# Evidence-enhanced disinformation detection for Crypto News

## Introduction

The world of cryptocurrencies is a fascinating and dynamic space, with a current market size of over 1.08 trillion dollars and predicted growth in the future. However, with the abundance of news and information available, it can be challenging to differentiate between reliable sources and disinformation. This is where my project comes in - I aim to develop a tool that can help users identify whether a piece of crypto news can be trusted or not by providing a relevancy score with evidence as source urls.

Disinformation is a major problem in today's world, and the rise of cryptocurrencies has only made it easier for individuals to spread false information and manipulate the market. Detecting disinformation is a difficult task, especially in the crypto space where there is a lot of noise and misinformation. Our project seeks to address this issue by leveraging machine-learning algorithms that can detect disinformation in crypto news articles by gathering evidence from multiple sources.

The outcome of this project can have significant impacts on the crypto community by providing users with the tools to make informed decisions and avoid being misled by false information. Additionally, the project can be extended and modified for different news categories to address disinformation more broadly.


## Problem Description

Detecting disinformation in the world of cryptocurrencies can be a daunting task, as the abundance of news and information available can be overwhelming. As an individual interested in the crypto market, I understand the importance of differentiating between reliable sources and disinformation to make informed decisions. With the current market size of over 1.08 trillion dollars and predicted growth in the future, it is becoming increasingly crucial to address this issue. To solve this problem, I am undertaking a project to develop a tool that can help me and other users identify whether a piece of crypto news can be trusted or not.

While several approaches, such as natural language processing (NLP) techniques, graph-based methods, social network analysis, and machine learning methods, have been proposed to tackle the problem of disinformation detection, these methods often rely on a single source of evidence and may not take into account the complex interplay between different sources. My proposed approach seeks to leverage evidence from multiple sources, including social media, news articles, and user comments, to detect disinformation more effectively. For the project purpose I have kept my scope small and have used an online available dataset, but one can easily create a script to link multiple sources with an ETL pipeline to extract key required information such as article text, source, and URL.

The impact of disinformation in the crypto community can be significant, as false information can easily manipulate the market and mislead users. Therefore, my project's outcomes can provide users with the tools to make informed decisions and avoid being misled by false information. Furthermore, the project can be modified and extended to address disinformation in different news categories and industries.

To achieve my goals, I plan to use machine-learning algorithms that can detect disinformation in crypto news articles by gathering evidence from multiple sources. The system will analyze the text of a document to extract relevant features such as sentiment, topics, and entities. It will also use network structures to model the relationships between documents, users, and other entities to identify patterns of disinformation propagation.

The goal of disinformation detection is not only to detect disinformation but also to understand why certain content is classified as disinformation. Therefore, the system will provide clear and understandable explanations of generated results with evidence urls or sources and scores. This will help build trust and confidence in the system among users.

## What has been done so far

This is the so far progress on my ML project:
- The final report in latex format has been started.
- Github repo is updated with partial final report and code to date
- Final Dataset and ipynb file prepared for ETL and cleanup
- NLP-based algorithm 1 was completed with analysis on test data.
- SVM-based algorithm 2 completed in-progress of testing on test data.

