

Bitcoin Transactions Network

Questions asked:

Is it possible to represent bitcoin transactions as a network?

Can you infer ownership of bitcoin addresses by analyzing transactions?

Can the tools offered to us by Network Analysis help us in this task?

The first task that was tackled in this project was data acquisition.

The bitcoin blockchain is by definition a huge database containing all transaction that were ever validated by the blockchain. Therefore, if we wish to gather data about transactions, we must query the blockchain for it.

To guarantee the integrity of the data, I decided to install the reference bitcoin node implementation, known as bitcoin-core, and to run a full node, in order to be able to directly query other nodes for blockchain data.

Bitcoin-core is a open-source version of the bitcoin software, and currently it's the most commonly used software by other nodes. If a node does not run bitcoin-core , it probably runs a modified version of it.

The software is available on github at

<https://github.com/bitcoin/bitcoin>

Once the software is installed and running, it will begin a (very slow) process of finding and querying other nodes for blockchain data. If one does not have a dedicated connection to a full node that is able to provide blockchain data, bitcoin-core will search for nodes at random.

This process might take a very long time.

This is crucial because in order for the blockchain to be accepted and for the data to be reliable, every node must process each block and validate it independently. Only by doing this we can be sure that the data we have is reliable and that it represents correctly the transactions that have occurred on the blockchain.

Data Acquisition

Once bitcoin-core is running, it's possible to query other nodes for blockchain data. A block is returned by the blockchain as a JSON object containing information about the block, such as height, hash, time it was mined, and most importantly, a list of transactions contained in that block

First big problem

While querying for blocks does not take a substantial amount of time, querying the blockchain for transactions is a LOT slower. Moreover, each block contains thousands of transactions.

This is why i decided that querying for transaction data had to be done through a websocket API.

Gathering block data from the blockchain and transaction data through the API allowed me to develop a simple and moderately fast pipeline for data extraction.

Typical transaction hex-dump

Each transaction comes in a raw hex-like format, like so:

```
02000000013066b6fd51e0aea6e456aac1b894a076a50634a1517c3c81219acee3
41892a600f00000006a47304402200223ed42107279f13dc4e35af13a235557b1aa7
b8c0ba5125026d47310241b1d0220081bd0ff5ee2629cd1e51d8aef120170df2421
682e580a641efd7c6b26618d3a012102d46e0710b40c285aec6c9e3d0271d13be5
2acb216f8c5b81d8875abcd5e35860ffffffff012681d0060000000001976a91442cdfb0
3341f83724902dc14f159f489d79beb0488ac00000000
```

This format encodes all the information stored in a transaction.

Edge extraction

From the hex dump we extract the information we need for our networks, such as:

- Sender(s) address
- Receiver(s) address
- Amount sent to each receiver

After extracting edge information, we create csv files for each block. This way we are able to feed R with our data and plot our networks

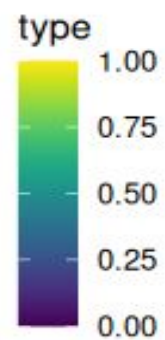
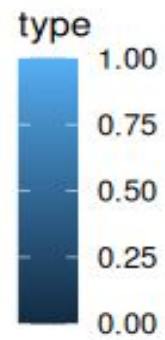
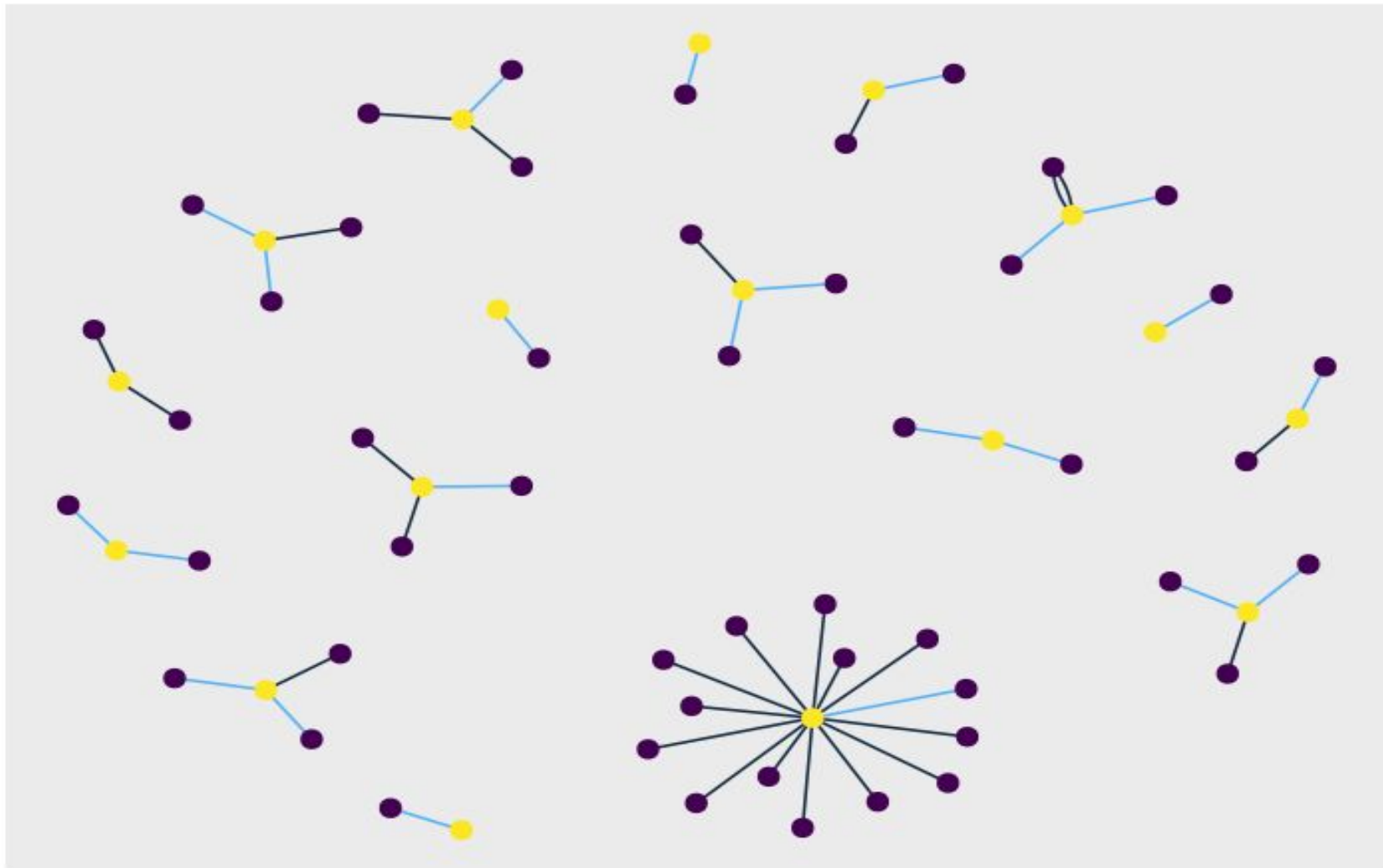
Transaction network

