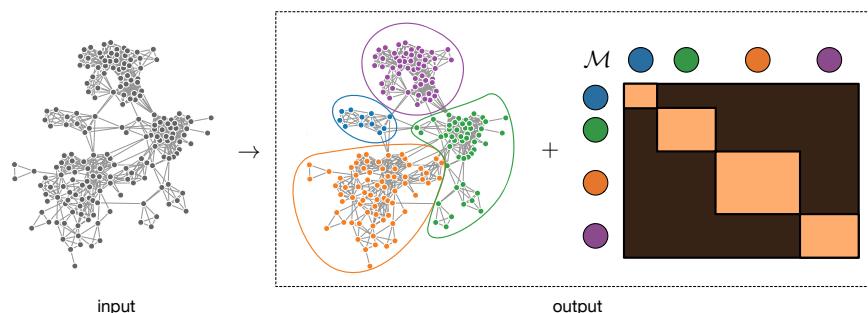


1 Detecting community structure

The process of decomposing a particular network G into “clusters” of nodes and “bundles” of edges is called **community detection**. By identifying groups of nodes that connect to other groups in similar ways, community detection allows us to obtain a description of the network’s structure that falls between the level of local (node-level) and global (network-level) statistics.¹ This kind of “graph clustering” can be used for many downstream purposes, from describing a network’s architecture in terms of its modules to better predict missing information.

The input to community detection is always a network, its nodes and their connections. In some cases, we can also use various forms of auxiliary information, e.g., node attributes or edge weights, to produce a better decomposition of a network into groups. The output is always a partition \bar{z} of the network, i.e., an assignment of nodes into communities (or modules, or compartments, or just groups). Given that partition, we can then derive a community mixing matrix \mathcal{M} that tells us how these communities fit together to construct the whole.



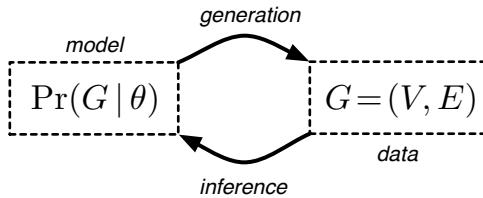
There are now hundreds of different algorithms for detecting communities in networks. At the end of this lecture, we will discuss some of their differences, and what it means that different algorithms can produce different partitions of the same network. In the meantime, we will learn about one of the more powerful and statistically principled approaches for community detection, and see how to use the stochastic block model (SBM) to detect communities in real networks.

2 Statistical inference

Each of the random graph models in our toolbox is a *probabilistic generative model*, because it defines a parameterized probability distribution over graphs $\Pr(G | \theta)$, where θ encodes the structural assumptions that define the particular model.

¹This wording recapitulates our definition of a “community” and generalizes many notions of group-level structure, including assortative or disassortative communities (edges inside or between groups), ordered communities (linear hierarchy of groups), core-periphery structure (dense core and sparse periphery), and mixtures thereof. See Lecture 5.

If we ourselves choose some value of θ , we can then *generate* a network instance G by flipping a set of coins whose biases are derived from θ . For instance, recall how we generated graphs from the Erdős-Rényi model with $\theta = p$, the edge density (Lecture 3), or from the stochastic block model, with $\theta = (c, \vec{z}, \mathcal{M})$, the community partition and mixing matrix (Lecture 5). The reverse of this process is called *inference*, in which we ask, for some observed network G and model $\Pr(G | \theta)$, what is the best choice of θ for generating G ? Together, generation and inference form a virtuous circle for modeling network data:



Given some data G , we can infer a model² $\Pr(G | \hat{\theta})$, and from that model, we can generate new (synthetic) data. If the parameters are interpretable, i.e., they have some understandable meaning for how edges are generated, inferring them can give us direct insight about the network's organization and why some edges exist but others do not. For instance, an inferred mixing matrix $\hat{\mathcal{M}}$ provides this kind of insight about how groups interact. A parameterized model also allows us to generate synthetic networks, which we can use as a substrate for simulation experiments (as in Lecture 5), or to check the validity of the inferred model, e.g., by comparing network measures of the generated networks against those of the empirical network, as in a null model (see Lecture 3).

Generative models have a variety of attractive features that partly explain why they represent the state-of-the-art in community detection. For instance, they make our assumptions about modular structure explicit, and their parameters θ typically relate to specific structural hypotheses (e.g., degrees and randomness, etc.).³ And, they make it easier to compare competing models of a network's structure, because each can be phrased as a probability distribution, and statistics provides many ways of comparing two distributions. And finally, they are easily adapted to tasks of missing attribute and link prediction, based on a partial or past set of observations. These benefits do come at some cost, as either specifying a new model or fitting it to data can seem more complicated than less principled heuristics. In practice, however, the benefits for science are very much worth the cost. In the rest of this lecture, we will focus on developing the tools of statistical inference for the stochastic block model (SBM) and its degree-corrected variant (the DC-SBM).

²A “hatted” symbol $\hat{\theta}$ means a parameter value we estimated with the benefit of data. An unhatted symbol θ means one we choose without data, i.e., we chose it ourselves to create synthetic data.

³Most generative models for networks are what you might call “edge generative,” meaning that they do not consider networks with different numbers of vertices, only networks of fixed size with different patterns of edges.

Optimize or sample? Selecting the best choice of θ for a model requires writing down some function $f(G, \theta)$ of both the data and the parameters that scores the quality of each choice. We then would either (i) *search* the parameter space of this function for a single high-scoring choice $\hat{\theta}$, ideally the global optimum but often we have to settle for merely a local one,⁴ or (ii) *sample* the parameter space for a set of high-scoring choices, in much the same way we used a sampler to generate configuration model random graphs (in Lecture 3). Both approaches are useful, depending on the circumstances. Below, we mainly consider optimization approaches.

2.1 The likelihood of a SBM

Recall that to specify a stochastic block model, we need to choose the number of communities c , the partition of nodes into groups \vec{z} , and the stochastic block matrix \mathcal{M} . From these, we can write down an explicit function of the probability or **likelihood** \mathcal{L} that we generate exactly (i) the set of observed edges $(i, j) \in E$ and (ii) the set of observed non-edges $(i, j) \notin E$ of a network G :

$$\begin{aligned} \mathcal{L}(G | z, \mathcal{M}) &= \prod_{i,j} \Pr(i \rightarrow j | z, \mathcal{M}) \\ &= \left(\prod_{(i,j) \in E} \Pr(i \rightarrow j | z, \mathcal{M}) \right) \left(\prod_{(i,j) \notin E} 1 - \Pr(i \rightarrow j | z, \mathcal{M}) \right) \\ &= \left(\prod_{(i,j) \in E} \mathcal{M}_{z_i, z_j} \right) \left(\prod_{(i,j) \notin E} 1 - \mathcal{M}_{z_i, z_j} \right), \end{aligned} \quad (1)$$

where we drop the c parameter for notational convenience.