1. I have started using ACM SIG Latex Template to prepare my final report

Fig 1: Overleaf Final Report

2. Uploaded my progress to date on the GitHub repo:
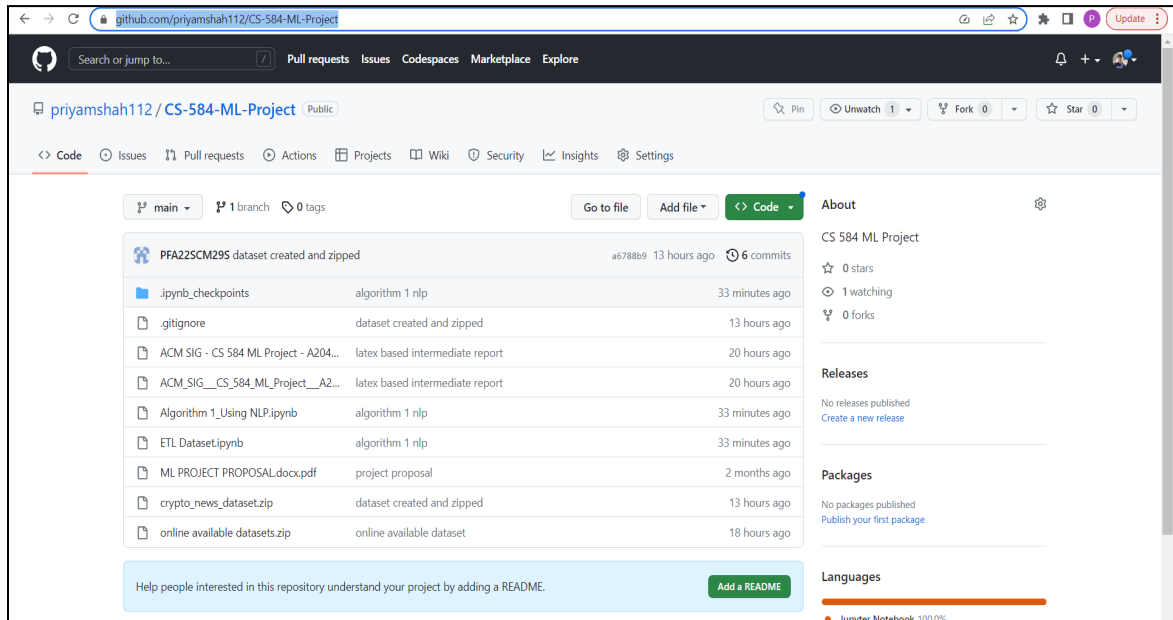   https://github.com/priyamshah112/CS-584-ML-Project



Fig 2: GitHub Snapshot of till date progress

3. For my analysis, I have found a few online datasets over Kaggle and figshare for
   cryptocurrency related news which have actual text, source, and URL available. I have
   created my own crypto_news_dataset.csv using these 3 available datasets. I have created
   my ETL Dataset.ipynb file which performs extraction of required column data and loads them
   in my csv dataset. Also, I have cleaned a few columns and rows which had garbage values.

   Below is the detailed dataset overview:



Fig 3: Online available datasets

The final dataset size and information is

Fig 4: crypto_news_dataset.csv [54572 rows x 5 columns]

4.  After reading the below-referenced research paper, I have planned to use NLP and SVM as 2 algorithms to start with my experiment.
5.  I have 2 ipynb files which are tentatively ready, I am still working on them to improvise and to capture results on more test data
6.  **The first one is Algorithm1_Using NLP.ipynb**



Fig 5: Algorithm1_using NLP.ipynb

This algorithm is an implementation of a content-based recommendation system using NLP for crypto news articles. It uses a dataset of crypto news articles in the form of a CSV file with columns for year, title, text, source, and URL. It takes a user text as an input and provides relevant [ text, source, url] with relevancy score to justify how authentic or relevant that news article can be. It also provides analysis parameters to assess the summarized results.

```python
query = "Walmart and Litecoin Payment News Debunked by Walmart Spokesperson, LTC Prices Shudder from Fake News"
similar_articles, analysis_params = get_similar_articles(query, df)
# print list of similar articles
for article in similar_articles:
    text, source, url, score = article
    print(f"Text: {text}\nSource: {source}\nURL: {url}\nRelevancy Score: {score}\n")

print(f"Number of similar articles: {analysis_params['num_similar_articles']}")
print(f"Mean sentiment score: {analysis_params['mean_sentiment_score']}")
print(f"Standard deviation of sentiment score: {analysis_params['std_sentiment_score']}")
print(f"Mean polarity score: {analysis_params['mean_polarity_score']}")
print(f"Standard deviation of polarity score: {analysis_params['std_polarity_score']}")
```

```
would allow the vehicl to access those restrict area addit that distribut network could serv as a way to track and authent th
e good themselv as the applic explain when a custom interact with a product the custom is permit to do so via a privat or pub
lic authent key in respon new block may be ad to subsequ root block which will contain inform relat to the date and time a pr
oduct deliv by the agv wa access as well as the authent key that access the product walmart wrote in the applic that it propo
system could increa custom loyalti due to the signif conveni it would provid likewi it ad that the system would like cut cost
becau of it autom natur at the expen of deliveri driver who might otherwi serv in that capac howev ind the applic further cem
ent the impress that wal mart is look at blockchain as a possibl tool for improv the autom of it servic as coindesk report in
march the retail which is separ research blockchain as a mean to track food shipment file a patent applic aim at creat a smar
t courier system walmart trailer imag via shutterstockth leader in blockchain news coindesk is a media outlet that strive for
the highest journalist standard and abid by a strict set of editori polici
Source: news
URL: https://www.coindesk.com/walmart-envisions-blockchain-powered-delivery-fleets/
Relevancy Score: 0.2562907813280027

Number of similar articles: 10
Mean sentiment score: 0.06501299226163085
Standard deviation of sentiment score: 0.060099302742572205
Mean polarity score: 0.40925166040176963
Standard deviation of polarity score: 0.09215631159297351
```

