

# Integration of Back-Propagation Neural Network to Classify of Cybercriminal Entities in Blockchain



Rohit Saxena, Deepak Arora, and Vishal Nagar

**Abstract** Bitcoin is a decentralized, pseudonymous cryptocurrency that has become one among the most demanded digital assets to date. Because of its uncontrolled nature and users' inherent anonymity, it has seen a significant surge in its use for illegal operations. As a result, numerous systems for characterizing diversified entities across the Bitcoin network must be developed. In this work, we offer a way for breaking Bitcoin anonymity using a revolutionary cascade machine learning model that only utilizes a few features taken straight from Bitcoin blockchain data. We gathered approximately 29 million samples from diverse sources and generated data for four different entities: exchanges, gambling, pools, and services. On a dataset balanced using SMOTE and weight of the entities, the back-propagation neural network (BPNN) model was trained and tested. Cross-validation accuracy has been utilized to evaluate the model's accuracy. On the dataset balanced using the weight of the entities, the BPNN model classified the entities with 71.51%, while with SMOTE, the accuracy of classification is 71.22%.

**Keywords** Bitcoin · Blockchain · Back propagation neural network

## 1 Introduction

Bitcoin was coined in 2008, and it has since risen to become the most successful cryptographic money among multiple competitors, boosting the economy by the huge sum of money within only a few years. Bitcoin is a type of cryptocurrency that is made up of a series of computer codes that have a monetary value. All transactions

---

The original version of this chapter was revised: The affiliation of the volume editors has been updated. A correction to this chapter is available at [https://doi.org/10.1007/978-981-16-8826-3\\_57](https://doi.org/10.1007/978-981-16-8826-3_57)

R. Saxena (✉) · D. Arora  
Amity University Uttar Pradesh, Lucknow, India

V. Nagar  
Pranveer Singh Institute of Technology, Kanpur, Uttar Pradesh, India

© The Author(s), under exclusive license to Springer Nature Singapore Pte Ltd. 2022, corrected publication 2022

M. S. Kaiser et al. (eds.), *Proceedings of Trends in Electronics and Health Informatics*, Lecture Notes in Networks and Systems 376,  
[https://doi.org/10.1007/978-981-16-8826-3\\_45](https://doi.org/10.1007/978-981-16-8826-3_45)

and payments are completed over the Internet in this setting. Bitcoin differs from traditional Internet transactions in that it is based on a peer-to-peer (P2P) network that is not linked to a centralized third-party organization, such as an e-bank, a notary, or any other traditional online financial service provider that supervises, monitors, and approves electronic payment transactions. Instead, Bitcoin users have absolute power over what they want to do with their own money because they may freely order how and when to utilize digital money. Bitcoin has grown in popularity, attracting an increasing number of customers who want to use it as a payment option in numerous organizations. Bitcoin is frequently described as being “quick,” “convenient,” “tax-free,” and “revolutionary.”

Bitcoin was, and still is, referred to as an anonymous currency in some contexts. Although Bitcoin addresses, unlike traditional financial systems, has nothing to do with any real-world identity at the protocol level, this fact does not guarantee great anonymity. Bitcoin transactions are broadcast in cleartext over a peer-to-peer network and then stored in a massively replicated shared database after being confirmed by miners. A common technique to improve Bitcoin privacy is to use a different address for each transaction.

Illegal laundering of money and financing terrorism have been linked to Bitcoin and other cryptocurrencies [1, 2], as well as the online drug trade [3]. On the online black marketplaces like Silk Road 3, Alpha Bay, and Valhalla, cryptocurrencies have been linked to cybercriminal behaviors. Consumers can use Bitcoins to buy cybercrime-as-a-service, hacking tools, malware, stolen credit card information, and compromised login and password combinations, for example.

Bitcoin is frequently used to propagate ransomware around the world. WannaCry, a ransomware attack that began in May 2017, spread fast around the world. In a few hours, WannaCry labeled “the worst ransomware outbreak in history,” infected over 300,000 machines in 150 countries. This attack was extremely destructive since it was a worm that looked for new computers and systems to infect, rather than just a ransomware application. WannaCry encrypts all files on infected devices, rendering them unreachable to the victim until the culprit is paid at least \$300 in Bitcoin (s). The ransomware payments were received using three hardcoded Bitcoin addresses/wallets. A total of 335 payments, totaling 51.91182371 Bitcoin or US\$144,010.54, has been sent into the three Bitcoin wallets as of June 20, 2017.