When G is a simple graph (undirected edges, no self-loops), Eq. (1) gives the probability of observing exactly the set of $\binom{n}{2}$ values in the adjacency matrix, both 1s and 0s. The second line above separates the product over pairs i, j into two parts, the total probability of the 1s and the total probability of the 0s. In this way, every pair i, j still appears in the likelihood function exactly once, and the function contains $\Theta(n^2)$ terms.

Eq. (1) is sufficient for scoring different choices of θ for the SBM, but it is not a particularly efficient function. A key feature of the SBM is that two nodes i and j in the same community $z_i = z_j$ are *stochastically equivalent*, meaning that each of their pairwise probabilities are equal, i.e., $\Pr(i \rightarrow k) = \Pr(j \rightarrow k)$. These symmetries imply that many of the terms in Eq. (1) are the same, and we can group them together to write down a simpler likelihood function.

⁴Because it's usually hard to guarantee that we have found the global maximum of $f(G, \theta)$. It's easier to guarantee that for some choice θ , the nearby points in the parameter space are no better than $f(G, \theta)$.

Before we write down such a function, we first define two auxiliary variables

$$e_{rs} = \sum_{i < j}^n A_{ij} \delta_{r,z_i} \delta_{s,z_j} \quad \text{and} \quad n_{rs} = \sum_{i < j}^n \delta_{r,z_i} \delta_{s,z_j} , \quad (2)$$

which count the actual number of edges that connect nodes in group r to nodes in group s , and the number of possible such connections,⁵⁶ respectively, where δ_{ab} is the Kronecker delta function. The sum indices here run $1 \leq i < j \leq n$ because G is undirected. If we allow self-loops, then we would have them run $1 \leq i \leq j \leq n$ instead.

With these, we can rewrite the likelihood of G as a product series over pairs of groups r, s :

$$\mathcal{L}(G | z, \mathcal{M}) = \prod_{r,s} (\mathcal{M}_{rs})^{e_{rs}} (1 - \mathcal{M}_{rs})^{n_{rs} - e_{rs}} , \quad (3)$$

where the first term $(\mathcal{M}_{rs})^{e_{rs}}$ is the likelihood of generating exactly e_{rs} edges (the 1s in the adjacency matrix), each with probability \mathcal{M}_{rs} , and the second term is the likelihood of generating exactly $n_{rs} - e_{rs}$ non-edges (0s in the adjacency matrix), each with probability $1 - \mathcal{M}_{rs}$. Crucially, just like in Eq. (1), every pair i, j is accounted for in Eq. (3) (do you see why?). Notably, the form of Eq. (3) is similar to that of a Binomial distribution, where the bias on a coin depends on the group labels r and s , e_{rs} specifies the number of “successes” and n_{rs} is the number of attempts.⁷

This version of the SBM likelihood function takes $\Theta(c^2)$ time to compute, a significant improvement, if we have precomputed the edge bundle matrix e_{rs} and group size array n_r . Given a partition z , both variables can be computed in $\Theta(n^2)$ time. If we had to recompute them for each choice of θ , Eq. (3) would not be more efficient than Eq. (1). Hence, we can only be more efficient if we pay that cost once, at initialization, and then efficiently update them for each subsequent choice of θ .

2.1.1 The maximum likelihood choice

A common strategy for sorting among different parameterizations is to prefer the ones that maximize the likelihood function. Subject to mild regularity conditions, maximum likelihood is asymptotically *consistent*, meaning that if we generate some data using a parameter θ , the *maximum likelihood estimate* (MLE) of that parameter $\hat{\theta}$ will converge on θ in the limit of large data sets.

⁵The matrix e in the SBM is similar to the mixing matrix ω in the DC-SBM. They differ, however, in how they count the within-group elements e_{rr} and ω_{rr} . In the SBM, e_{rr} counts within-group edges, while ω_{rr} counts within-group edge “stubs,” meaning that each within-group edge is counted twice, once for each endpoint (stub).

⁶Note that in order to make G a simple graph, when $r = s$, $n_{rr} = \binom{n_r}{2}$, and otherwise $n_{rs} = n_r n_s$. This implies that e_{rr} should count edges, rather than stubs.

⁷If one wants to do link prediction with the SBM, one simply sets $\text{score}(i, j) = \mathcal{M}_{z_i, z_j}$, which is exactly the probability that i, j are connected under the (fitted) SBM.

The stochastic block model is a model where this property holds.⁸ Hence, because the probability that i, j are connected in G is a Bernoulli trial with parameter $\mathcal{M}_{z_i z_j}$, the MLE of \mathcal{M}_{rs} is

$$\hat{\mathcal{M}}_{rs} = e_{rs} / n_{rs} . \quad (4)$$

That is, the best choice for a mixing parameter is the observed fraction of groupwise connections.

Plugging Eq. (4) into Eq. (3), and then taking the log of both sides yields a *log-likelihood function*:

$$\ln \mathcal{L} = \sum_{r,s} e_{rs} \ln \frac{e_{rs}}{n_{rs}} + (n_{rs} - e_{rs}) \ln \left(\frac{n_{rs} - e_{rs}}{n_{rs}} \right) . \quad (5)$$

which depends only on the actual e_{rs} and possible n_{rs} number of edges in the r, s edge bundle, both of which can be derived using only the partition choice z and graph G .⁹ Log-likelihood functions are significantly easier to work with on computers because the logarithm converts multiplication (okay) into addition (easy), and allows us to represent extremely small numbers without worrying about underflow errors. And, crucially, the parameter choice θ that maximizes \mathcal{L} also maximizes $\ln \mathcal{L}$ (do you see why?).

2.1.2 An example

To see how the SBM likelihood function sorts different partitions, consider a simple $n = 6$ node and $m = 7$ edge network composed of two triangles linked by a single edge, and two possible $c = 2$ partitions.

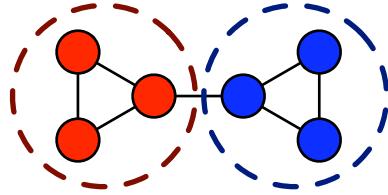
In the “good” partition, each triangle is placed into its own group, and $n_{\text{red}} = n_{\text{blue}} = 3$. In the “bad” partition, we place both nodes on either end of the “bridge” edge into the same group as one of the triangles. If Eq. (5) is to sort partitions well, then clearly the good partition should have a higher likelihood than the bad partition.

The \mathcal{M} stochastic block matrix in each case has been populated with the corresponding MLE fractions e_{rs}/n_{rs} . Then, applying Eq. (5) to the corresponding matrices yields:

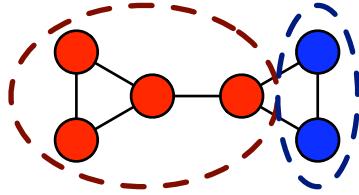
$$\begin{aligned} \ln \mathcal{L}_{\text{good}} &= \left(3 \ln \frac{3}{3} \right) + \left(1 \ln \frac{1}{9} + 8 \ln \frac{8}{9} \right) + \left(3 \ln \frac{3}{3} \right) = -3.1395 \dots \\ \ln \mathcal{L}_{\text{bad}} &= \left(4 \ln \frac{4}{6} + 2 \ln \frac{2}{6} \right) + \left(2 \ln \frac{2}{8} + 6 \ln \frac{6}{8} \right) + \left(1 \ln \frac{1}{1} \right) = -8.3178 \dots . \end{aligned}$$

⁸Somewhat amazingly, not all network models have this property. For network parameter estimates to be consistent, edges must be conditionally independent, which is true in the SBM and DC-SBM: edges are independent, conditioned on their group labels. For more information about network models that lack this property, see Shalizi & Rinaldo, *Ann. Statist.* 41(2), 508–535 (2013), <https://projecteuclid.org/euclid-aos/1366980556>.

