# Bitcoin De-Anonymization

ChengMing 517030910287

June 21, 2020

## 1  Introduction

Bitcoin uses a pseudonym system. Although all transactions of bitcoin are public, what is publicly displayed is only an account code of the transaction user, that is, "address". This address can be generated at will, so it can't be connected with real life identity. It objectively provides anonymity, but it is also used by lawbreakers. However, bitcoin can not guarantee complete anonymity. When bitcoin transactions have contact with the real world, it may disclose real information (such as using bitcoin for shopping, or using bitcoin to exchange for legal currency). Because all transactions on the block are public, this provides a breakthrough for the anonymization of bitcoin.

There are two ways to de-anonymize bitcoin network. One is de-anonymize on the network layer. Using the broadcast of bitcoin transaction in P2P network, we can determine the broadcast source by manipulating multiple nodes and other means, so as to associate the user address with the IP address and realize de-anonymity. The second is the application layer de anonymization, which mainly relies on the transaction map analysis, builds the transaction network in the blockchain, and matches the real transaction information obtained in various ways, so as to realize the de anonymization. In this project, I did some research on the application layer de anonymization

## 2  Basic idea

In the specific implementation process, we first need to obtain the transaction information on the bitcoin blockchain, then extract useful node and edge information from the transaction information, so as to build the bitcoin transaction network. Secondly, we need to obtain the information about the connection between the bitcoin transaction and the real world. I selected the record of transactions leaked on bitcoin transaction platform Mt.Gox in 2014. The account in this transaction record is registered with real world identity. If this part of transaction record can be matched with the transaction record of bitcoin blockchain, then our de anonymization work is successful

## 3  Data Preparation

### 3.1  bitcoin blockchain

There are two way to get the full data of bitcoin blockchain:

1. download bitcoin wallet, then it will download all blockchain automatically. The data we download is in the .dat file. The file has a very standard format. The specific format can be seen in the official document. There are multiple blocks in a .dat file. Each part of the block occupies the corresponding number of bytes, which can be divided easily and read in order.

| 数据项 | 字节 | 字段 | 说明 |
|---|---|---|---|
| Magic NO | 4 | 魔数 | 常数0xD9B4BEF9 |
| Blocksize | 4 | 区块大小 | 用字节表示的该字段之后的区块大小 |
| Blockheader | 80 | 区块头 | 组成区块头的几个字段 |
| Transaction counter | 1-9 | 交易计数器 | 该区块包含的交易数量，包含coinbase交易 |
| Transactions | 不定 | 交易 | 记录在区块里的交易信息，使用原生的交易信息格式，并且交易在数据流中的位置必须与Merkle树的叶子节点顺序一致 |

Figure 1: .dat file format

2. get block data or transaction data by requesting API, like https://api.blockcypher.com/v1/btc/main/txs/txhash. The data get from API is in json format. Here is an example:



Figure 2: API json format

## 3.2 Mt.Gox leaked DataSet

Mt.Gox leaked DataSet can be download in many dataset website, such as kaggle and github. here is one source: https://www.kaggle.com/xblock/mtgox-leaked-transaction
The data I need in this project is in complet_edge_v2.csv, the format of this file is shown below.

| | Source | Target | Trade_Id | Bitcoins | Money | Money_rate | Date | label |
|---|---|---|---|---|---|---|---|---|
| 1 | 895 | 3931 | 35372 | 23.02 | 18.061 | 0.784578627 | 2011/4/1 0:28:... | 0 |
| 2 | 895 | 722 | 35373 | 10 | 7.8 | 0.78 | 2011/4/1 0:28:... | 0 |
| 3 | 895 | 3605 | 35374 | 10 | 27.3 | 0.78 | 2011/4/1 0:28:... | 0 |

Figure 3: API json format

- **Source:** bitcoin seller

- **Target:** bitcoin buyer/receiver

- **Trade_Id:** Id of this transaction

- **Bitcoins:** the amount of bitcoins in this transction

- **Money:** the amount of dollars for buy bitcoin

- **Money_rate:** Unit bitcoin price

- **Date:** Time of transaction

- **label:** the type of user

# 4 Process data

In order to get useful information, we need to process the data.

## 4.1  bitcoin blockchain

our goal is to obtain particular information of every transction in blockchain, include input addresses, input value, output addresses and output value. Also, we need to obtain the time of each block.

The time of each block is in the blockheader, so can obtain it easily. The output value and public key is in output transaction part. we can obtain them too. and the address can be generate by the public key. But the address and public key of input is not in corresponding transction. we can only get the previous hash, which point to the source of bitcoin. We should search transaction whoes hash is same as the previous hash, and get the address and public key of its output.

After getting all these information, we should save them in json form. Some more details such as how to generate address by public key and how to read .dat file in certain order will not be mentioned here. You can see the detail of this part is in DataProcess.py



Figure 4: blockchain after process

## 4.2  Mt.Gox leaked DataSet

The Mt.Gox leaked DataSet don't need to change too much, just change the time into the form of seconds since 1970.1.1 which is easily to compare with other time. Then extract useful parts of it.

# 5  Network De-Anonymization

## 5.1  Algorithm

For each transaction (source, target) in the Mt.Gox data set, looking for a transaction (input, output) in the revenue and expenditure network with the same amount of the bitcoin, while the difference between the transaction time is less than 6 block times. When such a transaction is found, it means that source and input, target and output might belong to the same user, and their confidence degree $C_{ij}$ should plus 1.

Traverse all transactions in Mt.Gox data set, calculate in the above way, the confidence between each account and each address is obtained. For any bitcoin address, find out which Mt.Gox account has the highest level of confidence, then we think the address belongs to this account. If there is no match or there are more then one account with the highest confidence, it will not be assigned, or we say it is not de-anonymized.

## 5.2 Result

The whole data set is too large for my computer to handle. So I choose all the transaction in 2011.7.19 for experiment. The final result is shown below.



(a) Match Number

(b) De-Anonymization rate

Figure 5: Result

As you can see, 40% of the transactions in the data set match the only transactions in the bitcoin network. In addition, 30% of transactions match more than one transaction, probably because there are many integer values that often appear in the transaction (such as: 1, 10,100, etc.)

After complete the algorithm, nearly 60% of the addresses are de anonymous. The rest 40% of unsuccessful de-anonymization is because that the transactions in the dataset do not cover all transactions in blockchain and there are transactions that do not match the address.

# 6    Future Work

Based on the number of transactions between different addresses, clustering algorithms such as K-means can be used to cluster the addresses. By observing whether the addresses belonging to the same user in the de-anonymization result are in the same cluster, we can analyze the effectiveness of the de-anonymization result.