# MODELLING CRYPTOCURRENCY ACTIVITIES' REPORTAGE USING CORRELATED TOPIC MODEL

**T. O Maku [1][*]; M. O. Adenomon, and M. U. Adehi**

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

2, 3. Department of Statistics, Faculty of Science, Nasarawa State University, Kefi, Nigeria

**\*Corresponding Author:** makuto@fuotuoke.edu.ng; +2347060988998

## Abstract

This study aimed at unveiling the connections amidst cryptocurrency activities by users and how traders read these activities online to aid their trading capabilities. Therefore, it examined the daily actions of Bitcoin traders who may have learned about the market from news items. The study employed the latent Dirichlet Allocation (LDA) model and its version, and the Correlated Topic Model (CTM) to examine news stories that were scraped from the market section website of CNBC. Using this machine learning process, topics with their terminology and words that were hidden in the documents and articles were found. The Document-Term-matrix was used to create coherence and perplexity graphs, which were then used to determine the number of Topics (K) to prevent the CTM from being overfitted or underfitted. The topics were determined by estimating the proportions for each document, and the CTM was used to perform correlation tests between the themes that were found. Similarly, document-word proportions and topic-word proportions were also measured. The posterior covariance matrix produced by the CTM was used to create a dense topic graph, which was then fed straight into a network analysis model to show positive, negative, and no relationships between the topics. Additionally, papers were searched using the "Hellinger distance" approach to determine the relationship between two or more documents/articles. The outcome demonstrates that there were more positive correlations between the topics that were found and between the different documents and topics that might have included the collective knowledge underlying the cryptocurrency traders' actions.

**Keywords:** Topic modeling, Latent Dirichlet Allocation, Correlated Topic Model, Consumer News and Business Channel's market section website

## 1. INTRODUCTION

As early as 2200 BC, payments were made using several kinds of money before the invention of cryptocurrency. But since its inception, money has undergone constant transformation, taking on new forms while still performing the same essential functions. When money initially started to circulate, it was usually in the form of commodities, or goods with intrinsic value that were exchanged between people, such as produce or seeds. Some academics claim that the earliest money was made of cowry shells (Sehra *et al*, 2018).

As civilization developed, various forms of money, such as gold and silver, were produced and eventually took the role of barter as a medium of exchange.

Our perception of money started to change with the introduction of the World Wide Web, which eventually led to the development of credit cards, smartphones that could access bank

accounts, and other financial advancements. Simultaneously, technological improvements were also impacting other facets of our existence.

Because of the cryptocurrencies' recent development and increase in market capitalization, more people are opting to invest in platforms like Litecoin, Ethereum, Ripple, and Bitcoin. The price and policy uncertainty in the cryptocurrency market have an impact on the daily return pattern of the market (Sehra *et al*, 2018). One should be aware of the considerable danger that the decentralized nature of cryptocurrencies brings before making an investment or dealing with them. The uncertainty indices for investors and many users estimate the degree of unpredictability for the future direction of the financial market, spanning from numerous themes, the market factors (internal & external), and the trade friction among all the digital currencies.

We have decided to focus our analysis on Bitcoin out of all the various cryptocurrencies. According to a website (coindesk.com), the current price of one bitcoin is more than 42,978.125 US dollars or 30,480,937.80 million naira.

The legalization of Bitcoin in countries like Australia and Japan has significantly increased demand for the cryptocurrency and positively impacted the market. More countries throughout the world are gradually accepting Bitcoin. As per "The Africa Report," an online publication, the use of cryptocurrencies in Africa surged by 120% between July 2000 and June 2021. The behavior of individuals who possess a significant quantity of Bitcoins, together with societal opinions, political beliefs, and emotional states, all influence the price of Bitcoin.

News articles are a good source of knowledge on the previously stated components because they often provide the public with first-hand information. Through several mobile applications, people who are interested in Bitcoin transactions can access the most recent news, price movements, and their causes at any time and from any location. They could utilize this knowledge to help them decide whether to buy or sell. In this study, we investigated the influence of daily news on the price of bitcoin. The investigation shows that crowd wisdom is a powerful predictor of key events affecting cryptocurrency that investors and users use and can connect to data from news articles.

By evaluating and assigning a numerical value to these resources (textual data), we attempt to make use of the result of the investigation carried out. The most useful type of information is probably the textual information seen in news items, which is the views of traders and investors. It plays several significant roles. Text (in natural language form) is the most natural way to store and encode human information. Text is the most expressive form of information since it can explain other media kinds like images and videos. For example, to identify images that

match a user's textual input keyword query, the Google image search engine often matches a set of text of photos. (Zhai and Massung, 2016).

A substantial portion of the public gets at least some information via websites, applications, or social networking sites, according to several surveys. The web channel is currently the second most important source of information behind television among those who prefer reading news articles over watching or listening to the news. (Kaya and Karsligil, 2010). Given the exponential growth of digital text data, there is a great demand for methods and tools for handling and leveraging big text data in this environment. Text data is largely created by humans for communication, and as such, it is usually rich in semantic content and contains significant knowledge, opinions, and preferences. These factors all contribute to the demand for text data. However, most machine learning algorithms cannot deal with raw texts directly. News articles in plain text are converted into numerical features that can be processed by machine learning systems. The field of natural language processing (NLP) has made some text analysis algorithms and techniques easily accessible. Among the numerous text representation methods are distributional semantics, N-grams, Bag of Words, One-Hot Encoding, Vector Semantics (tf-idf), etc. The recommended approach made use of the Latent Dirichlet Allocation (LDA) and its version correlated Topic Model (CTM) which are part of the Topic model family.

Topic models that learn a limited range of topics, such as Latent Dirichlet Allocation (LDA), are said to be able to capture the co-occurrence of terms in discrete count data. (Blei and Jordan, 2003). Given themes, any article can be represented as a distribution over topics, and any downstream supervised job can then use the weights of this distribution as input. Topic models reduce the complexity of the data, which improves the interpretability of forecasts, particularly when the topics are interpretable. This renders topic models valuable in domains like criminal justice (Da *et al,* 2017) and health care (Michael *et al,* 2017) where interpretability is a need for downstream validation. Since it makes it simpler to evaluate this massive volume of textual data to learn about various information and preferences, including how much news affects the Bitcoin market, text/opinion mining has become quite popular (Zhai and Massung, 2016). Sources of textual data include both independently developed and corporately produced materials. Sources generated by the company, such as quarterly and annual reports, can provide a rich language structure that, when carefully examined, can forecast the company's future performance (Kloptchenko *et al,* 2004). Topic models are generative models with a probabilistic basis that are utilized in natural language processing and machine learning (Liu *et al,* 2016). "Topics" describes the undefined, hazy connections between words in a vocabulary and how they are used in writing. A document is conceptualized as an assortment of topics.

Topic models unearth the hidden themes in the collection and annotate the papers according to those themes. It is believed that each word comes from one of those topics. Lastly, a distribution of document coverage of topics is generated, providing a new way to analyze the topics' viewpoints data.

Blei and Jordan (2003) introduced the LDA, which is an even more extensive probabilistic generative model (unsupervised machine learning) that builds upon Probability latent semantic allocation (PLSA). Since the introduction of LDA (Blei and Jordan, 2003), the methodology has served as a foundation for several topic models. Although the LDA might be most appropriate for learning static or fixed set themes in a corpus, it can also be customized in some ways to learn dynamic topics. For example, to determine the links between topics in a corpus, Bollen *et al*, (2010) developed the Correlated Topic Model (CTM), which is based on LDA. The more adaptable Logistic Normal distribution was adopted by the CTM in place of the Dirichlet distribution for the latent variable that represents the percentage of each topic in a document. This allows for an examination of the covariance/correlation structure of the model's component elements. These provide a model where the existence of one latent topic may be related to another, and they more faithfully capture the latent topical structure of the data. The primary computational hurdle for Topic modeling using the LDA and its variants is approximating posterior distribution, and many of these problems are unsolvable. Numerous approximation strategies have been proposed, primarily categorized into two approaches: Variational Inference and Markov Chain Monte Carlo (MCMC). The MCMC approximation approach (Griffiths *et al,* 2004) uses collapse Gibbs sampling, while the variational inference approximation (Blei and Lafferty, 2006) uses mean field variational inference (Griffiths *et al,* 2004) collapse variational inference (Teh *et al,* 2006), and expectation propagation.

Stenqvist and L¨onn¨o, (2017), collected 2.27 million messages regarding Bitcoin from Twitter and utilized the dataset to infer the short-term price change. At intervals of four minutes to four hours, aggregated sentiment changes were performed, and values were advanced one to four times to correspond to corresponding price changes. The testing results show that their ideal parameters yield a 79% accuracy. Because their research is based on data from a single month, it is not representative.

400 distinct kinds of chain-lets were extracted from the Bitcoin blockchain by Akcora *et al*, (2018), who then used cosine similarity to cluster them. A random forest model was employed to predict the price of Bitcoin, and Granger Causality was utilized to show the dataset's predictive power. Certain forms of chain lets were shown to have the largest predictive

influence on both investment risk and the price of Bitcoin when they examined the Granger predictive causality of chain lets.

McNally *et al,* (2018) predicted the price of Bitcoin using mining difficulty, hash rate, and historical pricing data using a machine learning algorithm. They achieved the best classification accuracy of 52% and RMSE of 8% for price data, mining difficulty, and hash rate with the use of machine learning. Their RMSE of 8% and classification accuracy score of 52% were the best.