A transaction network is composed of two different kind of nodes:

Addresses and Transaction

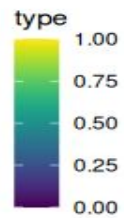
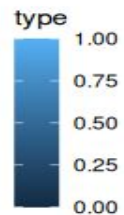
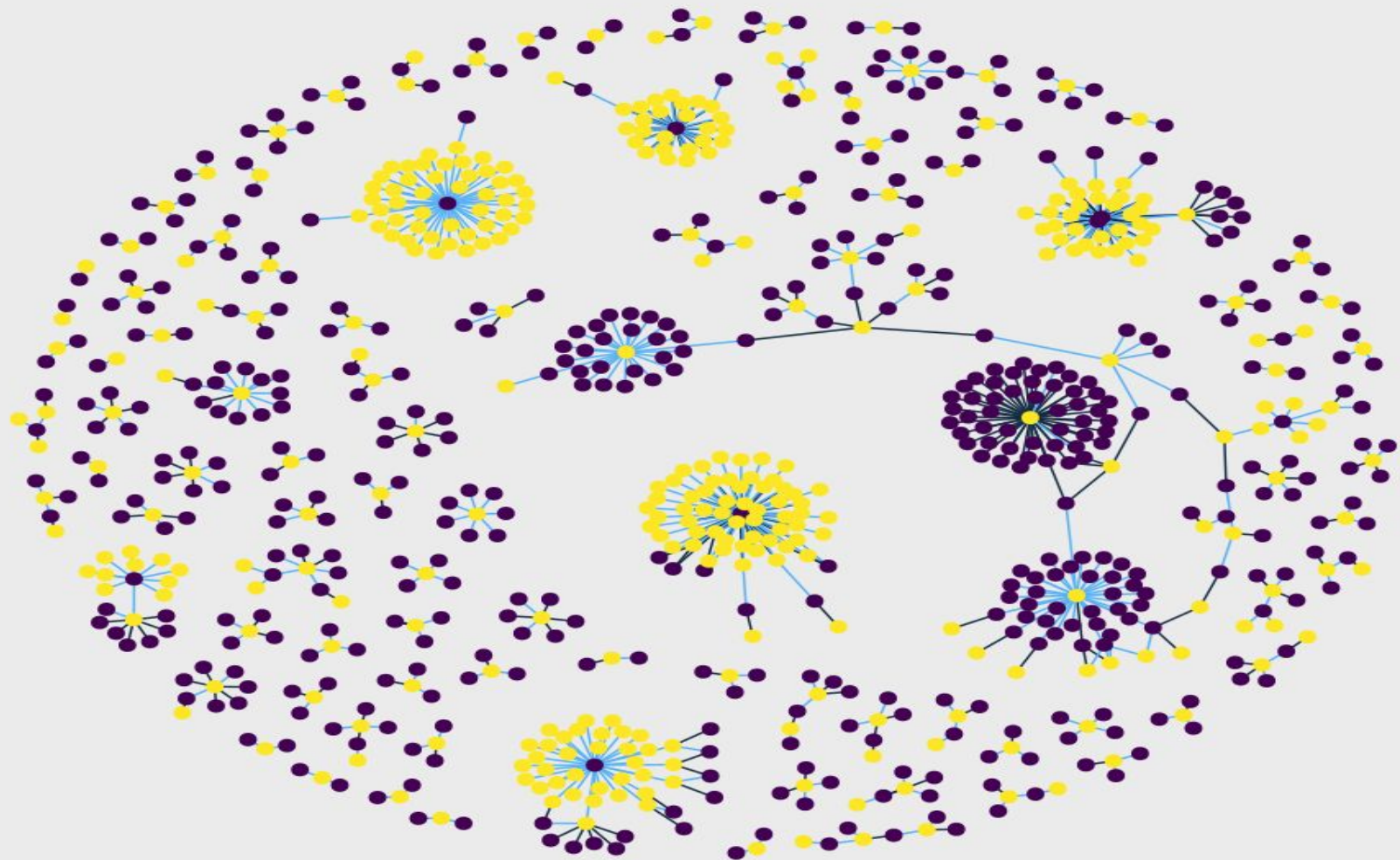
Addresses can function both as input or output to a transaction. Yellow nodes are transactions, and purple nodes are addresses. There are two different kind of edges. From address to transaction (input, type 0) and from transaction to address (output, type 1).

