

Modelling procedure while assessing the impact of news articles on cryptocurrency (Bitcoin) market movement

T.O. Maku¹, Monday Osagie Adenomon^{2*}, Mary U Adehi²

1. Department of Statistics, Federal University, Otuoke, Nigeria.

2. Department of Statistics, Nasarawa State University, Kefi, Nigeria.

(*Corresponding author: adenomonmo@nsuk.edu.ng, 
<https://orcid.org/0000-0002-9523-8032>)

Article Info	Abstract
<p>Original Article</p> <p>Main Object: Business & Economics, Data science, Statistics, computer science</p> <p>Received: 09 April 2025 Revised: 02 June 2025 Accepted: 02 June 2025 Published online: 13 July 2025</p> <p>Keywords: Bitcoin, CNBC's market section website, LDA, prediction, sLDA, Topic modelling.</p>	<p>Background: Cryptocurrencies have a variety of unique qualities, from cutting-edge technology to highly secure architecture. Additionally, the ability to invest in cryptocurrency, as an asset or a function of its prosperity has made crypto-currencies attractive to venture capitalists, computer scientists, and statisticians.</p> <p>Aims: In this study, we concentrated on a collection of documents web-scraped from the market section of CNBC, where each document is associated with a response variable.</p> <p>Methodology: These documents contain preprocessed words/terms of day-to-day reportage on cryptocurrency (Bitcoin). The corresponding response variables are the daily opening and closing price of Bitcoin prices. The Supervised Latent Dirichlet Allocation (sLDA), a statistical model of labeled documents, was used to analyze the textual data alongside their corresponding response variables, since our study aims to predict the response variable for unlabeled new documents.</p> <p>Results: Hidden Topics with their unique terms from the preprocessed articles were exposed through a Natural language processor. Mean absolute error (MAE), Mean absolute percentage error (MAPE), and Root mean square error (RMSE) graphs were constructed for the sLDA models with 'k = 3,10,20,30,50,75,100 and 200 Topics' values where the model with the best evaluation metric, was selected for prediction purpose.</p> <p>Conclusion: It was discovered that the sLDA model with k = 20. A posterior covariance matrix which shows the proportion of terms from the documents, making up a Topic. Coefficient values were generated in order to graphically visualize how important the discovered topics are and how they affect the market trend. Finally, the prediction of new labels (numeric-decoded closing prices) for the unlabeled documents was done and comparisons were made; the predicted labels follow a similar pattern to that of the time series closing price trend.</p>

Cite this article: Maku TO, Adenomon MO, Adehi MU. (2026). "Modelling procedure while assessing the impact of news articles on cryptocurrency (Bitcoin) market movement". *Cyberspace Studies*. 10(1): 23-42. doi: <https://doi.org/10.22059/jcss.2025.393177.1139>.



Creative Commons Attribution-NonCommercial 4.0 International License
Website: <https://jcss.ut.ac.ir/> | Email: jcss@ut.ac.ir |
EISSN: 2588-5502
Publisher: University of Tehran

1. Introduction

Cryptocurrencies have a variety of unique qualities, from cutting-edge technology to highly secure architecture. Additionally, the ability to invest in cryptocurrency, as an asset or a function of its prosperity has made crypto-currencies attractive to venture capitalists, computer scientists, and statisticians. Aside from all these aspects, there are other troubling issues such as the absence of financial institution regulation (Buchholz et al., 2012). Therefore, sentiments and ideologies have an impact on the movement of Bitcoin's price. The way that large-scale Bitcoin owners behave, as well as society norms, political beliefs, and emotional states, all influence the price of Bitcoin. Financial analysts and scholars have become increasingly interested in digital currency because of its widespread acceptance. Financial analysts and researchers enjoy the challenge of market price extrapolation, but the developing cryptocurrency marketplaces lack in-depth study (Yao et al., 2019).

The lack of this in-depth study prompted us to investigate other hidden factors contributing to the market price trend. In clear terms, we identified a factor which is textual crowd-trading-knowledge. Textual crowd-trading-knowledge is believed to have a strong root in articles, especially News articles. Through several mobile applications, those interested in Bitcoin transactions can obtain the latest information about price movements, causes, and news at any time and location. They could use this knowledge to decide whether to buy or sell. It has been shown by numerous studies that a significant portion of the public uses' websites, applications, or social networking sites to receive information in some capacity. According to Kaya and Karsligil (2010), web channels are now the second most important source of information for people who prefer reading news articles over watching or listening to the news. Television is the most important source of information.

The above resource is what we have taken advantage of to analyze and give it a meaning numerically. The way Bitcoin market price trend gets affected by the daily News articles on traders' activities was investigated. Eklund & Bejerholm (2004) postulated that sources of textual data include both independently developed and corporately produced materials. Sources generated by the organization, such as quarterly and yearly reports, can provide a rich language structure; when carefully examined, can forecast future performance.

Our chosen model in analyzing this phenomenon is a supervised Topic model known as the Supervised Latent Dirichlet Allocation (sLDA), simply because, our target is prediction. According to Blei et al. (2017), the sLDA is a variant of the Latent Dirichlet Allocation (LDA).

A supervised topic model is required for document collections that contain response variables where the objective is to predict the response variable given the document in the collection. The Supervised Latent

Dirichlet Allocation (sLDA), created by Blei et al. (2017), aims to infer latent themes that predict the response variable. First, they used a publicly available dataset of newspaper movie reviews introduced by Pang & Lee (2008), which contains movie reviews paired with a certain number of stars (ratings). The analysis was treated as a regression problem rather than a classification task. Next, they applied the sLDA to the analysis of two real-world tasks. The second application was to study amendment texts from the 109th and 110th U.S. senates. The discriminating parameters from an ideal-point analysis based on voting history which is utilized in quantitative political science in this case to map senators to a real-valued point on a political spectrum, make up the response variables (Clinton et al., 2004). The authors observed that the sLDA offered a better prediction on all datasets when compared to the results of linear regression and the Least Absolute Shrinkage and Selection Operator (LASSO) regression of the unsupervised LDA model.

Zang and Kjellström (2014) researched on the ways to supervise a topic model. The research suggested two factorized supervised topic models that factorize the topics into signal and noise. They also offered a detailed analysis of the behaviour of supervised topic models using Supervised Latent Dirichlet Allocation (SLDA). The findings of experiments conducted on synthetic and real-world data for computer vision tasks indicate that factorized topic models can improve performance and that increased supervision is necessary for optimal outcomes.

Survival-supervised latent Dirichlet allocation models were developed by Dawson and Kendziorski (2012) for the genomic study of time-to-event outcomes. With the use of their approach, groups of clinical and genomic characteristics that are common to patient subgroups can be efficiently identified, and each patient is then uniquely defined by these qualities. A useful patient subgroup was discovered by applying survLDA to The Cancer Genome Atlas (TCGA) ovarian research. These patient subgroups are distinguished by varying propensities to display aberrant mRNA expression and methylations, which correlate to varying rates of survival from first therapy. The study demonstrates how technological developments such as supervised topic modeling continue to improve the ease and precision of measuring the genome and phenome; as a result, genomic-based studies of the disease frequently involve the collection of a wide variety of data types from large patient populations.

Wilcox et al. (2021) researched on Supervised Latent Dirichlet Allocation with Covariates which is a Bayesian Structural and Measurement Model of Text and Covariates. They proposed a novel statistical model, supervised latent Dirichlet allocation with covariates (SLDAX) that jointly incorporates a latent variable measurement model of text and a structural regression model to allow the latent topics and

other manifest variables to serve as predictors of an outcome. Using a simulation study with data characteristics consistent with psychological text data, they found that SLDAX estimates were generally more accurate and more efficient. To illustrate the application of SLDAX and a two-stage approach, they provided an empirical clinical application to compare the application of both the two-stage and SLDAX approaches.

