

TOPICAL REVIEW

Cluster Variation Method in Statistical Physics and Probabilistic Graphical Models

Alessandro Pelizzola

Dipartimento di Fisica, Politecnico di Torino, c. Duca degli Abruzzi 24, 10129 Torino, Italy and INFN, Sezione di Torino

Abstract.

The cluster variation method (CVM) is a hierarchy of approximate variational techniques for discrete (Ising-like) models in equilibrium statistical mechanics, improving on the mean-field approximation and the Bethe-Peierls approximation, which can be regarded as the lowest level of the CVM. In recent years it has been applied both in statistical physics and to inference and optimization problems formulated in terms of probabilistic graphical models.

The foundations of the CVM are briefly reviewed, and the relations with similar techniques are discussed. The main properties of the method are considered, with emphasis on its exactness for particular models and on its asymptotic properties.

The problem of the minimization of the variational free energy, which arises in the CVM, is also addressed, and recent results about both provably convergent and message-passing algorithms are discussed.

PACS numbers: 05.10.-a, 05.50.+q, 89.70.+c

Submitted to: *J. Phys. A: Math. Gen.*

E-mail: alessandro.pelizzola@polito.it

1. Introduction

The Cluster Variation Method (CVM) was introduced by Kikuchi [1] in 1951, as an approximation technique for the equilibrium statistical mechanics of lattice (Ising-like) models, generalizing the Bethe–Peierls [2, 3] and Kramers–Wannier [4, 5] approximations, an account of which can be found in several textbooks [6, 7]. Apart from rederiving these methods, Kikuchi proposed a combinatorial derivation of what today we can call the cube (respectively triangle, tetrahedron) approximation of the CVM for the Ising model on the simple cubic (respectively triangular, face centered cubic) lattice.

After the first proposal, many reformulations and applications, mainly to the computation of phase diagram of lattice models in statistical physics and material science, appeared, and have been reviewed in [8]. The main line of activity has dealt with homogeneous, translation-invariant lattice models with classical, discrete degrees of freedom, but several other directions have been followed, including for instance models with continuous degrees of freedom [9], free surfaces [10, 11], models of polymers [12, 13] and quantum models [14, 15]. Out of equilibrium properties have also been studied, in the framework of the path probability method [16, 17, 18], which is the dynamical version of the CVM. Despite the CVM predicts mean-field like critical behaviour, the problem of extracting critical behaviour from sequences of CVM approximations has also been considered by means of different approaches [19, 20, 21, 22, 23].

A line of research which is particularly relevant to the present discussion has considered heterogeneous and random models. Much work has been devoted in the 80's to applications of the CVM to models with quenched random interactions (see e.g. [24] and refs. therein), mainly aiming to the phase diagram, and related equilibrium properties, of Ising-like models of spin glasses in the average case. The most common approach was based on the distribution of the effective fields, and population dynamics algorithms were developed and studied for the corresponding integral equations. All this effort was however limited at the replica-symmetric level. Approaches taking into account the first step of replica symmetry breaking have been developed only recently [25], at the level of the Bethe–Peierls approximation, in its cavity method formulation, for models on random graphs in both the single instance and average case. These approaches have been particularly successful in their application to combinatorial optimization problems, like satisfiability [26] and graph coloring [27]. Another interesting approach going in a similar direction has been proposed recently [28], which relies on the analysis of the time evolution of message-passing algorithms for the Bethe–Peierls approximation.

Prompted by the interest in optimization and, more generally, inference problems, a lot of work on the CVM has been done in recent years also by researchers working on probabilistic graphical models [29], since the relation between the Bethe–Peierls approximation and the belief propagation method [30] was recognized [31]. The interaction between the two communities of researchers working on statistical physics and optimization and inference algorithms then led to the discovery of several new algorithms for the CVM variational problem, and to a deeper understanding of the method itself. There have been applications in the fields of image restoration [32, 33, 34, 35], computer vision [36], interference in two-dimensional channels [37], decoding of error-correcting codes [38, 39, 40], diagnosis [41], unwrapping of phase images [42], bioinformatics [43, 44, 45], language processing [46, 47].

The purpose of the present paper is to give a short account of recent advances on methodological aspects, and therefore applications will not be considered in detail. It is not meant to be exhaustive and the material included reflects in some way the interests of the author. The plan of the paper is as follows. In Section 2 the basic definitions for statistical mechanics and probabilistic graphical models are given, and notation is established. In Section 3 the CVM is introduced in its modern formulation, and in Section 4 it is compared with related approximation techniques. Its properties are then discussed, with particular emphasis on exact results, in Section 5. Finally, the use of the CVM as an approximation and the algorithms which can be used to solve the CVM variational problem are illustrated in Section 6. Conclusions are drawn in Section 7.

2. Statistical mechanical models and probabilistic graphical models

We are interested in dealing with models with discrete degrees of freedom which will be denoted by $\mathbf{s} = \{s_1, s_2, \dots, s_N\}$. For instance, variables s_i could take values in the set $\{0, 1\}$ (binary variables), $\{-1, +1\}$ (Ising spins), or $\{1, 2, \dots, q\}$ (Potts variables).

Statistical mechanical models are defined through an energy function, usually called Hamiltonian, $H = H(\mathbf{s})$, and the corresponding probability distribution at thermal equilibrium is the Boltzmann distribution

$$p(\mathbf{s}) = \frac{1}{Z} \exp[-H(\mathbf{s})], \quad (1)$$

where the inverse temperature $\beta = (k_B T)^{-1}$ has been absorbed into the Hamiltonian and

$$Z \equiv \exp(-F) = \sum_{\mathbf{s}} \exp[-H(\mathbf{s})] \quad (2)$$

is the partition function, with F the free energy.

The Hamiltonian is typically a sum of terms, each involving a small number of variables. A useful representation is given by the *factor graph* [48]. A factor graph is a bipartite graph made of variable nodes i, j, \dots , one for each variable, and *function nodes* a, b, \dots , one for each term of the Hamiltonian. An edge joins a variable node i and a function node a if and only if $i \in a$, that is the variable s_i appears in H_a , the term of the Hamiltonian associated to a . The Hamiltonian can then be written as

$$H = \sum_a H_a(\mathbf{s}_a), \quad \mathbf{s}_a = \{s_i, i \in a\}. \quad (3)$$

A simple example of a factor graph is reported in Figure 1, and the corresponding Hamiltonian is written as

$$H(s_1, s_2, s_3, s_4, s_5, s_6) = H_a(s_1, s_2) + H_b(s_2, s_3, s_4) + H_c(s_3, s_4, s_5, s_6). \quad (4)$$

The factor graph representation is particularly useful for models with non-pairwise interactions. If the Hamiltonian contains only 1-variable and 2-variable terms, as in the Ising model

$$H = - \sum_i h_i s_i - \sum_{(i,j)} J_{ij} s_i s_j, \quad (5)$$

then it is customary to draw a simpler graph, where only variable nodes appear, and edges are drawn between pairs of interacting spins (i, j) . In physical models the

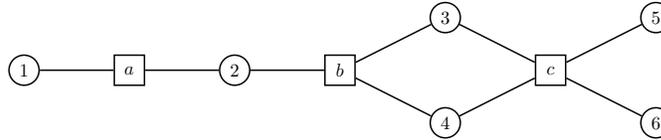


Figure 1. An example of a factor graph: variable and function nodes are denoted by circles and squares, respectively

interaction strength J_{ij} can depend on the distance between spins, and interaction is often restricted to nearest neighbours (NNs), which are denoted by $\langle i, j \rangle$.

In combinatorial optimization problems, the Hamiltonian plays the role of a cost function, and one is interested in the low-temperature limit $T \rightarrow 0$, where only minimal energy states (ground states) have a non-vanishing probability.

Probabilistic graphical models [29, 49] are usually defined in a slightly different way. In the case of *Markov random fields*, also called *Markov networks*, the joint distribution over all variables is given by

$$p(\mathbf{s}) = \frac{1}{Z} \prod_a \psi_a(\mathbf{s}_a), \quad (6)$$

where ψ_a is called *potential* (potentials involving only one variable are often called *evidences*) and

$$Z = \sum_{\mathbf{s}} \prod_a \psi_a(\mathbf{s}_a). \quad (7)$$

Of course, a statistical mechanical model described by the Hamiltonian (3) corresponds to a probabilistic graphical models with potentials $\psi_a = \exp(-H_a)$. On the other hand, *Bayesian networks*, which we will not consider here in detail, are defined in terms of directed graphs and conditional probabilities. It must be noted, however, that a Bayesian network can always be mapped onto a Markov network [29].

3. Fundamentals of the Cluster Variation Method

The original proposal by Kikuchi [1] was based on an approximation for the number of configurations of a lattice model with assigned local expectation values. The formalism was rather involved to deal with in the general case, and since then many reformulations came. A first important step was taken by Barker [50], who derived a computationally useful expression for the entropy approximation. This was then rewritten as a cumulant expansion by Morita [51, 52], and Schlijper [53] noticed that this expansion could have been written in terms of a Möbius inversion. A clear and simple formulation was then eventually set up by An [54], and this is the one we shall follow below.

The CVM can be derived from the variational principle of equilibrium statistical mechanics, where the free energy is given by

$$F = -\ln Z = \min_p \mathcal{F}(p) = \min_p \sum_{\mathbf{s}} [p(\mathbf{s})H(\mathbf{s}) + p(\mathbf{s}) \ln p(\mathbf{s})] \quad (8)$$

subject to the normalization constraint

$$\sum_{\mathbf{s}} p(\mathbf{s}) = 1. \quad (9)$$

It is easily verified that the minimum is obtained for the Boltzmann distribution

$$\hat{p}(\mathbf{s}) = \frac{1}{Z} \exp[-H(\mathbf{s})] = \arg \min \mathcal{F} \quad (10)$$

and that the variational free energy can be written in the form of a Kullback–Leibler distance

$$\mathcal{F}(p) = F + \sum_{\mathbf{s}} p(\mathbf{s}) \ln \frac{p(\mathbf{s})}{\hat{p}(\mathbf{s})}. \quad (11)$$

The basic idea underlying the CVM is to treat exactly the first term (energy) of the variational free energy $\mathcal{F}(p)$ in Equation (8) and to approximate the second one (entropy) by means of a truncated cumulant expansion.

We first define a *cluster* α as a subset of the factor graph such that if a factor node belongs to α , then all the variable nodes $i \in a$ also belong to α (while the converse needs not to be true, otherwise the only legitimate clusters would be the connected components of the factor graph). Given a cluster we can define its energy

$$H_\alpha(\mathbf{s}_\alpha) = \sum_{a \in \alpha} H_a(\mathbf{s}_a), \quad (12)$$

probability distribution

$$p_\alpha(\mathbf{s}_\alpha) = \sum_{\mathbf{s} \setminus \mathbf{s}_\alpha} p(\mathbf{s}) \quad (13)$$

and entropy

$$S_\alpha = - \sum_{\mathbf{s}_\alpha} p_\alpha(\mathbf{s}_\alpha) \ln p_\alpha(\mathbf{s}_\alpha). \quad (14)$$

Then the entropy cumulants are defined by

$$S_\alpha = \sum_{\beta \subseteq \alpha} \tilde{S}_\beta, \quad (15)$$

which can be solved with respect to the cumulants by means of a Möbius inversion, which yields

$$\tilde{S}_\beta = \sum_{\alpha \subseteq \beta} (-1)^{n_\alpha - n_\beta} S_\alpha, \quad (16)$$

where n_α denotes the number of variables in cluster α . The variational free energy can then be written as

$$\mathcal{F}(p) = \sum_{\mathbf{s}} p(\mathbf{s}) H(\mathbf{s}) - \sum_{\beta} \tilde{S}_\beta, \quad (17)$$

where the second summation is over all possible clusters.

The above equation is still an exact one, and here the approximation enters. A set R of clusters, made of maximal clusters and all their subclusters, is selected, and the cumulant expansion of the entropy is truncated retaining only terms corresponding to clusters in R . In order to treat the energy term exactly it is necessary that each function node is contained in at least one maximal cluster. One gets

$$\sum_{\beta} \tilde{S}_{\beta} \simeq \sum_{\beta \in R} \tilde{S}_{\beta} = \sum_{\alpha \in R} a_{\alpha} S_{\alpha}, \quad (18)$$

where the coefficients a_{α} , sometimes called Möbius numbers, satisfy [54]

$$\sum_{\beta \subseteq \alpha \in R} a_{\beta} = 1 \quad \forall \alpha \in R. \quad (19)$$

The above condition means that every subcluster must be counted exactly once in the entropy expansion and allows to rewrite also the energy term as a sum of cluster energies, yielding the approximate variational free energy

$$\mathcal{F}(\{p_{\alpha}, \alpha \in R\}) = \sum_{\alpha \in R} a_{\alpha} \mathcal{F}_{\alpha}(p_{\alpha}), \quad (20)$$

where the cluster free energies are given by

$$\mathcal{F}_{\alpha}(p_{\alpha}) = \sum_{\mathbf{s}_{\alpha}} [p_{\alpha}(\mathbf{s}_{\alpha}) H_{\alpha}(\mathbf{s}_{\alpha}) + p_{\alpha}(\mathbf{s}_{\alpha}) \ln p_{\alpha}(\mathbf{s}_{\alpha})]. \quad (21)$$

The CVM then amounts to the minimization of the above variational free energy with respect to the cluster probability distributions, subject to the normalization

$$\sum_{\mathbf{s}_{\alpha}} p_{\alpha}(\mathbf{s}_{\alpha}) = 1 \quad \forall \alpha \in R \quad (22)$$

and compatibility constraints

$$p_{\beta}(\mathbf{s}_{\beta}) = \sum_{\mathbf{s}_{\alpha \setminus \beta}} p_{\alpha}(\mathbf{s}_{\alpha}) \quad \forall \beta \subset \alpha \in R. \quad (23)$$

It is of great importance to observe that the above constraint set is approximate, in the sense that there are sets of cluster probability distributions that satisfy these constraints and nevertheless cannot be obtained as marginals of a joint probability distribution. An explicit example will be given in Section 5.

