# HIGH-FREQUENCY CRYPTOCURRENCY TRADING STRATEGY USING TWEET SENTIMENT ANALYSIS

Zhijun Chen

Department of Financial Engineering, SUSTech University, Shen Zhen, China

## ABSTRACT

*Sentiments are extracted from tweets with the hashtag of cryptocurrencies to predict the price and sentiment prediction model generates the parameters for optimization procedure to make decision and re-allocate the portfolio in the further step. Moreover, after the process of prediction, the evaluation, which is conducted with RMSE, MAE and R2, select the KNN and CART model for the prediction of Bitcoin and Ethereum respectively. During the process of portfolio optimization, this project is trying to use predictive prescription to robust the uncertainty and meanwhile take full advantages of auxiliary data such as sentiments. For the outcome of optimization, the portfolio allocation and returns fluctuate acutely as the illustration of figure.*

## KEYWORDS

*Cryptocurrency Trading Portfolio, Sentiment Analysis, Machine Learning, Predictive Prescription, Robust Optimization Portfolio.*

## 1. INTRODUCTION

As a decentralized digital asset, cryptocurrency does not exist as physical entity like paper money but secures transaction, controls creation by using strong cryptography [1]. The security and peer-to-peer benefits give bitcoin and many other different types of cryptocurrencies popularity and their markets quickly prosper. Started from 2008, a paper titled "Bitcoin: A Peer-to-Peer Electronic Cash System" marked its inception [2]. By 2017, the price of a single Bitcoin soared 2000% from $863 to $17,550 [3]. Although bubble exists [4], nowadays, the two largest cryptocurrencies, Bitcoin and Ethereum, had a 287.2 billion dollars market capitalization by October 2020, and Bitcoin alone shared 240.6 billion dollars [5].

In this paper, the way of predicting the price of cryptocurrency is by the result of sentiment analysis from tweets which are the short messages in a concise format published in social media platform – Twitter. On average, five hundred million tweets are sent each and every day, and around two hundred billion tweets per year [6]. Such a massive dataset helps this project for the sentiment analysis a lot.

Sentiment analysis or opinion mining combines the usage of natural language processing (NLP), text analysis, computational linguistics, and biometrics to analyse public emotion or preference in the area ranging from marketing to customer survey to recommender system [7][8][9][10]. One of trending task for sentiment analysis is classifying peoples' emotion into positive, negative or neural to analyse public views towards different topics on social network [11]. In this project,

sentiment analysis is utilized to generate prediction for the price of cryptocurrency, which then becomes the input of optimization model to work out the capital allocation strategy.

As for the decision, portfolio optimization is another important part in this paper, which is conducted in a prescriptive method. The prescriptions ensure the ability that allocates capital in the robust decision and accommodates the uncertainty in the real world by integrating operations research and management science with machine learning and utilizing both the auxiliary data and the data predicted by machine learning in the process of optimization [12].

## 2. LITERATURE REVIEW

Two main process of this paper – prediction via sentiment analysis and portfolio optimization both have a wide range of related former topics and researches in the financial field.

Emotion affects individual capital allocation strategy according to Amos Tversky and Daniel Kahneman as early as 1979 [13]. After decades of exploration of many great behavioural economists, Paul Tetlock concludes the negative correlation between pessimistic cognition and activeness in stock market [14]. Later, Galen Thomas Panger's research in 'Emotion in Social Media' further bridges the standpoints of behavioural economists and network platform such as social media [15]. Thus, extracting public sentiments for cryptocurrency from tweets in Twitter by data mining and sentiment analysis helps to predict traders' decision and then to predict the price. According to the continuity and time limited of public sentiment, a trading strategy with 2.5 million tweets is established by Hong Kee Sul, Alan R Dennis, and Lingyao Ivy Yuan and produces 11-15% annual returns with a good prospect [16]. Moreover, Y. B. Kim discovers the potential price fluctuation for Bitcoin due to user sentiment [17]. Given by these researches, the topic sentiment analysis for cryptocurrency price prediction is effectively practical.