Perotte et al. (2011) presented a model for hierarchically and multiple labeled bag-of-word data called hierarchically supervised latent Dirichlet allocation (HSLDA). Their work focused mostly on out-of-sample label prediction, but it was also interesting to see better lower-dimensional representations of the bag-of-word data. Using large-scale data from retail product classification and healthcare document labeling tasks, they showcased HSLDA. They demonstrated that, as compared to models that do not use the structure from hierarchical labels, out-of-sample label prediction is much improved.

Mohan et al. (2019) gathered a significant quantity of time series data and used deep learning models to analyze it in connection to linked news stories, improving the accuracy of stock price predictions. The study shows that there is difficulty in forecasting stock values due to their extremely volatile character, which is influenced by a wide range of political and economic issues, shifts in leadership, investor attitude, and numerous other factors. The study was able to reveal a high association between the movement of stock prices and the release of news stories with the use of topic modeling. The quantity of training data supplied determined how accurate deep learning algorithms were. However, by compiling a sizable amount of time series data and applying deep learning models to analyze it for relevant news stories, they were able to increase the accuracy of stock price predictions.

Yap et al. (2012) used computer methods to forecast stock values using financial data. They relied on intricate mathematical models and historical market data, which was limited to evaluating data that was available to them. Because of this flaw, they were unable to respond to unforeseen circumstances that deviated from previous patterns. Weighted terms were assigned to a new piece by the study's prediction to ascertain its expected direction of movement. To their credit, these more straightforward techniques have demonstrated a limited but significant ability to forecast price direction but not actual price. The popularity of Quantitative funds, or Quants, has grown in the last few years. Quants automatically sort through financial data using numbers and select stocks. These systems vary in the degree of trading control they possess, from basic stock recommenders to transaction executors, despite being built on proprietary technology.

Sharma (2020) compared stock price prediction models using news articles, using global indices, exchange rates, historical stock prices, global news, and technical indicators. They utilized an ensemble of Long Short-Term Memory (LSTM) models. To compare the results,

their investigation was also expanded to include some benchmark categorization models. According to their results, an ensemble of three LSTM models predicted growing and falling trends equally well and steadily, yielding an accuracy of 60% with the best recall and true negative score.

Sahut et al. (2024) evaluated five models based on cutting-edge machine-learning techniques. These models were picked from the literature on crude oil forecasting to evaluate the impact of news-based sentiment on crude oil price prediction. The COVID-19 pandemic era and the years from 1990 to the start of the pandemic were included in the results for each approach. This made it possible for them to investigate how news-based mood functions in various stages of economic growth and disaster. Compared to other eras, such as the 2008–2009 financial crisis, a notable impact of news-based sentiment was seen on the forecasting performance of machine learning techniques during the Covid-19 period.

Loughran et al. (2019) investigated if investors respond to news stories about oil that disclose supply and demand information in a timely and logical manner. To facilitate investors' and researchers' assessment of the information value of oil tales, their study generated a unique keyword list of 130 terms and modifiers connected to the oil industry. They discovered a notable overreaction in the short run to the news stories on oil connected to the Dow Jones. Lower oil prices the next trading day were linked to phrases in delayed news items like output cut, production reduced, scarcity, and demand rising. The data supports the theory that oil traders exaggerate their reactions to news articles that are extensively read.

The major challenge of using structured textual data to predict traders' interest is because of its extremely volatile character of such data, which is influenced by a wide range of political and economic issues, shifts in leadership, investor attitude, and numerous other factors (Mohan et al., 2019; Fataliyev et al., 2021).

The baseline models for textual data are social media content, News content, official company announcements, traditional textual representation techniques, deep learning based advanced NLP techniques, statistical models, machine learning techniques and deep learning models (Frank et al., 2017; Li et al., 2018; Shah et al., 2019; Kumar et al., 2020; Thakkar & Chaudhari, 2021; Fataliyev et al., 2021).

Hence, our suggested text-feature extraction of the sLDA a supervised Latent Dirichlet Allocation aims to predict traders' interest in cryptocurrency and raise traders' awareness when they rely on news stories for trading information.

This paper contributes to existing literatures on the potentials and capabilities of sLDA (supervised Latent Dirichlet Allocation) to predict traders' interest in cryptocurrency.

2. Materials and Method

News articles on cryptocurrency-related activities published in foreign media between 2016 and 2022 were the research population's focus. The Consumer News and Business Channel is where these news stories are found (CNBC). Because Bitcoin is so popular compared to other cryptocurrencies, it will be the only cryptocurrency sampled in this study.

As secondary data, the corpus of news items serves as the data set for this study. From 2016 until 2022, every one of the more than 6,000 news pieces or articles was written in English. Leveraging a custom Python script called "Beautiful Soap" created with the Jupyter Notebook, quick scrapping of the text data was accomplished by leveraging relevant meta-data from the previously mentioned source. The following meta-data was included in the pages, which were kept in a comma-separated (CSV) format: (i) article summary, (ii) article section, (iii) article link, (iv) article date, (v) article summary, (vi) article body, (vii) opening price, (viii) closing price.

We used the query "Daily Bitcoin reports". Furthermore, we made use of the articles' corresponding closing prices as our response variable.

The bag-of-words document representation is assumed by topic models (Blei et al., 2017). Each document is represented here as a bag of its terms/words, with no respect for word order or grammar. Many Natural Language Processing and Information Retrieval algorithms use this simplified model. The NLP stage consisted primarily of four broad steps: (1) loading the input data (News articles), (2) pre-processing the data, (3) transforming texts into bag-of-words vectors, and (4) training the sLDA models (Mckinney, 2010).

The news items (which will now be referred to as documents) cannot be simply fed into the model as raw data or free-text but must instead be transformed into a suitable form for the modeling framework. Normalization, tokenization, stemming/lemmatization, and stop-word removal are common text data pre-processing techniques. Following the collection and collation of articles, the text will be pre-processed in Python using the SpaCy, Gensim, and Pandas modules (Mckinney, 2010). Pre-processing is required before performing NLP on the text. The texts of the articles were subsequently normalized by making them lowercase. Then, word elongations and foreign characters that weren't words, such as punctuation and other non-ASCII characters, were eliminated. Next, non-informative stop-words that regularly occur, such as "the", "is", "I", and "did" were eliminated (using stop-words provided in the gensim module of Python). Following that, token words were lemmatized in Python using the gensim module. Lemmatization is a type of text normalization that involves grouping inflected forms of words into their base or dictionary root terms, known as lemma. For example, lemmatizing the terms 'trouble', 'troubling', and 'troubled'

yields the lemma 'trouble'. The traditional stemming of tokens will be avoided, based on Schofield et al. (2017)'s recommendation that topic coherence is rarely enhanced between the pre-stemming and post-stemming Topic models. Finally, whitespaces were removed to make the document more compact. Documents containing fewer than 50 words were eliminated. Words that appeared in less than 70% of the corpus were pruned as well.

2.1. Specifications and Estimation of the sLDA

For response-document pairs, the supervised latent Dirichlet allocation model (sLDA) shows better strength in executing such a plan.