⁹In both the SBM here, and the DC-SBM described below, we define $0^0 = 1$. This choice prevents the numerical calculation from failing when either $e_{rs} = 0$ or $n_{rs} - e_{rs} = 0$.



$$\mathcal{L}_{\text{good}} = 0.043304\dots \\ \ln \mathcal{L}_{\text{good}} = -3.1395\dots$$



$$\mathcal{L}_{\text{bad}} = 0.000244\dots \\ \ln \mathcal{L}_{\text{bad}} = -8.3178\dots$$

e_{rs}/n_{rs}	red	blue
red	3/3	1/9
blue	1/9	3/3

e_{rs}/n_{rs}	red	blue
red	4/6	2/8
blue	2/8	1/1

We can measure how much better the “good” partition is than the “bad” one by computing $e^{\Delta \ln \mathcal{L}}$, which is their likelihood ratio (do you see why?). Plugging in our results shows that the good partition is $\exp(\ln \mathcal{L}_{\text{good}} - \ln \mathcal{L}_{\text{bad}}) = 177$ times more likely to generate the observed data than the “bad” partition. In other words, the good partition is a much better model of the data.

2.2 The likelihood of a DC-SBM

Recall that to specify a degree-corrected SBM, we need to choose the number of communities c , the partition of nodes into groups \vec{z} , the expected degree sequence \vec{k} , and the mixing matrix \mathcal{M} .

Under the DC-SBM, the value we assign to any particular adjacency A_{ij} is a Poisson-distributed random variable with mean $\gamma_i \gamma_j \mathcal{M}_{z_i z_j}$, where each γ_i is a node-specific model parameter that quantifies the fraction of the group z_i ’s total degree that belongs to node i (see Lecture 5). Hence, the likelihood function is

$$\begin{aligned} \mathcal{L}(G | z, \gamma, \mathcal{M}) &= \prod_{i,j} \text{Poisson}(\gamma_i \gamma_j \mathcal{M}_{z_i z_j}) \\ &= \prod_{i < j} \frac{(\gamma_i \gamma_j \mathcal{M}_{z_i z_j})^{A_{ij}}}{A_{ij}!} \exp(-\gamma_i \gamma_j \mathcal{M}_{z_i z_j}) \times \prod_i \frac{\left(\frac{1}{2} \gamma_i^2 \mathcal{M}_{z_i z_i}\right)^{A_{ii}/2}}{(A_{ii}/2)!} \exp\left(-\frac{1}{2} \gamma_i^2 \mathcal{M}_{z_i z_i}\right) \quad (6) \end{aligned}$$

where the second part gives the likelihood of self-loops, and the factors of $1/2$ prevent double counting. The first part is only over $i < j$ because we assume an undirected network, and in the DC-SBM, we count stubs within groups but edges between groups.¹⁰

¹⁰If one wants to do link prediction with the DC-SBM, one sets $\text{score}(i, j) = \gamma_i \gamma_j \mathcal{M}_{z_i z_j}$, which is proportional to the probability that i, j are connected under the (fitted) DC-SBM.

Although it may seem rather complicated, we can substantially simplify Eq. (6) by first factoring out and collecting the various constants that depend only on the adjacency matrix,¹¹ and then substituting maximum likelihood estimators for γ and \mathcal{M} :

$$\hat{\gamma}_i = \frac{k_i}{\kappa_{z_i}} \quad \text{and} \quad \hat{\mathcal{M}}_{rs} = \omega_{rs} , \quad (7)$$

where

$$\omega_{rs} = \sum_{i,j=1}^n A_{ij} \delta_{z_i,r} \delta_{z_j,s} \quad \text{and} \quad \kappa_r = \sum_s \omega_{rs} = \sum_i k_i \delta_{z_i,r} \quad (8)$$

are the mixing matrix (counting stubs now, rather than edges) and the group “degrees,” respectively.

Taking the logarithm and then dropping the collected constants, which do not depend on the parameter choice, yields a nicely compact expression that depends only on the edge bundle counts induced by the partition choice z on the graph G :

$$\ln \mathcal{L} = \sum_{rs} \omega_{rs} \ln \frac{\omega_{rs}}{\kappa_r \kappa_s} , \quad (9)$$

where the sum runs over all the elements of the ω_{rs} matrix in order to account for both ends of every edge. Notice that the form of Eq. (9) is similar to that of the SBM’s log-likelihood function in Eq. (5), but where we replace the node-based count n_{rs} with an edge-based count $\kappa_r \kappa_s$.

2.2.1 An example

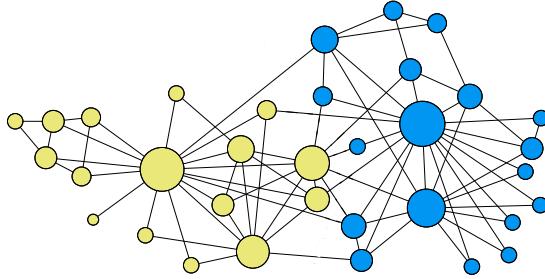
The canonical example for how the SBM and DC-SBM can produce different good partitions of the same network is in their application to the Zachary karate club social network.¹² This network represents the $m = 78$ friendships among $n = 34$ members of a karate club at a U.S. university in the 1970s. During the course of Zachary’s study, the club split into $c = 2$ similar-sized groups, centered around the club’s president and instructor (dashed line in the figures below). Hence, this network is simple graph in which nodes are annotated with a categorical attribute that represents a “social partition,” i.e., which of the two groups the person went with after the split.

The $c = 2$ partition that maximizes the SBM log-likelihood in Eq. (5) places the five highest-degree nodes into one group (including both the president and the instructor), and the other 29 nodes into the other. The tables below give the corresponding e_{rs} and n_{rs} matrices, and log-likelihoods, for both partitions. This partition is very different from the social partition, but it’s not less

¹¹Which we may collect in a term $1/C = \prod_{i < j} A_{ij}! \prod_i 2^{A_{ii}/2} (A_{ii}/2)!$.

¹²See W. W. Zachary, *J. Anthropol. Res.* **33**, 452–473 (1977), <http://www.jstor.org/stable/3629752>.