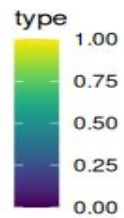
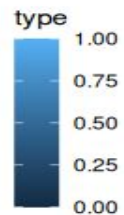
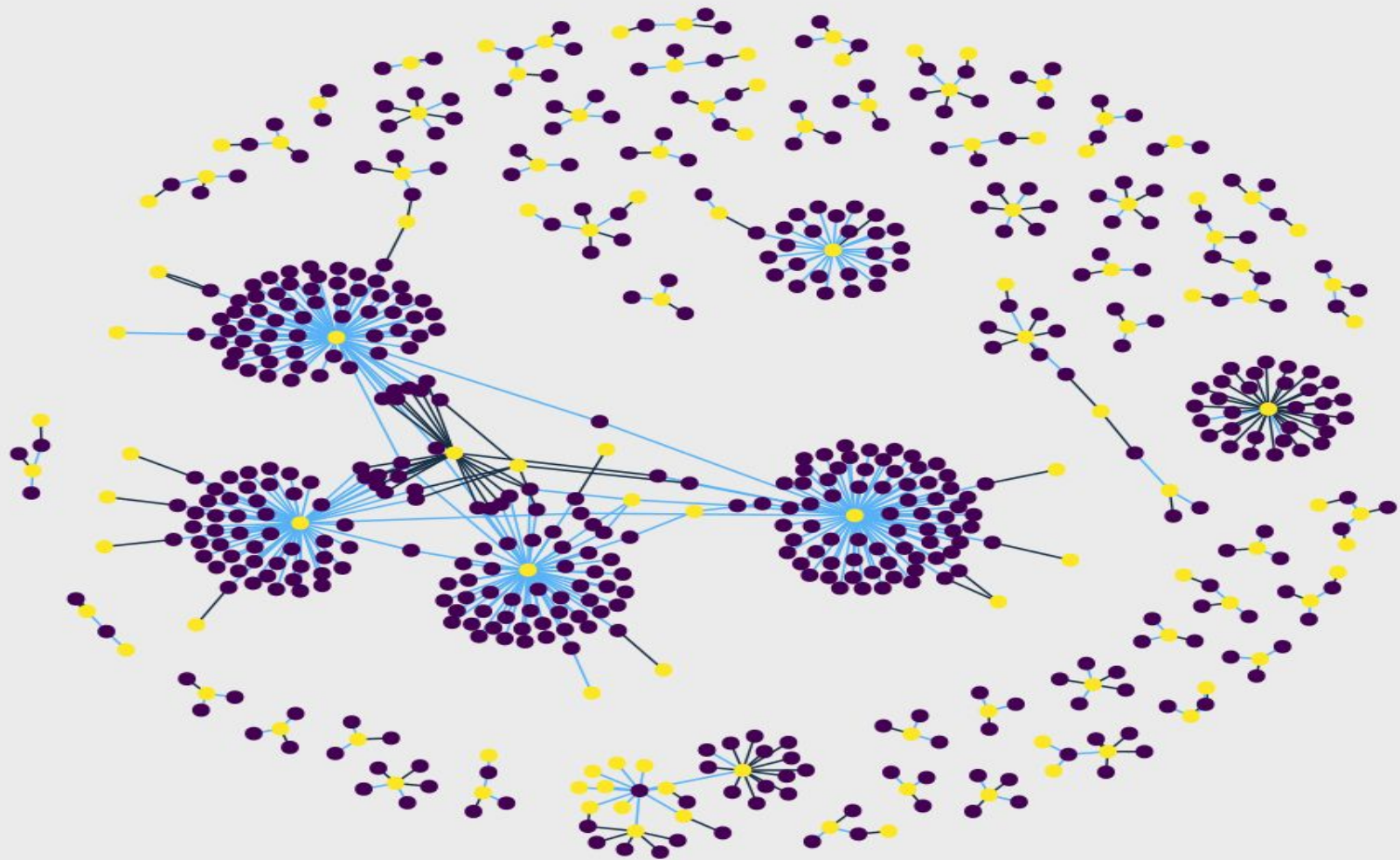
A common subset of transactions in a block might look like this

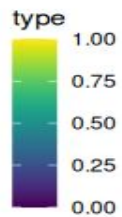
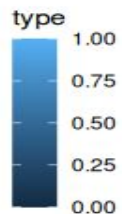
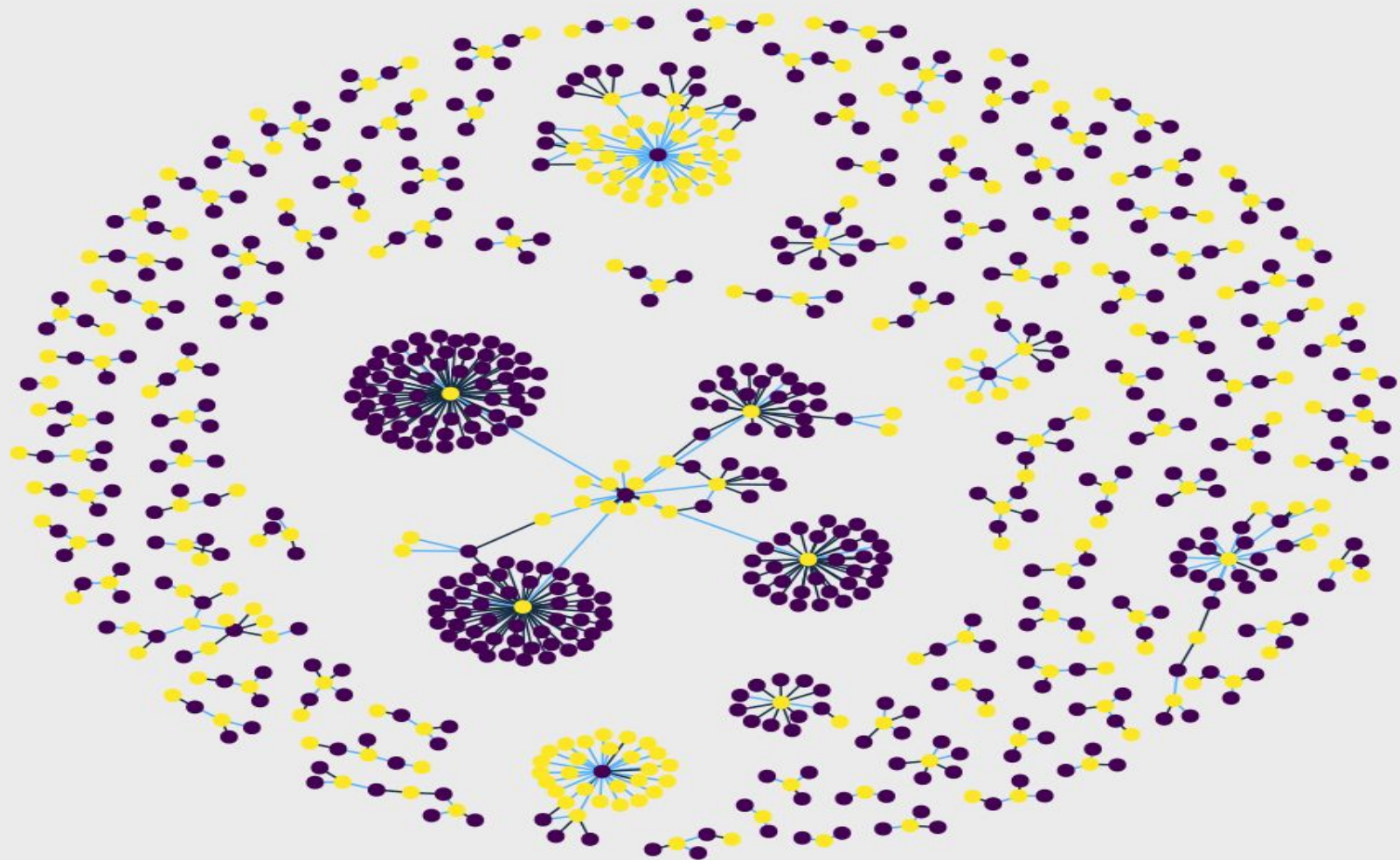


