Limited resolution in complex network community detection with Potts model approach

Jussi M. Kumpula¹, Jari Saramäki¹, Kimmo Kaski¹, and János Kertész^{1,2}

- Laboratory of Computational Engineering, Helsinki University of Technology, P.O. Box 9203, FIN-02015 HUT, Finland
- ² Department of Theoretical Physics, Budapest University of Technology and Economics, Budapest, Hungary

February 4, 2008

Abstract According to Fortunato and Barthélemy, modularity-based community detection algorithms have a resolution threshold such that small communities in a large network are invisible. Here we generalize their work and show that the q-state Potts community detection method introduced by Reichardt and Bornholdt also has a resolution threshold. The model contains a parameter by which this threshold can be tuned, but no a priori principle is known to select the proper value. Single global optimization criteria do not seem capable for detecting all communities if their size distribution is broad.

PACS. 89.75.-k Complex systems – 89.75.Hc Networks and genealogical trees – 89.75.Fb Structures and organization in complex systems – 89.65.-s Social and economic systems

1 Introduction

Networks are an efficient way to represent a variety of complex systems, including technological, biological and social systems [1,2]. Many networks have substructures called communities, which are, loosely speaking, groups of nodes that are densely interconnected but only sparsely connected with the rest of the network [3,4,5,6]. Detecting such communities is of interest, because they may provide valuable information of the substructure and functionality of the network, e.g., functional modules in metabolic networks, communities of individuals interacting with each other, etc. This analysis can also be extended to more complex properties, including networks of communities [7], roles of nodes inside and between communities [6], and the effect of communities on the dynamics of for example information flow through the network [8].

A large number of algorithms have been developed for detecting the communities, for reviews see [9,10]. A particularly popular method is based on the concept of modularity Q introduced by Newman and Girvan [11]:

$$Q = \sum_{s} e_{ss} - a_s^2, \tag{1}$$

where e_{rs} is the fraction of links that fall between nodes in communities r and s and $a_s = \sum_r e_{rs}$. Detecting communities is equivalent to optimizing the modularity of the network, where optimization is computationally demanding, especially for large networks, but solvable with vari-

ous approximate methods [12,13,11,14,15]. Modularity optimization has been shown to perform well for many test networks [9,16].

Recently, Fortunato and Barthélemy showed that modularity optimization fails to find small communities in large networks, indicating that it is favorable to combine small communities into larger ones [17]. In a network which has L links, there is a characteristic number of links, such that communities with less than $\sqrt{L/2}$ links are not visible. Earlier Reichardt and Bornholdt (RB) had introduced a general framework for community detection. which includes the modularity optimization as a special case [18,19]. Starting from a q-state Potts Hamiltonian, they show that community detection can be interpreted as finding the ground state of an infinite-range spin-glass. Potts spins are assigned to the nodes of the network and the communities can be identified as clusters of aligned spins in the ground state. The model is based on a comparison of the investigated network to a null model which can be arbitrarily chosen. In addition, the method contains a tunable parameter γ for detecting community structures at different hierarchical levels. The Newman-Girvan modularity optimization method is a special case in this general framework, where the null model is the configuration model [20] and $\gamma = 1$. The question arrises whether the more general RB spin-glass-based community detection method is able to overcome the limitations of the modularity optimization. Our paper addresses this question.

We analyze the effect of γ on community detection and consider how to design a network with optimal community structure, study the resolution limit and its estimates by

using a general null model, and, finally, demonstrate the consequences of our findings in certain example cases.

