

Profiling Cryptocurrency Influencers using Few-shot Learning

Notebook for the PAN Lab at CLEF 2023

Hamna Muslihuddeen¹, Pallapothula Sathvika¹, Shalaka Sankar¹, Shreya Ostwal¹ and Anand Kumar M¹

¹*Dept. of Information Technology, National Institute of Technology Karnataka, Surathkal, India 575025*

Abstract

This research provides a novel method for identifying cryptocurrency influencers on social media in a low-resource environment. The analysis focuses on English-language Twitter messages and divides influencers into impact categories ranging from minimal to massive. With a maximum of 10 English tweets per user, the dataset consists of 32 people per category. By comparing the suggested approach to two baseline models—Usercharacter Logistic Regression and t5-large (bi-encoders) using zero-shot and label tuning few-shot methods—the proposed system is evaluated using the Macro F1 measure. The findings show that the suggested approach operates effectively in low-resource environments and has the potential to be used to further in-depth studies of influencer profiling.

Keywords

low-resource, cryptocurrency, few-shot, zero-shot,

1. Introduction

Cryptography is used by cryptocurrencies, which are digital or virtual tokens, to safeguard their transactions and limit the generation of new tokens. They operate independently of a central authority or middleman, such as a bank or the government, because they are decentralised. In addition to being saved in digital wallets, cryptocurrency is frequently exchanged on online exchanges. They have no government or physical backing, and the market forces of supply and demand determine their price.

The rapidly rising ubiquity and dissemination of online information such as social media text and news improve user accessibility towards financial markets, however, modeling these vast streams of irregular, temporal data poses a challenge. (Ramit Sawhney, Shivam Agarwal, Megh Thakkar, Arnav Wadhwa, and Rajiv Ratn Shah. 2021.) Here, the challenge of effectively modeling large volumes of online information, such as social media text and news, which have

CLEF 2023: Conference and Labs of the Evaluation Forum, September 18–21, 2023, Thessaloniki, Greece

✉ hamnamuslihuddeen.201it221@nitk.edu.in (H. Muslihuddeen); psathvika.201it141@nitk.edu.in (P. Sathvika); shalakasankar.201it158@nitk.edu.in (S. Sankar); shreyaostwal.201it160@nitk.edu.in (S. Ostwal); m1_anandkumar@nitk.edu.in (A. K. M)

🌐 <https://github.com/hamna2905> (H. Muslihuddeen); <https://github.com/SathvikaP> (P. Sathvika); <https://github.com/shalakasankar> (S. Sankar); <https://github.com/Shreya41102> (S. Ostwal)



© 2023 Copyright for this paper by its authors. Use permitted under Creative Commons License Attribution 4.0 International (CC BY 4.0).

CEUR Workshop Proceedings (CEUR-WS.org)

irregular patterns and evolve over time. The authors introduce a novel model, HTLSTM, that uses hyperbolic geometry to better capture the unique characteristics of online information streams, especially in the context of finance. With cryptocurrencies becoming more and more popular, it is very common to find people or organisations with sizable online followings who are able to influence the thoughts and behaviours of their followers with regard to cryptocurrencies. These people or organisations are known as cryptocurrency influencers. These influencers may include traders, analysts, investors, journalists, or cryptocurrency specialists. Since their followers' buying and selling decisions are influenced by their thoughts and suggestions, crypto influencers can have a big impact on the acceptance and value of cryptocurrencies. Many cryptocurrency influencers express their views, analyses, and opinions about various cryptocurrencies and blockchain-related projects on social media sites like Twitter, YouTube, and Instagram.

While some influencers are renowned for their precise market predictions and analyses, others are renowned for their outspoken and divisive viewpoints. To advertise their goods and services, some influencers also work with cryptocurrency initiatives and businesses. But it's vital to remember that not all cryptocurrency influencers are reliable or trustworthy, and some can even participate in dishonest or deceptive behaviour. As a result, it's crucial for people to conduct their own research and use caution when acting on cryptocurrency influencers' advice. But since not everyone can afford this, a solution that can profile crypto influencers in real-time in a matter of milliseconds must be developed. This needs processing as little data as possible in order to get fast and accurate results.

Making use of the Few Shot Learning method is one strategy for solving this issue. A form of machine learning called few-shot learning entails teaching a model to recognise new classes of data with a very small sample size. When discussing cryptocurrency influencers, this refers to training a model to recognise influencers based on a sparse sample of their tweets.

