

Belief Propagation for the (physicist) layman

Florent Krzakala
ESPCI Paristech
Laboratoire PCT, UMR Gulliver CNRS-ESPCI 7083,
10 rue Vauquelin, 75231 Paris, France
fk@espci.fr
<http://www.pct.espci.fr/~florent/>
(Dated: January 30, 2011)

These lecture notes have been prepared for a series of course in a master in Lyon (France), and other teaching in Beijing (China) and Tokyo (Japan). It consists in a short introduction to message passing algorithms and in particular Belief propagation, using a statistical physics formulation. It is intended as a first introduction to such ideas which start to be now widely used in both physics, constraint optimization, Bayesian inference and quantitative biology. The reader interested to pursue her journey beyond these notes is refer to [1].

I. WHY STATISTICAL PHYSICS ?

Why should we speak about statistical physics in this lecture about complex systems?

A. Physics

First of all, because of physics itself. Statistical physics is a fundamental tool and one of the most formidable success of modern physics. Here, we shall however use it only as a probabilist toolbox for other problems. It is perhaps better to do a recap before we start.

In statistical physics, we consider N discrete (not always, actually, but for the matter let us assume discrete) variables σ_i and a cost function $\mathcal{H}(\{\sigma\})$ which we call the Hamiltonian.

We introduce the temperature T (and the inverse temperature $\beta = 1/T$ and the idea is that each configurations, or assignments, of the variable $(\{\sigma\})$ has a weight $e^{-\beta\mathcal{H}(\{\sigma\})}$, which is called the Boltzmann weight.

It is convenient to introduce the following function, called the partition sum:

$$Z(\beta) = \sum_{\{\sigma\}} e^{-\beta\mathcal{H}(\{\sigma\})} \quad (1)$$

One can compute average quantities of basically any quantity $A(\{\sigma\})$ that depends on the $(\{\sigma\})$ with respect to the Boltzmann measure as:

$$\langle A(\beta) \rangle = \frac{1}{Z(\beta)} \sum_{\{\sigma\}} A(\{\sigma\}) e^{-\beta\mathcal{H}(\{\sigma\})} \quad (2)$$

Of particular importance is the average of the Hamiltonian itself, which is called the energy

$$E(\beta) = \langle \mathcal{H} \rangle = \frac{1}{Z(\beta)} \sum_{\{\sigma\}} \mathcal{H}(\{\sigma\}) e^{-\beta\mathcal{H}(\{\sigma\})} \quad (3)$$

It is customary to introduce the so-called free energy

$$F(\beta) = -\frac{1}{\beta} \log Z(\beta) \quad (4)$$

Indeed, one can obtain the energy directly from the free energy as

$$E(\beta) = \frac{\partial \beta F(\beta)}{\partial \beta} \quad (5)$$

Another important quantity is the entropy which reads.

$$S(\beta) = -\beta(F(\beta) - E(\beta)) = \frac{\partial F(T)}{\partial T} \quad (6)$$

These quantities called free energy, energy, entropy, etc... have an interest of their own in statistical physics. However, we shall see that they are interesting as well in other contexts. For instance, it is important to notice that the zero temperature limit (when $\beta \rightarrow \infty$) of the energy give the value of the minimum of the Hamiltonian (which is called the ground state energy in physics):

$$\lim_{\beta \rightarrow \infty} E(\beta) = \min_{\{\sigma\}} \mathcal{H}(\{\sigma\}) \quad (7)$$

while the entropy at zero temperature give us the (logarithm of the) number of configuration Ω of $\{\sigma\}$ with this precise value of the cost function.

$$\lim_{\beta \rightarrow \infty} S(\beta) = \log \Omega \quad (8)$$

B. Constraint Satisfaction Problems

The last statement allow to see the connection with combinatorial optimization: finding the minimal value of a function of discrete variable, and the corresponding number of assignment means simply computing the ground-state energy (the zero temperature energy) and entropy of the associated statistical physics model! If we know how to solve the statistical physics model, we know the value of the best assignment!

Let us consider an example: the Coloring problem. We need to color a graph with q colors such that two neighboring nodes have a different color. The coloring problem was introduced by Francis Guthrie in 1852 when he was working on coloring of map, and notice that all maps (that is, all planar graphs!) can be colored with 4 colors. The corresponding 4-color theorem, was however only proven in 1976 by Kenneth and Appel [2, 3].

For a generic graph \mathcal{G} (that is, a non planar one) however, deciding if it is colorable with q color (that is, if the chromatic number $\chi \mathcal{G} \leq q$) is a NP-complete problem as soon as $q > 2$.

Let us now make the connection with statistical physics. The problem is equivalent to finding the ground state of the so called *Potts anti-ferromagnet* on this graph, a problem with Hamiltonian (or cost function):

$$\mathcal{H} = \sum_{ij \in \mathcal{G}} \delta s_i, s_j \quad \text{with} \quad s_i = 1, \dots, q \quad (9)$$

If one can solve the statistical physics problem, and compute $Z(\beta)$, then if the zero temperature energy is zero, the graph is colorable, and the entropy allows to obtain the number of such solutions... of course, it is usually not simple to compute the partition sum (in fact, it is generically a problem which is much harder than NP complete).

Another important point is that one can also use some method borrowed from statistical physics: simulated annealing [4] is for instance one of the most cited paper in computer science. Of course the famous K-SAT is another example where this approach can be applied [5].

C. Bayesian inference

Another example of application of the ideas of statistical physics is given by the field of Bayesian inference, or machine learning [6]. It is a method of statistical inference in which some observations are used to calculate the probability that an hypothesis may be true, or else to update its previously-calculated probability and to compute its most probable value.

Consider the following example: We have a disease which is connected with some 3 symptoms. A contaminated individual may display these symptoms with probability p_{11}, p_{12}, p_{13} while a healthy one will display them with probability p_{01}, p_{02}, p_{03} . We have a list of patient and we know if they display these symptom or not, and the goal is to estimate these set of probabilities $\{p\}$.

A simple way to do would be to maximize the probability of the parameter given that we observe such a graph of connections (see fig.1). So we want to consider $P(\{p\}|\mathcal{G})$. Unfortunately, we do not know any expression to do so. However, we could use Bayes theorem and write

$$P(\{p\}|\mathcal{G}) = P(\mathcal{G}|\{p\}) \frac{P(\{p\})}{P(\mathcal{G})} = \sum_{\{S\}} \frac{P(\{p\})}{P(\mathcal{G})} P(\mathcal{G}, \{S\}|\{p\}) \quad (10)$$

where the sum over $\{S\}$ is over all the possible assignments for each individual ($S_1 = 1$ means the individual 1 is sane, while $S_1 = 0$ means he is sick). It is often assumed that $P(\mathcal{G})$ and $P(\{p\})$ are *a priori* unknown, so that we should

use them basically as normalization constant. In fact the prior probability $P(\{p\})$ includes all graph independent information about the values of the parameters. With our definition of the problem, that is inference of parameters from the topology of the graph, we have no such prior information. The prior probability $P(\mathcal{G})$ does not depend on the values of parameters, so we realize that maximizing $P(\{p\}|\mathcal{G})$ over $\{p\}$ is equivalent to maximizing the following *partition sum* (which is called the *evidence* in Bayesian inference):

$$Z(\{p\}) = \sum_{\{S\}} P(\mathcal{G}, \{S\}|\{p\}) \quad (11)$$

At this point the connection with statistical physics is evident! It is even stronger when one write explicitly the $P(\mathcal{G}, \{S\}|\{p\})$, since indeed:

$$P(\mathcal{G}, \{S\}|\{p\}) = \prod_{ij} \left[p_{s_i,j}^{\mathcal{G}_{ij}} (1 - p_{s_i,j})^{1-\mathcal{G}_{ij}} \right] \quad (12)$$

where the indice i is for all individual, and the value j for all symptoms, and where $\mathcal{G}_{ij} = 1$ if there is indeed a link between i and j , and zero otherwise. One can rewrite the weight as follow to obtain the connection with statistical physics with an equivalent Hamiltonian

$$P(\mathcal{G}, \{S\}|\{p\}) = \exp \left(\sum_{ij} (\mathcal{G}_{ij} \log(p_{s_i,j}) + (1 - \mathcal{G}_{ij}) \log(1 - p_{s_i,j})) \right) \quad (13)$$

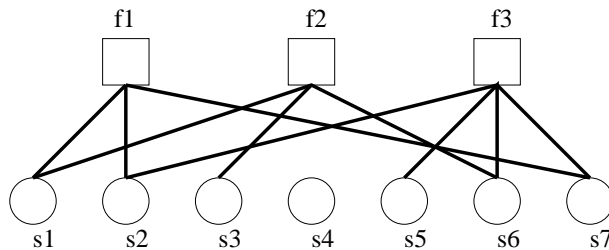


FIG. 1: Example of an inference problem represented by a factor graph.

Again, we see that the ability to compute a partition sum is very useful, as it allows to do Bayesian inference (where here we should use $\beta = 1$ in order to recover the statistical physics setting).

II. RECURSIONS ON A FINITE TREE: THE FERROMAGNETIC EXAMPLE

Computing a partition sum is hard! With N variable, we have a sum of exponentially many elements. However, pretty much any model or Hamiltonian can be solve exactly on tree, in a time linear in N . This is actually used in many different fields of science with various names (Bethe-Peierls in physics, or cavity method, Belief-Propagation in machine learning and artificial intelligence, or Sum-Product in coding theory). We shall here stick to the physicists notation, but refer to the algorithm as *Belief Propagation* [7].