2 Optimal number of communities in the RB model

For detecting the communities in a network, Reichardt and Bornholdt proposed the following Hamiltonian:

$$\mathcal{H} = -\sum_{i \neq j} (A_{ij} - \gamma p_{ij}) \,\delta(\sigma_i, \sigma_j), \tag{2}$$

where A_{ij} denotes the adjacency matrix of the graph with $A_{ij} = 1$ if an edge is present and zero otherwise, $\sigma_i \in$ $\{1, 2, \ldots, q\}$ denotes the group index of node i, γ is a parameter of the model, and p_{ij} denotes the link probability between nodes i and j according to the null model. The null model reflects the connection probability between nodes in a network having no apparent community structure. Possible choices for the null model are, for example, $p_{ij} = p$ and $p_{ij} = \frac{1}{2L}k_ik_j$, where k_i is the degree of node i and L is the number of links in the network. The former null model corresponds to the Erdős-Rényi network [21], whereas the latter one is closely related to the configuration model. The Hamiltonian (2) rewards existing links inside communities, but the reward is reduced if p_{ij} is large. Furthermore, the penalty of a missing link inside a community is proportional to its probability. The modularity Q of Eq. (1) is related to Eq. (2) as $Q = -\mathcal{H}/L$, provided that $\gamma = 1$ and $p_{ij} = \frac{1}{2L}k_ik_j$.

In order to gain some insight to the model given by Eq. (2), we consider two limits of γ . First, when $\gamma \to 0$ each link inside a community comes as a "surprise", while the missing links are not increasing the energy as they are not expected to exist. Thus, in the limit $\gamma = 0$ the minimum energy is obtained when all nodes are assigned into the same community, and the minimum energy is $\mathcal{H} = -2L$. Second, when $\gamma \gg 1$ communities are broken into smaller pieces because the penalty from missing links is large and all existing links are considered to be extremely likely. When γ exceeds the inverse of the minimum of non-zero p_{ij} :s, the terms $A_{ij} - \gamma p_{ij}$ in (2) become all negative, and the minimum energy is obtained when each node is regarded as a separate community, resulting in $\mathcal{H}=0$. This demonstrates that for small values of γ , one can expect to find large community structures, whereas for large values of γ only small community structures are found. The total amount of energy that can possibly be contributed by links and non-links is equal for $\gamma = 1$, which can be regarded as a natural choice. Later we show. however, that optimizing the energy with $\gamma = 1$ does not necessarily yield the obvious and most natural community structure even in a simple test case.

Following the steps in [17], we next consider how to design a connected network with N nodes and L links such that the energy (2) is minimized. In particular, we are interested in the optimal number of communities as a function of L and γ . Therefore, we study a network which has \hat{n} fully connected subgraphs (or cliques) of equal size,

being interconnected with \hat{n} links and arranged in a ring-like structure, see Fig. 1(A). This network has by construction \hat{n} communities, namely the cliques, i.e., the links inside the cliques are intra-community, while those connecting them are inter-community links. The minimization of (2) should reflect this structure providing the \hat{n} equal size communities. Moreover, for such an obvious structure this result should be robust against changing γ or even the null model.

Equation (2) can be rewritten as

$$\mathcal{H} = -\sum_{s=1}^{n} \left(l^s - \gamma [l]_{p_{ij}}^s \right), \tag{3}$$

where l^s is the number of links inside community s and $[l]_{p_{ij}}^s$ is the expected number of links in that community given the link distribution p_{ij} and the current assignment of nodes into communities [19]. In order to be compatible with the calculations in [17], we choose first to use $p_{ij} = \frac{1}{2L}k_ik_j$, i.e., our reference system is the configuration model. In this case, $[l]_{p_{ij}}^s = \frac{1}{4L}K_s^2$, where K_s is the sum of degrees of nodes in community s. It is straightforward to show that Eq. (3) is minimized when each community has L/n-1 links. Then, the energy is

$$\mathcal{H}_{min}(n,\gamma,L) = -(L - n - \gamma \frac{L}{n}). \tag{4}$$

The optimal number of communities, n^* , is obtained as the zero of the derivative $d\mathcal{H}_{min}(n,\gamma,L)/dn$. This yields $n^*=\sqrt{\gamma L}$, which in turn gives back the result of [17] for $\gamma=1$. If the null model is $p_{ij}=p$, i.e., an Erdős-Rényi graph, a similar calculation shows that the energy minimum is obtained when each community has an equal number of nodes. In this case, the optimal number of communities is $n^*=\sqrt{\gamma L\frac{N}{N-1}}$.