2. Literature Survey

In [1], it's about author profiling, which is the process of identifying characteristics of an author based on their writing, such as gender, age, native language, personality type, etc. It specifically focuses on the gender and age aspects of author profiling in social media, using everyday language to reflect basic social and personality processes. Author profiling involves applying computational tools and linguistic analysis to analyse written materials in order to predict and identify features about the authors, such as their demographics, personality traits, and behaviour. It can be used in areas like marketing, social media analysis, and forensic linguistics.

In [2], discuss few-shot and zero-shot learning in the context of author profiling. They explain that few-shot learning aims to train classifiers with little training data, while zero-shot learning does not use any labelled data. They also describe how the entailment approach can be used

for zero-shot text classification, which relies on neural language models such as BERT trained on large NLI datasets. The authors assess their framework using two tests that determine an author's gender and age based on their written work. They contrast their strategy with a number of industry standards and show that their framework produces competitive performance, especially in situations where only a small amount of labeled data is available.

Overall, the study makes a significant contribution to the field of author profiling and emphasises how zero-shot and few-shot learning have great potential for this task. By utilising these methods, models for author profiling can be created more accurately and more effectively while overcoming the problem of scarce labelled data.

In [3], An approach for active few-shot learning is suggested by the authors. The authors suggest a technique known as FASL (Fast Active Selection for Labelling), which combines few-shot learning and active learning. In order to enhance the performance of the few-shot learning model, FASL seeks to choose the most instructive examples for labelling. The process entails selecting the most instructive examples for labelling from a huge pool of unlabeled examples by first training a few-shot learning model on a limited set of labelled examples. The few-shot learning model is then retrained using the labelled examples that were chosen earlier. On a number of few-shot learning datasets, the authors assess their approach and compare it to other cutting-edge techniques. They show that FASL performs competitively on some datasets and outperforms other approaches on others. They also offer a thorough examination of the influence of various selection procedures as well as the efficacy of the active selection approach.

In [4], adopts machine learning model based on text analysis, using tf-idf for feature extraction of data, and then uses logistic regression model for data training. The idea of linear regression is to fit a straight line through historical data and use this line to predict new data. The objective of logistic regression is to calculate the likelihood that an observation belongs to a given class. A logistic (sigmoid) function is used in the logistic regression model to convert the linear regression equation and limit the output to a range between 0 and 1. This enables us to translate the result into a probability.

In[5], the application of active learning with random forest, a cutting-edge multi-class classifier. The suggested strategy uses an effective active learning algorithm to maximise the combined entropy of a group of samples while minimising information redundancy. The technique performs better than the basic batch mode of active learning when used to adaptively classify undersea mines.

In[6], it mentions that adequate hyper-parameter tuning is crucial for the effective use of SVM classifiers. For this issue, a number of techniques have been employed, including grid search, random search, estimation of distribution algorithms (EDAs), and bio-inspired metaheuristics. The conclusion is backed by experimental findings, and according to the set standards, EDAs are the best techniques for optimising SVM classifier hyperparameter settings. It is crucial to keep in mind that the effectiveness of the remaining algorithms depends on the precise values of the user-defined parameters that are used to control them.

In[7],The term "term frequency inverse document frequency" (TF-IDF) is used to examine the applicability of key terms to corpus documents. The application of the algorithm to various numbers of documents is the main topic of the study. To start, the actions that should be taken for TF-IDF implementation are explained along with their functioning principle. The results are then provided, and the strengths and shortcomings of the TD-IDF algorithm are contrasted in order to verify the conclusions drawn from using the algorithm.

In[8], examines the use of Word2Vec to identify implicit linkages in multi-participant Computer-Supported Collaborative Learning chat sessions. Word2Vec is a potent and one of the most recent Natural Language Processing semantic models used to determine text cohesion and document similarity. The intensity of the semantic ties between two utterances is measured by cohesion scores in this study; the higher the score, the more similar the two utterances are to one another. With Word2Vec, the context before and after each word occurrence in the training dataset is used to compute each embedding. As a result, words that frequently appear together in comparable contexts are represented closer together in the embedded space, while words that do not frequently occur together in similar situations are represented in various areas of this space.

3. Problem Statement

It is important to recognise that not all cryptocurrency influencers are trustworthy or honest, and that some may engage in misleading or manipulative behaviour. People should therefore do their homework and exercise caution while adopting the advice of these influencers.