Distributions over document collections are known as topic models, in which each document is represented by a set of discrete random variables, $W_{1:n}$, which are its words. The words in a document are treated as emerging from a set of latent themes in topic models, which are a set of unknown distributions over the vocabulary. Each document in a corpus uses a mixture of subjects with topic proportions that are specific to it, but all documents in the corpus share the same K topics. Different from traditional document mixing models, which link every document to a single, unidentified topic, topic models isolate each document. Erosheva et al. (2004) describe them as mixed-membership models. Each document has a corresponding response as covariates which are jointly modelled for prediction when determining labels for unlabeled new documents. By allowing a response to be associated with each document and jointly modeling the response variable and the corpus of documents, the LDA model is extended to a supervised learning environment using Supervised Latent Dirichlet Allocation (sLDA). According to Blei et al. (2017), this enables it to identify the latent topics that are most predictive of the response variables in the training set and even to make predictions about future unlabeled documents. Following the LDA model's notation, let y represent a response variable from a generalized linear model with parameters η and δ . Should we consider the subsequent fixed; $\beta_{1:K}$: the k topics with each β_k a vector of term probability, η and δ and the Dirichlet hyperparameter for the per-document topic proportion θ . For every document and response variable, the generative process assumed by the sLDA is as follows:

1. Draw topic proportion, $\theta | \alpha \sim \text{Di}(\alpha)$;
2. for each word,
 - a) Draw a topic assignment $Z_n | \theta \sim \text{Mult}(\theta)$
 - b) Draw word $w_n | Z_n, \beta_{1:K} \sim \text{Mult}(\beta_{z_n})$;
3. Draw a response variable $y | z_{1:N}, \eta, \delta \sim \text{GLM}(\bar{z}, \eta, \delta)$, where $\bar{z} = \frac{1}{N} \sum_{n=1}^N z_n$.

For sLDA, the Dirichlet distribution (Di) is used as a prior for the topic distributions. Each document is assumed to have a distribution

over topics drawn from a Dirichlet distribution. This helps in capturing the variability in topic proportions across different documents. Each topic is represented by a multinomial distribution over the vocabulary, indicating the probability of each word appearing in that topic. This allows the model to not only discover Topics but also predict the response variable based on the document's topic distribution.

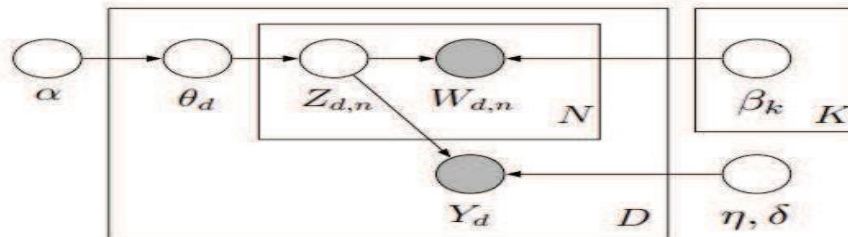


Figure 1. Graphical Model representation of Supervised Latent Dirichlet Allocation (sLDA)

2.2. Fitting the model

The method used for fitting the model is the variational expectation-maximization (VEM) algorithm. The sLDA model was equipped with parameters such as the “Document-Term matrix”, “K” (the number of topics), “vocab” (vocabulary words associated with word indices used in documents), “e.iteration”, “m.iteration”, “alpha”, “eta”, “var” (variance of the response variable), “annotation” (response variable) and various parameters was used to adjust the speed of convergence of the algorithm. Though there are no standards for the choice of parameters the parameters were used to set the convergence tolerance for the variance and E-M algorithms, respectively. Also, the maximum number of iterations for the conjugate gradient algorithm was set by one of the parameters which iterates between the E-step and M-step in a bid to maximize the likelihood of the corpus. The procedure finds the maximum bound about the latent variables (the topic proportions and the topic assignments Z) in the E-step, and the M-step, finds the maximum bound for the model parameters (the topics and the multivariate normal parameters).

Application of variational EM was done until the relative change in the likelihood bound was achieved. The iterations convergence ranged between 1hr to 6hrs 45minutes on a core i7 HP laptop of 2.7GHz with 8 GB RAM for each of the K values.

A data split was done on the 5000 plus preprocessed documents/articles in the proportion of 70% training data and 30% test data. The training data set was used to train our sLDA model while the test data set was used as a test for prediction of our response variable. The response variable is the Timeseries closing price of Bitcoin corresponding to the documents. Additionally, we categorized our response variable into a numeric classification of “low= 1”, “fairly

low= 2”, “fairly high= 3” and “high= 4” with a threshold of “less than or equal to 10000”, “less than or equal to 20000”, “less than or equal to 40000” and “less than or equal to 60000”, respectively.

Using the subsequent initializations: η to a grid on $[-1, 1]$ with increments of $1/K$, and β_k to randomly perturbed uniform topics. Different sLDA models with varying numbers of topics as 3, 10, 20, 30, 50, 75, 100 and 200 were trained so that the variational EM method was executed for the per-document ELBO at the E-steps as well as until a relative change in the corpus-like likelihood bound was less than 0.021% (Blei et al., 2017).

3. Result and Discussion

The experiment below was conducted to assess the quality of prediction using the sLDA model. We carried out a sensitivity analysis juxtaposing a latent Dirichlet allocation (LDA) with a regression model with our proposed Supervised latent Dirichlet allocation (sLDA) model. This experiment was carried out using two datasets of an equal number of documents. One dataset was from the CNBC market section (as stated in section 2), and the other was from the Bitcoinist website. Both datasets went through the same preprocessing stages (as seen in section 2). We will refer to the CNBC dataset as ‘Dataset(1)’ and the Bitcoinist dataset as ‘Dataset(2)’.

We compared the behaviors of sLDA and LDA + Regression models while trying to investigate the extrapolative strengths of the two models on the two different datasets. The LDA on its own cannot predict the Bitcoin price category. The usual practice when LDA is to be used for prediction has always been joint modeling of the LDA with a regression model (Clinton et al., 2004; Blei et al., 2017). To better understand the experiment, Figures 2 and 3 show some metric evaluations used to adjudge the better model in the topic range of 3 to 200. This explains how our proposed model (sLDA) for prediction purposes behaved against the conventional model (LDA+ regression) for prediction. In a quest for the optimal model between the LDA and the sLDA models, model evaluation metrics were employed to pick an optimal model that best fits the prediction task. The choice of our model evaluation was born out of the type of classification we used, i.e. the multinomial classification techniques. The evaluation metrics are Mean absolute Error (MAE), Mean absolute percentage Error (MAPE) Root mean squared Error (RSME), and the coefficient of Determination (R^2).

As seen in Figure 2, the sLDA model across the ‘K’ values, performed better with K=20 and K= 30 from Figure 3. Likewise, from Figure 2, the LDA + regression model across the ‘K’ values, performed better with K= 75 and K=100. Figure 3 shows that the sLDA model with K=20 has the lowest MAE, MAPE, and RSME values but has the highest R^2 value when compared with the LDA model. Also, the sLDA model with K=30 has the lowest MAE, MAPE, and RSME values and

the highest R^2 value when compared with the LDA model. The result in Figure 2 shows that LDA at $K=75$ has an R^2 value of 0.3246 while sLDA at $K= 20$ has an R^2 value of 0.6262. The sLDA model performance has an improvement of 90% when we consider the R^2 values for both models. From Figure 3, LDA at $k =100$ has an R^2 value of 0.4462 while sLDA at $K= 30$ has an R^2 value of 0.736. The sLDA model performance also had an improvement of 64%. According to Zang and Kjellström (2014), the improvement in classification results is always significant through different noise levels when the number of topics is small.

The subsequent subsections show the practical application of the above experiment using ‘Dataset(1)’.