Let us suppose that, given N and L, we have constructed a ring-like network as described above, having more than $\sqrt{\gamma L}$ cliques. Previous analysis shows, counterintuitively, that when each clique is considered as a separate community the energy (3) is not minimized. Instead, it is better to relabel the communities so that small communities are merged to form larger ones. On the other hand, if the number of communities is much smaller than $\sqrt{\gamma L}$ it might be advantageous to split large communities into smaller ones. Therefore, the original, well defined communities are not necessarily found by optimizing the quality function (3). In particular, small communities may remain unresolved.

3 Resolution threshold with a general null model

The previous section suggests that the most common null models, $p_{ij} = \frac{1}{2L} k_i k_j$ and $p_{ij} = p$, lead to merging of small communities in large networks. In this section we investigate the case of a general null model and the effect of γ on the resolution. Hence, we consider a general undirected,

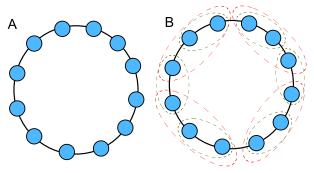


Figure 1. A: a ring-like network of n cliques joined by n links. B: Consecutive cliques can be merged to form larger communities. The optimal configuration depends on the network parameters and γ .

unweighted network having N nodes and L links. Let us suppose that the nodes have somehow been assigned to communities. We take two communities, labeled s and r, each having l^s and l^r links inside and $l^{s \leftrightarrow r}$ links between them. The question is, when should the communities be merged? At first, the energy (3) reads as follows

$$E_1 = \sum_{t \neq s,r} (-l^t + \gamma[l]_{p_{ij}}^t) + (-l^r + \gamma[l]_{p_{ij}}^r) + (-l^s + \gamma[l]_{p_{ij}}^s)$$

whereas after combining the communities the energy is

$$E_2 = \sum_{t \neq s,r} (-l^t + \gamma[l]_{p_{ij}}^t) + \left[-(l^r + l^s + l^{s \leftrightarrow r}) + \gamma[l]_{p_{ij}}^{s+r} \right],$$
(6)

where $[l]_{p_{ij}}^{s+r}$ is the expected number of links in the combined community and $l^{s \leftrightarrow r}$ is the number of links between the communities s and r. The communities should be combined if

$$\Delta E = E_2 - E_1 = -l^{s \leftrightarrow r} + \gamma \left([l]_{p_{ij}}^{s+r} - [l]_{p_{ij}}^r - [l]_{p_{ij}}^s \right) < 0.$$
(7)

But $[l]_{p_{ij}}^{s+r} - [l]_{p_{ij}}^r - [l]_{p_{ij}}^s \equiv [l]_{p_{ij}}^{s \leftrightarrow r}$ is the expected number of links between the communities and equation (7) reduces to

$$\gamma[l]_{p_{ij}}^{s \leftrightarrow r} < l^{s \leftrightarrow r}. \tag{8}$$

As the communities have n_s and n_r nodes each, the maximum number of links between the communities is $n_s n_r$. In a large network, the average probability for connecting two nodes has to be of the order of N^{-1} , regardless of the null model. Therefore, the expected number of links between the communities, $[l]_{p_{ij}}^{s \leftrightarrow r}$, is on average of the order of $n_s n_r/N$. Using this estimate in Eq. (8) suggests that even a single link between small communities may trigger merging if the communities are small, i.e., $n_s, n_r \ll N$. In particular, communities of approximately the same size are merged if

$$n_s \approx n_r \lesssim \sqrt{N l^{s \leftrightarrow r} / \gamma}.$$
 (9)

Now, let us suppose the communities are loosely connected to each other, that is, $l^{s \leftrightarrow r} \sim 1$. When this is applied in