Fig 6: Output of NLP-based algorithm on test data=" Walmart and Litecoin Payment News Debunked by Walmart Spokesperson, LTC Prices Shudder from Fake News"

**Output Result Understanding**

The output analysis parameters suggest the following:

The number of similar articles: 10: This indicates that the algorithm has retrieved 10 most similar articles from the dataset for the given query text.

Mean sentiment score: 0.06501299226163085: This is the average sentiment score of the retrieved articles. The sentiment score ranges from -1 (most negative) to 1 (most positive), with 0 being neutral. The mean sentiment score of 0.065 indicates that the retrieved articles are slightly positive.

Standard deviation of sentiment score: 0.060099302742572205: This value indicates the variability in the sentiment scores of the retrieved articles. A lower standard deviation indicates that the sentiment scores are clustered closely around the mean, while a higher standard deviation indicates greater variability in the scores. The relatively low standard

deviation of 0.060 suggests that the sentiment scores of the retrieved articles are relatively consistent.

Mean polarity score: 0.40925166040176963: This is the average polarity score of the retrieved articles. The polarity score ranges from 0 (most objective) to 1 (most subjective). The mean polarity score of 0.409 indicates that the retrieved articles are slightly subjective.

Standard deviation of polarity score: 0.09215631159297351: This value indicates the variability in the polarity scores of the retrieved articles. A lower standard deviation indicates that the polarity scores are clustered closely around the mean, while a higher standard deviation indicates greater variability in the scores. The relatively high standard deviation of 0.092 suggests that the polarity scores of the retrieved articles are relatively diverse.

**Code Explanation**
The function get_similar_articles takes a query string, a dataframe df containing the news articles, and an optional parameter n for the number of similar articles to return (default 10). The function first preprocesses the query string and the text of the news articles using a set of text preprocessing techniques such as removing non-alphabetic characters, converting to lowercase, tokenizing, and stemming.

Then it computes the TF-IDF (term frequency-inverse document frequency) matrix, which is a numerical statistic that reflects how important a word is to a document in a collection. It then computes the cosine similarity between the query and all the documents in the dataset using the TF-IDF matrix. Cosine similarity is a measure of similarity between two non-zero vectors of an inner product space.

It then returns a list of n tuples, where each tuple contains the text of the article, its source, its URL, and its relevancy score, sorted by the relevancy score in descending order. The function also calculates sentiment and polarity scores for each article using the TextBlob library, which is a Python library for processing textual data. The sentiment score measures the overall positive or negative sentiment of the article, while the polarity score measures the degree of subjectivity of the article.

Finally, the function returns a dictionary with various analysis parameters, such as the number of similar articles, the mean sentiment score, the standard deviation of the sentiment score, the mean polarity score, and the standard deviation of the polarity score.

In the main code, a sample query is provided and passed to the get_similar_articles function, which returns a list of similar articles and analysis parameters. The list of similar articles is printed to the console along with their source, URL, and relevancy score. The analysis parameters are also printed to the console.

Overall, this algorithm implements a content-based recommendation system for crypto news articles that uses a combination of text preprocessing techniques, TF-IDF matrix, and cosine similarity to recommend similar articles based on a query string. It also provides sentiment and polarity scores for each recommended article.