As for the related works in the field of portfolio optimization, in early 1952, Harry Markowitz formulates the well-known portfolio selection model or mean-variance model [18]. In this model, the portfolio return and risk get measured by expected value and variance respectively. Then, the strategy selection problem is converted into an optimization problem. Based on the key idea of mean-variance model, the Capital Asset Pricing Model (CAPM) is created by William Sharpe in 1964 [19]. These models have a very profound impact on robust portfolio optimization framework, which assumes the worst-casebehaviour faced with unknown parameters or market perturbations [20]. Nowadays, robust portfolio optimization is a wise choice to challenge with parameter uncertainty and estimation errors in portfolio management and to find the optimal portfolio over a basket of cryptocurrencies or the other financial securities with a limited risk.

## 3. RESEARCH METHODOLOGY

In this project, prediction and decision are conducted for Bitcoin and Ethereum, which are two popular cryptocurrencies currently. Six sections comprise the integrity of the trading strategy: Dataset compilation, Sentiment Analysis, Prediction, Measurement, Optimization and Visualization.

### 3.1. Dataset Compilation

Cryptocurrency price dataset or Open High Low Close Volume (OHLCV) is collected from CoinAPI via coinapi-sdk, while tweet dataset collected from tweets with the hash tag of the certain cryptocurrency is compiled from Tweet Archivist started from December twelfth, 2019 to February twelfth, 2020. Coordinate these two types of data set together with respect to the same

time interval in hours. That means each row item in the combined data set has both OHLCV and all tweets published in one hour correspondingly.

## 3.2. Sentiment Analysis and Imputation

Sentiment analysis and imputation comprise the further data wrangling.

Firstly, sentiment analysis is to figure out the emotional standpoint of tweets. Tweet text attribute is converted to a bunch of floating-number attributes named 'pos', 'neg', 'neu', and 'compound' which respectively represent the positive, negative, neutral viewpoints or emotional tendencies towards certain kind of cryptocurrency. During this process, one general sentiment package nltk.sentiment.vader is utilized. In the package, the model Sentiment Intensity Analyzer is already trained and performs well in daily dialog dataset. A more specialized model for financial or cryptocurrency is regarded to the future work for the project. However, to some extent, the general model fits the cryptocurrency when it comes to casual communication tweets. For example, the sentence "Bitcoin is awesome" has very high compound value, which is 0.6249, while "BAD NEWS FOR #BITCOIN" is totally contrary with a -0.5423-compound value. And a question "I am wondering how people trade on bitcoin" remains neutral, that is 0.0 for compound value.

Next, all the sentiment attributes for tweets are calculated the mean values contributing to the sentiment features for certain hour period which is one of the entries in the final dataset.

Secondly, due to the case that there are some time intervals without tweets for example midnight, it needs imputation for those sentiment attributes to avoid the empty entry. By the assumption of consistency emotion of society, the previous sentiment is pasted to fill the empty row.

Finally, the attribute information of the modified dataset is as follows:

Table 1.  Attribute Information for Dataset.

| Field Name | Count | Non-Null | Data Type |
|---|---|---|---|
| Period start | 24 | Non-null | Object |
| Period end | 24 | Non-null | Object |
| Time open | 24 | Non-null | Object |
| Time close | 24 | Non-null | Object |
| Price open | 24 | Non-null | Float64 |
| Price close | 24 | Non-null | Float64 |
| Price low | 24 | Non-null | Float64 |
| Price high | 24 | Non-null | Float64 |
| Volume traded | 24 | Non-null | Float64 |
| Trades count | 24 | Non-null | Int64 |
| Change percentage | 24 | Non-null | Float64 |
| Compound | 24 | Non-null | Float64 |
| Pos | 24 | Non-null | Float64 |
| Neg | 24 | Non-null | Float64 |
| Neu | 24 | Non-null | Float64 |

### 3.3. Prediction via Regressor Models

Split the dataset into train set and test set, where regression models are trained and then produce the prediction. During the process of prediction, grid search is used to select the best parameters for predicting the 'Change percentage' target, while the other attributes are the inputs of the prediction models. Four prediction models are used, which are random forest, k-nearest neighbours (KNN), classification and regression tree (CART), and Lasso regressor.