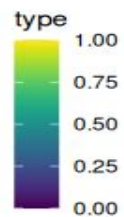
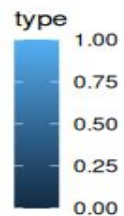
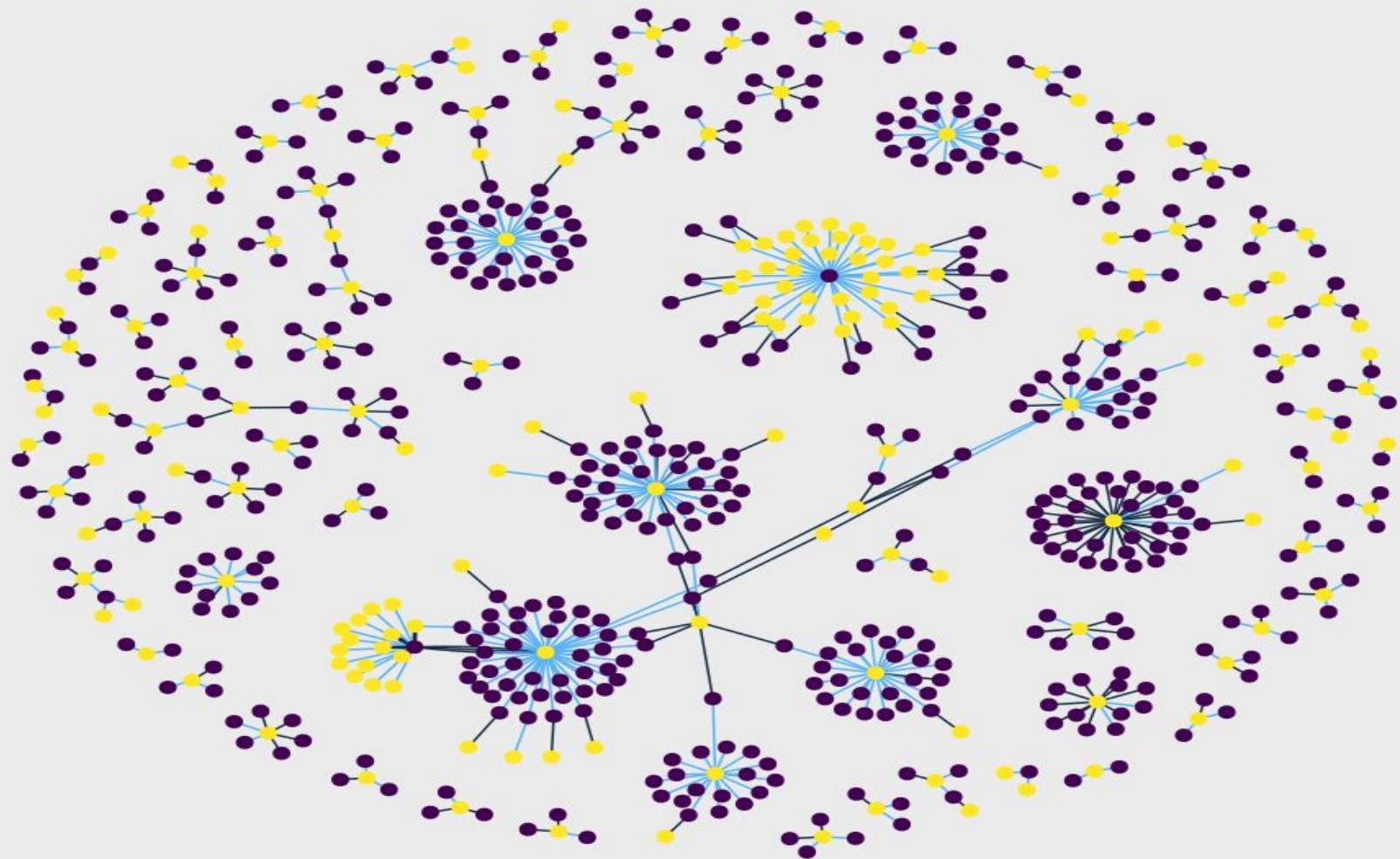
What we are looking for are large connected components inside a graph of a full block (or possibly more blocks). Once we try to plot more than 500 transactions, visualization becomes hard (and also meaningless), thus the approach i took was to extract a subset (or a slice) of unique nodes from the edgelist, and then extract again from the whole graph all edges containing at least one of the unique nodes from our “slice”. In this way we are able to determine if one node is involved in more than one transaction, or is part of a larger connected component of other nodes.