The simplest example is the pair approximation for a model with pairwise interactions, like the Ising model (5). The maximal clusters are the pairs of interacting variables, and the other clusters appearing in R are the variable nodes. The pairs have Möbius number 1, while for the variable nodes $a_i = 1 - d_i$, where d_i is the *degree* of node i , that is, in the factor graph representation, the number of function nodes it belongs to.

The quality of the approximation (18) depends on the value of the neglected cumulants. In the applications to lattice systems it is typically assumed that, since cumulants are related to correlations, they vanish quickly for clusters larger than the correlation length of the model. In Figure 2 the first cumulants, relative to the site (single variable) entropy, are shown for the homogeneous ($J_{ij} = J$), zero field ($h_i = 0$), square lattice Ising model, in the square approximation of the CVM.

It can be seen that the cumulants peak at the (approximate) critical point and decrease as the cluster size increases. This property is not however completely general, it may depend on the interaction range. It has been shown [55] that this does not

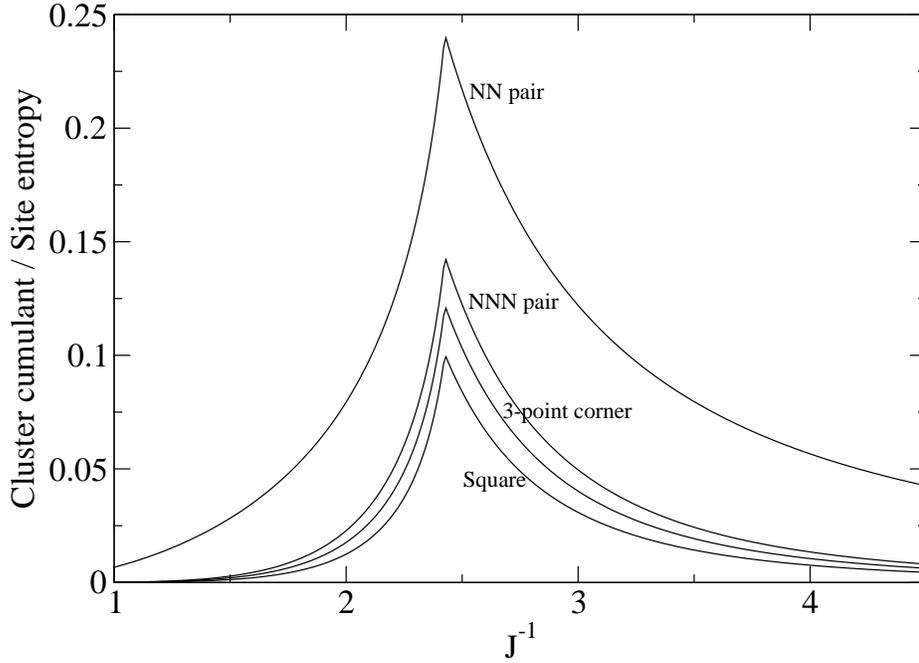


Figure 2. Cumulants for the square lattice Ising model

hold for finite instances of the Sherrington–Kirkpatrick spin–glass model, which is a fully connected model.

The meaning of cumulants as a measure of correlation can be easily understood by considering a pair of weakly correlated variables and writing their joint distribution as

$$p_{12}(s_1, s_2) = p_1(s_1)p_2(s_2) [1 + \varepsilon q(s_1, s_2)], \quad \varepsilon \ll 1. \quad (24)$$

The corresponding cumulant is then

$$\tilde{S}_{12} = S_{12} - S_1 - S_2 = -\langle \ln [1 + \varepsilon q(s_1, s_2)] \rangle = O(\varepsilon). \quad (25)$$

4. Region–based free energy approximations

The idea of *region–based free energy approximations*, put forward by Yedidia [56], is quite useful to elucidate some of the characteristics of the method, and its relations to other techniques. A region–based free energy approximation is formally similar to the CVM, and can be defined through equations (20) and (21), but the requirements on the coefficients a_α are weaker. The single counting condition is imposed only on variable and function nodes, instead of all subclusters:

$$\sum_{\alpha \in R, a \in \alpha} a_\alpha = 1 \quad \forall a, \quad (26)$$

$$\sum_{\alpha \in R, i \in \alpha} a_\alpha = 1 \quad \forall i. \quad (27)$$

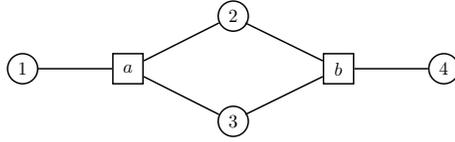


Figure 3. Factor graph of a model for which the Bethe–Peierls approximation is not a special case of the CVM

Interesting particular cases are obtained if R contains only two types of regions, *large regions* and *small regions*. The *junction graph* method [56, 57] is obtained if they form a directed graph, with edges from large to small regions, such that:

- (i) every edge connects a large region with a small region which is a subset of the former;
- (ii) the subgraph of the regions containing a given node is a connected tree.

On the other hand, the *Bethe–Peierls approximation*, in its most general formulation, is obtained by taking function nodes (with the associated variable nodes) as large regions and variable nodes as small regions. This reduces to the usual statistical physics formulation in the case of pairwise interactions.

The CVM is a special region–based free energy approximation, with the property that R is closed under intersection. Indeed, one could define R for the CVM as the set made of the maximal clusters and all the clusters which can be obtained by taking all the possible intersections of (any number of) maximal clusters.

It is easy to verify that the Bethe–Peierls approximation is a special case of CVM only if no function node shares more than one variable node with another function node. If this is not the case, one should be careful when applying the Bethe–Peierls approximation. Consider a model with the factor graph depicted in Figure 3, where $s_i = \pm 1$ ($i = 1, 2, 3, 4$), $H = H_a + H_b$ and

$$H_a(s_1, s_2, s_3) = -h_0 s_1 - \frac{h}{2}(s_2 + s_3) - J s_1 s_2 s_3, \quad (28)$$

$$H_b(s_2, s_3, s_4) = -h_0 s_4 - \frac{h}{2}(s_2 + s_3) - J s_2 s_3 s_4. \quad (29)$$

The CVM, with function nodes as maximal clusters, is exact (notice that it coincides with the junction graph method), and the corresponding exact cumulant expansion for the entropy is

$$S = S_a + S_b - S_{23}, \quad (30)$$

while the Bethe–Peierls entropy is

$$S_{\text{BP}} = S_a + S_b - S_2 - S_3. \quad (31)$$

The two entropies differ by the cumulant $\tilde{S}_{23} = S_{23} - S_2 - S_3$, and hence correlations between variable nodes 2 and 3 cannot be captured by the Bethe–Peierls approximation. In Figure 4 it is clearly illustrated how the Bethe–Peierls

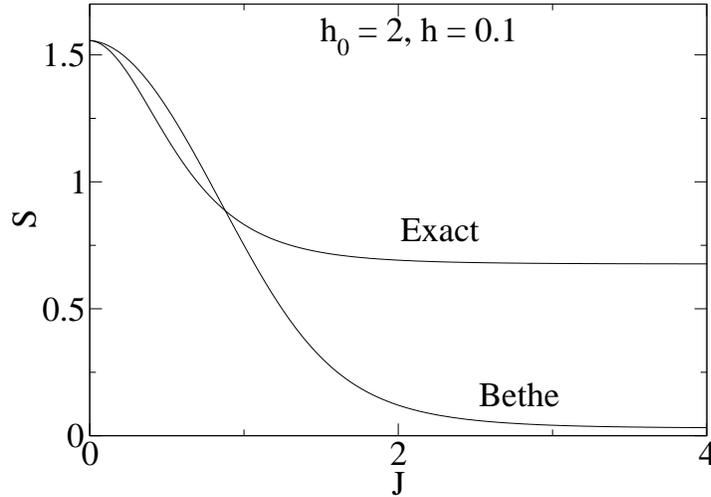


Figure 4. Entropy of the Bethe–Peierls approximation vs the exact one for a model for which the Bethe–Peierls approximation is not a special case of the CVM

approximation can fail. At large enough J the exact entropy is larger (by roughly $\ln 2$) than the Bethe–Peierls one.

5. Exactly solvable cases

The CVM is known to be exact in several cases, due to the topology of the underlying graph, or to the special form of the Hamiltonian. In the present section we shall first consider cases in which the CVM is exact due to the graph topology, then proceed to the issue of realizability and consider cases where the form of the Hamiltonian makes an exact solution feasible with the CVM.

5.1. Tree-like graphs

It is well known that the CVM is exact for models defined on tree-like graphs. This statement can be made more precise by referring to the concept of *junction tree* [58, 59], which we shall actually use in its generalized form given by Yedidia, Freeman and Weiss [56]. A junction tree is a tree-like junction graph. The corresponding large regions are often called *cliques*, and the small regions *separators*. With reference to Figure 3 it is easy to check that the CVM, as described in the previous section, corresponds to a junction tree with cliques $(a123)$ and $(b234)$ and separator (23) , while the junction graph corresponding to the Bethe–Peierls approximation is not a tree.

For a model defined on a junction tree the joint probability distribution factors [56, 60] according to

$$p(\mathbf{s}) = \frac{\prod_{\alpha \in R_L} p_{\alpha}(\mathbf{s}_{\alpha})}{\prod_{\beta \in R_S} p_{\beta}^{d_{\beta}-1}(\mathbf{s}_{\beta})}, \quad (32)$$

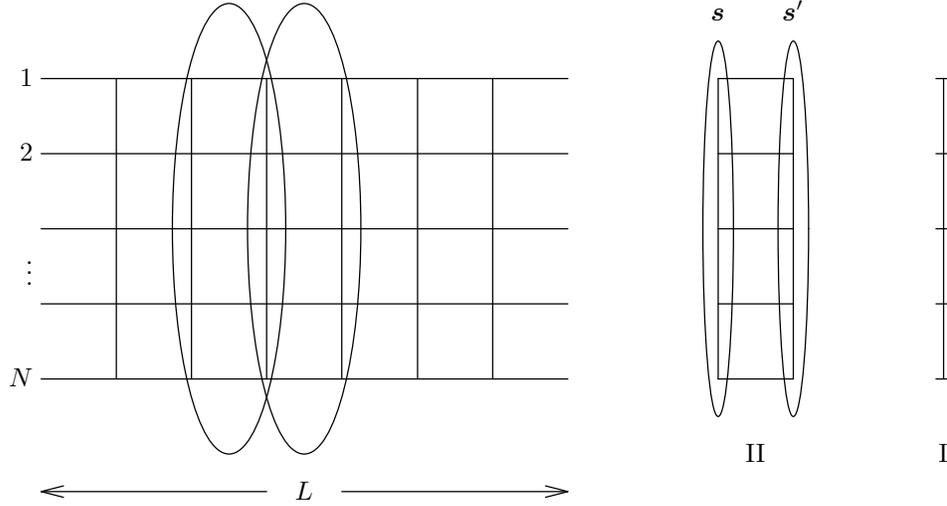


Figure 5. A one-dimensional strip and the clusters used to solve a pairwise model on it

where R_L and R_S denote the sets of large and small regions, respectively, and d_β is the degree of the small region β in the junction tree. Notice that no normalization is needed.

The above factorization of the probability leads to the exact cumulant expansion

$$S = \sum_{\alpha \in R_L} S_\alpha - \sum_{\beta \in R_S} (d_\beta - 1) S_\beta, \quad (33)$$

and therefore the CVM with $R = R_L \cup R_S$ is exact.

As a first example, consider a particular subset of the square lattice, the strip depicted in Figure 5, with open boundary conditions in the horizontal direction, and define on it a model with pairwise interactions (we do not use the factor graph representation here).

According to the junction tree property, the joint probability factors as follows:

$$p(\mathbf{s}) = \frac{\prod_{\alpha \in II} p_\alpha(\mathbf{s}_\alpha)}{\prod_{\beta \in I} p_\beta(\mathbf{s}_\beta)}, \quad (34)$$

where I and II denote the sets of chains (except boundary ones) and ladders, respectively, shown in Figure 5. As a consequence, the cumulant expansion

$$S = \sum_{\alpha \in II} S_\alpha - \sum_{\beta \in I} S_\beta \quad (35)$$

of the entropy is also exact, and the cluster variation method with $R = II \cup I$ is exact. For strip width $N = 1$ we obtain the well-known statement that the Bethe–Peierls approximation is exact for a one-dimensional chain. Rigorous proofs of this statement have been given by Brascamp [61] and Percus [62]. More generally, Schlijper has shown [63] that the equilibrium probability of a d -dimensional statistical mechanical model with finite range interactions is completely determined by its restrictions (marginals) to $d - 1$ -dimensional slices of width at least equal to the interaction range.

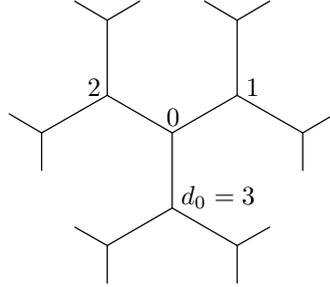


Figure 6. A small portion of a tree

In the infinite length limit $L \rightarrow \infty$ translational invariance is recovered

$$\frac{\mathcal{F}}{L} = \sum_{\mathbf{s}, \mathbf{s}'} ([p_{\text{II}}(\mathbf{s}, \mathbf{s}') H_{\text{II}}(\mathbf{s}, \mathbf{s}') + p_{\text{II}}(\mathbf{s}, \mathbf{s}') \ln p_{\text{II}}(\mathbf{s}, \mathbf{s}')] - \sum_{\mathbf{s}} p_{\text{I}}(\mathbf{s}) \ln p_{\text{I}}(\mathbf{s}) \quad (36)$$

and solving for p_{II} we obtain the transfer matrix formalism:

$$\frac{F}{L} = - \ln \max_{p_{\text{I}}} \left\{ \sum_{\mathbf{s}, \mathbf{s}'} p_{\text{I}}^{1/2}(\mathbf{s}) \exp[-H_{\text{II}}(\mathbf{s}, \mathbf{s}')] p_{\text{I}}^{1/2}(\mathbf{s}') \right\} \quad (37)$$

$$\sum_{\mathbf{s}} p_{\text{I}}(\mathbf{s}) = 1 \quad (38)$$

The natural iteration method (see section 6.3) in this case reduces to the power method for finding the largest eigenvalue of the transfer matrix.