Such obstacles and challenges, unfortunately, still exist today. Due to its features, particularly its pseudo-anonymity, Bitcoin has become the preferred payment system for illicit activities. Bitcoin transactions are linked to public keys or addresses rather than real-world identities, and the latter does not require any verified personal data to create. While some enterprises openly provide their addresses when it is essential to delivering their services, others hide their spending habits by using privacy-enhancing payment channels or mixing services. This is a regular occurrence in entities associated with tor markets, ransom payments, frauds, and thievery.

The purpose of this study is to investigate if an uncovered cluster can be classified into one of the categories of exchange, gambling, hosted wallet, merchant services, mining pool, mixing, ransomware, fraud, tor market, or others. This study will increase the transparency of the ecosystem, encouraging businesses and

consumers to utilize Bitcoin as a payment method and expand the economy without resorting to criminal actions. This research works reviewed in this paper focus on the classification of Bitcoin addresses based on their activity.

The remainder of the paper contributes as follows: Sect. 2 reviews the related work, followed by a discussion of conceptual overview of the research in Sect. 3. Section 4 gives an overview of the data preparation procedure; Sect. 5 suggests the methodology adopted in the study. The results obtained are outlined in Sect. 6, followed by the conclusion and future scope in Sect. 7.

## 2 Related Work

Several publications have attempted to disprove Bitcoin's alleged pseudo-anonymity. The first way of breaking down Bitcoin's anonymity used network analytical techniques over addresses mixed with open-source data from Wikileaks, demonstrating that Bitcoin user addresses can be connected [4]. A subsequent technique required direct interaction with the network by submitting transactions and clustering public keys using co-spend algorithms, which led to the discovery of 1.9 million Bitcoin addresses related to real entities or pseudo-identities [5]. Using an open-source framework, the Bitcoin blockchain was processed; public keys were clustered; clusters were tagged, and the network was shown. The system was able to identify an address holding 111,114 BTC pertaining to a Silk Road cold wallet and precisely estimate ransoms delivered to CryptoLocker using only an address provided by a victim on a forum as a lead [6]. Another technique was to use statistical analysis to figure out how its users behaved when sending, receiving, and keeping money. Unlike large portions of transactions moving minimal quantities of coins and the specific subject of study, countless numbers of transactions sending more than 50,000 BTC all at once, this approach discovered that the vast majority of coins remain hidden in addresses that haven't been associated in outgoing transactions [7]. The k-means algorithm was used to cluster a portion of the Bitcoin blockchain with the purpose of detecting anomalous behavior, uncovering anomalous transactions, and detecting abnormal behavior from some clients suspected of money laundering [8]. In [9], the authors developed a method for correlating Bitcoin users' pseudonyms behind NAT with the public IP address of the host where the transaction is generated. The attack's goal was to use one octet of outbound connections to identify each client. Even if they connect to the Bitcoin network via Tor, the approach outlined in [10] linked the sessions of unreachable nodes. The authors do this by employing a novel method that organizes block requests made by nodes in a Bitcoin session graph. The modified Bitcoin client is likewise vulnerable to this attack [10]. A transaction clustering approach was presented and executed based on the analysis of propagation times on four popular cryptocurrencies: Bitcoin, Zcash, Dash, and Monero [11]. Biryukov and Tikhomirov [12] was the first to analyze the prevalence of cybercriminal entities in the Bitcoin ecosystem. They trained unsupervised machine learning classifiers

using the dataset provided by the data provider, which used three methods of clustering Bitcoin transactions to categorize entities: co-spend, intelligence-based, and behavior-based. A multiclass classification on Bitcoin blockchain clusters was carried out with the objective of seeing if supervised machine learning algorithms could be helpful to predict the category of an undiscovered cluster given a set of previously recognized clusters as training data [13]. On the basis of transaction history summarization, a multiclass service identification technique in Bitcoin was presented. For improved identification, the suggested technique gets transaction history and analyses the working of the retrieved transactions [14]. Transaction history summary for Bitcoin addresses and entity classification was included as new features. To create a model for machine learning classification for detecting abnormalities of Bitcoin network addresses, the transaction history summary was composed of basic statistics, supplementary statistics, and transaction moments [15]. A cascade machine learning model combined with a sufficient collection of input attributes directly extracted from Bitcoin blockchain data was used to demonstrate a technique to challenge Bitcoin anonymity through entity characterization [16].