7. **The second one is Algorithm 2_SVM.ipynb**

   Since in this dataset we assume that the text column has authentic data from mentioned sources and urls we can classify them as positive/true labeled data. Using a Support Vector Machine algorithm can help in identifying if the user-provided text can be trusted, not trusted or uncertain to be said this time based on the knowledge base of our captured dataset.

**SVM to provide evidence based dis-information detection in crypto news space**

```python
In [*]: import pandas as pd
        import re
        from sklearn.feature_extraction.text import TfidfVectorizer
        from sklearn.metrics.pairwise import cosine_similarity
        from sklearn.model_selection import train_test_split
        from sklearn.metrics import accuracy_score, precision_score, recall_score, f1_score
        from nltk.stem import PorterStemmer
        from nltk.sentiment.vader import SentimentIntensityAnalyzer

        def preprocess_text(text):
            text = re.sub('[^a-zA-Z]', ' ', str(text))
            text = text.lower()
            text = text.split()
            ps = PorterStemmer()
            text = [ps.stem(word) for word in text]
            text = ' '.join(text)
            return text

        def get_polarity_score(text):
            sid = SentimentIntensityAnalyzer()
            polarity_score = sid.polarity_scores(text)
            return polarity_score['compound']

        def get_similar_articles(user_text, n_similar=10):
            df = pd.read_csv('crypto_news_dataset.csv')
            df.dropna(inplace=True)
            df['processed_text'] = df['text'].apply(preprocess_text)
            tfidf = TfidfVectorizer(max_features=10000)
            tfidf.fit(df['processed_text'])
            user_text_processed = preprocess_text(user_text)
            user_tfidf = tfidf.transform([user_text_processed])
            similarity_scores = cosine_similarity(user_tfidf, tfidf.transform(df['processed_text']))
            df['similarity'] = similarity_scores[0]
            df.sort_values(by='similarity', ascending=False, inplace=True)
            similar_articles = df[['text', 'source', 'url']].head(n_similar).values.tolist()
            return similar_articles
```

Fig 7: Algorithm 2_SVM.ipynb

Fig 8: Algorithm 2 SVM Testing

## What remains to be done

The things that are remaining to be done are
- Testing SVM algorithm for my problem statement
- Testing KNN algorithm for my problem statement
- Generating a comparative study between all used algorithms and results
- Completing Final Report and Presentation

## References

1. Hamdi T., Slimi H., Bounhas I., Slimani Y. (2020) A Hybrid Approach for Fake News Detection in Twitter Based on User Features and Graph Embedding. In: Hung D., D´Souza M. (eds) Distributed Computing and Internet Technology. ICDCIT 2020. Lecture Notes in Computer Science, vol 11969. Springer, Cham. https://doi.org/10.1007/978-3-030-36987-3_17

2. Bharadwaj, Pranav and Shao, Zongru, Fake News Detection with Semantic Features and Text Mining (July 24, 2019). International Journal on Natural Language Computing (IJNLC) Vol.8, No.3, June 2019, Available at SSRN: https://ssrn.com/abstract=3425828.

3. Zhang, Y., Yang, D., He, D., & Huang, J. (2020). GTC: A Graph-based Text Classification Method for Detecting Disinformation. IEEE Transactions on Knowledge and Data Engineering.

4. Liu, Z., Zhang, Y., He, D., & Huang, J. (2019). A graph-based approach for detecting disinformation in social media. Journal of Computer Science and Technology.

5. Arroyo-Fernández, I., Fernández-Vilas, A., & Anido-Rifón, L. E. (2021). Disinformation Detection in Twitter through Topic Modeling and Sentiment Analysis. International Journal of Interactive Multimedia and Artificial Intelligence.

6. Feng, H., Zhang, J., Gao, L., & Zhang, J. (2020). Multi-modal disinformation detection based on transfer learning. IEEE Access.

7. Chen, W., & Wang, Y. (2019). Disinformation Detection on Social Media via Propagation Path Based Learning. Proceedings of the 28th ACM International Conference on Information and Knowledge Management.

8. CoinMarketCap. (2023, February 26). Cryptocurrency Prices, Charts And Market Capitalizations. Retrieved from https://coinmarketcap.com/