As a second example, consider a tree, like the one depicted in Figure 6, and a model with pairwise interactions defined on it.

In this case the probability factors according to

$$p(\mathbf{s}) = \frac{\prod_{\langle ij \rangle} p_{ij}(s_i, s_j)}{\prod_i p_i^{d_i-1}(s_i)}, \quad (39)$$

where $\langle ij \rangle$ denotes a pair of adjacent nodes. The cumulant expansion of the entropy is therefore

$$S = \sum_{\langle ij \rangle} S_{ij} - \sum_i (d_i - 1) S_i, \quad (40)$$

and the pair approximation of the CVM (coinciding with Bethe–Peierls and junction graph) is exact. Recently this property has been exploited to study models on finite connectivity random graphs, which strictly speaking are not tree-like: loops are present, but in the thermodynamic limit their typical length scales like $\ln N$ [64].

As a final example, consider the so-called (square) cactus lattice (the interior of a Husimi tree), depicted in Figure 7.

Here the probability factors according to

$$p(\mathbf{s}) = \frac{\prod_{\square} p_{\square}(\mathbf{s}_{\square})}{\prod_i p_i(s_i)}, \quad (41)$$

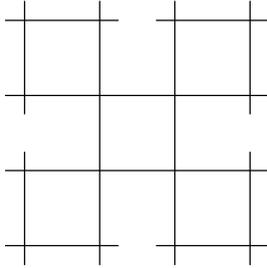


Figure 7. A small portion of a square cactus lattice

the entropy cumulant expansion takes the form

$$S = \sum_{\square} S_{\square} - \sum_i S_i, \quad (42)$$

and the CVM with R made of square plaquettes and sites is exact. Again, this coincides with the junction graph method and, if function nodes are associated to square plaquettes (so that the corresponding factor graph is tree-like), with Bethe–Peierls.

5.2. Realizability

We have seen that when the probability factors in a suitable way, the CVM can be used to find an exact solution. By analogy, we could ask whether, as in the mean field approximation, CVM approximations can yield an estimate of the joint probability distribution as a function of the cluster distributions, in a factorized form. In the general case, the answer is negative. One cannot, using a trial factorized form like

$$\prod_{\alpha} [p_{\alpha}(\mathbf{s}_{\alpha})]^{a_{\alpha}} \quad (43)$$

(which would lead to a free energy like that in Eqs. 20-21), obtain a joint probability distribution which marginalizes down to the cluster probability distributions used as a starting point. As a consequence, we have no guarantee that the CVM free energy is an upper bound to the exact free energy. Moreover, in sufficiently frustrated problems, the cluster probability distributions cannot even be regarded as marginals of a joint probability distribution [65].

It can be easily checked that Equation (43) is not, in the general case, a probability distribution. It is not normalized and therefore its marginals do not coincide with the p_{α} 's used to build it. At best, one can show that

$$\prod_{\alpha} [p_{\alpha}(\mathbf{s}_{\alpha})]^{a_{\alpha}} \propto \exp[-H(\mathbf{s})], \quad (44)$$

but the normalization constant is unknown. This has been proven in [66] at the Bethe–Peierls level, and the proof can be easily generalized to any CVM approximation.

Let us now focus on a very simple example. Consider three Ising variables, $s_i = \pm 1$, $i = 1, 2, 3$, with the following node and pair probabilities:

$$p_i(s_i) = 1/2 \quad i = 1, 2, 3 \quad (45)$$

$$p_{ij}(s_i, s_j) = \frac{1 + cs_i s_j}{4}, \quad -1 \leq c \leq 1, \quad i < j. \quad (46)$$

A joint $p(s_1, s_2, s_3)$ marginalizing to the above probabilities exists for $-1/3 \leq c \leq 1$, which shows clearly that the constraint set Equation (23) is approximate, and in particular it can be too loose. For instance, in [67] it has been shown that due to this problem the Bethe–Peierls approximation for the triangular Ising antiferromagnet predicts, at low temperature, unphysical results for the correlations and a negative entropy.

Moreover, the joint probability $p(s_1, s_2, s_3)$ is given by the CVM–like factorized form

$$\frac{p_{12}(s_1, s_2)p_{13}(s_1, s_3)p_{23}(s_2, s_3)}{p_1(s_1)p_2(s_2)p_3(s_3)} \quad (47)$$

only if $c = 0$, that is if the variables are completely uncorrelated.

As a more interesting case, we shall consider in the next subsection the square lattice Ising model. In this case it has been shown [68, 69] that requiring realizability yields an exactly solvable case.

5.3. Disorder points

For a homogeneous (translation–invariant) model defined on a square lattice, the square approximation of the CVM, that is the approximation obtained by taking the elementary square plaquettes as maximal clusters, entails the following approximate entropy expansion:

$$S \simeq \sum_{\square} S_{\square} - \sum_{\langle ij \rangle} S_{ij} + \sum_i S_i. \quad (48)$$

The corresponding factorization

$$\prod_{\square} p_{\square}(s_{\square}) \prod_{\langle ij \rangle} p_{ij}^{-1}(s_i, s_j) \prod_i p_i(s_i) \quad (49)$$

for the probability does not, in general, give an approximation to the exact equilibrium distribution. Indeed, it does not marginalize to the cluster distributions and is not even normalized.

One could, however, try to impose that the joint probability given by the above factorization marginalizes to the cluster distributions. It turns out that it is sufficient to impose such a condition on the probability distribution of a 3×3 square, like the one depicted in Figure 8. It is easy to check that for an Ising model the CVM–like function

$$\frac{p_4(s_1, s_2, s_5, s_4)p_4(s_2, s_3, s_6, s_5)p_4(s_4, s_5, s_8, s_7)p_4(s_5, s_6, s_9, s_8)p_1(s_5)}{p_2(s_2, s_5)p_2(s_5, s_8)p_2(s_4, s_5)p_2(s_5, s_6)} \quad (50)$$

marginalizes to the square, pair and site distributions (p_4 , p_2 and p_1 respectively) only if odd expectation values vanish and

$$\langle s_i s_k \rangle_{\langle\langle ik \rangle\rangle} = \langle s_i s_j \rangle_{\langle ij \rangle}^2, \quad (51)$$

where the l.h.s. is the next nearest neighbour correlation, while the r.h.s. is the square of the nearest neighbour correlation.

Leaving apart the trivial non–interacting case, the above condition is satisfied by an Ising model with nearest neighbour, next nearest neighbour and plaquette interactions, described by the Hamiltonian

$$H = -J_1 \sum_{\langle ij \rangle} s_i s_j - J_2 \sum_{\langle\langle ij \rangle\rangle} s_i s_j - J_4 \sum_{\square} s_i s_j s_k s_l, \quad (52)$$

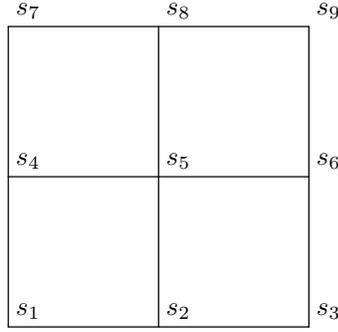


Figure 8. A 3×3 square on the square lattice

if the couplings satisfy the *disorder* condition (see [68] and refs. therein)

$$\cosh(2J_1) = \frac{e^{4J_2+2J_4} + e^{-4J_2+2J_4} + 2e^{-2J_2}}{2(e^{2J_2} + e^{2J_4})}. \quad (53)$$

This defines a variety in the parameter space, lying in the disordered phase of the model, and in particular in the region where nearest neighbour and next nearest neighbour interactions compete. In this case the square approximation of the CVM yields the exact solution, including the exact free energy density

$$f = -\ln[\exp(-J_4) + \exp(J_4 - 2J_2)], \quad (54)$$

and the nearest neighbour correlation

$$g = \langle s_i s_j \rangle_{\langle ij \rangle} = \frac{\exp(-4J_2) - \cosh(2J_1)}{\sinh(2J_1)}. \quad (55)$$

Higher order correlations can be derived from the joint probability Equation (49), for example the two-body correlation function $\Gamma(x, y) = \langle s(x_0, y_0) s(x_0 + x, y_0 + y) \rangle$ (where spin variables have been identified by their coordinates on the lattice), which simply reduces to a power of the nearest neighbour correlation: $\Gamma(x, y) = g^{|x|+|y|}$. For this reason a line of disorder points is often referred to as a one-dimensional line. Or the plaquette correlation:

$$q = \langle s_i s_j s_k s_l \rangle_{\square} = \frac{e^{4J_4} (1 - e^{8J_2}) + 4e^{2J_2} (e^{2J_4} - e^{2J_2})}{e^{4J_4} (1 - e^{8J_2}) + 4e^{2J_2} (e^{2J_4} + e^{2J_2})}. \quad (56)$$

Finally, since all the pair correlations are given simply as powers of the nearest-neighbour correlation we can easily calculate the momentum space correlation function, or structure factor. We first write $\Gamma(x, y) = \exp\left(-\frac{|x|+|y|}{\xi}\right)$, where $\xi = -(\ln g)^{-1}$. After a Fourier transform one finds $S(p_x, p_y) = S_1(p_x) S_1(p_y)$, where

$$S_1(p) = \frac{\sinh(1/\xi)}{\cosh(1/\xi) - \cos p}. \quad (57)$$

It can be verified that the structure factors calculated by Sanchez [70] and (except for a misprint) Cirillo and coworkers [71] reduce to the above expression on the disorder line.

	j									
	1	2	3	4	5	6	7	8	9	10
1	○	○	○	○	○	○	○	○	○	○
2		●	●	●	●	○	○	○	○	○
3			●	●	●	○	○	○	○	○
4				●	●	○	○	○	○	○
5					●	○	○	○	○	○
6						○	○	○	○	○
7							○	○	○	○
8								●	●	●
9									●	●
10										●

Figure 9. A typical configuration of the Muñoz–Eaton model. An empty (resp. filled) circle at row i and column j represents the variable $x_{i,j}$ taking value 0 (resp. 1).

5.4. Wako–Saitô–Muñoz–Eaton model of protein folding

There is at least another case in which the probability factors at the level of square plaquettes, and the CVM yields the exact solution. It is the Wako–Saitô–Muñoz–Eaton model of protein folding [72, 73, 74, 75, 76, 77, 78, 79]. Here we will not delve into the details of the model, giving only its Hamiltonian in the form

$$H = \sum_{i=1}^L \sum_{j=i}^L h_{i,j} x_{i,j}, \quad x_{i,j} = \prod_{k=i}^j x_k, \quad x_k = 0, 1. \quad (58)$$

It is a one–dimensional model with arbitrary range multivariable interactions, but the particular form of these interactions makes an exact solution possible. A crucial step in this solution was the mapping to a two–dimensional model [77], where the statistical variables are the $x_{i,j}$ ’s (see Figure 9 for an illustration). In terms of these variables the Hamiltonian is local, and the only price one has to pay is to take into account the constraints

$$x_{i,j} = x_{i+1,j} x_{i,j-1}, \quad (59)$$

which can be viewed as local interactions.

In order to derive the factorization of the probability [79], we need first to exploit the locality of interactions, which allows us to write

$$p(\{x_{i,j}\}) = \frac{p^{(1,2)} p^{(2,3)} \dots p^{(L-1,L)}}{p^{(2)} \dots p^{(L-1)}}, \quad (60)$$

where $p^{(j)}$ denotes the probability of the j th row in Figure 9 and $p^{(j,j+1)}$ denotes the joint probability of rows j and $j+1$.

As a second step, consider the effect of the constraints. This is best understood looking at the following example:

$$\begin{aligned} p^{(j)}(0, \dots, 0_i, 1_{i+1}, \dots, 1) &= p_{i,i+1}^{(j)}(0, 1) = \\ &= \frac{p_{1,2}^{(j)}(0, 0) \cdots p_{i,i+1}^{(j)}(0, 1) \cdots p_{j-1,j}^{(j)}(1, 1)}{p_2^{(j)}(0) \cdots p_i^{(j)}(0) p_{i+1}^{(j)}(1) \cdots p_{j-1}^{(j)}(1)}. \end{aligned} \quad (61)$$

The CVM-like factorization is possible since every factor in the numerator, except $p_{i,i+1}^{(j)}(0, 1)$, cancels with a factor in the denominator. A similar result can be obtained for the joint probability of two adjacent rows, and substituting into (60) one eventually gets

$$p(\{x_{i,j}\}) = \prod_{\alpha \in R} p_\alpha(x_\alpha)^{a_\alpha}, \quad (62)$$

where $R = \{\text{square plaquettes, corners (on the diagonal), and their subclusters}\}$ and a_α is the CVM Möbius number for cluster α .

6. Cluster Variation Method as an approximation

In most applications the CVM does not yield exact results, and hence it is worth investigating its properties as an approximation.

An important issue is the choice of maximal clusters, and in particular the existence of sequence of approximations (that is, sequence of choices of maximal clusters) with some property of convergence to the exact results. This has been long studied in the literature regarding applications to lattice, translation invariant, systems and will be the subject of subsection 6.1. In particular, rigorous results concerning sequences of approximations which converge to the exact solution have been derived by Schlijper [53, 63, 80], providing a sound theoretical basis for the earlier investigations by Kikuchi and Brush [81].

Another important issue is related to the practical determination of the minima of the CVM variational free energy. In the variational formulation of statistical mechanics the free energy is convex, but this property here is lost due to the presence of negative a_α coefficients in the entropy expansion. More precisely, it has been shown [82] that the CVM variational free energy is convex if

$$\forall S \subseteq R \quad \sum_{\alpha \in R_S} a_\alpha \geq 0 \quad R_S = \{\alpha \in R \mid \exists \beta \subseteq \alpha, \beta \in S\}. \quad (63)$$

Similar conditions have been obtained by McEliece and Yildirim [83] and Heskes, Albers and Kappen [84].

