# Detecting fuzzy community structures in complex networks with a Potts model

Jörg Reichardt[1,2] and Stefan Bornholdt[1,2]

[1] *Interdisciplinary Center for Bioinformatics, University of Leipzig, Kreuzstr. 7b, D-04103 Leipzig, Germany*
[2] *Institute for Theoretical Physics, University of Bremen, D-28359 Bremen, Germany (present address)*
(Dated: September 10, 2018)

A fast community detection algorithm based on a q-state Potts model is presented. Communities in networks (groups of densely interconnected nodes that are only loosely connected to the rest of the network) are found to coincide with the domains of equal spin value in the minima of a modified Potts spin glass Hamiltonian. Comparing global and local minima of the Hamiltonian allows for the detection of overlapping ("fuzzy") communities and quantifying the association of nodes to multiple communities as well as the robustness of a community. No prior knowledge of the number of communities has to be assumed.

Finding groups of alike elements in data is of great interest in all quantitative sciences. For multivariate data, where the objects are characterized by a vector of attributes, a number of efficient and well understood clustering algorithms exist [1]. They allow to find clusters of similar objects based on a metric between the attribute vectors. If, however, the data is of relational form as, e.g., a network or graph $G(V, E)$ consisting of a set $V$ of $N$ nodes and a set $E$ of $M$ links or edges connecting them and representing some relation between the nodes, the problem of finding alike elements corresponds to discovering communities: sets of nodes interconnected more densely among themselves than with the rest of the network (for a recent review see ref. [2]). For any induced subgraph $g(v, e)$ of the graph $G(V, E)$ with $n$ nodes and $m$ internal edges and $m_{nN}$ edges connecting the $n$ nodes to the $N - n$ remaining nodes of the graph, this can be formalized as:

$$\frac{2m}{n(n-1)} > \frac{2M}{N(N-1)} > \frac{m_{nN}}{n(N-n)}. \tag{1}$$

In other words, the inner link density should be higher than the average link density in the network which again should be higher than the outer density of the community. As a community structure we thus define a set of induced subgraphs $g(v, e)$ that covers $G(N, E)$ and that fulfills (1). Note that the problem of community detection is different from that of minimal cut graph partitioning, as for $g(v, e)$ to be a community it is not necessary that the number of external edges is a global minimum. Rather, it only needs to be smaller than a certain threshold that depends on $G(V, E)$ and the size of $g(v, e)$. We see from (1) that the presence of communities is bound to the presence of inhomogeneities in the link distribution of a graph. Furthermore, it is understood that a community structure is not defined uniquely on a network. Rather, a number of community structures different in size and number of communities may exists that all fulfill the inequalities (1). Certain nodes may belong to the same community in one realization and may be assigned to a different community in another realization. The differences and similarities of these realizations yield valuable information about the robustness of a particular community structure. Furthermore, the nodes which can be assigned into more than one community represent an overlap of possible community structures that cannot be interpreted as a hierarchy of communities, since the overlap may only be partial. Here, we introduce a new algorithm that can rapidly detect a community structure and allows for a quantitative assessment of the individual realizations.

In this paper, we combine the early idea by Fu and Anderson for graph bi-partitioning with a modified Ising Hamiltonian [3] and the recent Potts model clustering of multivariate data by Blatt et al. [4]. This will allow us to map the communities of a network onto the magnetic domains in the ground state or in local minima of a suitable Hamiltonian. For this purpose we alter a q-state Potts Hamiltonian by adding a global constraint that forces the spins into communities according to (1):

$$\mathcal{H} = -J \sum_{(i,j)\in E} \delta_{\sigma_i,\sigma_j} + \gamma \sum_{s=1}^{q} \frac{n_s(n_s - 1)}{2}. \tag{2}$$

Here, $\sigma_i, i = 1...N$ denote the individual spins which are allowed to take $q$ values $1...q$, $n_s$ denotes the number of spins that have spin $s$ such that $\sum_{s=1}^{q} n_s = N$, $J$ is the ferromagnetic interaction strength, $\gamma$ is a positive parameter, and $\delta$ is the Kronecker delta. The first sum is the standard ferromagnetic Potts term for nodes connected by an edge in the network, and is minimized by $\mathcal{H}_{\text{ferr}} = -JM$. It favors a homogeneous distribution of spins over the network. Diversity, on the other hand, is introduced by the second term which sums up all possible pairs of spins which have equal value. It counter-balances the first sum and increases the energy with increasing homogeneity of the spin configuration. It represents a global anti-ferromagnetic interaction being maximal when all nodes have the same spin, and minimal when all possible spin values are evenly distributed over all nodes.

