# Community Detection in Complex Networks using Genetic Algorithms

Simon Lehnerer[1]

**Abstract:** Detecting the community structure is of great interest when analyzing the topology of a network, however it is not a trivial problem. In this article a genetic algorithm is proposed which finds the community structure of a network based on the maximization of a quality function called modularity. Tests using several sample networks show that it reliably finds the community structure. However it does not resolve sufficiently small communities as intuitively expected due to an effect known as *resolution limit*.

**Keywords:** complex networks; community detection; genetic algorithm
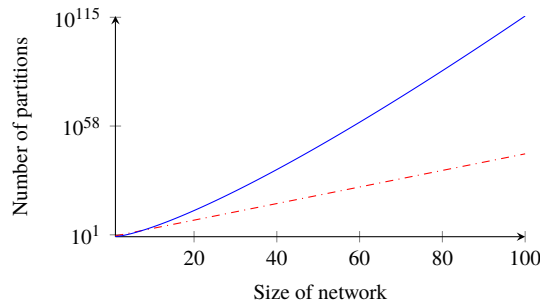
## 1 Introduction

Broadly speaking, the *community structure* of a network is the division of the network into groups of nodes called *communities*, which are internally densely connected and loosely connected to nodes from other communities. There is no unique formalized definition of the term *community* (*structure*) [Fo10]. In fact, many different approaches to define and detect communities in a network have been developed in the past [FH16; Fo10; Sc07]. We will give a few examples in the following: A simple method to find communities is to divide the nodes of a network in $g$ groups of predefined size such that the total number of edges connecting different groups in minimized [Po97]. For networks with a known hierarchical structure, e.g. often found in social networks, hierarchical clustering algorithms have been developed. Based on a distance measure between nodes they either start by considering the whole network as one single community and then iteratively divide the community into smaller communities (*top-down*), or they start by considering each node as a single community and iteratively merge them into larger communities (*bottom-up*) [HTF09]. Another popular ansatz is to determine communities using the eigenvalues of the Laplacian matrix of the network [Lu07]. Methods based on statistical inference like the *stochastic block model*, a generative model for random graphs, are also often used to study the community structure of a network [KN11]. Several information theoretic approaches try to uncover the community structure by investigating the trace of random walkers on a network and utilize the fact that the walkers are "trapped" inside communities [RB08]. Because the community problem is ill-defined there is no superior approach. It rather depends on the specific network that is

---

[1] ETH Zürich, Rämistrasse 101, 8092 Zürich, Switzerland, lehsimon@ethz.ch

to be studied and the notion of *community* that is used by a specific method. For certain networks some approaches might be more useful than others [FH16].

In this work we will focus on a quality function called *modularity* [FC12]. The modularity of a community partition is a measure that is the higher the "better" a partition indicates the community structure of the network. Thus the community structure is given by the community partition with maximal modularity.
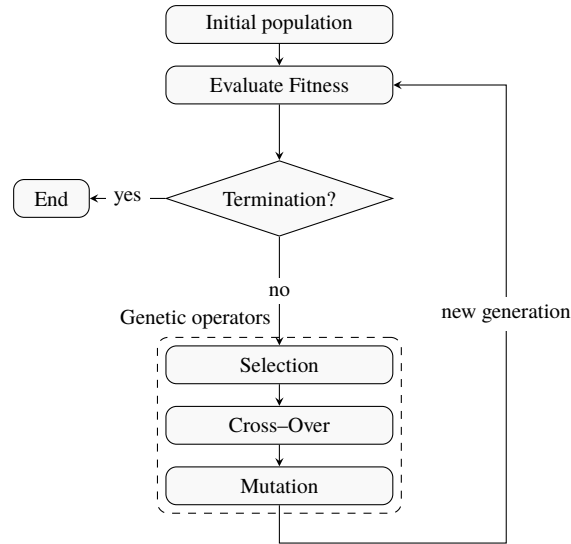


**Fig. 1:** Number of possible partitions (blue) compared with exponential curve (red). For a network consisting of twelve nodes there already exist more than one million possible community partitions.