If this is not the case multiple minima can appear, and their determination can be nontrivial. Several algorithms have been developed to deal with this problem, falling mainly in two classes: message-passing algorithms, which will be discussed in subsection 6.2 and variational, provably convergent algorithms, which will be discussed in subsection 6.3.

6.1. Asymptotic behaviour

The first studies on the asymptotic behaviour of sequences of CVM approximations are due to Schlijper [53, 63, 80]. He showed that it is possible to build sequences of CVM approximations (that is, sequences of sets of maximal clusters) such that the corresponding sequence of free energies converge to the exact one, for a translation-invariant model in the thermodynamic limit. The most interesting result, related to the transfer matrix idea, is that for a d -dimensional model the maximal clusters considered have to increase in $d - 1$ dimensions only.

These results provided a theoretical justification for the series of approximations developed by Kikuchi and Brush [81], who introduced the B_{2L} series of approximations for translation-invariant models on the two-dimensional square lattice, based on zig-zag maximal clusters, as shown in Figure 10, and applied it to the zero field Ising model. Based solely on the results from the B_2 (which is equivalent to the plaquette approximation) and B_4 approximations, Kikuchi and Brush postulated a linear behaviour for the estimate of the critical temperature as a function of $(2L+1)^{-1}$.

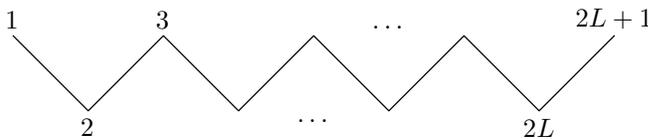


Figure 10. Maximal cluster for the B_{2L} approximation

In Figure 11 we have reported the inverse critical temperature as a function of $(2L+1)^{-1}$ for $L = 1$ to 6. The extrapolated inverse critical temperature is $\beta_c \simeq 0.4397$, to be compared with the exactly known $\beta_c = \frac{1}{2} \ln(1 + \sqrt{2}) \simeq 0.4407$.

It is not our purpose here to make a complete finite size scaling analysis, in the spirit of the coherent anomaly method (see below), of the CVM approximation series. We limit ourselves to show the finite size behaviour of the critical magnetization. More precisely, we have computed the magnetization of the zero field Ising model on the square lattice at the exactly known inverse critical temperature, again for $L = 1$ to 6. Figure 12 shows that the critical magnetization vanishes as $(2L + 1)^{\beta/\nu}$, and the fit gives a very good estimate for the exponent, consistent with the exact result $\beta/\nu = 1/8$.

As a further illustration of the asymptotic properties of the B_{2L} series we report in Figure 13 the zero temperature entropy (actually the difference between the extrapolated entropy density and the B_{2L} estimate) of the Ising triangular antiferromagnet as a function of $1/L$ [67]. It is clearly seen that asymptotically $s_L = s_0 - aL^{-\psi}$, and the fit yields the numerical results $s_0 \approx 0.323126$ (the exact value being $s \approx 0.323066$) and $\psi \approx 1.7512$ (remarkably close to $7/4$).

An attempt to extract non-classical critical behaviour from high precision low and high temperature results from CVM was made by the present author [19, 20, 21, 22], using Padé and Adler approximants. This work has led to the development of an 18 ($3 \times 3 \times 2$) site cluster approximation for the simple cubic lattice Ising model [22], which is probably the largest cluster ever considered. The results obtained for the Ising model are still compatible with the most recent estimates [85], although of lower precision.

It has also been considered the possibility of extracting non-classical critical behaviour from CVM results by means of the coherent anomaly method, which applies

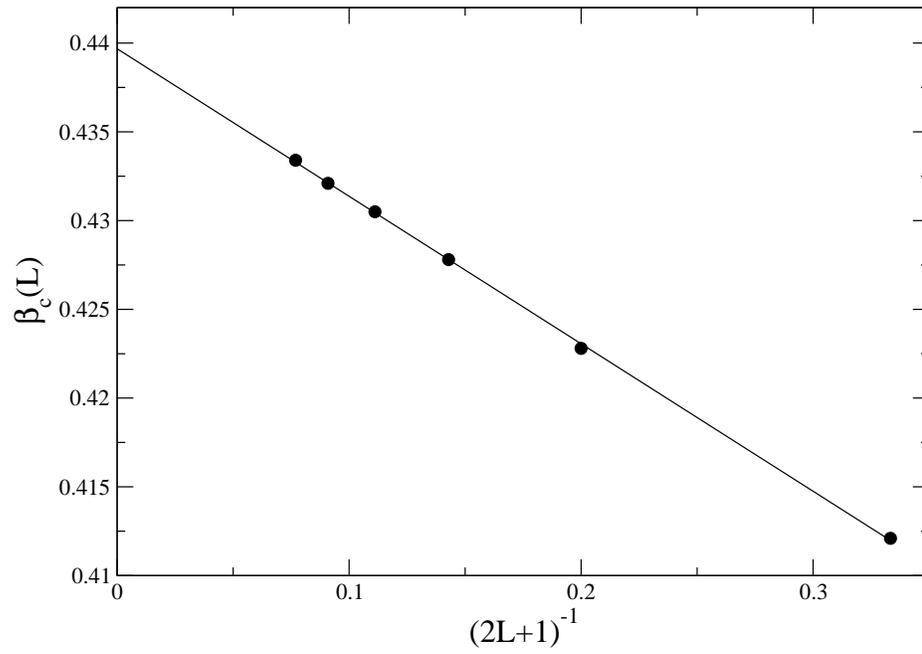


Figure 11. Inverse critical temperature of the B_{2L} approximation series

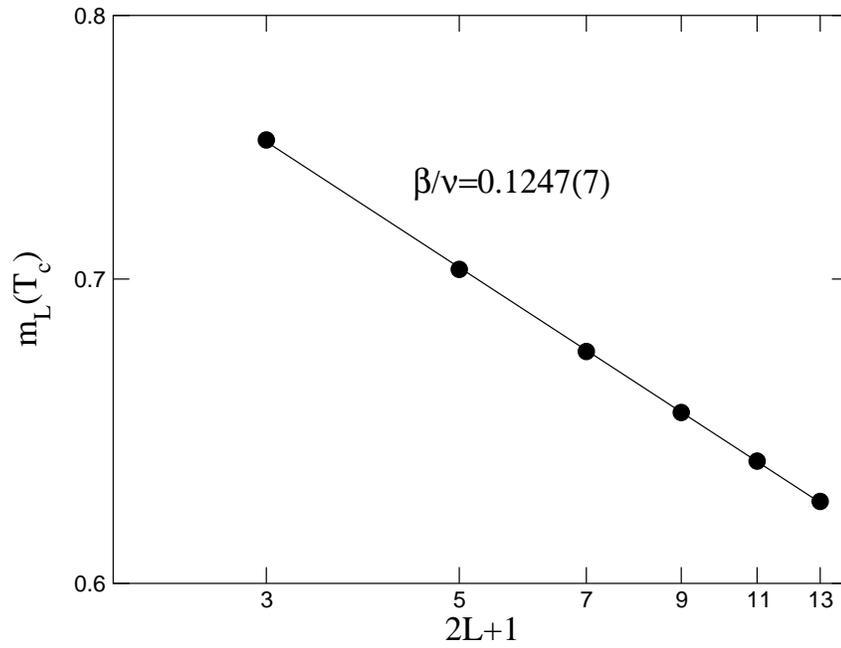


Figure 12. Critical temperature of the B_{2L} approximation series

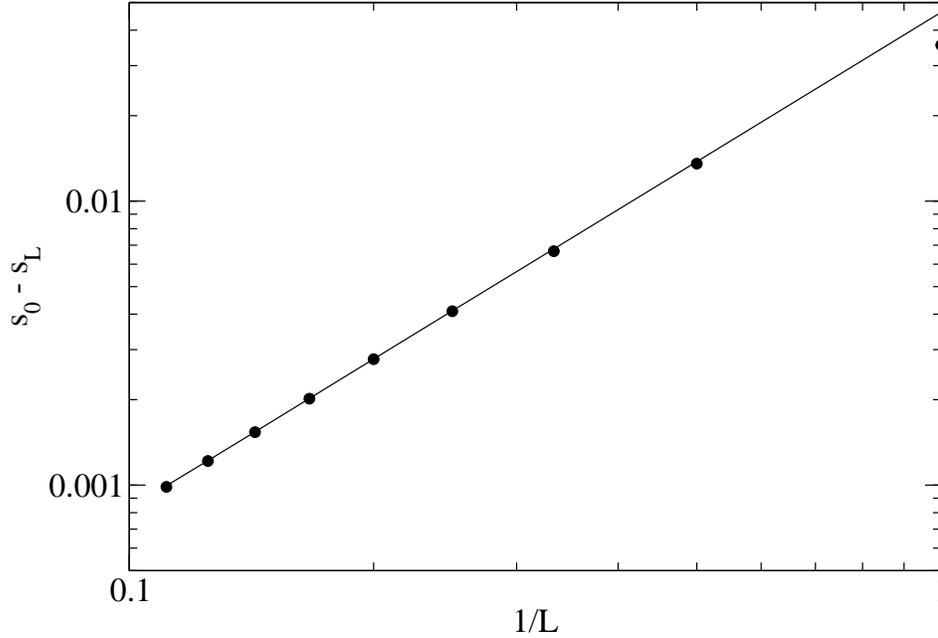


Figure 13. Zero temperature entropy of the triangular Ising antiferromagnet in the B_{2L} approximation series

finite size scaling ideas to series of mean-field-like approximations. A review of these results can be found in [23].

6.2. Message-passing algorithms

In order to describe this class of algorithms it is useful to start with the Bethe-Peierls approximation (pair approximation of the CVM) free energy for the Ising model Equation (5):

$$\begin{aligned}
\mathcal{F} = & - \sum_i h_i \sum_{s_i} s_i p_i(s_i) - \sum_{\langle ij \rangle} J_{ij} \sum_{s_i, s_j} s_i s_j p_{ij}(s_i, s_j) + \\
& + \sum_{\langle ij \rangle} \sum_{s_i, s_j} p_{ij}(s_i, s_j) \ln p_{ij}(s_i, s_j) - \sum_i (d_i - 1) \sum_{s_i} p_i(s_i) \ln p_i(s_i) \\
& + \sum_i \lambda_i \left(\sum_{s_i} p_i(s_i) - 1 \right) + \sum_{\langle ij \rangle} \lambda_{ij} \left(\sum_{s_i, s_j} p_{ij}(s_i, s_j) - 1 \right) + \\
& + \sum_{\langle ij \rangle} \left[\nu_{i,j} \left(\sum_{s_i} s_i p_i(s_i) - \sum_{s_i, s_j} s_i p_{ij}(s_i, s_j) \right) + \right. \\
& \left. + \nu_{j,i} \left(\sum_{s_j} s_j p_j(s_j) - \sum_{s_i, s_j} s_j p_{ij}(s_i, s_j) \right) \right]. \tag{64}
\end{aligned}$$

One can easily recognize the energy terms, the pair entropy, the site entropy (with a Möbius number $-(d_i - 1)$, where d_i is the degree of node i), and the Lagrange terms corresponding to the normalization and pair-site compatibility constraints. Observe that, due the presence of normalization constraints, it is enough to impose the consistency between the spin expectation values given by the site and pair probabilities.

A physical way of deriving message-passing algorithms for the determination of the stationary points of the above free energy is through the introduction of the effective field representation. Consider the interaction J_{ik} and assume that, whenever this is not taken into account exactly, its effect on variable s_i can be replaced by an effective field $h_{i,k}$. This can be made rigorous by observing that the stationarity conditions

$$\begin{aligned} \frac{\partial \mathcal{F}}{\partial p_i(s_i)} &= 0 \\ \frac{\partial \mathcal{F}}{\partial p_{ij}(s_i, s_j)} &= 0 \end{aligned} \quad (65)$$

can be solved by writing the probabilities as

$$p_i(s_i) = \exp \left[F_i + \left(h_i + \sum_{k \text{ NN } i} h_{i,k} \right) s_i \right] \quad (66)$$

$$p_{ij}(s_i, s_j) = \exp \left[F_{ij} + \left(h_i + \sum_{k \text{ NN } i}^{k \neq j} h_{i,k} \right) s_i + \left(h_j + \sum_{k \text{ NN } j}^{k \neq i} h_{j,k} \right) s_j + J_{ij} s_i s_j \right], \quad (67)$$

where the effective fields, and the site and pair free energies F_i and F_{ij} , are related to the Lagrange multipliers through

$$\begin{aligned} \lambda_i &= (d_i - 1)(1 + F_i) \\ \lambda_{ij} &= -1 - F_{ij} \\ \nu_{i,j} &= h_i + \sum_{k \text{ NN } i}^{k \neq j} h_{i,k}. \end{aligned} \quad (68)$$

F_i and F_{ij} are determined by the normalization, but first of all the effective fields must be computed by imposing the corresponding compatibility constraints, which can be cast into the form

$$h_{i,j} = \tanh^{-1} \left[\tanh \left(h_j + \sum_{k \text{ NN } j}^{k \neq i} h_{j,k} \right) \tanh J_{ij} \right]. \quad (69)$$

This is a set of coupled nonlinear equations which is often solved by iteration, that is an initial guess is made for the $h_{i,j}$'s, plugged into the r.h.s. of Equation (69) which returns a new estimate, and the procedure is then repeated until a fixed point is (hopefully) reached. The values of the effective fields at the fixed point can then be used to compute the probabilities according to Equation (67).