There’s also this tumblr: <http://networkkarate.tumblr.com>.



meaningful—instead, divides the network into what we might call leaders and followers. Under the SBM, the leader-follower partition is an astounding $\exp(196.29 - 179.39) \approx 10^{7.3}$ more likely to generate the observed network than is the social partition.

e_{rs}/n_{rs}	A (16)	B (18)	e_{rs}/n_{rs}	A (5)	B (29)
A (16)	33/120	10/288	A (5)	5/10	54/145
B (18)	—	35/153	B (29)	—	19/406

$\ln \mathcal{L}_{\text{SBM}} = -196.29$ (social) $\ln \mathcal{L}_{\text{SBM}} = -179.39$ (leader-follower)

In contrast, the $c = 2$ partition that maximizes the DC-SBM log-likelihood in Eq. (9) presents a very different picture: the social partition is now a far better model of the edges than the leader-follower partition. The tables below give the corresponding ω_{rs} matrix and κ_r values, and the log-likelihoods, for the two partitions under this model.

ω_{rs}	A (16)	B (18)	κ_r	ω_{rs}	A (5)	B (29)	κ_r
A (16)	66	10	76	A (5)	10	54	64
B (18)	10	70	80	B (29)	54	38	92

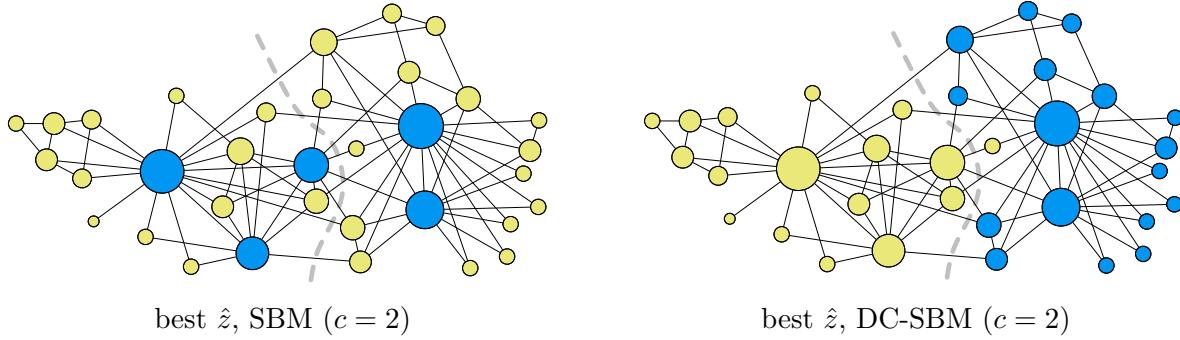
$\ln \mathcal{L}_{\text{DC-SBM}} = -739.43$ (social) $\ln \mathcal{L}_{\text{DC-SBM}} = -772.28$ (leader-follower)

Although close, the social partition is not the globally best DC-SBM partition for this network. That partition moves one node between the groups, and yields a log-likelihood of $\ln \mathcal{L} = -739.39$, which about 4% more likely to generate the observed network than the social partition.¹³

2.3 Finding good partitions

The preceding sections simplified the task of finding a good decomposition of a network into communities. First, we reduced the problem to estimating parameters via statistical inference. Then, we reduced it further to searching over partitions all $\{\vec{z}\}$ to find a partition \vec{z} that maximizes the

¹³Here's a fun, possibly apocryphal story about the one “misclassified” node, which has an equal number of connections to each group. This person apparently had a karate exam coming up soon, and when the club split, they chose to go with the instructor's faction, instead of the president's, in order to be better prepared for the exam.



SBM or DC-SBM log-likelihood of a network.

There are many ways we could perform such a search over partitions. These include powerful methods like Markov chain Monte Carlo, expectation-maximization (the “EM” algorithm), belief propagation (also called “BP”), among many others. Here, we will learn about a **locally greedy heuristic**, which is a kind of generalization of the Kernighan-Lin algorithm developed in 1970 for solving the minimum-cut graph partitioning problem.¹⁴ Like Kernighan-Lin, the locally greedy heuristic constructs a long sequence of partitions $z_i, z_{i+1}, z_{i+2}, \dots$, in a series of phases. This process continues until an entire phase passes in which every partition in that phase is no better than z_* , the partition with the best log-likelihood \mathcal{L}_* prior to the current phase.¹⁵

Below we give pseudocode for the heuristic. In words, the heuristic works as follows.

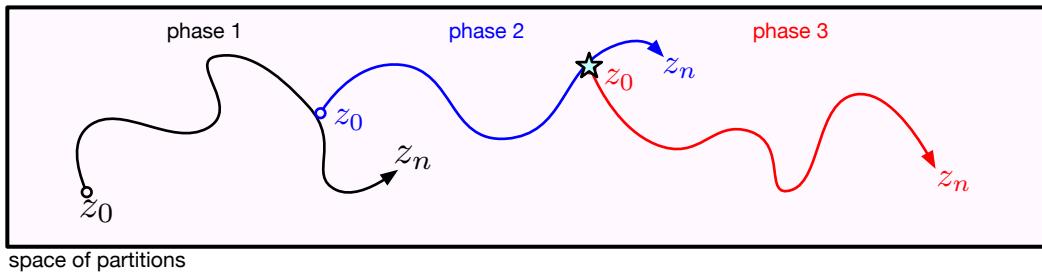
Initially, the heuristic begins with a “random” partition z_0 , in which each node is assigned one of the c labels, uniformly at random, which initializes the first phase.

Within a phase, the heuristic will “move” every node to a different group exactly once, in a greedy fashion. Once a node has been moved, we say it is “frozen,” meaning that it cannot be moved again in this phase. Let z_t denote the partition in this phase with $t \in [0, n]$ frozen nodes. For each value of t , we calculate the change in log-likelihood for all $(c - 1)(n - t)$ possible label changes among the unfrozen nodes. (Do you see why there are that many?) From among this set, we make the greedy choice, changing the label of the node j that yields the most positive (least negative) change in log-likelihood, producing the partition z_{t+1} . Then, we freeze node j . Because each value of t results in freezing one more node, this process produces exactly n new partitions in a phase. Although each move is greedy, a move may result in the log-likelihood *decreasing*; this is fine.

¹⁴Kernighan & Lin, *Bell System Technical Journal* **49**, 291 (1970), <https://archive.org/details/bstj49-2-291>.

¹⁵See Karrer & Newman, *Phys. Rev. E* **83**, 016107 (2011), <https://arxiv.org/abs/1008.3926>.

At the end of a phase, once all n nodes are frozen, the heuristic then re-examines the partitions z_0, z_1, \dots, z_n and chooses the partition z_* that has the highest log-likelihood among them. If z_* 's log-likelihood \mathcal{L}_* is better than the log-likelihood of the phase's initial partition z_0 , the heuristic seeds a new phase with z_* as the initial partition, and unfreezes all nodes; otherwise, if none of the new partitions was better, z_0 must be a local optimum, and the heuristic halts and returns z_0 as a good partition of the network.