The research by Yao *et al,* (2019) is closely related to ours; in this study, the researchers investigated if news stories have an impact on the price of bitcoin by applying machine learning and the Senti-Graph. Sentiment analysis is used by Senti-Graph to turn a news article into a graph. In comparison to previous feature extraction strategies, their trial findings indicate that the strategy had a good forecast accuracy, which also illustrates the impact of news items on the price of bitcoin.

Bollen *et al,* (2010), attempted to predict the stock market. From daily tweets, Opinion Finder and GPOMS gathered seven features. Before training self-organizing fuzzy neural networks, these traits were investigated via Granger causality. They only use news articles in their scenarios, which are different from instant communications. The news contains more thorough information and is relatively longer. The length of the tweets utilized in their research was limited to 140 characters, whereas the news does not have this restriction.

Nagar and Hahsler, (2012), searched the news and eliminated articles of specific stock categories. Articles, sentences, paragraphs, and titles are a few examples. The number of variances between positive and negative words indicates the polarity of an event. The percentage of positive polarity incidences in a corpus determines its score. They found a strong relationship between the corpora scores and the firm price.

Mueller and Rauh, (2016), provided an innovative method for topic-modelling newspaper text summaries (LDA). They predicted the beginning of the bloody fight one or two years in advance. When making their prediction, they distinguished between the possibility of violence between nations and the occurrence of conflict within each nation at a given time. Their analysis shows that the within-country variation, or their generated news data, has a comparative advantage in wartime prediction. They proposed that this is a particularly useful addition since it incorporates textual data into the most often used conflict predictors in the study.

Geletta *et al,* (2020), used natural language processing (NLP) and machine learning (ML) techniques to identify recurring patterns in narrative research materials to distinguish between successful and failed studies. Their goal was to identify the factors that led to unsuccessful research endeavours.

Colianni *et al,* (2015), looked into whether supervised machine learning algorithms may be used to develop profitable bitcoin trading techniques utilizing cryptocurrency-related Twitter data. They provided an overview of many machine-learning techniques for identifying changes in the Bitcoin market. The well-known alternative currency that is examined in this study is called Bitcoin (BTC). To ensure that exact inputs were used across the entire model, a thorough error analysis was conducted. Their analysis produced an accuracy gain of 25% on average.

Blei and Lafferty, (2006), developed the correlated topic model (CTM), which displays the correlation between topic proportions using the logistic normal distribution. They created an approximation posterior inference method based on mean-field variational inference for this model. Their analysis revealed that CTM yields a better match than LDA when applied to a batch of OCRed articles from the journal Science. They also showed how the CTM may be used to intuitively view and analyze this and other unstructured data sets.

Utilizing an interpolated multi-model approach, Wolk, (2020), examined the impact of social media on cryptocurrency prices. They illustrated how the extremely speculative valuations of cryptocurrencies were substantially impacted by the mental and behavioral attitudes of the general public.

Parekh *et al,* (2022), proposed DL-Gues, a robust and hybrid framework that considered both market sentiment and the interdependence of each coin for predicting cryptocurrency prices. They considered the Dash price projection, which was calculated for various validation loss functions based on price history and tweets from Dash, Litecoin, and Bitcoin. To evaluate the applicability of DL-Gues on other cryptocurrencies, they also analyzed the price history and tweets of Bitcoin-Cash, Litecoin, and Bitcoin to deduce findings for price prediction of Bitcoin-Cash.

Card *et al,* (2017), provided a broad neural framework based on topic models and built upon recent advancements in variational inference techniques to enable the flexible and rapid insertion of metadata and the study of alternative models. Their approach yielded good results with a fair trade-off between ambiguity, coherence, and sparsity. Finally, to demonstrate the potential of their method, they looked at a corpus of papers about immigration to the US.

Valencia *et al,* (2019), proposed predicting fluctuations in the price of Litecoin, Ethereum, Bitcoin, and Ripple cryptocurrency marketplaces by utilizing publicly accessible social media data and conventional machine learning techniques. They compared the use of neural networks (NN), support vector machines (SVM), and random forests (RF) using components from Twitter and market data as input features. The results showed that NN outperformed the other models, that machine learning and sentiment analysis may be used to predict cryptocurrency prices, and that some cryptocurrencies may be extrapolated just from Twitter data.

Previous studies like Yao *et al,* (2019) and Wolk, (2020), have shown that the impact of investor opinion on Bitcoin returns can be measured using unstructured language, such as sentiments on Twitter, which may or may not include meaningful and sufficient information. Our study will use structured textual data to examine how the sentiment of Bitcoin users affects Bitcoin returns. The goal of our suggested text-feature extraction of unsupervised Latent Dirichlet Allocation, or the CTM, is to raise traders' awareness when they rely on news items for trade expertise.

## 2  MATERIALS AND METHOD

The study population focused on news stories about cryptocurrency-related activity published in foreign media between 2016 and 2022. The Consumer News and Business Channel (CNBC) is the source of these news pieces, and Appendix IV has hyperlinks to each document. Because of the widespread use of Bitcoin among other cryptocurrencies, it is used in this study.

Every one of the more than 6,000 news pieces published between 2016 and 2022 was written in English. With a custom Python script called "beautiful Soap" created with the Jupyter Notebook, the text data was quickly scraped from the source mentioned above together with the accompanying meta-data. The query 'Daily Bitcoin reports was used. The pages were stored in a comma-separated (CSV) format with the following meta-data:

    i.    Article Headline

    ii.    Article section

    iii.    Article link

    iv.    Article date

    v.    Article summary

    vi.    Article body

    vii.    Opening Price

    viii.    Closing Price

## 2.1    Specifications and Estimation of the CTM

To overcome a limitation, the Correlated Topic Model (CTM) introduces a more flexible distribution that permits investigating the covariance structure across components (Blei and Lafferty, 2006). This was accomplished by replacing the Dirichlet distribution with the Logistic Normal distribution (Aitchison and Shen, 2019), which is conjugated to the topic assignments and so more computationally convenient. The CTM yields a more accurate and realistic fit of the latent topical structure by using the logistic normal distribution in its hierarchical model to describe the hidden composition of topics linked with each document (Aitchison, 2019, Blei and Lafferty, 2006).



**Figure 2.1**: **The CTM model**

CTM assumes the following generative process and produces *N-word* documents.

1. For each topic $k \in \{1,\ldots,K\}$

   a. Draw a distribution over words $\beta_k \sim Dir\,(\emptyset)$

2. For each document $d \in \{1,\ldots,D\}$,

   a. Draw  $\eta_d|(\mu,\Sigma) \sim N\,(\mu, \Sigma\,)$

   b. Obtain a topic distribution $\theta_d = (\theta_{d1},\ldots,\theta_{dK})$ as $\theta_{dk} = f(\eta_d)$

   c. Let $\eta_d$ be the word in document d, $\forall n \in \{1,\ldots,N_d\}$

      i. Draw topic assignment $z_{d,n}|\eta \sim Mult(f\,(\eta d))$

      ii. Draw word   $w_{d,n}|\,(z_{d,n}, \beta_{1:K}) \sim Mult(\,\beta_{z_{d,n}})$

where  $f(\eta_d)$ is a function that maps the real vector natural parameter, η to the simplex

$$\theta = f(\eta) = \frac{\exp(\eta_i)}{\Sigma_j \exp(\eta_i)}$$

The full joint distribution is as follows:

$$p(w, z, \beta, \eta|\emptyset, \mu, \Sigma) = p(\beta|\emptyset)p(\eta|\mu, \Sigma)p(z|\eta)\, p(w|\beta, z)$$

$$=\prod_k \frac{\Gamma(\Sigma_v \emptyset_v)}{\prod_v \Gamma(\emptyset_v)} \beta_{kv}^{\emptyset_w - 1} \cdot \prod_d \frac{1}{2\pi^{\frac{k}{2}}|\Sigma|^{\frac{1}{2}}} \exp\left\{-\frac{1}{2}(\eta_d - \mu)^T \Sigma^{-1}(\eta_d - \mu)\right\} \cdot \prod_d \prod_n \frac{\exp(\eta_{dz_{dj}})}{\Sigma_k \exp(\eta_{dk})} \beta_{z_{dn}w_{dn}} \quad (1)$$

Given a collection of topics and distribution over topic proportions $\{\beta_{1:K}, \mu, \Sigma\}$, for a document $w_d$, the posterior distribution of the latent variables conditioned on the words of a document is;

$$p(\eta, z | w, \beta_{1:K}, \mu, \Sigma) = \frac{p(\eta|\mu,\Sigma)\prod_{n=1}^{N}p(Z_n|\eta)p(W_n|Z_n,\beta_{1:K})}{\int p(\eta|\mu,\Sigma)\prod_{n=1}^{N}\Sigma_{Z_n}^{k}p(Z_n|\eta)p(W_n|Z_n,\beta_{1:K})d\eta} \tag{2}$$

We estimated the CTM parameters using variational expectation maximization and given a collection of documents by attempting to maximize the likelihood of the corpus as a function of the topic $\beta_{1:K}$ and the multivariate Gaussian $(\mu, \Sigma)$.