A note of warning with respect to what can be find in the literature: We will for the moment restrict to consider models where the clauses have connectivity $k = 2$: for these models, one does not need a factor graph representation, since interactions can be represented by a standard graph with vertices $i = 1, \dots, N$ representing variables σ_i and links $\langle i, j \rangle$ representing interactions $\psi_{ij}(\sigma_i, \sigma_j)$. We shall repeatedly use in the following the same notation we used for factor graphs, specialized to the $k = 2$ case: therefore we use ∂i for the set of vertices adjacent to a given vertex i , i.e. for the sites which interact with i , and $\partial i \setminus j$ for those vertices around i distinct from j .

We shall consider for pedagogical purpose what is called *ferromagnetic model* in the physics literature.

1. A single spin in a field



FIG. 2: One spin in a field.

We consider a single spin in a field. The Hamiltonian is

$$\mathcal{H} = -h_0 S_0 \quad (14)$$

Let us solve the problem. We denote the total partition sum $Z(\beta)$, which simply reads

$$Z(\beta) = \sum_{S_0} e^{\beta h_0 S_0} = 2 \cosh(\beta h_0) \quad (15)$$

It is in fact interesting to define a partition sum when S_0 is fixed. In this case

$$Z_0(\beta, S_0) = e^{\beta h_0 S_0} \quad (16)$$

$$\text{and } Z = Z_0(+) + Z_0(-) \quad (17)$$

The energy is given by $E(\beta) = -h_0 \tanh \beta h_0$ and it gives at zero temperature $E(T = 0) = -h_0$ at it should. The average value of the spin is given by

$$m(\beta) = \frac{1}{Z(\beta)} \sum_{S_0} S_0 e^{\beta h_0 S_0} = \frac{Z_0(+)-Z_0(-)}{Z_0(+)+Z_0(-)} = \tanh(\beta h_0) \quad (18)$$

From this, we know the probability to be plus ($\eta_+ = \frac{1+m}{2}$) and to be minus ($\eta_- = \frac{1-m}{2}$).

This is a trivial example but there is an important comment to be made here. Note that for such a single random variable, one can specify either m , η_{\pm} or h_0 : this is completely equivalent! Of course, it is more physicists who like to use the field notation, but it is actually often very practical for computations.

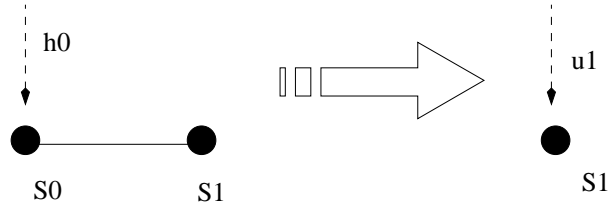


FIG. 3: Two spins in a field...

2. Adding one spin

We now add a new spin and we obtain the new Hamiltonian:

$$\mathcal{H} = -h_0 S_0 - J S_0 S_1 \quad (19)$$

How to compute Z for this new problem? We again introduce a partial partition sum $Z_1(S_1)$ which tell us what is the partition sum of the total system when S_1 is fixed). We find:

$$Z_1(+)=\sum_{S_0} e^{-\beta \mathcal{H}}=e^{\beta h_0} e^{\beta J}+e^{-\beta h_0} e^{-\beta J} \quad (20)$$

$$Z_1(-)=\sum_{S_0} e^{-\beta \mathcal{H}}=e^{\beta h_0} e^{-\beta J}+e^{-\beta h_0} e^{\beta J} \quad (21)$$

$$(22)$$

Of course, $Z=Z_1(+)+Z_1(-)$. Interestingly, we recover here a factor $e^{\pm \beta h_0}$ which was the partial partition sum of the spin S_0 before we added the new spin S_1 ! In the context of this new problem, with Hamiltonian 19 we shall call this partial partition sum where S_1 is not present $Z_{0 \rightarrow 1}(S_0)$ (and we shall now reserve the notation $Z_0(S_0)$ for the partial partition sum of the variable (S_0) with the new Hamiltonian 19. Note that now $Z_0(S_0)$ is different from $2 \cosh(\beta h_0)$ because of the interaction with the new spin.). Denoting therefore $Z_{0 \rightarrow 1}(S_0)=e^{\beta h_0 S_0}$ we that

$$Z_1(S_1)=\sum_{S_0} Z_{0 \rightarrow 1}(S_0) e^{\beta J S_0 S_1} \quad (23)$$

$$=e^{\beta h_0} e^{\beta J S_1}+e^{-\beta h_0} e^{-\beta J S_1} \quad (24)$$

We see that the partial partition sum of the new spin S_1 can be express as a simple function of the partition sum of S_0 before S_1 was introduced! This is typical of the kind of equation we shall see in the following. This recursion is in fact the root of the BP equation.

In fact, a different —but equivalent— recursion can be written if one is interested in the average values of variables instead of the partition sum. What is the average value of the new spin S_0 ? As before, we want to use a physicist notation and we wrote

$$m_1=\langle S_1 \rangle=\tanh(\beta u_1) \quad (25)$$

This is a *definition* of u_1 . This is to say that the spin 1 acts as if it is *alone* in a field u_1 . We have, using either 22 or 24

$$\tanh(\beta u_1)=\frac{Z_1(+)-Z_1(-)}{Z_1(+)+Z_1(-)}=\tanh \beta J \tanh \beta h_0. \quad (26)$$

At this point something miraculous has happened: The new spin is acting as if it is alone, without any interaction with the other spin 0, expect that it is in an effective field given by the recursion 26. As we see, a recursion can be written for the field. Going back to the average value of the variable, we obtain

$$m_1=\tanh \beta J m_{0 \rightarrow 1} \quad (27)$$

where $m_{0 \rightarrow 1}$ the average magnetization of the variable 0 when variable 1 is not (yet) present. Note that $m_{0 \rightarrow 1}$ might be different from m_0 in the two spins problems, so one should not confuse the two! Finally, for the more probabilistic reader, we could note the recursion as:

$$\eta_1(+)=2 * \tanh (\beta J)\left(1+\eta_{0 \rightarrow 1}\right)-1 \quad (28)$$

In the physicist jargon, we denote the quantities with indice $0 \rightarrow 1$ that denote averages when the link between 0 and 1 is not present, as *cavity quantity*, as we are somehow making a cavity in the system and consider only a sub-system. $m_{0 \rightarrow 1}$ is for instance the cavity magnetization.

The very basic point of BP is to write a recursion for these cavity variables.

3. A third spin...

We add a again new spin called S_2 and we have the new Hamiltonian:

$$\mathcal{H} = -h_0 S_0 - JS_0 S_1 - JS_1 S_2 \quad (29)$$

Let us proceed as before: we have just computed the value of the partition sum of spin 1 when S_2 was not (yet) present and now denote this $Z_{0 \rightarrow 2}(S_0)$ in the new context of Hamiltonian 29. Clearly, the new partition sum when S_2 is fixed is simply given by

$$Z(S_2) = \sum_{S_1} Z_{1 \rightarrow 2}(S_1) e^{\beta J S_1 S_2} \quad (30)$$

and it can be computed right away as soon as $Z_{1 \rightarrow 2}(S_1)$ is known! We see that any time we add a new spin, only a single operation is needed to compute the new partition sum, and this is faster than the use of the original formula.

In terms of field, one can repeat the previous argument, and argue that the new field S_2 is basically acting as if it was isolated, but in a field u_2 given by

$$\tanh(\beta u_2) = \tanh \beta J \tanh \beta u_{1 \rightarrow 2} = (\tanh \beta J)^2 \tanh \beta h_0 \quad (31)$$

4. The 1D Ising chain

Now, we can have a first idea of what these iterations can do for us: let us consider a very long chain of N spins, with Hamiltonian:

$$\mathcal{H} = -h_1 S_1 - J \sum_{i=1}^{N-1} S_i S_{i+1} \quad (32)$$

The partition sum when the last spin is fixed is simply given by

$$Z(S_N) = \sum_{S_{N-1}} Z_{N-1 \rightarrow N}(S_{N-1}) e^{\beta J S_{N-1} S_N} \quad (33)$$

so that one just need to compute a new operation anytime one add a new spin to compute the partition sum. This can be done in $O(N)$ operation, to be compare with the costly 2^N evaluation if one uses the definition of the partition sum: for such a chain, we went from exponential to linear!

The same can be done if we focus on the average value as well. Repeating the previous argument, we can obtain that the last variable behave as if in a field u_N where

$$\tanh \beta u_N = \tanh \beta h_1 (\tanh \beta J)^{N-1} \quad (34)$$

Or equivalently

$$m_N = \tanh \beta h_1 (\tanh \beta J)^{N-1} \quad (35)$$

In particular, since $\tanh \beta J \leq 1$ at finite temperature, we immediately see that the magnetization will be zero far enough from the first spin, so that there is no ordered phase (in accord with the well known fact that there is no ordered phase at finite temperature in one dimension).

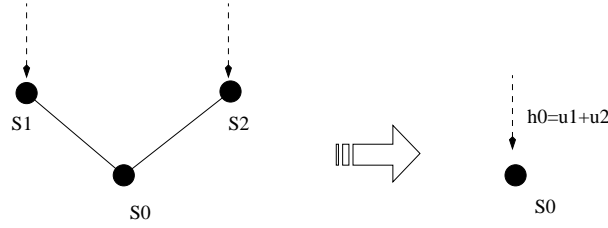


FIG. 4: The simplest tree ...

