

Adaptive and Selfaveraging Thouless-Anderson-Palmer Mean Field Theory for Probabilistic Modelling

Manfred Opper

Neural Computing Research Group,
Department of Computer Science and Applied Mathematics,
Aston University, Birmingham B4 7ET, United Kingdom

Ole Winther

Center for Biological Sequence Analysis,
Department of Biotechnology,
Technical University of Denmark, B208, 2800 Lyngby, Denmark.
Informatics and Mathematical Modelling,
Technical University of Denmark, B321, 2800 Lyngby, Denmark.
Theoretical Physics, Lund University, Sölvegatan 14 A,
223 62 Lund, Sweden.

May 8, 2001

Abstract

We develop a generalization of the TAP mean field approach of disorder physics which makes the method applicable to the computation of approximate averages in probabilistic models for real data. In contrast to the conventional TAP approach, where the knowledge of the distribution of couplings between the random variables is required, our method adapts to the concrete set of couplings. We show the significance of approach in two ways: Our approach reproduces replica symmetric results for a wide class of toy models (assuming a nonglassy phase) with given disorder distributions in the thermodynamic limit. On the other hand, simulations on a real data model demonstrate that the method achieves more accurate predictions as compared to conventional TAP approaches.

1 Introduction

Models of statistical physics with random infinite ranged interactions have the nice feature that they can be treated exactly by mean field methods. To compute average properties of the system, one may choose two possible alternative but equivalent approaches. The first one is based on the replica method in which the replicated system is first averaged over the disorder and the resulting nonrandom system is decoupled by saddlepoint methods which leads to exact average case mean field equations [21]. In the second approach one derives a mean field theory for a *fixed set of random couplings* which becomes exact in the thermodynamic limit for almost all realizations of the randomness. This type of mean field theory is traditionally called the TAP approach after *Thouless, Anderson and Palmer* [41] who developed it first for the *Sherrington Kirkpatrick* model [37] of spin glasses. In a final step, one may average TAP mean field equations over the couplings to achieve the same result as in the replica approach [6].

Besides the importance of the TAP approach in the theory of disordered systems there is a recent interest in this method which comes from a more applied area of research dealing with adaptive probabilistic data models (for a review see e.g. [10]). The goal of such models is to explain complex observed data by a set of unobserved, *hidden* random variables based on the joint distribution of both sets of variables. A few popular examples are: Bayes belief networks [32] (used as trainable expert systems), independent component analysis [16] (abbreviated ICA, which detects independent sources in nonlinear signal processing), Gaussian process models [29] (modelling hidden spatial structures by random fields) and Boltzmann machines [8](the Ising version of the random fields).

The price that a modeler has to pay for the high degree of flexibility of these models is the vast increase in computational complexity when the number of hidden variables is large. Both the statistical inference about the hidden variables and the learning of the model parameters requires the computation of marginal distributions of the hidden variables/the observed data, i.e. the evaluation of high dimensional sums or integrals. Since similar types of calculations are ubiquitous in computing thermal averages, e.g. for finding local magnetizations and free energies, there is a great deal of interest in adopting approximation techniques from statistical physics. Already the simple (often called naive) mean field (MF) method, which neglect all correlations of random variables has been applied successfully to a variety of probabilistic data models. At present, there is a growing research activity in the field of probabilistic models trying to overcome the limitations of the simple MF method by partly including the dependencies of variables but still keeping the approximation tractable. In the case where the individual dependencies are weak but their total effect cannot be neglected, the TAP method is a natural candidate for such an improved approximation. TAP approaches for different probabilistic models have already been discussed for neural networks [19, 44, 28], Boltzman machines [11, 39], Gaussian process models

for classification [29], error correcting codes [14], for a review see also [27].

Unfortunately, the TAP mean field approach shows a characteristic difference to the naive MF method which makes its straightforward application to models for real data nontrivial. While the simple MF equations are expressed in terms of the concrete couplings (which encode observed data in applications), the *Onsager correction* to the naive MF theory (for models with extensive connectivities) provided by the TAP approach will explicitly depend on the distribution from which these couplings were generated at random. Two models with the same connectivities but different distributions for the couplings, like e.g. the SK model and the Hopfield model [9] have different expressions for the Onsager corrections (see e.g. [21], chapter XIII). While in the models of disorder physics this distribution is given in advance, such a knowledge is obviously not available for models of real data. Simply taking results from a theory that *assumes* a specific distribution may lead to suboptimal performance. To define the ‘correct’ TAP approach which is valid in these more general situations truncated perturbative expansions for the free energy or for marginal distributions have been considered [33, 2, 1, 13]. While such a finite order truncation becomes exact for the SK model, this is, in general not to be expected.

In this paper we present a new approach to this problem.¹ Our criteria for a valid TAP method are twofold: We require that the lack of knowledge of the underlying distribution of the couplings must be compensated by a selfconsistent computation which *adapts* the Onsager correction to the *concrete* set of couplings. Second, when applied to a set of interactions which was randomly generated from a *known* distribution (for which the mean field assumption is valid), a suitable average of the adaptive TAP method should reproduce the *correct* average case results known from the replica approach. We achieve the first goal by combining the cavity approach [21] with a simple linear response technique which yields a second set of TAP equations for the Onsager correction. This method fulfils the second requirement sofar only for the class of extensively connected models with non-glassy behaviour. These models can be described by a single ergodic phase which is correctly described by a finite set of order parameters (unlike the sparsely connected models [21]) in replica symmetry. Our experience with average case studies of neural networks makes us expect that the assumption of replica symmetry may well describe practical situations when the models are sufficiently matched to the data.

Our approach is most naturally developed for models with pairwise interactions between variables S_i , $i = 1, \dots, N$

$$P(\mathbf{S}) = \frac{\rho(\mathbf{S})}{Z(\boldsymbol{\theta}, \mathbf{J})} \exp \left[\sum_{i < j} S_i J_{ij} S_j + \sum_i S_i \theta_i \right]. \quad (1)$$

Here $\mathbf{S} = (S_1, \dots, S_n)$ and we have set $J_{ii} = 0$. All self-interactions are contained

¹A shorter presentation of the main results of this paper can be found in [30].

in the factorizing distribution $\rho(\mathbf{S}) \equiv \prod_j \rho_j(S_j)$ which also contains all single variable constraints of the variables S_i like their *range*, their *discreteness* etc. Examples of models that are included in this framework are Ising models (like the SK model, the Hopfield model, the Boltzmann machine in the neural computation context), the finite temperature versions of the matching and travelling salesman problems [20, 21], Gaussian process models [29] and the ICA model of [17]. However, many other interesting data models are of a more complicated form such as

$$P(\mathbf{S}) \propto \rho(\mathbf{S}) \exp \left[\sum_{i < j} S_i J_{ij} S_j + \sum_i S_i \theta_i \right] \prod_{k=1}^m F \left(\sum_{i=1}^N \hat{J}_{ki} S_i \right), \quad (2)$$

which includes a variety of popular ‘network’ models like perceptrons (see section 4), sigmoid belief networks [23, 30, 35] and combinatorial optimization problems with inequality constraints, e.g. the knapsack problem [24].