The heuristic works by building a sequence of partitions that “walk” through the space of all possible partitions $\{z\}$, starting a random initial location z_0 . The greedy choice ensures that a phase explores only the best one-node change from each partition within a phase. The phase structure ensures that each of the partitions that begins a phase is a strictly better than all the partitions that came before it, until the sequence finds a local optimum (the blue star in the figure above). Within a phase, the sequence can explore a string of “bad” moves—those that decrease the log-likelihood—which may allow it to find a different, better local optimum later in the phase.

We can make this verbal description more concrete by writing it out in pseudocode:

```
function local-greedy( $G$ ):
    1. initialization
        (a) create  $z_0$  by assigning  $z_{0,i} = \text{Uniform}(\{1, \dots, c\})$ ,
        (b) record its log-likelihood  $L_0 = \ln \mathcal{L}(z_0)$ 
        (c) set  $t = 1$  and set all nodes to be “unfrozen”
    2. run a phase
        while  $t < n$  (some node is unfrozen):
            (a)  $z_t = z_{t-1}$  (copy the last partition)
            (b) try all the moves
                for each of the  $(n - t)$  unfrozen nodes:
                    i. let  $i$  be an unfrozen node, and let  $s = z_{t,i}$  ( $i$ 's current group is  $s$ )
```

- ii. for each choice of a group $r \in \{1, \dots, c\} - s$, set $z_{t,i} = r$, compute the corresponding log-likelihood value, and store it as a tuple $(i, \ln \mathcal{L}, r)$
 - iii. of these tuples, keep the one with the most positive log-likelihood $\ln \mathcal{L}$
 - iv. reset $z_{t,i} = s$
 - (c) a **greedy choice**
 of the $(n-t)$ tuples $(i, \ln \mathcal{L}, r)$, one for each unfrozen node, select the one with the most positive $\ln \mathcal{L}$
 - (d) **freeze that node**
 set $z_{t,i} = r$ (move one node)
 set $L_t = \ln \mathcal{L}$ (record its log-likelihood)
 freeze node i
 $t = t + 1$ (next partition in the sequence)
3. **evaluate the phase**
- (a) among the $(n+1)$ log-likelihoods L_0, \dots, L_n , identify the largest L_* and corresponding partition z_*
 - (b) (halt) if $L_* \leq L_0$, then halt and return z_0
 - (c) (new phase) else, set $z_0 = z_*$, $L_0 = L_*$, $t = 1$, unfreeze all nodes, and run a new phase

This locally greedy heuristic is only guaranteed to converge on a local optimum, not the global optimum. A standard approach to mitigating this problem is to simply run the heuristic a number of times; each time will start the optimization from a different random initial partition, and we then take the best output across many of these independent runs.

2.3.1 An example

To illustrate how this locally greedy heuristic works in practice, we'll apply it to a synthetic network with $n = 30$ nodes, connected according to a DC-SBM planted partition of $c = 3$ equal-sized and assortative groups (see Lecture 5), where every node has the same expected degree. In this experiment, we'll track both the log-likelihood (given by Eq. (9)), and a measure of the distance between the current and planted partitions called the variation of information (VI).¹⁶

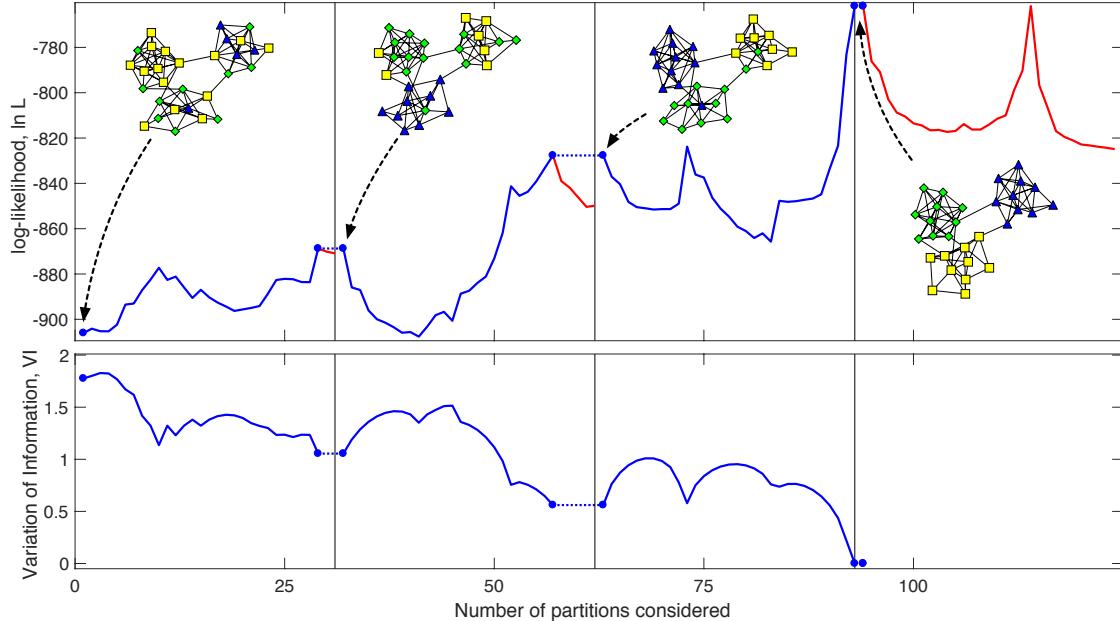
Before showing the results of the experiment, let's first define the variation of information (VI), which measures the “distance” between two partitions X and Y . The VI has the nice property that it obeys the triangle inequality, which means its a distance metric. For our purposes, we'll define X to be the planted partition we used to generate the network, which means the VI provides a complementary measure for quantifying how well the locally greedy heuristic finds good partitions.

For two partitions X and Y , the **variation of information** is defined as

$$\text{VI}(X, Y) = H(X) + H(Y) - 2I(X, Y) \quad (10)$$

$$= -\frac{1}{n} \sum_{r,s} d_{rs} \left[\log\left(\frac{d_{rs}}{p_r}\right) + \log\left(\frac{d_{rs}}{q_s}\right) \right], \quad (11)$$

¹⁶Introduced in Meilă, *J. Multivariate Analysis* **98**, 873–895 (2007), <https://bit.ly/2qdunK1>.



where Eq. (10) is expressed in terms of the standard entropy and mutual information functions $H(\cdot)$ and $I(\cdot, \cdot)$, while Eq. (11) is expressed in terms of simple counts

$$d_{rs} = |X_r \cap Y_s| \quad \text{and} \quad p_r = |X_r| \quad q_r = |Y_r| , \quad (12)$$

where d_{rs} counts the number of nodes that are in group r in partition X and in group s in partition Y , and p_r and q_r are the group sizes in X and Y .¹⁷ In the latter formulation, calculating the VI takes time $\Theta(n)$ from the partitions X and Y . In addition to being a metric, the variation of information has the property $VI(X, X) = 0$ (do you see why?), and is bounded $VI(X, Y) \leq 2 \log c$. In our experiment, we will set X equal to the planted partition, and Y will be the current partition being considered. If the locally greedy heuristic is working, the log-likelihood will go up, and the VI will go down.

