - np.max(wealthRelative)) / np.max(wealthRelative) *
222           100)
223     print("Maximum drawdown benchmark [%]:",
224           (np.min(wealthBenchmark) - np.max(wealthBenchmark)) / np.max(wealthBenchmark) *
225           100)
226
227     plt.gca().xaxis.set_major_formatter(mdates.DateFormatter('%Y-%m-%d'))

```

```

215 plt.gca().axis.set_major_locator(mdates.DayLocator(interval=round((len(dates) -
    lenOutSample) / 7)))
216
217 plt.plot(x, (wealthAbsolute - 1) * 100, label="Absolute")
218 plt.plot(x, (wealthRelative - 1) * 100, label="Relative")
219 plt.plot(x, (benchmark(lenInSample, len(dates)) - 1) * 100, label="Benchmark")
220 plt.axhline(y=0, color='black', linestyle='--')
221
222 # for i in range(folds):
223 #     plt.axvline(x=x[lenOutSample * i], color='b')
224
225 plt.legend(loc="lower left")
226 plt.xlabel('Date')
227 plt.ylabel('Revenue [%]')
228
229 plt.tight_layout()
230 plt.gcf().autofmt_xdate()
231 plt.savefig('plots/revenue.pdf', bbox_inches='tight')
232 plt.show()

```

Listing A.9: Backtesting