To solve this problem, we aim to develop a low-resource model that can categorise cryptocurrency influencers on social media into five different groups based on their level of influence: null, nano, micro, macro, and mega. Our concentration is on English-language Twitter messages, and our goal is to create a strong model that, using the Few Short Learning technique[2], can precisely profile and categorise cryptocurrency influencers on social media. By doing this, we intend to give people a useful tool that will help them decide wisely when interacting with social media bitcoin influencers[1].

4. Methodology

The project involves 2 major parts: data processing and developing the model. Under data processing we perform feature extraction to obtain maximum information possible from the limited dataset. Following feature extraction we proceed to develop the model based on few short learning which makes use of active learning.

4.1. Dataset

The dataset used in our few-shot learning task consists of two JSON files: train_text.json and train_truth.json. The train_text file contains 160 JSON objects, each of which contains the Twitter user ID and the corresponding user tweets. The number of tweets per user varies

```
{"twitter user id": "05ca545f2f700d0d5c916657251d010b", "texts": [{"text": "I got $20 on Boston winning tonight, who trying to be"}]
{"twitter user id": "062492818c984febba843b650a4a602e", "texts": [{"text": "@1inch, my favorite aggregator has a sweet booth this"}]}
```

```
{"twitter user id": "0003d5772f14b3147659f37b5aa4399e", "class": "no influencer"}
{"twitter user id": "00230caa0289b84a7a077457435d26b8", "class": "macro"}
```

Figure 1: Train_text and Train_truth

between 2 to 12 tweets. The train_truth file contains 160 JSON objects, each of which contains the Twitter user ID and the profiled class label. The class labels used are null, nano, micro, macro, and mega.

The dataset is relatively small, with only 32 users under each of the five class labels, resulting in a total of 160 entries. This presents a challenge for the few-shot learning task, as the model must learn to recognize and classify users based on a limited amount of training data. Therefore, it is important to carefully select and pre-process the data to ensure that the model can effectively learn and generalize from it.

4.2. Data Pre-processing

Data cleaning and feature extraction are important steps in preparing the dataset for few-shot learning. In our approach, we perform the following steps to preprocess the data:

Combine tweets: To simplify the data and create a single input sequence for each user, we combine all the tweets of a particular user into a single sentence.

Remove punctuation: We remove all punctuation marks from the text, as they do not provide meaningful information for our task.

Convert emojis and emoticons: Emojis or emoticons are often used to communicate thoughts and can be perceived to be more effective than text in social media. Therefore we convert all emojis and emoticons into their corresponding text representations to get more information and to ensure consistency in the data.

Replace hyperlinks: Some tweets contain hyperlinks to other websites. For valid links, we replace the hyperlink with the data scraped from the website or the title of the website. For invalid links, we replace the hyperlink with a blank space.

Overall, these preprocessing steps help to standardize the data and remove irrelevant information, allowing the model to focus on the key features that are important for classifying cryptocurrency influencers. By cleaning and processing the data in this way, we can improve the model's performance and accuracy.

4.3. Feature extraction

Feature extraction is a crucial step in data preprocessing that aims to reduce the dimensionality of raw data by selecting and transforming relevant features to improve the accuracy and

	twitter user id	texts	tweet ids	class	encoded_text
0	0037a672f0ed64b3231bac64853a278d	[{"text": "rt aroundmycitys ape ape give some ..."}]	[{"tweet id": "686f570816c7f81571173988773e755..."}]	nano	3
1	03eaa72711143b521c073d9ac5745923	[{"text": "rt quasimondo now we just need to k..."}]	[{"tweet id": "82923a7577ee8e06030bd2eeaff0c41..."}]	nano	3
2	0409fe210a0edfe258d21e3404e1ce05	[{"text": "rt solanadaily top ten solana colle..."}]	[{"tweet id": "dfe28ca24005b2412580d9ab68e045f..."}]	micro	2
3	05ca545f2f700d0d5c916657251d010b	[{"text": "winkingface httpstcocw3uyna1do"}, {...}]	[{"tweet id": "65408feeb147b509e4bc47280c062e1..."}]	mega	1
4	062492818c984febba843b650a4a602e	[{"text": "takes two matic or one hundred bank..."}]	[{"tweet id": "6539be1639225a7b362e71dba7dcf18..."}]	nano	3
5	0d3700fa5c7c3fce6fd1e1ffd5282f50	[{"text": "data shows bitcoin could be on the ..."}]	[{"tweet id": "d0a8bc3061e9b39ae052b8c91f6e0a3..."}]	macro	0
6	0e1d2c43b93e8e80dc8eb6b29d48b2c1	[{"text": "I'm sending a bunch of eth out to r..."}]	[{"tweet id": "ff58910b19727d93235b3cb28a268a4..."}]	no influencer	4
7	0ed9637249db91cb2c256ec156ce1977	[{"text": "updated price for qtum bnb trxqtum ..."}]	[{"tweet id": "b595be76a45eabea695d37951194515..."}]	no influencer	4
8	0f0942696ae8bcadf0db494cce7333e0	[{"text": "rt avagyuusuzi you'll receive eight..."}]	[{"tweet id": "88a2912ac05afda90c379f59ce8607e..."}]	nano	3