5. A basic recursion on a proto-tree

This is all nice and well, but chains are pretty limited. Can we do better? In fact we can, let us move to the following situation with 3 spins:

$$\mathcal{H} = -h_1 S_1 - h_2 S_2 - JS_0 S_1 - JS_0 S_2 \quad (36)$$

We shall apply the very same strategy and write the partition sum for Z_0 . We find:

$$Z_0(S_0) = \sum_{S_1, S_2} e^{-\beta \mathcal{H}(S_0)} = \sum_{S_1, S_2} e^{\beta h_1 S_1 + \beta h_2 S_2 + \beta JS_0 S_1 + \beta JS_0 S_2} \quad (37)$$

$$= \prod_{i=1,2} \left[\sum_{S_i} e^{\beta h_i S_i + \beta JS_0 S_i} \right] \quad (38)$$

$$= \prod_{i=1,2} \left[\sum_{S_i} Z_{i \rightarrow 0}(S_i) e^{\beta JS_0 S_i} \right] = \prod_{i=1,2} Z_{0i}(S_0) \quad (39)$$

This factorization we just did is really convenient, and in particular we again see that we could write the recursion using or notation $Z_{i \rightarrow 0}(S_i)$, just as before! This is the root of the BP recursion we are about to write!

The expression is also very convenient in order to write the effective field h_0 on S_0 such that $S_0 = \tanh \beta h_0$ so that

$$\langle S_0 \rangle = \tanh(\beta h_0) = \frac{e^{\beta h_0} - e^{-\beta h_0}}{e^{\beta h_0} + e^{-\beta h_0}} = \frac{Z_0(+)-Z_0(-)}{Z_0(+)+Z_0(-)} \quad (40)$$

If we define u_i by $h_0 = \sum_i u_i$ we see immediately that a solution is given by

$$\frac{e^{\sum_i \beta u_i} - e^{-\sum_i \beta u_i}}{e^{\sum_i \beta u_i} + e^{-\sum_i \beta u_i}} = \frac{\prod_i e^{\beta u_i} - \prod_i e^{-\beta u_i}}{\prod_i e^{\beta u_i} + \prod_i e^{-\beta u_i}} = \frac{\prod_i Z_{0i}(+) - \prod_i Z_{0i}(-)}{\prod_i Z_{0i}(+) + \prod_i Z_{0i}(-)} \quad (41)$$

Given what we have done for one spin with one neighbors, we thus must satisfy

$$e^{\pm \beta u_i} \propto \sum_{S_i} e^{\pm \beta h_i} e^{\pm \beta JS_i} \quad (42)$$

We know the solution to this problem, since it is the one with two spins! It is simply

$$\tanh(\beta u_i) = \tanh \beta J \tanh \beta h_i. \quad (43)$$

And we thus find that the following result, which in fact holds for more than two spins:

$$\beta h_0 = \sum_i \operatorname{atanh}(\tanh \beta J \tanh \beta h_i) \quad (44)$$

This can be interpreted as a message passing algorithm on any tree, if one interpret, as before, this Z_0 as being in fact a $Z_{0 \rightarrow 4}$ upon the addition of a new variable 4.

6. *Application: the solution of the ferromagnet on an infinite tree*

A first application is the solution of the ferromagnetic Ising model on an infinite tree with connectivity k (this is an approximation [8] of the finite dimensional problem in physics where $2d = k + 1$.) We have simply the fixed point:

$$\beta h = k \tanh(\tanh \beta J \tanh \beta h) \quad (45)$$

$$m = \tanh[k \tanh(\tanh \beta J m)] \quad (46)$$

$$(47)$$

Let us see how it works: these are the value it predicts for β_c compares with the correct one:

TABLE I: Gaussian couplings: critical exponent in MK and in 2D

d	Bethe tree	A lattice in finite dimension
1	0	0
2	0.346	0.44
3	0.203	0.221
4	0.144	0.149
5	0.112	0.114

TABLE II: Critical β

It is really a good approximation! We see also that we can interpret these equations as message passing on trees!

III. GENERIC EQUATIONS FOR BELIEF PROPAGATION

We are now ready to write the BP recursion [1]. Let us do it for any kind of interaction involving two variables, where the interaction graph is a tree. Following what we have done, we define the quantity $Z_{i \rightarrow j}(\sigma_i)$, for two adjacent sites i and j , as the partial partition function for the subtree rooted at i , excluding the branch directed towards j , with a fixed value of the spin variable on the site i . We also introduce $Z_i(\sigma_i)$, the partition function of the whole tree with a fixed value of σ_i . These quantities can be computed according to the following recursion rules, see Fig. 5 for an example,

$$Z_{i \rightarrow j}(\sigma_i) = \prod_{k \in \partial i \setminus j} \left(\sum_{\sigma_k} Z_{k \rightarrow i}(\sigma_k) \psi_{ik}(\sigma_i, \sigma_k) \right), \quad Z_i(\sigma_i) = \prod_{j \in \partial i} \left(\sum_{\sigma_j} Z_{j \rightarrow i}(\sigma_j) \psi_{ij}(\sigma_i, \sigma_j) \right). \quad (48)$$

This is the cavity recursion for partition sums.

We saw that it was also possible to write recursion in terms of probabilities. Let us thus rewrite these equations in terms of normalized quantities which can be indeed interpreted as probability laws for the random variable σ_i , namely $\eta_{i \rightarrow j}(\sigma_i) = Z_{i \rightarrow j}(\sigma_i) / \sum_{\sigma'} Z_{i \rightarrow j}(\sigma')$ and $\eta_i(\sigma_i) = Z_i(\sigma_i) / \sum_{\sigma'} Z_i(\sigma')$. The recursion equations read in these notations

$$\eta_{i \rightarrow j}(\sigma_i) = \frac{1}{z_{i \rightarrow j}} \prod_{k \in \partial i \setminus j} \left(\sum_{\sigma_k} \eta_{k \rightarrow i}(\sigma_k) \psi_{ik}(\sigma_i, \sigma_k) \right), \quad \eta_i(\sigma_i) = \frac{1}{z_i} \prod_{j \in \partial i} \left(\sum_{\sigma_j} \eta_{j \rightarrow i}(\sigma_j) \psi_{ij}(\sigma_i, \sigma_j) \right), \quad (49)$$

where $z_{i \rightarrow j}$ and z_i are normalization constants:

$$z_{i \rightarrow j} = \sum_{\sigma_i} \prod_{k \in \partial i \setminus j} \left(\sum_{\sigma_k} \eta_{k \rightarrow i}(\sigma_k) \psi_{ik}(\sigma_i, \sigma_k) \right), \quad z_i = \sum_{\sigma_i} \prod_{j \in \partial i} \left(\sum_{\sigma_j} \eta_{j \rightarrow i}(\sigma_j) \psi_{ij}(\sigma_i, \sigma_j) \right), \quad (50)$$

Since the leaves are isolated when the link connecting them is removed, one has $Z_{i \rightarrow j}(\sigma_i) = \text{const.}$ and $\eta_{i \rightarrow j}(\sigma_i) = \text{const.}$ for leaves. However, one can also choose to put an arbitrary $\eta_{i \rightarrow j}(\sigma_i)$ on the leaves: this might represent an external field acting on them, or the effect of a given boundary condition. Moreover the quantity $\eta_i(\sigma_i)$ is exactly the marginal probability law of the Gibbs-Boltzmann distribution, hence the local magnetizations can be computed as $m_i = \langle \sigma_i \rangle = \sum_{\sigma} \eta_i(\sigma) \sigma$. Then, on a given tree the recursion equations on the $\eta_{i \rightarrow j}$ for all directed edges of the graph have a single solution, easily found by propagating the recursion from the leaves of the graph.

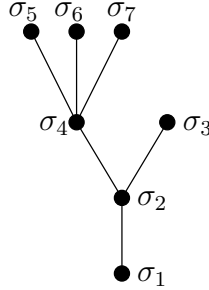


FIG. 5: Example of an Ising tree model on 7 vertices. The definition of $Z_{2 \rightarrow 1}$ and its recursive computation reads here: $Z_{2 \rightarrow 1}(\sigma_2) = \sum_{\sigma_3, \dots, \sigma_7} \psi_{23}(\sigma_2, \sigma_3) \psi_{24}(\sigma_2, \sigma_4) \psi_{45}(\sigma_4, \sigma_5) \psi_{46}(\sigma_4, \sigma_6) \psi_{47}(\sigma_4, \sigma_7) = \sum_{\sigma_3, \sigma_4} Z_{3 \rightarrow 2}(\sigma_3) Z_{4 \rightarrow 2}(\sigma_4) \psi_{23}(\sigma_2, \sigma_3) \psi_{24}(\sigma_2, \sigma_4)$.