Eq. (9), we obtain that it is beneficial to combine communities smaller than $\sim \sqrt{N/\gamma}$. This is the lower limit for the community size that the method is able to detect. Large values of γ decrease this resolution threshold, but rather inefficiently. When the communities are more densely interconnected, the resolution threshold increases. In the extreme (unphysical) limit, when the communities are connected with $l^{s \leftrightarrow r} \sim L$ links, Eq. (9) indicates that even communities whose size is comparable to the whole network may remain unresolved. Similar results for the resolution thresholds were obtained in Ref. [17] for the case $\gamma = 1$, $p_{ij} = k_i k_j / 2L$: Two tightly connected communities may be merged if each has less than L/4 links, whereas the lower limit is $\sqrt{L/2}$ for communities connected with a single link.

The community structure found by the RB model corresponds to the global minimum of (3). It should be noted that the previous calculations do not prove that the particular communities s and r will be in the same community for the global minimum. The calculations show, however, that the global minimum does not contain connected communities smaller than the above mentioned size limits, because by combining them a lower energy would be achieved.

Equation (8) shows also that cliques are stable against splitting for any reasonable γ . Suppose that a clique is split into two parts each having n_s and n_r nodes. The parts have the maximum number $l^{s \leftrightarrow r} = n_s n_r$ of connecting links. Substituting this and $[l]_{p_{ij}}^{s \leftrightarrow r} \sim n_s n_r/N$ into Eq. (8) shows that it is beneficial to split a clique only when $\gamma \sim N$. Such high value of γ does not, however, make sense because it would lead to splitting the network into individual nodes for the following reason. In this case the average value of links from a node according to the null model would exceed the maximum possible number of links from a node, and the communities would be split into individual nodes. We conclude that when $\gamma \ll N$ cliques and almost complete cliques are not split.

4 Examples

We illustrate the consequences of the above results in three example cases. Let us first consider the simplest possible case of community detection [17]: the network consists of a ring of complete cliques joined by single links, Fig. 1(A). There are n cliques and each clique has m nodes and m(m-1)/2 links. Figure 1(B) shows a case where r consecutive cliques are merged to form a single community. A straighforward calculation shows that in this case, the energy is given by

$$\mathcal{H}_{\frac{n}{r}}(\gamma) = -n\left(\frac{m(m-1)}{2} + \frac{r-1}{r}\right) + \gamma \frac{rn}{4L}m^2(m-1)^2,\tag{10}$$

when $p_{ij} = \frac{1}{2L}k_ik_j$. By joining cliques, we get a "bonus" from the links joining the cliques, i.e., term (r-1)/r, but in large communities the expected number of links inside the communities is increasing faster than in small communities. Thus, for small γ the merged cliques have low

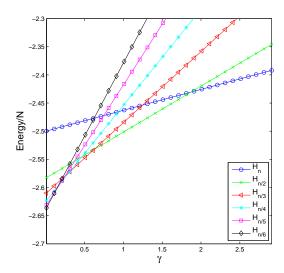


Figure 2. Energy (10) as a function of γ for a system where n = 60, m = 6 and $r = 1, \ldots, 6$. The optimal configuration depends on γ and the natural communities are found only when $\gamma > 1.875$, c.f. Eq. (11).

energy, but as γ increases the energy is growing quite fast as illustrated in Fig. 2. The optimal configuration found by optimizing Eq. (2) is the configuration that has the lowest energy for the given values of n, m and γ . Especially, it can be shown that the natural communities are found only if

$$m(m-1) + 2 > \frac{n}{\gamma}.\tag{11}$$

When the link probability is $p_{ij} = p$, we obtain the same result with a correction term of the order of $(\gamma m)^{-1}$.