In the E-step, we will use variational inference to maximize the (coordinate ascent) objective function constrained concerning the variational parameters for each document.

$$L(\mu, \Sigma, \beta_{1:K} | w_{1:D}) \geq \sum_{d=1}^{D} \mathbb{E}_q[\ln p(\eta_d z_d w_d | \mu, \Sigma, \beta_{1:K})] + H(q_d) \tag{3}$$

We maximized the bound concerning the model parameters in the M-step. This included a maximum likelihood estimate of the themes and multivariate Gaussian using assumed adequate statistics, with the expectation being taken about the variational distributions derived in the E-step.

$$\hat{\beta}_i \; \alpha \; \sum_d \tau_{d,i} n_d \tag{4}$$

$$\hat{\mu} = \frac{1}{D}\sum_d \lambda_d \tag{5}$$

$$\hat{\Sigma} = \frac{1}{D}\sum Iv_d^2 + (\lambda_d - \hat{\mu})(\lambda_d - \hat{\mu})^\tau \tag{6}$$

where $n_d$ is the vector of word counts for document *d*.


## 2.2 Fitting the Model

The variational expectation-maximization (VEM) approach is employed to fit the models. Parameters like the "Document-Term matrix," "K" (the number of topics), and "control" (the correlated Topic Model (CTM)-VEM) were used to fit the correlated Topic Model. The algorithm's rate of convergence was modified by varying the settings of the "control" (CTM-VEM). The convergence tolerance for the variance and E-M algorithms was established using the parameters. Additionally, one of the settings for the conjugate gradient algorithm which alternates between the E-step and M-step to optimize the likelihood of the corpus sets the maximum number of iterations that may be performed.

In the M-step, the technique determines the maximum bound for the model parameters (the topics and the multivariate normal parameters), and in the E-step, it determines the maximum bound for the latent variables (the topic proportions and the topic assignments Z). Variational inference was applied until the relative change in the probability was less than $10^{-6}$, and variational EM was applied until the relative change in the likelihood bound was less than

$10^{-4}$. The iterations on a 2.7GHz core i7 HP laptop with 8 GB RAM reached convergence in precisely 6 hours and 45 minutes.

### 3   RESULTS AND DISCUSSION



**Figure 3.1: Bar plot of documents from various sections on the CNBC website.**

Figure 3.1 shows that most of the documents out of the 5819 documents, came from the technology section. This simply means that cryptocurrency activities are driven by technology.



**Figure 3.2: Bar plot of the top 40 words in the corpus**

The top 40 words according to Figure 3.2 above show words like "bitcoin", "year", "company", and "market", e.t.c has a high frequency of occurrence in the corpus.



**Figure 3.3: Wordcloud of the first 400 words in the corpus**

Regarding Figure 3.3, the word cloud is a broader view of word occurrence frequencies when compared with the barplot in Figure 3.2 above. The bigger the word, the higher the frequencies.

**Table 3.1: Tabular Summary of the corpus**

| Number of documents | Total word count | Total number of unique words |
|---|---|---|
| 5,819 | 1,119,023 | 17,616 |

Table 3.1 shows a summary of the entire corpus, given that we analyzed 5,819 articles with a total word count of 1,119,023 and unique terms totaling 17,616.



**Figure 3.2: Graph of coherence scores for the corpus**

We used the elbow approach to examine the graph in Figure 3.2 and determine which number of topics is the most correct. When taking the coherence score into account, the elbow with the highest frequency (k = 13) has the optimal number of topics. We chose the second elbow with the highest frequency, k = 35, to prevent the model from being under-fitted. This seems to be a good way to represent the corpus without going overboard.

**Table 3.2: Table showing the first Seven (7) Topics with Ten (10) words/terms from the CTM output**

| Topic 1 | Topic 2 | Topic 3 | Topic 4 | Topic 5 | Topic 6 | Topic 7 |
|---|---|---|---|---|---|---|
| tokenize | playing | first | developer | staid | polling | buildup |
| honig | netherlands | tripled | fond | mediumsized | interfax | affirmed |
| imply | feel | thankful | tape | firstout | piloted | superrich |
| maxim | marked | regret | doug | costing | crow | iaccino |
| Aaound | sloan | stifel | chip | oncepopular | accrued | nonaccredited |
| Cryptic | bashed | supervision | perjury | sample | distributor | premature |
| incorporates | faceless | export | takeover | wikipedia | becky | mood |
| safety | supplychain | celebrate | screenshot | overexuberance | bancorp | faint |
| supplydemand | permissioned | vaynerchuk | natixis | iaccino | private | mile |
| cybercrime | devastating | gear | chosen | silver | buzzy | threat |

From the above table, Topic 1 talks about the News articles' safety measures, when considering adding bitcoin and other cryptocurrencies to one's retirement savings. Topic 2 exposes the News article's concern about the losses made in the bitcoin market, the sharp tumbling of bitcoin price, and its recovery trend. Awareness of scammers' activities in the Bitcoin and cryptocurrency market, education on safety guilds against falling victim, and the proportion of Bitcoin investors adopting artificial intelligence while trading were all captured in Topic 3. Topic 4 captures how Investors fled risk assets after the Federal Reserve of the US affirmed its commitment to an aggressive tightening path. Also, a measure of the activity of bitcoin miners' possibilities of giving investors a clue as to where the digital currency is headed next is captured in this Topic. Topic 5 encapsulates how Bitcoin could be having a "watershed" moment, as the cryptocurrency continued to rally off of lows. Sports trading card activities, fetching record prices in the millions, and collectors lining up online to buy virtual basketball "moments" that use blockchain were also captured in Topic 5. Activities as to how the crypto market was hit by several factors, ranging from the collapse of stable-coin terra-USD to questions of solvency at crypto lender Celsius are explained by Topic 6.

The activities of Draper, an early backer of Bitcoin and its underlying blockchain technology, participating in a so-called "initial coin offering" of Tezos was revealed by Topic 7. The Abu

Dhabi Financial Services Regulatory Authority's declaration on bitcoin having the same status as "commodities", is also captured in this same topic. The dollar falling to two-week lows on choppy trading, led by losses against the yen and euro, was recorded to be at risk sentiment improved in the afternoon session amid stock market gains and as U.S. Treasury yields rose. However, this led bitcoin price surging overnight, which is captured in Topic 8.

There are other 22 topics uncovered by the CTM, which can be seen in Appendix (i).

**Table 3.4 shows the CTM Per Topic-word proportion output.**

### TERMS/WORDS

| | | according | affirmed | aggressive | algorithmic | alltime | analyst |
|---|---|---|---|---|---|---|---|
| | 1 | 5.34E-05 | 5.37E-05 | 9.16E-05 | 9.83E-05 | 3.78E-05 | 5.79E-05 |
| | 2 | 8.25E-05 | 1.87E-05 | 8.29E-05 | 3.18E-05 | 6.71E-05 | 9.18E-05 |
| TOPICS | 3 | 3.45E-05 | 3.68E-05 | 0.00011 | 9.49E-05 | 9.66E-05 | 4.39E-05 |
| | 4 | 2.24E-05 | 6.25E-05 | 7.63E-05 | 4.18E-06 | 6.52E-05 | 9.3E-05 |
| | 5 | 3.97E-05 | 2.97E-05 | 6.81E-05 | 3.14E-05 | 8.86E-05 | 0.000101 |
| | 6 | 9.11E-05 | 0.000102 | 1.15E-05 | 6.1E-05 | 6.82E-05 | 7.19E-05 |
| | 7 | 1.63E-05 | 0.000112 | 5.11E-05 | 9.27E-05 | 4.96E-05 | 5.38E-05 |
| | 8 | 9.5E-05 | 6.33E-05 | 6.01E-05 | 9.87E-05 | 7.13E-05 | 2.7E-05 |
| | 9 | 9.83E-05 | 2.99E-05 | 0.000105 | 2.02E-05 | 6.14E-05 | 7.46E-05 |
| | 10 | 8.85E-05 | 7.42E-05 | 5.7E-05 | 2.11E-05 | 8.64E-05 | 2.94E-05 |

Table 3.4 above, shows the proportions or percentages of terms/words in each of the 35 Topics unveiled by the CTM. Ten (10) Topics and Six (6) words/terms were randomly selected and displayed in the table above while the full information can be seen in the appendix. The result reveals that 0.00534% of Topic 1 is formed by the word/term "according", 0.00825% of Topic 2 is formed by the word/term "according", 0.00345% of Topic 3 is formed by the word/term "according" and Same theory applies to other Topics and words/terms.

**Table 3.5 shows the CTM Per Document-Topic proportion output.**

|  | Topic 1 | Topic 2 | Topic 3 | Topic 4 | Topic 5 | Topic 6 | Topic 7 |
|---|---|---|---|---|---|---|---|
| document 0 | 0.014426 | 0.011609 | 0.017581 | 0.591978 | 0.012079 | 0.012051 | 0.011941 |
| document 1 | 0.0066 | 0.007897 | 0.008036 | 0.313197 | 0.012198 | 0.229338 | 0.007285 |
| document 2 | 0.011883 | 0.010888 | 0.030879 | 0.015107 | 0.012577 | 0.220742 | 0.009185 |
| document 3 | 0.009424 | 0.013925 | 0.011353 | 0.361039 | 0.010979 | 0.013062 | 0.008897 |
| document 4 | 0.006524 | 0.007213 | 0.007856 | 0.009461 | 0.007212 | 0.007661 | 0.006423 |
| document 5 | 0.010468 | 0.019943 | 0.01357 | 0.017734 | 0.015375 | 0.011094 | 0.010197 |
| document 6 | 0.013044 | 0.017754 | 0.01755 | 0.022268 | 0.018448 | 0.016679 | 0.015714 |
| document 7 | 0.015493 | 0.021778 | 0.02131 | 0.042051 | 0.024368 | 0.022564 | 0.015066 |
| document 8 | 0.010995 | 0.011721 | 0.01078 | 0.301849 | 0.012202 | 0.009947 | 0.011227 |

According to Table 3.5 above, there is a 5819(row) by 35(column) matrix which shows that the proportions of terms/words in the documents/articles, came from the CTM Topics. However, because of the volume of this result and for clarity's sake, just the eight (8) documents/articles were selected and displayed in the table above why the full information can be seen in the appendix. The result reveals that 1.4% of the words in Document 0 are from Topic 1, 0.66% of the words in Document 1 are from Topic 1, 1.18% of the words in Document 2 are from Topic 1 and this goes on till the 35th topic and 5819th documents are exhausted.