The above equations, and their generalizations at the CVM level, have been intensively used in the 80's for studying the average behaviour of models with quenched random interactions, like Ising spin glass models. This work was started by a paper by Morita [86], where an integral equation for the probability distribution of the effective

field, given the probability distributions of the interactions and fields, was derived. In the general case this integral equation takes the form

$$p_{i,j}(h_{i,j}) = \int \delta \left(h_{i,j} - \tanh^{-1} \left[\tanh \left(h_j + \sum_{\substack{k \neq i \\ k \text{ NN } j}} h_{j,k} \right) \tanh J_{ij} \right] \right) \times \\ \times P_j(h_j) dh_j P_{ij}(J_{ij}) dJ_{ij} \prod_{\substack{k \neq i \\ k \text{ NN } j}} p_{j,k}(h_{j,k}) dh_{j,k}, \quad (70)$$

with simplifications occurring if the probability distributions can be assumed to be site-independent, which is the most studied case. Typical calculations concerned: the determination of elements of the phase diagrams of Ising spin glass models, through the calculation of the instability loci of the paramagnetic solution; results in the zero temperature limit, where solutions with a discrete support are found; iterative numerical solutions of the integral equation. A review of this line of research until 1986 can be found in [87]. It is important to notice that the results obtained by this approach are equivalent to those by the replica method, at the replica symmetric level.

The effective field approach is reminiscent of the message-passing procedure at the heart of the belief propagation (BP) algorithm, and indeed the messages appearing in this algorithm are related, in the Ising case, to the effective fields by $m_{\langle ij \rangle \rightarrow i}(s_i) = \exp(h_{i,j} s_i)$, where $m_{\langle ij \rangle \rightarrow i}(s_i)$ denotes a message going from the NN pair $\langle ij \rangle$ to node i .

In order to derive the BP algorithm consider the Bethe–Peierls approximation for a model with variable nodes i and factor nodes a . The variables s_i need not be limited to two states and the Hamiltonian is written in the general form Equation (3).

The CVM free energy, with the normalization and compatibility constraints, can then be written as

$$\mathcal{F} = - \sum_a \sum_{\mathbf{s}_a} H_a(\mathbf{s}_a) p_a(\mathbf{s}_a) + \\ + \sum_a \sum_{\mathbf{s}_a} p_a(\mathbf{s}_a) \ln p_a(\mathbf{s}_a) - \sum_i (d_i - 1) \sum_{s_i} p_i(s_i) \ln p_i(s_i) + \\ + \sum_i \lambda_i \left(\sum_{s_i} p_i(s_i) - 1 \right) + \sum_a \lambda_a \left(\sum_a \sum_{\mathbf{s}_a} p_a(\mathbf{s}_a) - 1 \right) + \\ + \sum_a \sum_{i \in a} \sum_{s_i} \mu_{a,i}(s_i) \left(p_i(s_i) - \sum_{\mathbf{s}_a \setminus i} p_a(\mathbf{s}_a) \right), \quad (71)$$

where $\mathbf{s}_a \setminus i$ denotes the set of variables entering factor node a , except i .

The stationarity conditions

$$\frac{\partial \mathcal{F}}{\partial p_i(s_i)} = 0 \\ \frac{\partial \mathcal{F}}{\partial p_a(\mathbf{s}_a)} = 0 \quad (72)$$

can then be solved, in analogy with Equation (67), by

$$p_i(s_i) = \frac{1}{Z_i} \prod_{i \in a} m_{a \rightarrow i}(s_i)$$

$$p_a(\mathbf{s}_a) = \frac{1}{Z_a} \psi_a(\mathbf{s}_a) \prod_{k \in a} \prod_{k \in b}^{b \neq a} m_{b \rightarrow k}(s_k). \quad (73)$$

In the particular case of an Ising model with pairwise interactions, the previously mentioned relationship between messages and effective fields is evident from the above equation.

Now Z_i and Z_a take care of normalization, and the messages $m_{a \rightarrow i}(s_i)$ are related to the Lagrange multipliers by

$$\mu_{a,k}(s_k) = \sum_{k \in b}^{b \neq a} \ln m_{b \rightarrow k}(s_k). \quad (74)$$

Notice that the messages can be regarded as exponentials of a new set of Lagrange multipliers, with the constraints rewritten as in the following free energy

$$\begin{aligned} \mathcal{F} = & - \sum_a \sum_{\mathbf{s}_a} H_a(\mathbf{s}_a) p_a(\mathbf{s}_a) + \\ & + \sum_a \sum_{\mathbf{s}_a} p_a(\mathbf{s}_a) \ln p_a(\mathbf{s}_a) - \sum_i (d_i - 1) \sum_{s_i} p_i(s_i) \ln p_i(s_i) + \\ & + \sum_i \lambda_i \left(\sum_{s_i} p_i(s_i) - 1 \right) + \sum_a \lambda_a \left(\sum_a \sum_{\mathbf{s}_a} p_a(\mathbf{s}_a) - 1 \right) + \\ & + \sum_a \sum_{i \in a} \sum_{s_i} \ln m_{a \rightarrow i}(s_i) \left((d_i - 1) p_i(s_i) - \sum_{i \in b} \sum_{\mathbf{s}_b \setminus i}^{b \neq a} p_b(\mathbf{s}_b) \right). \quad (75) \end{aligned}$$

Again, imposing compatibility between variable nodes and factor nodes, one gets a set of coupled equations for the messages which, leaving apart normalization, read

$$m_{a \rightarrow i}(s_i) \propto \sum_{\mathbf{s}_a \setminus i} \psi_a(\mathbf{s}_a) \prod_{k \in a} \prod_{k \in b}^{k \neq i, b \neq a} m_{b \rightarrow k}(s_k). \quad (76)$$

The above equations, and their iterative solution, are the core of the BP algorithm. Also, their structure justifies the name ‘‘Sum-Product’’ [48], which is often given them in the literature on probabilistic graphical models, and the corresponding term ‘‘Max-Product’’ for their zero temperature limit.

There are several issues which must be considered when discussing the property of an iterative algorithm based on Equation (76). First of all, one could ask whether messages have to be updated sequentially or in parallel. This degree of freedom does not affect the fixed points of the algorithm, but it affects the dynamics. This issue has been considered in some depth by Kfir and Kanter [88] in the context of the decoding of error-correcting codes. In that case they showed that the sequential update results in twice faster convergence with respect to the parallel update.

Convergence is however not guaranteed if the underlying graph is not tree-like, that is if the pair approximation of the CVM is not exact. This issue has been investigated theoretically by Tatikonda and Jordan [89], Mooij and Kappen [90], Ihler et al [91], who derived sufficient conditions for convergence, and by Heskes [92], who derived sufficient conditions for the uniqueness of the fixed point. In practice it is typically observed that the BP algorithm converges if the frustration due to competitive interactions, like those characteristic of spin-glass or constraint satisfaction models, is not too large. In some cases, the trick of damping, or inertia,

can help extending the convergence domain. The trick consists in taking the updated message equal to a weighted (possibly geometrical) average of the old message and the new one given by Equation (76). The convergence domain of the BP algorithm has been determined for several problems, like satisfiability [26], graph colouring [27], error correcting codes [40] and spin glasses [93]. Within its convergence domain, the BP algorithm is indeed very fast, and this is its real strength. See the next subsection for some performance tests and a comparison with provably convergent algorithms.

Once a fixed point has been obtained it is worth asking whether this corresponds to a minimum of the free energy or not. This has been partially solved by Heskes [94], who has shown that stable fixed points of the belief propagation algorithm are minima of the CVM pair approximation free energy, but the converse is not necessarily true. Actually, examples can be found of minima of the free energy which correspond to unstable fixed points of the belief propagation algorithm.

An important advancement in this topic is the *generalized belief propagation* (GBP) algorithm by Yedidia and coworkers [31]. The fixed points of the GBP algorithm for a certain choice of clusters correspond to stationary points of the CVM free energy at the approximation level corresponding by the same choice of clusters or, more generally, of a region graph free energy. Actually, for a given choice of clusters, different GBP algorithms can be devised. Here only the so-called *parent to child* GBP algorithm [56] will be considered. Other choices are described in [56].

In order to better understand this algorithm, notice a few characteristics of the belief propagation algorithm. First of all, looking at the probabilities Equation (73) one can say that a variable node receives messages from all the factor nodes it belongs to, while a factor node a receives messages from all the other factor nodes to which its variable nodes $i \in a$ belong. In addition, the constraint corresponding to the message $m_{a \rightarrow i}(s_i)$ (see Equation (75)) can be written as

$$\sum_{\mathbf{s}_{\mathbf{a} \setminus i}} p_{\mathbf{a}}(\mathbf{s}_{\mathbf{a}}) = \sum_{i \in b} \sum_{\mathbf{s}_{b \setminus i}} p_b(\mathbf{s}_b) - (d_i - 1)p_i(s_i). \quad (77)$$

The parent to child GBP algorithm generalizes these characteristics in a rather straightforward way. First of all, messages $m_{\alpha \rightarrow \beta}(\mathbf{s}_{\beta})$ ($\beta \subset \alpha$) are introduced from regions (parent regions) to subregions (child regions). Then, the probability of a region takes into account messages coming from outer regions to itself and its subregions. Finally, exploiting the property Equation (19) of the Möbius numbers, the constraint corresponding to $m_{\alpha \rightarrow \beta}(\mathbf{s}_{\beta})$ is written in the form

$$\sum_{\alpha \subseteq \gamma \in R} a_{\gamma} \sum_{\mathbf{s}_{\gamma \setminus \beta}} p_{\gamma}(\mathbf{s}_{\gamma}) = \sum_{\beta \subseteq \gamma \in R} a_{\gamma} \sum_{\mathbf{s}_{\gamma \setminus \beta}} p_{\gamma}(\mathbf{s}_{\gamma}). \quad (78)$$

It can be shown [56] that this new set of constraints is equivalent to the original one.

To make this more rigorous, consider the free energy given by Equations (20) and (21), with the above compatibility constraints (with Lagrange multipliers $\ln m_{\alpha \rightarrow \beta}(\mathbf{s}_{\beta})$) and the usual normalization constraints (with multipliers λ_{α}). One obtains

$$\begin{aligned} \mathcal{F} = & \sum_{\gamma \in R} a_{\gamma} \sum_{\mathbf{s}_{\gamma}} [p_{\gamma}(\mathbf{s}_{\gamma}) H_{\gamma}(\mathbf{s}_{\gamma}) + p_{\gamma}(\mathbf{s}_{\gamma}) \ln p_{\gamma}(\mathbf{s}_{\gamma})] + \sum_{\gamma \in R} \lambda_{\gamma} \left[\sum_{\mathbf{s}_{\gamma}} p_{\gamma}(\mathbf{s}_{\gamma}) - 1 \right] + \\ & + \sum_{\beta \subset \alpha \in R} \sum_{\mathbf{s}_{\beta}} \ln m_{\alpha \rightarrow \beta}(\mathbf{s}_{\beta}) \left[\sum_{\alpha \subseteq \gamma \in R} a_{\gamma} \sum_{\mathbf{s}_{\gamma \setminus \beta}} p_{\gamma}(\mathbf{s}_{\gamma}) - \sum_{\beta \subseteq \gamma \in R} a_{\gamma} \sum_{\mathbf{s}_{\gamma \setminus \beta}} p_{\gamma}(\mathbf{s}_{\gamma}) \right], \quad (79) \end{aligned}$$

where it is not necessary to put all the possible $\alpha \rightarrow \beta$ compatibility constraints, but it is enough to put those which satisfy $a_\alpha \neq 0$, $a_\beta \neq 0$ and β is a direct subregion of α , that is there is no region γ with $a_\gamma \neq 0$ such that $\beta \subset \gamma \subset \alpha$. Notice also that the Lagrange term corresponding to the $\alpha \rightarrow \beta$ constraint can be written as

$$-\ln m_{\alpha \rightarrow \beta}(\mathbf{s}_\beta) \sum_{\beta \subseteq \gamma \in R}^{\alpha \not\subseteq \gamma} a_\gamma \sum_{\mathbf{s}_{\gamma \setminus \beta}} p_\gamma(\mathbf{s}_\gamma). \quad (80)$$

The stationarity conditions

$$\frac{\partial \mathcal{F}}{\partial p_\gamma(\mathbf{s}_\gamma)} = 0 \quad (81)$$

can then be solved, leaving apart normalization, by

$$p_\gamma(\mathbf{s}_\gamma) \propto \exp[-H_\gamma(\mathbf{s}_\gamma)] \prod_{\beta \subseteq \gamma}^{\alpha \not\subseteq \gamma} \prod_{\beta \subset \alpha \in R} m_{\alpha \rightarrow \beta}(\mathbf{s}_\beta), \quad (82)$$

where \mathbf{s}_β denotes the restriction of \mathbf{s}_γ to subregion β .

Finally, message update rules can be derived again by the compatibility constraints, though some care is needed, since in the general case these constraints are not immediately solved with respect to the (updated) messages, as it occurs in the derivation of Equation (76). Here one obtains a coupled set of equations in the updated messages, which can be solved starting from the constraints involving the smallest clusters.