The reader might complain that with these new notation, one could not compute the partition sum. In fact, it is perfectly possible: we can write the free energy of the system as a function of the cavity probabilities. For this it is useful to first define the object

$$z_{ij} = \sum_{\sigma_i, \sigma_j} \eta_{j \rightarrow i}(\sigma_j) \eta_{i \rightarrow j}(\sigma_i) \psi_{ij}(\sigma_i, \sigma_j) = \frac{z_j}{z_{j \rightarrow i}} = \frac{z_i}{z_{i \rightarrow j}}, \quad (51)$$

where the last two equalities are easily derived using Eqs. 49. Indeed, the total partition function is

$$Z = z_i \prod_{j \in \partial i} z_{j \rightarrow i} \prod_{k \in \partial j \setminus i} z_{k \rightarrow j} \cdots = z_i \prod_{j \in \partial i} \frac{z_j}{z_{ij}} \prod_{k \in \partial j \setminus i} \frac{z_k}{z_{jk}} \cdots = \frac{\prod_i z_i}{\prod_{\langle i, j \rangle} z_{ij}} \quad (52)$$

and the free energy is

$$\begin{aligned} F &= -T \log Z = \sum_i f_i - \sum_{\langle i, j \rangle} f_{ij} , \\ f_i &= -T \log z_i , \\ f_{ij} &= -T \log z_{ij} . \end{aligned} \quad (53)$$

This is the full form of BP for any graph with two body interactions.

IV. RANDOM GRAPHS AND HYPER-GRAPHS

The common notion of computational complexity, being based on worst-case considerations, could overlook the possibility that “most” of the instances of an NP problem are in fact easy and that the difficult cases are very rare. Random ensembles of problems have thus been introduced in order to give a quantitative content to this notion of typical instances; a property of a problem will be considered as typical if its probability (with respect to the random choice of the instance) goes to one in the limit of large problem sizes. Of course the choice of a distribution over the instances is arbitrary and could not reflect the properties of the instances that relevant for a given practical application. Still, the introduction of a distribution over the instances allows to formulate the problem in terms of the statistical mechanics of a spin-glass-like model. We will see that this formulation provides important insight in the properties of difficult instances of these problems.

An instance of a random CSP is defined by two objects: the underlying factor graph and, as in fully connected models, the set of *couplings* J appearing in the constraints (e.g. the right hand side of an equation in XORSAT). Both the factor graph and the couplings can be taken as random variables to define a probability distribution over instances. Recall that we have N variable nodes and M constraint nodes. In the statistical mechanics approach we will be interested in the thermodynamic limit of large instances where N and M both diverge with a fixed ratio $\alpha = M/N$.

Among many possible ensembles of graphs, two have been investigated in great detail:

- *Random regular graphs* (or fixed connectivity): each constraint involves k distinct variables (k is a free parameter for (XOR)SAT and $k = 2$ for COL), and each variable enters in *exactly* c different constraints. Uniform probability is given to all graphs satisfying this property. Note that one must have $Mk = Nc$, i.e. $c = kM/N = k\alpha$.
- *Erdős-Rényi random graphs*: For each of the M clauses a a $k(\geq 2)$ -uplet of distinct variable indices (i_a^1, \dots, i_a^k) is chosen uniformly at random among the $\binom{N}{k}$ possible ones. For large N, M the degree of a variable node of the factor graph converges to a Poisson law¹ of average αk . To compare with regular graphs we shall use the notation $c = k\alpha$ for the average connectivity.

In principle one might allow also the connectivity of the constraints to be a random variable but we do not discuss such case here. Note that the limit $c \rightarrow \infty$ with a proper scaling of the couplings gives back the fully connected model. In this limit, the fluctuations of c in Erdős-Rényi graphs can be neglected and the two ensembles of graph are equivalent.

Random (hyper)graphs have many interesting properties in this limit. In particular, picking at random one variable node i and isolating the subgraph induced by the variable nodes at a graph distance smaller than a given constant L yields, with a probability going to one in the thermodynamic limit, a (random) tree. This tree can be described by a Galton-Watson branching process: the root i belongs to l constraints, where l is a Poisson random variable of parameter αk (or $l = \alpha k$ in the fixed connectivity case). The variable nodes adjacent to i give themselves birth to new constraints, in numbers which are independently Poisson distributed with the same parameter. This reproduction process is iterated on L generations, until the variable nodes at graph distance L from the initial root i have been generated.

The algorithm to construct Erdős-Rényi graphs is trivial, since it is given by the definition. Fixed connectivity graphs can be constructed as follows: first one attaches to each variable node a number c of links, thus obtaining Nc links. These links have to be connected to the $Mk = Nc$ links attached to constraint nodes. To do this, one simply numbers the links from 1 to Nc and then pick up a random permutation of these numbers in order to decide which of the variable links has to be attached to a given constraint link. The resulting graph however might have variables that are connected twice or more to the same node. In this case the permutation is discarded and a new one is picked until a good graph is reached. In practice the probability of such event is small if N is large and c not too large, so that the procedure converges fast to a good graph.

The key properties of random graphs that is exploited by the cavity method is that loops are very large in the thermodynamic limit, as we discussed in the introduction to this section. Therefore, locally random graphs looks like trees. Before studying the cavity method, we must understand how to solve statistical mechanics models defined on trees.

¹ This can be seen as follow: the probability that a given site is connected to another one at each steps is $p = 2/N$. Therefore, the distribution probability of the number of connection is given by $Q(k) = p^k(1-p)^{M-k}\binom{M}{k}$. In the limit when $M \gg k$ this goes to $Q(k) = \frac{M^k}{k!}e^{-Mp}p^k = \frac{1}{k!}e^{-c}c^k$.

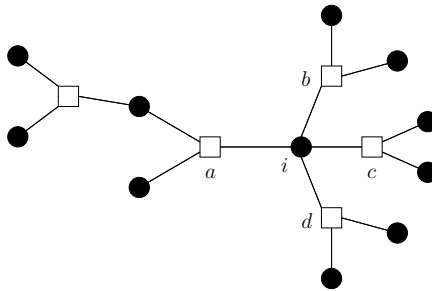


FIG. 6: An example of factor graph. The neighborhoods are for instance $\partial i = \{a, b, c, d\}$ and $\partial i \setminus a = \{b, c, d\}$

V. GENERIC EQUATIONS FOR BELIEF PROPAGATION ON HYPER-GRAPHS

What if we have a problem with more than 3-body interaction, such as the one in the inference problem in the first section? Luckily, our equations generalize easily to this case and we can write the recursion using a ψ_{ijk} just as easily as we did for a ψ_{ij} . In order to do so, we shall the graphical models and factor graphs.

Factor graphs provide an useful representation of a CSP. These graphs (see Fig. 6 for an example) have two kind of nodes. Variable nodes (filled circles on the figure) are associated to the degrees of freedom σ_i , while constraint nodes (empty squares) represent the clauses ψ_a . An edge between constraint a and variable i is drawn whenever ψ_a depends on σ_i . The neighborhood ∂a of a constraint node is the set of variable nodes (i_1, \dots, i_{K_a}) that appear in a . Conversely we will denote by ∂i the set of all constraints (a_1, a_2, \dots) in which variable i is involved. We shall conventionally use the indices i, j, \dots for the variable nodes, a, b, \dots for the constraints, and denote \setminus the subtraction from a set. The graph distance between two variable nodes i and j is the number of constraint nodes encountered on a shortest path linking i and j (formally infinite if the two variables are not in the same connected component of the graph).

Let us do it for any kind of interaction involving two variables, again where the interaction graph is a tree. We refer the reader to [1] for a detailed derivation, and we shall here directly gives the results in term of (cavity) probabilities.

We saw that it was also possible to write recursion in terms of probabilities. The recursion equations read

$$\eta_{j \rightarrow a}(\sigma_j) \propto \prod_{b \in \partial j \setminus a} \nu_{b \rightarrow j}(\sigma_j) \quad (54)$$

with

$$\nu_{a \rightarrow j}(\sigma_j) \propto \sum_{\sigma \in \partial a \setminus j} \psi_a \prod_{k \in \partial a \setminus j} \eta_{k \rightarrow a}(\sigma_k) \quad (55)$$

VI. APPLICATION: THE ISING FERROMAGNET ON A RANDOM GRAPH

We can now solve a first non trivial problem: the Ising ferromagnet on a Random Graph. For simplicity, we shall consider the regular case with a fixed value of c .

In the high temperature phase, we find that $\eta_{i \rightarrow j} = 1/2$. A simple computation leads to the free energy:

$$-\beta f = \log 2 + \frac{c}{2} \log \cosh \beta \quad (56)$$

Below the critical temperature, we must now use the value of η that is the fixed point of

$$\eta_c = \frac{[\eta_c e^\beta + (1 - \eta_c) e^{-\beta}]^{c-1}}{[\eta_c e^\beta + (1 - \eta_c) e^{-\beta}]^{c-1} + [\eta_c e^{-\beta} + (1 - \eta_c) e^\beta]^{c-1}} \quad (57)$$

which is fully equivalent to the previous equation if one sets $\eta_c = \frac{1 + \tanh \beta m}{2}$. From then, we can determine the total value of the magnetization, and the value of the free energy.

Note that the validity of this approach was proven rigorously by Montanari and Dembo [9]. However, we know this is not generic and that things can be (much) more complicated for “glassy” or “complex” systems: in particular, there could be many (exponentially many) solution to those equations [10].

VII. RANDOM CONSTRAINT SATISFACTION PROBLEMS: THE COLORING PROBLEM

We shall now discuss a constraint satisfaction problem, and start by one which is particularly convenient: the coloring problem.

The Graph Coloring problem (COL) is a very basic and famous problems in combinatorics [11] and in statistical physics[12]. Given a graph, or a lattice, and given a number q of available colors, the problem consists in finding a coloring of vertices such that no two neighboring vertices have the same color. The minimally needed number of colors is the chromatic number of the graph. For planar graphs there exists a famous theorem [2] showing that four colors are sufficient, and that a coloring can be found by an efficient algorithm. On the contrary, for general graphs the problem is computationally hard to solve: already in 1972 it was shown that Graph Coloring is NP-complete [13] which means, roughly speaking, that the time required for determining the existence of a proper coloring grows exponentially with the graph size.