Starting with a random initial partition of our $n = 30$ node network, the figure below illustrates the locally greedy heuristic's trajectory through the partition space. Within each phase (demarcated by a black vertical line), the highest log-likelihood partition is marked, and a dotted line connects it to the initialization of the next phase. In the upper panel, the lower-quality partitions explored after the best-in-phase partition is found are indicated in a different color, which reflect subsequent node moves that worsen the quality of the partition. In this run of the algorithm, the heuristic finds

¹⁷As in both the SBM and DC-SBM log-likelihoods, we must define $0^0 = 1$ in practice.

the planted partition at the end of the third phase, at which point the distance drops to $VI = 0$ in the lower panel.¹⁸

2.4 A diversity of algorithms, and choosing the number of groups c

The SBM and DC-SBM are just two of many different algorithms for finding community structure in networks. Among the most popular is to simply maximize the modularity function Q (see previous lecture) over all possible partitions, and there are many others. Some are based on compressing short random walks over the network (InfoMap), others consider entirely other criteria. Most techniques implicitly or explicitly search for assortative community structure, while the SBM and DC-SBM are more flexible, allowing for the discovery of other types of large-scale patterns.

One attractive feature of other community detection algorithms is that many automatically choose the number of groups c . The formulation of the SBM and DC-SBM above do not include a way to do this, but there are many ways to add such a behavior to the model. The key step to doing so is to control, somehow, for the additional model flexibility that comes with increasing c —having an extra group significantly increases the model’s ability to fit the observed network because the mixing matrix embodies $\Theta(c^2)$ parameters. The standard solution is *regularization*, which penalizes a model’s fit to the data in a way that grows with the model’s flexibility. Hence, more flexible models have to fit the data *much* better than less flexible models, in order to be worth it. But, just as there are many different ways to detect communities, there are many different ways to regularize a model.¹⁹

In a 2019 paper on community detection,²⁰ Ghasemian et al. compared the results of applying 16 popular community detection algorithms—all of which had some way of selecting the number of communities c —to 572 real-world networks. A key takeaway from this exercise is that even on the same data, different algorithms can find both substantially different communities and very different numbers of communities. Figure 1 shows the number of communities c found by these algorithms as a function of network size, either the number of nodes n or the number of edges m . Beyond

¹⁸This run was particularly lucky in finding a short path from the random initial partition to the planted partition. More typically, the heuristic gets stuck in a local optimum, which is why the right strategy is to run the heuristic many times, each from a different random initial partition, and take the best output across these runs.

¹⁹The process of choosing one out many models is called *model selection* or *model comparison*. (Not to be confused with a goodness-of-fit test, which asks whether a single model is a good fit to the data.) There are many approaches, each of which comes with a set of philosophical assumptions about *what makes a better model*. The most common is a kind of penalized likelihood function, in which we add to the log-likelihood a penalty function that grows with the “size” of the model. The AIC and BIC are examples of this idea, and MDL is a less well-known example. Another approach is to consider likelihood ratios between smaller and larger models, which can allow one to say that neither model is a better fit than the other; or to integrate the Bayesian posterior distribution over all choices of c . Etc. Which approach a scientist favors often comes down to the scientific culture and history of their field.

²⁰Ghasemian et al., “Evaluating overfit and underfit in models of network community structure.” *IEEE Trans. Knowledge and Data Engineering* **32**, 1722–1735 (2019)

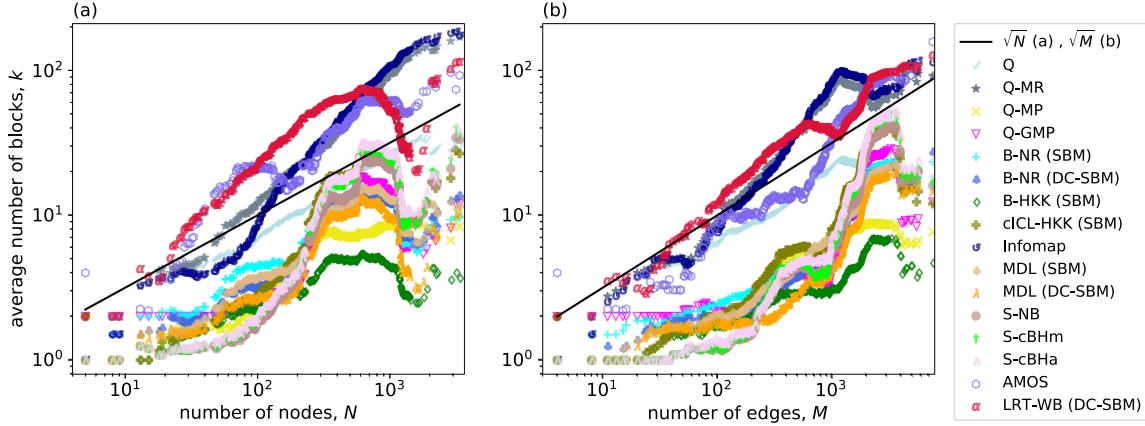


Figure 1: For 16 popular community detection algorithms applied to a diverse corpus of 572 real-world networks, the inferred number of communities as a function of the size of the network, either number of nodes n or number of edges m . In both cases, the number of communities varies substantially depending on which method is used, but the trend of $O(\sqrt{m})$ is a reasonable overall rule-of-thumb. (Reproduced from Ghasemian et al., *TKDE* 2019.))

some algorithms finding 10 times more or less communities than other algorithms, the overall trend is very roughly $O(\sqrt{m})$. That is, generally speaking, the algorithms agree that as a network gets larger, it takes more parameters (communities) to effectively capture its community structure.

The particular communities found, however, correlate with the particular type of algorithm used to find them. Figure 2 below shows the overall pairwise similarity of the partitions found by algorithm \mathcal{A} and those found by algorithm \mathcal{B} . The resulting similarity matrix exhibits strong clustering itself, but now the clustering represents clustering of algorithms.

These results illustrate that both how we define what communities are, and how we go about finding them in a network, have an enormous influence on what communities we find. If there was a fundamental or “true” community partition to be found in a network, we would expect to see far more agreement between the methods. Instead, we find evidence that there are many different, and perhaps roughly equally good partitions, each of which highlights different aspects of the large-scale structure in a network.