An example can be helpful here. Consider a model defined on a regular square lattice, with periodic boundary conditions, and the CVM square approximation, that is the approximation obtained by taking the elementary square plaquettes as maximal clusters. The entropy expansion contains only terms for square plaquettes (with Möbius numbers 1), NN pairs (Möbius numbers -1) and single nodes (Möbius numbers 1), as in Equation (48). A minimal set of compatibility constraints includes node-pair and pair-square constraints, and one has therefore to deal with square-to-pair and pair-to-node messages, which will be denoted by $m_{ij,kl}(s_i, s_j)$ and $m_{i,j}(s_i)$ respectively. With reference to the portion of the lattice depicted in Figure 14 the probabilities, according to Equation (82), can be written as

$$\begin{aligned} p_i(s_i) &\propto \exp[-H_i(s_i)] m_{i,a}(s_i) m_{i,j}(s_i) m_{i,l}(s_i) m_{i,h}(s_i), \\ p_{ij}(s_i, s_j) &\propto \exp[-H_{ij}(s_i, s_j)] m_{i,a}(s_i) m_{i,l}(s_i) m_{i,h}(s_i) \times \\ &\quad \times m_{j,b}(s_j) m_{j,c}(s_j) m_{j,k}(s_j) m_{ij,ab}(s_i, s_j) m_{ij,lk}(s_i, s_j), \\ p_{ijkl}(s_i, s_j, s_k, s_l) &\propto \exp[-H_{ijkl}(s_i, s_j, s_k, s_l)] m_{i,a}(s_i) m_{i,h}(s_i) m_{j,b}(s_j) m_{j,c}(s_j) \times \\ &\quad \times m_{k,d}(s_k) m_{k,e}(s_k) m_{l,f}(s_l) m_{l,g}(s_l) \times \\ &\quad \times m_{ij,ab}(s_i, s_j) m_{jk,cd}(s_j, s_k) m_{kl,ef}(s_k, s_l) m_{lj,gh}(s_l, s_j). \end{aligned} \quad (83)$$

Imposing node-pair and pair-square constraints one gets equations like

$$\begin{aligned} \exp[-H_i(s_i)] m_{i,j}(s_i) &\propto \sum_{s_j} \exp[-H_{ij}(s_i, s_j)] \times \\ &\quad \times m_{j,b}(s_j) m_{j,c}(s_j) m_{j,k}(s_j) m_{ij,ab}(s_i, s_j) m_{ij,lk}(s_i, s_j), \\ \exp[-H_{ij}(s_i, s_j)] m_{i,f}(s_i) m_{j,k}(s_j) m_{ij,lk}(s_i, s_j) &\propto \sum_{s_k, s_l} \exp[-H_{ijkl}(s_i, s_j, s_k, s_l)] \times \\ &\quad \times m_{k,d}(s_k) m_{k,e}(s_k) m_{l,f}(s_l) m_{l,g}(s_l) \times \\ &\quad \times m_{jk,cd}(s_j, s_k) m_{kl,ef}(s_k, s_l) m_{lj,gh}(s_l, s_j). \end{aligned} \quad (84)$$

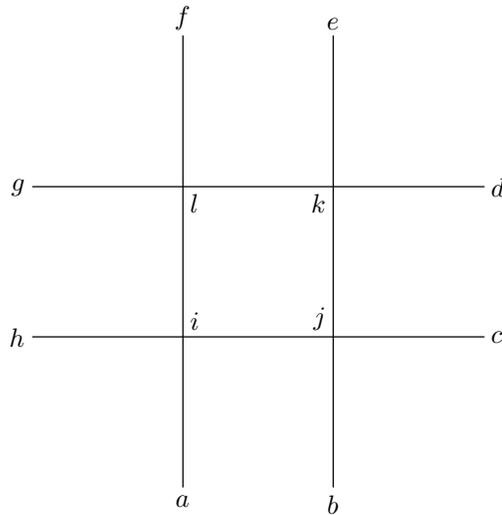


Figure 14. A portion of the square lattice

The above equations can be viewed as a set of equations in the updated messages at iteration $t + 1$, appearing in the l.h.s., given the messages at iteration t , appearing in the r.h.s.. It is clear that one has first to calculate the updated pair-to-site messages according to the first equation, and then the updated square-to-pair messages according to the second one, using in the l.h.s. the updated pair-to-site messages just obtained.

GBP (possibly with damping) typically exhibits better convergence properties (and greater accuracy) than BP, but the empirical rule that a sufficient amount of frustration can make it not convergent is valid also for GBP. It is therefore fundamental to look for provably convergent algorithms, which will be discussed in the next subsection. A variation of the BP algorithm, the conditioned probability (CP) algorithm, with improved convergence properties, has recently been introduced [95]. The extension of this algorithm beyond the BP level is however not straightforward.

We conclude the present subsection by mentioning that techniques like the Thouless-Anderson-Palmer equations, or the cavity method, both widely used in the statistical physics of spin glasses, are strictly related to the Bethe-Peierls approximation.

The Thouless-Anderson-Palmer [96] equations can be derived from the Bethe-Peierls free energy for the Ising model, through the so-called Plefka expansion [97]. One has first to write the free energy as a function of magnetizations and nearest-neighbour correlations through the parameterization

$$p_i(s_i) = \frac{1 + s_i m_i}{2} \quad p_{ij}(s_i, s_j) = \frac{1 + s_i m_i + s_j m_j + s_i s_j c_{ij}}{4}, \quad (85)$$

then to solve analytically the stationarity conditions with respect to the c_{ij} 's and finally to expand to second order in the inverse temperature.

Finally, the cavity method [98, 99, 100] is particularly important since it allows to deal with replica symmetry breaking. The cavity method, though historically derived in a different way, can be regarded as an alternative choice of messages and effective fields in the Bethe-Peierls approximation. With reference to Equation (73), introduce

messages $m_{k \rightarrow a}(s_k)$ from variable nodes to factor nodes according to

$$m_{k \rightarrow a}(s_k) = \prod_{\substack{b \neq a \\ k \in b}} m_{b \rightarrow k}(s_k). \quad (86)$$

Then the probabilities Equation (73) become

$$\begin{aligned} p_i(s_i) &= \frac{1}{Z_i} \prod_{i \in a} m_{a \rightarrow i}(s_i) \\ p_a(\mathbf{s}_a) &= \frac{1}{Z_a} \psi_a(\mathbf{s}_a) \prod_{k \in a} m_{k \rightarrow a}(s_k), \end{aligned} \quad (87)$$

and the message update equations (76) become

$$m_{a \rightarrow i}(s_i) \propto \sum_{\mathbf{s}_a \setminus i} \psi_a(\mathbf{s}_a) \prod_{\substack{k \neq i \\ k \in a}} m_{k \rightarrow a}(s_k). \quad (88)$$

The effective fields corresponding to the factor-to-variable messages $m_{a \rightarrow i}(s_i)$ are usually called cavity biases, while those corresponding to the variable-to-factor messages $m_{i \rightarrow a}(s_i)$ are called cavity fields. In the Ising example above a factor node is just a pair of NNs and cavity biases reduce to effective fields $h_{i,j}$, while cavity fields

take the form $\sum_{\substack{k \neq j \\ k \text{ NN } i}} h_{i,k}$.

The cavity method admits an extension to cases where one step of replica symmetry breaking occurs [100, 101]. In such a case one assumes that there exist many states characterized by different values of the cavity biases and fields, and introduces the probability distributions of cavity biases and fields over the states. From the above message update rules one can then derive integral equations, similar to Equation (70), for the distributions. These integral equations can in principle be solved by iterative population dynamics algorithms, but most often one restricts to the zero temperature case, where these distributions have a discrete support.

The zero temperature case is particularly relevant for hard combinatorial optimization problems, where 1-step replica symmetry breaking corresponds to clustering of solutions. Clustering means that the space of solutions becomes disconnected, made of subspaces which cannot be reached from one another by means of local moves, and hence all local algorithms, like BP or GBP, are bound to fail. The cavity method has been used to solve these kind of problems in the framework of the survey propagation algorithm [25], which has been shown to be a very powerful tool for constraint satisfaction problems like satisfiability [26] and colouring [27] defined on finite connectivity random graphs. These graphs are locally tree-like and therefore all the analysis can be carried out at the Bethe-Peierls level. A sort of generalized survey propagation capable of dealing with short loops would really be welcome, but it seems that realizability issues are crucial here and replica symmetry breaking can only be introduced when CVM gives an exact solution.

A different approach, still aimed to generalize the BP algorithm to situations where replica symmetry breaking occurs, has been suggested by van Mourik [28], and is based on the analysis of the time evolution of the BP algorithm.

6.3. Variational algorithms

In the present subsection we discuss algorithms which update probabilities instead of messages. At every iteration a new estimate of probabilities, and hence of the free energy, is obtained. These algorithms are typically provably convergent, and the proof is based on showing that the free energy decreases at each iteration. This is of course not possible with BP and GBP algorithms, where the probabilities and the free energy can be evaluated only at the fixed point. The price one has to pay is that in variational algorithms one has to solve the compatibility constraints at every iteration, and therefore these are double loop algorithms, where the outer loop is used to update probabilities and the inner loop is used to solve the constraints.

The natural iteration method (NIM) [102, 103] is the oldest algorithm specifically designed to minimize the CVM variational free energy. It was originally introduced [102] in the context of homogeneous models, for the pair and tetrahedron (for the fcc lattice) approximations. In such cases the compatibility constraints are trivial. Later [103] it was generalized to cases where the compatibility constraints cannot be solved trivially. An improved version of the algorithm, with tunable convergence properties, appeared in [104] and its application is described in some detail also in [105], where higher order approximations are considered.

The algorithm is based on a double loop scheme, where the inner loop is used to solve the compatibility constraints, so that at each iteration of the outer loop a set of cluster probabilities which satisfy the constraints is obtained.

Proof of convergence, based on showing that the free energy decreases at every outer loop iteration, exist in many cases, but it has also been shown that there are non-convergent cases, like the four-dimensional Ising model [106] in the hypercube approximation.

We do not discuss in detail this algorithm since it is rather slow, and better alternatives have been recently developed.

A first step in this direction was the *concave-convex procedure* (CCCP) by Yuille [107], who started from the observation that the non-convergence problems of message-passing algorithms arise from concave terms in the variational free energy, that is from the entropy of clusters with negative Möbius numbers. His idea was then to split the CVM free energy into a convex and a concave part,

$$\mathcal{F}(\{p_\alpha\}) = \mathcal{F}_{\text{vex}}(\{p_\alpha\}) + \mathcal{F}_{\text{cave}}(\{p_\alpha\}), \quad (89)$$

and to write the update equations to be iterated to a fixed point as

$$\nabla \mathcal{F}_{\text{vex}}(\{p_\alpha^{(t+1)}\}) = -\nabla \mathcal{F}_{\text{cave}}(\{p_\alpha^{(t)}\}), \quad (90)$$

where $p_\alpha^{(t)}$ and $p_\alpha^{(t+1)}$ are successive iterates. In order to solve the compatibility constraints, at each iteration of Equation (90), the Lagrange multipliers enforcing the constraints are determined by another iterative algorithm where one solves for one multiplier at a time, and it can be shown that the free energy decreases at each outer loop iteration. Therefore we have another double loop algorithm, which is provably convergent, faster than NIM (as we shall see below), and allows some freedom in the splitting between convex and concave parts.

A more general and elegant formalism, which will be described in the following, has however been put forward by Heskes, Albers and Kappen (HAK) [84]. Their basic idea is to consider a sequence of convex variational free energies such that the sequence of the corresponding minima tends to the minimum of the CVM free energy. More

precisely, if the CVM free energy $\mathcal{F}(\{p_\alpha, \alpha \in R\})$ is denoted for simplicity by $\mathcal{F}(p)$, they consider a function $\mathcal{F}_{\text{conv}}(p, p')$, convex in p , with the properties

$$\begin{aligned}\mathcal{F}_{\text{conv}}(p, p') &\geq \mathcal{F}(p), \\ \mathcal{F}_{\text{conv}}(p, p) &= \mathcal{F}(p).\end{aligned}\quad (91)$$

The algorithm is then defined by the update rule for the probabilities

$$p^{(t+1)} = \arg \min_p \mathcal{F}_{\text{conv}}(p, p^{(t)}), \quad (92)$$

and it is easily proved that the free energy decreases at each iteration and that a minimum of the CVM free energy is recovered at the fixed point.

A lot of freedom is left in the definition of $\mathcal{F}_{\text{conv}}$, and strategies of varying complexity and speed can be obtained. NIM (when convergent) and CCCP can also be recovered as special cases. The general framework is based on the following three properties.

(i) If $\beta \subset \alpha$, then

$$-S_\alpha + S_\beta = \sum_{\mathbf{s}_\alpha} p_\alpha(\mathbf{s}_\alpha) \ln p_\alpha(\mathbf{s}_\alpha) - \sum_{\mathbf{s}_\beta} p_\beta(\mathbf{s}_\beta) \ln p_\beta(\mathbf{s}_\beta) \quad (93)$$

is convex over the constraint set, i.e. it is a convex function of p_α and p_β if these satisfy the compatibility constraint Equation (23).

(ii) The linear bound

$$S_\beta = - \sum_{\mathbf{s}_\beta} p_\beta(\mathbf{s}_\beta) \ln p_\beta(\mathbf{s}_\beta) \leq - \sum_{\mathbf{s}_\beta} p_\beta(\mathbf{s}_\beta) \ln p'_\beta(\mathbf{s}_\beta) = S'_\beta \quad (94)$$

holds, with equality only for $p'_\beta = p_\beta$

(iii) If $\gamma \subset \beta$, and p_β and p_γ (p'_β and p'_γ) satisfy the compatibility constraints, the bound

$$\begin{aligned}S_\beta - S_\gamma &= - \sum_{\mathbf{s}_\beta} p_\beta(\mathbf{s}_\beta) \ln p_\beta(\mathbf{s}_\beta) + \sum_{\mathbf{s}_\gamma} p_\gamma(\mathbf{s}_\gamma) \ln p_\gamma(\mathbf{s}_\gamma) \leq \\ &\leq - \sum_{\mathbf{s}_\beta} p_\beta(\mathbf{s}_\beta) \ln p'_\beta(\mathbf{s}_\beta) + \sum_{\mathbf{s}_\gamma} p_\gamma(\mathbf{s}_\gamma) \ln p'_\gamma(\mathbf{s}_\gamma) = S'_\beta - S'_\gamma\end{aligned}\quad (95)$$

holds, and it is tighter than the previous bound. A tighter bound typically entail faster convergence.

In order to give an example, consider again the CVM square approximation for a model on a regular square lattice with periodic boundary conditions and focus on the entropy part of the free energy, which according to the entropy expansion Equation (48) has the form

$$- \sum_{\square} S_{\square} + \sum_{\langle ij \rangle} S_{ij} - \sum_i S_i = \sum_{\square} p_{\square} \ln p_{\square} - \sum_{\langle ij \rangle} p_{ij} \ln p_{ij} + \sum_i p_i \ln p_i. \quad (96)$$

This contains both convex (from square and site entropy) and concave terms (from pair entropy). Notice that the numbers of plaquettes is the same as the number of sites, while there are two pairs (e.g. horizontal and vertical) per site. This implies that the free energy is not convex over the constraint set.