### 3 Conceptual Overview

#### 3.1 *Cryptocurrencies*

Cryptocurrency is conceivably the most secured form of digital money exchange, as it is built on a decentralized network via the Internet and uses cryptography to perform financial transactions. Although the conceptual theory behind the cryptocurrency technology was proposed in 1991 [17], it was only introduced to the world in 2009 as bitcoin [18]. Unlike traditional money exchange systems, bitcoin relies on a dispersed network of nodes known as miners to keep track of transactions. Each mining node maintains a list of transactions in blocks, each of which contains the preceding block's SHA 256 cryptographic hash. Because this process is persistent, and transaction blocks are stored in all nodes, it is almost difficult to change it. Each node checks all transaction history before committing any transaction to validate that the amount to be sent is correct. Following validation, each transaction in a block conducts repeated hashing procedures using the sender's public key and the preceding transaction's hash value.

#### 3.2 *Anonymity*

Anonymity is defined as a way of obtaining “freedom from identification, concealment, and lack of distinction” [19], and it can also be defined as a phenomenon in which one can conceal one's identity from others [20]. With the advent of the

Internet and subsequent advancements in electronic commerce, communications, and social media, as well as developments like Web 2.0, there is indeed a growing discussion regarding anonymity, particularly in online environments. The proponents of anonymity see online anonymity as an essential tool for preserving information privacy by shielding personal data from untrustworthy platforms and parties [20]. Anonymity, on the other hand, is regularly exploited, creating a climate conducive to hate speech and libelous remarks by those who act irresponsibly with impunity [20, 21]. Communications over the Internet can be formed with a high degree of certainty, hiding the identity of the communication's source, thanks to the development of public-key cryptography and software agents in the 1990s, such as anonymous remailer servers [22]. These techniques cleared the ground for the formation of pseudonymous entities in Internet communication, which can send and receive messages while keeping the originator's identity hidden. Pseudonymity varies from anonymity in that anonymity demands the complete elimination of all identification information, whereas pseudonymity allows for the construction and maintenance of a pseudo/alternate identity, allowing for partial concealing of the true identity information [19, 23].

### 3.3 *Entity Categorization*

One of the most active kinds of enterprises is the exchange, which is a global digital marketplace, wherein traders can trade cryptocurrencies using different fiat (money made legal tender by a government edict) or other digital currencies. As stated in [24], exchanges serve as the “front and exit doors” to the bitcoin realm and are perfect for concealing unlawful activities. The **darknet market** is another option. These are online marketplaces where users can buy narcotics, ammunition, and other commodities and services that are prohibited in most nations. These crypto-markets promote legal and illegal transactions among their customers by using electronic currencies [25]. Furthermore, as described in [26], so-called **mixers** are services that allow users to obfuscate processes. Mixed transactions, on the other hand, boost user privacy and can be used to launder unlawful payments. A few more categories are listed in [12] which are **gambling**, an entity that provides Bitcoin-accepting gambling services, such as Lucky Games and Nitrogen Sports. In exchange for a reward proportional to their contribution to a block's solution, **mining pools** are made up of distributed miners that pool their processing power over a mining network. AntPool and BTC Top are two examples.

## 4 Data Preparation

The dataset for this study was gathered from the Blockchair and WalletExplorer repositories. Samples for three months were chosen for the study, namely December 2020, November 2020, and January 2021.

### 4.1 Raw Dataset

Blockchair is a blockchain explorer which has the dumps for the cryptocurrencies such as Bitcoin, Ethereum, Bitcoin Cash, Dogecoin, Litecoin, Bitcoin-SV, and ZCash. The transaction dataset collected from the Blockchair comprises of `block_id`, `hash`, `time`, `size`, `weight`, `version`, `lock_time`, `is_coinbase`, `has_witness`, `input_count`, `output_count`, `input_total`, `output_total`, `input_total_usd`, `output_total_usd`, `fee`, `fee_used`, `fee_per_kb`, `fee_per_kb_usd`, `fee_kb_kwu`, `fee_per_kwu_usd`, `cdd_total`. WalletExplorer is a bitcoin explorer with address grouping and wallet labeling. The transaction dataset collected from this repository contains `date`, `received from`, `received amount`, `sent amount`, `sent to`, `balance`, and `transaction`.