Figure 2: Dataset after preprocessing.

efficiency of machine learning algorithms. In our task of profiling cryptocurrency influencers in social media, we adopt various feature extraction techniques to capture the most important characteristics of the dataset and enable effective analysis, modeling, and decision-making.

We first encode the preprocessed tweet of each Twitter user into a vector using TF-IDF [7]. In natural language processing, the TF-IDF (term frequency-inverse document frequency) sentence transformation approach is used to assess the significance of a given sentence within a corpus or document. It determines each word's relevance in a sentence by comparing its frequency in the sentence to its rarity across the entire document or corpus. By doing so, the original sentences can be changed into vector representations that reflect the semantic significance of the words included within them. This approach makes it simpler to assess how similar two sentences are.

In addition to TF-IDF encoding, we also extract various other features that are unique to our problem statement and help in improving the model's performance. We keep count of the number of tweets per user, the number of hyperlinks mentioned in the tweets, and the number of valid and invalid hyperlinks. We also count the number of cryptocurrency-related terms used in the text, which is a critical aspect of our task, given that we are profiling cryptocurrency influencers. Finally, we also include the word2vec [8] embedding of size 200 of the names of the most popular cryptocurrencies to capture the relationships and similarities between the various cryptocurrencies.

When we combine all the selected features, our feature matrix consists of 973 columns and 160 rows in this case, with each row, representing a unique Twitter user. In general the number of columns in the feature matrix will be equal to the sum of total number of words in the dataset and the additional columns we have added and the number of rows will be the total number of twitter users given in the dataset. This feature matrix serves as the input to our machine learning algorithm, which use it to train and make predictions on new and unseen data. By adopting a comprehensive feature extraction approach, we can capture the most critical and relevant information from the dataset, which helps us to understand better and analyze the dynamics of cryptocurrency influencers on social media.

4.4. Classification model using Few-Shot Learning

In order to execute few shot learning in our project, we employed active learning. Selecting the most informative samples from a dataset for expert annotation or for active learning by