Luckily, the models of the type eq. (2) can be easily represented in the standard form eq. (1) by the ‘field theoretic’ trick of introducing the fields $\sum_{i=1}^N \hat{J}_{ki} S_i$, $k = 1, \dots, m$ as new variables by using Dirac δ -functions and their exponential representations. Denoting the purely imaginary conjugate variables by $\hat{\mathbf{S}} = (\hat{S}_1, \dots, \hat{S}_m)$, the space of variables is augmented to the set $(\mathbf{S}, \hat{\mathbf{S}})$ where the ‘prior distribution’ for the hatted variables is given by

$$\hat{\rho}(\hat{S}) = \int \frac{d\hat{h}}{2\pi i} e^{-\hat{S}\hat{h}} F(\hat{h}) \quad (3)$$

and the augmented coupling matrix is

$$\mathbf{J}_{\text{aug}} = \begin{pmatrix} \mathbf{J} & \hat{\mathbf{J}}^T \\ \hat{\mathbf{J}} & 0 \end{pmatrix}. \quad (4)$$

Since all the subsequent manipulations are of the analytic type, ie. they are based on certain formal expansions rather than on probabilistic arguments (in the sense of assuming positive normalized measures), we expect that our use of nonpositive and even complex measures will not be problematic.

The paper is organized as follows: The adaptive TAP equations are derived in section 2 (with a summary given in section 2.5). In section 3, we show how the adaptive theory reproduces the correct average case results when applied to a fairly general class of distributions for the couplings. We derive ‘selfaveraging’ TAP equations, replica results and the stability condition for the mean field solution (the AT condition). In section 4, we apply our results to the SK-model, the Hopfield network and the simple perceptron. Finally, we present an outlook in section 5.

2 Adaptive TAP Approach

In this section, we will derive both an adaptive TAP approximation for the marginal distribution $P_i(S) \equiv \int \prod_{j \neq i} dS_j P(\mathbf{S})$ (section 2.1-2.3) and the free energy $F(\mathbf{J}, \boldsymbol{\theta}) = -\ln Z(\mathbf{J}, \boldsymbol{\theta})$ (section 2.4). The free energy corresponds to the negative log probability of the observed data which can be used as a yardstick for deciding which model best fits the data.

Our derivation will be based on the cavity approach introduced by [21]. We will assume that we are not dealing with a glassy system with its many ergodic components, but that all averages are for a single state. This is as discussed above usually expected to hold when the probabilistic model is well matched to the data. We expect that the adaptive TAP approximation can be extended to glassy systems along the line of chapter V in [21].

2.1 The Marginal Distribution

The starting point of our derivation is the following exact equation for the marginal distribution of the variable S_i

$$P_i(S_i) = \frac{\int \prod_{j, j \neq i} dS_j \rho(S_i) e^{S_i(\sum_j J_{ij} S_j + \theta_i)} P(\mathbf{S} \setminus S_i)}{\int \prod_j dS_j \rho(S_i) e^{S_i(\sum_j J_{ij} S_j + \theta_i)} P(\mathbf{S} \setminus S_i)}, \quad (5)$$

where $P(\mathbf{S} \setminus S_i)$ is the distribution of all variables for a system where the i th variable is absent. We see from eq. (5) that S_i interacts with the remaining variables only through the field $h_i = \sum_j J_{ij} S_j$. Hence, we introduce its ‘cavity’ distribution, i.e. the distribution of the field at the ‘position’ of the ‘empty’ site i by

$$P(h_i \setminus S_i) = \int \prod_{j \neq i} dS_j \delta\left(h_i - \sum_j J_{ij} S_j\right) P(\mathbf{S} \setminus S_i) \quad (6)$$

and rewrite (5) as

$$P_i(S_i) = \frac{\rho_i(S_i)}{Z_0^{(i)}} e^{-H_i(S_i)}, \quad (7)$$

where we have introduced an effective single variable Hamiltonian

$$-H_i(S) = \ln \left\langle e^{Sh_i} \right\rangle_{\setminus i}, \quad (8)$$

and the brackets $\langle \dots \rangle_{\setminus i}$ denote an average with respect to $P(h_i \setminus S_i)$, eq. (6). The corresponding partition function is

$$Z_0^{(i)} = \int dS \rho_i(S) e^{-H_i(S)}. \quad (9)$$

The complete knowledge of $H_i(S)$ would provide us with the ability to compute averages of functions of a single variable S_i like e.g.

$$\langle S_i \rangle = \frac{\partial}{\partial \theta_i} \ln Z_0^{(i)}. \quad (10)$$

Second, by using appropriate derivatives with respect to the external fields θ_j we could compute correlation functions. E.g., the connected correlation function of two variables is expressed by the *linear response relation* as

$$\chi_{ij} \equiv \langle S_i S_j \rangle - \langle S_i \rangle \langle S_j \rangle = \frac{\partial \langle S_i \rangle}{\partial \theta_j} = \frac{\partial^2 \ln Z_0^{(i)}}{\partial \theta_i \partial \theta_j}. \quad (11)$$

Moreover, from the definition (8), we also realize that the distribution of h_i can be reconstructed using derivatives with respect to S . A corresponding result which will be useful in the following is

$$\langle h_i \rangle = \frac{1}{Z_0^{(i)}} \int dS \rho_i(S) \frac{\partial}{\partial S} e^{-H_i(S)}. \quad (12)$$

2.2 The Cavity Approach

In general, we can express eq. (8) by the cumulants $\kappa_k^{(i)}$ of the *cavity* distribution

$$-H_i(S) = \sum_{k=1}^{\infty} \frac{\kappa_k^{(i)}}{k!} S^k \quad (13)$$

with $\kappa_1^{(i)} = \langle h_i \rangle_{\setminus i}$, $\kappa_2^{(i)} = \langle h_i^2 \rangle_{\setminus i} - \langle h_i \rangle_{\setminus i}^2$ etc.