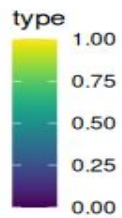
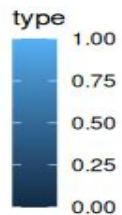
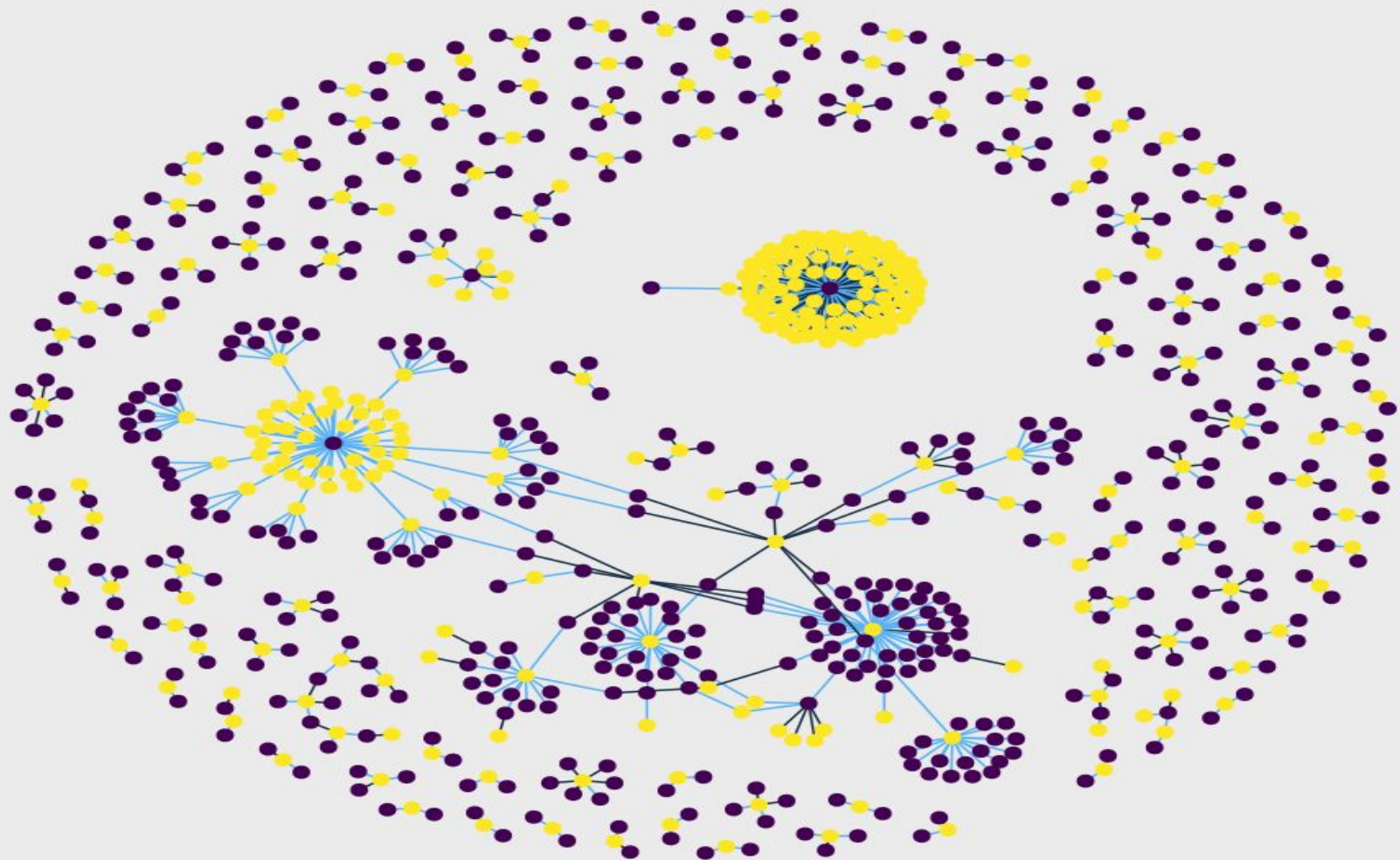
Here are some of the results achieved with this method.

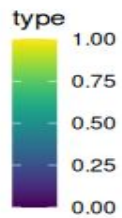
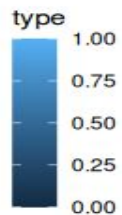
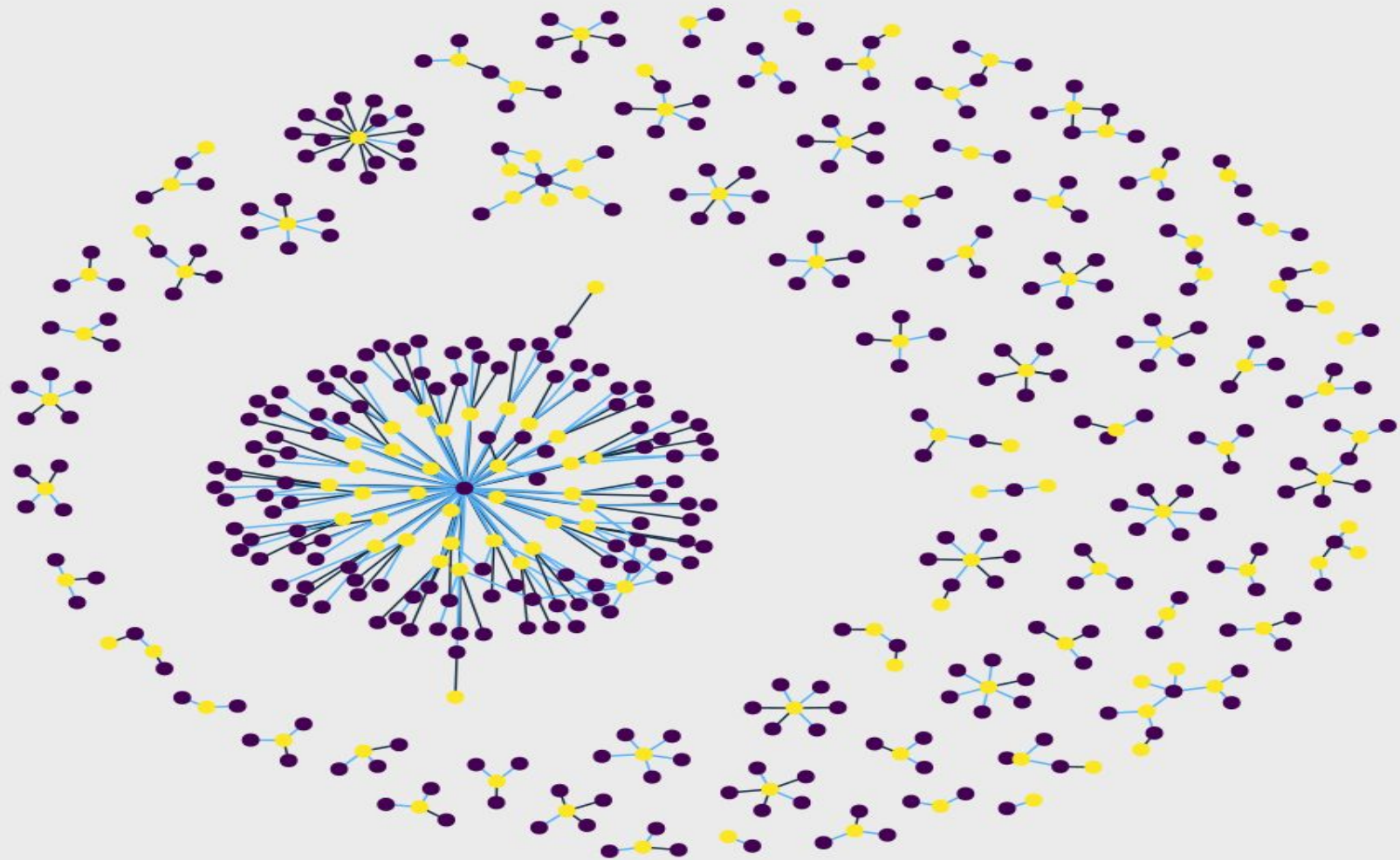