Based on the idea of ensemble learning by training a batch of decision trees [21][22], the pseudo code for random forests or random decision forests is as below. By taking the mode or mean value from the set of trees, the overfitting problem gets released due to avoiding only focusing on individual decision tree model [23].

---

Algorithm 1: Pseudo code for the random forest algorithm **Error! Reference source not found.**

---

To generate $c$ classifiers:
**for** $i = 1$ to $c$ **do**
  Randomly sample the training data $D$ with replacement to produce $D_i$
  Create a root node, $N_i$ containing $D_i$
  Call BuildTree($N_i$)
**end for**

**BuildTree(N):**
**if** $N$ contains instances of only one class **then**
**return**
**else**
  Randomly select x% of the possible splitting features in $N$
  Select the feature $F$ with the highest information gain to split on
  Create f child nodes of $N$, $N_1, …, N_f$, where $F$ has $f$ possible values ($F_1, …, F_f$)
**for** $i = 1$ to $f$ **do**
     Set the contents of $N_i$ to $D_i$, where $D_i$ is all instances in $N$ that match
$$F_i$$
     Call BuildTree($N_i$,)
**end for**
**end if**

---

As a non-parametric method introduced by Thomas Cover in the field of pattern recognition [25], KNN pseudo code is as follow. And the output is regarded as the average of k closet training samples [26].

---

Algorithm 2: Pseudo code for the KNN algorithm **Error! Reference source not found.**

---

Classify ($X,Y, x$) // $X$: training data, $Y$: class labels of $X$, $x$: unknown sample
**for** $i = 1$ to $m$ **do**
Compute distance $d(X_i,x)$
**end for**
Compute set $I$ containing indices for the $k$ smallest distances $d(X_i,x)$
**return** majority label for {$Y_i$ where $i \in I$}

---

Due to its precision and explicitness, which are demonstrated in the below pseudo code, decision tree or CART is one of the most widely used machine learning algorithm [28][29].

---

Algorithm 3: Pseudo code for the CART algorithm **Error! Reference source not found.**

---

d=0, endtree=0
Node(0)=1, Node(1)=0, Node(2)=0
**while** endtree<1
**if** Node($2^d$-1) + Node($2^d$) +…+ Node($2^{d+1}$-2) = 2-$2^{d+1}$
       endtree=1
**else**
       **do** i=$2^d$-1, $2^d$ ,…,$2^{d+1}$-2
       **if** Node(i)>-1
       **Split tree**
       **else**
          Node(2i+1)=-1
          Node(2i+2)=-1
       **end if**
       **end do**
 **end if**
d = d + 1
**end while**

---

Aimed to increase the accuracy of the statistical model, lasso regression works both on constructing model by a subset of relevant features and adding information to prevent overfitting and was firstly used in geophysics literature in 1986 [30].