Our second example is a random network, which has often been used as a test network for community detection algorithms [11]. The network consists of n communities each having m nodes. Each node has on average $\langle k \rangle$ links of which $\langle k_{in} \rangle$ go to random nodes in the same group and $\langle k_{out} \rangle = \langle k \rangle - \langle k_{in} \rangle$ links lead randomly to nodes in other communities. Let us now calculate when, on the average, it is beneficial to merge two designed communities. We obtain that the average number of observed links between the communities is

$$\langle l^{s \leftrightarrow r} \rangle = \frac{m}{n-1} \langle k_{out} \rangle,$$
 (12)

where the averaging is done over all the realizations of the network. Note that if $m/(n-1)\langle k_{out}\rangle < 1$ we have to set $\langle l^{s \leftrightarrow r} \rangle = 1$ because we are considering only communities which are connected by at least one link. The null model is again $p_{ij} = \frac{1}{2L} k_i k_j$. According to the null model the expected number of links between communities is

$$[l]_{p_{ij}}^{s \leftrightarrow r} = \frac{1}{2L} (m\langle k \rangle)^2 = \frac{m\langle k \rangle}{n},$$
 (13)

when averaged over the realizations of the network. Now Eqs. (8), (12) and (13) give that the communities are

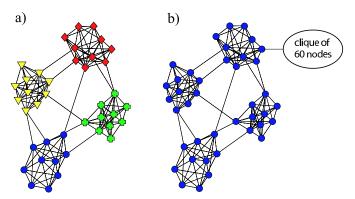


Figure 3. An example of the effect of network size on the resolution of the Potts method. Symbols correspond to communities. See text for details.

merged if

$$\gamma < \begin{cases} \frac{\langle k_{out} \rangle}{\langle k \rangle} \frac{n}{n-1} & \text{for } n < 1 + m \langle k_{out} \rangle \\ \frac{n}{m \langle k \rangle} & \text{for } n > 1 + m \langle k_{out} \rangle. \end{cases}$$
(14)

For typical values $n=4,\ m=32,\ \langle k\rangle=16$ and $\langle k_{out}\rangle=1\dots 8$ we find that $\gamma=1$ should give the correct communities. Thus, it is not surprising that community detection based modularity optimization (1) performs well for this network. We point out that it is possible to choose the parameters $n,m,\langle k\rangle$ and $\langle k_{out}\rangle$ in such a way that modularity optimization with $\gamma=1$ does not give the designed communities.

As a third example we note that the Potts Hamiltonian (2) can be generalized to weighted networks by using a weighted adjacency matrix W_{ij} . A simple way to do this is to define

$$\mathcal{H}_w = -\sum_{i \neq j} (W_{ij} - \gamma \overline{w_{ij}} p_{ij}) \delta(\sigma_i, \sigma_j), \tag{15}$$

where $\overline{w_{ij}}$ is the average link weight. In this way, strong links inside communities lower the energy greatly, while missing links are assumed to be of average weight. Using weights does not, however, resolve the underlying problem that in a large network even a single link easily exceeds the expected weight between the communities.

Finally, in Figure 3 we demonstrate the effect of network size on the resolution of the Potts method. Panel a) shows a network of four groups of 10 nodes. We have compared the energies (3) for two community divisions using $\gamma = 1$ and the configuration null model. $E_1 = 0$ is the energy for the case when all four groups are assigned to a single community, whereas $E_4 = -100.2$ is the energy when the four groups are each assigned to a different community. In this case, $E_4 < E_1$, i.e. the groups are properly identified as communities. However, if the original network of panel a) is modified such that an additional 60-clique community is connected to it via a single link, the situation is changed. All nodes of this new 60-clique are assigned to a single community. Now, $E'_1 = -271.17$ is the energy when the original four groups are merged into a Potts community, and $E'_4 = -269.73$ the energy when they are assigned to separate communities. Hence $E_1' < E_4'$, i.e., the energy for merged groups is lower. This is unphysical, since connecting the new clique via a single link does not alter the original four-group topology.