The choice of $\gamma$ determines how strongly the minimum of the combined Hamiltonian depends on the topology of the network. Consider a network of two communities $g_1(n_1, m_1)$ and $g_2(n_2, m_2)$ with $m_{12}$ edges connecting them. For the ground state to be composed of these two communities, $\gamma^*$ has to obey a simple condition

$$\mathcal{H}_{\text{homogeneous}} \geq \mathcal{H}_{\text{diverse}} \qquad (3)$$

$$-J(m_1 + m_2 + m_{12}) + \gamma^* \frac{(n_1 + n_2)(n_1 + n_2 - 1)}{2} \geq$$
$$-J(m_1 + m_2) + \gamma^*(\frac{n_1(n_1 - 1)}{2} + \frac{n_2(n_2 - 1)}{2})$$

$$\gamma^* \geq J\frac{m_{12}}{n_1 n_2}. \qquad (4)$$

Comparing with (1) we see that, apart from the ferro-magnetic coupling $J$, $\gamma^*$ is just the outer link density of community $g_1(n_1, m_1)$. Thus, with the parameter $\gamma$ we enforce a ground state of the system such that all groups of nodes with equal spin have a an outer link density smaller than $\gamma$. Setting $J = 1$ and $\gamma$ to be the average connection probability of the network $p = \frac{2M}{N(N-1)}$ (or $\gamma^* = p\langle J_{ij}\rangle$ for weighted networks), we thus satisfy the second inequality in (1). The first inequality in (1) is satisfied implicitly, because high inner link densities are energetically favored by the Hamiltonian. Different local minima of the Hamiltonian then correspond to different possible assignments of community structures. It is instructive to write the Hamiltonian (2) in the form

$$\mathcal{H} = \sum_{i<j} \delta(\sigma_i, \sigma_j)(\gamma - J_{ij}) \qquad (5)$$

with $J_{ij}$ as the (weighted) adjacency matrix of the graph. The ground state structure of this spin glass Hamiltonian corresponds to the community structure of the network. Fortunately, finding the ground state is difficult only for random networks which usually do not exhibit any clear community structure, and where the ambiguous community assignment corresponds to a typical spin glass situation of multiple local energy minima. Relevant examples of networks with non-random community structure, however, usually correspond to Hamiltonians with prominent ground states in large basins of attraction which makes our approach particularly practicable.

The number of possible communities $q$ is not a critical parameter in the algorithm: it only needs to be chosen large enough to accommodate for all possible communities. If the number of communities is smaller than $q$, the remaining spin states will not be populated. However, since the runtime of the algorithm is linear in $q$, a reasonable value should be chosen ($q < 100$ was sufficient in our case).

It remains to define a measure of the statistical significance of the communities found. Given the number of nodes in a community $n$, the number of inner links $l_{in}$, and the number of outer links $l_{out}$ we can calculate the expected number of possible equivalent communities $E(n, l_{in}, l_{out})$ in a random network of the same size $(N, M)$ and connection probability $p = \frac{2M}{N(N-1)}$:

$$E(n, l_{in}, l_{out}) = \binom{N}{n} \binom{\frac{n(n-1)}{2}}{l_{in}} \binom{n(N-n)}{l_{out}} \times$$
$$p^{l_{in}}(1 - p)^{\frac{n(n-1)}{2} - l_{out}} p^{l_{out}}(1 - p)^{n(N-n) - l_{out}} \qquad (6)$$

If $E(n, l_{in}, l_{out})$ is larger than 1, we can expect to find such a community in a random network of the same size, marking the border of statistical significance.