Two main problems: the $q \log q$ and the $2q \log q$ problems. Hard/Easy transition and Possible/Impossible transition. The same happens in SAT.

So, what should we do? Let us first do the annealed computation.

Then, let us apply cavity. We find the same !!! In fact, this is the rigorous results at low connectivity (Bandyopadhyay and Gamarnik). So there is a phase transition somewhere

$$-\beta f = c \log q + \frac{c}{2} \log 1 - \frac{1 - e^{-\beta}}{q} \quad (58)$$

$$s = \log q + \frac{c}{2} \log 1 - \frac{1}{q} \quad (59)$$

Bounds...

There is a glassy phase

Appearance of many many solution

A simple argument in the large connectivity limit

We recover $q \log q$ and the $2q \log q$.

Glass transition. In the Glassy phase we use Survey propagation instead of Belief Propagation (replica symmetric breaking).

Note that BP can be use on a single graph as solver (decimation) and as a way to estimate the entropy. Difference between average cases and a given graph.

VIII. WHERE DO WE STAND?: A REVIEW

In the last section, we perform a kind of review on the knowledge we have now of hard random constraint satisfaction problems. This chapter can be skipped by most readers...

A. Introduction

Spin glass theory has a large and probably initially unexpected impact on some problems far from condensed matter physics and one example of such spectacular outcome is the application of statistical physics ideas to combinatorial optimization and of the concept of phase transitions to the probabilistic analysis of Constraint Satisfaction Problems (CSPs) [14–16]. Given a set of N discrete variables subject to a set of M constraints, a CSP consists in deciding if there exists an assignment of the variables satisfying all the constraints. This is a generic setting that is currently used to tackle problems as diverse as, among others, error-correcting codes, register allocation in compilers or genetic regulatory networks. The class of NP-complete problems [11], for which no algorithm is known that guarantees to decide satisfiability in a time polynomial in N , is particularly interesting. Well-studied examples of such problems are the satisfiability of boolean formulas (SAT), and the q -coloring problem (q -COL, see figure 7) that we shall discuss here. Given a graph with N vertices and M edges connecting certain pairs of them, and given q colors, can we color the vertices so that no two connected vertices have the same color?

Crucial empirical observations were made when considering the ensemble of random graphs with a given average vertex connectivity c : while below a critical value c_s a proper q -coloring of the graph exists with a probability going to one in the large size limit, it was found that beyond c_s no proper q -coloring exists asymptotically. This sharp threshold (which appears in other CSPs such as K-SAT and whose existence is partially proved in [17]) is an example of a phase transition arising in random CSPs. It was also observed empirically [18, 19] that deciding colorability becomes on average much harder near to the coloring threshold c_s than far away from it. It is therefore natural to ask ourselves: Can the value of the colorable/uncolorable (COL/UNCOL) phase transition be computed? Can the number of all possible colorings be also computed? Are there other interesting phase transitions? Can these transitions explain the fact that solutions are sometimes very hard to find? Can this knowledge help us in designing new algorithms? These questions, and their answers, are at the roots of the interest of the statistical physics community in optimization problems [15, 16].

B. A Potts anti-ferromagnet on random graphs

It is immediate to realize that the q -coloring problem is equivalent to the question of determining if the ground-state energy of a Potts anti-ferromagnet on a random graph is zero or not [20]. Consider indeed a graph $G = (\mathcal{V}, \mathcal{E})$ defined by its vertices $\mathcal{V} = \{1, \dots, N\}$ and edges $(i, j) \in \mathcal{E}$ which connect pairs of vertices $i, j \in \mathcal{V}$; and the Hamiltonian

$$\mathcal{H}(\{s\}) = \sum_{(i,j) \in \mathcal{E}} \delta(s_i, s_j). \quad (60)$$

With this choice there is no energy contribution for neighbors with different colors, but a positive contribution otherwise. The ground state energy (the energy at zero temperature) is thus zero *if and only if* the graph is q -colorable. This transforms the coloring problem into a well-defined statistical physics model. Usually, two types of random graphs are considered: in the c -regular ensemble all points are connected to exactly c neighbors, while in the Erdős-Rényi case the connectivity has a Poisson distribution.

C. Cavity method: Warnings, Beliefs and Surveys

Over the last few years, a number of studies have investigated CSPs following the adaptation of the so-called cavity method [14] to random graphs [10, 16]. It is a powerful heuristic tool —whose exactness is widely accepted but has still to be rigorously demonstrated— equivalent to the replica method of disordered systems [14]. Its main idea lies in the fact that a large random graph is locally tree-like, and that an iterative procedure known in physics as the Bethe-Peirls method can solve exactly any model on a tree (such models are often qualified as “mean field” in physics). Interestingly, it was realized [6] that an equivalent formalism has been developed independently in computer science [7], where it is called Belief Propagation (BP, which conveniently enough, may also stand for Bethe-Peirls).

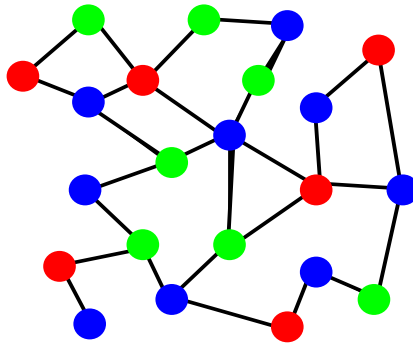


FIG. 7: Example: a proper 3-coloring of a small graph.

Defining $\psi_c^{i \rightarrow j}$ ($c = 1, \dots, q$) as the probability that the spin i has color c in absence of the spin j (the “belief” that the spin j has on the properties of the spin i), BP reads

$$\psi_c^{i \rightarrow j} = \frac{1}{Z_0^{i \rightarrow j}} \prod_{k \in N(i) \setminus \{j\}} (1 - \psi_c^{k \rightarrow i}) \quad (61)$$

where $Z_0^{i \rightarrow j}$ is a normalization constant and the notation $k \in N(i) \setminus \{j\}$ means the set of neighbors of i except j . From a fixed point of these equations, the complete beliefs in presence of all spins can be also computed. They give, for each vertex, the probability of each color from which other quantities, as for instance the number of solutions, can be computed. A simpler formalism, called Warning Propagation (WP), restricts itself to frozen variables (*i.e.* to variables for which only one color can satisfy the constraints). However, WP does not allow to compute the number of solutions, only their existence, but is definitely simpler to handle.

It was soon realized, however, that these methods developed for trees could not be used straightforwardly on all random graphs because of a non-trivial phenomenon called clustering [10, 21] (for which rigorous results are now available, see [22]). Indeed, while for graphs with very low connectivities all solutions are “connected” —in the sense that it is easy with a local dynamics to move from one solution to another— they regroup into a large number of disconnected clusters for larger connectivities. It can be argued that each of these clusters corresponds to a different fixed point of the BP equations, so that a survey over the whole set of the fixed points should be performed. This can be done in the cavity method by the now famous Survey Propagation (SP) equations [16] which, in the physics language, correspond to the Parisi’s one-step Replica Symmetry Breaking (RSB) scheme [14]. Within this formalism, the number of clusters (which behaves as $\mathcal{N} = e^{N\Sigma}$, where Σ is a fundamental quantity called the *complexity*) and their sizes (the number of proper colorings inside the cluster) can be determined.

This formalism has been applied on the SAT [16] and COL [23–26] problems in the limit of infinitely large graphs. These cutting edge studies were however restricted to SP applied to the clusters corresponding to fixed points of WP and not to those of BP. Although this already allowed the correct computation of the COL/UNCOL transition and the development of a powerful algorithm [16], it meant that the description of the clustered phase was only partial and this resulted in a number of problems and inconsistencies that stayed unanswered until very recently. These issues have been today clarified [25, 27, 28] and we shall now discuss this new understanding.

D. The phase diagram of the coloring problem on a random graph

Consider that we have $q \geq 4$ colors (the $q = 3$ case being a bit particular [25, 26], as we shall see) and a large random graph whose connectivity c we shall increase. Different phases are encountered that we will now describe (and enumerate) in order of appearance (the corresponding phase diagram is depicted in figure 8).