### 4.2 Data Preprocessing

When considered separately, the features in both the dataset were insufficient for both training and testing. To make the best use of the collected dataset, the features of both datasets were merged on the basis of transaction hash. After merging, 29,228,184 feature-rich samples were obtained. Out of 29,228,184 samples, 427,625 samples were be labeled. The available labeled entities were *exchange*, *gambling*, *mining pool*, *mixing services*.

### 4.3 Data Cleaning

The labeled dataset thus obtained was cleaned by performing the following tasks:

- a **.Handling null and infinite values:** Deep learning models do not handle the null and infinite values; therefore, such samples were dropped.
- b **.Encoding string values to integer:** Using the encoder library of scikit-learn, string values were encoded to an integer to make them adaptable for deep learning models.

#### 4.4 Selection of Features

The transaction hash is generated randomly corresponding to the transaction performed by the Bitcoin users. Because of this random nature, they were dropped. Moreover, few features such as `block_id` were of low co-relation and were not participating in the training and testing of a classification model. Hence, they were also removed. The features which got selected are as follows: `label`, `size`, `weight`, `version`, `lock_time`, `is_coinbase`, `has_witness`, `input_count`, `output_count`, `inout_total_usd`, `output_total_usd`, `fee_usd`, `fee_per_kb_usd`, `fee_per_kwu_usd`, `cdd_total`, `rec/sent`, `amount`. At this step, the dataset is now ready for training and testing.

### 5 Methodology

To classify and predict the accuracy of the classification, the dataset obtained was trained and tested over the back-propagation neural network (BPNN) and evaluated on cross-validation (CV) accuracy.

#### 5.1 Balancing of Classes/Entities

In the dataset of 427,625 samples, the entity exchanged had the majority of samples 335,847, i.e., 78.53%, while the entity gambling had merely 2254 samples which were 0.53% of the total samples. This led to the issue of class imbalance. To handle this issue, SMOTE [27] and the weight of the entities are employed. SMOTE stands for **s**ynthetic **m**inority **o**versampling **t**echnique which produces synthetic samples for the minority classes so that the utilization of imbalanced classes can be enhanced. The other technique, i.e., the weight of the entities oversamples the minority classes and under-samples majority classes so that there's a uniform distribution of samples of all four entities. The weight of an entity was calculated using the following formulae:

$$w_i = n / (k * n_i)$$

where.

- $w_i$     the weight of the entity  $i$ ,
- $n$       number of dataset samples,
- $k$       number of dataset entities,
- $n_i$     number of dataset samples of entity  $i$

With this approach, samples of mining pool, gambling, and mixing services were oversampled while that of exchange was under-sampled.

5.2 Training and Testing of Model

In this study, the back-propagation neural network (BPNN) was validated using the SMOTE and weight of entities techniques on balanced dataset samples. The y-axis had the labels (exchange, pool, services, and gaming), whereas the X-axis contains the remaining attributes. Samples from the training and testing datasets are distributed 60% and 40%, respectively. Both approaches were trained and tested using the TensorFlow libraries on Google’s colaboratory.

6 Result

The BPNN was trained and tested to classify the entities and predict the CV accuracy. The model was trained for the dataset samples prepared using SMOTE and the weight of the entities. The classification of entities is shown in Fig. 1, while the results obtained on training and testing of the model are as shown in Table 1.

From Table 1, it is evident that the class balancing carried out using the weight of the entities is giving a slightly better prediction than SMOTE strategy with the exception that the number of samples trained using SMOTE is more than that of the weight of the entities.