To practically find or approximate the ground state of our system we employ a simple Monte-Carlo heat-bath algorithm with simulated annealing [5]. Starting from a temperature with an acceptance ratio of $> 95\%$, the system is subsequently cooled down using a decrement function for the temperature of the form $T_{k+1} = \alpha T_k$ with $\alpha = 0.99$ or similar values for the $k^{th}$ step, until it reaches a configuration where no more than a given number of spin flips are accepted during a certain number of sweeps over the network. In one such run, one reaches the ground state or another low lying local minimum that corresponds to a community structure of the network. With a set of several runs, we are able to evaluate the robustness of the community classification by sampling the local minima of the energy landscape of the Hamiltonian. The number of co-appearances of nodes in one community are then binned in an $N \times N$ matrix. We then order the rows and columns of this matrix according to the assignment of communities from a single simulated annealing run. Well defined community structures then appear as blocks of high co-appearance along the diagonal. Off-diagonal instances of high co-appearance indicate a possible overlap between clusters.

Let us first check our algorithm by applying it to a number of computer-generated random test networks with known community structure as suggested in [6]. Nodes are assigned to communities and are randomly connected to members of the same community by an average of $\langle k_{in}\rangle$ and to members of different communities by an average of $\langle k_{out}\rangle$ links. Fixing the average degree of all nodes to $\langle k \rangle = \langle k_{in}\rangle + \langle k_{out}\rangle = 16$, it becomes more and more difficult for any algorithm to detect the communities as $\langle k_{in}\rangle$ decreases on the expense of $\langle k_{out}\rangle$. Sensitivity and specificity are benchmarked over all possible pairs of nodes. As true positive (negative) we count a pair of nodes that is in the same (a different) community by design and is classified accordingly by the algorithm. We tested two sets of networks. The first is composed of four equally sized communities of 32 nodes each and the second is composed of four communities of 128, 96, 64 and 32 nodes respectively. Performance of our algorithm
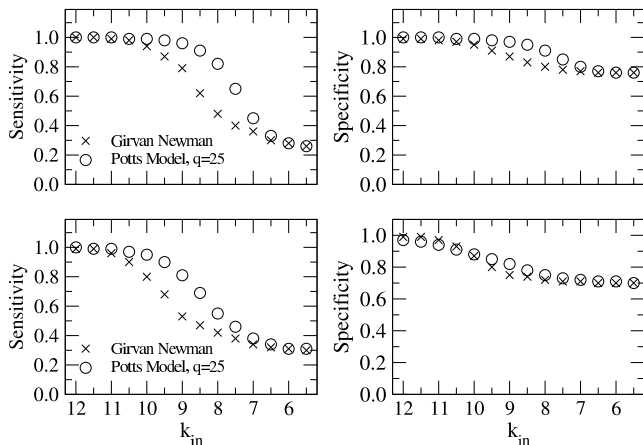
FIG. 1: Benchmark of the algorithm for networks with known community structure and comparison with Girvan and Newman. Top row: 4 communities of 32 nodes each, bottom row: 4 communities of 128, 96, 64 and 32 nodes respectively. Symbol size corresponds to error bars.
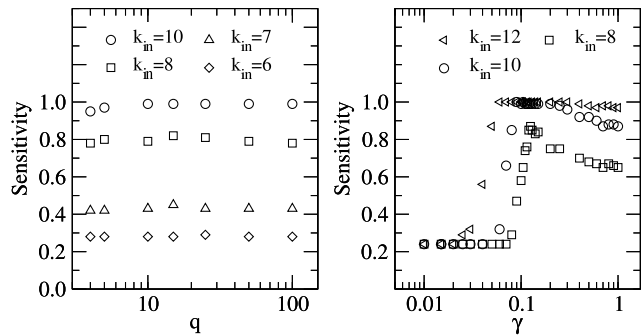