Four machine learning models are used in the project, which are imported from sklearn package. They are Random Forest Regressor, KNeighbors Regressor, Decision Tree Regressor and linear_model.Lasso(). Regressor models, combination with grid searched parameters and cross validated training dataset, get trained by training dataset and then produce the prediction results for test dataset as shown below.
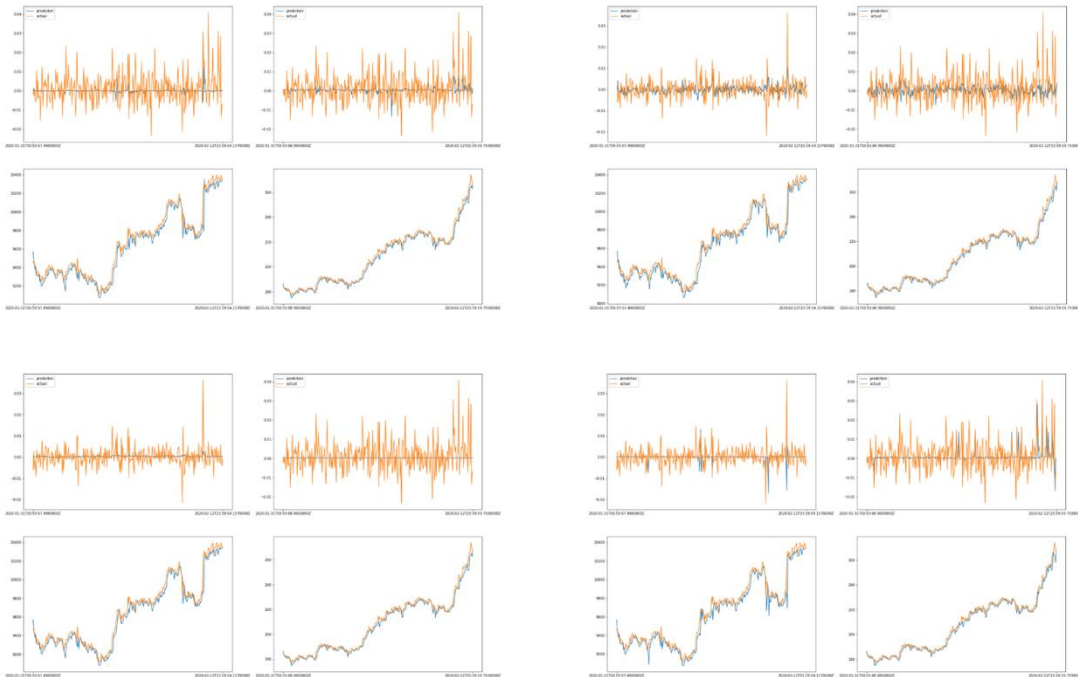
Figure 1.  Predictions of change percentage and close price

The sixteen charts in Figure 1 can be divided into four groups. With the writing direction of letter "z", from the position of top left to bottom right, there are the prediction results of random forests, KNN, CART, Lasso regressor respectively. Inside each group, with the same direction, there are the prediction of change percentage for Bitcoin and for Ethereum, the prediction of close price for Bitcoin and for Ethereum respectively. Inside each chart, prediction result, which is the blue line, and actual value, which is the orange one, are plotted in the time interval from January thirty-first, 2020 to February twelfth, 2020.

## 3.4. Measurement for Prediction

To measure the performance of these models and select the most accurate one except the others. Three types of measurements are taken into account, which are root-mean-square error (RMSE), mean absolute error (MAE), and coefficient of determination (denoted by $R^2$).

RMSE is the scale-dependent measurement which calculates the square root of the average of squared errors [31]. And the formula is

$$\text{RMSE} = \sqrt{\frac{\sum_{t=1}^{T}(\widehat{y_t}-y_t)^2}{T}},$$

Where $\widehat{y_t}$ denotes the predicted results for times t, $y_t$ represents the observed values for the same times, and their difference, also called error, is calculated quartic sum over T times [32].

MAE is the arithmetic average for the absolute error and commonly used in time series analysis [34]. With the above denotation, MAE is given by

$$\text{MAE} = \frac{\sum_{t=1}^{T}|\widehat{y_t}-y_t|}{T} \text{ Error! Reference source not found.[34].}$$

$R^2$ measures the degree of replication by the proportion of variation [35][36] and is equivalent to the explained sum of squares over the total sum of squares, that is

$$R^2 = \frac{SS_{reg}}{SS_{tot}},$$
$$SS_{reg} = \sum_{t=1}^{T} (\hat{y}_t - y_t)^2,$$
$$SS_{tot} = \sum_{t=1}^{T} (y_t - \bar{y})^2,$$

where $SS_{reg}$ is the quadratic sum of the difference between prediction and observed data, and $SS_{tot}$ is the squared sum of the difference between observed data and mean value [36].