Fig. 1 Classification of entities

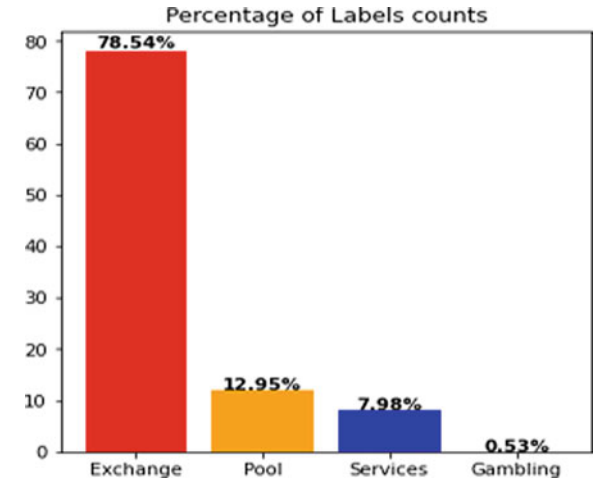


Table 1 Comparison of approaches based on CV accuracy

Class balancing strategies	CV accuracy (%)
SMOTE	71.22
Weight of entities	71.51



## 7 Conclusion and Future Work

Back propagation neural networks were trained on dataset samples of Bitcoin transactions obtained from the Blockchair and WalletExplorer repositories for this study. The transaction hash was then used to merge it. The issue of class imbalance was overcome by carrying out necessary oversampling using SMOTE and oversampling and under-sampling as required using the weight of the entities. The model was evaluated for both the approaches of class balancing on the basis of cross-validation accuracy. The strategy of the weight of the entities seems to be more accurate in comparison to SMOTE.

The accuracy of the model can be enhanced by developing a hybrid approach employing the heuristics on the existing model. Moreover, the availability of cluster is also the limitation of the research; therefore, available datasets can be clustered that can widen the scope of the research.

## References

1. Irwin ASM, Milad G (2016) The use of crypto-currencies in funding violent jihad. *J Money Laundering Control* 19(4):407–425. <https://doi.org/10.1108/JMLC-01-2016-0003>
2. Pflaum I, Hateley E (2014) A bit of a problem: national and extraterritorial regulation of virtual currency in the age of financial disintermediation. *Georgetown J Int Law* 45(4):1169–1215
3. Martin J (2014) Lost on the silk road: online drug distribution and the ‘cryptomarket.’ *Criminol Crim Just* 14(3):351–367. <https://doi.org/10.1177/1748895813505234>
4. Reid F, Harrigan M (2013) An analysis of anonymity in the bitcoin system. In: Altshuler Y, Elovici Y, Cremers A, Aharony N, Pentland A (eds) *Security and privacy in social networks*. Springer, New York. [https://doi.org/10.1007/978-1-4614-4139-7\\_10](https://doi.org/10.1007/978-1-4614-4139-7_10)
5. Meiklejohn S, Pomarole M, Jordan G, Levchenko K, McCoy D, Voelker GM, Savage S (2016) A fistful of bitcoins: characterizing payments among men with no names. *Commun. ACM* 59, 4 (April 2016):86–93. <https://doi.org/10.1145/2896384>
6. Spagnuolo M, Maggi F, Zanero S (2014) BitIodine: extracting intelligence from the bitcoin network. In: Christin N, Safavi-Naini R (eds) *Financial cryptography and data security*. FC 2014. *Lecture Notes in Computer Science*, vol 8437. Springer, Berlin, Heidelberg. [https://doi.org/10.1007/978-3-662-45472-5\\_29](https://doi.org/10.1007/978-3-662-45472-5_29)
7. Ron D, Shamir A (2013) Quantitative analysis of the full bitcoin transaction graph. In: Sadeghi AR (ed) *Financial cryptography and data security*. FC 2013. *Lecture Notes in Computer Science*, vol 7859. Springer, Berlin, Heidelberg. [https://doi.org/10.1007/978-3-642-39884-1\\_2](https://doi.org/10.1007/978-3-642-39884-1_2)
8. Hirshman J, Huang Y, Macke S (2013) Unsupervised approaches to detecting anomalous behavior in the bitcoin transaction network, 3rd ed. Technical report, Stanford University
9. Biryukov A, Khovratovich D, Pustogarov I (2014) Deanonymisation of clients in bitcoin P2P network. In: *Proceedings of the 2014 ACM SIGSAC conference on computer and communications security (CCS ‘14)*. Association for Computing Machinery, New York, 15–29. <https://doi.org/10.1145/2660267.2660379>
10. Mastan ID, Paul S (2017) A new approach to deanonymization of unreachable bitcoin nodes. In: *International conference on cryptology and network security*. Springer, Cham, pp 277–298
11. Biryukov A, Tikhomirov S (2019) Deanonymization and linkability of cryptocurrency transactions based on network analysis. *2019 IEEE European symposium on security and privacy (EuroS&P)*, pp 172–184. <https://doi.org/10.1109/EuroSP.2019.00022>