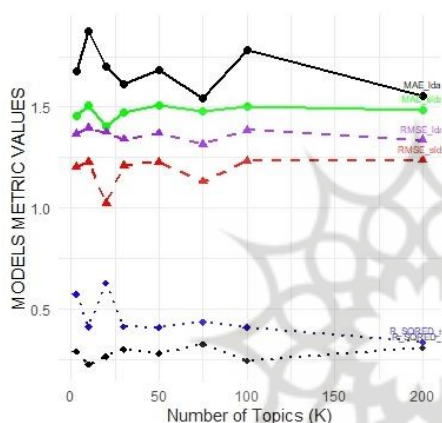


Figure 2. sLDA vs LDA behaviour with Dataset(1)

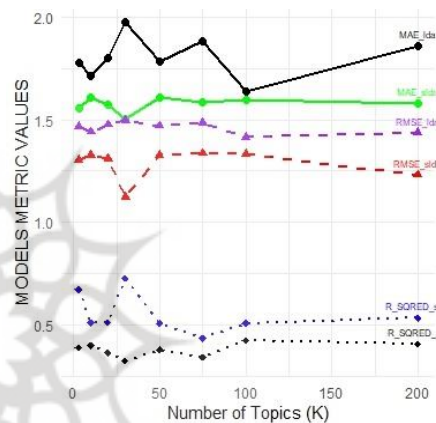


Figure 3. sLDA vs LDA behaviour with Dataset(2)

Table 1. Tabular summary of the split corpus

Number of documents	Total word count	Total number of unique words
4073(training data)	1,151,424	17,330
1746(testing data)	494,535	15,591

3.1. sLDA model evaluation

In a quest for the optimal sLDA model, model evaluation metrics were employed to pick an optimal model that best fits our prediction task on the first dataset. The evaluation metrics are Mean absolute error, mean absolute percentage error, and root mean squared error.

Table 2 shows the sLDA model with its numeric evaluation values was picked. A good look across the table shows that the sLDA model with $K= 3$ has the lowest numeric values across the three metrics but choosing that of $K= 3$ may lead to underfitting our model. We settled for the sLDA model with the next lowest numeric value, which is the $K= 20$ sLDA model, as the optimal model.

Exploring Figure 4, gives an insight into the topics that are predictive of the reported News from each of the documents. Our Topics ranges

from Minner’s investment, Crypto ransom, Economic policies, Scandals, exchange rates etc. We will briefly discuss the first four Topics. Topic 1 depicts simply, talked about the community of bitcoin miners in Texas who thinks that investing in the energy sector can aid in resolving the state's (Texas) power grid issues. This news was heavily reported by the media and appeared in some of the web scraped documents at different dates (Figure 4).

Table 2. Table showing numeric values of three sLDA model evaluation metric

K	MAE	RMSE	MAPE
3	1.05841	1.202434	0.6576
10	1.11	1.228	0.72566
20	1.070072	1.223507	0.639
30	1.076113	1.220401	0.6587
50	1.11068	1.227168	0.7296
75	1.073	1.231	0.645
100	1.11445	1.2332	0.7268
200	1.082776	1.235816	0.6508

Topic 1	mining	bitcoin	miner	power	energy	mine	state	electricity	cost	texas
Topic 2	attack	account	hacker	security	data	wallet	computer	criminal	informati	ransomwa
Topic 3	crypto	cryptocurr	token	musk	digital	dogecoin	ether	asset	market	tesla
Topic 4	china	country	chinese	governme	bank	south	korea	yuan	world	north
Topic 5	stock	percent	market	year	growth	price	investor	average	analyst	equity
Topic 6	blockchai	technolog	network	transactio	token	ripple	ledger	project	applicatio	coin
Topic 7	card	credit	account	income	gain	cash	inflation	youre	return	rate
Topic 8	president	trump	house	russia	bill	senate	russian	governme	state	federal
Topic 9	exchange	security	company	million	commissio	offering	agency	comment	customer	statement
Topic 10	dollar	rate	inflation	bank	market	index	week	high	yield	reserve
Topic 11	bank	currency	digital	financial	central	cryptocurr	payment	regulator	regulator	libra
Topic 12	investor	fund	asset	investmer	bitcoin	investing	market	invest	financial	advisor
Topic 13	people	money	student	home	life	woman	family	inflation	school	university
Topic 14	long	trader	kelly	brian	short	fast	call	buyer	spread	adami
Topic 15	bitcoin	currency	bitcoins	gold	value	digital	transactio	price	read	virtual
Topic 16	bitcoin	price	cryptocurr	percent	trading	exchange	digital	currency	market	according
Topic 17	company	coinbase	million	square	venture	startup	billion	customer	product	business
Topic 18	inflation	going	money	cramer	really	thats	thing	want	people	know
Topic 19	share	stock	premarke	revenue	cnbc	earnings	company	quarter	tesla	estimate
Topic 20	buffett	apple	read	berkshire	warren	google	amazon	food	book	billionaire

Figure 4. The twenty probable (20) Topics with Ten (10) words/terms from the sLDA model with K= 20

Topic 2 is all about (crypto security) such as ransomware attacks posing a severe risk, thus companies in the financial services, healthcare & pharmaceuticals, and employee personally identifiable information (PII) and HR sectors need to implement strong security measures to safeguard their data and systems. The news was also reported by the media and appeared in some of the web scraped documents at different dates, either by reference to it or just a few mentions.

Topic 3 shows deep thought about Bitcoin’s market competition

amidst other cryptocurrencies. That is, despite its humorous beginnings, Dogecoin has grown into an intriguing aspect of the cryptocurrency world, mostly because of Elon Musk's lighthearted engagement. Observing Musk's tweets as the market moves forward may be just as important as examining technical charts.

Topic 4 is about how some countries' economic circumstances affect the market trend of Bitcoin. All the topics had words that were either heavily mentioned, partially mentioned or not mentioned at all in the news encapsulated in our various published documents.

Figure 5 explains how important the topics are to the documents. Topic 3, shows a high level of importance to the documents while Topic 20, shows the lowest level of importance to the documents.

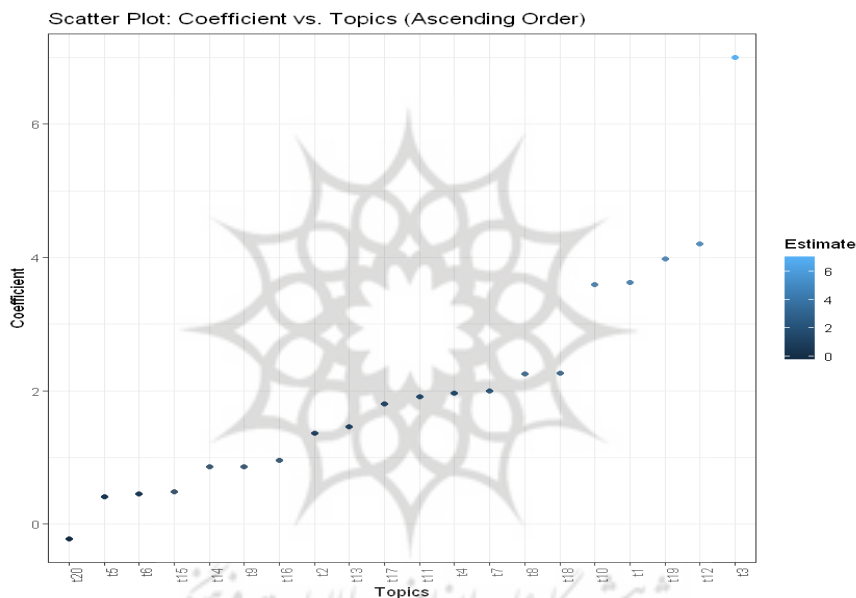


Figure 5. Graph of the coefficients of the sLDA model's probable Topics

Table 3 is a 4073 (row) by 20 (column) matrix which shows the proportions of terms/words from the documents/articles that make up the sLDA Topics. It is a posterior covariance matrix. However, because of the volume of this result and for clarity's sake, just the first eight Topics and the first Fourteen documents were selected and displayed in Table 4, while the full information can be seen in the appendix.

A critical look into the first four documents, Table 4 shows how the topic-proportion in the News known as a document and the response of the market price trend, 24 hours after publication. With reference to Table 4, Topic 1, 6, 7 and 8 from Document 0 had its impact on the market trend while Topic 2, 3, 4, 5 had their proportional impact on the market trend from Document 0. Topic 6 had no impact on the market trend from Document 1 while Topic 1, 2, 3, 4, 5, 7 and 8 had their proportional impact on the market trend from Document 1.