After comparing the three measurements, RMSE, MAE and $R^2$ for four models used to predict the price respectively. Conclusion can be drawn to the case that KNN performs best for the price prediction of Bitcoin while CART predicts the price of Ethereum most accurately, which both have lowest RMSE, MAE and highest $R^2$ (see Table 2 and Table 3)

Table 2. Bitcoin Prediction Measurement.

|  | RMSE | MAE | $R^2$ |
|---|---|---|---|
| Random forests | 0.008249 | 0.005964 | -0.05151 |
| KNN | 0.008025 | 0.006007 | 0.004652 |
| CART | 0.008293 | 0.005971 | -0.06280 |
| Lasso regression | 0.008070 | 0.005960 | -0.006503 |

Table 3. Ethereum Prediction Measurement.

|  | RMSE | MAE | $R^2$ |
|---|---|---|---|
| Random forests | 0.007977 | 0.005954 | 0.01649 |
| KNN | 0.008320 | 0.006271 | -0.06975 |
| CART | 0.007717 | 0.005747 | 0.07954 |
| Lasso regression | 0.008105 | 0.005967 | -0.0152 |

As data shows, some of the prediction methods does not perform effectively with respect to their negative $R^2$ scores. These may be due to the high volatility of the cryptocurrency market and frequent fluctuation of the cryptocurrency price which decrease the accuracy of prediction.

## 3.5. Portfolio Optimization

Before formulating the strategy, covariance matrix between the history price of Bitcoin and Ethereum needs to be introduced. According to

$$\Sigma = \begin{bmatrix} var(B) & cov(B,E) \\ cov(E,B) & var(E) \end{bmatrix} [37],$$

where $B$ is for the value of Bitcoin, $E$ is for the one of Ethereum, $var$ calculates the variance and $cov$ calculates the covariance.

As for the objective function and constraints, according to Bertsimas work in "From Predictive to Prescriptive Analytics", objective function or predictive prescription is given by

$$\hat{z}_N(x) \in \arg\min_{z \in Z} \sum_{i=1}^{N} w_N^i(x) c(z; y^i) [12],$$

while the constraints – robust maximum return formulation in the worst case are

$$\text{maximize} \min_{\{\mu \in S_m\}} \mathrm{E}[r_\phi],$$
$$\text{subject to} \max_{\{V \in S_v, D \in S_d\}} \mathrm{Var}[r_\phi] \leq \lambda,$$
$$1^T \phi = 1 \ [20].$$

In this project, risk $\lambda$ is taken from the max one of variance of Bitcoin and variance of Ethereum. Transaction cost, which is the extra manipulation cost for re-allocation the capital, is 0.5%. The perception allocation percentage is half and half, that is holding Bitcoin and Ethereum equal amount at the beginning time.