The usual argument in deriving TAP equations is based on the assumption of weak dependencies between the random variables \mathbf{S} which is expressed in the so-called *clustering hypothesis* [21]

$$\frac{1}{N^2} \sum_{ij} (\langle S_i S_j \rangle - \langle S_i \rangle \langle S_j \rangle)^2 \rightarrow 0 \quad (14)$$

and similar relations for higher connected correlation functions. One concludes that if the variable S_i is removed from the system, the effect of the correlations between the S'_j s on the distribution of the field h_i is so weak that (as in proofs of the central limit theorem by characteristic functions) one can neglect all cumulants of order greater than 2, i.e.,

$$\kappa_k^{(i)} = 0 \quad \text{for} \quad k > 2. \quad (15)$$

If \mathbf{S} is a real random vector this is equivalent to approximating (6) by a Gaussian distribution [21, 29], setting

$$P(h_i \setminus S_i) \approx \frac{1}{\sqrt{2\pi V_i}} \exp\left(-\frac{(h_i - \langle h_i \rangle_{\setminus i})^2}{2V_i}\right), \quad (16)$$

where $V_i \equiv \kappa_2^{(i)} = \langle h_i^2 \rangle_{\setminus i} - \langle h_i \rangle_{\setminus i}^2$. From eq. (15), we immediately get the marginal distribution

$$P_i(S) = \frac{1}{Z_0^{(i)}} \rho_i(S) \exp \left(S(\langle h_i \rangle_{\setminus i} + \theta_i) + \frac{V_i}{2} S^2 \right), \quad (17)$$

and the single variable partition function

$$Z_0^{(i)} = \int dS \rho_i(S) \exp \left(S(\langle h_i \rangle_{\setminus i} + \theta_i) + \frac{V_i}{2} S^2 \right). \quad (18)$$

In the following, we will assume the validity of eq. (15) and the resulting eq. (17) also in the cases, where S is a complex variable with nonreal prior distribution $\rho(S)$.

All that remains to derive the TAP equations is to compute the sets of the first two cumulants $\langle h_i \rangle_{\setminus i}$ and V_i for $i = 1, \dots, N$ selfconsistently. The first cumulant is easily found from eq. (12) using eq. (17) as

$$\langle h_i \rangle = \langle h_i \rangle_{\setminus i} + V_i \langle S_i \rangle. \quad (19)$$

Hence, we can eliminate $\langle h_i \rangle_{\setminus i}$ in favour of the mean field variables $\langle S_i \rangle$ and V_i , $i = 1, \dots, N$ via

$$\langle h_i \rangle_{\setminus i} = \sum_j J_{ij} \langle S_j \rangle - V_i \langle S_i \rangle. \quad (20)$$

This has the well known structure of a mean field corrected by a so called *Onsager reaction* term which accounts for the nontrivial correlations between variables that are neglected in a naive mean field approach.

So far, this approach is well known. The new aspect of our paper is in the way we compute the V_i 's. A naive computation would lead to

$$V_i = \sum_{j,k} J_{ij} J_{ik} \left(\langle S_j S_k \rangle_{\setminus i} - \langle S_j \rangle_{\setminus i} \langle S_k \rangle_{\setminus i} \right) \approx \sum_j J_{ij}^2 \left(\langle S_j^2 \rangle_{\setminus i} - \langle S_j \rangle_{\setminus i}^2 \right), \quad (21)$$

neglecting the nondiagonal correlations. This turns out to be correct in the thermodynamic limit $N \rightarrow \infty$ for models (like the SK model) where different pairs of couplings J_{ij} and J_{ik} are drawn *independently* at random. However, it fails when the J_{ij} 's become weakly correlated. Consider e.g. the simple Gaussian model $\rho(\mathbf{S}) \propto e^{-\frac{1}{2}|\mathbf{S}|^2}$ with a 'Hopfield type' coupling matrix defined by $J_{ij} = \frac{1}{N} \sum_{k=1}^{\alpha N} x_i^k x_j^k$, where x_i^k iid random variables with zero mean and unit variance. The corresponding covariance

$$\langle S_i S_j \rangle - \langle S_i \rangle \langle S_j \rangle = (\mathbf{I} - \mathbf{J})_{ij}^{-1} \quad (22)$$

would still fulfil the clustering condition eq. (14), but non-diagonal contributions to V_i do not vanish because of higher order correlation, e.g., $\sum_{jk} J_{ij} J_{ik} J_{jk} \approx \alpha$ do not vanish.

While it is well known how to derive TAP equations for this type of coupling matrix \mathbf{J} (see e.g. [21, 15]), these approaches require the explicit knowledge of the statistics of the J_{ij} 's. In the following section, we aim at a computation of the cavity field variances V_i 's which does not assume such a knowledge.

2.3 Computing the Second Moments

Our derivation is based on a selfconsistent computation of the matrix of susceptibilities $\chi_{ij} = \frac{\partial \langle S_i \rangle}{\partial \theta_j} = \langle S_i S_j \rangle - \langle S_i \rangle \langle S_j \rangle$ based on the mean eq. (10). Hence, the diagonal elements χ_{ii} , $i = 1, \dots, N$ are expressed both by linear response and by the explicit result $\chi_{ii} = \langle S_i^2 \rangle - \langle S_i \rangle^2$ which can be evaluated using the marginal distribution (17). Equating the two expressions we obtain implicit equations for the variances $V_i = \langle h_i^2 \rangle_{\setminus i} - \langle h_i \rangle_{\setminus i}^2$, $i = 1, \dots, N$. Selfinteractions $V_i \langle S_i \rangle$ determined by the the linear response method have also been introduced in [12] as a heuristics to correct the naive MF equations for Boltzmann machines.