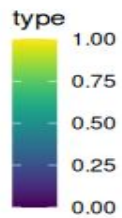
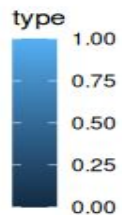
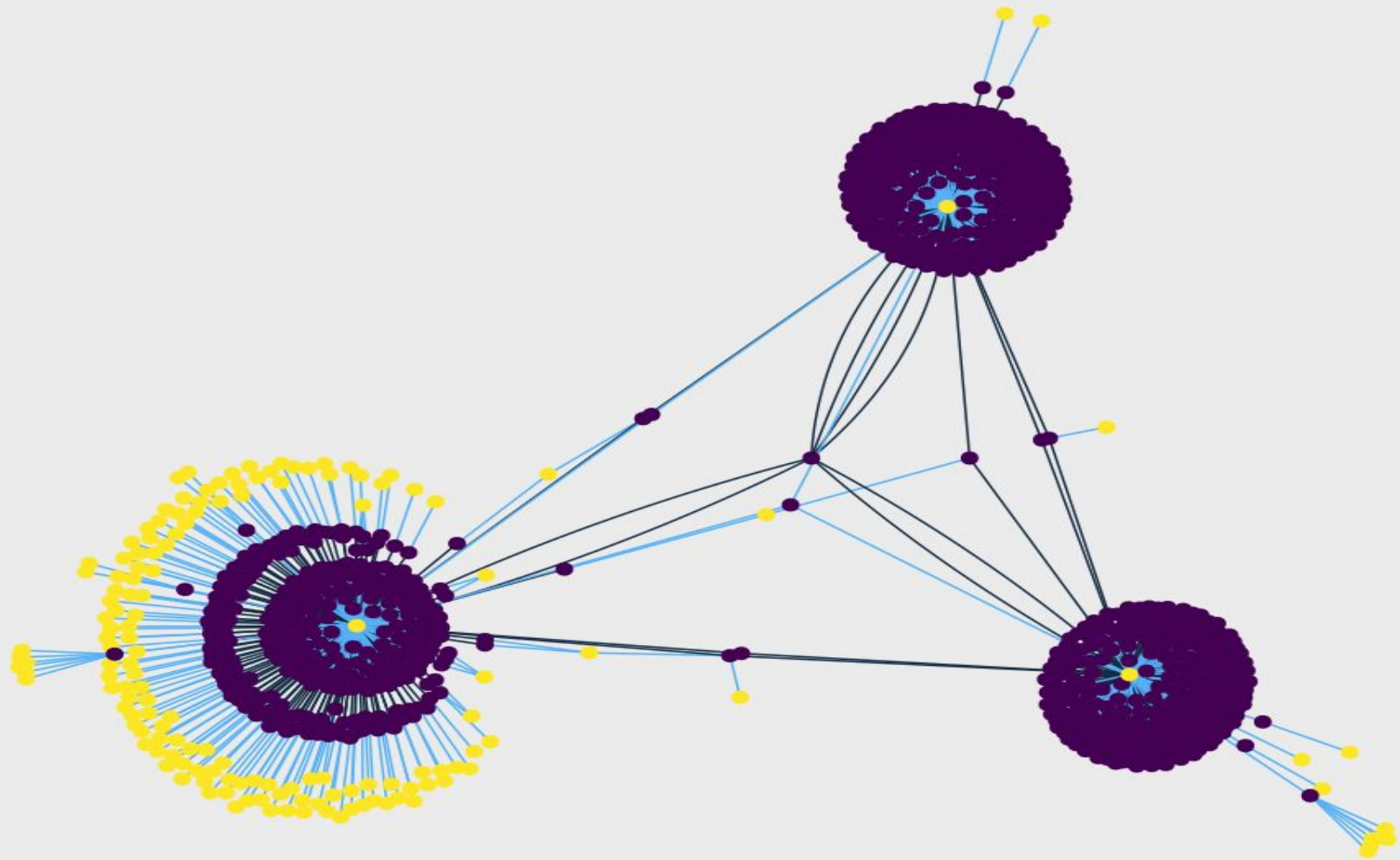












These results come from analyzing and plotting just one single block. It is clear that while the majority of transactions follow a simple pattern of one address sending to two addresses (the receiver and possibly another address owned by the same sender), sometimes we are able to identify clusters of addresses being part of the same transaction. This is crucial information, especially in the input side of the transaction.

If a cluster of addresses all function as input to one transaction, it means that either these addresses belong to the same entity (remember that each wallet can generate a seemingly infinite amount of addresses), or the sender has access to the private keys of the addresses involved in the sending side of the transaction. This is a crucial security issue in the so called “anonymity” of Bitcoin.

Largest graph

Upon realizing that our pipeline worked, i focused on trying to analyze a larger graph. The data i acquired consists of (almost) all the transactions in blocks between 500503 (07:36:28 AM 22-12-2017) and 500599 (11:10:12 PM 22-12-2017). This is almost 16 hours, in which around 100 blocks were mined. The network generated from these blocks contains over 1 Million edges.

Some calculations on this network are meaningless to try (such as diameter), but some of them are still doable.