## 3.6. Visualization of Strategy

The results show that all the capital is allocated to buy Bitcoin or Ethereum at every time period as the Figure 2 shows. It seems the capital basket is too monotonous in this case. And the returns for the portfolio tend to soar at some time periods as the illustration of Figure 3.
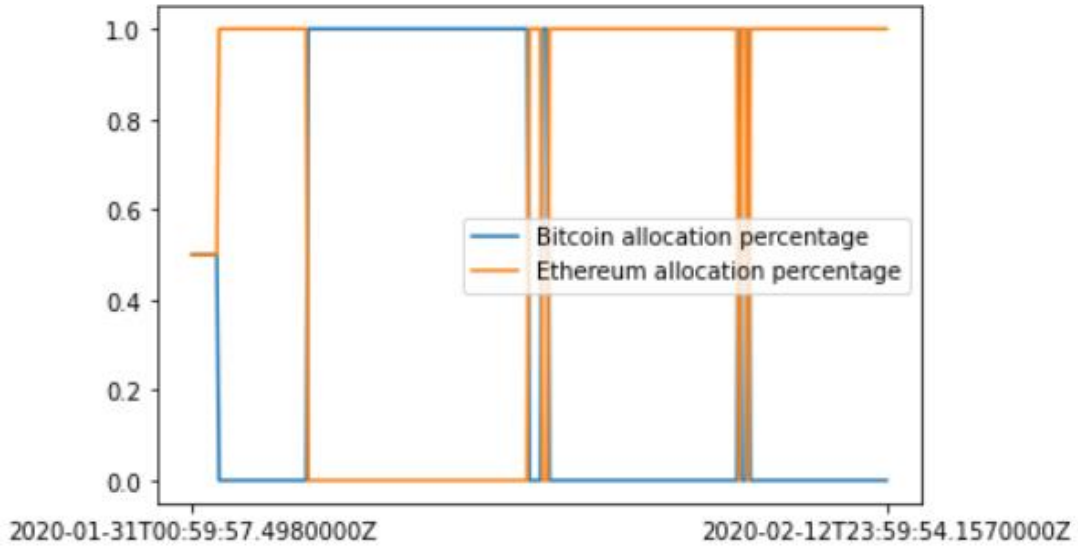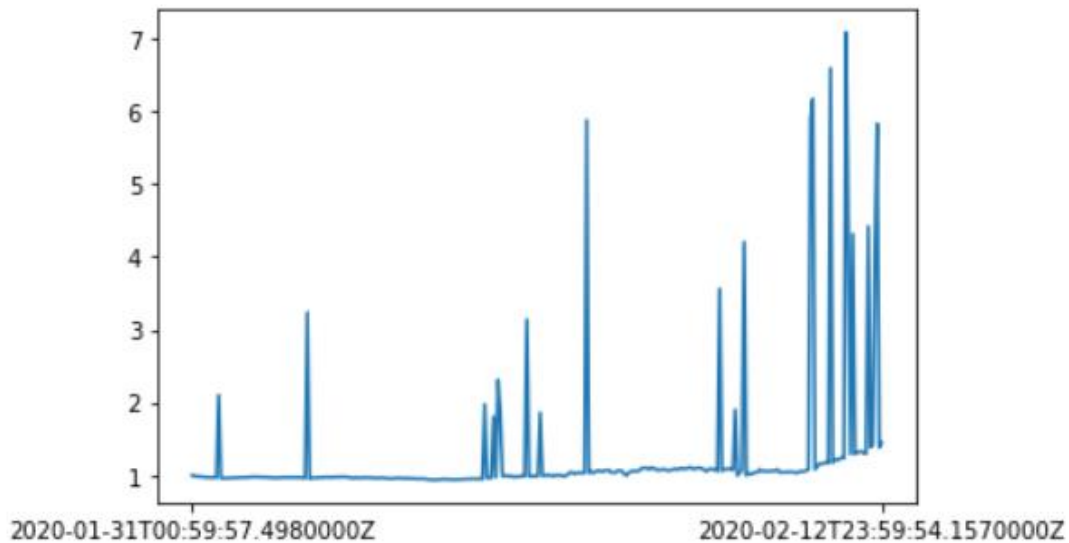


Figure 2.  Capital allocation for portfolio

Figure 3. Returns for portfolio

Within the timespan of half month, returns fluctuated acutely with even greater than 7% at some points, but coming down to no more than 1.5% at the end.

## 4. DISCUSSION AND CONCLUSION

The combination of sentiment features and predictive prescription helps to collect auxiliary data from social network and produce an uncertainty-robust portfolio strategy to help capital allocation in the cryptocurrency market.

More work can be done in the future both in the predictive part and prescriptive analysis part, such as replace the general sentiment model by a specific financial-target one, change the risk to a wider range, select more perception allocation percentage, and try different transaction cost, which are all the factors that can have an impact on the results of portfolio construction and final returns of this strategy. Moreover, more trading rules can be included such as short position.

REFERENCES

[1]    Greenberg, A. (2011) "CRYPTO CURRENCY-Money you can't trace". Forbes, 40.
[2]    Nakamoto, S. (2008)"Bitcoin: a peer-to-peer electronic cash system". Retrieved from https://bitcoin.org/bitcoin.pdf (accessed April 30, 2018).
[3]    Abraham, Jethin & Higdon, Daniel& Nelson, John& and Ibarra, Juan (2018) "Cryptocurrency Price Prediction Using Tweet Volumes and Sentiment Analysis," SMU Data Science Review: Vol. 1 : No. 3 , Article 1. Available at: https://scholar.smu.edu/datasciencereview/vol1/iss3/1
[4]    R. C. Phillips &D. Gorse, (2017) "Predicting cryptocurrency price bubbles using social media data and epidemic modelling", 2017 IEEE Symposium Series on Computational Intelligence (SSCI), pp. 1–7.
[5]    "Bit Info Charts" (2013). Available online at: https://bitinfocharts.com (accessed October 24, 2020).

