

COMPUTATIONAL COMPLEXITY, PHASE TRANSITIONS, AND MESSAGE-PASSING FOR COMMUNITY DETECTION

AURÉLIEN DECELLE, JANINA HÜTTEL, ALAA SAADE, CRISTOPHER MOORE

These are notes from the lecture of Cristopher Moore given at the autumn school “Statistical Physics, Optimization, Inference, and Message-Passing Algorithms”, that took place in Les Houches, France from Monday September 30th, 2013, till Friday October 11th, 2013. The school was organized by Florent Krzakala from UPMC & ENS Paris, Federico Ricci-Tersenghi from La Sapienza Roma, Lenka Zdeborová from CEA Saclay & CNRS, and Riccardo Zecchina from Politecnico Torino. Much of this material is covered in Chapters 4, 5, 13, and 14 of [5].

CONTENTS

1. Computational complexity	1
2. Hardness: P, NP or EXP?	5
2.1. Examples	5
3. Random Graphs	8
3.1. Erdős-Rényi Random Graphs	8
3.2. The Giant Component	9
3.3. Giant component and configuration model	11
3.4. The k -core	12
4. Random k -SAT	13
4.1. Easy Upper Bound	14
4.2. Lower Bounds from Differential Equations and the Analysis of Algorithms	14
4.3. Lower bounds from the second moment method	16
5. Community detection	19
5.1. The stochastic block model	19
5.2. Spectral methods	26
6. Appendix	28
6.1. Definition of Perfect Matching	28
6.2. Definition of Max Flow	28
6.3. Definition of k -SAT	28
References	29

1. COMPUTATIONAL COMPLEXITY

Computational complexity is a branch of complex systems which has the advantages of being very well defined and presenting a lot of rigorous results. It allows qualitative distinctions between different kinds of “hardness” in computational problems, such as the difference between polynomial and exponential time. Amongst the historical examples of combinatorial problems, the “Bridges of Königsberg” problem

is a famous one solved by Euler, where the premises of computational complexity can be already sketched. This problem deals with finding a path through all parts of the city of Königsberg using each bridge only once. It was solved quickly by Euler who, by using the dual graph, showed that a solution exists only if the degree of each node is even (except from the starting and ending node). This perspective is a profound change of how to view the problem since by Euler's argument, the verification of the existence of such a path can be done very quickly (in a polynomial time since it is enough to check the degree of all edges) as opposed to an exhaustive search through all possible paths.

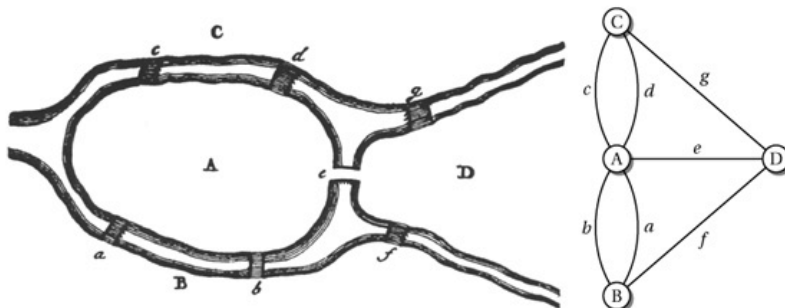


FIGURE 1. from [5]

Hamilton came later with a problem which looks similar. It is about the possibility of finding a path in a graph that visits each vertex exactly once. Although it looks similar, the fundamental difference with the previous problem is the impossibility (as far as we know) to verify the existence of such a path by a simple (and quick) argument. Therefore finding a Hamiltonian path seems to be possible only by an exhaustive search (with an exponential number of possibilities). On the other hand, checking if a given path is a Hamiltonian cycle or not is an easy problem. One simply goes through the path, node by node, to check whether a path is a solution of the problem or not.

To go from these first simple examples to the theory of computational complexity, we first need a few definitions.

Definition 1 (Problem). *A problem is given by defining an input (or instance), for example in the case of the Hamilton cycle, the input is: “a graph G ”, and a question on that input, “does there exist a Hamiltonian path on this graph? (yes or no)”*

Definition 2 (The class P). *P is the class of problems that can be “solved” in polynomial time $\text{poly}(x) = x^c$ for some constant c and where x is the size of the instance. In the example above, the Euler Path belongs to P as, given a graph G , we can answer the question: “does there exist an Eulerian path on G ?” in polynomial time.*

Definition 3 (The class NP). *NP is a class of decision problems, which have a yes/no answer. If the answer is “yes”, there is a proof that the answer is “yes” that can be checked in polynomial time.*

Let's look at the Hamiltonian case. Consider the question: “Is there a Hamiltonian path?” We do not know any polynomial algorithm answer, but if someone provides us

such a path, it is a proof that the answer is yes, and it can be checked in polynomial time.

Note that the definition of NP problems is not symmetrical. If we claim “Prove that there is not a Hamiltonian path!”, it does not seem easy at all to prove that such a path does not exist: as far as we know, in general all possible paths have to be examined.

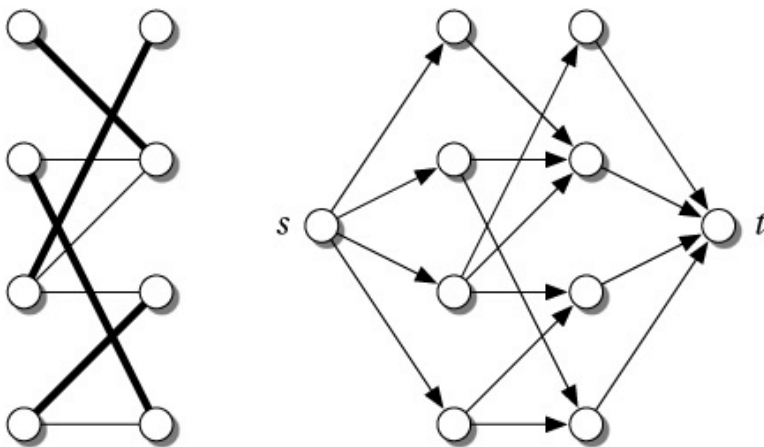
Definition 4 (The class NP-complete). *A problem is called NP-complete if it “somehow” encompasses any problem in NP.*

To be more precise we need to introduce the notion of “reduction” of a problem into another. A problem (A) can be reduced to (B) if, given an instance x of A, there exists a function f so that $f(x)$ is an instance of (B) and f is in P. In addition, $f(x)$ is a “yes”-instance of (B) if and only if x is a “yes”-instance of (A). Such a reduction is written as $A \leq B$. In practice, it means that if I have an algorithm to solve (B), I can use it in some way to solve (A). It also implies that if (B) is in P, (A) is in P since f is a polynomial function. And as can be understood by the definition, if any problem of the NP-complete class is in P, then $P=NP$.

Here’s an example of a reduction between two problems in P.

Example 1 (Reduction— Perfect bipartite matching \leq max-flow).

The problem of telling whether a bipartite graph has a perfect matching can be mapped on the max-flow problem by adding the nodes s and t to the bipartite graph as illustrated below. The node s is linked to the left part of the bipartite graph by directed edges, at the center, the edges are replaced by directed edges going from left to right and the right part is linked to the node t by directed edges. Each edge has a weight of one. The intuition is that, if there is a max-flow of value m , then there is a matching consisting of m edges. By the direction of the arrow and the weight of the edges in the max flow problem, the principle of the proof can be understood.



It can also be demonstrated that max-flow can be reduced to the weighted min-cut problem and vice versa.

The reduction of one problem into another implies also that, if $A \leq B$ then, if A is not in P then B is not in P (which justified the symbol “less or equal”). Now we can define more clearly the class of NP-complete problems:

Definition 5 (NP-completeness). *B is NP-complete if*

- $B \in NP$
- $A \leq B, \forall A \text{ in } NP$

This definition implies that, if B is in P then $P=NP$, so it is enough to find a polynomial solution to an NP-complete problem to demonstrate the equality. It can be shown that the problem of finding a Hamiltonian path is NP-complete. We introduce now some examples of NP-complete problems.

Example 2 (NP-complete problems).

- program SAT: A program Π and a time t given in unary. Does there exist an instance x such that $\Pi(x) = \text{“yes”}$ (in t steps or less)?

This problem is NP-complete because it reproduces the exact structure of any NP problem: a general “yes/no” question that should be checkable in a time at most of size t . Remark: t is given in unary (i.e. a string of t 1s), instead of in binary. Otherwise t would be exponential as a function of the size of the input, measured in bits.

- circuit SAT: Boolean circuits are a set of source nodes (bits) connected to logical gates (AND, OR, NOT) on a directed acyclic graph. The output consists of a set of sink nodes. A circuit is called satisfiable if there exists an assignment of the input nodes such that the output is true. The problem can be phrased as the following:

Input: A circuit C

Question: Does there exist an x so that $C(x)$ is true?

This problem is NP-complete because the program SAT can be reduced to a circuit-SAT instance. This can be understood as any algorithm can be stated in term of a circuit-SAT form (it is enough to think that a computer is only an ensemble of nodes and logical gates, and that for-loops can be unfolded to produce more layers of these gates). Therefore, any program can be mapped to a circuit SAT instance and it is an NP-complete problem.

- 3-SAT: the 3-SAT problem is a constraint satisfaction problem (CSP) which is defined by an ensemble of variable nodes (true/false). Then a 3-SAT instance is a formula, which consists of a set of clauses all connected by an OR operator. A clause contains three variables linked by an AND operator and each variable can be negated or not. We can prove the following: $\text{circuit-SAT} \leq 3\text{-SAT}$.

This reduction can be proven by first demonstrating that any circuit-SAT can be mapped into a SAT formula, and then showing that any SAT formula can be mapped in a 3-SAT formula.

Circuit-SAT to SAT formulas: To map a circuit to a formula, we add variables for the internal values on the wires, and then transform each logical gate into a set of SAT formulas. Then the output of the gate is used as a new variable that is propagated throughout the rest of the circuit.

– 1) **AND**: $y = x_1 \wedge x_2 \Leftrightarrow (x_1 \vee \bar{y}) \wedge (x_2 \vee \bar{y}) \wedge (\bar{x}_1 \vee \bar{x}_2 \vee y)$.

