

MidTerm Peer2Peer Systems and blockchains

Analysing the Kademia DHT

Lorenzo Bellomo, 531423

Project Usage, Parameters and Attached Files

The project can be executed both using command line arguments, or through the interactive shell. The three parameters to be passed for the execution are m , n and k , which are respectively the system wide parameters modelling the number of bits of the identifier space, the number of joining nodes in the network and the size of the routing table buckets.

The only fixed system wide parameter is α , which is fixed as 3 and is used in the Kademlia lookup procedure.

Together with this report and the source files (6 *.java* files, inside the *src* folder), a folder (*networkData*) with a set of files will be attached:

- *programAnalysis.xlsx*: it is an *Excel* spreadsheet with the parameters of used in the section , which analyses the statistics computed by the program. It contains the output of 5 runs for 3 different network topologies.
- *network1Dump.txt* and *network2Dump.txt*: Those files contain the full output of the runs analysed in depth in section , and using *Cytoscape* to analyse the result.

Class Overview

The project comes organized in 6 classes

- *Coordinator*: This class models the *Coordinator* entity (which has a global view of the network). It keeps a collection of *NodeDescriptor* instances (one for each node joining the network), together with some statistics to be displayed at the end of the execution.

The Coordinator is the only entity which is spawned by the Main class, and it provides as only public method the create network. This method, which is the one invoked by the main, calls two private methods: *initialize* (which resembles the *2a* phase in the MidTerm text) and *generateNewNode* (which instead models the points *2b* and *2b*).

The two other main functions of the Coordinator are:

- At the end of the creation of the network some *statistics* are computed and displayed through the command line
 - It acts as a *RPC registry* in the Kademlia network. In this way nodes that want to execute a *remote procedure call* on another node, have to ask the instance to the coordinator.
- *Main*: This class is the one with the main function, which simply asks for the three system wide parameters and spawns the Coordinator life cycle.
 - *Node*: This class models the node entity, as a triple of *IP address* (modelled as a *String* in order to avoid exceptions relative to non-existing IP

addresses), *identifier* (a *BigInteger* instance) and *UDP port* (which is not used in this simulation).

- *RoutingTable*: This class models a routing table of the Kademlia network. It consists of a kind-of Map structure, where each node keeps m bucket, each one of size k. This class provides methods to update the routing table, finding the best entries currently present in it. The procedure adopted for routing table update is the one proposed by Kademlia, in which the least recently seen node is pinged. Nevertheless, the network to be created is a *churn-free* one, which means that the ping method is a fake one, that always acts as if the nodes were alive.
- *NodeDescriptor*: This is the class that models the view that the Coordinator has of each node. This class basically models an actual node in the network, as a couple <Node, RoutingTable>. This class is the one over which the *Remote Procedure Calls* are instantiated. The main methods that are provided by this class are:
 - *findNode*: This method is the one exposed in the Kademlia original paper. The only main difference resides in the fact that the parameter passed to this method, in addition to the id, is the list of travelled nodes (needed in order to allow a new node to become known in the network). It was chosen to pass this list as a parameter in order to avoid the burden of modelling RPC messages in this simulation.
 - *nodeLookup*: This is the Kademlia lookup procedure. It finds in the network the closest k nodes to the target id passed as parameter. It doesn't explicitly exploit the bootstrap node (the one the Coordinator generates for him), but it just uses the best known nodes to get to fill the routing table k-buckets.
- *Utils*: This class simply provides some utility function (mainly to be used by the class *NodeDescriptor*) in the form of public static methods. This is the class that provides a function to generate a new node, and which is exploited by the *Coordinator* in order to populate the network.

It also provides a method for generating an identifier fitting one specific bucket of the routing table. This method (*generateIDInRightBucket*) is exploited to generate identifiers for each bucket of a node in order to initiate the lookup procedure on this specific id.

Data Structures

The main data structure to be presented in this report is the one used in the routing tables. It has to model 160 buckets, each one of size at most k. The final choice was to use an ArrayList of 160 positions, and for each entry, using a Queue (a LinkedList) of Node classes. In order to access the right bucket out of the 160, whenever an update to the routing table has to be issued, then a

findBucketIndex(id). This method computes the bucket index as $\log_2 d$, where $d = myId \oplus targetId$.