**Figure 3.3: Topic Correlation graph**

Visualizing the clear correlative relationship between the identified Topics using the Network analysis. [10], which relied on the edges (relationships) and the centrality of the Nodes (Topics) to determine the correlative measure of the Network Analysis graph. It is evident from Figure 3.3 above that the connecting lines, or edges, vary in thickness, distance, and color. The stronger the correlation, the more ticker the edges are; the blue and red colors indicate positive and negative correlations, respectively, between the Topics. The graph's density allows it to be successful in capturing the intricate correlation between the Topics. The hierarchy of Topic relevance is indicated by the nodes' (Topics') centrality.

Topic 35 occupies a central position in the graph, with thicker edges indicating a few Topics that are highly and positively connected with it. This only indicates that Topic 35, Topic 1, Topic 12, and Topic 21 have strong centrality and are hence highly predictive of traders.



**Figure 3.4: Topics' strength centrality correlation**

Figure 3.4 provides further information on the graph in Figure 3.3, which illustrates the relationship between the subjects. The subjects are arranged in ascending order based on their significance and strength; correlations are shown by the black squares, whereas none are shown by the grey squares. Upon closely examining the subjects presented in Table 3.2, Subject 1 is related to subjects 21, 28, 8, 12, 34, 1, and 20. The topics 31, 32, 29, 10, 2, 33, 16, 22, 14, 17, 24, and 18 are linked to topic 2. There is a relationship between subjects 11,9,27,5,4,23,3,13,19,30,15,6,31 and topic 3. Topics 11, 9,27,5,25,4,23,3,13,19,30,15,6,31,32,29, and 10 are linked to subject 4. Correlates with subjects 11, 9, 27, 25, 4, 23, 3, 13, and 19 are topic 5. For the sake of concision, the degree of correlation between the correlated topics is displayed

in Figure 3.3's dense sparse graph. Topic 1 is one instance of this, as it has favorable relationships with subjects 21, 28, 8, 12, 34, 1, and 20. The sparse graphs, however, indicate that it has a strong correlation with subject 12 initially, followed by topic 21 and then topic 20.



**Figure 3.5: Similarity of a queried document and closest documents**

One randomly chosen document was queried using the Hellinger distance approach to determine which two documents were closest to the requested document. The graphs below illustrate how the approach calculates the distances between the two nearest documents and the document being queried, accordingly. The graphs in Figure 3.5 above show that, among the two nearest documents and the questioned document, Topics 4 and 6 have the largest proportions. Every topic follows a similar pattern of appearance.

## 4. CONCLUSION

The significance and application of LDA and its version, the CTM, in textual data mining and analysis have been investigated to uncover latent themes, inferred topics, and topics' associations among documents. Relationships between the Hidden Topics in our textual data were also disclosed. Similarly, the relationships between the papers, Topic, and Word/terms were also disclosed. linkages between things like "safe investment plans by traders," "retirement saving incentives," "security of investors brokerage," "cryptocurrency theft," "enthusiasm among investors and traders," "excitement and interest of traders," "crypto mining and the US health care system," etc. are examples of these linkages. Lastly, the relationship between the documents was revealed thanks in large part to the "Hellinger distance."

## REFERENCES

Akcora C., Dey A. K., Gel Y. R., and Kantarcioglu M., *"Forecasting bitcoin price with graph chainlets,"* in Pacific-Asia Conference on Knowledge Discovery and Data Mining. Springer, 2018, pp. 765–776.

Aitchison J. and Shen S. (2019*). Logistic normal distributions: Some properties and uses.* Biometrika 67 261–272.

Blei D. and Jordan M. *Modeling annotated data.* In Proceedings of the 26th annual International ACM SIGIR Conference on Research and Development in Information Retrieval, pages 127–134. ACM Press, 2003.

Blei D. and Lafferty J. (2006), *Dynamic topic models.* In Proceedings of the 23rd International Conference on Machine Learning, pages 113–120.

Blei D. and Lafferty J. (2006). *Correlated topic models.* In Weiss, Y., Schölkopf, B., and Platt, J., editors, Advances in Neural Information Processing Systems 18. MIT Press, Cambridge, MA.

Bollen J., Mao H., and Zeng X., *"Twitter mood predicts the stock market,"* CoRR, vol. abs/1010.3003, 2010. [Online]. Available: http://arxiv.org/abs/1010.3003

Card, D., Tan, C., & Smith, N. A. (2017). *Neural models for documents with metadata*. arXiv preprint arXiv:1705.09296

Colianni, S., Rosales, S., &Signorotti, M. (2015). *Algorithmic trading of cryptocurrency based on Twitter sentiment analysis.* CS229 Project,1(5), 1-4.

Da Kuang, P. Jeffrey Brantingham, and Andrea L. Bertozzi. *Crime topic modeling. Crime Science*, 6:1–20, 2017.

Epskamp S, Borsboom D and Fried  E (2017) *"Estimating psychological networks and their accuracy"* Behav Res (2018) 50:195–212 DOI 10.3758/s13428-017-0862-1

Geletta S., Follett L., and Laugerman M. *"Latent Dirichlet Allocation in Predicting Clinical Trial Terminations"Department of Public Health,* Des Moines University, 169 Ryan Hall, 3200 Grand Ave, Des Moines, IA, USA

Griffiths, Thomas &Steyvers, Mark. (2004). *Finding Scientific Topics.* Proceedings of the National Academy of Sciences of the United States of America. 101 Suppl 1. 5228-35. 10.1073/pnas.0307752101.

Kaya M. Y. and Karsligil M. E., *"Stock price prediction using financial news articles,"* in 2010 2nd IEEE International Conference on Information and Financial Engineering. IEEE, 2010, pp. 478–482.

Kloptchenko, A., T. Eklund, et al. (2004). *"Combining Data and Text Mining Techniques for Analysing Financial Reports."* Intelligent Systems in Accounting, Finance & Management 12(1): 29-41.

Liu, L., Tang, L., Dong, W. *et al. An overview of topic modeling and its current applications in bio-informatics. Springer Plus* **5**, 1608 (2016). https://doi.org/10.1186/s40064-016-3252-8

Nagar A. and Hahsler M., *"Using text and data mining techniques to extract stock market sentiment from live news streams,"* in International Conference on Computer Technology and Science (ICCTS 2012), IACSIT Press, Singapore, 2012.

McNally S., Roche J., and Caton S., *"Predicting the price of bitcoin using machine learning,"* in 2018 26th Euro micro International Conference on Parallel, Distributed and Network-based Processing (PDP). IEEE,2018, pp. 339–343.

Michael C Hughes, Leah Weiner, Gabriel Hope, Thomas H McCoy Jr, Roy H Perlis, Erik B Sudderth, and Finale Doshi-Velez. *Prediction-constrained training for semi-supervised mixture and topic models. arXiv preprint arXiv:1707.07341*, 2017b.

Mueller H., Rauh C., (2016), University of Cambridge, INET Institute) *"Reading Between the Lines: Prediction of Political Violence Using Newspaper Text"*

Parekh, R., Patel, N. P., Thakkar, N., Gupta, R., Tanwar, S., Sharma, G. & Sharma, R. (2022). DL-GuesS: *Deep learning and sentiment analysis-based cryptocurrency price prediction.* IEEE Access, 10, 35398-35409.

Sehra, Avtar & Cohen, Richard &Arulchandran, Vic. (2018). *On cryptocurrencies, digital assets and private money*. Journal of Payments Strategy and Systems. 12. 13-32.

Stenqvist E. and L¨onn¨o J., "Predicting bitcoin price fluctuation with Twitter sentiment analysis," 2017.

Teh, Yee & Jordan, Michael & Beal, Matthew & Blei, David. (2006*). Hierarchical Dirichlet Processes. Machine Learning.* 1-30. 10.1198/016214506000000302.

Wolk, K. (2020). *Advanced social media sentiment analysis for short-term cryptocurrency price prediction.* Expert Systems, 37(2), e12493

Valencia, F., GÃ³mez-Espinosa, A., &ValdÃs-Aguirre, B. (2019). *Price movement prediction of cryptocurrencies using sentiment analysis and machine learning.* Entropy, 21(6), 589.

Yao, Wenbing& Xu, Ke & Li, Qi. (2019). *Exploring the Influence of News Articles on Bitcoin Price with Machine Learning.* 10.1109/ISCC47284.2019.8969596.