– 2) **OR**: $y = x_1 \vee x_2 \Leftrightarrow (\bar{x}_1 \vee y) \wedge (\bar{x}_2 \vee y) \wedge (x_1 \vee x_2 \vee \bar{y})$.

– 3) **NOT**: $y = \bar{x} \Leftrightarrow (x \vee y) \wedge (\bar{x} \vee \bar{y})$

Any k -SAT clause can be written as a 3-SAT formula. We refer to [5] for the case $k > 3$ and show the cases $k = 1, 2$ here so that the reduction from a circuit-SAT to 3-SAT is complete.

- A single variable $x \Leftrightarrow (x \vee z_1 \vee z_2) \wedge (x \vee z_1 \vee \bar{z}_2) \wedge (x \vee \bar{z}_1 \vee z_2) \wedge (x \vee \bar{z}_1 \vee \bar{z}_2)$. In the 3-SAT formula above, whatever the values of z_1 and z_2 are, the expression is satisfiable if x is true and unsatisfiable otherwise.
- Two variables $(x \vee y) \Leftrightarrow (x \vee y \vee z) \wedge (x \vee y \vee \bar{z})$. Same as above, whatever the value z takes, the expression is satisfiable only if x OR y is true.

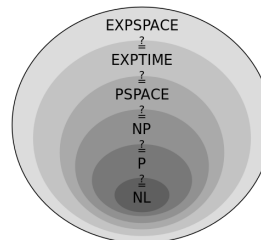
With these two pieces, it is now easy to map any instance of a circuit SAT toward a 3-SAT problem and therefore $\text{Circuit SAT} \leq 3\text{-SAT}$.

- **NAE-SAT (NonAllEqual-SAT):** *This problem is very similar to a SAT problem with the small difference that all the variables in one clause cannot be all true (or false) at the same time. We can observe that any NAE-SAT solution is symmetric (unlike the k-SAT case). For this problem we have the reduction $3\text{-SAT} \leq \text{NAE-3-SAT}$. To prove this reduction, it is easier to prove first that $3\text{-SAT} \leq \text{NAE-4-SAT}$ and then to show that $\text{NAE-3-SAT} \leq \text{NAE-4-SAT}$. The first point can be proven by adding a variable to each clause of a 3-SAT instance. So any formula $\phi = (x_1 \vee x_2 \vee x_3) \wedge \dots$ becomes $\phi' = (x_1, x_2, x_3, s) \wedge \dots$, with the same variable s added to every clause. If ϕ is satisfiable, we can satisfy ϕ' by setting s to false. Now take a SAT-instance of ϕ' . Because the problem is symmetric, the symmetric version of that instance will exist, and amongst those two one where s is false. Therefore the rest of the variables satisfy the 3-SAT formula as well. Then it remains to convert any NAE-4-SAT formula into a NAE-3-SAT formula. This can be done easily by using another variable z and seeing that $(x_i, x_j, x_k, x_l) = (x_i, x_j, z) \wedge (\bar{z}, x_k, x_l)$.*

2. HARDNESS: P, NP OR EXP?

We already discussed the hardness of a problem in the previous lecture. Now we want to discuss counting problems. We already know that problems in NP ask whether an object with a certain property exists. We now define #P, pronounced “sharp P”, as the class of problems that ask how many such objects exist. As in NP, we require that this property can be checked in polynomial time.

Definition 6. #P is the class of functions $A(x)$ of the form $A(x) = \#\{w \mid B(x, w)\}$, where $B(x, w)$ is a property that can be checked in polynomial time, and where $|w| = \text{poly}(|x|)$ for all w such that $B(x, w)$ holds.

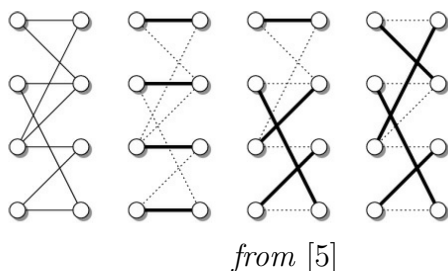


from [5]

2.1. Examples.

Example 3 (Perfect Matching). *The problem of deciding if there exists a perfect matching in a given graph G , is, as we already mentioned, in P. On the other hand the problem of counting these perfect matchings is in #P. Suppose we have a bipartite graph G with $2n$ vertices. We can represent it as a $(n \times n)$ matrix A , where $A_{ij} = 1$ if the i th vertex on the left is connected to the j th vertex on the right, and $A_{ij} = 0$*

otherwise.



As an example we can write down the matrix A for the first graph on the left:

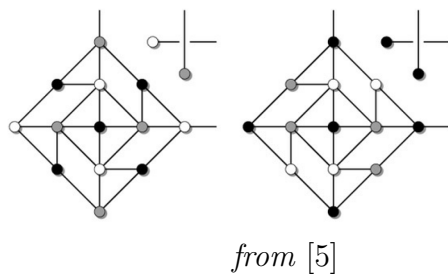
$$A = \begin{pmatrix} 1 & 1 & 0 & 0 \\ 0 & 1 & 0 & 1 \\ 1 & 1 & 1 & 0 \\ 0 & 0 & 1 & 1 \end{pmatrix}$$

But how can we express the number of perfect matching in terms of A ? Each perfect matching is a permutation that maps the vertices on the left to their partners on the right. Each permutation σ corresponds to a matching if and only if for each i there is an edge from i to $\sigma(i)$, i.e., if $A_{i\sigma(i)} = 1$. Therefore, the number of matchings is given by the following quantity, which is called the **permanent of A** :

$$\text{perm}(A) = \sum_{\sigma \in S_n} \prod_{i=1}^n A_{i\sigma(i)}$$

Note that this is just like the determinant, except that it doesn't have the parity $(-1)^\sigma$. Ironically, this makes it harder to compute.

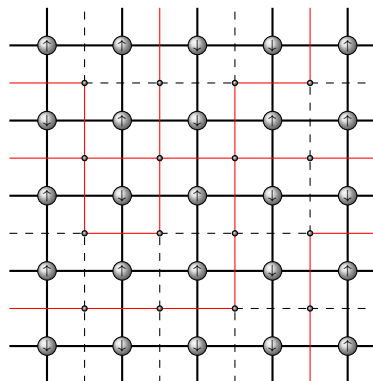
Example 4 (Graph 3-colorability). We want to give a hint on how it is possible to reduce graph 3-colorability to planar graph 3-colorability. In other words how we can transform an arbitrary graph G to a planar graph G' , such that G' is 3-colorable if and only if G is. We can easily see that the reverse reduction is trivial. Thus planar graph 3-colorability is just as hard as graph 3-colorability in general, and the two problems have the same complexity.



This crossover gadget allows us to convert an arbitrary graph into a planar one while preserving the property of 3-colorability. The gadget can be colored in two different ways up to symmetry, depending on whether the transmitted colors are different (left) or the same (right).

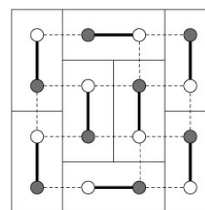
Example 5 (2-dimensional Ising model). We want to think about the Ising model as a 2-dimensional lattice with n spins which can take the values ± 1 . Now we can consider the dual lattice, whose edges cross the edges of our Ising model lattice.

An edge in the dual lattice is colored in red if the two spins that are connected through its crossing edge have different values. It is clear that one vertex of the dual lattice is always adjacent to an even number of red edges. Now we can think about a ground state of the Ising model as a minimum weight perfect matching in a decorated version of the dual graph (see [5] for the gadgets and details). We know that every spin assignment corresponds to a single coloring of the dual lattice.

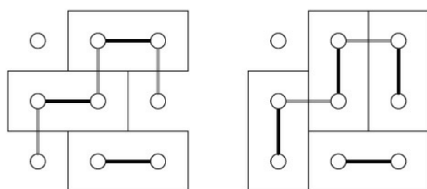


Example 6 (Domino tiling).

Our next task is to reduce domino tiling (where we want to cover a $n \times n$ -chess board C with dominos) to finding a perfect matching in a bipartite graph. We therefore define a graph G where each vertex corresponds to a square in C and is connected to the four vertices corresponding to the neighboring squares. A domino tiling of C is just a perfect matching of G . So domino tiling is in P .



from [5]



from [5]

We can improve a partial domino tiling, or equivalently a partial matching, by flipping the edges along an alternating path. Now if we want to compute the number of domino tilings for the matrix C we have to compute $\text{perm}(C)$ as above.

Turning Permanents into Determinants Again we have a planar, bipartite graph G with n vertices of each type. Now we color the two types of vertices in black and white (see above). As in example 3 we define a $n \times n$ -matrix A . We know already that each perfect matching corresponds to a permutation $\sigma \in S_n$, which maps each black vertex to its white partner. We saw in example 3 that $\text{perm}(A)$ counts all of these permutations. But the determinant of A counts them weighted by their parities, $\det(A) = \sum_{\text{matchings } \sigma} (-1)^\sigma$. The idea is to compensate the parity weights of the determinant in order to obtain the correct count of perfect matching. To do this, we place weights $w_{ij} = \pm 1$ on the edges of G . It defines a quantity $\tilde{A}_{ij} = w_{ij}$. Now, each matching σ of \tilde{A} has a weight

$$w(\sigma) = \prod_i w_{i\sigma(i)}$$

and the determinant of \tilde{A} is given by

$$\det(\tilde{A}) = \sum_{\text{matchings } \sigma} (-1)^\sigma w(\sigma)$$

Now we would like to write the permanent of A as the determinant of \tilde{A} . To do this, we should choose the weights w_{ij} such that $(-1)^\sigma w(\sigma)$ has the same sign for all σ . It means that the matching σ should change the weight $w(\sigma)$ by a factor -1 when its parity changes. Then we would have that $|\det(\tilde{A})| = \text{perm}(A) = \#$ of perfect matchings. The whole trick is to find a proper set of weights. For instance, in the particular case of the chess board (see the figure above), one can decide to put a weight i on all vertical edges. Then, by changing two horizontal dimers to two vertical ones, the weight changes by $i^2 = -1$, and so does the parity. It should be emphasized that a set of such weights can always be found for planar graphs. In addition, an analogous approach allows us to solve the 2-d Ising model (using a trick to count perfect matchings on planar graph that are not bipartite). For further details we refer to [5]. We finally notice that some problems (like Graph 3-Coloring) are just as hard when limited to the planar case, while others (like 4-Coloring or counting matchings) get much easier.