There is no method known that can determine this community partition in a simple and straightforward way. That is why modularity based approaches are in fact optimization approaches that try to maximize modularity [FH16]. The number of possible community partitions of a network with $n$ nodes is given by the Bell number $B_n$, which has factorial growth (see fig. 1). A "brute-force" approach to find the community structure by comparing all possible partitions clearly becomes impractical for larger networks. In addition to that, the modularity function features a high degeneracy which makes it difficult to find its global maximum [Ba16]. For that reason we need an algorithm that can find the community partition with maximal modularity both efficiently and effectively. For this purpose a genetic algorithm (see fig. 2) is proposed in the following.
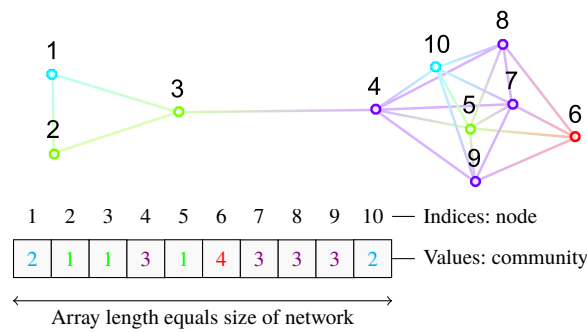
## 2   The algorithm

The algorithm [Le18] described here is aimed at finding the community structure of an undirected, unweighted network. A *chromosome* represents a possible community partition of the network. Technically it is a list where the $i$-th entry contains the community label $c_i$ of node $i$ (see fig. 3). Thus two nodes $i$ and $j$ are in the same community iff $c_i = c_j$. The chromosomes of the initial population are generated randomly.

**Fig. 2:** Schematic of a genetic algorithm

Genetic algorithms, inspired by processes of evolution, are stochastic optimization methods used to find an optimal solution in a large solution space. They use a *population* of *individuals* where each *chromosome* of an individual represents a possible solution. A *fitness* value reflects the quality of the solution of an individual and is used to stochastically perform *selection*, *crossover* and *mutation* operations to form a new generation of individuals (also called *children*). Ideally, the quality of the solutions improves with every new generation. The process is repeated until a reasonable solution is found.
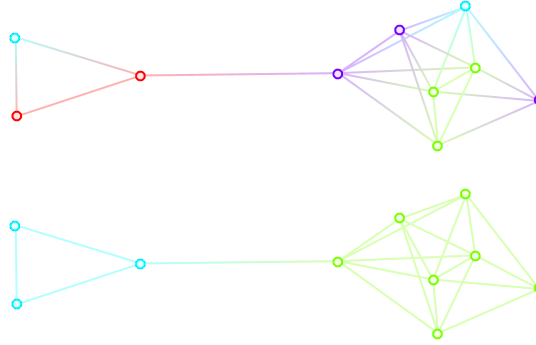


**Fig. 3:** Example of a chromosome for a network of size 10. It represents a possible community partition of the network (the communities are color-coded so that two nodes are in the same community iff their color is the same).

The fitness is calculated from the modularity of the community partition, which is a global measure of the network and given by

$$Q = \frac{1}{2m} \sum_{i,j} \left( A_{ij} - \frac{d_i d_j}{2m} \right) \delta(c_i, c_j) \quad \in [-0.5, 1]$$
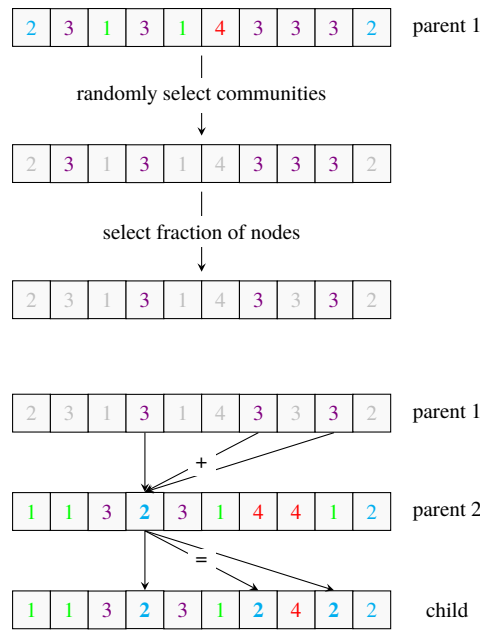
Here $m$ is the number of total links, $A_{ij}$ is the adjacency matrix (which is 1 if there is a connection between nodes $i$ and $j$, otherwise 0), $d_k$ is the degree of node $k$, $c_k = l$ if node $k$ is in community $l$, and $\delta(c_i, c_j) = 1$ if $c_i = c_j$, otherwise 0 [Ne06]. The formula compares the number of internal links in a community with the expected number of links of a randomized graph and therefore implies a definition of the term *community*: A community is a group of nodes that is more tightly connected than expected at random [FC12].