The only other data structure worth mentioning is the *HashMap* kept by the coordinator and containing all the *NodeDescriptor* instances.

Stats Display

At the end of the creation of the whole network, some statistics are displayed through the command line. Their purpose is mainly to show what happens with different network wide parameters, and also compute an expected value for the number of edges.

The expected value of edges is computed as follows:

- A number of random identifiers is generated (i.e. 20), and for each of them, the number of nodes “falling” in every *k*-bucket is computed.
- At the end of this process, an array of *m* position is filled, where in position *i* the average number of nodes falling in this *k*-bucket for a generic node is stored. This way I can estimate the number of edges for each bucket *i* in this way:
 - *if the number is greater than k* then this bucket is probably getting filled, so the process adds *k* to the total number of edges
 - *otherwise* this bucket has a lot less edges to offer to the whole network, so the process adds the value computer averaging on the sample generated in advance.
- Throughout this process, by summing the temporary values, it is possible to compute the average number of outgoing edges from a generic routing table.
- The total number of estimated outgoing edges is computed by simply multiplying this value with the total number of nodes.

This value should be considered valid with reasonable certainty because the Kademlia process is, by construction, expected to converge with a few number of steps, and get to the *best solution* in a few steps. In order to validate this hypothesis, another stats comes helping: the *recursive lookup depth*. For each lookup call, the recursive depth of given call is computed by counting the times that the loop in the lookup process is executed. At the end of the construction of the network, for each recursive depth, the number of times it was reached is displayed. This value, in order to respect the *fast convergence hypothesis* of the Kademlia protocol, should be not only low on average, but with also relatively low variance. Other stats that are displayed at the end of the network creations are:

- The *time* needed to generate the network.

- The *expected number of nodes* for each distance (bucket) for a generic node.
- The number of registered *collisions* for each hash function.
- The value *maxTheoretical*, which is the number of edges a theoretical network with every k-bucket of every routing table being filled would have. This value has to be considered an upper bound to the number of total edges in the network, as not enough nodes for each bucket can be present in the network, and it is only displayed for comparison purposes. This value is computed by summing up k for each bucket, except for the first ones (for example bucket 1 accounts for 1, bucket 2 accounts for 2...).
- The top 10 and bottom 10 nodes according to the in-degree, together with their order of generation. This value will be useful in Section *Tool Analysis* when analysing the in-degree and which nodes have an higher and lower value for this stats.

Node Join

A node join is organized as follows:

1. the Coordinator entity creates a new node with a random IP address. The identifier is then built by using the *SHA1* function on given IP address, and computing this value modulo 2^m . The identifier value is stored using the Java class *BigInteger*, which belongs to the *java.crypto* package. This class provides a number of methods for special operations of common use in cryptography, like XOR(\oplus) or modulo(%), which are useful also for the Kademlia protocol. The *BigInteger* class also provides also allows the use of long identifiers, without bounds to the length (for example 2^{64} overflows for longs in Java). Nevertheless generating a new node does not prevent collisions like *SHA1* does, since the identifier space is shrank to $[0, 2^m)$. The Coordinator checks if the newly generated id is already in the network, and in case of a collision, it simply increments the collision counter and generates a new node
2. the Coordinator finds a bootstrap node for the newly generated node, so that the new node has an outgoing edge towards the network, useful for populating its routing table. The new node then issues one lookup procedure for its own identifier,
3. after the first lookup, the new node issues m lookups, one for each *k-bucket*.
4. at this point the node join procedure is finally complete. Every node which received a findNode call originated in the new node, had a chance to know its new identifier, and, if possible adds this new contact to its routing table. This only happens “if possible” since the algorithm applied to update the routing table resembles the original one, which by default

keeps the “stablest” nodes, so the ones added to the routing table in the farthest moment in time. The last nodes have a considerably lower chance of being added, in this *churn free* environment, to any routing table. In conclusion, a higher *in-degree* should be expected for the first nodes that were generated, with respect to an actual working Kademlia network.

Topology Analysis

The following section addresses the network topology analysis part of the assignment. It will be organised in two parts:

- *Program Analysis*: a first, more lightweight part, in which only the program generated parameters are analysed (so without the use of external tools for graph analysis). This part is useful in order to ensure the correctness of the procedure and to see some small details. The analysis is in section Program Analysis.
- *Tool Analysis*: a second part, where a small number of networks are analysed in depth. This part is in section Tool Analysis.