Degree Distribution of large network

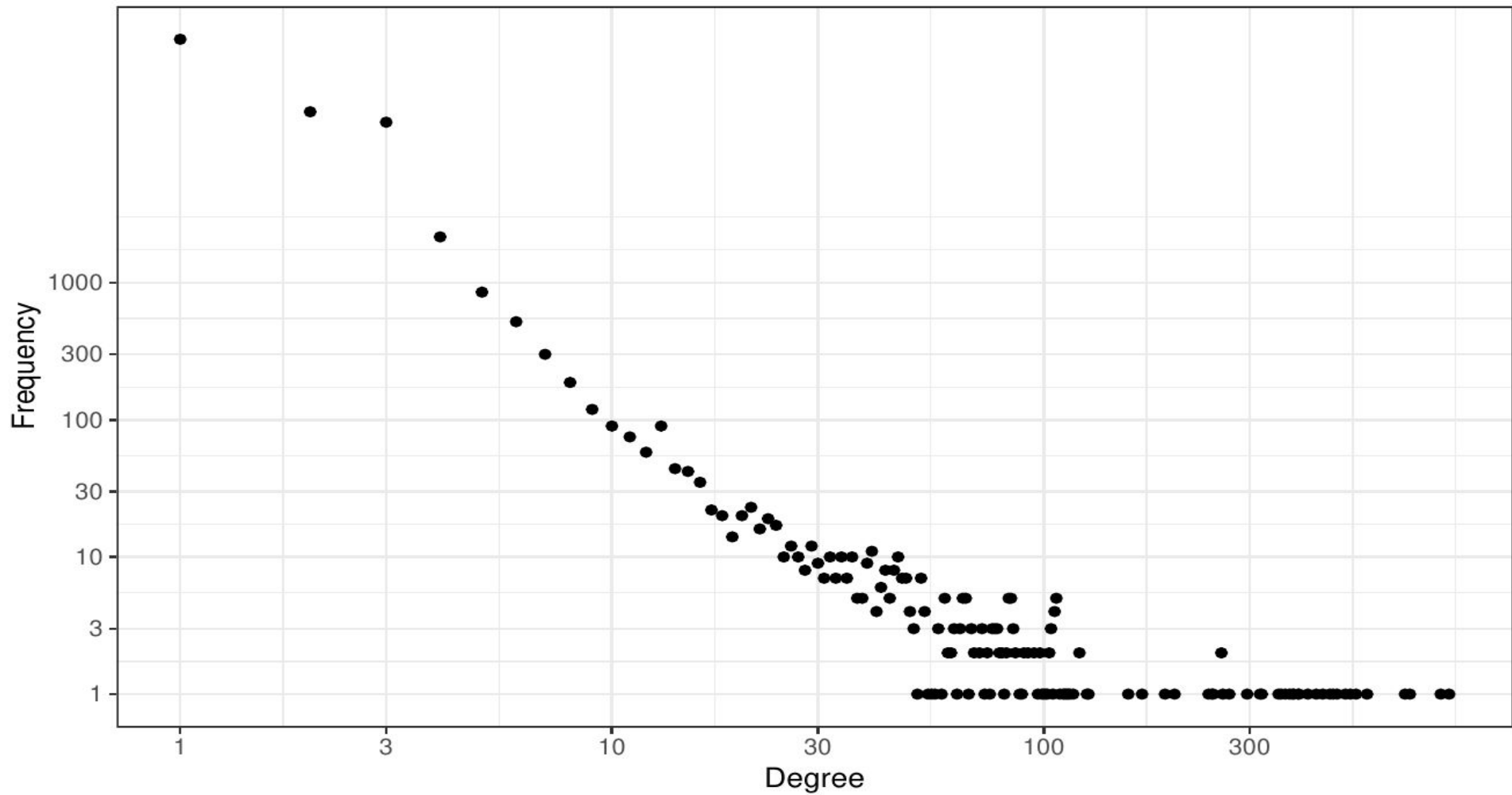
The log-log degree distribution of samples of 100'000, 500'000, and the full network all suggest power-law properties for this network. Most nodes have a low degree, while the number of nodes with high degree (hubs) decreases exponentially as the degree increases. Because of the size of the graph, for analysis we turn our attention to the giant component

Giant component in this graph contains 704k nodes, 856k edges, with max degree equal to 4908.

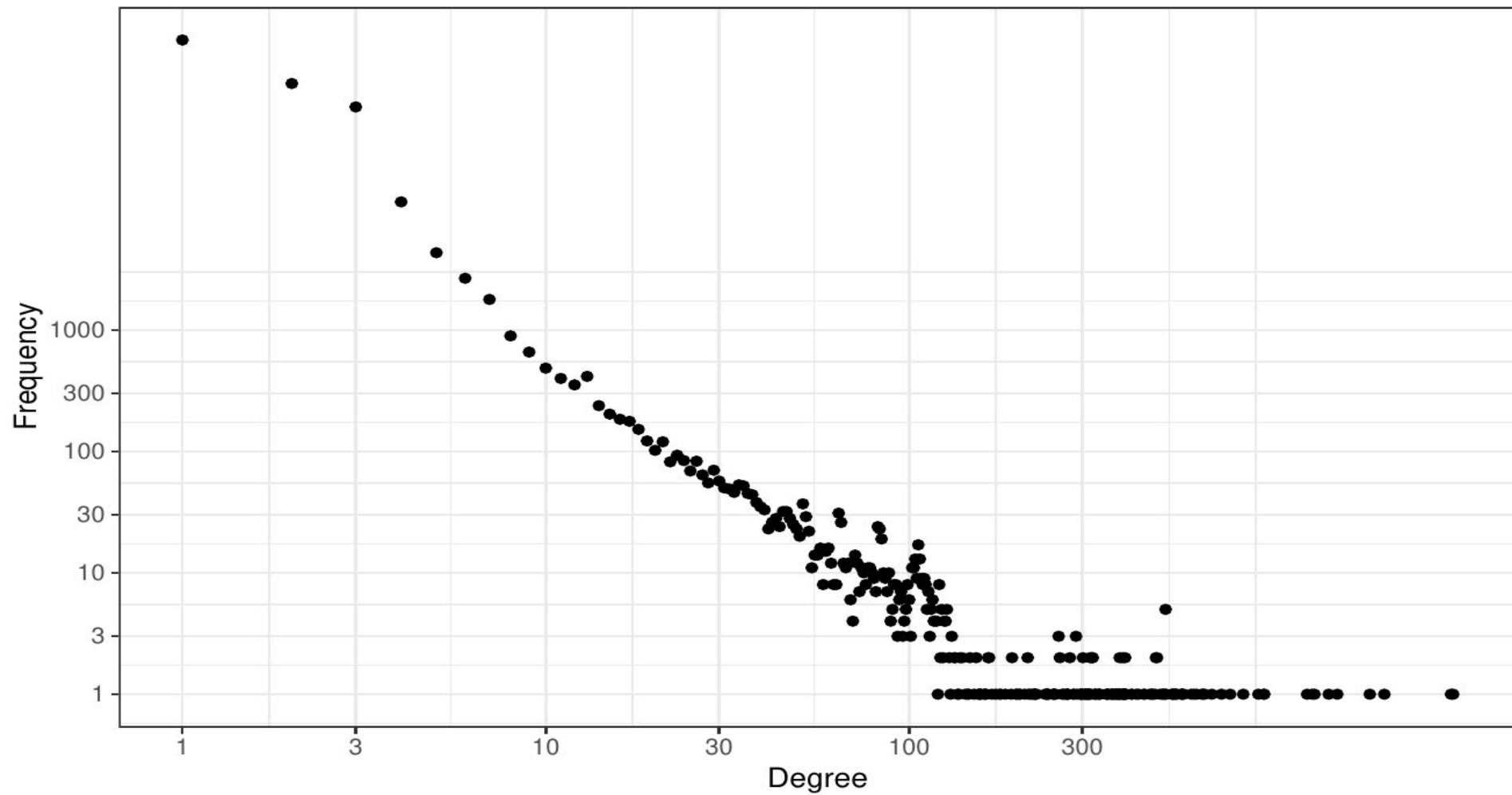

```
> summary(G.giant_component)
IGRAPH 79b06c2 DN-B 704023 856455 --
+ attr: name (v/c), type (v/n), amount (e/n), type (e/n)
> max(degree(G.giant_component))
[1] 4908
>
> mean(degree(G.giant_component))
[1] 2.433031
Warning message:
In degree(G.giant_component) :
  At vendor/cigraph/src/cliques/maximal_cliques_template.h:219 : Edge directions
are ignored for maximal clique calculation.

> vcount(G.giant_component)
[1] 704023
> ecoun(G.giant_component)
[1] 856455
```