**Fig. 4:** The figure shows two possible community partitions of a network. While the upper partition has a low modularity of about −0.08, the partition below it has a modularity of about 0.22 which is maximal among all possible partitions and thus shows the community structure of the network.
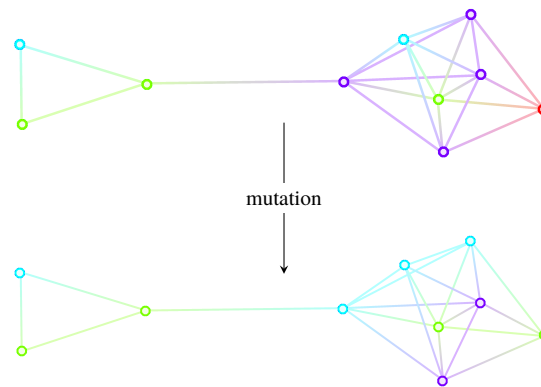
The algorithm provides two methods of selection. The first, in the following called *Elitism*, only keeps the fittest individual. All children are unaltered copies of this single individual. The second method, in the following called *Mating*, retains a certain fraction of the fittest individuals (*parents*), and performs cross-over operations to create new children (see figure 5 for a detailed description). Both methods are indeed a variant of *elitism*, where only the fittest individuals are preserved. Important to note is that the community partition is a relational property, meaning that the labeling of the communities is ambiguous. Two individuals could actually refer to the same community but use different community labels. This is considered by the algorithm by "translating" the community label when doing cross-over operations.

The mutation process consists of two independent parts and is applied to each child by a certain chance. The first mutation process randomly alters parts of the chromosome (see fig. 6a). In the second mutation process a random sample of nodes will inherit the community that is most dominant in their respective neighborhood (the nodes they are directly connected with).

| 2 | 3 | 1 | 3 | 1 | 4 | 3 | 3 | 3 | 2 | parent 1

randomly select communities

| 2 | 3 | 1 | 3 | 1 | 4 | 3 | 3 | 3 | 2 |

select fraction of nodes

| 2 | 3 | 1 | 3 | 1 | 4 | 3 | 3 | 3 | 2 |

| 2 | 3 | 1 | 3 | 1 | 4 | 3 | 3 | 3 | 2 | parent 1

+

| 1 | 1 | 3 | 2 | 3 | 1 | 4 | 4 | 1 | 2 | parent 2

=

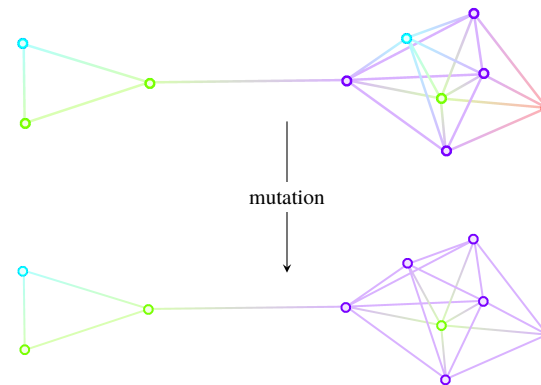| 1 | 1 | 3 | 2 | 3 | 1 | 2 | 4 | 2 | 2 | child

**Fig. 5:** Cross-over process

The aim of the cross-over process is to combine the chromosomes of two individuals (called *parents*) in order to create a new chromosome (called *child*). This is done in the following way: At first a number of communities of *parent 1* are randomly selected. From the selected communities only a certain fraction of the nodes are then chosen. These nodes are then transferred to *parent 2*, i.e. for those nodes we replace their respective community label of *parent 2* by the respective community label of *parent 1*. Since the community labelling is a relational property, we "translate" the community label of *parent 1*. For each community to be transferred it uses the corresponding community label of the **first node** in *parent 2* as community label (indicated by the arrows which all lead to the first node of the community). The resulting chromosome gives the chromosome of the child.