We now turn to random graphs, and give some of their basic properties. These properties will be useful when discussing some other problems in more details.

3. RANDOM GRAPHS

3.1. Erdős-Rényi Random Graphs. First we want to give the definition of a random graph.

Definition 7. A graph $G(n, p)$ with n vertices and edge probability p is one model of a **random graph**. Edge probability p means that between any two vertices v and w there exists an edge with probability p . For every pair of vertices this event is independent from each other pair. An **Erdős-Rényi graph** $G(n, m)$ has n vertices and is uniformly chosen from all graphs with exactly m edges.

Our first task is to compute the expected degree of one vertex and the expected number of edges in the $G(n, p)$ model. We want to investigate the sparse case $pn = c$ where c is a constant independent of n . We use \mathbb{E} to denote the expectation.

$$(3.1) \quad \mathbb{E}[m] = p \binom{n}{2} \approx \frac{pn^2}{2} = \frac{cn}{2}$$

$$(3.2) \quad \mathbb{E}[deg] = p(n-1) \approx pn = c$$

Now we can compute some other interesting quantities in our model. We compute the expected number of triangles and bicycles in $G(n, p)$,

$$(3.3) \quad \mathbb{E}[\#\Delta] = \binom{n}{3} p^3 \approx \frac{c^3}{6}$$

$$(3.4) \quad \mathbb{E}[\#\text{ bicycles}] = \text{const} \times \binom{n}{k} p^{k+1} \approx n^k \frac{c^{k+1}}{n^{k+1}} \xrightarrow{n \rightarrow \infty} 0$$

Bicycles are subgraphs containing k vertices and $k + 1$ edges. Thus they contain two loops, connected by a path or by an edge belonging to both loops. The constant depends on k but not on n . For any fixed k there are probably no bicycles as $n \rightarrow \infty$, so most vertices have a treelike neighborhood.

Our next task is to compute the probability that a vertex v has degree k .

$$(3.5) \quad \mathbb{P}[\deg(v) = k] = \binom{n-1}{k} p^k (1-p)^{n-1-k} \approx \frac{n^k c^k}{k! n^k} \left(1 - \frac{c}{n}\right)^n = \frac{e^{-c} c^k}{k!}$$

This is called the **degree distribution**. As n goes to infinity, it becomes a Poisson distribution with mean c .

3.2. The Giant Component. In this section we want to describe the most basic phase transition occurring in $G(n, p)$ for $pn = c$. For very small c , G consists of small components isolated from each other, and almost all of them are trees. For c bigger and bigger we add more edges and these trees grow and connect with each other. Suddenly, at $c = 1$, many of them come together to form a giant component. We want to observe this phase transition later on in this chapter. First we study the expected component size of a vertex v . We start at this vertex and go further to its neighbors and their neighbors and so on. We can think of those vertices as v 's children and the descendants with distance k from our starting vertex v are the so called k^{th} **generation** of v . By doing that we can develop the whole component of v . For large n we can assume that every child of v has again Poisson(c) children. That is, we can approximate the process of exploring v 's component as a branching process. We only have to count the descendants of v to get the component size.

$$(3.6) \quad \mathbb{E}[\text{Component size of } v] = 1 + c + c^2 + c^3 + \dots = \frac{1}{1-c} \text{ for } c < 1$$

For $c < 1$ this sum converges and the expected component size is finite. For $c > 1$, on the other hand, the sum diverges and v has, in expectation, an infinite number of descendants; in fact, the number of descendants is infinite with positive probability. At this point the branching process is no longer an accurate model to compute the component size of v . For further computations we call γ the fraction of all vertices that lies in the giant component of $G(n, p)$. When $c > 1$, with high probability there is a unique giant component.

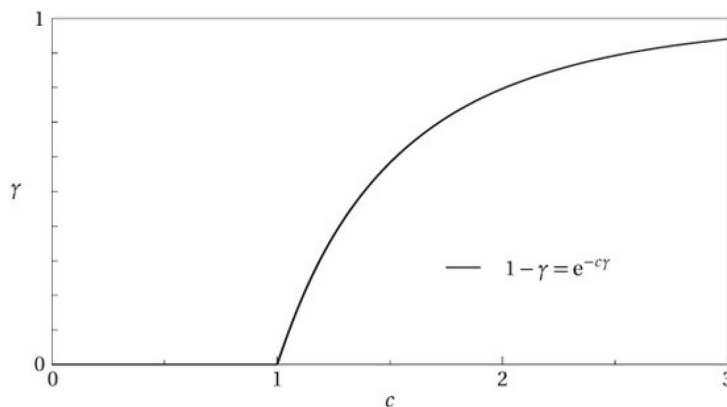
We will see two different methods to estimate the value of γ .

3.2.1. Method #1. The probability that a vertex v is not part of the giant component is the sum over all k of the probability that v has degree k and none of its children are part of the giant component.

$$(3.7) \quad 1 - \gamma = \sum_k \mathbb{P}[\deg(v) = k] (1 - \gamma)^k = \sum_k k \frac{e^{-c} c^k}{k!} (1 - \gamma)^k = e^{-c} e^{c(1-\gamma)} = e^{-c\gamma}$$

So we get a transcendental equation for the size of the giant component:

$$(3.8) \quad \gamma = 1 - e^{-c\gamma}$$



from [5]

3.2.2. *Method #2.* We compute γ again but now with a system of differential equations. We start again with a vertex v and explore its connected component, now by using an algorithm. At each point in time, a vertex is labeled **Reached**, **Boundary**, or **Untouched**.

- “Reached” means that a vertex lies in v ’s component and its neighborhood has been explored
- “Boundary” means that a vertex lies in v ’s component but its neighbors have not been explored yet
- “Untouched” means that it is not yet known if the vertex is in v ’s component.

Algorithm 1 Cluster expansion

Require: a vertex v

Ensure: a connected-graph

- 1: label v Boundary
 - 2: label all other vertices Untouched
 - 3: **while** there are Boundary vertices **do**
 - 4: choose a Boundary vertex
 - 5: label each of v ’s Untouched neighbors Boundary
 - 6: label v Reached
 - 7: **end while**
-

This algorithm explores G one vertex at a time, until there are no Boundary vertices left and v ’s component has been labeled Reached. Let R , B , and U denote the number of vertices of each type at a given point in time. At each step of the algorithm, R increases by 1, and the expected change in B and U is

$$(3.9) \quad \Delta R = 1$$

$$(3.10) \quad \mathbb{E}[\Delta B] = -1 + pU = \frac{c}{n}U - 1$$

$$(3.11) \quad \mathbb{E}[\Delta U] = -pU = -\frac{c}{n}U$$

We have an expected change from U to B of pU , since each Untouched vertex is connected to v with probability p . By changing the chosen vertex v from Boundary

to Reached we get the term -1 . When n is large we can rescale these stochastic difference equations, so they become a system of differential equations:

$$(3.12) \quad \frac{dr}{dt} = 1$$

$$(3.13) \quad \frac{db}{dt} = cu - 1$$

$$(3.14) \quad \frac{du}{dt} = -cu$$

We solve these differential equations with the initial conditions $b(0) = 0$ and $u(0) = 1$:

$$(3.15) \quad u(t) = e^{-ct} \text{ and } b(t) = 1 - t - e^{-ct}$$

The fraction γ of vertices in the giant component is the value of t at which no Boundary vertices remain and the algorithm stops. This gives the same equation as before,

$$(3.16) \quad b(\gamma) = 1 - \gamma - e^{-c\gamma} = 0$$

Remark: The differential equation approach to the size of the giant component in random graphs is similar to the dynamics for the SIR model, where p is the transmission rate and γ is the fraction of the population that eventually becomes infected.

Another interesting quantity we can consider is the expected number of components/trees of G with k vertices. We approximate this with the number of trees:

$$(3.17) \quad \mathbb{E}[\text{Components/trees of size } k] = \binom{n}{k} p^{k-1} (1-p)^{k-2} (1-p)^{kn-k}$$

Here k^{k-2} comes from Cayley's formula for the number of labeled trees with k vertices. For $p = c/n$ and $c = 1$, the expression above becomes a power law by applying Stirling's formula:

$$(3.18) \quad \mathbb{E}[\text{Components/trees of size } k] \approx \frac{1}{\sqrt{2\pi}} \frac{n}{k^{5/2}} (1 + O(k^2/n))$$

3.3. Giant component and configuration model. One problem of random graphs is that the degree distribution converges toward the Poisson distribution which is unrealistic in many concrete examples (social networks, biological networks, etc.). The configuration model can be used to remedy this problem. This model deals with a sequence of nodes and degrees. The nodes are chosen randomly before being connected by an edge. During this process, each node is chosen with a probability proportional to the number of unmatched edges it has. Equivalently, each vertex with degree d has d "stubs" or half-edges. The total number of stubs is $2m$ where m is the number of edges, and we choose a uniformly random matching of these stubs with each other.

This procedure may create some self-loops or multiple edges, so strictly speaking this model produces random multigraphs. However, for reasonable degree distributions (with bounded mean and variance) the resulting graph is simple with constant probability.

The next question is to determine where, in this new setting, the giant component appears. As long as the graph is sparse (i.e. the average degree is constant) the

Algorithm 2 Configurational model

Require: degree sequence**Ensure:** a graph

- 1: **while** there remain unmatched edges **do**
 - 2: choose two unmatched edges uniformly and randomly
 - 3: put an edge between them
 - 4: **end while**
-

appearance of the giant component can be analyzed by a branching process as before. When we follow a link to a vertex of degree k , there are $k - 1$ “children” consisting of new edges that come out of it. Thus the average branching ratio is

$$(3.19) \quad \lambda = \sum_k (k-1) \frac{ka_k}{\sum_j ja_j} = \frac{\mathbb{E}[k(k-1)]}{\mathbb{E}[k]}$$

where a_k is the fraction of nodes in the degree sequence that have degree k . A giant component appears when $\lambda > 1$, or equivalently

$$(3.20) \quad \mathbb{E}[k^2] > 2\mathbb{E}[k]$$

Thus both the first and second moments of the degree distribution matter.