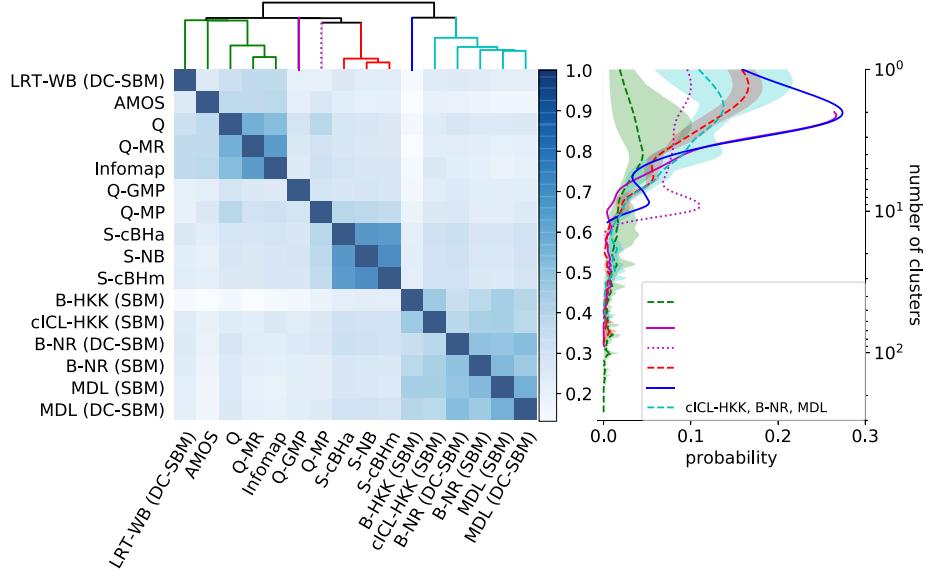


Figure 2: For the same algorithms and data in Fig. 1, the pairwise similarity of the resulting network partition \vec{z} , organized using a hierarchical clustering algorithm. Similar algorithms produce similar clusterings, mostly. (Reproduced from Ghasemian et al., *TKDE* 2019.))

2.5 Caveats and limitations of community detection

Generally, the task of community detection is the network analog of clustering data in a vector space. While clustering aims to find groups of points that clump together, community detection aims to find groups of nodes in a network with similar connectivity patterns.

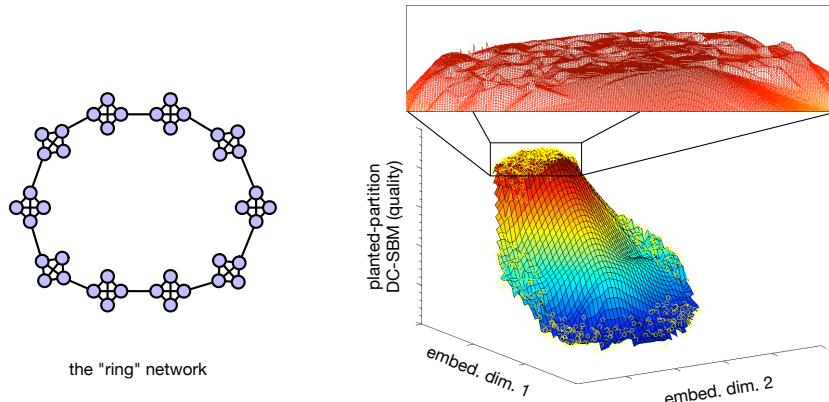
The general problem of clustering is notoriously slippery, and cannot, in fact, be solved universally.²¹ Ultimately, *which* clustering is optimal depends on what we intend to do with it later. Crucially, a similar fact holds for community detection. It's a remarkably useful tool in our toolbox, but its output should always be interpreted with some skepticism. Good community detection methods, like the SBM and the DC-SBM, can be powerful tools for exploratory data analysis, able to uncover a wide variety of such patterns in real networks. They can also be used to predict missing links or missing node attributes, as null models, or to test specific network hypotheses. But, the particular output of any algorithm may not tell us the full story about a network's modular structure, as we

²¹For a truly excellent explanation of this fact, see von Luxburg et al. *J. Mach. Learn. Res.* **27**, 65–79 (2012). <http://proceedings.mlr.press/v27/luxburg12a/luxburg12a.pdf>, and Kleinberg *Adv. Neural Inf. Process. Syst.* 463–470 (2003), <https://www.cs.cornell.edu/home/kleinber/nips15.pdf>.

will see below.

2.5.1 Competitive local optima

A common feature across many of the most popular and most powerful ways of decomposing a network into its constituent communities is that the *objective function* (the likelihood function above) is not convex. If it were, then it would be relatively easy to pursue some kind of “gradient ascent” strategy in order to reliably find the globally best partition under a given model. Instead, these objective functions are often very “rugged”, with a large number of local optima in which search algorithms tend to get stuck. For instance, below is a visualization of the landscape of the likelihood function (right-hand figure) for a planted-partition version of the DC-SBM, applied to a “ring” network of n/d cliques, each of size d , arranged in a ring (left-hand figure). The zoomed-in portion illustrates just how rugged the surface function is—the global maximum is in there somewhere, but one would be hard pressed to tell which peak it is.²²



For this reason, the output of nearly any community detection algorithms cannot be considered the *best possible* partition under a given model. Rather, it is likely to be merely a *pretty good* solution. Furthermore, as the figure’s inset illustrates, many of the local optima are also *competitive* in quality with the global optimum, and mathematically, there tend to be an exponential number of them.²³ In a very real sense, the prevalence of these merely good solutions is why so many different algorithmic strategies for optimizing this or that objective function work relatively well—it’s not hard to find a decent local optimum, and, in truth, finding even one can be useful if we are mainly

²²Reproduced from Good et al., *Phys. Rev. E* **81**, 046106 (2010), <http://arxiv.org/abs/0910.0165>. Embedding dimensions are constructed by computing pairwise VI distances among a set of sampled partitions, and then embedding those points in a 2-dimensional space that preserves those pairwise distances, .

²³That might sound like a big number, but they represent a vanishingly small fraction of the total space of all partitions, which is super-exponential in size.

doing data exploration.

However, if we have a scientific application in mind, e.g., we want to use the output partition \tilde{z} to conduct additional analyses, modeling, or hypothesis testing, then we must ask whether it matters which local optimum an algorithm converges on. If most good solutions are similar to each other in how they decompose the network, then any one will be representative of the set and we need not worry which we found. But, if good solutions are structurally dissimilar to each other, then there's no guarantee that our algorithm of choice has found the one that is scientifically meaningful. What are we to do in such situations? A reasonable strategy is to switch from optimization to sampling to instead obtain a set of good partitions $\{\tilde{z}\}$ (e.g., running the algorithm of Section 2.3 many times) and then check whether they are similar or different before proceeding with any downstream analyses.

2.5.2 Detectability limits

In addition to not necessarily finding the best partition possible, community detection is not guaranteed to find the actual communities in a network, if the boundaries between them are insufficiently sharp. This is true not just on real-world networks, where the underlying community structure is not known, but also true on synthetic networks with known structure, and it holds for all algorithms.