Program Analysis

The main point of interest for this part is the correctness of the implemented Kademlia procedure. Regarding that, two important focus point have to be analysed: *number of edges* and *recursive depth*. The first one accounts for the correctness of the joining procedure, together with the fact that “old” nodes get a chance to know new nodes. The second point instead focuses on the *fast convergence* of the algorithm. Of course those values have to be coherent in different networks (i.e. filled ones, and almost empty ones).

All the following analysis is the result of the average of 5 runs for each network. The single values of those runs, and their averages, are present in the excel file *programAnalysis.xlsx*, added as an attachment to the project.

There follows a naive analysis of some networks:

1. $\langle m, n, k \rangle = \langle 10, 1000, 10 \rangle$: This network is characterized by an almost full identifier space (1000 nodes in a 1024 range identifier space). Something to expect from such a filled network is a high number of collisions in the node generation phase, and a really high number of edges, according to the fact that a lot of buckets (also those ones relative to close neighbours) should be filled for the majority of nodes.

The average values are:

- The upper bound for the number of edges is 75 000.
- the number of edges is 74 650, while the expected theoretical value is 75 670. The percentage difference (so the variation) between those values is around 1,50%. It is possible to notice that the expected

value is actually higher than the maximum possible theoretical value. This should not be a problem, since the estimate is probabilistic based. The important thing is that the actual number of edges stays below the threshold.

- The average time to generate the network is 6,58 sec.
- The average recursive lookup depth is 1,70 (so most of the times, the loop in the lookup is executed 1 or 2 times), while the maximum depth reached is 5.
- The number of collisions with the hash function is 2 835.

The two highlights from this process are the values of average recursive depth, which is very low (so the procedure actually has a fast convergence, when not considering churn), and the number of edges, which comes really close to the expected one. The filling factor of this network (where $fillingFactor = numEdges/upperBound$) has been computed to be 99,53%. An important thing about the filling factor has to be clarified: this value is not the network filling factor, with respect to $numberOfNodes/numberOfIdentifiers$, but it is the filling factor of the routing table, with respect to the upper bound computed.

This shows the fact that the simulation works well in a filled environment.

2. $\langle m, n, k \rangle = \langle 20, 500, 5 \rangle$: This network differs from the precedent because it has a significantly lower filling factor, as the range of identifier is $[0, 2^{20})$ with only 500 nodes inside. This network should have a significantly lower number of edges with respect to the upper bound.

The average values are:

- The upper bound for the number of edges is 46 000.
- the number of edges is 18 646, while the expected theoretical value is 18 925. The percentage difference (so the variation) between those values is around 1,41%.
- The average time to generate the network is 2,04 sec.
- The average recursive lookup depth is 2,11, while the maximum depth is 5.
- The number of collisions with the hash function is 0.

As expected, the number of collisions is dramatically lower (actually, 0 over 5 runs). The recursive depth is a little bit increased, but the real important piece of information relies around the number of edges. The upper bound is far from reached, (filling factor is 40,45%), but the theoretical estimate is matched within 1,41% accuracy.

3. $\langle m, n, k \rangle = \langle 64, 10\,000, 20 \rangle$: This one is an example of a huge network. This is used to assert the goodness of the implemented procedure even in big networks.

The average values are:

- The upper bound for the number of edges is 12 110 000.
- the number of edges is 1 975 067, while the expected theoretical value is 1 989 000. The percentage difference (so the variation) between those values is around 0,71%.
- The average time to generate the network is around 30 min.
- The average recursive lookup depth is 2,38, while the maximum depth is 6.
- The number of collisions with the hash function is 0.

The purpose of this simulation is to show that, even for a huge network, the recursive depth reached by the algorithm is very low. As expected, the lookup depth did not grow almost at all, and the number of edges respects the theoretical assumptions.

Tool Analysis

The tool used for graph analysis is *Cytoscape*.

The two Networks that will be analysed are the following:

1. $\langle m, n, k \rangle = \langle 15, 1\,000, 10 \rangle$: Network1
2. $\langle m, n, k \rangle = \langle 10, 200, 1 \rangle$: Network2

Network1

Network Parameters: $\langle m, n, k \rangle = \langle 15, 1\,000, 10 \rangle$

Now follows an analysis of the various parameters of the network

Degree The main point of focus is the one regarding the difference between *in-degree* and *out-degree* (Figure 1 and 2).