3.4. The k -core.

Definition 8. *The k -core of a graph G is the largest sub-graph where each vertex has minimal degree k . Equivalently, it is the graph remaining after removing vertices with $d < k$ neighbors (iteratively).*

This object is related to the k -colorable property of a graph: if there is no k -core, the graph is k -colorable [9] (note that the converse is not necessarily true). We will again use a branching process to characterize the k -core: at which α it appears and what fraction of the system is in it at that point. Unlike the giant component, the k -core appears suddenly for $k \geq 3$: that is, it includes a constant fraction of vertices when it first appears.

Branching process. We consider only $k > 2$ since for any $c > 0$, $G(n, p = c/n)$ typically contains a 2-core as a loop of size $\mathcal{O}(\log n)$ or even a triangle is present with constant probability. To describe the problem with a branching process, we need to consider first a root node v , again treating its neighbors as its children, their neighbors as its grandchildren, and so on. Let us (recursively) say that a node in this tree is *well-connected* if it has at least $k - 1$ well-connected children. Such a node will survive the process that deletes nodes with degree less than k , if we think of this process as moving up the tree from the leaves.

If the fraction of well-connected children is q , the number of well-connected children a given node has is Poisson distributed with mean cq . The probability that the number of well-connected children is less than k is then given by

$$(3.21) \quad Q_k = \sum_{j=0}^{k-1} \frac{e^{-cq} (cq)^j}{j!}.$$

This gives us the fixed-point equation

$$(3.22) \quad q = 1 - Q_{k-1}.$$

With this equation, we can find the value of c at which the k -core appears, i.e., the smallest c such that this equation has a positive root. For $k = 3$, for instance, we have $c_3^{\text{CORE}} = 3.351$.

The probability that a given node is in the core is a little tricky. Here we treat the node as the root. In order to survive the deletion process, it has to have at least k well-connected children. Thus the fraction of nodes in the core is

$$(3.23) \quad \gamma_k = 1 - Q_k.$$

where q is the root of (3.22). For $k = 3$ this gives $\gamma_3 = 0.268$.

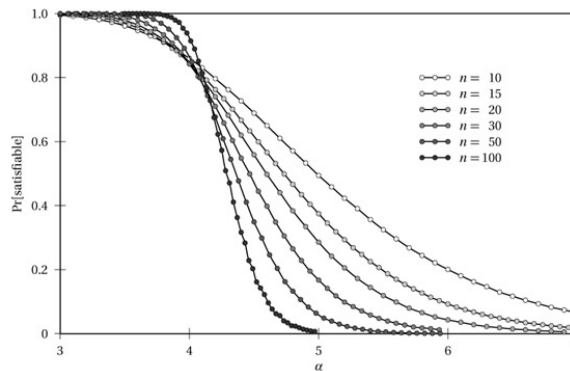
4. RANDOM k -SAT

In analogy with random graphs $G(n, m)$ that have n nodes and m edges, we can define random k -SAT formulas $F_k(n, m)$ with n variables and m clauses. We choose each clause uniformly and independently from the $2^k \binom{n}{k}$ possible clauses. That is, we choose a k -tuple of variables, and then for each one independently negate it with probability $1/2$.

We focus on the sparse case where the number of clauses scales as $m = \alpha n$ where α is a constant. In the following analysis we will not care about rare events such as two variables appearing in the same clause, or the same clause appearing twice in the formula.

Exercise: show that in the sparse case, the probability of either of these happening tends to zero as $n \rightarrow \infty$.

Phase Diagram: It is now commonly accepted that the random k -SAT (here $k = 3$) problem undergoes a phase transition at a critical value α_k . This phase transition is characterized by the probability of having a satisfying assignment converging to 1 below α_k and to 0 above.



from [5]

Mathematically speaking, this picture is still a conjecture which is known as:

Conjecture 1 (Threshold conjecture). $\forall k \geq 3, \exists \alpha_k$ such that $\forall \epsilon > 0,$

$$\lim_{n \rightarrow \infty} \mathbb{P}[F_k(n, \alpha m) \text{ is satisfiable}] = \begin{cases} 1 & \text{if } \alpha < (1 - \epsilon)\alpha_k \\ 0 & \text{if } \alpha > (1 + \epsilon)\alpha_k \end{cases}$$

The closest rigorous result for this conjecture came from Friedgut, who showed that there is a function $\alpha_k(n)$ for which it is true. But we don't know if $\alpha_k(n)$ converges to a constant as $n \rightarrow \infty$, or whether it continues to fluctuate in some way. Physically, this would be like the freezing point of water depending on the number of water

molecules, which seems very unlikely, but we still lack a rigorous proof. For the rest of this section, we will assume that the threshold conjecture is true.

4.1. Easy Upper Bound. It is easy in a first approach to derive a (not so good) upper bound for the position of the k -SAT threshold. We will use the “first moment” method: for any random variable that takes values $0, 1, 2, \dots$, such as the number of some object, we have (exercise)

$$\mathbb{P}[x > 0] \leq \mathbb{E}[x].$$

Let’s define x as the number of satisfying assignments for a given formula. It is easy to compute the average of x as the clauses are independent. Let write x as the sum over all truth assignments or “configurations”, using an indicator function to count only the satisfying ones:

$$(4.1) \quad x = \sum_{\sigma \in \{0,1\}^n} \mathbb{1}_\sigma$$

$$(4.2) \quad \mathbb{E}[x] = \sum_{\sigma} \mathbb{E}[\mathbb{1}_\sigma]$$

$$(4.3) \quad = \sum_{\sigma} \mathbb{P}(\sigma \text{ satisfies } \phi) = \sum_{\sigma} \prod_{c \in \phi} \mathbb{P}(\sigma \text{ satisfies } c)$$

$$(4.4) \quad = 2^n (1 - 2^{-k})^m.$$

For $k = 3$ in particular, this gives

$$\mathbb{E}[x] = \left[2 \left(\frac{7}{8} \right)^\alpha \right]^n.$$

Therefore we can compute an upper bound on α_3 by finding α such that $2\left(\frac{7}{8}\right)^\alpha = 1$. This gives $\alpha_3 \leq 5.19$. More generally, we have

$$\alpha_k < \frac{\ln 2}{\ln(1 - 2^{-k})} < 2^k \ln 2.$$

An interesting question, which was open until fairly recently, is whether this bound is asymptotically tight. How close to the truth is this simple counting argument? For that matter, why is it not exactly correct? For $k = 3$, in the range $4.27 < \alpha < 5.19$, the expected number of satisfying assignments is exponentially large, and yet most of the time there aren’t any at all. What’s going on?

4.2. Lower Bounds from Differential Equations and the Analysis of Algorithms. We will go through two proofs for the lower bound of the k -SAT threshold. The first (easy) one is based on the analysis of a very simple algorithm called “Unit Clause propagation” (UC). This algorithm deals with unit clauses which are clauses containing a single variable (x or \bar{x}). How could unit clauses appear in the k -SAT problem? Well, if one fixes a variable of the problem, three different options can happen concerning the clauses in which it appears. If it satisfies a clause then the clause disappears. If the clause remains unsatisfied, then the variable is removed from the clause. In that case, a 3-clause becomes a 2-clause, and a 2-clause becomes a unit clause. If the clause in which the variable appears is a unit clause and does not get satisfied when fixing the variable, then a contradiction appears.

The principle of UC is the following. Whenever there exists a unit clause, the algorithm chooses a unit clause uniformly at random and satisfies it by permanently fixing the variable of the clause. If no more unit clauses are present, the algorithm chooses a variable uniformly at random from the unset variables and fixes it randomly to 0 or 1 with probability 1/2. When all the clauses have been satisfied, the algorithm returns “satisfiable”. If the process encounters a contradiction (unsatisfied clause with all the variables of the clause fixed), it returns “contradiction” and the algorithm stops.

Algorithm 3 Unit Clause propagation (UC)

Require: a k-SAT instance

Ensure: “satisfiable” or “contradiction”

```

1: while (there is no contradiction) AND (there exists unsatisfied clauses) do
2:   if there exists a unit clause then
3:     (forced step) choose one uniformly at random and satisfy it
4:   else
5:     (free step) choose  $x$  uniformly from all the unset variables  $\rightarrow$  set  $x = 0, 1$ 
6:   end if
7:   if there exists an empty clause (or unsat clause) then
8:     return “contradiction”
9:   end if
10: end while
11: return “satisfiable”

```

At all times, the remaining formula is uniformly random conditioned on the number S_3, S_2, S_1 of unsatisfied 3-clauses, 2-clauses, and unit clauses. This comes from the fact that, since variables are removed from the formula whenever they are set, we know nothing at all about how the unset variables appear in the remaining clauses. Moreover, both moves effectively give a random variable a random value, since the (free) move is completely random, and the forced one satisfies a random unit clause. This property makes it easier to deal with the algorithm.

Let’s imagine now that we are at a time t where there remains $n - T$ variables. When dealing with 3-clauses, the probability that a chosen variable appears in it is $3/(n - T)$. When fixing this variable, there is half a chance that it will satisfy the clause and half a chance that the 3-clause became a 2-clause. The same calculation can be done for the 2-clause, a variable having $2/(n - T)$ chance to be in it. If we write $S_3 = s_3 n$, $S_2 = s_2 n$, and $T = tn$, and assume that the change in these rescaled variables equals its expectation, we can write a set of differential equations for s_2 and s_3 :

$$(4.5) \quad \frac{ds_3}{dt} = -\frac{3s_3}{1-t}$$

$$(4.6) \quad \frac{ds_2}{dt} = -\frac{2s_2}{1-t} + \frac{3}{2} \frac{s_3}{1-t}.$$

The solution of these equations is

$$(4.7) \quad s_3 = \alpha(1-t)^3$$

$$(4.8) \quad s_2 = \frac{3}{2}\alpha t(1-t)^2$$

These equations don't tell us how many unit clauses there are. Since there are $O(1)$ of these as opposed to $O(n)$ of them (indeed, there is a contradiction with high probability as soon as there are $O(\sqrt{n})$ of them) we will model them in a different way, using a branching process rather than a differential equation.