Our crucial approximation in the linear response calculation is that it is sufficient to include a perturbation of the means of the cavity fields $\langle h_i \rangle_{\setminus i}$ whereas the variances V_i are kept unchanged. This is consistent with the fact that the V_i 's become selfaveraging quantities in a suitable thermodynamic limit framework where the mean field method becomes exact. Hence, by differentiating eq. (10) with respect to the external field θ_j eq. (18), we get

$$\chi_{ij} = \frac{\partial \langle S_i \rangle}{\partial \theta_i} \left(\delta_{ij} + \frac{\partial \langle h_i \rangle_{\setminus i}}{\partial \theta_j} \right),$$

where $\frac{\partial \langle S_i \rangle}{\partial \theta_i}$ is the explicit derivative of eq. (10). Further differentiating eq. (20), we finally get

$$\chi_{ij} = \frac{\partial \langle S_i \rangle}{\partial \theta_i} \left(\delta_{ij} + \sum_k (J_{ik} - V_k \delta_{ik}) \chi_{kj} \right) \quad (23)$$

which can be solved with respect to $\chi = \{\chi_{ij}\}$ and yields

$$\chi = (\mathbf{\Lambda} - \mathbf{J})^{-1}. \quad (24)$$

Here we have introduced the diagonal matrix

$$\mathbf{\Lambda} = \text{diag}(\Lambda_1, \dots, \Lambda_N) \quad \Lambda_i \equiv V_i + 1/\chi_{ii}. \quad (25)$$

Specializing to the diagonal elements $\chi_{ii} = \langle S_i^2 \rangle - \langle S_i \rangle^2$ determines V_i implicitly via

$$\chi_{ii} = \left[(\mathbf{\Lambda} - \mathbf{J})^{-1} \right]_{ii}. \quad (26)$$

The requirement that the susceptibility matrix (i.e. the matrix of covariances) must be positive definite can be used to test whether the mean field solution is consistent. In the thermodynamic limit this requirement leads to a criterion which is equivalent to the well known *de Almeida–Thouless* stability condition [21] (see section 3.4).

Two important remarks should be made at this point: Although all approximations are expected to be exact in general only in a suitable thermodynamic limit framework (see section 3), the final result (26) is correct for a Gaussian

model for *arbitrary* N (see appendix A). Secondly, we note that the functional relationship between the χ_{ii} and the V_i , eq. (26) is independent of the specific single spin measure $\rho(S)$. This argument can be used for a derivation of the free energy different from the one of the following section.

2.4 The Adaptive TAP Free Energy

In this section, we will derive a TAP approximation to the free energy

$$F(\mathbf{J}, \boldsymbol{\theta}) = -\ln Z(\mathbf{J}, \boldsymbol{\theta})$$

using the adaptive form of the Onsager term given by (26). For this purpose, it is useful to generalize the model eq. (1) to a one parameter class of models where the interaction \mathbf{J} is replaced by $l\mathbf{J}$ with $0 \leq l \leq 1$, ie.

$$Z(l\mathbf{J}, \boldsymbol{\theta}) = \int d\mathbf{S} \prod_i \rho(S_i) \exp\left(\frac{l}{2}\mathbf{S}^T \mathbf{J} \mathbf{S} + \mathbf{S}^T \boldsymbol{\theta}\right). \quad (27)$$

Since the solutions of the TAP equations provide us with the moments $m_i \equiv \langle S_i \rangle$ and M_i for $i = 1, \dots, N$, we will work with the *Gibbs free energy*, i.e. the free energy for *fixed* m_i and M_i which is defined by a Legendre transform using external fields γ_i and λ_i conjugate to m_i and $M_i \equiv \langle S_i^2 \rangle$, i.e.

$$\Phi_l(\mathbf{m}, \mathbf{M}) = \underset{\boldsymbol{\lambda}, \boldsymbol{\gamma}}{\text{extr}} \Psi(\boldsymbol{\lambda}, \boldsymbol{\gamma}, \mathbf{m}, \mathbf{M}) \quad (28)$$

where

$$\Psi_l(\boldsymbol{\lambda}, \boldsymbol{\gamma}, \mathbf{m}, \mathbf{M}) = -\ln Z(l\mathbf{J} + \boldsymbol{\lambda}, \boldsymbol{\theta} + \boldsymbol{\gamma}) + \sum_i \gamma_i m_i + \sum_i \frac{\lambda_i}{2} M_i. \quad (29)$$

$\boldsymbol{\lambda}$ is a diagonal matrix with entries λ_i . The values \mathbf{m}^e and \mathbf{M}^e which make Φ_l stationary, ie. for which $\partial_{\mathbf{m}} \Phi_l = \partial_{\mathbf{M}} \Phi_l = 0$ determine the correct equilibrium expectation values: $\langle S_i \rangle_l = m_i^e$ and $\langle S_i^2 \rangle_l = M_i^e$ (where the index indicates that the expectation is taken with parameter l). Furthermore, by inserting these values back into Φ_1 , the original free energy $F(\mathbf{J}, \boldsymbol{\theta}) = -\ln Z(\mathbf{J}, \boldsymbol{\theta}) = \Phi_1(\mathbf{m}^e, \mathbf{M}^e)$ is recovered.

We compute the TAP approximation to Φ_1 from the relation

$$\Phi_1 = \Phi_0 + \int_0^1 dl \frac{\partial \Phi_l}{\partial l} = \Phi_0 - \frac{1}{2} \int_0^1 dl \left\{ \sum_{i,j} m_i J_{ij} m_j + \text{Tr}(\boldsymbol{\chi}_l \mathbf{J}) \right\} \quad (30)$$

with $\chi_{l,ij} = \langle S_i S_j \rangle_l - \langle S_i \rangle_l \langle S_j \rangle_l$. Note, that the derivatives of λ_i and γ_i with respect to l disappear from (30) because of $\partial_{\gamma_i} \Psi = \partial_{\lambda_i} \Psi = 0$.

We next insert our TAP approximation $\boldsymbol{\chi}_l = (\boldsymbol{\Lambda}_l - l\mathbf{J})^{-1}$ (24) with $\Lambda_{l,i} = V_{l,i} + 1/\chi_{ii}$, into (30) and integrate with respect to l . Note, that $\chi_{ii} = M_i - m_i^2$ is a fixed quantity which does not depend on l . Using

$$\begin{aligned} \text{Tr}(\boldsymbol{\chi}_l \mathbf{J}) &= \text{Tr}((\boldsymbol{\Lambda}_l - l\mathbf{J})^{-1} \mathbf{J}) = -\frac{d}{dl} \text{Tr} \ln(\boldsymbol{\Lambda}_l - l\mathbf{J}) + \text{Tr}(\boldsymbol{\chi}_l \frac{\partial \boldsymbol{\Lambda}_l}{\partial l}) \\ &= -\frac{d}{dl} \text{Tr} \ln(\boldsymbol{\Lambda}_l - l\mathbf{J}) + \sum_i \chi_{ii} \frac{\partial \Lambda_{l,i}}{\partial l}. \end{aligned}$$

and noting that $V_{0,i} = 0$, we obtain

$$\begin{aligned} \Phi_1 &= \Phi_0 - \frac{1}{2} \sum_{ij} m_i J_{ij} m_j + \Delta\Phi \\ \Delta\Phi &= \frac{1}{2} \ln \det(\boldsymbol{\Lambda} - \mathbf{J}) - \frac{1}{2} \sum_i V_i \chi_{ii} + \frac{1}{2} \sum_i \ln \chi_{ii}. \end{aligned} \quad (31)$$