The out degree behaves in a very regular way. This value ranges from a minimum of 67 outgoing edges to a maximum of 82, with no significant variance. This is expected, as at the end of the process, generically every node should have converged to his best possible solution. This means that, all the last k-buckets should be filled, while the close ones should contain as many close contacts as possible (though usually not quite enough to fill a k-bucket of 10 elements).

The same regularity does not apply to the in-degree, for which values range from around 20 to over 800. Why is that? A possible reason for this variance lies in the protocol, which, in case of a filled bucket, tends to keep older nodes (considered more stable). This gets only amplified in a *churn free* environment like the one in the simulation.

To verify this fact, the output of the network creation comes useful (provided as attachment with the file *network1Dump.txt*). As it is possible to see in Figure 3, the 10 nodes with highest in-degree all belong to the first 20 to be generated, while the lowest ones are all well over 900. In fact, the process that drives new nodes to gain knowledge of other nodes is strongly driven by the first k-buckets

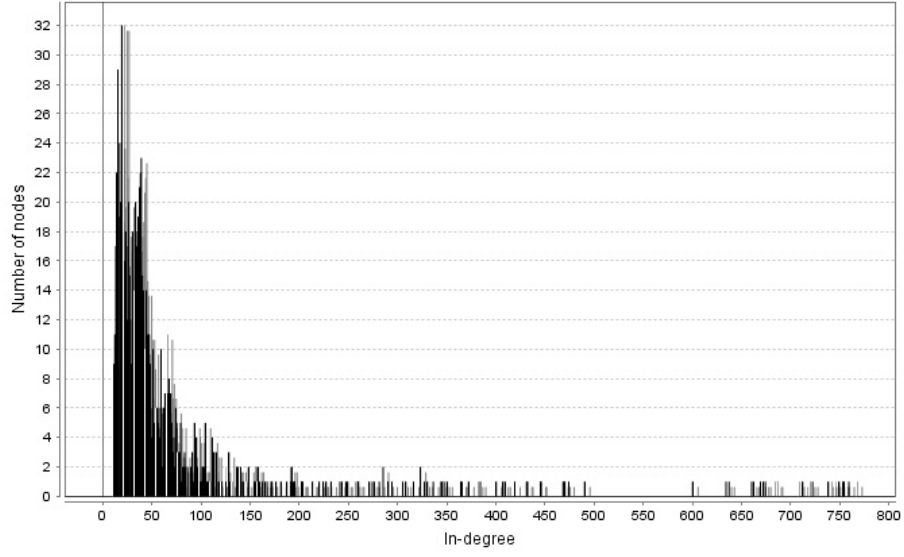


Figure 1: in-degree distribution, Network 1

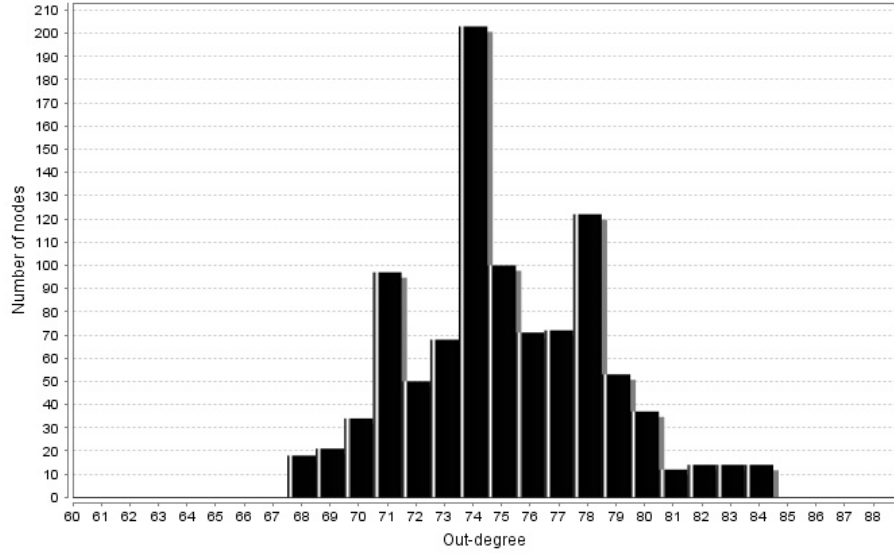


Figure 2: out-degree distribution, Network 1

that get filled in the network. The last nodes have a chance of entering in other nodes routing table only if their relative k-bucket is empty, which usually happens only with respect to close neighbours. In conclusion, the last nodes

```

Now printing 10 lowest in-degrees
ID: 28757 generated 979th and has inDegree 10
ID: 2176 generated 987th and has inDegree 10