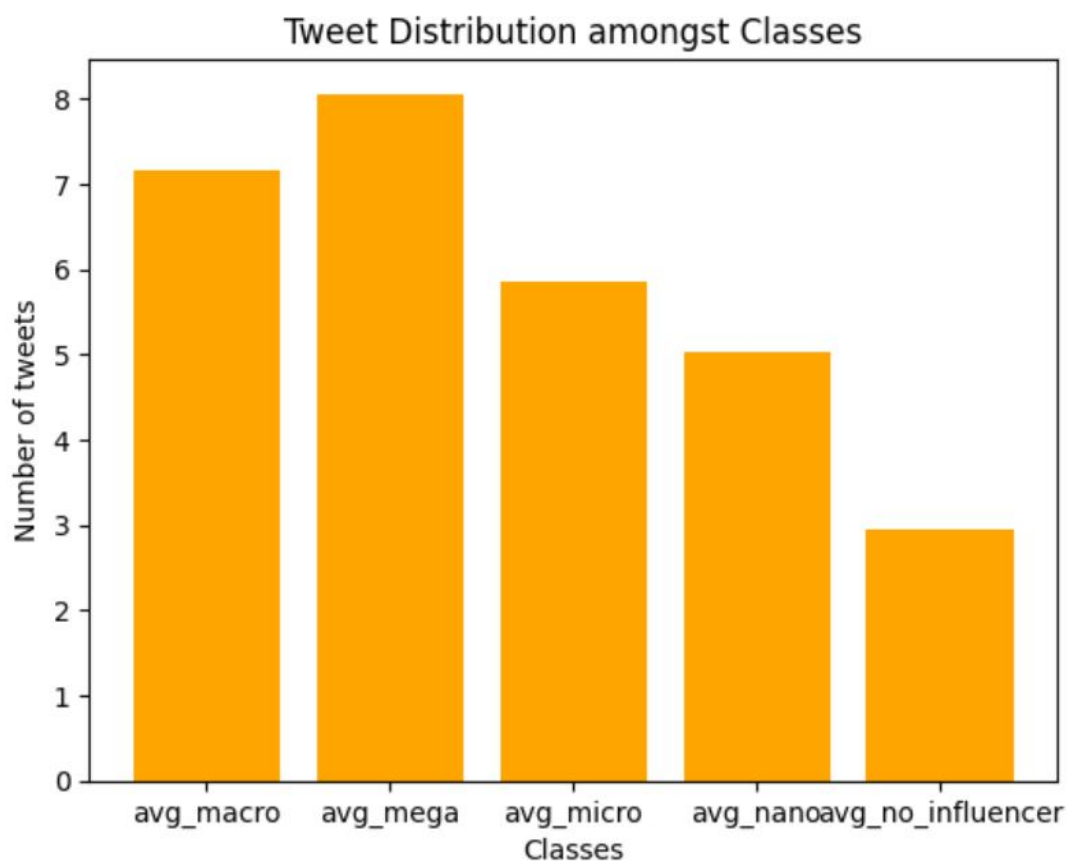


Figure 3: Distribution of tweets under each class of users.

a machine learning algorithm is a key component of the active learning machine learning approach[3]. For labels on particular samples that are most likely to increase the model's accuracy, the algorithm actively asks the user or expert in active learning. There are several sampling techniques available to choose the most illuminating samples for training. Only a few of the strategies are effective at picking samples from a dataset this huge. As follows:

- **Uncertainty Sampling:** According to uncertainty sampling, retraining the model using the most uncertain data points would improve the model's accuracy.
- **Diversity Sampling:** Using a diverse sample The representative test data are gathered using this technique from unsupervised learning. Predicted results for such samples are provided as training data to the algorithm.

In order to perform active learning, we used diversity sampling. K means clustering is the unsupervised learning method we used to choose the samples. We have divided the data into two parts, namely the initial training data and test data, to imitate the real-life scenario where the test data is unsupervised. On the first set of training data, the model is trained. Multinomial Logistic regression[4] is chosen as the base model. As there are 5 classes, the test data is divided into 5

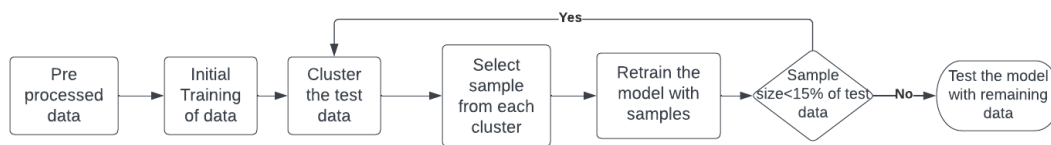


Figure 4: Working of the model.

clusters .By calculating the nearest data point to the training dataset, we choose a representative sample from each cluster. The data point is added to the training data along with the model’s prediction for it, while the test data is eliminated. We keep doing this until we have gathered the required number of representative samples. The steps involved are as follows:

1. Initial splitting of dataset into training and splitting.The desired splitting should contain more testing data. We have used 2 different approaches to split the data into testing and training sets. Firstly we used the conventional method of using 80%training data and 20% testing data. The second method involves finding the median of the number of tweets in each class of user and classifying the data such that users whose number of tweets is more than or equal to the median are classified into the training set and the rest are categorized to the test set.
2. Iterative clustering is performed on the testing data to collect the samples. The number of samples should at least be 15% of the data.
3. The collected samples are appended to the training data and are used to retrain the model
4. The remaining data is used to test the model.

Along with active learning, we also implemented transfer learning to compare and evaluate and choose the better model. To implement transfer learning, we have used logistic regression as both a pre-trained model and a learning model. Logistic regression is the most suitable classification model since it allows to extract coefficients for the learning model. First, we have pre-trained the logistic regression model with 12.5% of the dataset to extract the weights. The weights of the first layer are then set to that of the pre-trained model. The layer is frozen to fix those pre-defined weights. We then retrain the model with another 12.5% of the dataset. The remaining dataset is used for testing.