Several bounding schemes are possible to define $\mathcal{F}_{\text{conv}}$. For instance, one can obtain a function which is just convex over the constraint set by applying property (iii) to the site terms and half the pair terms, with the result

$$-\sum_{\square} S_{\square} + \sum_{\langle ij \rangle} S_{ij} - \sum_i S_i \leq -\sum_{\square} S_{\square} + \frac{1}{2} \sum_{\langle ij \rangle} S_{ij} + \frac{1}{2} \sum_{\langle ij \rangle} S'_{ij} - \sum_i S'_i. \quad (97)$$

In the following the HAK algorithm will always be used with this bounding scheme.

The NIM can be obtained if, starting from the above expression, one applies property (ii) to the not yet bounded pair terms, with the result

$$-\sum_{\square} S_{\square} + \sum_{\langle ij \rangle} S_{ij} - \sum_i S_i \leq -\sum_{\square} S_{\square} + \sum_{\langle ij \rangle} S'_{ij} - \sum_i S'_i. \quad (98)$$

This is clearly a looser bound than the previous one, and hence it leads to a (much) slower algorithm. In the general case, the NIM (which of course was formulated in a different way) can be obtained by bounding all entropy terms except those corresponding to the maximal clusters. This choice does not always lead to a convex bound (though in most practically relevant cases this happens) and hence convergence is not always guaranteed.

The CCCP recipe corresponds to bounding every convex ($a_{\beta} < 0$) term by

$$-a_{\beta} S_{\beta} \leq -S_{\beta} + (1 - a_{\beta}) S'_{\beta}, \quad (99)$$

using property (ii). In the present case this gives

$$-\sum_{\square} S_{\square} + \sum_{\langle ij \rangle} S_{ij} - \sum_i S_i \leq -\sum_{\square} S_{\square} - \sum_{\langle ij \rangle} S_{ij} + 2 \sum_{\langle ij \rangle} S'_{ij} - \sum_i S_i, \quad (100)$$

which is convex independently of the constraints, and hence the bound is again looser than Equation (97)

In all cases one is left with a double loop algorithm, the outer loop being defined by the update rule for probabilities, and the inner loop being used for the minimization involved in Equation (92). This minimization is simpler than the original problem, since the function to be minimized is convex. In each of the above schemes a particular technique was proposed for the convex minimization in the inner loop, and here these will not be covered in detail.

A point which is important to notice here is that the bounding operation gives a new free energy which is structurally different from a CVM free energy. It must be minimized with respect to p at fixed p' and, viewed as a function of p , it contains an entropy expansion with coefficients \tilde{a}_{β} which do not satisfy anymore the Möbius relation (19) (for instance, in the “just convex over the constraint set” scheme, we have $a_{\square} = 1$, $a_{ij} = -1/2$ and $a_i = 0$). This means that a message-passing algorithm like parent-to-child GBP, which relies on the Möbius property, cannot be applied. In [84] a different message-passing algorithm, which can still be viewed as a GBP algorithm, is suggested.

Observe also that there are entropy-like terms S'_{β} which are actually linear in p_{β} and must therefore be absorbed in the energy terms.

The main reason for investigating these double loop, provably convergent algorithms, is the non-convergence of BP and GBP in frustrated cases. Since BP and GBP, when they converge, are the fastest algorithms for the determination of the minima of the CVM free energy, it is worth making some performance tests to evaluate the speed of the various algorithms. The CPU times reported below refer to an Intel Pentium 4 processor at 3.06 GHz, using g77 under GNU/Linux.

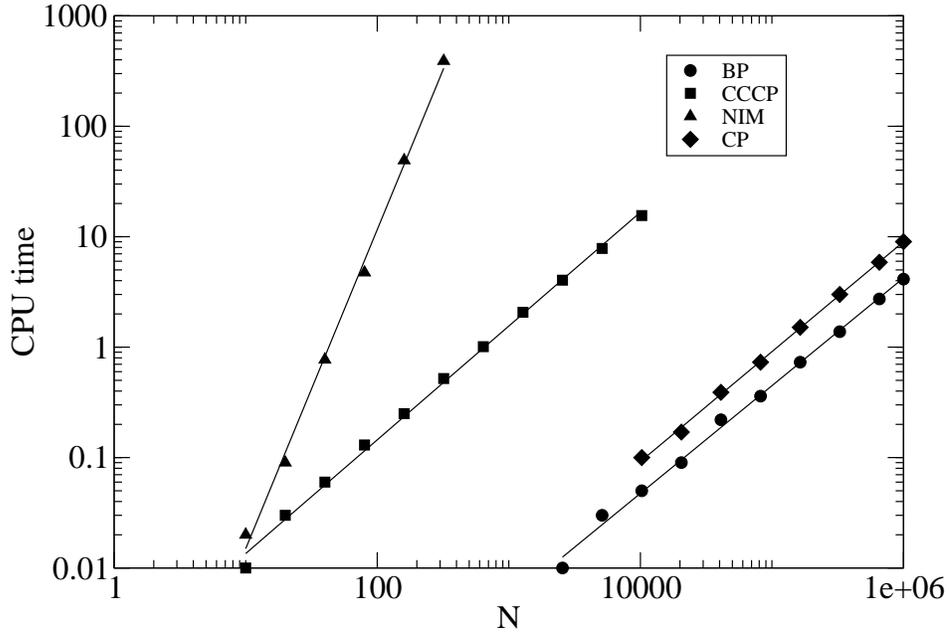


Figure 15. CPU times (seconds) for the 1d Ising chain with random fields

Consider first a chain of N Ising spins, with ferromagnetic interactions $J > 0$ and random bimodal fields h_i independently drawn from the distribution

$$p(h_i) = \frac{1}{2}\delta(h_i - h_0) + \frac{1}{2}\delta(h_i + h_0). \quad (101)$$

The boundary conditions are open, and the model is exactly solved by the CVM pair approximation. The various algorithms described are run from a disordered, uncorrelated state and stopped when the distance between two successive iterations, defined as the sum of the squared variations of the messages (or the probabilities, or the Lagrange multipliers, depending on the algorithm and the loop – outer or inner – considered). Figure 15 reports the CPU times obtained with several algorithms, for the case $J = 0.1$, $h_0 = 1$. The HAK algorithm is not reported since it reduces to BP due to the convexity of the free energy. It is seen that the CPU time grows linearly with N for all algorithms except NIM, in which case it goes like N^3 . Despite the common linear behaviour, there are order of magnitude differences between the various algorithms. While BP and CP converges in 4 and 9 seconds respectively for $N = 10^6$, CCCP takes 15 seconds for $N = 10^4$. For NIM, finally, the fixed point is reached in 12 seconds for $N = 10^2$.

As a further test, consider, again at the level of the pair approximation, the two-dimensional Edwards–Anderson spin glass model, defined by the Hamiltonian Equation (5) with $h_i = 0$ and random bimodal interactions J_{ij} independently drawn from the distribution

$$p(J_{ij}) = (1 - p)\delta(J_{ij} - J) + p\delta(J_{ij} + J). \quad (102)$$

Here the frustration effects are even more important and the non-convergence problem of BP becomes evident. As a rule of thumb, when the temperature, measured by J^{-1} ,

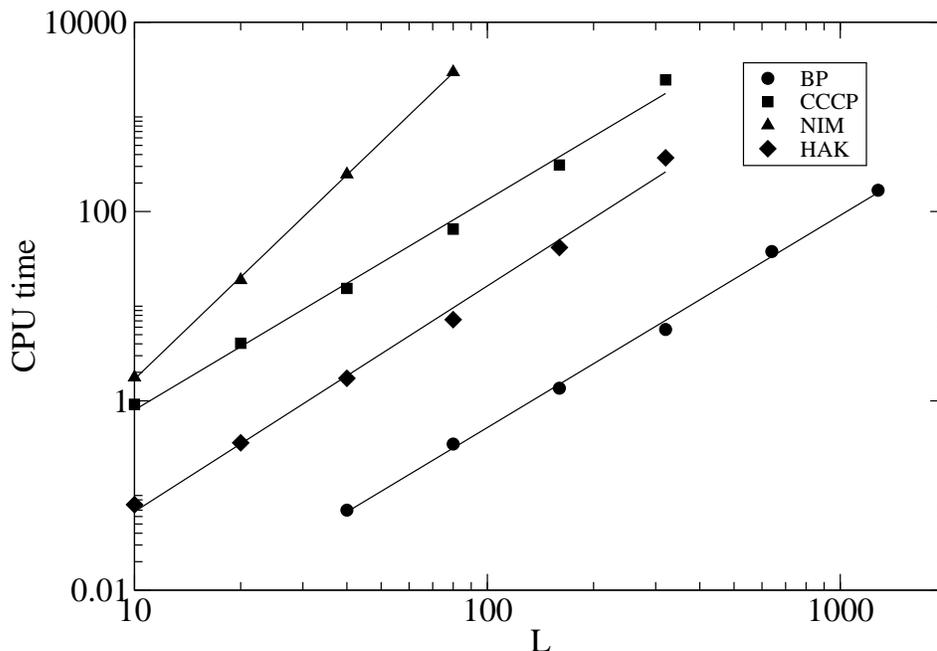


Figure 16. CPU times (seconds) for the 2d Edwards–Anderson model in the paramagnetic phase

is small enough and p (the fraction of antiferromagnetic bonds) is large enough, the BP algorithm stops converging. The condition for the instability of the BP fixed point has been computed, in the average case, for Ising spin glass models with pairwise interactions [93]. In order to compare algorithm performances, Figure 16 reports CPU times vs L for $N = L^2$ lattices with periodic boundary conditions, $J = 0.2$ and $p = 1/2$, that is well into the paramagnetic phase of the model. The initial guess is a ferromagnetic state with $m_i = 0.9, \forall i$. It is seen that the CPU times scale roughly as $N^{1.1}$ for all the algorithms considered except NIM, which goes like $N^{1.8}$. Again the algorithms with linear behaviour are separated by orders of magnitude. For $L = 320$ BP converges in 6 seconds, HAK in 370 seconds and CCCP in 2460 seconds.

CP has not been considered in the present and the following tests, although empirically it is seen that its behaviour is rather close to the HAK algorithm. Its performance is however severely limited as soon as one considers variable with more than two states, due to a sum over the configurations of the neighbourhood of a NN pair.

A similar comparison can be made in the ferromagnetic phase, setting $J = 0.5$ and $p = 0.1$. Here the CPU times for the BP algorithm exhibit large fluctuations for different realizations of the disorder, and the data reported are obtained by averaging over 30 such realizations. Now all algorithms exhibit comparable scaling properties, with CPU times growing like $N^{1.5} \div N^{1.7}$. As far as absolute values are concerned, for $L = 50$ convergence is reached in 4, 44, 680 and 1535 seconds by BP, HAK, CCCP and NIM respectively.

A similar scaling analysis was not possible into the glassy phase (which is unphysically predicted by the pair approximation), due to non-convergence of BP

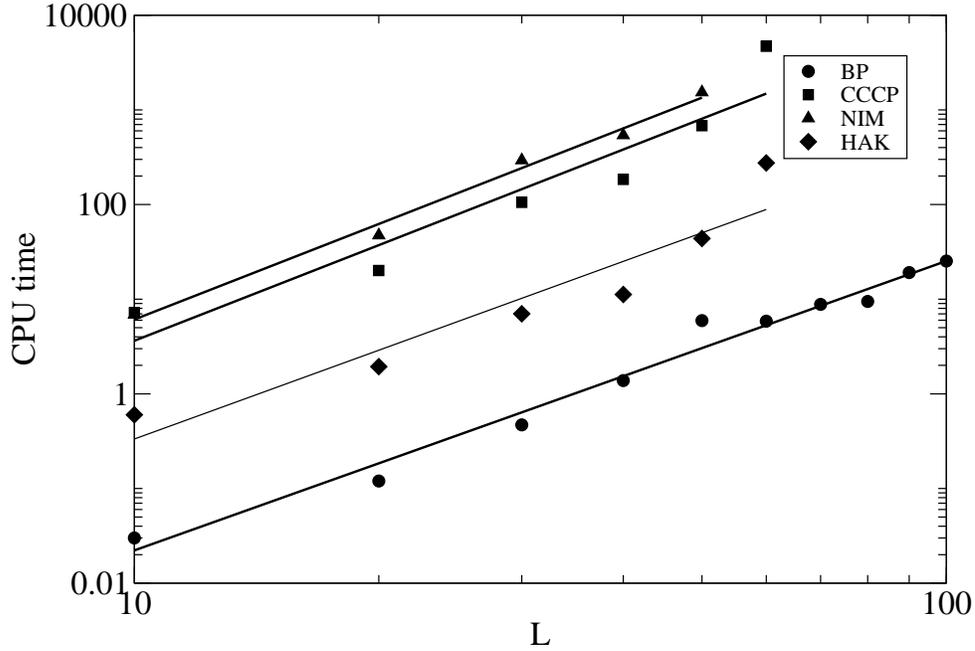


Figure 17. CPU times (seconds) for the 2d Edwards–Anderson model in the ferromagnetic phase

and too large fluctuations of the convergence time of the other algorithms.

As a general remark we observe that BP is the fastest algorithm available whenever it converges. Among the provably convergent algorithms, the fastest one turns out to be HAK, at least in the “just convex over the constraints” [84] scheme which was used here.

7. Conclusions

Some aspects of the cluster variation method have been briefly reviewed. The emphasis was on recent developments, not yet covered by the 1994 special issue of Progress of Theoretical Supplement [8], and the focus was on the methodological aspects rather than on the applications.

The discussion has been based on what can be considered the modern formulation of the CVM, due to An [54], based on a truncation of the cumulant expansion of the entropy in the variational principle of equilibrium statistical mechanics.

The advancements in this last decade were often due to the interaction between two communities of researchers, working on statistical physics and, in a broad sense, probabilistic graphical models for inference and optimization problems. The interest of both communities is currently on heterogeneous problems, while in the past the CVM was most often applied to translation invariant lattice models (in this topic, the only new advancements discussed have been the attempts to extract information about critical behaviour from CVM results). The more general point of view that has to be adopted in studying heterogeneous problems has been crucial to achieve many of the results discussed.