ID: 403 generated 986th and has inDegree 10
ID: 5233 generated 982th and has inDegree 10
ID: 8344 generated 976th and has inDegree 11
ID: 14610 generated 921th and has inDegree 11
ID: 14640 generated 932th and has inDegree 11
ID: 29316 generated 966th and has inDegree 11
ID: 15140 generated 978th and has inDegree 11
ID: 17271 generated 991th and has inDegree 11
Now printing 10 highest in-degrees
ID: 32702 generated 3th and has inDegree 767
ID: 18827 generated 12th and has inDegree 759
ID: 26161 generated 2th and has inDegree 754
ID: 25089 generated 6th and has inDegree 753
ID: 7599 generated 17th and has inDegree 749
ID: 18676 generated 1th and has inDegree 745
ID: 4870 generated 9th and has inDegree 741
ID: 1115 generated 16th and has inDegree 738
ID: 12278 generated 15th and has inDegree 722
ID: 9555 generated 22th and has inDegree 719

```

Figure 3: The in-degree output of the 1st network, Network 1

getting generated only get into a few routing tables, and the ones of really close neighbours (close with respect to their XOR distance).

Network Diameter The network diameter is 4. It must be noticed, as it is possible to see from Figure 4, that only few nodes actually have a characteristic shortest path of 4, but the majority can reach any destination in 2 or 3 hops.

Clustering Coefficient The clustering coefficient distribution is shown in Figure 5, while on average $CC = 0,396$.

Network2

Network Parameters: $\langle m, n, k \rangle = \langle 9, 500, 2 \rangle$

The network that follows is one which is smaller but significantly more full with respect to the identifier range. In Figure 6 a snapshot representing a part of the network is displayed. The degree is significantly lower with respect to Network1.