Each time we satisfy a unit clause, we might create some new ones by shortening a 2-clause. Thus the branching ratio, i.e. the expected number of new unit clauses, is $\lambda = s_2/(1-t)$. If $\lambda > 1$, this branching process explodes, leading to a contradiction. On the other hand, it can be shown that if $\lambda < 1$ throughout the algorithm, then it succeeds in satisfying the entire formula with positive probability. Because of Friedgut's theorem, proving positive probability is enough to prove probability 1. Therefore we should determine the maximum value α can take such that $\max_t \lambda \leq 1$. We find that

$$(4.9) \quad \lambda_{\max} = \frac{3}{8}\alpha,$$

and therefore

$$\alpha_3 > \frac{3}{8}.$$

This analysis can be generalized for any k , giving

$$(4.10) \quad \alpha_k \gtrsim \frac{2^k}{k}.$$

Note that this is a factor of k below the first-moment upper bound. Next, we will see how to close this gap.

4.3. Lower bounds from the second moment method. We now use the so-called second moment method to narrow the gap between the two bounds on α_k previously discussed. The second moment method relies on the inequality

$$(4.11) \quad \mathbb{P}[X > 0] \geq \frac{\mathbb{E}[X]^2}{\mathbb{E}[X^2]}.$$

X denotes the number of solutions to the formula ϕ . Note that by Friedgut's theorem, we will automatically have $\mathbb{P}[X > 0] = 1$ by simply proving that the RHS is strictly positive. We will now prove the following lower bound on α_k :

$$(4.12) \quad \alpha_k \geq 2^{k-1} \ln 2 - O(1)$$

This will reduce the width of the gap between upper and lower bound, from a factor of k to a factor of 2, independent of k . To do so, the new challenge is to compute $\mathbb{E}[X^2]$. Using the indicator function $\mathbb{1}_\sigma$ introduced previously — which is equal 1 if σ satisfies ϕ and 0 otherwise — we have the following:

$$(4.13) \quad X^2 = \left(\sum_{\sigma} \mathbb{1}_{\sigma} \right)^2 = \sum_{\sigma, \tau \in \{0,1\}^n} \mathbb{1}_{\sigma} \mathbb{1}_{\tau},$$

so that

$$(4.14) \quad \mathbb{E}[X^2] = \sum_{\sigma, \tau \in \{0,1\}^n} \mathbb{P}[\sigma, \tau \text{ both satisfy } \phi]$$

While using the first moment method previously, we had to compute $\mathbb{E}[X]$, which could be expressed in terms of $\mathbb{P}[\sigma \text{ satisfies } \phi]$, which was independent of σ due to the definition of $F_k(n, m)$. On the other hand, $\mathbb{P}[\sigma, \tau \text{ both satisfy } \phi]$ now depends on the Hamming distance z between the truth assignments σ and τ . More precisely, by independence of the clauses, we have

$$\begin{aligned} \mathbb{E}[X^2] &= \sum_{\sigma, \tau \in \{0,1\}^n} \prod_{c \in \phi} \mathbb{P}[\sigma, \tau \text{ both satisfy } c] \\ &= \sum_{\sigma, \tau \in \{0,1\}^n} \left(\mathbb{P}[\sigma, \tau \text{ both satisfy a random } c] \right)^m \end{aligned}$$

By the inclusion-exclusion principle, we then have:

$$\begin{aligned} \mathbb{P}[\sigma, \tau \text{ both satisfy a random } c] &= 1 - \mathbb{P}[\sigma \text{ doesn't satisfy } c] \\ &\quad - \mathbb{P}[\tau \text{ doesn't satisfy } c] + \mathbb{P}[\sigma, \tau \text{ both don't satisfy } c] \\ &= 1 - 2^{-k} - 2^{-k} + (z/n)^k 2^{-k} \\ &= f\left(\frac{z}{n}\right) \end{aligned}$$

To compute $\mathbb{P}[\sigma, \tau \text{ both don't satisfy } c]$, we used the fact that

$$\begin{aligned} \mathbb{P}[\sigma, \tau \text{ both don't satisfy } c] &= \mathbb{P}[\tau \text{ doesn't satisfy } c] \mathbb{P}[\sigma \text{ doesn't satisfy } c \mid \tau \text{ doesn't satisfy } c] \\ &= 2^{-k} (z/n)^k \end{aligned}$$

because σ won't satisfy c if and only if none of the variables that differ between σ and τ fall in the set of variables concerned by the clause c . This happens with probability $(z/n)^k$.

We therefore have

$$\mathbb{E}[X^2] = \sum_{z=0}^n 2^n \binom{n}{z} f\left(\frac{z}{n}\right)^m$$

where $2^n \binom{n}{z}$ is the number of couples of truth assignments σ and τ such that they differ by z bits. Note that when $\frac{z}{n} = \frac{1}{2}$, which is the most likely overlap between two truth assignments, then $f\left(\frac{z}{n}\right) = f\left(\frac{1}{2}\right) = (1 - 2^{-k})^2$, as if σ and τ were independent. If this term dominates the sum in $\mathbb{E}[X^2]$, then we have $\mathbb{E}[X^2] = 4^n (1 - 2^{-k})^{2m} = \mathbb{E}[X]^2$, so that there exists a solution to the k -SAT problem. We now need to check for which values of α this approximation holds. To do so, we write in the limit $n \rightarrow \infty$

$$\mathbb{E}[X^2] = 2^n n \int_0^1 d\zeta \binom{n}{\zeta n} f(\zeta)^m$$

By Stirling's formula,

$$\binom{n}{\zeta n} \sim \frac{1}{\sqrt{n}} e^{nh(\zeta)}$$

where $h(\zeta) = -\zeta \ln(\zeta) - (1 - \zeta) \ln(1 - \zeta)$, so that

$$\mathbb{E}[X^2] = 2^n \sqrt{n} \int_0^1 d\zeta e^{n\phi(\zeta)}$$

with $\phi(\zeta) = h(\zeta) + \alpha \ln f(\zeta)$. Using the Laplace method, we can write:

$$\mathbb{E}[X^2] \sim 2^n \sqrt{n} \sqrt{\frac{2\pi}{n|\phi''|}} e^{n\phi^{\max}}$$

When $\phi^{\max} = \phi(\frac{1}{2})$, corresponding to $e^{\phi(\frac{1}{2})} = 2(1 - 2^{-k})^{2\alpha}$, we recover that $\mathbb{E}[X^2] \sim 4^n (1 - 2^{-k})^{2m} = \mathbb{E}[X]^2$, consistently with a previous remark. For values of α where this holds, there exists a solution to the k -SAT problem. But we know that can't be true for all values of α , because we know there is a transition. We therefore need to compare the actual maximum of ϕ with ϕ^{\max} . But because h is symmetric with respect to the $\zeta = \frac{1}{2}$ axis, and $\ln f$ is increasing, it is clear that whenever $\alpha > 0$ also $\phi^{\max} > \phi(\frac{1}{2})$, so that in the limit where n goes to infinity, the second moment inequality (1) just tells us that $\mathbb{P}[X > 0] \geq 0$.

The way out is to consider another problem for which the function f is also symmetric with respect to the $\zeta = \frac{1}{2}$ axis. Consider the NAE k -SAT problem, where we ask that at least one literal in each clause be false. It is clear that a truth assignment that satisfies an NAE k -SAT formula is also a solution to the corresponding k -SAT problem, so that a lower bound for α in NAE also holds in k -SAT. It turns out the function f in NAE is given by

$$f(\zeta) = 1 - 2^{1-k} + \zeta^k 2^{1-k} + (1 - \zeta)^k 2^{1-k}$$

and therefore ϕ is symmetric around $\frac{1}{2}$. It is then straightforward to show that the maximum of ϕ is in $\frac{1}{2}$, up to some value of $\alpha = \frac{1}{2} 2^k \ln 2 - O(1)$. We therefore have $\alpha_k \geq \alpha_k^{\text{NAE}} \geq \frac{1}{2} 2^k \ln 2 - O(1)$, which is the lower bound we wanted. This lower bound can in fact be improved again to $2^k \ln 2 - O(1)$ by considering another random variable

$$X = \sum_{\sigma \text{ satis. } \phi} \eta^{\# \text{ true literals}}$$

for some η carefully chosen.

We conclude this section with some physical considerations on what the second moment inequality (4.11) means. It is interesting to interpret this inequality in terms of the planted ensemble. In the planted ensemble, we start by choosing at random a truth assignment σ , and then sample formulas ϕ such that σ satisfies ϕ . The expectation value of X^2 can be expressed in terms of the planted average in the following way:

$$\begin{aligned} \mathbb{E}[X^2] &= \sum_{\sigma, \tau} \mathbb{P}[\sigma, \tau \text{ sol.}] \\ &= \sum_{\sigma} \mathbb{P}[\sigma \text{ sol.}] \sum_{\tau} \mathbb{P}[\tau \text{ sol.} \mid \sigma \text{ sol.}] \\ &= \mathbb{E}[X] \mathbb{E}[X \mid \sigma] \\ &= \mathbb{E}[X] \mathbb{E}_{\text{planted}}[X] \end{aligned}$$

so that equation (4.11) can be rewritten

$$\mathbb{P}[X > 0] \geq \frac{\mathbb{E}[X]}{\mathbb{E}_{\text{planted}}[X]}$$

By construction, $\mathbb{E}_{\text{planted}}[X]$ is bigger than $\mathbb{E}[X]$, because a formula constructed in the planted ensemble has more solutions on average than a completely random instance. But we see here that if the number of solutions is not too different, i.e. if the planted expectation is not too different from the exact expectation in the large n limit, the second moment inequality tells us something new.

5. COMMUNITY DETECTION

5.1. The stochastic block model. It has been a general trend lately to consider everything as a network, either in physics, sociology, biology, finance. . . The question that naturally arises is whether this is justified, that is if seeing a system as a network actually allows to answer questions about the system.