Table 3. The sLDA model per topic-document proportion output

	Topic 1	Topic 2	Topic 3	Topic 4	Topic 5	Topic 6	Topic 7	Topic 8
Doc 0	0	0.011905	0.02381	0.011905	0.154762	0	0	0
Doc 1	0.013699	0.034247	0.013699	0.136986	0.068493	0	0.006849	0.020548
Doc 2	0.008	0.054	0.002	0.02	0.022	0.286	0.018	0.008
Doc 3	0.003236	0.038835	0.074434	0.116505	0	0.223301	0.029126	0.016181
Doc 4	0	0.00463	0.125	0.023148	0.101852	0.041667	0	0.013889
Doc 5	0.048421	0.065263	0.008421	0.012632	0.176842	0.002105	0.027368	0
Doc 6	0.021277	0.035461	0.014184	0.070922	0.060284	0.049645	0.035461	0.053191
Doc 7	0	0	0.22	0	0.28	0.02	0	0
Doc 8	0.601027	0.005137	0.02226	0.119863	0.025685	0.003425	0.001712	0.035959
Doc 9	0.008086	0	0.040431	0.016173	0.040431	0.002695	0.016173	0.105121
Doc 10	0.012346	0.018519	0.074074	0.018519	0.098765	0.006173	0	0.030864
Doc 11	0	0.018127	0.003021	0.039275	0.009063	0.069486	0.012085	0.003021
Doc 12	0.002899	0.005797	0.023188	0.156522	0.034783	0.008696	0.017391	0
Doc 13	0	0.003876	0.003876	0	0.027132	0.003876	0.023256	0.027132
Doc 14	0.009479	0.009479	0	0.004739	0.530806	0.080569	0.018957	0.004739

Table 4. The price market trend for the few selected published documents (see Appendix)

Doc	Date published	open	High	Low	Close	Volume	% Change
Doc 0	28/8/2022	20032.73	20133.36	19606	19606	4.16E+10	-2.13%
Doc1	24/8/2022	21513.00	21826.33	21172.0	21378.1	5.06E+10	-0.63%
Doc 2	25/8/2022	21367.72	21773.15	21330.7	21591.28	5.04E+10	1.05%
Doc 3	18/8/2022	23338.393	23575.61	23140.4	23206.1	5.1E+10	-0.57%

All the Topics had their proportional effects on the market trend from Document 2. Also, all the topics had their proportional effect on the market trend from document 3 except for Topic 5. The effects were either positive or negative change in the market trend. While there are many factors that can cause price fluctuations in financial markets, the function of the media, is a consistent source of information and sentiment.

Figure 6 is a graph of an in-sample forecast which we did to predict response variables for a Document; i.e., a document could be read and inferences made as to how it is going to affect the market price. This could be made useful by a would-be miner or cryptocurrency analysts. The sLDA prediction model was tested using the testing data which is assumed to be unlabeled documents/articles. The predicted classification in line with the earlier mentioned classification in section (2.2), was juxtaposed with the testing data's initial response variable (classification). For the case of our voluminous data, 200 documents of the testing data and its response variable, alongside the corresponding predicted response, were captured. The result shows that the two trends appear in a near-similar pattern.

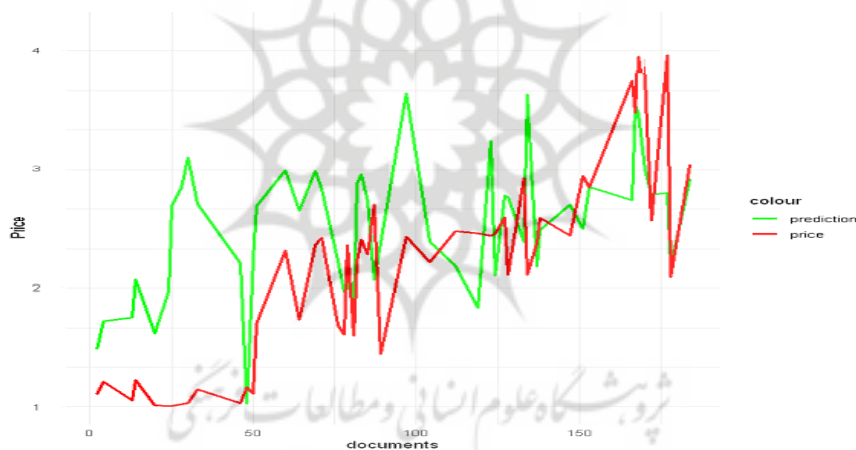


Figure 6. The initial price classification and the sLDA predicted classification

4. Conclusion

A helpful technique for examining huge text corpora is topic modeling. Topic modeling has applications outside of text mining, linguistics, and language modeling; it has been utilized in computer vision, population genetics analysis, social sciences, and humanities research.

Topic modeling, by revealing the topic structure buried in the text collection, can offer a higher-level exploration strategy to have a solid grasp of big text corpora when combined with other text-mining and machine-learning technologies. Additionally, the outputs of topic models can strengthen tasks like information retrieval, collaborative filtering, classification/categorization, and recommendation systems.

To understand how latent crowd knowledge affects the bitcoin market and to further leverage its predictive power, the research has looked at the significance and use of the sLDA-VEM and sLDA-predict for textual data mining and analysis. Hidden predictive topics used by traders in our textual data were exposed. Similarly, the links between Documents/ articles and topics were disclosed. Finally, we investigated the sLDA model's predictive power by extrapolating traders' interest or response variables for unlabeled documents.

Conflict of interest

The authors declared no conflicts of interest.

Ethical considerations

The authors have completely considered ethical issues, including informed consent, plagiarism, data fabrication, misconduct, and/or falsification, double publication and/or redundancy, submission, etc. This article was not authored by artificial intelligence.

Authors Contributions

TM, MOA and MUA: Developed the conceptual framework; TM-Source for the data and analyze the data; MOA and MUA-Supervised the work; TM-Came up with the first draft; TM, MOA and MUA- came up with the final draft; MOA and MUA- Do the editing of the workData availability.

Funding

This research did not receive any grant from funding agencies in the public, commercial, or non-profit sectors.

Acknowledgements

We wish to thank the effort of the reviewers and the Editor.

Appendix 1. The sLDA model Per Topic-Document proportion output (doc0-doc101)

Table with columns for Topic 1 through Topic 20 and rows for documents doc0 through doc101. Each cell contains a numerical value representing the proportion of a topic within a document.

Appendix 2. The sLDA model Per Topic-Document proportion output (doc3970-doc4072)

Table with columns for document ID (doc3970 to doc4072) and 30 numerical values representing proportions for different topics. Each row corresponds to a document and contains 31 columns of data.