12. Sun Yin H, Vatraru R (2017) A first estimation of the proportion of cybercriminal entities in the bitcoin ecosystem using supervised machine learning. In: 2017 IEEE international conference on big data (Big Data), pp 3690–3699. <https://doi.org/10.1109/BigData.2017.8258365>
13. Harlev MA, Sun Yin H, Langenheldt KC, Mukkamala R, Vatraru R (2018) Breaking bad: de-anonymizing entity types on the bitcoin blockchain using supervised machine learning. In: Proceedings of the 51st Hawaii international conference on system sciences. Hawaii international conference on system sciences. <https://doi.org/10.24251/hicss.2018.443>
14. Toyoda K, Ohtsuki T, Mathiopoulos PT (2018) Multi-class bitcoin-enabled service identification based on transaction history summarization. In: 2018 IEEE international conference on internet of things (iThings) and IEEE green computing and communications (GreenCom) and IEEE cyber, physical and social computing (CPSCom) and IEEE smart data (SmartData), pp 1153–1160. [https://doi.org/10.1109/Cybermatics\\_2018.2018.00208](https://doi.org/10.1109/Cybermatics_2018.2018.00208).
15. Lin YJ, Wu PW, Hsu CH, Tu IP, Liao SW (2019) An evaluation of bitcoin address classification based on transaction history summarization. In: 2019 IEEE international conference on blockchain and cryptocurrency (ICBC). IEEE, pp 302–310
16. Zola F, Eguimendia M, Bruse JL, Urrutia RO (2019) Cascading machine learning to attack bitcoin anonymity. In: 2019 IEEE international conference on blockchain (Blockchain). IEEE, pp 10–17
17. Haber S, Stornetta WS (1991) How to time-stamp a digital document. *J Cryptol* 3:99–111. <https://doi.org/10.1007/BF0019679199111>
18. Nakamoto S (2008) Bitcoin: a peer-to-peer electronic cash system. *Decentralized Bus Rev* 21260.
19. Scott SV, Orlikowski WJ (2014) Entanglements in practice: performing anonymity through social media. *MIS Q* 38(3):873–893
20. Brazier F, Oskamp A, Prins C et al (2004) Anonymity and software agents: an interdisciplinary challenge. *Artif Intell Law* 12:137–157. <https://doi.org/10.1007/s10506-004-6488-5>
21. Levmore S (2010) The internet’s anonymity problem. In: Lemore S, Nussbaum M (eds) *The offensive internet: speech, privacy, and reputation*. Harvard University Press, Cambridge
22. Michael FA (1995) Anonymity and its enmities (1995) 1 *J Online Law* art. 4. Available at SRN: <https://ssrn.com/abstract=2715621>
23. Michael Froomkin A (1999) legal issues in anonymity and pseudonymity. *Inf Soc* 15(2):113–127. <https://doi.org/10.1080/019722499128574>
24. Moore T, Christin N (2013) Beware the middleman: empirical analysis of bitcoin-exchange risk. In: Sadeghi AR (ed) *Financial cryptography and data security*. FC 2013. Lecture Notes in Computer Science, vol 7859. Springer, Berlin, Heidelberg. [https://doi.org/10.1007/978-3-642-39884-1\\_3](https://doi.org/10.1007/978-3-642-39884-1_3)
25. Christin N (2013) Traveling the silk road: a measurement analysis of a large anonymous online marketplace. In: Proceedings of the 22nd international conference on World Wide Web (WWW ‘13). Association for computing machinery, New York, 213–224. <https://doi.org/10.1145/2483888.2488408>
26. Moser M (2013) Anonymity of bitcoin transactions
27. Chawla NV, Bowyer KW, Hall LO, Kegelmeyer WP (2002) SMOTE: synthetic minority over-sampling technique. *J Artif Intell Res* 16:321–357