5 Conclusions

In the light of the above considerations it is clear that the problem of the resolution limit is not restricted to the Newman-Girvan method of modularity optimization. Rather, it is a flaw which seems to be present in any community detection scheme based on global optimization of intra- and extra-community links and on a comparison to any null model. The limited resolution rises from the fact that in a large network the expected number of links between two small sets of nodes is small and even a single link between the sets is enough to merge them. The null model uses the global probability of connecting nodes while the small communities should be considered on a more local level. We agree with the conclusion of Ref. [17] that presently, in large networks, local community detection methods like [6] seem to perform better from the point of view of resolution. An alternative solution to this problem could be to iteratively change the parameter γ when looking for smaller communities in a large network.

Our results indicate that when the community structure is not known beforehand, there is no simple way to decide which γ gives the most relevant communities. Moreover, if the size distribution of the communities is broad, like in collaboration networks [6] or school friendship networks [22], there is no single proper value of γ for the optimal resolution. The hierarchical structure can be examined to some extent by using several values of γ [19], but this method may find too much hierarchy in the network as it tends to artificially merge communities. Because of this tendency, one should always carefully investigate the structure of the found communities.

Acknowledgements: JK thanks Santo Fortunato for inspiring discussions at ISI, Torino. This work was partially supported by OTKA K60456 and the Academy of Finland (Center of Excellence program 2006-2011).

References

- 1. R. Albert, A.L. Barabási, Rev. Mod. Phys. **74**, 47 (2002)
- S. Boccaletti, V. Latora, Y. Moreno, M. Chavez, D.U. Hwang, Phys. Rep. 424, 175 (2006)
- R. Guimerá, S. Mossa, A. Turtschi, L.A.N. Amaral, PNAS 102, 7794 (2005)
- A. Arenas, L. Danon, A. Díaz-Guilera, P.M. Gleiser, R. Guimerá, Eur. Phys. J. B 38, 373 (2004)
- 5. R. Guimera, L.A.N. Amaral, Nature 433, 895 (2005)
- G. Palla, I. Derényi, I. Farkas, T. Vicsek, Nature 435, 814 (2005)
- P. Pollner, G. Palla, T. Vicsek, Europhys. Lett. 73, 478 (2006)
- 8. J.P. Onnela, J. Saramäki, J. Hyvönen, G. Szabo, D. Lazer, K. Kaski, J. Kertész, A.L. Barabási, Proc. Natl. Acad. Sci. (USA), in press, e-print physics/0610104

- L. Danon, A. Díaz-Guilera, J. Duch, A. Arenas, J. Stat. Mech. 2005, P09008 (2005)
- 10. M.E.J. Newman, Eur. Phys. J. B 38, 321 (2004)
- M.E.J. Newman, M. Girvan, Phys. Rev. E. 69, 026113 (2004)
- 12. M.E.J. Newman, Phys. Rev. E **74**, 036104 (19) (2006)
- A. Clauset, M.E.J. Newman, C. Moore, Phys. Rev. E 70, 66111 (2004)
- 14. J. Duch, A. Arenas, Phys. Rev. E 72, 027104 (2005)
- 15. M.E.J. Newman, Phys. Rev. E 69, 066133 (2004)
- M. Gustafsson, M. Hornquist, A. Lombardi, Physica A 367, 559 (2006)
- S. Fortunato, M. Barthélemy, Proc. Natl. Acad. Sci. USA 104, 36-41 (2007), e-print physics/0607100
- J. Reichardt, S. Bornholdt, Phys. Rev. Lett. 93, 218701 (2004)
- 19. J. Reichardt, S. Bornholdt, Phys. Rev. E 74, 016110 (2006)
- 20. M.E.J. Newman, SIAM Review 45, 167 (2003)
- 21. P. Erdös, A. Rényi, Publ.Math.Debrecen ${\bf 6},$ 290 (1959)
- M.C. Gonzales, H.J. Herrmann, J. Kertész, T. Vicsek, Physica A (in press), physics/0611268