Zhai C. and Massung S., (2016). *Text Data Management and Analysis,* ACM Books series, #12, 3

**APPENDICES**

**(I)**

**Table showing the thirty-five (35) Topics with Ten (10) words/terms from the CTM output**

| | | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|
| Topic 1 | tokenize | honig | imply | maxim | abound | cryptic | incorporat | safety | supplyder | cybercrim |
| Topic 2 | playing | netherlan | feel | marked | sloan | bashed | faceless | supplycha | permissio | devastatir |
| Topic 3 | first | tripled | thankful | regret | stifel | supervisic | export | celebrate | vaynerchu | gear |
| Topic 4 | develope | fond | tape | doug | chip | perjury | takeover | screensho | natixis | chosen |
| Topic 5 | staid | mediumsi | firstout | costing | oncepopu | sample | wikipedia | overexub | iaccino | silver |
| Topic 6 | polling | interfax | piloted | crow | accrued | distributo | becky | bancorp | private | buzzy |
| Topic 7 | buildup | affirmed | superrich | iaccino | nonaccrec | premature | mood | faint | mile | threat |
| Topic 8 | koomey | utah | ascendanc | cryptofun | rwyo | nvidias | yeah | digitize | salesforce | cochief |
| Topic 9 | course | baseball | tattoo | monetize | moscow | summons | farreachin | stave | assign | somewhat |
| Topic 10 | underesti | belshe | rieder | stopgap | divorced | suspicious | iphone | banker | citizen | dialing |
| Topic 11 | temptatio | instead | scramblin | realized | chess | global | bentsen | reining | urge | wasted |
| Topic 12 | bloodclot | seller | processed | softened | diversifica | system | spec | runoff | tallied | exited |
| Topic 13 | bread | bolder | soccer | switzerlar | highlyanti | centraliza | parity | chicago | mahaney | richer |
| Topic 14 | loginov | quicken | jefferies | everyone | misrepres | plummeti | emin | barclays | bitpanda | worked |
| Topic 15 | ballooned | ethereal | bitcoincha | soybean | oversight | arising | caixin | microsoft | monitor | knowyour |
| Topic 16 | disparate | observe | analytics | versa | hoped | materially | squarespa | rigid | scrutinizir | download |
| Topic 17 | kaleido | unlike | monaco | zscaler | gearing | regional | diamondb | ottersec | obie | netscape |
| Topic 18 | parcel | favoring | jackie | breyer | blowback | lang | nlnk | module | tendency | calling |
| Topic 19 | chappelle | opecled | benefittin | microlend | resulting | facebookt | yhoo | saudi | redesign | arose |
| Topic 20 | memorab | neighborh | commissic | morningst | elder | patino | encourage | niche | whistlebl | raced |
| Topic 21 | wayfair | moderate | fyre | johnson | counterof | designing | offensive | considere | quits | happier |
| Topic 22 | oneal | bead | document | gigabyte | shortsight | distressec | warmed | consists | talking | tradeoff |
| Topic 23 | victory | kenyan | spate | cyberint | expiring | pandemic | asus | ebrokers | tianjin | draft |
| Topic 24 | hileman | keurig | garland | reply | gamers | esma | nakamoto | pivoted | till | assisting |
| Topic 25 | likelihood | opioids | tailspin | fixing | trump | talented | misunders | singapore | line | physic |
| Topic 26 | explosion | easymone | inverse | attributio | disclosure | fired | castle | vcbacked | alleviated | everest |
| Topic 27 | tested | buzzy | reside | unsealed | standoff | reader | specialty | twoshot | assign | trailed |
| Topic 28 | amend | content | gamble | explainer | streeters | shelled | surfside | mitsubish | namesake | pedal |
| Topic 29 | ziop | sledgehar | virgin | religion | soccer | based | contact | amending | corn | uptake |
| Topic 30 | russia | panteras | witnessin | winkler | financebit | buyback | imperfect | harleydav | spun | tucked |
| Topic 31 | hopefully | patch | commerci | copeland | solves | forfeited | guggenhe | stratosphe | roubini | sarah |
| Topic 32 | vehicle | rainy | overheate | clearly | poland | contends | warns | zavery | transactio | piloted |
| Topic 33 | grmn | tragedy | assume | argues | nice | hayek | body | sort | cornell | specially |
| Topic 34 | clothes | twohour | anecdotal | radically | garlinghou | respective | ingenuity | housing | jeep | nbas |
| Topic 35 | button | supporter | slim | springing | subset | unnecessa | ritholtz | sopori | helsinki | mobius |

(II)