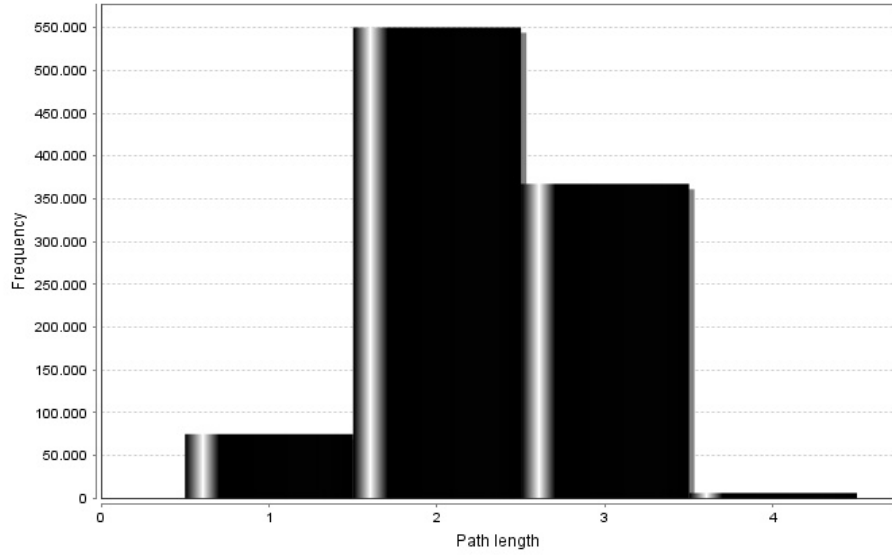


Figure 4: shortest path distribution, Network 1

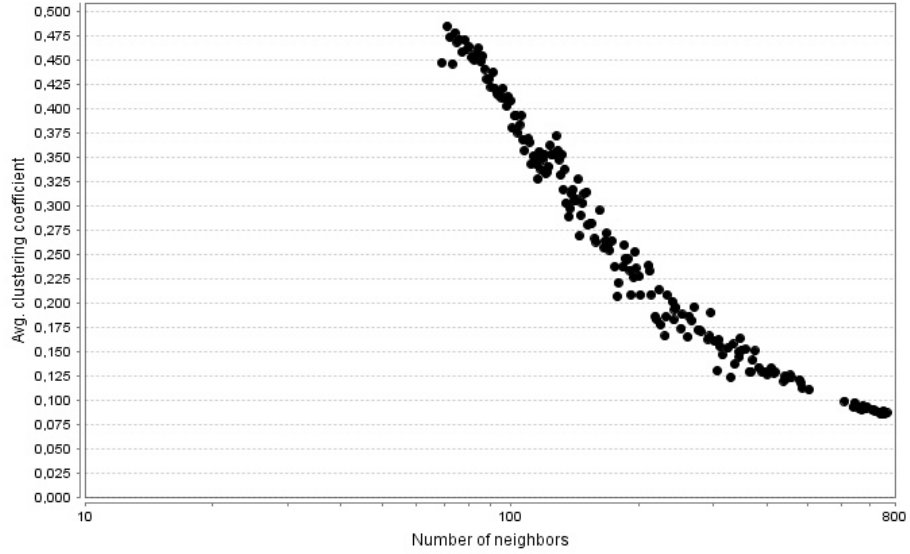


Figure 5: average clustering coefficient, Network 1

Degree Figure 7 and Figure 8 display the in-degree and out-degree of the network. The result obtained is coherent with respect to Network1, so the analysis will follow the same logical link.



Figure 6: snapshot of a portion of Network 2

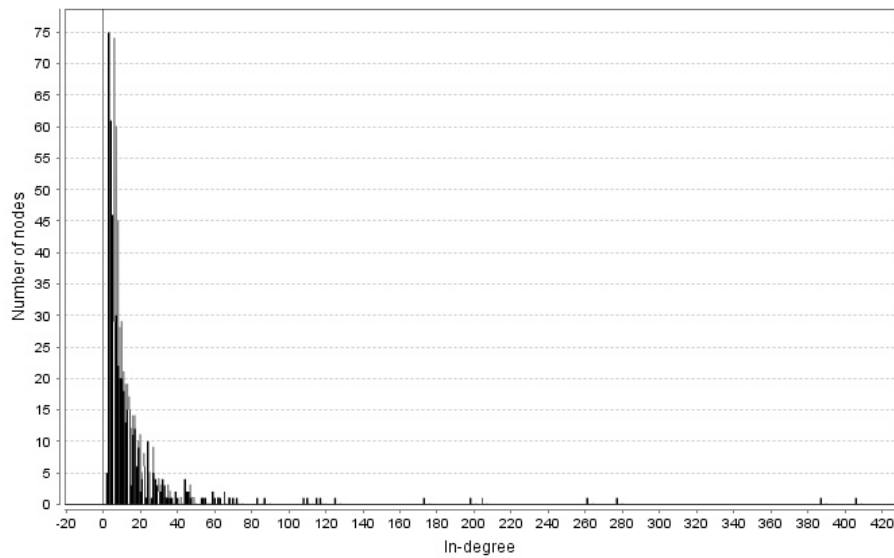


Figure 7: in-degree distribution, Network 2

Network Diameter The network diameter is 5, but the average path length is around 3 (as shown in Figure 9).

Clustering Coefficient The clustering coefficient distribution is shown in Figure 10, while the average clustering coefficient is $CC = 0,317$.