The first two terms constitute the naive mean field approximation to Φ and the last term $\Delta\Phi$ is the Onsager correction. Note, that this result is not equivalent to a truncation of a power series expansion of Φ to second order in l (often termed the Plefka expansion [33]) but contains terms of all orders. It is easy to see that we also recover the TAP equations from the equilibrium conditions $\partial_{\mathbf{m}} \Phi_s = \partial_{\mathbf{M}} \Phi_s = 0$. We finally identify the conjugate fields as the mean and the variance of the cavity field via $\gamma_i = \langle h_i \rangle_{\chi_i}$ and $\lambda_i = V_i$.

In section 3 we give a simplification of the free energy for the thermodynamic limit when the distribution of the couplings \mathbf{J} is explicitly given, i.e. for the conventional TAP approach. Assuming that the free energy is selfaveraging in this case, we show the equivalence of our TAP approach to the results of a replica calculation.

In appendix B, an alternative derivation of the TAP free energy is given. It is based on the observation that the functional form (as a function of m_i and M_i) of the Onsager term V_i in the TAP equations does not depend on the specific single variable densities $\rho(\mathbf{S})$. Hence, we may compute the corresponding universal form of $\Delta\Phi$ by calculating Φ for an exactly solvable model, i.e. for a Gaussian ρ and subtract the naive mean field part. This is the strategy used by Parisi and Potters [31] to derive TAP equations for a spin glass model with an orthogonal random matrix \mathbf{J} .

2.5 Summary of Adaptive TAP Equations

To summarize our results so far, we write the adaptive TAP equations for $\langle S_i \rangle$ and V_i , $i = 1, \dots, N$ as

$$\langle S_i \rangle = \frac{\partial}{\partial \theta_i} \ln Z_0^{(i)} \quad (32)$$

where the single variable partition function is

$$Z_0^{(i)} = \int dS \rho_i(S) e^{S(\sum_j J_{ij}\langle S_j \rangle - V_i\langle S_i \rangle + \theta_i) + \frac{V_i}{2}S^2} . \quad (33)$$

The second set of TAP equations for V_i is

$$\frac{\partial \langle S_i \rangle}{\partial \theta_i} = \frac{\partial^2}{\partial \theta_i^2} \ln Z_0^{(i)} = [(\mathbf{\Lambda} - \mathbf{J})^{-1}]_{ii} \quad (34)$$

where

$$\mathbf{\Lambda} = \text{diag}(\Lambda_1, \dots, \Lambda_N) \quad \Lambda_i \equiv V_i + \left(\frac{\partial \langle S_i \rangle}{\partial \theta_i} \right)^{-1} . \quad (35)$$

Note, that the partial derivatives are now taken with respect to the explicit θ_i dependence, i.e. all remaining arguments in $Z_0^{(i)}$ are *fixed*. Finally, the free energy is given by eq. (31) with $m_i = \langle S_i \rangle$, $M_i = \langle S_i^2 \rangle$, $\gamma_i = \langle h_i \rangle_{\setminus i}$ and $\lambda_i = V_i$.

2.6 The Generalized Model

The special structure of the augmented coupling matrix eq. (4) allows for a variety of simplifications in treating the model eq. (2). By introducing hatted variables $\langle \hat{\mathbf{S}} \rangle$, $\hat{\mathbf{V}}$ and $\hat{\mathbf{\Lambda}}$ explicitly, the previous results for the means eqs. (32) and (33) are $\langle S_i \rangle = \frac{\partial}{\partial \theta_i} \ln Z_0^{(i)}$ and $\langle \hat{S}_k \rangle = \frac{\partial}{\partial \hat{\theta}_k} \ln \hat{Z}_0^{(k)}$ with

$$Z_0^{(i)} = \int dS \rho_i(S) e^{S(\sum_j J_{ij}\langle S_j \rangle + \sum_k \hat{J}_{ki}\langle \hat{S}_k \rangle - V_i\langle S_i \rangle + \theta_i) + \frac{V_i}{2}S^2} \quad (36)$$

$$\hat{Z}_0^{(k)} = \int d\hat{S} \hat{\rho}_k(\hat{S}) e^{\hat{S}(\sum_i \hat{J}_{ki}\langle S_i \rangle - \hat{V}_k\langle \hat{S}_k \rangle + \hat{\theta}_k) + \frac{\hat{V}_k}{2}\hat{S}^2} \quad (37)$$

The augmented susceptibility matrix is given by

$$\boldsymbol{\chi}_{\text{aug}} = (\mathbf{\Lambda}_{\text{aug}} - \mathbf{J}_{\text{aug}})^{-1} = \begin{pmatrix} \mathbf{\Lambda} - \mathbf{J} & -\hat{\mathbf{J}}^T \\ -\hat{\mathbf{J}} & \hat{\mathbf{\Lambda}} \end{pmatrix}^{-1} = \begin{pmatrix} \boldsymbol{\chi} & \tilde{\boldsymbol{\chi}}^T \\ \tilde{\boldsymbol{\chi}} & \hat{\boldsymbol{\chi}} \end{pmatrix} ,$$

where $\tilde{\chi}_{ki} \equiv \frac{\partial \langle \hat{S}_k \rangle}{\partial \theta_i}$. It can be shown that although the vector $\hat{\mathbf{S}}$ is purely imaginary, its expectation $\langle \hat{\mathbf{S}} \rangle$ and the susceptibility matrix $\boldsymbol{\chi}_{\text{aug}}$ come out real.