In this section we consider the problem of community detection, i.e. the problem of identifying groups of nodes in a graph that share a common set of features. Community structure is called *assortative* if each group has a larger connectivity between its members than with other communities. The opposite situation is called *disassortative*. While trying to guess communities in a general setting, one doesn't know *a priori* if they are associative or disassortative.

To generate instances of this problem, we consider the Stochastic Block Model (SBM). It is a generalization of the Erdős-Renyí ensemble, with k types of vertices. These types can be thought of as colors, groups, spins, and so on. We start by supposing that k is known, and we encode the parameters of the model in a $k \times k$ affinity matrix p where $p_{r,s}$ is the probability that a given node from group r is connected with a given node from group s . Denoting the type of node i as t_i , we then construct the graph with the rule

$$\mathbb{P}[(i, j) \in E] = p_{t_i, t_j}.$$

Our goal is, given the graph G , to find the types of the nodes.

Let us begin by rewriting the problem as a statistical physics problem. The probability of a given graph G conditioned on the types of the nodes t and the affinity between communities p is given by

$$\mathbb{P}[G|t, p] = \prod_{(i,j) \in E} p_{t_i, t_j} \prod_{(i,j) \notin E} (1 - p_{t_i, t_j})$$

In the problem we are considering, t is unknown—it is what we are looking for. But we might know a priori probabilities q_r for $r \in \{1, \dots, k\}$ that a given node has

type r . We then have

$$\begin{aligned} \mathbb{P}[G|p, q] &= \sum_t \mathbb{P}[G|t, p] \mathbb{P}[t|q] \\ &= \sum_t \underset{(i,j) \in E}{\text{underset}(i, j) \in E} \prod p_{t_i, t_j} \prod_{(i,j) \notin E} (1 - p_{t_i, t_j}) \prod_i q_{t_i} \\ &= \sum_t e^{-H(t)}. \end{aligned}$$

This is the partition function of a physical system at inverse temperature $\beta = 1$, with Hamiltonian

$$(5.1) \quad H(t) = -\sum_{i \sim j} \ln p_{t_i, t_j} - \sum_{i \not\sim j} \ln (1 - p_{t_i, t_j}) - \sum_i \ln q_i$$

corresponding to a generalized Potts model with coupling constants $J_{ij} = \ln p_{t_i, t_j}$ and external fields $h_i = \ln q_i$. Note that it includes interactions between non-neighboring sites. In the sparse case ($p_{rs} = \frac{c_{rs}}{n}$), the coupling between two non-neighboring sites $i \not\sim j$ is of order $\frac{1}{n}$. However, we cannot simply get rid of these interactions, since the sum contains $O(n^2)$ such terms.

The Boltzmann distribution over assignments of types t is given, using Bayes' rule, by

$$\begin{aligned} \mathbb{P}[t|G, p, q] &= \frac{\mathbb{P}[G, t|p, q]}{\mathbb{P}[G|p, q]} \\ &= \frac{\mathbb{P}[G|t, p] \mathbb{P}[t|q]}{\mathbb{P}[G|p, q]} \\ &= \frac{e^{-H(t)}}{\mathbb{P}[G|p, q]} \end{aligned}$$

Now suppose we want to find the parameters p, q that maximize $\mathcal{Z} = \mathbb{P}[G|p, q]$; that is, we want to maximize the total probability of the network, summed over all type assignments. We can relate these optimal values of the parameters to thermodynamic quantities in the following way,

$$\begin{aligned} \frac{\partial P}{\partial p_{rs}} &= \sum_t \frac{\partial}{\partial p_{rs}} e^{-H(t)} \\ &= -\sum_t e^{-H(t)} \frac{\partial H(t)}{\partial p_{rs}} \\ &= -\sum_t e^{-H(t)} \left(-\sum_{\substack{i \sim j \\ t_i=r \\ t_j=s}} \frac{1}{p_{rs}} + \sum_{\substack{i \not\sim j \\ t_i=r \\ t_j=s}} \frac{1}{1 - p_{rs}} \right), \end{aligned}$$

where m_{rs} denotes the number of edges between group r and group s , and n_r denotes the number of nodes in group r . In a particular instance of the problem, we have

$$(5.2) \quad \begin{aligned} \frac{\partial P}{\partial p_{rs}} &= - \sum_t e^{-H(t)} \left(- \frac{m_{rs}}{p_{rs}} + \frac{n_r n_s - m_{rs}}{1 - p_{rs}} \right) \\ &\propto - \frac{\langle m_{rs} \rangle}{p_{rs}} + \frac{\langle n_r n_s \rangle - \langle m_{rs} \rangle}{1 - p_{rs}} \end{aligned}$$

where the brackets denote the average over the Boltzmann distribution.

We will assume $\langle n_r n_s \rangle = \langle n_r \rangle \langle n_s \rangle$; this holds, in particular, if both n_r and n_s are tightly concentrated with $O(\sqrt{n})$ fluctuations. In that case, (5.2) gives

$$p_{rs} = \frac{\langle m_{rs} \rangle}{\langle n_r \rangle \langle n_s \rangle}.$$

Similarly, by taking the derivative of P with respect to q_r , we find

$$q_r = \frac{\langle n_r \rangle}{n}.$$

Thus if we can estimate the averages with respect to the Boltzmann distribution, we can learn the parameters using the following Expectation-Maximization algorithm:

- E step: Compute the averages $\langle m_{rs} \rangle$ and $\langle n_r \rangle$ using the current values of p and q .
- M step: Update p and q to their most likely values given $\langle m_{rs} \rangle$ and $\langle n_r \rangle$.

We iterate until we reach a fixed point.

The only question that remains is how to estimate averages with respect to the Boltzmann distribution. This can be done by doing a Monte Carlo (MC) simulation, also known as Gibbs sampling. For instance, we can use the Heat Bath algorithm: at each step, a node in the graph is chosen and ‘‘thermalized,’’ which means that its group t_i is sampled from the marginal distribution imposed by its neighbors. More precisely, we set $t_i = s$ with probability proportional to $q_s \prod_{j \sim i} p_{s,t_j}$. It is straightforward that this

algorithm verifies detailed balance with respect to the Boltzmann distribution (5.1). After convergence, it will sample configurations with the correct weights. However, in order to compute averages like $\langle m_{rs} \rangle$, we need many independent samples, which forces us to run the algorithm for a large multiple of the autocorrelation time.

We can do better using Belief Propagation (BP). The idea of BP is that vertices pass each others estimates of marginals until consistency (a fixed point) is achieved. More precisely, we write

$$\mu_r^{i \rightarrow j} = \mathbb{P}[t_i = r \text{ if } j \text{ were absent}].$$

which is the cavity interpretation of the messages of BP. Our goal is to estimate the one-node and the two-node marginals:

$$\begin{aligned} \mu_r^i &= \mathbb{P}[t_i = r] \\ \mu_{rs}^{ij} &= \mathbb{P}[t_i = r, t_j = s] \end{aligned}$$

BP estimates of these marginals can be expressed in terms of the messages. In particular,

$$\mu_{rs}^{ij} \propto \mu_r^{i \rightarrow j} \mu_s^{j \rightarrow i} \begin{cases} p_{rs} & \text{if } i \sim j \\ 1 - p_{rs} & \text{otherwise} \end{cases}$$

The BP update rule is given by

$$(5.3) \quad \mu_r^{i \rightarrow j} = \frac{1}{\mathcal{Z}^{i \rightarrow j}} q_r \prod_{\substack{k \sim i \\ k \neq j}} \sum_s \mu_s^{k \rightarrow i} p_{rs} \prod_{\substack{k \not\sim i \\ k \neq j}} \sum_s \mu_s^{k \rightarrow i} (1 - p_{rs})$$

As usual in BP, this expression assumes that nodes other than i are independent conditioned on t_i . This is only approximately true if G is not a tree; we believe that it is approximately true when G is locally treelike, as in graphs generated by the sparse stochastic block model. Even in real networks, it works surprisingly well.

Note that (5.3) actually takes place on a fully connected graph, because of the interaction between non-neighboring sites. This yields $O(n^2)$ calculations at each update. To recover sparsity in this formula, we will assume that site i only feels the mean field of its non-neighboring sites, that is $\mu_r^{k \rightarrow i} = \mu_r^k$ for all $k \not\sim i$. With this assumption, one iteration of BP requires only $O(m)$ computations, where m is the number of edges: for sparse graphs, we have $m = O(n)$ rather than $O(n^2)$. This makes the algorithm far more scalable.

Once the parameters p, q are set correctly, we can run MC or BP to determine the most likely assignment of group types. But we get a much finer sense of what's really going on by taking a look at the whole distribution of group assignments. Zachary's Karate Club provides a cautionary tale about the risks of the procedure we just described. Wayne Zachary collected relationship data in a university karate club composed of 34 people in 1977. Due to an argument between the president of the club and the instructor, the club split up in two groups. Some followed the president, some others followed the instructor. Zachary asked whether it is possible, from the friendship network he collected before the split-up, to retrodict the composition of the two groups. This is easy to do except for one node, who was closer to the president, but ended in the instructor's group because he had a black belt exam three weeks later and didn't want to change instructor!

This story tells us that the graph cannot contain all the information relevant to assigning communities to nodes. On the other hand, by looking at the distribution of group assignments t , we can get a sense of how tightly each node is linked to its community. For instance, we could have computed that the president, the instructor, and their closest friends had a 99% chance of ending up in their respective groups, while some other nodes only had a 60% chance to end up in a given group. In this sense, determining the assignment t is rather a marginalization than a maximization problem.

A somewhat similar problem arises at the level of determining the parameters p and q . In the network literature, many authors vary the parameters p, q in order to minimize the ground state energy. That is, they minimize

$$\mathbb{P}[t|G, p, q] = \frac{e^{-H(t)}}{\mathcal{Z}}.$$

The problem with this approach is that by varying the parameters p, q , one can reach artificially low ground state energies. [11] showed that in a random 3-regular graph, one can always find a partition of the nodes in two groups such that only 11% of the edges cross from one group to the other. One might think that such a partition shows communities in the graph, although it was generated completely at random. The point is that there are exponentially many ways to partition the nodes in two groups, and it is hardly astonishing that there should exist one with few edges between the two groups: but this is really just overfitting, fitting the random noise in the graph rather than finding statistically significant communities. Instead of looking at the ground state energy, one should focus on the free energy $F = -\ln \mathcal{Z} = -\ln P[G|p, q]$. This is another reason we prefer BP to MC. In MC, computing the entropy is tricky because it requires us to integrate over different temperatures, and running MC for a long time at each temperature.