- (i) **A unique cluster exists:** For low enough connectivities, all the proper colorings are found in a single cluster, where it is easy to “move” from one solution to another. Only one possible —and trivial— fixed point of the BP equations exists at this stage (as can be proved rigorously in some cases [29]). The entropy can be computed and reads in the large graph size N limit

$$s = \frac{\log \mathcal{N}_{\text{sol}}}{N} = \log q + \frac{c}{2} \log \left(1 - \frac{1}{q} \right). \quad (62)$$

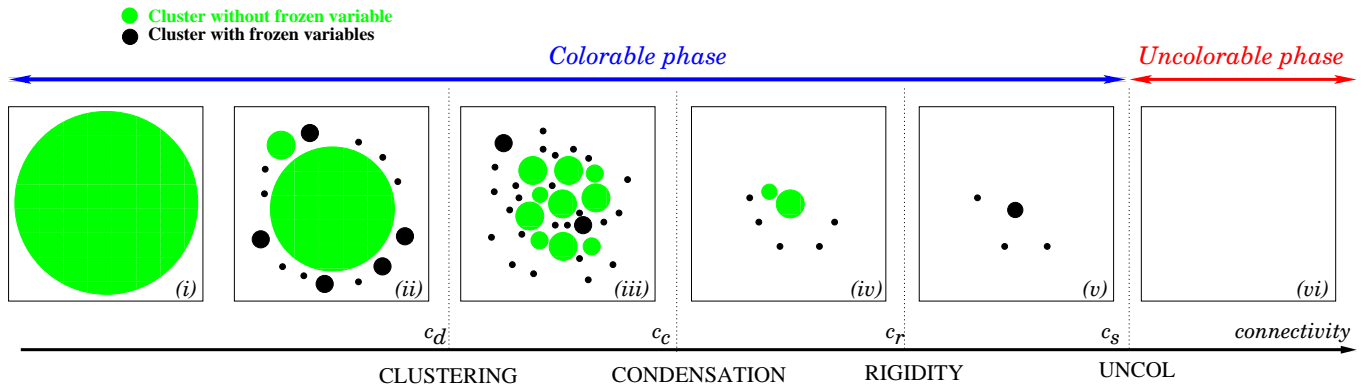


FIG. 8: Sketch of the space of solutions —colored points in this representation— in the q -coloring problem on random graphs when the connectivity c is increased. (i) At low c , all solutions belong to a single cluster. (ii) For larger c , other clusters of solutions appear but a giant cluster still contains almost all solutions. (iii) At the clustering transition c_d , it splits into an exponentially large number of clusters. (iv) At the condensation transition c_c , most colorings are found in the few largest of them. (v) The rigidity transition c_r ($c_r < c_c$ and $c_r > c_c$ are both possible depending on q) arises when typical solutions belong to clusters with frozen variables (that are allowed only one color in the cluster). (vi) No proper coloring exists beyond the COL/UNCOL threshold c_s .

- (ii) **Some (irrelevant) clusters appear:** As the connectivity is slightly increased, the phase space of solutions decomposes into an large (exponential) number of different clusters. It is tempting to identify that as the clustering transition, but it happens that all (but one) of these clusters contain relatively very few solutions —as compare to whole set— and that almost all proper colorings still belong to one single giant cluster. Clearly, this is not a proper clustering phenomenon and in fact, for all practical purpose, there is still only one single cluster. equation (62) still gives the correct entropy at this stage.
- (iii) **The clustered phase:** For larger connectivities, the large single cluster also decomposes into an exponential number of smaller ones: this now defines the genuine clustering threshold c_d^2 . Beyond this threshold, a local algorithm that tries to move in the space of solutions will remain prisoner of a cluster of solutions. Interestingly, it can be shown that the total number of solutions is still given by equation (62) in this phase. This is because, as is well known in the replica method, the free energy has no singularity at the dynamical transition (which is therefore not a true transition in the sense of Ehrenfest, but rather a dynamical or geometrical transition in the space of solutions).
- (iv) **The condensed phase:** As the connectivity is further increased, a new sharp phase transition arises at the condensation threshold c_c where most of the solutions are found in a finite number of the largest clusters. From this point, equation (62) is not valid anymore and becomes just an upper bound. The entropy is non-analytic at c_c therefore this is a genuine static phase transition.
- (v) **The rigid phase:** As mentioned in section VIII C, two different types of clusters exist: In the first type, that we shall call the *unfrozen* ones, all spins can take at least two different colors. In the second type, however, a finite fraction of spins is allowed only one color within the cluster and are thus “frozen” into this color. These *frozen* clusters actually correspond to non-trivial fixed points of BP *and* WP, while the first kind are non-trivial fixed points of BP *only*. It follows that a transition exists, that we call *rigidity*, when frozen variables appear inside the dominant clusters (those that contains most colorings). If one takes a proper coloring at random beyond c_r , it will belong to a cluster where a finite fraction of variables is frozen into the same color. Depending on the value of q , this transition may arise before or after the condensation transition (see table III).
- (vi) **The UNCOL phase:** Eventually, the connectivity c_s is reached beyond which no more solutions exist. The ground state energy (sketched in figure 9) is zero for $c < c_s$ and then grows continuously for $c > c_s$. The values

² It is important to point out that the location of the clustering transitions was therefore not computed correctly when the dependence on the size of clusters was not taken into account. Also, different results were obtained previously depending on whether or not unfrozen fixed points were explicitly considered.

q	c_d	c_r	c_c	c_s
3	5^+	-	6	6
4	9	-	10	10
5	14	14	14	15
6	18	19	19	20
7	23	-	25	25
8	29	30	31	31
9	34	36	37	37
10	39	42	43	44

q	c_d	c_r	c_c	c_s
3	4	4.66(1)	4	4.687(2)
4	8.353(3)	8.83(2)	8.46(1)	8.901(2)
5	12.837(3)	13.55(2)	13.23(1)	13.669(2)
6	17.645(5)	18.68(2)	18.44(1)	18.880(2)
7	22.705(5)	24.16(2)	24.01(1)	24.455(5)
8	27.95(5)	29.93(3)	29.90(1)	30.335(5)
9	33.45(5)	35.658	36.08(5)	36.490(5)
10	39.0(1)	41.508	42.50(5)	42.93(1)

TABLE III: Threshold connectivities c_d (dynamical/clustering) [25, 26, 32, 33], c_r (rigidity/freezing) [26, 34], c_c (condensation/Kauzmann) [25, 26] and c_s (COL/UNCOL) [23, 24] for regular (left) and Erdős-Rényi (right) random graphs. In the large q -limit, one finds in both cases that [25, 26]: $c_r = q[\log q + \log \log q + 1 + o(1)]$, $c_c = 2q \log q - \log q - 2 \log 2 + o(1)$ and [24]: $c_s = 2q \log q - \log q - 1 + o(1)$.

c_s computed within the cavity formalism are in perfect agreement with the rigorous bounds [30] derived using probabilistic methods and are widely believed to be exact (although they remains to be rigorously proven, but see [31] for a proof that they are at least rigorous upper bounds).

We report the values of the threshold connectivities corresponding to all these transitions in table III for the regular and the Poissonian (*i.e.* Erdős-Rényi) random graphs ensembles. Notice that the 3-coloring is peculiar because $c_d = c_c$ so that the clustered phase is always condensed in this case. In view of this rich phase diagram, it is important to get an intuition on the meaning and the properties of these different phases and, in this respect, it is interesting before entering the algorithmic implications to discuss the analogies with the glass transition.

E. A detour into the ideal glass transition phenomenology

To those familiar with the replica theory and the mean field theory of glasses, the phenomenology depicted in the former section should look familiar: these successive transitions are indeed very well known in the picture of the ideal glass transition [35]. This striking analogy is in fact quite natural since, despite the fact that there is no disorder in the interactions in Hamiltonian (60), the frustration due to the loops in the random graph makes the model behaving like a disordered “anti-ferromagnetic” Potts spin glass [36] and such models are known to display the glassy phenomenology [20, 35].

The phase diagram obtained on Poissonian random graphs with average connectivity c for $q \geq 4$ is sketched in figure 9 (the $q = 3$ model is slightly different as, again, $T_d = T_c$). At high temperature the system behaves as a liquid (or a paramagnet in the language of magnetic systems). Below a temperature T_d a first transition —called “dynamical”— happens and the system falls out of equilibrium. For $T < T_d$ it is not possible for a physical dynamics to equilibrate the system and the ergodicity is broken: this is due to the appearance of exponentially many different states. However, the would-be equilibrium properties of the problem remain similar (and in particular, the free-energy has no singularity at this temperature). Only at temperature T_c the free energy is non analytic and a true “static” glass transition happens, called the Kauzmann transition [35]. In this phase only a finite number of states does matter at a given connectivity. Finally, for larger connectivities, a third phenomenon is observed as the temperature is further lowered, called the Gardner transition [37]. It is a transition towards a more complicated phase, similar to the one found in the celebrated solution of the Sherrington-Kirkpatrick model [14]. The fact that the Gardner transition arises for connectivities *larger* the COL/UNCOL one is very important in this respect: it shows that the study of this phase, that requires a more involved cavity formalism (and probably further RSB), does not seem to be needed in the colorable phase³. We also now recognize that the “clustering” and the “condensation” transitions in the coloring problem are just the zero temperature relics of the dynamic and Kauzmann transitions at finite temperatures.

³ The expert reader might find this puzzling, as many papers stated that the simple “one replica symmetry breaking” [10, 16] solution was unstable towards a more complex solution in some region of the COL/SAT phase [24, 38]. However, these results were obtained neglecting the role of the sizes of the clusters: while in some cases *most* clusters are indeed unstable, our studies [36] indicate that the *relevant ones* seem always stable in the COL phase (although the cases of 3-COL and 3-SAT might be problematic, see [26, 36]).

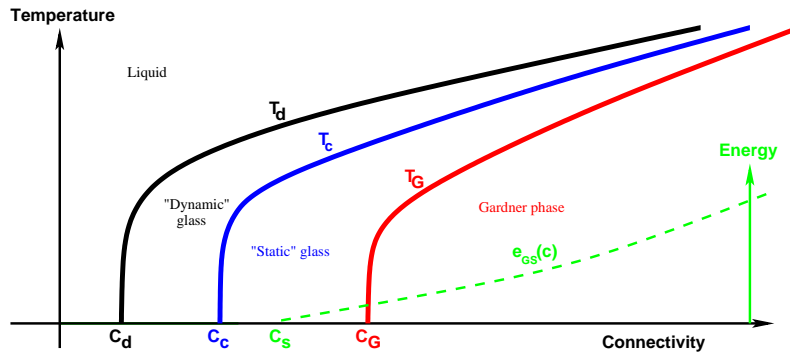


FIG. 9: Sketch of the phase diagram in the coloring problem at finite temperature (from [36]). At T_d , the system falls out of equilibrium (“dynamic” transition). At T_c the system undergoes a “static” glass transition. Finally, at T_g , a Gardner transition appears. e_{GS} represents the ground state energy.