The TAP equations for the V_i 's can be simplified using identities for the inverse and the determinant of partitioned matrices [34]:

$$\boldsymbol{\chi} = (\mathbf{\Lambda} - \mathbf{J} - \hat{\mathbf{J}}^T \hat{\mathbf{\Lambda}}^{-1} \hat{\mathbf{J}})^{-1} \quad (38)$$

and $\det(\mathbf{\Lambda}_{\text{aug}} - \mathbf{J}_{\text{aug}}) = \det \hat{\mathbf{\Lambda}} \det(\mathbf{\Lambda} - \mathbf{J} - \hat{\mathbf{J}}^T \hat{\mathbf{\Lambda}}^{-1} \hat{\mathbf{J}})$ which leads to

$$\hat{\chi}_{kk} = \frac{\partial}{\partial \hat{\Lambda}_k} \ln \det(\mathbf{\Lambda}_{\text{aug}} - \mathbf{J}_{\text{aug}}) = \frac{1}{\hat{\Lambda}_k} + \frac{1}{\hat{\Lambda}_k^2} \sum_{ij} \hat{J}_{ki} \hat{J}_{kj} \chi_{ij} , \quad (39)$$

showing that the hatted covariances can be explicitly computed from the non-hatted ones. The variances V_i and \hat{V}_k are obtained from a straight forward generalization of eq. (34): $\frac{\partial \langle S_i \rangle}{\partial \theta_i} = \left[\left(\mathbf{\Lambda} - \mathbf{J} - \hat{\mathbf{J}}^T \hat{\mathbf{\Lambda}}^{-1} \hat{\mathbf{J}} \right)^{-1} \right]_{ii}$ and $\frac{\partial \langle \hat{S}_k \rangle}{\partial \theta_k} = \hat{\chi}_{kk}$. Finally, we can use the identity for partitioned determinants to show that in the Onsager term eq. (32) a few terms cancel and we can write

$$\begin{aligned} \Phi_{\text{Ons}} &= \frac{1}{2} \ln \det \left(\mathbf{\Lambda} - \mathbf{J} - \hat{\mathbf{J}}^T \hat{\mathbf{\Lambda}}^{-1} \hat{\mathbf{J}} \right) - \frac{1}{2} \sum_i V_i \chi_{ii} + \frac{1}{2} \sum_i \ln \chi_{ii} \\ &\quad - \frac{1}{2} \sum_k \hat{V}_k \hat{\chi}_{kk} + \frac{1}{2} \sum_k \ln(1 + \hat{V}_k \hat{\chi}_{kk}) . \end{aligned} \quad (40)$$

Finally, we note that for the consistency of the TAP equations, only the positive definiteness of the submatrix χ for the original real random variables \mathbf{S} (and *not* χ_{aug}) is required.

3 Thermodynamic Limit and Selfaveraging Theory

For specific choices of the distribution of disorder where different sites i appear in a symmetric way, we can expect that quantities like V_i and the free energy become selfaveraging in the thermodynamic limit $N \rightarrow \infty$, i.e. $V_i = V$ independent of the specific realization of the disorder. Such quantities can be computed by suitable quenched averages using the statistics of the disorder variables J_{ij} .

As usual, quenched averages over the distribution of the coupling matrix \mathbf{J} are performed within the replica framework using $[\ln Z]_{\mathbf{J}} = \lim_{n \rightarrow 0} \frac{1}{n} \ln [Z^n]_{\mathbf{J}}$, where for integer n , we have

$$[Z^n]_{\mathbf{J}} = \int d\mathbf{S}^n \prod_{ia} \rho(S_{ia}) \left[e^{\frac{1}{2} \sum_a \sum_{ij} S_{ia} J_{ij} S_{ja}} \right]_{\mathbf{J}} . \quad (41)$$

$a = 1, \dots, n$ are replica indices and the brackets $[\dots]_{\mathbf{J}}$ denote the average over the J'_{ij} 's.

If all matrix elements J_{ij} for $i < j$ are assumed to be iid Gaussian random variables (as for the SK model) the average over \mathbf{J} is easily carried out. The simplest way to generalize this *Gaussian orthogonal ensemble* and allow for correlations between J_{ij} 's is to keep the orthogonality of the ensemble but make it non Gaussian. Such distributions which also lead to a well defined thermodynamic limit (introduced by Ref. [18] and subsequently used by Ref. [31]) are defined by generating functions of the type

$$\left[e^{\frac{1}{2} \text{Tr} \mathbf{A} \mathbf{J}} \right]_{\mathbf{J}} = e^{N \text{Tr} G(\mathbf{A}/N)} , \quad (42)$$

with the function G fully specifying the ensemble. In appendix D, it is shown how G is related to the spectrum of the matrix \mathbf{J} .² The Gaussian ensemble is recovered by setting $G(x) \propto x^2$. Note, that the scaling with $1/N$ inside of G keeps the trace of order one for $N \rightarrow \infty$ when the elements of the matrix \mathbf{A} are of order one.

Distributions with generating functions eq. (42) have the nice feature that the average eq. (41) depends only on a single set of quadratic orderparameters given by $q_{ab} \equiv \frac{1}{N} \sum_i S_{ia} S_{ib}$.³ This can be seen by we applying eq. (42) to the matrix \mathbf{A} $A_{ij} = \sum_{a=1}^n S_{ia} S_{ja}$ appearing in eq. (41). We note that \mathbf{A} has $N - n$ eigenvalues equal to zero, but in the space spanned by the n vectors \mathbf{S}_a we get

$$\frac{1}{N} \sum_j A_{ij} S_{ja} = \sum_a q_{ab} S_{ib}. \quad (43)$$

Hence, in this n dimensional subspace, the matrix \mathbf{A}/N acts as the matrix $\mathbf{q} = \{q_{ab}\}$.

In this way we can derive general results for the selfaveraging case and connect these with the previous results for the adaptive TAP theory. In sections 3.1, 3.2 and 3.3, we compute the Onsager term, the replica free energy and the average of the TAP free energy and show that both coincide. Finally, in section 3.4, we show how the condition of a positive definite susceptibility matrix χ translates to the AT stability condition.

3.1 The Onsager Term

In this section, we compute $V = [V_i]_{\mathbf{J}}$ and $\Delta\Phi^{\text{self}} = [\Delta\Phi]_{\mathbf{J}}$ from (31) for the model eq. (1). We will briefly sketch how to generalize these results to the model eq. (2) at the end of the section. To compute V , we write eq. (34) as

$$\chi_{ii} = [(\mathbf{\Lambda} - \mathbf{J})^{-1}]_{ii} = \frac{\partial}{\partial \Lambda_i} \ln \det(\mathbf{\Lambda} - \mathbf{J}) \quad (44)$$

and replace the right hand side by its average over the distribution of the matrix \mathbf{J} for $N \rightarrow \infty$. Since there is no spin glass ordering for a Gaussian model, it is sufficient to perform an annealed average sufficient and we can use the identity

$$\det^{-\frac{1}{2}}(\mathbf{\Lambda} - \mathbf{J}) = \int \frac{d\mathbf{z}}{(2\pi)^{N/2}} e^{-\frac{1}{2}\mathbf{z}^T(\mathbf{\Lambda} - \mathbf{J})\mathbf{z}}. \quad (45)$$

Applying eq. (42) to the matrix $\mathbf{A} = \mathbf{z}\mathbf{z}^T$ which has a single eigenvalue equal to $\mathbf{z}^T\mathbf{z}$ (with eigenvector \mathbf{z}) and a $N - 1$ fold degenerate eigenvalue 0 and using the

²In section 3.1, we discuss the extension of this assumption to the model eq. (2).