FIG. 2: Robustness of results for the test network with 4 communities of 32 nodes each. Left: Sensitivity vs. $q$ at the end of a Monte-Carlo optimization at different values of $\langle k_{in} \rangle$. Averaged over 50 graphs. Right: Sensitivity for $q = 25$ as a function of $\gamma$ for different values of $\langle k_{in} \rangle$. All results averaged over 10 graphs.
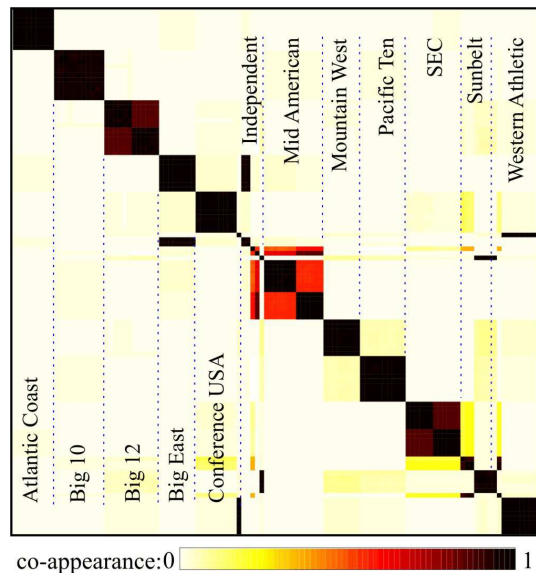


FIG. 3: Co-appearance matrix for the football network. $0.1p \le \gamma \le p$, matrix ordering taken from assignment of teams into conferences according to game schedule.

and, for comparison, the one by Girvan and Newman (GN) [7] is shown in Figure 1. Note the high sensitivity and specificity of our algorithm for both types of networks. When running our algorithm without simulated annealing, but simply relaxing the system at temperature zero from a random initial condition it is extremely fast, yet still performs as good as the GN method.

Figure 2 shows the dependence of the sensitivity of the algorithm on $q$ in the case of the test network with equally sized communities for four different values of $\langle k_{in} \rangle$. Note that results do not depend on $q$. For the same type of test networks, Figure 2 also shows the robustness of the sensitivity with respect to the choice of $\gamma$. The better the communities are defined (the larger $\langle k_{in} \rangle$), the more robust are the results. The maxima of the curves for all values of $\langle k_{in} \rangle$, however, coincide at $\gamma = p \simeq 0.125$ which again justifies this choice of parameter. The same statements apply to the specificity.

One real world example with known community structure is the College Football network from ref. [7]. It represents the game schedule of the 2000 season of Division I of the US college football league. The nodes in the network represent the 115 teams, while the links represent 613 different games played in the course of the year. The community structure of this network arises from the grouping into conferences of 8-12 teams each. On average, each team has 7 matches with members of its own conference and another 4 matches with members of different conferences. We perform a parameter variation in $\gamma$ at ten values between $0.1p \le \gamma \le p$. At each value of $\gamma$ we relax the system 50 times from a randomly assigned initial configuration at $T = 0$ using $q = 50$. Figure 3 shows the resulting $115 \times 115$ co-appearance matrix, normalized and color coded. The ordering of the matrix cor-

responds to the assignment of the teams into conferences according to the game schedule. The dashed blue lines indicate this. Apart from regaining almost exactly the known community structure, our algorithm is also able to detect inhomogeneities in the distribution of intra- and inter-conference games. For instance, we see a large overlap of the Pacific Ten and Mountain West conference and also a possible subdivision of the Mid American conference into two sub-conferences, one of which contains Ball State, Toledo, Central, Eastern, Northern and Western Michigan. This is due to the fact, that geographically close teams are more likely to play against each other as already pointed out in ref. [7].

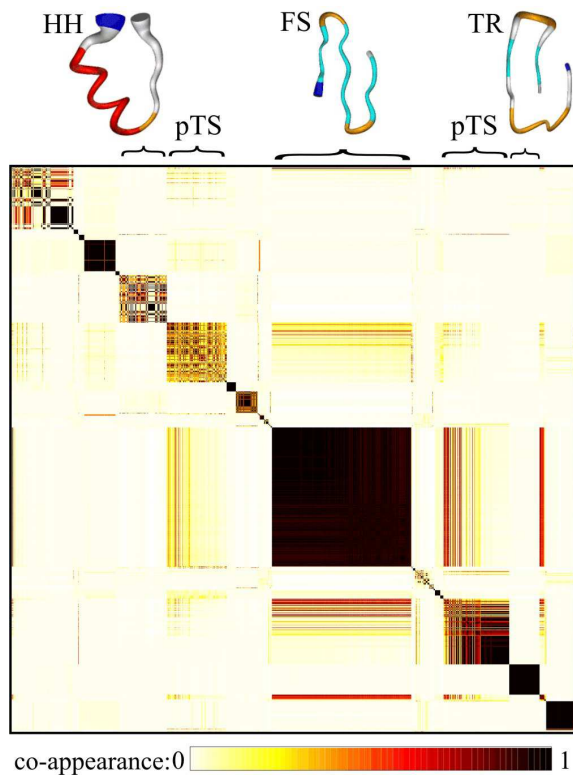Finally we consider a large real world example with