**Fig. 6a:** First mutation process
A random sample of nodes is assigned to random communities (e.g. the community of the red node in the above network is changed from red to green). Through this stochastic process we sample the solution space in an effective way. For example it enables the creation of new communities.



**Fig. 6b:** Second mutation process
A random sample of nodes inherits the most dominant community in its respective neighborhood (e.g. the red node in the above network inherits the purple community since three out of four neighbors are in the purple community). This improves the speed of convergence since it is likely that neighbors belong to the same community.

Finally the algorithm removes communities that consist of only one single node, because a proper community must contain at least two nodes. The whole process is repeated until a certain number of generations have been created. The algorithm then stops and returns the best found solution.

## 3   Experimental results

Since the algorithm finds solutions stochastically, its outcome might differ on every execution. To test the effectiveness of our algorithm two networks consisting of two respectively five communities were used (see fig. 7 and 8). For each trial the algorithm was run 100 times with the same parameters and the number of times where the algorithm found the correct community partition was counted.

We used a population size of 20 and, for each child, a 75% probability for the first mutation process and a 50% probability for the second mutation process. The algorithm was set to terminate after 5 generations for the first network and after 50 generations for the second network. For the *Mating*-method the 15% fittest were selected and all nodes from the selected communities were transferred (see fig. 5).

Results of the testing are shown in table 1. The two communities of the first network were correctly detected in all 100 runs within two generations on average. The communities of the second network were detected correctly in 95% of all cases for the *Elitism*-method and 99% for the *Mating*-method within eight respectively seven generations on average. In most of the failed cases the algorithm merged two communities into one single community.

The results show that our algorithm detects the community structure efficiently, as the average number of generations used to find the community structure was low. It also is effective, because the failure rate was very low as well (0% for the first network, 5% respectively 1% for the second network).



**Fig. 7:** First tested network consisting of 16 nodes, 49 edges and two communities. Number of possible community partitions $\approx 1 \cdot 10^{10}$

Also, the algorithm was tested using a large network with unknown community structure and its result were compared with different methods of *Mathematica*'s built-in function `FindGraphCommunities`. For our algorithm we used the *Mating*-method, the same parameters as for the other two networks and set it to terminate after 100 generations. In order to estimate

| method | fails | avg. gen. | max. gen. |
|--------|-------|-----------|-----------|
| Elitism | 0% | 2.13 ± 0.34 | 3 |
| Mating | 0% | 2.19 ± 0.39 | 3 |

**(a)** Results for first network

| method | fails | avg. gen. | max. gen. |
|--------|-------|-----------|-----------|
| Elitism | 5% | 8.56 ± 5.43 | 30 |
| Mating | 1% | 7.08 ± 2.53 | 27 |

**(b)** Results for second network

**Tab. 1:** Results of the testing
The columns show the used selection method, the number of runs where the algorithm did not find the correct community partition, and the average and maximum number of generations needed to find the best community partition.
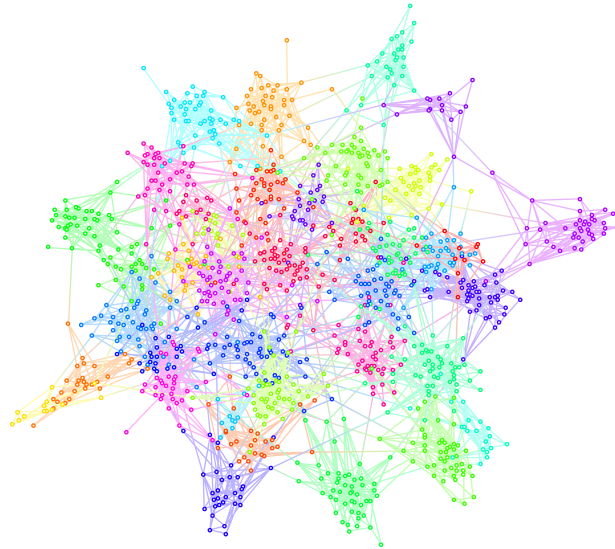