³Our ensembles do not apply to diluted models, for which usually an infinite sequence of orderparameters of any order has to be considered.

fact that $G(0) = 0$ we arrive at

$$\begin{aligned} [\det^{-\frac{1}{2}}(\mathbf{\Lambda} - \mathbf{J})]_{\mathbf{J}} &= \int \frac{d\mathbf{z}}{(2\pi)^{N/2}} e^{-\frac{1}{2} \sum_i \Lambda_i z_i^2 + NG(\frac{1}{N} \sum_i z_i^2)} \\ &= \int \frac{dr d\hat{r}}{4\pi i/N} \frac{d\mathbf{z}}{(2\pi)^{N/2}} e^{-\frac{1}{2} \sum_i \Lambda_i z_i^2 + \frac{1}{2} \hat{r} (\sum_i z_i^2 - Nr) + NG(r)} , \end{aligned}$$

where in the last line the orderparameter $r = \frac{1}{N} \sum_i z_i^2$ and its conjugate \hat{r} have been introduced. For $N \rightarrow \infty$, \hat{r} is found from the saddlepoint of

$$\ln \det(\mathbf{\Lambda} - \mathbf{J}) = \sum_i \ln(\Lambda_i - \hat{r}) + N\hat{r}r - 2NG(r) \quad (46)$$

yielding

$$r = \frac{1}{N} \sum_i \frac{1}{\Lambda_i - \hat{r}} \quad (47)$$

Hence, the averaged TAP equation (34) reads

$$[\chi_{ii}]_{\mathbf{J}} = \frac{1}{\Lambda_i - \hat{r}} \quad (48)$$

giving

$$r = \bar{\chi} \equiv \frac{1}{N} \sum_i [\chi_{ii}]_{\mathbf{J}} \quad (49)$$

and $V_i = V = \hat{r}$ (when taken together with the definition $\Lambda_i = V_i + 1/\chi_{ii}$). Finally, variation with respect to r yields $\hat{r} = 2G'(r)$. Summarizing, we find that in the thermodynamic limit,

$$V = 2G'(\bar{\chi}) . \quad (50)$$

As sketched in appendix C, the same result is obtained when the general Gaussian model used in the derivation of the Gibbs free energy is replaced by a *spherical* model, where only a single Lagrange parameter λ (or Λ) is coupled to $\sum_i S_i^2$.

Inserting the saddle point values into eq. (46), the expression for the Onsager term eq. (32) simplifies remarkably

$$\Delta\Phi^{\text{self}} = NG(\bar{\chi}) . \quad (51)$$

Turning to the generalized model, it can be seen from eqs. (38) and (40) that \mathbf{J} in the original model is replaced by $\mathbf{J} + \hat{\mathbf{J}}^T \hat{\mathbf{\Lambda}}^{-1} \hat{\mathbf{J}}$ in the generalized model. For this case, we extend the definition of orthogonal ensembles (42) to

$$\left[e^{\frac{1}{2} \text{Tr} \mathbf{A}(\mathbf{J} + \hat{\mathbf{J}}^T \hat{\mathbf{\Lambda}}^{-1} \hat{\mathbf{J}})} \right]_{\mathbf{J}, \hat{\mathbf{J}}} = e^{N \text{Tr} G_{\hat{\mathbf{\Lambda}}}(\mathbf{A}/N)} . \quad (52)$$

With this definition, the result eq. (50) for V remains valid, i.e. $V = 2G'_{\hat{\mathbf{\Lambda}}}(\bar{\chi})$. Selfaveraging results for \hat{V}_{kk} obtained from eq. (38) and the Onsager term eq. (40) are derived for neural network models in section 4.

3.2 Replica Free Energy

To get the average free energy, we use (41) and compute⁴

$$[Z^n]_{\mathbf{J}} = \int d\mathbf{S}^n \prod_{ia} \rho(S_{ia}) \left[e^{\frac{1}{2} \sum_a \sum_{ij} S_{ia} J_{ij} S_{ja}} \right]_{\mathbf{J}} = \int d\mathbf{S}^n \prod_{ia} \rho(S_{ia}) e^{N \text{Tr} G(\mathbf{q})} . \quad (53)$$

In replica symmetry, \mathbf{q} has only two types of eigenvalues: a nondegenerate one given by $(n-1)q + q_0$ and an $n-1$ fold degenerate eigenvalue equal to $q_0 - q$. Hence

$$\text{Tr} G(\mathbf{q}) = (n-1)G(q_0 - q) + G(nq + (q_0 - q)) \quad (54)$$

After introducing and eliminating conjugate parameters \hat{q} and \hat{q}_0 with a saddle-point method we find the replica symmetric free energy

$$\begin{aligned} -\frac{1}{N} [\ln Z]_{\mathbf{J}} &= -\frac{1}{N} \frac{d}{dn} \ln [Z^n]_{\mathbf{J}} \Big|_{n=0} \\ &= -\int Dz \ln \int dS \rho(S) \exp \left[\sqrt{2qG''(q)} zS + G'(q_0 - q) S^2 \right] \\ &\quad -G(q_0 - q) + (q_0 - q) (qG''(q_0 - q) + G'(q_0 - q)) , \end{aligned} \quad (55)$$

where $Dz \equiv dz e^{-z^2/2} / \sqrt{2\pi}$ and we must insert the values for q_0 and q which make the right hand side stationary.