Appendix 3. Preprocessed data randomly attached from Document 0 to Document 5818

	A	B	C	D	E	F	G	H	I	J	K
1	Unnamed	article_he article_se article_lin article_fir article_las article_su article_bo	Date	Open2	Close2						
2	doc 0	Bitcoin A Markets	https://w	2022-08-2	2022-08-2	Investors bitcoin bri	8/28/2022	20032.73	19606.74		
3	doc 1	A closely-Cryptocur	https://w	2022-08-2	2022-08-2	A measuri bitcoin pc	8/24/2022	21513	21378.1		
4	doc 2	Bitcoin ha Cryptocur	https://w	2022-08-2	2022-08-2	Bitcoin mi crypto wir	8/25/2022	21367.72	21591.28		
5	doc 3	Sudden er Crypto Wc	https://w	2022-08-1	2022-08-2	Bitcoin hit bitcoin fri	8/18/2022	23338.39	23206.19		
6	doc 4	Ether is u Technolo	https://w	2022-08-1	2022-08-1	Since find since find	8/18/2022	23338.39	23206.19		
7	doc 5	Skybridge Crypto	https://w	2022-08-1	2022-08-1	Bitcoin fui bitcoin fui	8/14/2022	24443.21	24323		
8	doc 6	Bitcoin toj Cryptocur	https://w	2022-08-1	2022-08-1	Bitcoin bri bitcoin bri	8/14/2022	24443.21	24323		
9	doc 7	BlackRock Markets	https://w	2022-08-1	2022-08-1	The large blackrock	8/10/2022	23145.39	23940.2		
10	doc 8	MicroStra Crypto Wc	https://w	2022-08-0	2022-08-0	MicroStra microstrat	8/1/2022	23311.56	23311.59		
11	doc 9	Bitcoin bri Cryptocur	https://w	2022-07-2	2022-07-2	Bitcoin's r bitcoin toj	7/28/2022	22955	23847.01		
12	doc 10	Tesla has Technolo	https://w	2022-07-2	2022-07-2	Tesla CEO early yea	7/19/2022	22517.1	23407.35		
13	doc 11	Bitcoin jui Markets	https://w	2022-07-2	2022-07-2	Bitcoin roi bitcoin roi	7/21/2022	21279.62	22943.12		
14	doc 12	Crypto mi Crypto Wc	https://w	2022-07-1	2022-07-1	New data data block	7/17/2022	21207.85	20769.4		
15	doc 13	Bitcoin toj Cryptocur	https://w	2022-07-1	2022-07-1	Bitcoin bo bitcoin bo	7/17/2022	21207.85	20769.4		
16	doc 14	Bitcoin toj Technolo	https://w	2022-07-2	2022-07-2	Bitcoin toj bitcoin toj	7/19/2022	22517.1	23407.35		
17	doc 15	Bitcoin Fa Crypto Wc	https://w	2022-07-0	2022-07-0	The 'Bitcc bitcoin fan	7/1/2022	19922.82	19336.1		
18	doc 16	World's -iETF Edge	https://w	2022-07-0	2022-07-0	Grayscale digital cur	7/6/2022	20209.98	20557.54		
19	doc 17	For bitcoi Technolo	https://w	2022-07-1	2022-07-1	Industry pi improven	7/13/2022	19330.03	20187.15		
20	doc 18	Coinbase Technolo	https://w	2022-07-1	2022-07-1	Coinbase i share coin	7/17/2022	21207.85	20769.4		
21	doc 19	El Salvado Crypto Wc	https://w	2022-06-2	2022-06-2	The gover salvador e	6/24/2022	21097.66	21279.53		
22	doc 20	Billions in Crypto Wc	https://w	2022-07-0	2022-07-1	The bad d crypto len	7/5/2022	20219.67	20185.96		
23	doc 21	Bitcoin he Markets	https://w	2022-07-0	2022-07-0	Bitcoin is i price bitcc	7/7/2022	20562.57	21647.17		
24	doc 22	Bitcoin jui Crypto Wc	https://w	2022-06-3	2022-06-3	Bitcoin jui bitcoin fin	6/29/2022	20278.9	20101.12		
25	doc 23	Five reasc Technolo	https://w	2022-07-0	2022-07-0	Bitcoin loi bitcoin wc	6/30/2022	20117.92	19816.75		
26	doc 24	Bitcoin po Cryptocur	https://w	2022-06-3	2022-06-3	Bitcoin ha bitcoin thi	6/29/2022	20278.9	20101.12		
27	doc 25	Bitcoin fal Cryptocur	https://w	2022-06-3	2022-06-3	Cypto hed bitcoin thi	6/29/2022	20278.9	20101.12		
28	doc 26	Charts sug Mad Moni	https://w	2022-06-2	2022-06-2	"Bitcoin ex cnbs cran	6/22/2022	20722.23	19976.05		
29	doc 27	ProShares Markets	https://w	2022-06-2	2022-06-2	Eight mon eight mor	6/19/2022	19010.48	20554.9		
30	doc 28	Bitcoin bri Crypto Wc	https://w	2022-06-2	2022-06-2	Bitcoin fel bitcoin bri	6/28/2022	20715.2	20293.45		
31	doc 29	Bitcoin bo Cryptocur	https://w	2022-06-2	2022-06-2	The cryptc bitcoin jui	6/19/2022	19010.48	20554.9		
32	doc 30	Bitcoin bri Crypto Wc	https://w	2022-06-1	2022-06-1	Bitcoin fel bitcoin pli	6/17/2022	20406.02	20468.21		
33	doc 31	Bitcoin's p Technolo	https://w	2022-06-1	2022-06-1	A \$4 billio lost billioi	6/13/2022	26599.66	22509.1		
34	doc 32	Bitcoin co Technolo	https://w	2022-06-2	2022-06-2	Ian Harnel cryptos pe	6/21/2022	20577.6	20739.94		
35	doc 33	Bitcoin ha Invest in Y	https://w	2022-06-1	2022-06-1	Bitcoin's r cryptocur	6/14/2022	22456.78	22209.52		
36	doc 34	New York Crypto Wc	https://w	2022-06-0	2022-06-0	Lawmakei following	6/2/2022	29795.17	30441.47		
37	doc 35	Bitcoin ca Technolo	https://w	2022-06-1	2022-06-1	Bitcoin's p selloff cry	6/14/2022	22456.78	22209.52		
38	doc 36	MicroStra Cryptocur	https://w	2022-06-1	2022-06-1	MicroStra aggressive	6/14/2022	22456.78	22209.52		
39	doc 37	Bitcoin rel Markets	https://w	2022-06-2	2022-06-2	Bitcoin bo bitcoin cli	6/20/2022	20580.04	20672.8		
40	doc 38	Cramer d Crypto	https://w	2022-06-1	2022-06-1	CNBC's Jin cnbs cran	6/12/2022	28415.54	26704.73		
41	doc 39	Bitcoin dr Crypto Wc	https://w	2022-06-1	2022-06-1	Factors in bitcoin tui	6/12/2022	28415.54	26704.73		
42	doc 40	Human rig Crypto Wc	https://w	2022-06-0	2022-06-0	Human rig washington	6/7/2022	31380.14	31193.07		
43	doc 41	Bitcoin dr Cryptocur	https://w	2022-06-1	2022-06-1	Bitcoin A t bitcoin bri	6/13/2022	26599.66	22509.1		
44	doc 42	Binance p Markets	https://w	2022-06-1	2022-06-1	Binance tr binance ri	6/12/2022	28415.54	26704.73		
45	doc 43	Bitcoin bri Crypto Wc	https://w	2022-06-0	2022-06-0	Bitcoin A f bitcoin sli	6/6/2022	29930.46	31367.42		
46	doc 44	Grayscale Technolo	https://w	2022-06-3	2022-06-3	The U.S. S grayscale	6/29/2022	20278.9	20101.12		
47	doc 45	Despite t Empowen	https://w	2022-06-1	2022-06-1	CNBC spol even cryp	6/12/2022	28415.54	26704.73		
48	doc 46	Bitcoin sir Technolo	https://w	2022-06-1	2022-06-1	Crypto inv bitcoin dri	6/17/2022	20406.02	20468.21		
49	doc 47	Bitcoin ris Crypto Wc	https://w	2022-06-0	2022-06-0	Bitcoin roi bitcoin mc	6/5/2022	29847.33	29909.21		
50	doc 48	Bitcoin fal Cryptocur	https://w	2022-06-1	2022-06-1	Bitcoin is i bitcoin fel	6/15/2022	22174.8	22490.59		
51	doc 49	MicroStra U.S. Marki	https://w	2022-05-1	2022-05-1	MicroStra microstrat	5/11/2022	31000.98	28984.54		
52	doc 50	Bitcoin wc Next Gen	https://w	2022-07-0	2022-07-0	Bitcoin wc bitcoin pli	7/7/2022	20562.57	21647.17		
53	doc 51	Bitcoin pr Technolo	https://w	2022-05-1	2022-05-1	New rese bitcoin mi	5/17/2022	29833.05	30452.62		
54	doc 52	Bitcoin bil Technolo	https://w	2022-06-2	2022-06-2	Bankman- central ba	6/21/2022	20577.6	20739.94		
55	doc 53	Coinbase Markets	https://w	2022-06-2	2022-06-2	Coinbase i coinbase i	6/21/2022	20577.6	20739.94		
56	doc 54	A \$3.5 bill Crypto Wc	https://w	2022-05-0	2022-05-1	South Kor multibillc	5/8/2022	34062.89	33933.54		
57	doc 55	New York Crypto Wc	https://w	2022-05-0	2022-05-0	Lawmakei state york	5/4/2022	34062.89	33933.54		
58	doc 56	Bitcoin bo Markets	https://w	2022-05-1	2022-05-1	The price price bitcc	5/12/2022	29044.26	29130.04		
59	doc 57	Bitcoin co Davos WE	https://w	2022-05-2	2022-05-2	Guggenhe bitcoin dr	5/22/2022	29445.06	30366.04		
60	doc 58	Jack Dorse Bitcoin	https://w	2022-05-1	2022-05-1	At Block's block exei	5/17/2022	29833.05	30452.62		
61	doc 59	Luna Four Crypto Wc	https://w	2022-05-0	2022-05-0	The Luna i luna founi	5/4/2022	34062.89	33933.54		
62	doc 60	This real e Crypto De	https://w	2022-04-2	2022-04-2	Compass i taing busi	4/22/2022	34062.89	33975.55		
63	doc 61	Bitcoin in Crypto Wc	https://w	2022-05-1	2022-05-1	Investors investor b	5/9/2022	34062.89	30363.87		
64	doc 62	Warren Bl 2022 Berk	https://w	2022-04-3	2022-05-0	Despite a bitcoin st	4/29/2022	34062.89	38592		
65	doc 63	\$3 billion Crypto Wc	https://w	2022-05-1	2022-05-1	Luna Four investor e	5/15/2022	30076.93	31333.77		
66	doc 64	Grayscale Finance	https://w	2022-05-1	2022-05-1	A spot-ba grayscale	5/10/2022	30106.43	31024.7		
67	doc 65	Crypto exi Technolo	https://w	2022-06-2	2022-06-2	CoinFlex i digital ass	6/28/2022	20715.2	20293.45		
68	doc 66	El Salvado Crypto Wc	https://w	2022-05-0	2022-06-2	El Salvado salvador a	5/8/2022	34062.89	33933.54		
69	doc 67	40% of bit Crypto Wc	https://w	2022-05-0	2022-05-1	Bitcoin is i bitcoin ne	5/8/2022	34062.89	33933.54		
70	doc 68	Bitcoin dr Cryptocur	https://w	2022-05-0	2022-05-0	Bitcoin cri bitcoin bri	5/5/2022	34062.89	36568.76		
71	doc 69	Fort World Crypto Wc	https://w	2022-04-2	2022-04-2	Three Bitr for worth	4/25/2022	34062.89	40475.28		
72	doc 70	Bitcoin dri Crypto Wc	https://w	2022-05-0	2022-05-0	Cryptocur bitcoin dri	5/8/2022	34062.89	33933.54		
73	doc 71	Michael S Crypto Wc	https://w	2022-04-2	2022-04-2	MicroStra miami mii	4/20/2022	34062.89	41363.5		
74	doc 72	Bitcoin dri Crypto Wc	https://w	2022-05-1	2022-05-1	The price bitcoin dri	5/10/2022	30106.43	31024.7		
75	doc 73	ItA -s eETF Edge	https://w	2022-04-1	2022-04-1	Crypto exj crypto exj	4/16/2022	34062.89	40386.51		
76	doc 74	What to ci invest in Y	https://w	2022-04-2	2022-04-2	Before ad investor s	4/25/2022	34062.89	40475.28		
77	doc 75	Bitcoin jui Crypto Wc	https://w	2022-05-0	2022-05-0	The cryptc price bitcc	5/3/2022	34062.89	37732.3		
78	doc 76	Bitcoin bri Crypto Wc	https://w	2022-05-1	2022-05-1	Bitcoin fel bitcoin dri	5/9/2022	34062.89	30363.87		
79	doc 77	Bitcoin dri Cryptocur	https://w	2022-05-0	2022-05-0	The cryptc bitcoin co	5/7/2022	34062.89	35552.56		
80	doc 78	Bitcoin Pli Davos WE	https://w	2022-05-2	2022-05-2	The indus davos swi	5/23/2022	30295.93	29114.08		