[6]    "Internet Live Stats" (2011). Available online at: https://www.internetlivestats.com/twitter-statistics/ (accessed October 24, 2020).

[7]    Pang, Bo&Lee, Lillian & Vaithyanathan, Shivakumar (2002). "Thumbs up? Sentiment Classification using Machine Learning Techniques". Proceedings of the Conference on Empirical Methods in Natural Language Processing (EMNLP). pp. 79–86.

[8]    Thelwall, Mike& Buckley, Kevan& Paltoglou, Georgios & Cai, Di; Kappas, Arvid (2010). "Sentiment strength detection in short informal text". Journal of the American Society for Information Science and Technology. 61 (12): 2544–2558. CiteSeerX 10.1.1.278.3863. doi:10.1002/asi.21416.

[9]    Turney, Peter (2002). "Thumbs Up or Thumbs Down? Semantic Orientation Applied to Unsupervised Classification of Reviews". Proceedings of the Association for Computational Linguistics. pp. 417–424.

[10]   Korkontzelos, Ioannis & Nikfarjam, Azadeh & Shardlow, Matthew & Sarker, Abeed & Ananiadou, Sophia & Gonzalez, Graciela H. (2016). "Analysis of the effect of sentiment analysis on extracting adverse drug reactions from tweets and forum posts". Journal of Biomedical Informatics. 62: 148–158. doi:10.1016/j.jbi.2016.06.007. PMC 4981644. PMID 27363901.

[11]   R. Khan& H. U. Khan& M. S. Faisal&K. Iqbal&M. S. I. Malik, (2016)"An Analysis of Twitter users of Pakistan" Int. J. Comput. Sci. Inf. Secur., vol. 14, no. 8, p. 855.

[12]   Bertsimas, D. &Kallus, N., (2014)"From predictive to prescriptive analytics". arXiv preprint arXiv:1402.5481.

[13]   Kahneman, D. & Tversky, A. (1979)"Prospect theory: An analysis of decision under risk." Econometrica 47(2), pp263-291.

[14]   Tetlock, P.C.(2007)"Giving content to invsotry sentiment: The role of media in the stock market." The Journal of Finance.

[15]   Panger, G.T. (2017)"Emotion in Social Media". PhD thesis, University of California, Berkeley.

[16]   Hong Kee Sul& Alan R Dennis&Lingyao Ivy Yuan. (2016)"Trading on twitter: Using social media sentiment to predict stock returns". Decision Sciences.

[17]   Y. B. Kim& J. Lee, N. Park& J. Choo& J.-H. Kim&C. H. Kim, (2017)"When Bitcoin encounters information in an online forum: Using text mining to analyse user opinions and predict value fluctuation" PLOS ONE, vol. 12, no. 5, p. e0177630.

[18]   Harry Markowitz (1952)"Portfolio selection". The Journal of Finance, 7(1):77–91.

[19]   Sharpe, W. (1964). "Capital asset prices: A theory of market equilibrium under conditions of risk" J. Finance 19(3) 425–442.

[20]   Md. Asadujjaman & Kais Zaman (2014) "Robust Portfolio Optimization under Data Uncertainty" 15th National Statistical Conference, Dhaka, Bangladesh.

[21]   Ho, Tin Kam (1995). "Random Decision Forests" (PDF). Proceedings of the 3rd International Conference on Document Analysis and Recognition, Montreal, QC, 14–16 August 1995. pp. 278–282. Archived from the original (PDF) on 17 April 2016. Retrieved 5 June 2016.