To see how BP provides us with an estimate of the free energy, we transform the problem of evaluating F into a variational problem. For simplicity we write $\mathbb{P}[\cdot]$ for $\mathbb{P}[\cdot | p, q]$. Let \mathbb{Q} be an arbitrary probability distribution over the possible assignments t . The free energy is

$$\begin{aligned} -F &= -\ln \mathcal{Z} = \ln \mathbb{P}[G] \\ &= \ln \sum_t \mathbb{P}[G, t] \\ &= \ln \sum_t \mathbb{Q}(t) \frac{\mathbb{P}[G, t]}{\mathbb{Q}(t)} \\ &= \ln \mathbb{E}_{t \sim \mathbb{Q}(t)} \frac{\mathbb{P}[G, t]}{\mathbb{Q}(t)} \end{aligned}$$

Jensen's inequality (i.e., the concavity of the logarithm) gives $\ln \mathbb{E}X \geq \mathbb{E} \ln X$. Then

$$\begin{aligned} -F &\geq \mathbb{E}_{t \sim \mathbb{Q}(t)} \ln \frac{\mathbb{P}[G, t]}{\mathbb{Q}(t)} \\ &\geq \mathbb{E}_{t \sim \mathbb{Q}(t)} \ln \mathbb{P}(t) - \sum_t \mathbb{Q}(t) \ln \mathbb{Q}(t) \end{aligned}$$

so that

$$F \leq E_{\mathbb{Q}} - S(\mathbb{Q})$$

where $E_{\mathbb{Q}}$ is the average energy if a configuration t has probability $\mathbb{Q}(t)$ instead of the Boltzmann probability \mathbb{P} , and $S(\mathbb{Q})$ is the entropy of the distribution \mathbb{Q} . Furthermore, this inequality is saturated if and only if $\mathbb{P} = \mathbb{Q}$.

This allows us to approximate the free energy by minimizing the Gibbs free energy $E_{\mathbb{Q}} - S(\mathbb{Q})$. To make the variational problem tractable, we can constrain \mathbb{Q} within a family probability distributions with a small (i.e., polynomial) number of parameters. A popular choice is the mean-field assumption, in which correlations between different sites are neglected. This assumes that \mathbb{Q} is simply a product distribution,

$$\mathbb{Q}(t) = \prod_i \mu_{t_i}^i.$$

BP does better than this rough assumption as it considers two-node correlations, in addition to single-site marginals. In fact, as was shown by [10], BP fixed points are the stationary points of the free energy within the family where \mathbb{Q} is of the form

$$\mathbb{Q}(t) = \frac{\prod_{i \sim j} \mu_{t_i, t_j}^{i, j}}{\prod_i \mu_{t_i}^{d_i - 1}}$$

Note that this form isn't even a distribution unless the underlying graph G is a tree. Plugging this form into the Gibbs free energy in general therefore only gives an approximation of the Helmholtz free energy, called the Bethe free energy.

Let us now see how BP performs on a concrete example. Again letting k denote the number of groups, we generate a graph using the following parameters of the Stochastic Block Model:

$$\begin{aligned} q_r &= \frac{1}{k} \\ p_{rs} &= \frac{c_{rs}}{n} \\ c_{rs} &= \begin{cases} c_{\text{in}} & \text{if } r = s \\ c_{\text{out}} & \text{if } r \neq s \end{cases} \end{aligned}$$

The average connectivity is

$$c = \frac{1}{k} c_{\text{in}} + \frac{k-1}{k} c_{\text{out}}.$$

Given the graph, the goal is to infer an assignment of the groups t that maximizes the overlap with the true values. Figure 2 shows the running time of BP as a function of the ratio $c_{\text{out}}/c_{\text{in}}$. We notice that it has hardly any dependence on the number of nodes n , and that at some particular value c^* a critical slowdown appears, typical of second-order phase transitions.

It is in fact possible to evaluate analytically at which value of the parameters BP begins to fail to detect the communities. All the nodes have the same average degree, so that the uniform distribution over all nodes is a fixed point of the BP equations. Let us study the linear stability of this fixed point. To do so, we introduce a small perturbation in the following way:

$$\mu_s^{i \rightarrow j} = \frac{1}{k} + \epsilon_s^{k \rightarrow i}$$

The linearization of the BP equations then takes the form:

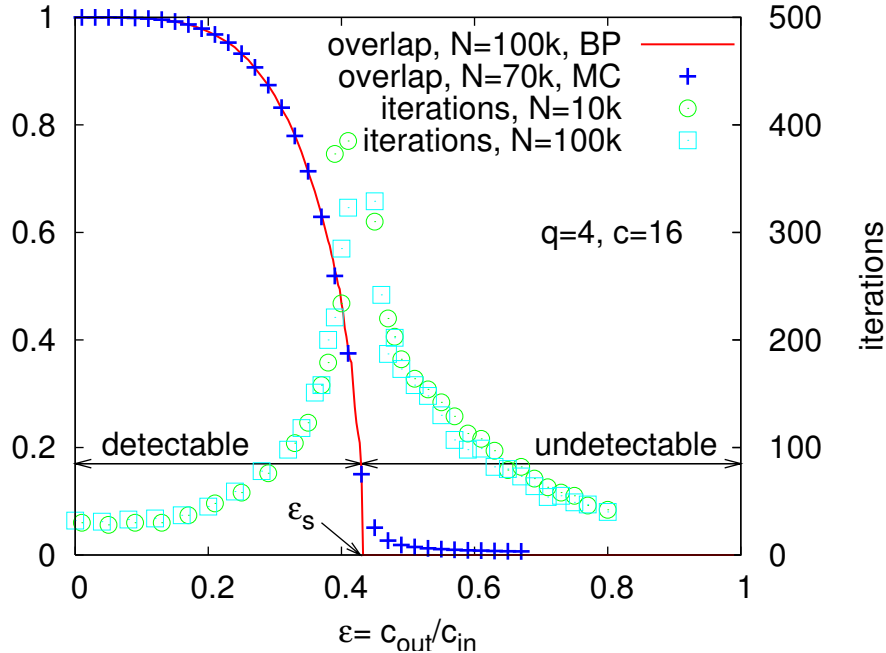
$$\epsilon^{i \rightarrow j} = \sum_{\substack{k \sim i \\ k \neq j}} T \epsilon^{k \rightarrow i}$$

where T is a $k \times k$ matrix with components

$$T_{rs} = \frac{1}{k} \left(\frac{c_{rs}}{c} - 1 \right).$$

The eigenvalues of this matrix are 0 and

$$\lambda = \frac{c_{\text{in}} - c_{\text{out}}}{kc}.$$

FIGURE 2. Second order phase transition for $k \leq 4$, from [1]

If λ is too small, the uniform distribution is a stable fixed point of BP, and we won't learn anything by running it: it will simply conclude that every node is equally likely to be in every group.

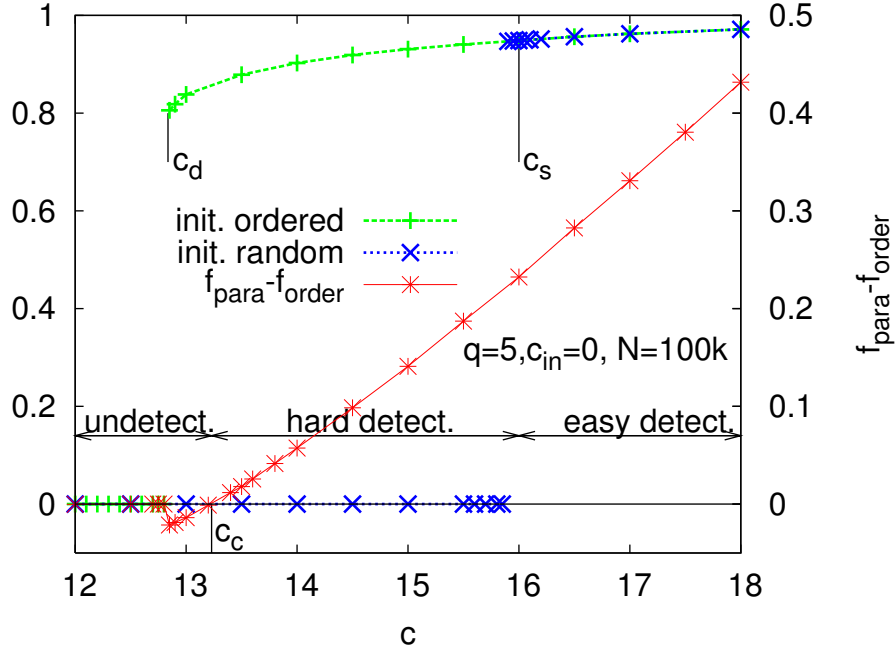
More precisely, if we assume that distant nodes are independent, then the perturbation ϵ gets multiplied by a factor $(\lambda\sqrt{c})^d$ at level d of the tree. The stability condition of the uniform distribution is therefore $\lambda\sqrt{c} < 1$, and the transition happens at

$$c_{in} - c_{out} = k\sqrt{c}.$$

Thus BP fails if $\lambda\sqrt{c} < 1$, labeling the nodes no better than chance. Our claim is that any other algorithm will also fail in this region: the randomness in the block model “washes away” the information of the underlying types.

Indeed, after this argument was presented by [1], It was shown rigorously and independently by [6] and [4] in the case of $k = 2$ groups of equal size that if $\lambda\sqrt{c} < 1$, the marginal distributions of the nodes approach the uniform distribution. In fact, the ensemble of graphs generated by the Stochastic Block Model is contiguous to the Erdős-Renyí ensemble, meaning that one graph is not enough to distinguish the two ensembles: there is no statistical test that determines, with probability approaching 1, whether communities exist or not. Conversely, if $\lambda\sqrt{c} > 1$, then a BP-like algorithm can indeed label the nodes better than chance.