FIG. 4: Co-appearance matrix for the reduced version of the protein folding network. $0.1p \leq \gamma \leq p$, matrix ordering taken from a simulated annealing run of the full network.

with 50 repetitions at each value of $\gamma$ and $q = 50$. For this, we used the reduced version of the folding network as in [8] that contains only nodes which are visited 20 times or more in the course of the simulation, resulting in 1287 nodes and 23948 links. Figure 4 shows the resulting $1287 \times 1287$ nodes co-appearance matrix. The rows and columns are ordered with respect to one single simulated annealing run at $\gamma = p$. Thus, we see how well the ground state is approximated by the local minima and how robust the assignment into communities is with respect to $\gamma$. Again we find a clear characterization of the FS and TR communities. The helical conformations (HH), however, do not occur in one community for all values of $\gamma$ which indicates many different possible assignments into communities and is an indication of their high entropy nature. Furthermore, a number of putative transition states (pTS) could be assigned, that mediate the folding from certain denatured configurations into the folded state.

In conclusion, we discuss a new algorithm for community detection in complex networks based on a modified q-state Potts model. Communities appear as domains of equal spin value near the ground state of the system. which is approximated through Monte-Carlo optimization. Only local information is used to update the spins which makes parallelization of the algorithm straightforward and allows the application to very large networks. On both, computer-generated and real world networks as studied here the algorithm performs fast, often considerably faster than current state-of-the-art algorithms. Without using prior knowledge it automatically detects the number of communities as the number of occupied spin states. As the algorithm is non-deterministic and non-hierarchical, it allows for the quantification of both, the stability of the communities, as well as the affiliation of a node to more than one community ("fuzzy communities").

only partially known community structure, a large protein folding network compiled by Rao and Caflisch [8]. This network represents the conformation space of a 20 amino acids peptide sampled by molecular dynamics at the melting temperature. $5 \times 10^5$ subsequent conformational snapshots were taken at time intervals of 20ps, resulting in 132168 different configurations sampled and 228972 observed transitions between two different conformations. These represent a network of conformations, where a link indicates that two conformations follow each other in time. Analysis of this network yields valuable information about the free energy landscape of the folding Hamiltonian without the need of projecting it onto arbitrarily chosen coordinates. Applying the algorithm to the complete unweighted network using $q = 50$ and $\gamma = p$ yields a largest community of 16,000 nodes, correctly corresponding to the folded state (FS). The statistical weight of the nodes in this community was found to be 55% of the total weight which confirms the expectation of the folded and denatured state being equally populated at the melting temperature. The characteristic conformations of the denatured state, the high enthalpy, high entropy conformations, such as the helical conformations (HH), as well as low entropy conformations such as the curl like trap (TR) are also recognized as communities. Again, $\gamma$ is varied between $0.1p \leq \gamma \leq p$ and $T = 0$

[1] L. Kaufman and P. Rousseeuw, *Finding Groups in Data: an introduction to cluster analysis* (Wiley-Interscience, 1990).
[2] M. E. J. Newman, Eur. Phys. J. B **38** (2004).
[3] Y. Fu and P. W. Anderson, J. Phys. A: Math. Gen. **19**, 1605 (1986).
[4] M. Blatt, S. Wiseman, and E. Domany, Phys. Rev. Lett. **76** (1996).
[5] S. Kirkpatrick, C. G. Jr., and M. Vecchi, Science **220**, 671 (1983).
[6] M. Newman, Phys. Rev. E. **69**, 066133 (2004).
[7] M. Newman and M. Girvan, Proc. Natl. Acad. Sci. **99**, 7821 (2003).
[8] F. Rao and A. Caflisch, J. Mol. Bio. (2004).