**Fig. 8:** Second tested network consisting of 125 nodes, 390 edges and five communities. Number of possible community partitions $\approx 2.5 \cdot 10^{153}$.

how far *Mathematica*'s results and the result of our algorithm were apart we calculated the partition distance, i. e. the number of nodes that were classified into different communities (this was transformed into an assignment problem and solved using *Mathematica*'s built-in function `FindIndependentEdgeSet`). Our algorithm detected 38 communities (see fig. 9). The community structure differs only slightly from the ones obtained by *Mathematica*'s *Centrality* and *Hierarchical* methods (see table 2).

The results of our algorithm yielded a higher modularity value than all different methods of *Mathematica*'s built-in function `FindGraphCommunities`. We are therefore tempted to claim that our algorithm is more effective. However, the outcome and hence effectiveness of our algorithm does not only strongly depend on the used parameters of our algorithm, but also the investigated network itself. Whether our algorithm outperforms *Mathematica*'s methods can hence not answered with certainty. In terms of speed our implementation of the algorithm can not compete with *Mathematica*'s function. Our implementation is almost $10^3$ times slower (time-wise) and therefore not suited for networks consisting of millions of nodes.



**Fig. 9:** Large network consisting of 1293 nodes and 4145 edges showing the community structure detected by our algorithm. About $6.7 \cdot 10^{2609}$ possible community partitions.

| method | communities | modularity | partition distance |
|--------|-------------|------------|--------------------|
| Centrality | 38 | 0.907235 | 3 |
| Hierarchical | 38 | 0.904835 | 8 |
| Modularity | 35 | 0.900042 | 83 |
| Spectral | 43 | 0.817467 | 143 |
| Our algorithm | 38 | 0.907236 | - |

**Tab. 2:** Comparison of the results of *Mathematica*'s built-in function `FindGraphCommunities` with the result of our algoritm. The columns show the method used by *Mathematica*'s built-in function, the number of detected communities, the modularity of the partition and the partition distance.
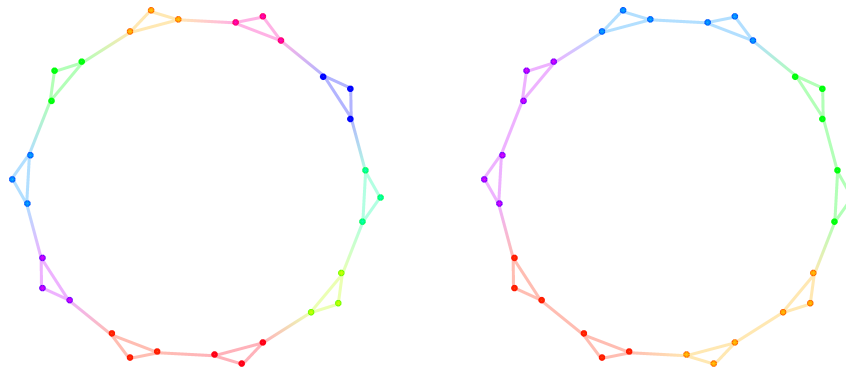
## 4   Analysis of the algorithm

Due to its stochastic way of finding solutions the algorithm does not always find the best solution, which is supported by the high degeneracy of the modularity function with a lot of local maxima close to the global maximum. However our simulations showed that it usually finds a solution that is at least close to the best solution.