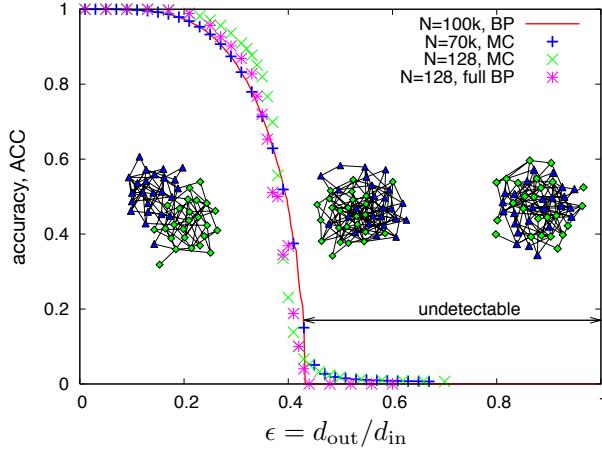
To illustrate this fact, recall the planted partition model for the SBM, in which we have $c=2$ equal-sized communities and where the sharpness of their boundary is governed by a single parameter ϵ . In seminal work, Decelle et al. (2011) proved mathematically that for this simple setting, there exists a regime of ϵ beyond which *no algorithm* can successfully recover the planted partition, and that this “undetectable” regime emerges well before the communities are formally equivalent. The figure below, from their paper, shows a sharp phase transition in the recoverability of the planted communities. (In their formulation, ϵ is defined as the ratio of out-of-group edges to in-group edges, but the model is otherwise the same as the one we saw in Lecture 5.)²⁴

These “detectability limits” have since been extended to networks that vary over time, to multi-layer networks, and to networks with node metadata that we might exploit to identify underlying communities (e.g., if the metadata labels correlate with the community structure pattern of edges). The implication of these limits is that we should be mildly skeptical of the output of any community detection algorithm, as it may not have been able to detect the actual underlying structure.

2.5.3 No ground truths and No Free Lunches

Beyond not finding the best partition and not necessarily detecting the underlying communities, there is no guarantee that community detection can recover the correct “ground truth” and, on

²⁴Reproduced in part from Decelle, et al., *Phys. Rev. E* **84**, 066106 (2011), <https://arxiv.org/abs/1109.3041>.



average, there is no “best” community detection algorithm.

Many networks have nodes annotated with metadata \vec{x} , and it’s common to treat the partition these labels induce as being the thing that a good community detection algorithm should find, using only the edges. But, node metadata are not really “ground truth” for network communities; they’re just more data. It’s possible, sometimes even likely, that \vec{x} correlates with the community structure of the network, e.g., if the probability that an edge (i, j) exists is conditioned in some way on x_i and x_j . But without a clear model of that relationship, metadata always have an uncertain relationship with ground truth.

Think of it this way: if an algorithm fails to find a good division \vec{z} that correlates with some node metadata \vec{x} , there are four possible reasons:

- (i) (no correlation) the metadata are irrelevant to the network’s structure,
- (ii) (they’re different) the detected communities and the metadata capture different aspects of the networks structure,
- (iii) (no correlation) the network contains no communities (as in a non-modular random graph) or its communities fall below the detectability threshold, or
- (iv) the community detection algorithm performed poorly.

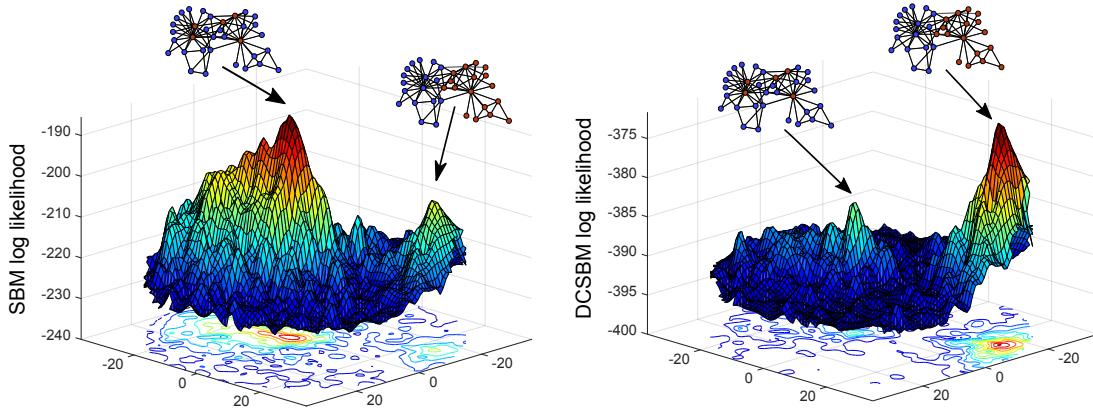
In practice, we often cannot tell which of these cases is more or less likely, and the presence of measurement errors or noise in the network’s connectivity or metadata only increases that uncertainty.

There is no ground truth. Formally, community detection is a kind of inverse problem. We imagine generating a network G by applying a generative process g to a “ground truth” set of communities \mathcal{T} , i.e., we compute $G = g(\mathcal{T})$. In community detection, we want an algorithm g^{-1} that

can reliably identify the \mathcal{T} that produced G , i.e., community detection aims to compute $\mathcal{T} = g^{-1}(G)$.

Now suppose we have a network G that can be produced by multiple, distinct generative processes, each with its own ground truth. That is, we have $G = g_1(\mathcal{T}_1) = g_2(\mathcal{T}_2)$, but where $(g_1, \mathcal{T}_1) \neq (g_2, \mathcal{T}_2)$. If this is possible, and it is, then no community detection algorithm method can uniquely solve the problem for every network, or even just for one arbitrary network. In other words, the community detection problem is not invertible because there is not a bijection between networks G and ground truths \mathcal{T} . Or, more succinctly, there is no (identifiable) ground truth.

We have already seen one such network: the Zachary karate club can be generated using either the SBM with the leader-follower partition, or the DS-SBM with something like the social partition. The figure below illustrates this point using the likelihood surfaces for both models, in which we have labeled the local maxima for both partitions.



There are No Free Lunches. In machine learning, the famous No Free Lunch (NFL) theorem proves that, subject to mild constraints on the evaluation criteria, no learning algorithm has an *a priori* advantage over any other across all possible inputs. In other words, for a set of cases that a particular method \mathcal{A}_1 outperforms \mathcal{A}_2 , the NFL theorem implies that there must exist a set of cases where \mathcal{A}_2 outperforms \mathcal{A}_1 in such a way that, on average, no algorithm performs better than any other. This theorem also holds for the task of community detection.²⁵

²⁵See Peel et al., *Science Advances*, **3**(5), e1602548 (2017). Formally, the NFL holds for community detection algorithms evaluated using the standard measure of adjusted mutual information (AMI), which is closely related to the variation of information in Eq. (10). Notably, there do exist areas of statistical learning where the NFL doesn't hold. Identifying these narrow settings is an area of active research.

The lack of generally better, or worse, performance by different community detection algorithms does not mean community detection itself is pointless. Rather, the theorem implies that if the inputs of interest represent a restricted subset of cases, e.g., all networks with realistic structure, or specific groups of networks like food webs or online social networks, then there may indeed be a method that out-performs others *within the confines of that subset*. In other words, an algorithm can be better on some groups of inputs only if it is worse on enough other inputs to compensate. Hence, if we intend to use community detection algorithms in practice, we are more likely to obtain good solutions if we match our beliefs about the underlying data generating process g with the particular assumptions of an algorithm \mathcal{A} . That is, we should favor specialized algorithms whenever we know something about what kind of network we're analyzing.

3 Supplemental readings

1. Chapters 14.1 and 14.4 in *Networks*