The formal properties of the CVM have been better understood by comparing it with other region-based approximations, like the junction graph method or the most general formulation of the Bethe–Peierls approximation (the lowest order of the CVM), which can treat also non-pairwise interactions. Studying realizability, that is the possibility of reconstructing a global probability distribution from the marginals predicted by the CVM, has led to the discovery of non-tree-like models for which the CVM gives the exact solution.

A very important step was made by understanding that belief propagation, a message-passing algorithm widely used in the literature on probabilistic graphical models, has fixed points which correspond to stationary points of the Bethe–Peierls approximation. The belief propagation can thus be regarded as a powerful algorithm to solve the CVM variational problem, that is to find minima of the approximate free energy, at the Bethe–Peierls level. This opened the way to the formulation of generalized belief propagation algorithms, whose fixed points correspond to stationary points of the CVM free energy, at higher level of approximation.

Belief propagation and generalized belief propagation are certainly the fastest available algorithms for the minimization of the CVM free energy, but they often fail to converge. Typically this happens when the problems under consideration are sufficiently frustrated. In order to overcome this difficulty double loop, provably convergent algorithms have been devised, for which the free energy can be shown to decrease at each iteration. These are similar in spirit to the old natural iteration method by Kikuchi, but orders of magnitude faster, though not as fast as BP and GBP.

When the frustration due to competitive interactions or constraints is very strong, like in spin-glass models in the glassy phase or in constraints satisfaction problems in the hard regime, even double loop algorithms become useless, since we are faced with the problem of replica symmetry breaking, corresponding to clustering of solutions. Very important advancements have been made in recent years by extending the belief propagation algorithm into this domain. These results are in a sense at the border of the CVM, since they are at present confined to the lowest order of the CVM approximation, that is the Bethe–Peierls approximation.

It will be of particular importance, in view of the applications to hard optimization problems with non-tree-like structure, to understand how to generalize these results to higher order approximations.

Acknowledgments

I warmly thank Pierpaolo Bruscolini, Carla Buzano and Marco Pretti, with whom I have had the opportunity to collaborate and to exchange ideas about the CVM, Marco Zamparo for a fruitful discussion about Equation (44), Riccardo Zecchina for many discussions about the survey propagation algorithm, and the organizers of the Lavin workshop “Optimization and inference in machine learning and physics” where I had the opportunity to discuss an early version of this work.

References

- [1] Kikuchi R 1951 *Phys. Rev.* **81** 988
- [2] Bethe H 1935 *Proc. R. Soc. A* **150** 552
- [3] Peierls R E 1936 *Proc. Cambridge Philos. Soc.* **32** 477

- [4] Kramers H A and Wannier G H 1941 *Phys. Rev.* **60** 252
- [5] Kramers H A and Wannier G H 1941 *Phys. Rev.* **60** 263
- [6] Plischke M and Bergersen B 1994 *Equilibrium Statistical Physics* (Singapore: World Scientific Publishing)
- [7] Lavis D A and Bell G M 1999 *Statistical Mechanics of Lattice Systems* (Berlin: Springer)
- [8] Morita T, Suzuki M, Wada K and Kaburagi M (eds) 1994 *Progr. Theor. Phys. Suppl.* **115** (special issue)
- [9] Kikuchi R and Masuda-Jindo K 1999 *Comput. Mater. Sci.* **14** 295
- [10] Díaz-Ortiz A, Sanchez J M and Morán-López J L 1998 *Phys. Rev. Lett.* **81** 1146
- [11] Buzano C and Pelizzola A 1995 *Physica A* **216** 158
- [12] Aguilera-Granja F and Kikuchi R 1994 *Progr. Theor. Phys. Suppl.* **115** 195
- [13] Lise S, Maritan A and Pelizzola A 1998 *Phys. Rev. E* **58** R5241
- [14] Morita T 1994 *Progr. Theor. Phys.* **92** 1081
- [15] Danani A and Pelizzola A 1993 *Mod. Phys. Lett. B* **7** 1761
- [16] Ishii T 1994 *Progr. Theor. Phys. Suppl.* **115** 243
- [17] Ducastelle F 1994 *Progr. Theor. Phys. Suppl.* **115** 255
- [18] Wada K and Kaburagi M 1994 *Progr. Theor. Phys. Suppl.* **115** 273
- [19] Pelizzola A 1994 *Phys. Rev. E* **49** R2503
- [20] Pelizzola A 1995 *J. Magn. Magn. Mat.* **140–144** 1491
- [21] Pelizzola A 1996 *Phys. Rev. E* **53** 5825
- [22] Pelizzola A 2000 *Phys. Rev. E* **61** 4915
- [23] Suzuki M et al 1995 *Coherent Anomaly Method: Mean Field, Fluctuations and Systematics* (Singapore: World Scientific)
- [24] Seino M and Katsura S 1994 *Progr. Theor. Phys. Suppl.* **115** 237
- [25] Mézard M, Parisi G and Zecchina R 2002 *Science* **297** 812
- [26] Mézard M and Zecchina R 2002 *Phys. Rev. E* **66** 056126
- [27] Braunstein A et al 2003 *Phys. Rev. E* **68** 036702
- [28] van Mourik J 2005 *Time averaged belief propagation*, preprint
- [29] Smyth P 1997 *Patt. Rec. Lett.* **18** 1261
- [30] Pearl J 1988 *Probabilistic Reasoning in Intelligent Systems* (San Francisco: Morgan Kaufmann)
- [31] Yedidia J S, Freeman W T and Weiss Y 2001 *Advances in Neural Information Processing Systems* ed T K Leen, T G Dietterich and V Tresp (Cambridge: MIT Press) p 689
- [32] Tanaka K and Morita T 1995 *Phys. Lett. A* **203** 122
- [33] Tanaka K 2002 *J. Phys. A* **35** R81
- [34] Tanaka K, Inoue J and Titterton D M 2003 *J. Phys. A* **36** 11023
- [35] Tanaka K et al 2004 *J. Phys. A* **37** 8675
- [36] Freeman W T, Pasztor E C and Carmichael O T 2000 *Int. J. Comp. Vision* **40** 25
- [37] Shental O, Shental N, Weiss A J and Weiss Y 2004 *Proc. IEEE Information Theory Workshop*
- [38] Gallager R G 1963 *Low-density parity check codes* (Cambridge: MIT Press)
- [39] McEliece R, MacKay D and Cheng J 1998 *IEEE J. Sel. Communication* **16** 140
- [40] Kabashima Y and Saad D 2004 *J. Phys. A* **37** R1
- [41] Kappen H J et al 1999 *Patt. Rec. Lett.* **20** 1231
- [42] Frey B J, Koetter R and Petrovic N 2002 *Advances in Neural Information Processing Systems 14* ed T Dietterich, S Becker and Z Ghahramani (Cambridge: MIT Press) p 737
- [43] Burge C B and Karlin S 1998 *Curr. Opin. Struct. Biol.* **8** 346
- [44] Durbin R et al 1998 *Biological sequence analysis* (Cambridge: Cambridge University Press)
- [45] Krogh A et al 2001 *J. Mol. Biol.* **305** 567
- [46] Huang X et al 2001 *Spoken Language Processing* (New York: Prentice Hall)
- [47] Manning C D and Schütze H 1999 *Foundations of Statistical Natural Language Processing* (Cambridge: MIT Press)
- [48] Kschischang F R, Frey B J and Loeliger H-A 2001 *IEEE Trans. Inf. Theory* **47** 498
- [49] Lauritzen S 1996 *Graphical models* (Oxford: Oxford University Press)
- [50] Barker J A 1953 *Proc. R. Soc. A* **216** 45
- [51] Morita T 1957 *J. Phys. Soc. Jpn.* **12** 753
- [52] Morita T 1972 *J. Math. Phys.* **13** 115
- [53] Schlijper A G 1983 *Phys. Rev. B* **27** 6841
- [54] An G 1988 *J. Stat. Phys.* **52** 727
- [55] Kappen H J and Wiegerinck W 2002 *Advances in Neural Information Processing Systems 14* ed T Dietterich, S Becker and Z Ghahramani (Cambridge: MIT Press) p 415
- [56] Yedidia J S, Freeman W T and Weiss Y 2004 *Constructing Free Energy Approximations and Generalized Belief Propagation Algorithms* (MERL TR-2004-040)

- [57] Aji S M and McEliece R J 2001 *Proceeding of the 41st Allerton Conference on Communication, Control and Computing* p 672
- [58] Lauritzen S and Spiegelhalter D 1988 *J. Royal Stat. Soc. B* **50** 157
- [59] Jensen F 1996 *An introduction to Bayesian networks* (London: UCL Press)
- [60] Cowell R 1998 *Advanced inference in Bayesian networks* in *Learning in graphical models* ed M I Jordan (Cambridge: MIT Press)
- [61] Brascamp H J 1971 *Commun. Math. Phys.* **21** 56
- [62] Percus J K 1977 *J. Stat. Phys.* **16** 299
- [63] Schlijper A G 1984 *J. Stat. Phys.* **35** 285
- [64] Bollobás B 1985 *Random Graphs* (New York: Academic Press)
- [65] Schlijper A G 1988 *J. Stat. Phys.* **50** 689
- [66] Wainwright M J, Jaakkola T S and Willsky A S 2003 *IEEE Trans. Inf. Theory* **49** 1120
- [67] Pelizzola A and Pretti M 1999 *Phys. Rev. B* **60** 10134
- [68] Pelizzola A 2000 *Phys. Rev. B* **61** 11510
- [69] Pelizzola A 2000 *Nucl. Phys. B (Proc. Suppl.)* **83-84** 706
- [70] Sanchez J M 1982 *Physica A* **111** 200
- [71] Cirillo E N M, Gonnella G, Troccoli M and Maritan A 1999 *J. Stat. Phys.* **94** 67
- [72] Wako H and Saitō N 1978 *J. Phys. Soc. Jpn.* **44** 1931
- [73] Wako H and Saitō N 1978 *J. Phys. Soc. Jpn.* **44** 1939
- [74] Muñoz V, Thompson P A, Hofrichter J and Eaton W A 1997 *Nature* **390** 196
- [75] Muñoz V, Henry E R, Hofrichter J and Eaton W A 1998 *Proc. Natl. Acad. Sci. U.S.A.* **95** 5872
- [76] Muñoz V and Eaton W A 1999 *Proc. Natl. Acad. Sci. U.S.A.* **96** 11311
- [77] Bruscolini P and Pelizzola A 2002 *Phys. Rev. Lett.* **88** 258101
- [78] Bruscolini P and Pelizzola A 2002 *Modeling of complex systems* ed P L Garrido and J Marro (New York: AIP) p 205
- [79] Pelizzola A 2005 *Exactness of the cluster variation method and factorization of the equilibrium probability for the Wako-Saitō-Muñoz-Eaton model of protein folding*, submitted to JSTAT
- [80] Schlijper A G 1985 *J. Stat. Phys.* **40** 1
- [81] Kikuchi R and Brush S G 1967 *J Chem Phys* **47** 195
- [82] Pakzad P and Anantharam V 2002 *Proc. Conference on Information Sciences and Systems* paper no. 225
- [83] McEliece R J and Yildirim M 2002 *Mathematical Systems Theory in Biology, Communication, Computation, and Finance* ed D Gilliam and J Rosenthal (Berlin: Springer) p 275
- [84] Heskes T, Albers K and Kappen B 2003 *Uncertainty in Artificial Intelligence: Proceedings of the 19th Conference (UAI-2003)* (San Francisco: Morgan Kaufmann) p 313
- [85] Pelissetto A and Vicari E 2002 *Phys. Rep.* **368** 549
- [86] Morita T 1979 *Physica A* **98** 566
- [87] Katsura S 1986 *Progr. Theor. Phys. Suppl.* **87** 139
- [88] Kfir H and Kanter I 2003 *Physica A* **330** 259
- [89] Tatikonda S and Jordan M I 2002 *Uncertainty in Artificial Intelligence: Proceedings of the 18th Conference (UAI-2002)* (San Francisco: Morgan Kaufmann) p 493
- [90] Mooij J and Kappen H 2005 *Sufficient conditions for convergence of Loopy Belief Propagation*, to appear in Proc. of the 21st Conference on Uncertainty and Artificial Intelligence (UAI 2005)
- [91] Ihler A T, Fisher J W and Willsky A S 2005 *J. Mach. Learn. Res.* **6** 905
- [92] Heskes T 2004 *Neur. Comp.* **16** 2379
- [93] Kabashima Y 2003 *J. Phys. Soc. Jpn.* **72** 1645
- [94] Heskes T 2003 *Advances in Neural Information Processing Systems 15* ed S Becker, S Thrun and K Obermayer (Cambridge: MIT Press) p 343
- [95] Pelizzola A and Pretti M 2003 *J. Phys. A* **36** 11201
- [96] Thouless D J, Anderson P W and Palmer R G 1977 *Phil. Mag.* **35** 593
- [97] Plefka T 1982 *J. Phys. A* **15** 1971
- [98] Mézard M and Parisi G 1986 *Europhys. Lett.* **1** 77
- [99] Mézard M and Parisi G 1987 *Europhys. Lett.* **3** 1067
- [100] Mézard M and Parisi G 2001 *Eur. Phys. J. B* **20** 217
- [101] Mézard M and Parisi G 2003 *J. Stat. Phys.* **111** 1
- [102] Kikuchi R 1974 *J. Chem. Phys.* **60** 1071
- [103] Kikuchi R 1976 *J. Chem. Phys.* **65** 4545
- [104] Kikuchi R, Kokubun H and Katsura S 1986 *J. Phys. Soc. Jpn.* **55** 1836
- [105] Pelizzola A 1994 *Physica A* **211** 107
- [106] Pretti M 2005 *J. Stat. Phys.* in press (cond-mat/0404654)

[107] Yuille A L 2002 *Neur. Comp.* **14** 1691