A similar connection with the physics of glassy system can also be obtained directly at zero temperature via the jamming phenomenology [39] where the density of constraints (in this case the volume of some non-overlapping spheres in a box of fixed volume) are increased and where a dynamical transition is first met while some authorized configurations exist much beyond this point. We thus see that the coloring problem on random graphs translates into a very general mean field model of a complex liquid. This convergence of interest between different disciplines is quite interesting in itself and allows to discuss a number of matters, as we shall now see.

F. Onset of hardness for local search algorithms

The properties of the phase diagram we just discussed are based on analytical computations through the cavity method. We would like to discuss now what are the implications of these different phases on the performance of simple local search algorithms that try to find a solution. This is, however, a much harder subject to handle analytically and we shall thus leave the field of analytical computations to enter the one of phenomenology. Still, the behavior of an algorithm trying to find a solution is reminiscent of the behavior of the physical dynamics in glassy systems, and we can at least exploit this analogy in order to get an intuition for the problem that we can later confirm with numerical simulations.

It is first tempting to identify the point c_d , where a physical Monte-Carlo dynamics gets trapped into a cluster, with the onset of computational hardness⁴. However, a second moment though indicates that this should not be the case: In the glass transition phenomenology, it is well known that, although the system falls out of equilibrium beyond T_d , its energy can be further lowered by lowering the temperature, or just by waiting a bit longer [39]. In short: the fact that dynamics is prisoner of a given region of the set of possible configurations does not mean that no solution can be found in this region. Although c_d is indeed a sharp transition for the Monte-Carlo sampling, there is no reason, a priori, to experience difficulties if one just want to find one solution beyond this point. This is particularly transparent in the analysis and the algorithm introduced in [39] (and directly inspired from the analogy with jamming):

1. Start with a graph of connectivity c and a proper coloring.
2. Increase the density of constraint by adding a link in the graph.
3. Use a simple algorithm in order to solve the contradiction introduced by the link. When it is done, go back to step (i).

By applying this strategy, starting from scratch (*i.e.* from a graph with N vertices and *no* link), the set of all proper colorings undergoes the successive transitions described in figure 8 as connectivity increases. When the dynamical transition is reached, one is trapped inside a cluster of solutions, but this is not really a problem as one is still free to move inside the cluster. As more links are added, the cluster size is decreasing continuously but while it still exists, the local algorithm should be in principle able to find solutions nearby. Only for larger connectivities, when the cluster gets frozen, it disappears and consequently the algorithm stops. It was shown in Ref. [39] through numerical simulations that this strategy, using the Walk-COL [26] algorithm for step (iii), is indeed efficient, and linear in N , much beyond the dynamical transition.

⁴ When the clustering phenomena was discovered in CSPs, it was indeed initially conjectured to be responsible for the onset of hardness for local search strategies [16, 21] and the clustered phase was named the “hard phase”. However, some local algorithms were found to easily beat the threshold (see [40] for SAT and [26] for COL).

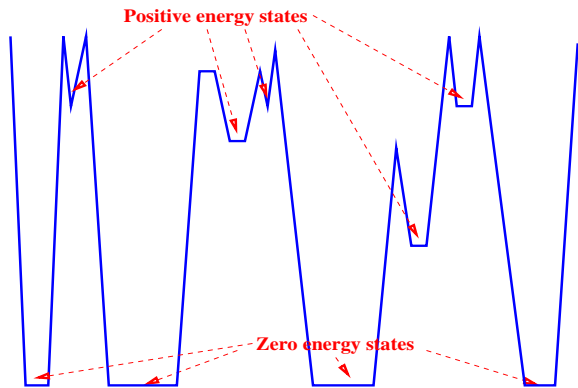


FIG. 10: Artist's view of the energy landscape for low connectivities $c > c_d$: a region dominated by canyons that reach the ground-states.

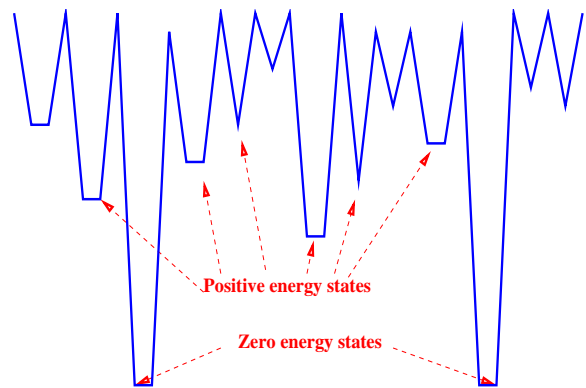


FIG. 11: Artist's view of the energy landscape for large connectivities $c > c_d$: a region dominated by high mountains and deep valleys.

The reason why this recursive strategy becomes inefficient when the cluster in which the dynamics is trapped freezes is the following: if a link is put between two vertices frozen in the same color, it is impossible to satisfy the constraints while remaining in the cluster. As opposed to the unfrozen clusters, the frozen clusters thus have a finite probability to disappear when a new link is added. A cavity-like analysis [34], confirmed by numerical data [39], actually shows that the number of changes that the algorithm must perform, in order to solve the contradictions imposed by the addition of new links, increases with the connectivity and diverges when the frozen variables appear. The source of difficulties is therefore not the clustering phenomenon in itself, but rather the appearance of frozen variables. This makes the analysis and the prediction of a EASY/HARD threshold much harder since (as one can see on figure 8) clusters of different sizes freeze at different connectivities, although a connectivity $c_* \geq c_r$ exists where all clusters are frozen, thus putting a strict bound to the efficiency of this procedure.

Interestingly, even non-incremental algorithms may also pass the c_d threshold [26, 40], and that might come as a surprise for those having in mind the “rugged” many-valley energy landscape picture of spin glasses. This apparent paradox can be clarified by the following considerations: It is possible that at lower connectivity $c > c_d$ the energy landscape is dominated by deep canyons (figure 10), where it is in principle easy to go down as one has just to jump ahead! At larger connectivities a more rugged region with many deep valleys and high mountains is found (figure 11) in which case, as any mountain-hiker will undoubtedly know, it takes some time to go to the deepest valley because many hills have to be climbed first. This difference in behavior might explain the “unreasonable efficiency” of local algorithm and the performance of the annealing procedure beyond c_d [32].

To further illustrate this point, consider the Walk-COL algorithm introduced in [26] (and adapted from a similar one in SAT [40]) defined by the following procedure

- (i) Randomly choose a spin that has the same color as at least one of its neighbors.
- (ii) Change randomly its color. Accept this change with probability one if the number of unsatisfied spins has been lowered, otherwise accept it with probability p (this is a parameter that has to be tuned for better efficiency).
- (iii) If there are unsatisfied vertices, go to step (i) unless the maximum running time is reached.

This algorithm can easily find colorings for large sizes in linear time beyond c_d [26], but certainly not too close to the UNCOL transition where it gets trapped at higher energies (see figure 12).

So far, there are few analytical results about the energy landscape in this problem and it is likely that this will be the subject of further studies. It is unfortunately very hard to say for which connectivities the landscape goes from canyons-dominated to mountains-dominated as this may not be a sharp transition and more a matter of —certainly algorithm-dependent— basins of attraction. The rigidity transition for typical clusters is certainly a good candidate as a crossover in behavior (as is the connectivity where *all* clusters are frozen).

To conclude, one sees that although the algorithmic issues are indeed more difficult to handle than the phase diagram, at least two important points can already be made: First, the dynamical transition does not correspond to the onset of hardness, and second, the source of difficulty seems more to be related with the appearance of frozen variables.

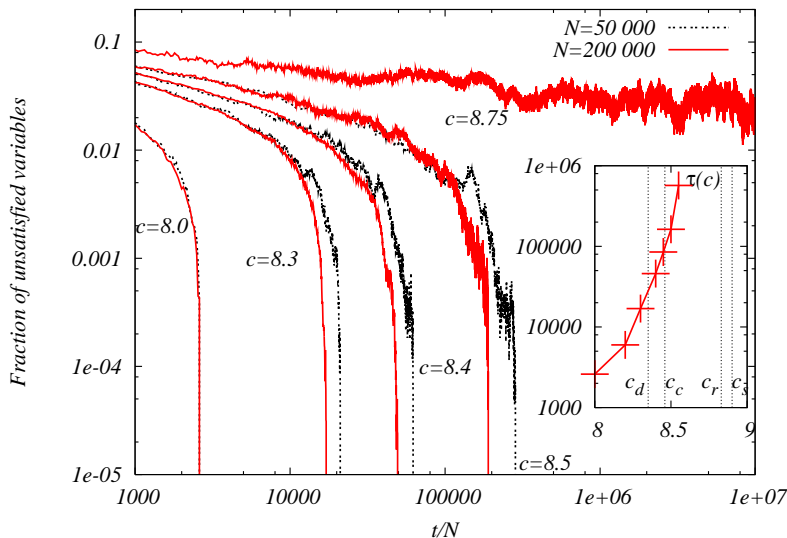


FIG. 12: Fraction of unsatisfied variables versus the number of attempted flips of the Walk-COL algorithm divided by the size of the graph in the $q = 4$ coloring: Walk-COL [26] algorithm is able to find some solutions in the clustered phase for low c (in the canyon-like region), but get trapped in the high energy valleys for larger c . Inset: estimated time $\tau = t/N$ needed to find a solution versus connectivity.