## Tables showing the CTM Per Topic-word proportion output (not all).

|  | according | affirmed | aggressive | algorithm | alltime | analyst | apiece | arrow | asset | average | bank | bankruptc | bitcoin | blink | briefly |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
|  |  |  |  |  |  | WORDS/TERMS |  |  |  |  |  |  |  |  |  |
| T1 | 5.34E-05 | 5.37E-05 | 9.16E-05 | 9.83E-05 | 3.78E-05 | 5.79E-05 | 1.25E-05 | 9.20E-05 | 3.85E-05 | 8.41E-05 | 3.45E-05 | 9.05E-05 | 6.67E-05 | 3.60E-06 | 7.48E-05 |
| T2 | 8.25E-05 | 1.87E-05 | 8.29E-05 | 3.18E-05 | 6.71E-05 | 9.18E-05 | 6.35E-05 | 1.16E-05 | 2.13E-05 | 8.28E-05 | 7.21E-05 | 5.07E-06 | 9.09E-05 | 7.18E-05 | 2.55E-05 |
| T3 | 3.45E-05 | 3.68E-05 | 0.00011 | 9.49E-05 | 9.66E-05 | 4.39E-05 | 6.81E-05 | 3.05E-05 | 7.87E-05 | 6.61E-05 | 3.27E-05 | 7.78E-05 | 6.65E-06 | 3.50E-05 | 6.28E-05 |
| T4 | 2.24E-05 | 6.25E-05 | 7.63E-05 | 4.18E-06 | 6.52E-05 | 9.30E-05 | 3.36E-05 | 7.80E-06 | 2.38E-05 | 5.99E-05 | 0.000111 | 2.79E-05 | 8.62E-05 | 8.42E-05 | 0.000109 |
| T5 | 3.97E-05 | 2.97E-05 | 6.81E-05 | 3.14E-05 | 8.86E-05 | 0.000101 | 2.32E-05 | 7.08E-05 | 7.20E-05 | 8.48E-05 | 4.28E-05 | 2.87E-05 | 0.000102 | 2.92E-05 | 8.53E-05 |
| T6 | 9.11E-05 | 0.000102 | 1.15E-05 | 6.10E-05 | 6.82E-05 | 7.19E-05 | 5.06E-06 | 4.33E-05 | 8.95E-05 | 8.00E-05 | 1.18E-05 | 6.06E-05 | 5.58E-05 | 0.000102 | 3.35E-05 |
| T7 | 1.63E-05 | 0.000112 | 5.11E-05 | 9.27E-05 | 4.96E-05 | 5.38E-05 | 5.26E-05 | 9.08E-05 | 6.44E-05 | 1.52E-05 | 2.73E-06 | 9.72E-05 | 4.92E-05 | 1.03E-05 | 1.88E-05 |
| T8 | 9.50E-05 | 6.33E-05 | 6.01E-05 | 9.87E-05 | 7.13E-05 | 2.70E-05 | 3.40E-05 | 4.02E-05 | 3.19E-05 | 3.48E-05 | 5.88E-05 | 8.84E-05 | 9.59E-05 | 4.70E-05 | 1.78E-06 |
| T9 | 9.83E-05 | 2.99E-05 | 0.000105 | 2.02E-05 | 6.14E-05 | 7.46E-05 | 8.27E-05 | 1.86E-05 | 7.59E-05 | 5.42E-05 | 3.46E-05 | 5.75E-05 | 4.15E-05 | 3.97E-05 | 9.13E-05 |
| T10 | 8.85E-05 | 7.42E-05 | 5.70E-05 | 2.11E-05 | 8.64E-05 | 2.94E-05 | 5.26E-05 | 5.97E-05 | 1.14E-05 | 1.51E-05 | 8.90E-05 | 9.40E-05 | 7.08E-05 | 4.06E-05 | 9.12E-05 |
| T11 | 4.51E-05 | 1.33E-05 | 7.75E-05 | 9.47E-05 | 2.67E-05 | 0.000106 | 1.53E-05 | 8.65E-05 | 4.19E-05 | 1.23E-05 | 9.29E-05 | 5.37E-05 | 0.000108 | 5.04E-05 | 6.15E-05 |
| T12 | 6.59E-05 | 7.74E-05 | 4.67E-05 | 9.49E-05 | 7.53E-05 | 1.28E-06 | 9.27E-05 | 2.43E-05 | 3.84E-05 | 4.29E-05 | 4.37E-05 | 1.04E-05 | 1.90E-05 | 2.91E-05 | 2.57E-05 |
| T13 | 6.12E-05 | 0.000105 | 4.97E-05 | 6.02E-05 | 4.63E-05 | 0.0001 | 3.24E-05 | 8.34E-06 | 9.54E-05 | 3.23E-05 | 0.000104 | 7.40E-05 | 0.000102 | 7.86E-05 | 4.20E-05 |
| T14 | 2.08E-05 | 6.29E-05 | 0.000105 | 8.05E-05 | 3.05E-05 | 2.17E-05 | 2.27E-05 | 5.36E-05 | 8.14E-05 | 3.21E-05 | 6.78E-05 | 9.33E-05 | 2.35E-05 | 7.72E-05 | 4.19E-05 |
| T15 | 9.13E-05 | 1.81E-05 | 6.62E-05 | 1.76E-05 | 2.08E-05 | 8.74E-05 | 7.70E-05 | 0.000107 | 2.86E-05 | 0.000104 | 2.73E-05 | 9.37E-05 | 8.35E-05 | 0.0001 | 2.75E-05 |
| T16 | 5.22E-06 | 3.49E-05 | 2.07E-05 | 2.99E-05 | 0.000103 | 2.77E-05 | 5.97E-05 | 5.57E-05 | 4.85E-05 | 0.000109 | 8.62E-05 | 0.00011 | 7.57E-05 | 1.66E-05 | 9.47E-06 |
| T17 | 5.00E-05 | 9.61E-05 | 1.05E-05 | 3.56E-05 | 7.49E-05 | 4.38E-05 | 0.000108 | 3.02E-05 | 1.61E-06 | 1.15E-05 | 2.18E-05 | 9.26E-05 | 8.87E-05 | 4.99E-05 | 7.10E-05 |
| T18 | 1.51E-05 | 6.86E-05 | 8.00E-06 | 2.06E-05 | 8.97E-05 | 4.03E-05 | 6.14E-05 | 5.04E-05 | 7.05E-05 | 3.41E-05 | 4.36E-05 | 1.58E-05 | 9.68E-05 | 3.03E-05 | 7.31E-05 |
| T19 | 4.74E-05 | 3.99E-05 | 9.67E-05 | 1.60E-06 | 6.69E-05 | 0.000108 | 1.69E-05 | 4.52E-06 | 5.60E-05 | 1.12E-05 | 4.95E-05 | 5.27E-05 | 8.32E-05 | 4.00E-05 | 4.25E-05 |
| T20 | 8.26E-05 | 2.61E-05 | 8.41E-05 | 6.42E-05 | 3.46E-05 | 2.63E-05 | 7.40E-06 | 3.07E-05 | 3.29E-05 | 5.69E-05 | 8.26E-05 | 5.53E-05 | 1.31E-05 | 3.15E-05 | 2.81E-06 |
| T21 | 5.93E-05 | 1.75E-05 | 3.18E-05 | 0.000103 | 7.00E-05 | 2.73E-05 | 1.01E-05 | 0.000111 | 9.04E-05 | 2.22E-05 | 0.000108 | 6.04E-05 | 2.52E-05 | 3.21E-05 | 4.49E-05 |
| T22 | 6.32E-05 | 9.13E-05 | 7.94E-05 | 4.22E-05 | 3.48E-05 | 0.000107 | 5.70E-05 | 1.60E-06 | 1.81E-05 | 0.000103 | 0.000111 | 3.90E-05 | 3.83E-05 | 7.64E-06 | 1.60E-05 |
| T23 | 9.12E-06 | 6.39E-05 | 0.00011 | 3.18E-05 | 3.23E-05 | 7.04E-05 | 5.03E-05 | 7.81E-05 | 8.16E-05 | 3.64E-05 | 7.24E-05 | 6.81E-05 | 4.11E-05 | 2.93E-05 | 3.21E-06 |
| T24 | 1.74E-05 | 3.92E-05 | 7.65E-06 | 6.28E-05 | 7.81E-05 | 9.72E-06 | 5.45E-05 | 6.51E-05 | 0.000106 | 9.68E-05 | 9.23E-05 | 7.64E-05 | 3.84E-05 | 1.49E-05 | 4.43E-05 |
| T25 | 9.95E-05 | 1.20E-05 | 1.92E-06 | 9.30E-05 | 6.37E-05 | 1.43E-05 | 8.30E-05 | 5.28E-05 | 4.49E-05 | 9.54E-05 | 6.18E-06 | 5.86E-05 | 7.58E-05 | 4.95E-05 | 3.03E-05 |
| T26 | 0.000112 | 9.68E-06 | 3.13E-05 | 2.65E-06 | 0.000109 | 2.19E-05 | 0.000107 | 2.94E-05 | 4.58E-05 | 2.93E-05 | 9.20E-05 | 9.22E-05 | 2.27E-05 | 5.61E-05 | 4.77E-05 |
| T27 | 9.82E-06 | 5.40E-05 | 3.39E-05 | 7.58E-05 | 1.75E-05 | 4.18E-05 | 8.03E-05 | 3.92E-05 | 5.56E-05 | 6.76E-05 | 8.32E-05 | 0.000102 | 9.73E-05 | 6.04E-05 | 5.89E-05 |
| T28 | 8.51E-05 | 4.57E-05 | 7.73E-05 | 2.97E-05 | 3.81E-05 | 9.60E-05 | 0.000107 | 8.82E-05 | 2.88E-05 | 2.11E-05 | 8.96E-05 | 1.99E-05 | 0.000112 | 5.31E-05 | 6.06E-05 |
| T29 | 8.54E-05 | 4.06E-05 | 0.000113 | 5.16E-05 | 7.70E-06 | 3.31E-05 | 5.67E-05 | 2.44E-05 | 2.79E-05 | 1.77E-05 | 9.11E-05 | 9.39E-05 | 8.73E-05 | 3.17E-05 | 9.03E-06 |
| T30 | 2.25E-05 | 8.98E-05 | 8.19E-05 | 7.82E-05 | 1.46E-05 | 6.93E-05 | 5.65E-05 | 9.92E-05 | 0.000101 | 7.65E-05 | 3.00E-05 | 7.22E-05 | 2.06E-05 | 4.93E-05 | 9.37E-05 |
| T31 | 4.17E-05 | 2.39E-05 | 3.41E-05 | 8.25E-05 | 8.40E-05 | 5.19E-05 | 7.22E-05 | 8.86E-05 | 4.80E-05 | 0.000111 | 6.10E-06 | 0.000101 | 0.000101 | 6.61E-06 | 5.86E-05 |
| T32 | 0.000103 | 2.89E-05 | 0.000109 | 2.10E-05 | 8.61E-05 | 0.000101 | 3.59E-05 | 7.87E-05 | 6.44E-05 | 1.69E-05 | 2.18E-05 | 0.000102 | 9.42E-05 | 5.17E-05 | 8.06E-05 |
| T33 | 1.36E-05 | 0.000106 | 3.90E-05 | 2.42E-05 | 5.59E-05 | 7.62E-05 | 1.55E-05 | 5.29E-05 | 0.000112 | 5.45E-05 | 5.63E-05 | 3.35E-05 | 4.70E-05 | 7.20E-05 | 5.64E-05 |
| T34 | 8.89E-05 | 3.13E-05 | 5.99E-05 | 4.01E-05 | 2.80E-05 | 2.45E-05 | 5.45E-06 | 6.84E-05 | 0.00011 | 9.80E-05 | 6.62E-05 | 9.51E-05 | 7.65E-05 | 7.18E-05 | 4.71E-05 |
| T35 | 7.05E-05 | 7.65E-05 | 7.93E-05 | 2.77E-05 | 6.13E-05 | 9.52E-05 | 5.46E-05 | 2.91E-05 | 4.85E-05 | 6.58E-05 | 8.94E-05 | 7.38E-05 | 6.89E-05 | 8.60E-05 | 1.34E-05 |