5792	doc 5790	5 signs you Money M	https://www.2018-01-21/2018-01-21/Here's how horrible h	1/25/2018	34062.89	11219.44
5793	doc 5791	If you invest The Begin	https://www.2018-01-1/2018-01-1/Google's search enj	1/10/2018	34062.89	15065
5794	doc 5792	Travis Kalani Money M	https://www.2018-01-01/2018-01-01/Wealth m travis kala	1/7/2018	34062.89	16515.29
5795	doc 5793	Meet the Careers N	https://www.2018-01-1/2018-01-1/Millennial werent da	1/9/2018	34062.89	14710.6
5796	doc 5794	The 10 most Careers N	https://www.2018-01-2/2018-01-2/Education though co	1/22/2018	34062.89	10917.31
5797	doc 5795	Legendary Entrepren	https://www.2018-01-2/2018-01-2/Bill Miller legendary	1/21/2018	34062.89	11481.03
5798	doc 5796	Men most Money M	https://www.2018-01-1/2018-01-1/A new GO financial r	1/15/2018	34062.89	13769.22
5799	doc 5797	Make it's Money M	https://www.2017-12-3/2017-12-3/These weird bitcoin cri	12/30/2017	34062.89	12732.36
5800	doc 5798	Self-made Money M	https://www.2018-01-01/2018-01-01/If you have want supe	1/4/2018	34062.89	15163.9
5801	doc 5799	GOP tax p Money M	https://www.2017-12-2/2017-12-2/More than middle cli	12/19/2017	34062.89	17700.78
5802	doc 5800	3 common Careers N	https://www.2017-12-2/2017-12-2/Your essay college ap	12/20/2017	34062.89	16599.69
5803	doc 5801	Mark Cubi Money M	https://www.2017-12-2/2017-12-2/Shark Tan year oppo	12/25/2017	34062.89	13965.21
5804	doc 5802	The 8 best Entrepren	https://www.2017-12-1/2017-12-1/Here are t heart shar	12/18/2017	34062.89	18931.2
5805	doc 5803	Vanguard Money M	https://www.2017-12-1/2017-12-1/I don't lik jack bogie	12/12/2017	34062.89	16798.03
5806	doc 5804	Here's how Careers N	https://www.2017-12-2/2017-12-2/The bill cc morning s	12/19/2017	34062.89	17700.78
5807	doc 5805	Here's how Money M	https://www.2017-11-1/2017-11-1/With Andri billionair	11/14/2017	34062.89	6647.346
5808	doc 5806	Departme Careers N	https://www.2017-10-2/2017-10-2/Criminals safe cyber	10/23/2017	34062.89	5946.927
5809	doc 5807	Mark Cubi Money M	https://www.2017-10-2/2017-12-2/Even billk come sma	10/20/2017	34062.89	6000.943
5810	doc 5808	Success hi Leadershi	https://www.2017-11-01/2017-11-01/Not every charles bu	11/1/2017	34062.89	6751.17
5811	doc 5809	Billionaire Entrepren	https://www.2017-10-2/2017-10-2/A "I woi peter thie	10/25/2017	34062.89	5725.302
5812	doc 5810	Mark Cubi Money M	https://www.2017-10-2/2017-10-2/It's "a far t prefer pa	10/24/2017	34062.89	5526.929
5813	doc 5811	The 9 sure Leadershi	https://www.2017-09-01/2017-09-01/Take any resume te	9/4/2017	34062.89	4224.22
5814	doc 5812	Looks like Money M	https://www.2017-06-1/2017-11-1/We live i rising hes	8/10/2017	34062.89	3391.23
5815	doc 5813	For the bu Careers N	https://www.2017-06-2/2017-06-2/Traveling corporate	6/20/2017	34062.89	2718.49
5816	doc 5814	This start- South by S	https://www.2017-03-1/2017-03-1/A former tmarjuana	3/14/2017	34062.89	1240.06
5817	doc 5815	Want to tr Money M	https://www.2017-03-2/2017-03-2/An easy, a brick farm	3/19/2017	34062.89	1038.05
5818	doc 5816	'Shark Tan Careers N	https://www.2017-02-2/2017-02-2/The cyber tempting	2/27/2017	34062.89	1179.03
5819	doc 5817	The Five F Money M	https://www.2016-10-2/2016-10-2/If you and bought so	10/26/2016	34062.89	678.486
5820	doc 5818	Turn your How I Mac	https://www.2016-05-07/2016-07-01/The Collis john collis	5/8/2016	34062.89	458.6