One may then ask what happens if we have more than two groups. It turns out that for $k > 4$, the situation is different, and the transition is first order (see figure 3). The purple curve indicates the probability of detection when BP starts from a random initial condition. This is called robust reconstruction, but it fails at an “easy/hard transition” well above the detectability threshold c_d . On the other hand, if we give BP a hint, starting from a configuration not too far from the true assignment, then

FIGURE 3. First order phase transition for $k > 4$, from [1]

it achieves the detectability threshold (green curve). In between the detectability transition and the easy/hard transition, there are two fixed points of BP: the uniform one, and an accurate one. The Bethe free energy of the accurate was is lower, so we would choose it if we knew it; however, its basin of attraction is exponentially small. Thus we have a “hard but detectable” regime where detection is information-theoretically possible, but (we believe) exponentially hard.

5.2. Spectral methods. We now turn to another family of algorithms for community detection. These are called spectral methods, and aim at inferring types of nodes by diagonalizing a matrix which encodes the graph we are given. Popular choices for this matrix include the adjacency matrix A , the random walk matrix $D^{-1}A$ (where D is the diagonal matrix where D_{ii} is the degree d_i of node i) the Laplacian $D - A$, and the symmetrically normalized Laplacian $D^{-1/2}(D - A)D^{-1/2}$.

Generically, the largest eigenvalues of these matrices will be correlated with the community structure—except for the Laplacian, where one has to look at the smallest positive eigenvalues. More precisely, the largest eigenvector depends on the node degree or other notions of “centrality” (see for instance [8]), while (for $k = 2$ groups) the second eigenvector will be positive in one group and negative in the other. However, for the matrices listed above, this program will work up to a value of $c_{\text{in}} - c_{\text{out}}$ slightly above the BP threshold. Can we understand why?

The problem is that the spectrum of sparse random matrices is more difficult to study than that of dense matrices. If a random matrix is dense enough, with the average degree growing at least like $\log n$, then the spectrum can be modeled as a discrete part related to the community structure, and a continuous “bulk” that follows Wigner’s semicircle law, lying in the interval $[-2\sqrt{c}, 2\sqrt{c}]$ (see [7]). However, if the average degree is $O(1)$, the distribution of eigenvalues differs from Wigner’s

semi-circle law. In particular, tails appear at the edges of the semi-circle due to the high-degree vertices, which might drown the informative eigenvalue.

To make this argument more precise, it is easy to see that the adjacency matrix A (for example) has an eigenvalue λ such that

$$|\lambda| \geq \sqrt{d_{\max}}$$

where d_{\max} is the largest degree of a node in the graph. To see this, consider a node i with degree d_{\max} , and let e_i be the vector where is 1 at i and zero elsewhere. Note that $(A^k)_{ij}$ is the number of paths from i to j in the graph of exactly k steps. Since there are $d_i = d_{\max}$ ways to leave i and return with two steps, we have $A^2_{i,i} = d_{\max}$. Taking inner products, we have

$$\langle e_i, A^2 e_i \rangle = d_{\max}.$$

Since A^2 is symmetric and e_i can be orthogonally decomposed as a linear combination of eigenvectors, it follows that A^2 has an eigenvector at least d_{\max} , and A has an eigenvector which is at least $\sqrt{d_{\max}}$ in absolute value.

In the Erdős-Renyí ensemble we have $d_{\max} = O(\log n / \log \log n)$. Thus for sufficiently large n , $|\lambda|$ becomes bigger than the edge of the Wigner semicircle. As a consequence, spectral methods based on A or L will tend to find localized eigenvectors, clustered around high-degree vertices, as opposed to finding those correlated with the community structure. The normalized versions of A and L have analogous problems, where they get stuck in the trees of size $O(\log n)$ that dangle from the giant component. Thus, in the sparse case where the average degree is $O(1)$, all these methods fail significantly above the detectability threshold.

We therefore see that if we want to achieve optimal community detection by spectral methods, we need to design a matrix that forbids going back and forth on the same edge. The simplest way to do this is to define a matrix B that describes a non-backtracking walk on the edges of the network,

$$B_{i \rightarrow j, k \rightarrow \ell}^{\text{non-backtracking}} = \delta_{jk}(1 - \delta_{i\ell})$$

[3] conjectured that even in the sparse case, all of B 's eigenvalues are inside a circle of radius \sqrt{c} in the complex plane, except for the leading eigenvalue and the ones correlated with community structure. This conjecture is supported by extensive numerical experiments (see Fig. 4). Moreover, B describes to first order how messages propagate in BP, and in particular the linear stability of the uniform fixed point discussed above. Thus, if this conjecture is true, spectral clustering with B succeeds all the way down to the detectability transition for $k = 2$, and to the easy/hard transition for larger k .

As a concluding remark, we should keep in mind that this discussion assumed that the network is generated by the stochastic block model. In fact, for many networks this model is inaccurate; for instance, it doesn't produce heavy-tailed degree distributions, which are common in real networks. The good news is that many of these techniques, including belief propagation, the Bethe free energy, and so on, are easy to extend to more elaborate models, such as the degree-corrected block model of [2].

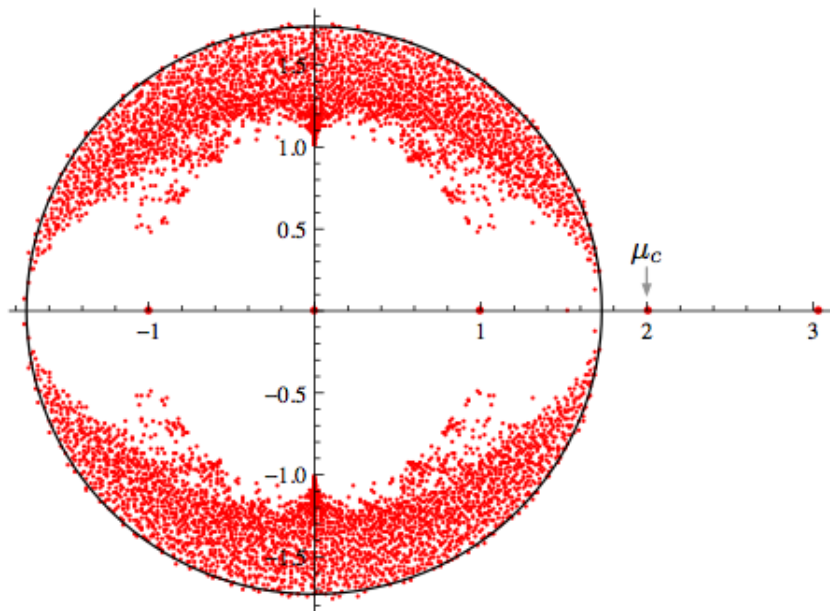


FIGURE 4. The eigenvalues of the non-backtracking matrix, from [3]

6. APPENDIX

6.1. Definition of Perfect Matching. The perfect matching problem consists in finding a set of matching between nodes such that it is compatible with the structure of the network. Let's consider the graph $G = (V, E)$ of vertices and edges. It aims to find a set of disjoint edges that covers as many nodes as possible, giving each of those nodes a unique partner. A perfect matching is a matching where all the nodes are covered.

6.2. Definition of Max Flow. Consider a directed graph $G = (V, E)$, with a source node s and a sink node t . Each edge e has a given capacity $c(e)$ which is a non-negative integer. A flow assigns a number $f(e)$ to each edge so that the total flow in and out of each node is conserved, except for the source and sink. We require that $0 \leq f(e) \leq c(e)$, and we seek to maximize the total flow, i.e., the sum of $f(e)$ over all the outgoing edges of s (or the incoming edges of t).

6.3. Definition of k-SAT. An instance of k -SAT is a formula involving n Boolean variables x_1, \dots, x_n . There are many variations such as NAESAT, 1-in- k SAT, and so on, but in the standard version of k -SAT the formula is in “conjunctive normal form.” That is, it is the AND of a set of clauses; each clause is the OR of a set of literals; and each literal is either a variable x_i or its negation \bar{x}_i . Thus a truth assignment, i.e., an assignment of true/false values to the variables, satisfies a clause if its makes at least one of its literals true, and it satisfies a formula if and only if it satisfies all its clauses. A formula is satisfiable if a satisfying assignment exists.

REFERENCES

- [1] Aurelien Decelle, Florent Krzakala, Cristopher Moore, and Lenka Zdeborová. Inference and phase transitions in the detection of modules in sparse networks. *Physical Review Letters*, 107(6):065701, 2011.
- [2] Brian Karrer and Mark EJ Newman. Stochastic blockmodels and community structure in networks. *Physical Review E*, 83(1):016107, 2011.
- [3] Florent Krzakala, Cristopher Moore, Elchanan Mossel, Joe Neeman, Allan Sly, Lenka Zdeborová, and Pan Zhang. Spectral redemption in clustering sparse networks. *Proceedings of the National Academy of Sciences*, 110(52):20935–20940, 2013.
- [4] Laurent Massoulié. Community detection thresholds and the weak ramanujan property. *arXiv preprint arXiv:1311.3085*, 2013.
- [5] Cristopher Moore and Stephan Mertens. *The nature of computation*. Oxford University Press, 2011.
- [6] Elchanan Mossel, Joe Neeman, and Allan Sly. A proof of the block model threshold conjecture. *arXiv preprint arXiv:1311.4115*, 2013.
- [7] Raj Rao Nadakuditi and Mark EJ Newman. Graph spectra and the detectability of community structure in networks. *Physical review letters*, 108(18):188701, 2012.
- [8] Mark Newman. *Networks: an introduction*. Oxford University Press, 2010.
- [9] Boris Pittel, Joel Spencer, and Nicholas Wormald. Sudden emergence of a giant k_i/k_j -core in a random graph. *Journal of Combinatorial Theory, Series B*, 67(1):111–151, 1996.
- [10] Jonathan S Yedidia, William T Freeman, Yair Weiss, et al. Generalized belief propagation. In *NIPS*, volume 13, pages 689–695, 2000.
- [11] Lenka Zdeborová and Stefan Boettcher. A conjecture on the maximum cut and bisection width in random regular graphs. *Journal of Statistical Mechanics: Theory and Experiment*, 2010(02):P02020, 2010.