Yet there might be a principal problem of the algorithm. Imagine $n$ identical complete graphs $K_m$ with $m$ nodes each, which are arranged on a ring lattice and connected via a single link. Intuitively we would say that each complete graph forms a community since all nodes within a complete graph are densely (in fact maximally) connected and share only two links to nodes from other communities. However it was proved that for a sufficiently large number of graphs $n$ modularity maximization will lead to a partition that merges pairs of communities into one single community [FB07]. Therefore modularity maximization can lead to a partition which seems counter-intuitive as it "neglects" small community structures, which is however the consequence of the implicit definition of a community using modularity (see fig. 10). In general, communities smaller than a certain threshold, which depends on the network size, will always be merged, which is known as *resolution limit*. A possible way to solve this problem is to check every community whether it is a merging of smaller communities.

## 5   Conclusion

In this article we proposed a genetic algorithm aimed at detecting the community structure of a network by maximization of a measure called *modularity*. The algorithm consists of different genetic operators, namely a selection process, a cross-over process and a mutation process. It provides two different methods for the selection process. Tests using several sample networks of known community structure showed that our algorithm detects communities both efficiently and effectively. For the investigated network with unknown

**Fig. 10:** Resolution Limit
Example network consisting of 10 complete graphs with three nodes each. For the left, intuitively expected community partition the modularity is 0.65, whereas in the right partition, where each time two communities are merged, the modularity is 0.675, which is higher than for the expected partition. Modularity maximization will therefore lead to a community partition that might be counter-intuitive.

community structure the algorithm even seems to outperform *Mathematica*'s built-in methods for finding communities. Hence we can conclude that our algorithm shows a very good performance in detecting the communities of complex networks and is suited for cases where a community partition close to the best community partition is required. However, our algorithm can fail to resolve sufficiently small communities as intuitively expected due to a problem known as *resolution limit*. From the investigation of the fitness function *modularity* the question arises if one could improve the definition of modularity such that it also resolves sufficiently small communities as intuitively expected.

# References

[Ba16]    Barabási, A.-L.: Network Science. Cambridge University Press, 2016.

[FB07]    Fortunato, S.; Barthélemy, M.: Resolution limit in community detection. PNAS volume 104 no. 1/, pp. 36–41, 2007.

[FC12]    Fortunato, S.; Castellano, C.: Community structure in graphs. Springer New York Computational Complexity: Theory, Techniques, and Applications/, pp. 490–512, 2012.

[FH16]    Fortunato, S.; Hric, D.: Community detection in networks: A user guide. Physics Reports 659/, Community detection in networks: A user guide, pp. 1–44, 2016.

[Fo10]   Fortunato, S.: Community detection in graphs. Physics Reports 486/3, pp. 75–174, 2010.

[HTF09]  Hastie, T.; Tibshirani, R.; Friedman, J.: The Elements of Statistical Learning: Data Mining, Inference, and Prediction, Second Edition. Springer New York, 2009.

[KN11]   Karrer, B.; Newman, M. E. J.: Stochastic blockmodels and community structure in networks. Phys. Rev. E 83/, p. 016107, Jan. 2011.

[Le18]   Lehnerer, S.: SimonNick/genetic-algorithm-community-detection: Initial release, 2018, URL: https://zenodo.org/badge/latestdoi/139072714.

[Lu07]   Luxburg, U.: A Tutorial on Spectral Clustering. Statistics and Computing 17/4, pp. 395–416, Dec. 2007.

[Ne06]   Newman, M. E. J.: Modularity and community structure in networks. Proceedings of the National Academy of Sciences 103/23, pp. 8577–8582, 2006.

[Po97]   Pothen, A.: Graph Partitioning Algorithms with Applications to Scientific Computing. In (Keyes, D. E.; Sameh, A.; Venkatakrishnan, V., eds.): Parallel Numerical Algorithms. Springer Netherlands, Dordrecht, pp. 323–368, 1997.

[RB08]   Rosvall, M.; Bergstrom, C. T.: Maps of random walks on complex networks reveal community structure. Proceedings of the National Academy of Sciences 105/4, pp. 1118–1123, 2008.

[Sc07]   Schaeffer, S. E.: Survey: Graph Clustering. Comput. Sci. Rev. 1/1, pp. 27–64, Aug. 2007.