References

- Blei, D.M.; Kucukelbir, A. & McAuliffe, J.D. (2017). "Variational inference: A review for statisticians". *Journal of the American Statistical Association*. 112(518): 859-877. <https://doi.org/10.1080/01621459.2017.1285773>.
- Buchholz, M.; Delaney, J.; Warren, J. & Parker, J. (2012). *Bits and Bets, Information, Price Volatility, and Demand for Bitcoin*. Economics. 312. <https://www.reed.edu/economics/parker/s12/312/finalproj/Bitcoin.pdf>.
- Clinton, J.; Jackman, S. & Rivers, D. (2004). "The statistical analysis of roll call data". *American Political Science Review*. 98(2): 355-370. <https://doi.org/10.1017/S0003055404001194>.
- Dawson, J. & Kendzioriski, C. (2012). "Survival-supervised latent Dirichlet allocation models for genomic analysis of time-to-event outcomes". *arXiv*. TR 225. <https://doi.org/10.48550/arXiv.1202.5999>.
- Eklund, M. & Bejerholm, U. (2004). "Time use and occupational performance among persons with schizophrenia". *Occupational Therapy in Mental Health*. 20: 27-47. https://doi.org/10.1300/J004v20n01_02.
- Erosheva, E.; Fienberg, S. & Lafferty, J. (2004). "Mixed-membership models of scientific publications". *Proceedings of the National Academy of Science*. 101(1): 5220-5227. <https://doi.org/10.1073/pnas.0307760101>.
- Fataliyev, K.; Chivukula, A.; Prasad, M. & Liu, W. (2021). "Text-based stock market analysis: A review". 1(1), July. <https://arxiv.org/pdf/2106.12985>.
- Frank, X.; Cambria, E. & Welsch, R.E. (2017). "Natural language based financial forecasting: A survey". *Artificial Intelligence Review*. 50(3): 49-73. <https://link.springer.com/article/10.1007/s10462-017-9588-9>.
- Kaya, M.Y. & Karşlıgil M.E. (2010). "Stock price prediction using financial news articles". *2nd IEEE International Conference on Information and Financial Engineering*. IEEE: 478-482.
- Kumar, G.; Jain, S. & Singh, U.P. (2020). "Stock market forecasting using computational intelligence: A survey". *Archives of Computational Methods in Engineering*. 28(6): 1-33. <http://dx.doi.org/10.1007/s11831-020-09413-5>.
- Li, Q.; Chen, Y.; Wang, Y.; Chen, Y. & Chen, H. (2018). "Web media and stock markets: A survey and future directions from a big data perspective". *IEEE Transactions on Knowledge and Data Engineering*. 30: 381-399. <http://dx.doi.org/10.1109/TKDE.2017.2763144>.
- Loughran, T.; McDonald, B. & Pragidis, I. (2019). "Assimilation of oil news into prices". *International Review of Financial Analysis*. 63. <https://doi.org/10.1016/j.irfa.2019.03.008>.

- Mckinney, W. (2010). "Data structures for statistical computing in Python". *Pyton in Science Conference*. <http://dx.doi.org/10.25080/Majora-92bf1922-00a>.
- Mohan, S.; Mullapudi, S.; Sammeta, S.; Vijayvergia, P. & Anastasiu, D. (2019). "Stock price prediction using news sentiment analysis". *IEEE Fifth International Conference on Big Data Computing Service and Applications (BigDataService)*. 205-208. <http://dx.doi.org/10.1109/BigDataService.2019.00035>.
- Pang, B. & Lee, L. (2008). "Opinion Mining and Sentiment Analysis". *Found. Trends Inf. Retr.* 2(1-2): 1-135. <http://dx.doi.org/10.1561/1500000011>.
- Perotte, A.; Bartlett, N.; Elhadad, N. & Wood, F. (2011). "Hierarchically supervised latent dirichlet allocation". *Neural Information Processing Systems*. 24. https://www.researchgate.net/publication/228449895_Hierarchically_Supervised_Latent_Dirichlet_Allocation.
- Sahut, J.; Hájek, P.; Olej, V. & Hikkerova, L. (2024). "The role of news-based sentiment in forecasting crude oil price during the Covid-19 pandemic". *Annals of Operations Research*. 345(2): 861-884. <http://dx.doi.org/10.1007/s10479-024-05821-z>.
- Schofield, A.; Magnusson, M.; Thompson, L. & Mimno, D. (2017). "Understanding text pre-processing for latent dirichlet allocation". *EMNLP*. <https://www.cs.cornell.edu/~xanda/winlp2017.pdf>.
- Shah, D.; Isah, H. & Zulkernine, F. (2019). "Stock market analysis: A review and taxonomy of prediction techniques". *International Journal of Financial Studies*. 7(2): 26. <https://doi.org/10.3390/ijfs7020026>.
- Sharma, R.K. (2020). "Comparison of stock price prediction models using news articles, currency exchange rates and global indicator performance". *Journal of Advanced Research in Dynamical and Control Systems*. 12(7). <http://doi.org/10.5373/JARDCS/V12SP7/20202273>.
- Thakkar, A. & Chaudhari, K. (2021). "Fusion in stock market prediction: A decade survey on the necessity, recent developments, and potential future directions". *Information Fusion*. 65: 95-107. <https://doi.org/10.1016/j.inffus.2020.08.019>.
- Wilcox, K.T.; Jacobucci, R.; Zhang, Z. & Ammerman, B.A. (2021). "Supervised latent dirichlet allocation with covariates: A bayesian structural and measurement model of text and covariates". *Psychological Methods*. 28: <http://dx.doi.org/10.31234/osf.io/62tc3>.
- Yao, W.; Xu, K. & Li, Q. (2019). *Exploring the influence of news articles on Bitcoin Price with Machine Learning*. Graduate School at Shenzhen, Tsinghua University, Shenzhen, China. <http://dx.doi.org/10.1109/ISCC47284.2019.8969596>.
- Yap, A.Y.; Schumaker, R. & Chen, H. (2012). "Predicting Stock price movement from financial news articles". *Information Systems for Global Financial Markets*. <http://dx.doi.org/10.4018/978-1-61350-162-7.ch006>.
- Zang, C. & Kjellström, H. (2012). "How to supervise topic models". *European Conference on Computer Vision*. pp. 500-515. Springer.