4.5. Hyperparameters

The performance of active learning was compared with respect to three classifiers: Logistic Regression, Support Vector machine classifier[6], and Random Forest classifier[5]. In Logistic Regression, the model has been initialized with the following parameters: a. “multiclass” which has been set to multinomial, and solver which has been set to ‘lbfgs’, which stands for Limited-memory Broyden-Fletcher-Goldfarb-Shanno. This solver is most suitable for multiclass problems with small to medium-sized datasets. In the Random Forest classifier, the only hyperparameter used was “n-estimators” which was set to 40. The value for this hyperparameter was obtained by the trial-and-error method.

Table 1
Predicted Class Results

Twitter_UserID	Class	Probability
0037a672f0ed64b3231bac64853a278d	Nano	0.453
03eaa72711143b521c073d9ac5745923	Nano	0.39
0409fe210a0edfe258d21e3404e1ce05	Micro	0.42
05ca545f2f700d0d5c916657251d010b	Macro	0.53
062492818c984febba843b650a4a602e	Micro	0.35
f8404b995e68fac31ac3f8318884a0a9	Micro	0.34
fa8d3942f7e1420a5a3d9d19b672f8f2	Nano	0.59
faed79f6d4a62e1c984a62d064f2c8d6	Mega	0.98
fb41227b22d727fbcff4fe780d849de6	Micro	0.49
fec182516cba4b665e2215094bbcc527	Nano	0.50

The hyperparameters defined in The SVM classifier[6] are the kernel, degree, random state, and gamma. the kernel used here was “poly” due to the higher dimensionality and linear separability. The degree for the poly kernel was set to 5. The random seed for various random processes within an algorithm, including random initialization, shuffling of data, or random sampling is set to 0 to ensure that the algorithm’s random processes produce the same results when the code is run again with the same random seed value. The gamma hyperparameter is set to ‘auto’, which means the value of gamma is automatically determined based on the training data. The ‘auto’ option calculates gamma as $1 / n\text{-features}$, where n-features is the number of features in the input data.

5. Results

We evaluated the performance of our created model using a different test dataset. We compared the performance of our model with the accuracy and F1 Macro scores attained by conducting logistic regression, SVM, and Random Forest, each of which is a well-established machine learning approach for classification problems, in order to guarantee the validity of our results.

We were able to validate our method and evaluate its effectiveness in identifying cryptocurrency influencers on social media by comparing the accuracy of our model with that of the logistic regression, random forest, and SVM. Through this comparison, we were able to ensure that our model is reliable and robust and that it can correctly classify influencers into the five categories of null, nano, micro, macro, and mega. Refer TABLE 2

This evaluation procedure provides a quantitative assessment of the accuracy and effectiveness of our developed model and demonstrates its ability to perform well in a low-resource setting when compared to a well-established machine-learning algorithm. In addition to accuracy, we also utilized the F1 score as an evaluation metric for our developed model. The F1 score is a commonly used performance measure in classification tasks, which considers both precision and recall of the model’s predictions.

The F1 scores obtained from our developed model as recorded in TABLE 2 can confidently conclude that our developed model of combining Logistic Regression with Active Learning out-

Table 2
Comparing Models

Model	Accuracy	F1 Macro Score
Random Forest	0.266	0.23
SVM	0.12	0.05
Logistic Regression	0.29	0.26
Random forest with Active Learning	0.29	0.24
SVM with Active Learning	0.14	0.08
Logistic Regression with Active learning	0.34	0.32
Logistic Regression with Transfer Learning	0.34	0.23

Table 3
Evaluation Metrics For Dataset split in Logistic Regression with Active Learning Model

Sno	Dataset split Method	Accuracy	F1 Macro Score
1	80-20 method	0.34	0.32
2	Median Method	0.18	0.23

Table 4
Evaluation Metrics of Logistic Regression with Active Learning

Sno	Precision	Recall	F1 score	Support
1	0.25	0.33	0.29	6
2	0.50	0.40	0.44	5
3	0.25	0.25	0.25	4
4	0.50	0.12	0.20	8
5	0.36	0.56	0.43	9
Accuracy			0.34	32
Macro Average	0.37	0.33	0.32	32
Weighted Average	0.38	0.34	0.33	32

performs logistic regression and other methods in a low-resource situation. This is a significant finding, as it indicates that our model is effective in profiling cryptocurrency influencers on social media using limited resources, and can be a valuable tool for researchers and individuals seeking to make informed decisions when engaging with influencers.