|  | bring | capital | cause | celsius | central | chain | chair | chairman | close | coin | coincided | collapse | commitm | continue | crypto |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
|  |  |  |  |  |  | WORDS/TERMS |  |  |  |  |  |  |  |  |  |
| T1 | 7.88E-05 | 7.09E-05 | 3.76E-05 | 6.86E-05 | 2.98E-05 | 6.22E-05 | 3.57E-05 | 7.97E-05 | 3.07E-05 | 4.16E-06 | 8.86E-05 | 9.18E-05 | 6.71E-05 | 2.11E-05 | 4.41E-05 |
| T2 | 2.65E-05 | 3.52E-05 | 4.50E-05 | 9.36E-05 | 4.61E-05 | 5.85E-05 | 5.75E-05 | 8.92E-05 | 5.24E-05 | 9.61E-05 | 8.86E-05 | 1.25E-05 | 2.79E-05 | 0.000109 | 7.26E-05 |
| T3 | 0.000103 | 3.05E-06 | 1.16E-05 | 5.45E-05 | 8.93E-05 | 0.000107 | 2.38E-05 | 9.69E-05 | 7.22E-05 | 9.38E-05 | 8.73E-05 | 9.10E-05 | 8.09E-05 | 4.88E-05 | 8.21E-05 |
| T4 | 6.64E-05 | 0.000111 | 8.19E-05 | 2.34E-05 | 8.71E-05 | 9.99E-05 | 7.76E-05 | 3.07E-05 | 3.06E-05 | 6.14E-05 | 9.45E-05 | 6.69E-05 | 5.14E-05 | 4.29E-05 | 0.000101 |
| T5 | 3.81E-05 | 6.31E-06 | 7.15E-05 | 5.99E-05 | 6.47E-05 | 7.67E-05 | 4.89E-05 | 2.50E-05 | 2.62E-05 | 4.40E-05 | 9.00E-06 | 4.72E-06 | 7.33E-05 | 6.77E-05 | 8.13E-05 |
| T6 | 7.92E-05 | 6.32E-05 | 3.19E-05 | 3.69E-06 | 0.000103 | 5.67E-05 | 5.55E-05 | 1.66E-05 | 4.76E-05 | 0.000108 | 4.26E-06 | 4.89E-05 | 6.18E-05 | 8.37E-05 | 0.000111 |
| T7 | 5.35E-05 | 0.000107 | 8.82E-06 | 4.69E-05 | 2.88E-05 | 4.52E-05 | 9.16E-05 | 6.45E-05 | 7.11E-05 | 0.000107 | 9.66E-06 | 0.00011 | 5.79E-06 | 0.0001 | 4.04E-05 |
| T8 | 4.96E-05 | 4.30E-05 | 8.60E-05 | 4.04E-05 | 9.92E-05 | 0.000111 | 5.70E-05 | 1.60E-05 | 3.47E-05 | 1.41E-05 | 5.17E-05 | 8.88E-05 | 5.82E-06 | 3.99E-05 | 4.98E-05 |
| T9 | 9.14E-05 | 5.50E-05 | 9.30E-05 | 0.00011 | 7.22E-05 | 4.16E-06 | 4.60E-05 | 8.04E-05 | 7.06E-05 | 8.91E-05 | 2.46E-05 | 5.71E-05 | 8.60E-05 | 6.50E-05 | 5.20E-05 |
| T10 | 4.79E-05 | 9.84E-06 | 6.49E-05 | 6.71E-05 | 0.000106 | 3.36E-05 | 6.66E-05 | 8.11E-05 | 9.56E-05 | 0.000108 | 1.09E-05 | 6.08E-05 | 2.56E-05 | 4.19E-05 | 2.56E-05 |
| T11 | 6.59E-05 | 4.59E-05 | 7.60E-05 | 8.60E-05 | 2.12E-05 | 4.29E-05 | 5.11E-05 | 1.40E-05 | 4.58E-05 | 0.000102 | 8.42E-05 | 7.89E-05 | 4.99E-05 | 2.21E-05 | 2.94E-05 |
| T12 | 0.000107 | 2.47E-05 | 5.90E-05 | 2.45E-05 | 3.30E-05 | 6.49E-05 | 0.000103 | 9.46E-05 | 4.70E-05 | 1.64E-06 | 9.50E-05 | 6.79E-05 | 3.65E-05 | 4.95E-05 | 4.30E-06 |
| T13 | 4.02E-05 | 1.10E-05 | 5.42E-05 | 7.40E-05 | 0.000104 | 0.000104 | 8.86E-05 | 7.21E-05 | 0.000104 | 4.23E-05 | 5.77E-05 | 1.42E-05 | 7.63E-05 | 2.67E-05 | 8.79E-05 |
| T14 | 9.47E-06 | 6.29E-05 | 6.88E-05 | 0.000109 | 1.32E-05 | 4.00E-05 | 5.05E-05 | 1.40E-05 | 0.000101 | 9.54E-05 | 8.35E-05 | 2.17E-05 | 5.44E-05 | 7.30E-05 | 4.51E-05 |
| T15 | 6.85E-05 | 9.99E-05 | 2.90E-05 | 1.78E-05 | 1.18E-05 | 9.57E-05 | 4.78E-05 | 3.53E-05 | 4.86E-05 | 1.11E-05 | 5.30E-05 | 0.00011 | 6.77E-05 | 0.000106 | 2.87E-05 |
| T16 | 7.51E-05 | 6.28E-05 | 7.22E-05 | 8.97E-05 | 7.60E-05 | 6.96E-05 | 5.20E-05 | 1.36E-05 | 1.64E-05 | 7.48E-05 | 7.44E-05 | 5.62E-05 | 1.04E-05 | 1.86E-06 | 3.67E-05 |
| T17 | 8.47E-05 | 0.000102 | 6.00E-05 | 7.00E-05 | 3.64E-05 | 8.93E-05 | 2.81E-06 | 5.15E-05 | 4.89E-05 | 3.46E-05 | 7.86E-05 | 1.23E-05 | 5.44E-05 | 0.0001 | 6.21E-05 |
| T18 | 0.000109 | 9.94E-05 | 0.000111 | 8.15E-05 | 1.68E-05 | 7.03E-05 | 6.96E-05 | 9.99E-06 | 7.45E-05 | 5.02E-05 | 9.45E-05 | 9.34E-05 | 5.75E-05 | 3.25E-05 | 1.42E-05 |
| T19 | 8.23E-05 | 2.64E-05 | 3.62E-05 | 2.48E-05 | 5.86E-05 | 6.28E-05 | 9.63E-05 | 1.70E-05 | 0.000106 | 9.76E-05 | 8.46E-05 | 2.58E-05 | 3.51E-05 | 3.70E-06 | 9.97E-05 |
| T20 | 1.17E-05 | 4.02E-05 | 8.11E-05 | 4.06E-05 | 0.000102 | 0.000112 | 9.09E-05 | 4.32E-05 | 2.51E-05 | 4.74E-05 | 6.14E-05 | 5.67E-05 | 2.41E-05 | 7.24E-05 | 2.48E-05 |
| T21 | 4.11E-05 | 0.000104 | 3.43E-06 | 8.89E-05 | 2.40E-06 | 4.33E-05 | 0.000105 | 3.82E-05 | 0.000108 | 5.62E-05 | 6.51E-05 | 7.32E-05 | 1.78E-05 | 5.52E-05 | 7.41E-05 |
| T22 | 4.48E-05 | 7.14E-05 | 7.98E-05 | 0.000101 | 4.34E-05 | 8.65E-05 | 3.54E-05 | 1.20E-05 | 2.03E-05 | 9.13E-06 | 6.64E-05 | 0.00011 | 4.50E-05 | 0.000111 | 3.61E-06 |
| T23 | 1.16E-05 | 5.76E-05 | 2.96E-05 | 0.000103 | 6.36E-06 | 4.65E-05 | 6.07E-05 | 3.00E-05 | 4.36E-05 | 3.93E-06 | 7.81E-06 | 9.89E-05 | 0.000108 | 7.60E-05 | 0.000102 |
| T24 | 7.34E-05 | 9.56E-05 | 2.80E-05 | 9.99E-05 | 1.69E-05 | 8.23E-05 | 9.79E-05 | 1.14E-05 | 3.36E-05 | 0.000105 | 0.000105 | 9.30E-05 | 6.57E-05 | 5.06E-05 | 3.39E-05 |
| T25 | 5.00E-05 | 2.95E-05 | 5.14E-05 | 2.46E-05 | 2.96E-05 | 9.61E-05 | 5.50E-05 | 3.11E-05 | 9.81E-05 | 8.25E-05 | 9.58E-05 | 4.52E-06 | 0.000103 | 1.61E-05 | 5.51E-05 |
| T26 | 1.43E-06 | 7.82E-05 | 1.50E-05 | 8.99E-05 | 6.86E-05 | 3.44E-05 | 8.30E-05 | 9.36E-05 | 3.03E-05 | 2.56E-05 | 4.94E-06 | 0.0001 | 0.000101 | 5.21E-05 | 5.64E-05 |
| T27 | 3.92E-05 | 4.75E-05 | 9.56E-06 | 9.79E-05 | 6.23E-05 | 0.000105 | 3.99E-05 | 0.000106 | 5.59E-05 | 5.68E-05 | 0.000109 | 0.000106 | 0.000102 | 0.000103 | 2.05E-05 |
| T28 | 7.39E-05 | 0.000106 | 2.56E-05 | 6.03E-05 | 5.03E-05 | 1.02E-05 | 3.79E-05 | 8.49E-05 | 7.12E-05 | 0.000112 | 8.84E-05 | 1.20E-05 | 5.37E-05 | 6.09E-05 | 2.08E-06 |
| T29 | 9.37E-05 | 1.46E-05 | 6.80E-05 | 4.81E-05 | 6.51E-05 | 5.24E-06 | 8.18E-05 | 8.44E-05 | 7.83E-05 | 0.000103 | 0.00011 | 4.28E-05 | 3.70E-06 | 8.73E-06 | 5.18E-05 |
| T30 | 1.11E-05 | 7.66E-05 | 8.34E-05 | 7.67E-05 | 7.79E-05 | 0.000107 | 2.74E-05 | 6.53E-05 | 6.91E-05 | 6.36E-05 | 8.66E-05 | 2.44E-05 | 3.49E-05 | 9.98E-05 | 8.06E-05 |
| T31 | 6.46E-05 | 5.32E-05 | 4.78E-05 | 4.14E-05 | 6.34E-05 | 6.50E-05 | 9.69E-05 | 1.11E-05 | 5.39E-06 | 9.99E-05 | 0.000108 | 9.27E-05 | 3.12E-05 | 5.17E-06 | 6.79E-05 |
| T32 | 8.76E-05 | 9.64E-05 | 1.49E-05 | 5.73E-05 | 2.97E-06 | 0.000102 | 9.07E-06 | 7.69E-05 | 8.97E-05 | 1.57E-05 | 7.55E-05 | 1.30E-05 | 0.000103 | 4.59E-05 | 6.50E-05 |
| T33 | 1.88E-05 | 2.53E-05 | 0.000108 | 8.10E-06 | 7.81E-05 | 8.98E-05 | 3.42E-05 | 8.09E-05 | 5.16E-06 | 6.33E-05 | 8.41E-05 | 8.71E-05 | 3.97E-05 | 0.000107 | 5.97E-05 |
| T34 | 2.47E-05 | 0.000106 | 1.00E-05 | 3.89E-05 | 0.000104 | 4.72E-05 | 6.66E-06 | 3.68E-05 | 6.00E-05 | 0.000109 | 8.51E-05 | 0.00011 | 9.68E-05 | 3.35E-05 | 5.82E-05 |
| T35 | 0.000104 | 8.43E-06 | 5.90E-05 | 4.85E-05 | 3.31E-06 | 9.33E-06 | 3.27E-05 | 9.50E-05 | 7.29E-05 | 4.25E-05 | 5.84E-05 | 6.51E-05 | 2.69E-05 | 6.96E-05 | 6.09E-05 |