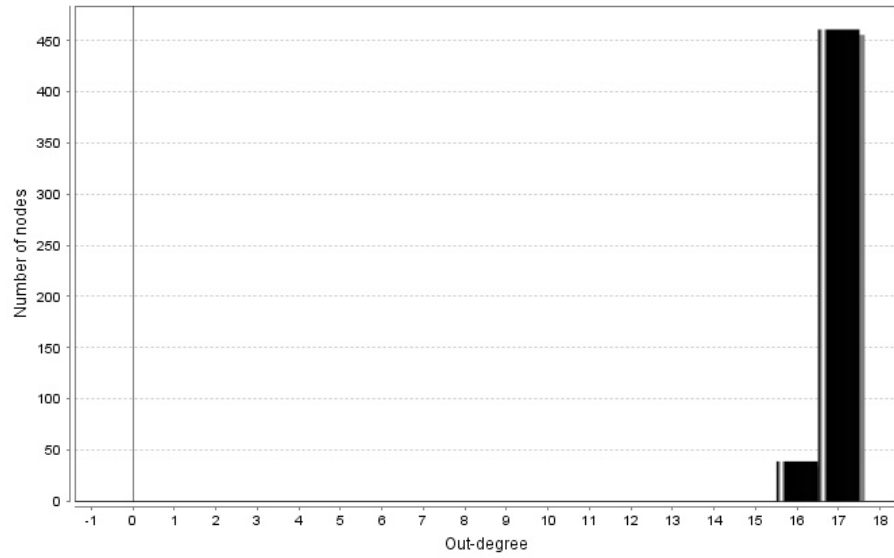


Figure 8: out-degree distribution, Network 2

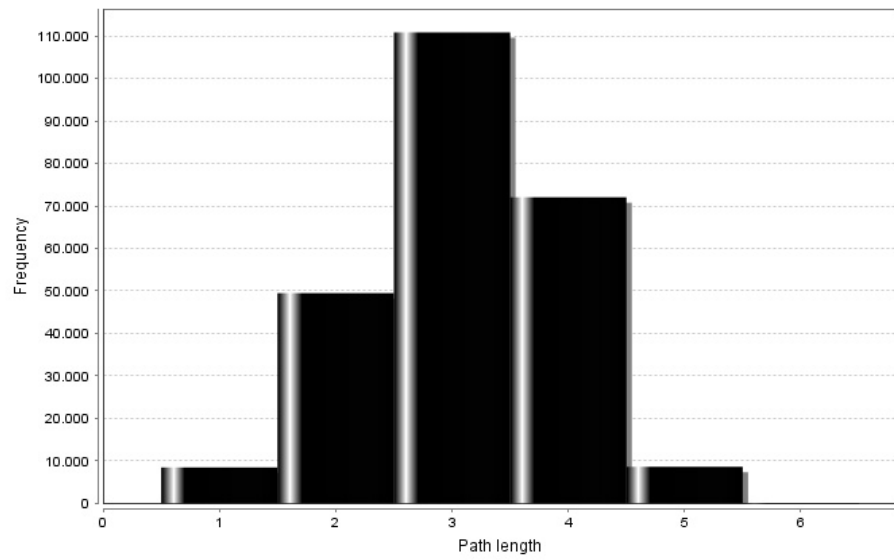


Figure 9: shortest path distribution, Network 2

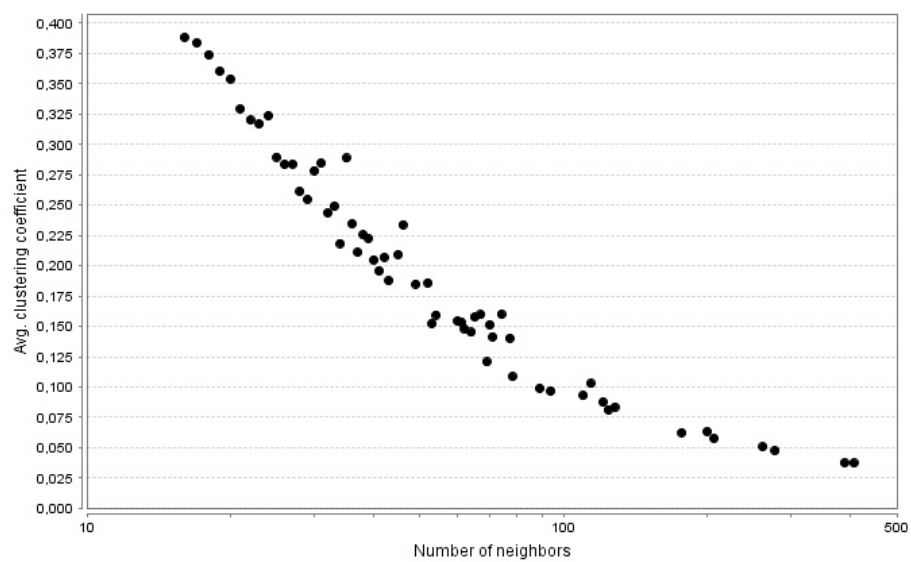


Figure 10: clustering coefficient distribution, Network 2