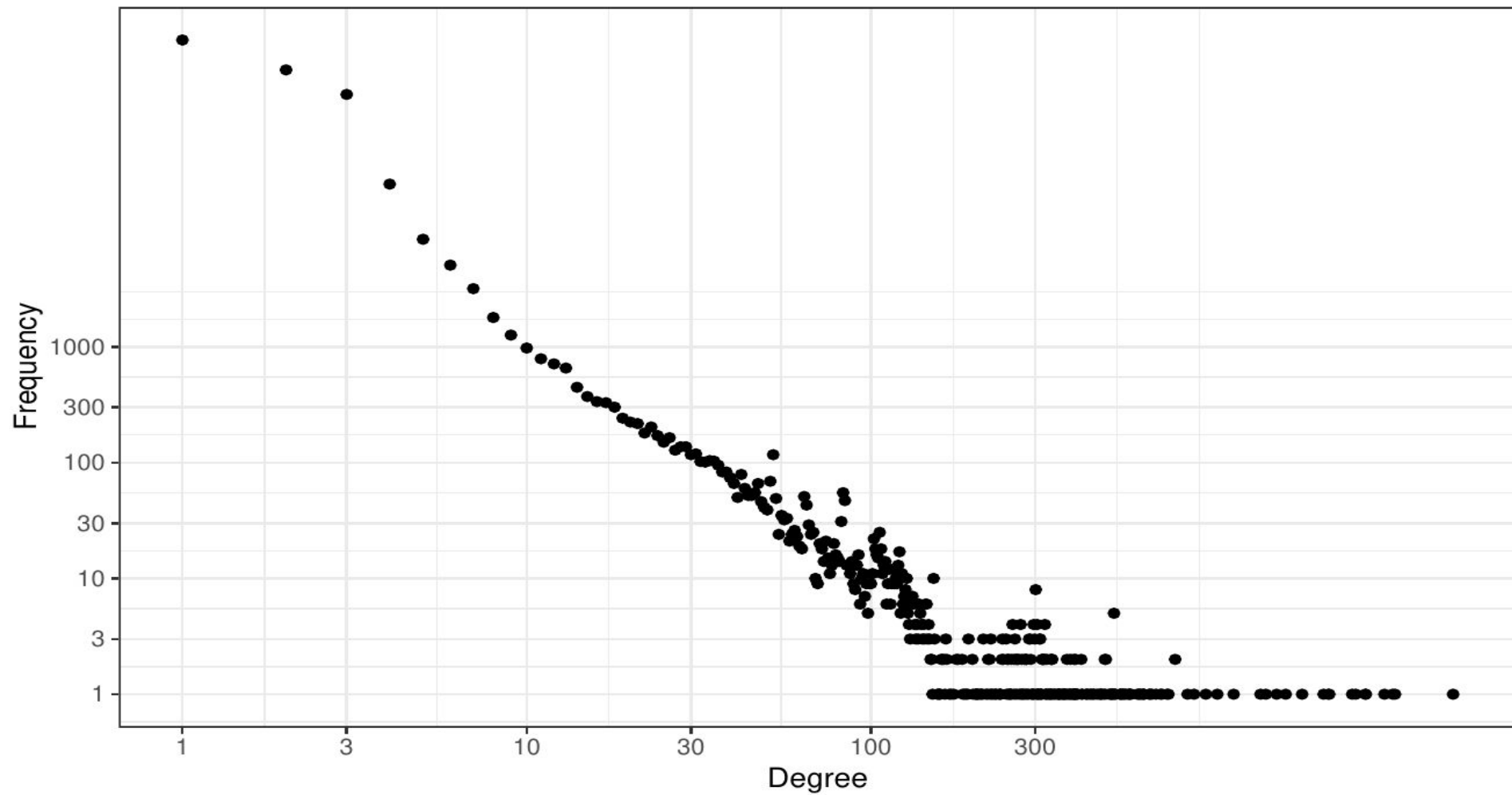
100k nodes log-log Degree Distribution



500k nodes log-log Degree Distribution



100 blocks log-log Degree Distribution



Conclusions

Although network analysis tools aid us in the description of this network, tracking ownership of addresses becomes easier with this kind of approach

References

- Ron, D., & Shamir, A (2013). Quantitative Analysis of the Full Bitcoin Transaction Graph
- Network Science (Albert-László Barabási)
- Programming Bitcoin Learn How to Program Bitcoin from Scratch (Jimmy Song)
- Mastering Bitcoin Programming the Open Blockchain (Andreas M. Antonopoulos)
- https://www.lambertleong.com/projects/bitcoin_network_analysis

Code

Full project code at <https://github.com/ricvigi/DMAU2-public>