G. Message Passing and Decimation

The class of local search algorithms is only one part of the story. A different class, where messages are exchanged through the nodes of the graph, was proven to be very efficient in computer science. BP and WP are examples of such procedures and it is thus interesting to use the information given by the cavity analysis to discuss their performance in estimating the marginals—or other informations—in the problem. A major outcome of the last years has also been the application of SP as a message passing (MP) [16].

From the information given by the fixed point of the MP, an algorithm can be defined in the following way [16]: (i) run the MP to obtain some information, and (ii) fix (decimate) some variables according to this information. This sequence is iterated until a solution is found, or until a contradiction is met. The two parts are quite independent and each of them can be changed separately (for instance a self consistent re-enforcement has been tried instead of the decimation with very good results [41]).

According to the cavity formalism, the iterative fixed point of BP correctly estimates the marginals until the condensation at c_c [25] (which, conveniently, is very close to c_s). Indeed the application of BP plus decimation is numerically efficient in finding solutions for both SAT and COL [25, 26, 42] and it would be interesting to see if this method allows to find “typical” solutions beyond c_d , thus bypassing the problems of Monte-Carlo algorithms. It seems that the strategy is efficient even *beyond* c_c in some cases [26], although the BP recursion does not always converge (and when it does, it does it slowly) on the decimated graph, in which case an imprecise and approximate information is obtained. These issues are thus far from being properly understood at the present time, and we hope that new works will be done in this direction.

A more powerful strategy in the clustered phase is to use SP to compute the probability that a given variable is frozen in a cluster [16]. Fixing the variables which are frozen in most clusters seems a good way to decimate the graph. Using this method in random 3-SAT allows to outperform any other known algorithm and to find solutions of huge instances rapidly even very close to the threshold [16, 41]. The method has been subsequently adapted for the 3-coloring in [23]. Interestingly, SP behaves better than BP on the decimated graph and has no problem of convergence. Other strategies are possible that have not yet been exploited [26].

The best way to extract the informations given by the fixed point of these message passing procedures is however not yet known, nor is the limit until where these strategies are efficient. However, this line of thoughts is certainly the most promising way in the direction of better solving and sampling algorithms.

H. Conclusions and perspectives

In this paper, we considered and reviewed some aspects of the q -coloring of large random graphs. We discussed the properties of the set of solutions and the different sharp phase transitions it undergoes when the average connectivity is varied. The problem translates in physics into a mean-field model with an ideal glass transition. The cavity method has been efficient in giving insight into the problem, although the important challenge of proving rigorously these results remains. It would be interesting in this respect to confirm these predictions by performing extensive

Monte-Carlo simulations of the phase diagram depicted in figure 9.

We also discussed the dynamical behavior of local search algorithms. We saw that although the “dynamical” transition has a direct meaning as the point where local Monte-Carlo algorithms get out of equilibrium, it is not directly connected with the onset of hardness in the problem which is rather due to the fact that clusters “freeze” as the connectivity increases. So far, it has not therefore been possible, despite initial hope, to have well-defined HARD and EASY phases. The precise role of frozen variables and the part played by the rigidity transition, or by the connectivity where all clusters becomes frozen, will undoubtedly be the subject of new research, both in numerical and theoretical directions, that will help to clarify these issues. A refined knowledge of the energy landscape would also be valuable.

Finally, we also discuss the major breakthrough in the algorithmic strategy that emerged from the application of the cavity solution on a single given graph. It is not clear at the present time what is the best way to use this approach and how efficient it can be, nor if it is possible to go arbitrarily close to the satisfiability/colorability threshold for any values of q , and it is likely that these questions will also trigger a lot of work in the future.

-
- [1] M. Mézard and A. Montanari, *Physics, Information, Computation* (Oxford Press, Oxford, 2009).
 - [2] K. Appel and W. Haken, *Illinois J. Math.* **21** (1977).
 - [3] K. Appel and W. Haken, *Illinois J. Math.* **21** (1977).
 - [4] S. Kirkpatrick, C. D. Gelatt Jr., and M. P. Vecchi, *Science* **220**, 671 (1983).
 - [5] S. Kirkpatrick and B. Selman, *Science* **264**, 1297 (1994).
 - [6] J. Yedidia, W. Freeman, and Y. Weiss, in *Exploring Artificial Intelligence in the New Millennium* (Science & Technology Books, 2003), pp. 239–236.
 - [7] J. Pearl, in *Proceedings American Association of Artificial Intelligence National Conference on AI* (Pittsburgh, PA, USA, 1982), pp. 133–136.
 - [8] H. A. Bethe, *Proc. Roy. Soc. London A* **150**, 552 (1935).
 - [9] A. Dembo and A. Montanari (2008), arXiv:0804.4726v2 [math.PR].
 - [10] M. Mézard and G. Parisi, *Eur. Phys. J. B* **20**, 217 (2001).
 - [11] M. Garey and D. Johnson, *Computers and intractability: A guide to the theory of NP-completeness* (Freeman, San Francisco, 1979).
 - [12] F. Y. Wu, *Rev. Mod. Phys.* **54**, 235 (1982).
 - [13] R. Karp, in *Complexity of Computer Computations*, edited by R. Miller and J. Thatcher (Plenum Press, New-York, 1972), pp. 85–103.
 - [14] M. Mézard, G. Parisi, and M. A. Virasoro, *Spin-Glass Theory and Beyond*, vol. 9 of *Lecture Notes in Physics* (World Scientific, Singapore, 1987).
 - [15] R. Monasson, R. Zecchina, S. Kirkpatrick, B. Selman, and L. Troyansky, *Nature* **400**, 133 (1999).
 - [16] M. Mézard, G. Parisi, and R. Zecchina, *Science* **297**, 812 (2002).
 - [17] E. Friedgut, *J. Amer. Math. Soc.* **12** (1999).
 - [18] P. Cheeseman, B. Kanefsky, and W. M. Taylor, in *Proc. 12th IJCAI* (Morgan Kaufmann, San Mateo, CA, USA, 1991), pp. 331–337.
 - [19] B. Selman, D. G. Mitchell, and H. J. Levesque, *Artif. Intell.* **81**, 17 (1996), ISSN 0004-3702.
 - [20] I. Kanter and H. Sompolinsky, *J. Phys. A: Math. Gen* **20**, L673 (1987).
 - [21] G. Biroli, R. Monasson, and M. Weigt, *Eur. Phys. J. B* **14**, 551 (2000).
 - [22] T. Mora, Ph.D. thesis, Université Paris-Sud (2007), <http://tel.archives-ouvertes.fr/tel-00175221/en/>.
 - [23] R. Mulet, A. Pagnani, M. Weigt, and R. Zecchina, *Phys. Rev. Lett.* **89**, 268701 (2002).
 - [24] F. Krzakala, A. Pagnani, and M. Weigt, *Phys. Rev. E* **70**, 046705 (2004).
 - [25] F. Krzakala, A. Montanari, F. Ricci-Tersenghi, G. Semerjian, and L. Zdeborová, *Proc. Natl. Acad. Sci. U.S.A* **104**, 10318 (2007).
 - [26] L. Zdeborová and F. Krzakala, *Phys. Rev. E* **76**, 031131 (2007).
 - [27] L. Zdeborová and M. Mézard, *Phys. Rev. Lett.* **101**, 078702 (2008).
 - [28] L. Zdeborová and M. Mézard, *J. Stat. Mech.* p. P12004 (2008).
 - [29] A. Bandyopadhyay and D. Gamarnik, in *Proc. of the 17th ACM-SIAM Symposium on Discrete Algorithms* (ACM Press, New York, USA, 2006), pp. 890 – 899.
 - [30] T. Luczak, *Combinatorica* **11**, 45 (1991).
 - [31] S. Franz and M. Leone, *J. Stat. Phys.* **3-4**, 535 (2003).
 - [32] J. van Mourik and D. Saad, *Phys. Rev. E* **66**, 056120 (2002).
 - [33] M. Mézard and A. Montanari, *J. Stat. Phys.* **124**, 1317 (2006).
 - [34] G. Semerjian, *J. Stat. Phys.* **130**, 251 (2008).
 - [35] D. Gross and M. Mézard, *Nucl. Phys. B* **240**, 431 (1984).
 - [36] F. Krzakala and L. Zdeborová, *Europhys. Lett.* **81**, 57005 (2008).
 - [37] E. Gardner, *Nuclear Physics B* **257**, 747 (1985).

- [38] A. Montanari, G. Parisi, and F. Ricci-Tersenghi, J. Phys. A **37**, 2073 (2004).
- [39] F. Krzakala and J. Kurchan, Phys. Rev. E **76**, 021122 (2007).
- [40] J. Ardelius and E. Aurell, Phys. Rev. E **74**, 037702 (2006).
- [41] J. Chavas, C. Furtlehner, M. Mézard, and R. Zecchina, J. Stat. Mech. p. P11016 (2005).
- [42] A. Montanari, F. Ricci-Tersenghi, and G. Semerjian, J. Stat. Mech. p. P04004 (2008).