**(III)**

**Table showing the CTM Per Document-Topic proportion output, For 34 documents out of 5819 documents.**



**(IV)**

**Table showing Preprocessed data randomly attached from Document 0 to Document 5818.**

(v)

## Table showing articles link.

| Unnamed | article_he | article_se | article_link | article_fir | article_las | article_su | article_bo | Date |
|---|---|---|---|---|---|---|---|---|
| 0 | BitcoinÃ L | Markets | https://www.cnbc.com/2022/08/29/bitcoin-drops-below-20000-to-lowest-level-since-mid-july-as-investors-dump-risk-assets.html | 2022-08-2 | 2022-08-2 | Investors | bitcoin bri | 8/28/2022 |
| 1 | A closely- | Cryptocur | https://www.cnbc.com/2022/08/26/bitcoin-btc-price-key-metric-flashes-bottom-for-the-crypto.html | 2022-08-2 | 2022-08-2 | A measur | bitcoin pc | 8/24/2022 |
| 2 | Bitcoin ha | Cryptocur | https://www.cnbc.com/2022/08/26/crypto-winter-is-coming-but-it-will-be-a-warm-winter-says-vc-firm.html | 2022-08-2 | 2022-08-2 | Bitcoin m | crypto wir | 8/25/2022 |
| 3 | Sudden cr | Crypto Wc | https://www.cnbc.com/2022/08/19/sudden-crypto-market-drop-sends-bitcoin-below-22000.html | 2022-08-1 | 2022-08-2 | Bitcoin hit | bitcoin fri | 8/18/2022 |
| 4 | Ether is up | Technolog | https://www.cnbc.com/2022/08/19/ether-eth-price-outpaces-bitcoin-btc-as-ethereum-merge-nears.html | 2022-08-1 | 2022-08-1 | Since find | since find | 8/18/2022 |
| 5 | Skybridge | Crypto | https://www.cnbc.com/2022/08/15/skybridges-scaramucci-on-two-things-that-will-prop-up-bitcoin-demand.html | 2022-08-1 | 2022-08-1 | Bitcoin fu | bitcoin fu | 8/14/2022 |
| 6 | Bitcoin to | Cryptocur | https://www.cnbc.com/2022/08/15/bitcoin-tops-25000-for-the-first-time-since-june-before-slipping.html | 2022-08-1 | 2022-08-1 | Bitcoin br | bitcoin bri | 8/14/2022 |
| 7 | BlackRock | Markets | https://www.cnbc.com/2022/08/11/blackrock-launches-a-private-trust-to-give-clients-exposure-to-spot-bitcoin.html | 2022-08-1 | 2022-08-1 | The larges | blackrock | 8/10/2022 |
| 8 | MicroStrat | Crypto Wc | https://www.cnbc.com/2022/08/02/microstrategy-ceo-saylor-moves-to-chairman-role-focusing-on-strategy-and-bitcoin.html | 2022-08-0 | 2022-08-0 | MicroStra | microstrat | 8/1/2022 |
| 9 | Bitcoin bri | Cryptocur | https://www.cnbc.com/2022/07/29/bitcoin-btc-price-rises-following-stocks-higher-in-a-post-fed-rally-.html | 2022-07-2 | 2022-07-2 | Bitcoin's r | bitcoin top | 7/28/2022 |
| 10 | Tesla has c | Technolog | https://www.cnbc.com/2022/07/20/tesla-converted-75percent-of-bitcoin-purchases-to-fiat-currency-in-q2-2022.html | 2022-07-2 | 2022-07-2 | Tesla CEO | early year | 7/19/2022 |
| 11 | Bitcoin jur | Markets | https://www.cnbc.com/2022/07/28/bitcoin-jumps-above-23000-after-federal-reserve-interest-rate-hike.html | 2022-07-2 | 2022-07-2 | Bitcoin ro | bitcoin ros | 7/27/2022 |
| 12 | Crypto mi | Crypto Wc | https://www.cnbc.com/2022/07/18/crypto-miners-moved-over-300-million-of-bitcoin-in-one-day.html | 2022-07-1 | 2022-07-1 | New data | data block | 7/17/2022 |
| 13 | Bitcoin to | Cryptocur | https://www.cnbc.com/2022/07/18/bitcoin-btc-tops-22000-ethereum-jumps-as-crypto-market-rallies.html | 2022-07-1 | 2022-07-1 | Bitcoin bo | bitcoin bo | 7/17/2022 |
| 14 | Bitcoin to | Technolog | https://www.cnbc.com/2022/07/20/crypto-prices-bitcoin-btc-climbs-above-23000.html | 2022-07-2 | 2022-07-2 | Bitcoin to | bitcoin br | 7/19/2022 |
| 15 | Bitcoin Fa | Crypto Wc | https://www.cnbc.com/2022/07/02/bitcoin-family-say-they-lost-1-million-in-value-this-year.html | 2022-07-0 | 2022-07-0 | The "Bitco | bitcoin far | 7/1/2022 |
| 16 | Worldâ„¢ | ETF Edge | https://www.cnbc.com/2022/07/07/worlds-largest-bitcoin-fund-sues-sec-over-crypto-etf-rejection.html | 2022-07-0 | 2022-07-0 | Grayscale | digital cur | 7/6/2022 |
| 17 | For bitcoir | Technolog | https://www.cnbc.com/2022/07/15/bitcoin-btc-price-bottom-heres-what-the-market-wants-to-see.html | 2022-07-1 | 2022-07-1 | Industry p | improvem | 7/13/2022 |
| 18 | Coinbase | Technolog | https://www.cnbc.com/2022/07/18/coinbase-stock-pops-17percent-as-cryptocurrencies-bitcoin-and-ether-rally.html | 2022-07-1 | 2022-07-1 | Coinbase | share coin | 7/17/2022 |
| 19 | El Salvado | Crypto Wc | https://www.cnbc.com/2022/06/25/el-salvador-bitcoin-experiment-not-saving-countrys-finances.html | 2022-06-2 | 2022-06-2 | The gover | salvador e | 6/24/2022 |
| 20 | Billions in | Crypto Wc | https://www.cnbc.com/2022/07/06/tax-writeoff-possible-for-bitcoin-lost-to-lending-sites-like-celsius.html | 2022-07-0 | 2022-07-1 | The bad d | crypto len | 7/5/2022 |
| 21 | Bitcoin he | Markets | https://www.cnbc.com/2022/07/08/bitcoin-heads-toward-best-week-since-october-as-crypto-collapse-stabilizes.html | 2022-07-0 | 2022-07-0 | Bitcoin is | price bitc | 7/7/2022 |
| 22 | Bitcoin jus | Crypto Wc | https://www.cnbc.com/2022/06/30/bitcoin-just-had-its-worst-month-in-more-than-a-year.html | 2022-06-3 | 2022-06-3 | Bitcoin jur | bitcoin fin | 6/29/2022 |
| 23 | Five reaso | Technolog | https://www.cnbc.com/2022/07/01/bitcoin-btc-posts-worst-quarter-in-more-than-a-decade-5-reasons-why.html | 2022-07-0 | 2022-07-0 | Bitcoin los | bitcoin wc | 6/30/2022 |
| 24 | Bitcoin po | Cryptocur | https://www.cnbc.com/2022/06/30/bitcoin-btc-on-track-for-its-worst-quarter-in-more-than-a-decade.html | 2022-06-3 | 2022-06-3 | Bitcoin ha | bitcoin th | 6/29/2022 |
| 25 | Bitcoin fal | Cryptocur | https://www.cnbc.com/2022/06/30/bitcoin-falls-below-19000-again-as-pressure-mounts-on-crypto-firms.html | 2022-06-3 | 2022-06-3 | Cypto hec | bitcoin the | 6/29/2022 |
| 26 | Charts sug | Mad Mone | https://www.cnbc.com/2022/06/22/charts-suggest-bitcoin-could-rally-over-the-next-few-months-but-likely-wont-reach-old-highs-jim-cramer-says.html | 2022-06-2 | 2022-06-2 | "Bitcoin c | cnbcs cran | 6/22/2022 |
| 27 | ProShares | Markets | https://www.cnbc.com/2022/06/20/proshares-is-launching-a-short-bitcoin-etf-this-week.html | 2022-06-2 | 2022-06-2 | Eight mon | eight mor | 6/19/2022 |
| 28 | Bitcoin bri | Crypto Wc | https://www.cnbc.com/2022/06/29/bitcoin-btc-briefly-drops-below-20000-as-pressure-mounts-on-crypto.html | 2022-06-2 | 2022-06-2 | Bitcoin fel | bitcoin bri | 6/28/2022 |
| 29 | Bitcoin bo | Cryptocur | https://www.cnbc.com/2022/06/20/bitcoin-btc-rebounds-but-struggles-to-hold-above-20000-.html | 2022-06-2 | 2022-06-2 | The crypto | bitcoin jur | 6/19/2022 |