Along with comparing the models TABLE 3 depicts the accuracy obtained for the different types of test-train dataset split. Upon observing it we can say that the convention 40-60 split is a better approach to split the dataset rather than taking the median as the threshold for the split.

Overall, the use of the F1 score provides a more comprehensive and robust evaluation of our developed model's performance and further confirms its superiority over logistic regression in a low-resource setting.

6. Conclusion

In conclusion, the task of profiling cryptocurrency influencers in social media and categorizing related aspects of their influence is a challenging one, especially when working with low-resource settings. However, by focusing on English Twitter posts and making use of Few shot learning, it is possible to extract valuable insights and information about these influencers. It is important to continue developing and refining techniques for analyzing social media data in order to better understand the influence of cryptocurrency influencers and their impact on the wider industry. By doing so, we can gain a deeper understanding of the dynamics of the cryptocurrency world and make more informed decisions about its future.

Acknowledgments

First and foremost, we would like to express our sincere gratitude to CLEF for giving us the chance to take part in this prestigious competition. We deeply value their support in making it possible for us to demonstrate our abilities on this platform. We would also like to take this opportunity to thank our professor and guide Dr Anand Kumar for their guidance, support, and encouragement throughout the entire process. Their mentorship and expertise were invaluable in helping us to shape the direction of our research and to bring our ideas to fruition.

References

- [1] Rangel, F.; Rosso, P.; Koppel, M.; Stamatatos, E.; Inches, G. (2013). Overview of the Author Profiling Task at PAN 2013. CLEF Conference on Multilingual and Multimodal Information Access Evaluation. 352-365.
- [2] Mara Chinea-Rios, Thomas Müller, Gretel Liz De La Peña Sarracén, Francisco Rangel, Marc Franco-Salvador. Zero and Few-Shot Learning for Author Profiling In: NLDB, pp. 333–344, 2022
- [3] Thomas Müller, Guillermo Pérez-Torró, Angelo Basile, Marc Franco-Salvador. Active Few-Shot Learning with FASL In: NLDB, pp. 323–333, 2022
- [4] T. Liu and L. Zhang, "Application of Logistic Regression in WEB Vulnerability Scanning," 2018 International Conference on Sensor Networks and Signal Processing (SNSP), Xi'an, China, 2018, pp. 486-490, doi: 10.1109/SNSP.2018.00097.
- [5] H. T. Nguyen, J. Yadegar, B. Kong and H. Wei, "Efficient batch-mode active learning of random forest," 2012 IEEE Statistical Signal Processing Workshop (SSP), Ann Arbor, MI, USA, 2012, pp. 596-599, doi: 10.1109/SSP.2012.6319769.
- [6] A. Rojas-Domínguez, L. C. Padierna, J. M. Carpio Valadez, H. J. Puga-Soberanes and H. J. Fraire, "Optimal Hyper-Parameter Tuning of SVM Classifiers With Application to Medical Diagnosis," in IEEE Access, vol. 6, pp. 7164-7176, 2018, doi: 10.1109/ACCESS.2017.2779794.
- [7] Qaiser, Shahzad & Ali, Ramsha. (2018). Text Mining: Use of TF-IDF to Examine the Relevance of Words to Documents. International Journal of Computer Applications. 181. 10.5120/ijca2018917395.
- [8] G. Gutu, M. Dascalu, S. Ruseti, T. Rebedea and S. Trausan-Matu, "Unlocking the Power

of Word2Vec for Identifying Implicit Links," 2017 IEEE 17th International Conference on Advanced Learning Technologies (ICALT), Timisoara, Romania, 2017, pp. 199-200, doi: 10.1109/ICALT.2017.120.

- [9] Ramit Sawhney, Shivam Agarwal, Megh Thakkar, Arnav Wadhwa, and Rajiv Ratn Shah. 2021. Hyperbolic Online Time Stream Modeling. In Proceedings of the 44th International ACM SIGIR Conference on Research and Development in Information Retrieval (SIGIR '21). Association for Computing Machinery, New York, NY, USA, 1682–1686. <https://doi.org/10.1145/3404835.3463119>