[22]   Ho TK (1998). "The Random Subspace Method for Constructing Decision Forests" (PDF). IEEE Transactions on Pattern Analysis and Machine Intelligence. 20 (8): 832–844. doi:10.1109/34.709601.

[23]   Hastie, Trevor& Tibshirani, Robert& Friedman, Jerome (2008). "The Elements of Statistical Learning (2nd ed.)". Springer. ISBN 0-387-95284-5.

[24]   Guo, Hongquan & Nguyen, Hoang & Vu, Diep-Anh & Bui, Xuan-Nam. (2019). "Forecasting mining capital cost for open-pit mining projects based on artificial neural network approach." Resources Policy. 10.1016/j.resourpol.2019.101474.

[25]   Altman, Naomi S. (1992). "An introduction to kernel and nearest-neighbor nonparametric regression" (PDF). The American Statistician. 46 (3): 175–185. doi:10.1080/00031305.1992.10475879. hdl:1813/31637.

[26]   Piryonesi S. Madeh & El-Diraby Tamer E. (2020). "Role of Data Analytics in Infrastructure Asset Management: Overcoming Data Size and Quality Problems". Journal of Transportation Engineering, Part B: Pavements. 146 (2): 04020022.

[27]   Tay B, Hyun JK, Oh S.(2014)"A machine learning approach for specification of spinal cord injuries using fractional anisotropy values obtained from diffusion tensor images". Comput Math Methods Med. 2014; 2014:276589. doi: 10.1155/2014/276589. Epub 2014 Jan 21. PMID: 24575150; PMCID: PMC3918356.

[28]   Wu, Xindong & Kumar, Vipin& Ross Quinlan, J.& Ghosh, Joydeep & Yang, Qiang & Motoda, Hiroshi & McLachlan, Geoffrey J.& Ng, Angus& Liu, Bing & Yu, Philip S. & Zhou, Zhi-Hua

(2008). "Top 10 algorithms in data mining". Knowledge and Information Systems. 14 (1): 1–37. doi:10.1007/s10115-007-0114-2. ISSN 0219-3116. S2CID 2367747.

[29] Piryonesi S. Madeh & El-Diraby Tamer E. (2020-03-01). "Data Analytics in Asset Management: Cost-Effective Prediction of the Pavement Condition Index". Journal of Infrastructure Systems. 26 (1): 04019036. doi:10.1061/(ASCE)IS.1943-555X.0000512.

[30] Santosa, Fadil& Symes, William W. (1986). "Linear inversion of band-limited reflection seismograms". SIAM Journal on Scientific and Statistical Computing. SIAM. 7 (4): 1307–1330. doi:10.1137/0907087.

[31] Hyndman, Rob J.& Koehler, Anne B. (2006). "Another look at measures of forecast accuracy". International Journal of Forecasting. 22 (4): 679–688. CiteSeerX 10.1.1.154.9771. doi:10.1016/j.ijforecast.2006.03.001.

[32] "Coastal Inlets Research Program (CIRP) Wiki - Statistics" (2015). Retrieved 4 February 2015.

[33] Hyndman, R. and Koehler A. (2005). "Another look at measures of forecast accuracy" [1]

[34] Willmott, Cort J.& Matsuura, Kenji (2005). "Advantages of the mean absolute error (MAE) over the root mean square error (RMSE) in assessing average model performance". Climate Research. 30: 79–82. doi:10.3354/cr030079.

[35] Steel, R. G. D.&Torrie, J. H. (1960). "Principles and Procedures of Statistics with Special Reference to the Biological Sciences". McGraw Hill.

[36] Glantz, Stanton A.& Slinker, B. K. (1990). "Primer of Applied Regression and Analysis of Variance". McGraw-Hill. ISBN 978-0-07-023407-9.

[37] Park, Kun Il (2018). "Fundamentals of Probability and Stochastic Processes with Applications to Communications". Springer. ISBN 978-3-319-68074-3.

## AUTHORS

A student majored in financial engineering in SUS Tech. Trying to undertake the way in Quant for stock market.