3.3 Averaging the TAP Free Energy

In order to show the consistency of our TAP approach with the replica theory, we will show that the replica result for the average free energy eq. (55) is also obtained by averaging over the TAP solutions. We will only give a brief description. Introducing the auxiliary partition function

$$Y = \int d\gamma d\mathbf{m} e^{-\beta \mathcal{L}} , \quad (56)$$

with the Lagrangean defined by

$$\begin{aligned} \mathcal{L} &= \sum_i \ln Z_0(\gamma_i, \lambda) + \sum_i m_i \gamma_i + \frac{\lambda}{2} q_0 N \\ &\quad - \frac{1}{2} \sum_{ij} m_i J_{ij} m_j - NG \left(q_0 - \frac{1}{N} \sum_i m_i^2 \right) \end{aligned}$$

where $m_j \equiv \langle S_j \rangle$ and $q_0 = \frac{1}{N} \sum_i \langle S_i^2 \rangle$. In the limits $N \rightarrow \infty$ and $\beta \rightarrow \infty$, choosing the paths of integration such that the integral converges and varying

⁴This analysis can be easily generalized to include cases with more than one orderparameter, e.g. for a neural network learning from a teacher.

with respect to λ and q_0 , the integral will be dominated by the values of m_i and γ_i which correspond to the solution of the TAP equations. We thus recover the TAP free energy *at equilibrium* from

$$\Phi = - \text{extr}_{\lambda, q_0} \lim_{N, \beta \rightarrow \infty} \frac{1}{\beta} \ln Y. \quad (57)$$

Variation with respect to q_0 yields $\lambda = 2G' \left(q_0 - \frac{1}{N} \sum_i m_i^2 \right)$. The calculation of $[\Phi]_{\mathbf{J}}$ proceeds by a straightforward replica calculation using the average

$$\frac{1}{N} \ln \left[e^{\sum_a \sum_{ij} \beta m_{ia} J_{ij} m_{ja}} \right]_{\mathbf{J}} = (n-1)G(\beta(q - \bar{q})) + G(\beta(n\bar{q} + q - \bar{q})). \quad (58)$$

For $\beta \rightarrow \infty$, the integrations over γ_i and m_i are decoupled and performed by the saddlepoint method yielding $m_i = 0$ and $\gamma_i = \sqrt{qG''(q - q_0)}z_i$ where z_i is a standard normal random variable showing the equivalence to eq. (55). In comparing both replica calculations it is useful to note that we can identify

$$\lim_{\beta \rightarrow \infty} \beta(q - \bar{q}) = \bar{\chi} = q - q_0, \quad (59)$$

which can be shown by a linear response argument. Putting everything together we find that $[\Phi]_{\mathbf{J}} = -[\ln Z]_{\mathbf{J}}$.

3.4 Stability and AT Condition

We will show next (by generalizing the arguments of [3]) how the positive definiteness of the susceptibility matrix (i.e. the matrix of covariances) χ eq. (24) (or eq. (38) for the model eq. (2)) translates into the *de Almeida–Thouless* stability condition well known from the replica theory.

By eq. (24), positive definiteness of χ is equivalent to the condition that $\mathbf{H} = \mathbf{\Lambda} - \mathbf{J}$ has only positive eigenvalues. Hence, in the thermodynamic limit, the eigenvalue density $\rho(\gamma) \equiv \lim_{N \rightarrow \infty} \frac{1}{N} \sum_{\mu} \delta(\mu - \gamma)$, where μ denote the eigenvalues of \mathbf{H} , must be exactly zero for small positive γ . Using a standard representation of δ -functions, we have

$$\begin{aligned} \rho(\gamma) &= \frac{1}{N\pi} \lim_{\delta \rightarrow 0^+} \text{Im} \sum_{\mu} \frac{1}{\mu - \gamma - i\delta} \\ &= \frac{1}{N\pi} \lim_{\delta \rightarrow 0^+} \text{Im} \text{Tr} [\mathbf{H} - (\gamma + i\delta)\mathbf{I}]^{-1} \\ &= -\frac{1}{N\pi} \lim_{\delta \rightarrow 0^+} \text{Im} \frac{\partial}{\partial \gamma} \ln \det [\mathbf{\Lambda} - \mathbf{J} - (\gamma + i\delta)\mathbf{I}]. \end{aligned}$$

Since we have already calculated $\ln \det(\mathbf{\Lambda} - \mathbf{J})$ in section 3.1, we can immediately write down the result as

$$\rho(\gamma) = \lim_{\delta \rightarrow 0^+} \text{Im} \frac{1}{\pi N} \sum_i \frac{1}{\Lambda_i - \hat{r}(\gamma) - (\gamma + i\delta)}. \quad (60)$$

In general it is hard to obtain a close form solution to the saddlepoint equation for $\hat{r}(\gamma)$. However for γ close to zero—the interesting region with regard to the stability— $\hat{r}(\gamma)$ is close to $\hat{r}(0) = V$ and we can easily get a solution for \hat{r} and r by expanding to second order in $\delta\hat{r} \equiv \hat{r}(\gamma) - \hat{r}(0)$ and $\delta r = r(\gamma) - r(0)$, $r(0) = \bar{\chi}$. We find that as long as the condition

$$1 - 2G''(\bar{\chi})\frac{1}{N}\sum_i[\chi_{ii}]_{\mathbf{J}}^2 > 0 \quad (61)$$

is satisfied, $\hat{r}(\gamma)$ is *real* and the density $\rho(\gamma)$ eq. (60) vanishes for small $\gamma > 0$.

On the other hand, if the left hand side of eq. (61)—the stability condition—is zero, then $\delta\hat{r}$ has an imaginary part

$$\delta\hat{r} = \sqrt{\frac{-\gamma}{2G''(\bar{\chi})\frac{1}{N}\sum_i[\chi_{ii}]_{\mathbf{J}}^3 + \frac{G'''(\bar{\chi})}{[2G''(\bar{\chi})]^2}}} \quad (62)$$

and the support of the density of eigenvalues eq. (60)

$$\begin{aligned} \rho(\gamma) &\simeq \frac{1}{\pi N} \sum_i [\chi_{ii}]_{\mathbf{J}}^2 \text{Im } \delta\hat{r} \\ &= \frac{1}{\pi} \sqrt{\frac{\gamma}{[2G''(\bar{\chi})]^3 \frac{1}{N} \sum_i [\chi_{ii}]_{\mathbf{J}}^3 + G'''(\bar{\chi})}}. \end{aligned} \quad (63)$$

extends to $\gamma = 0$ and the solutions of TAP equations are only marginally stable. As we will show for some examples in the next section, eq. (61) coincides with the AT stability condition of replica theory [21].

4 Applications

In this section we will give explicit examples for adaptive TAP equations for a few models that have been previously considered in the literature. These are the SK- and Hopfield-models and the perceptron. The latter is of the generalized model type eq. (2). Finally, in simulations, we investigate the effect of the choice the Onsager-term for some perceptron learning problems.

4.1 Ising Models

For Ising models we have a prior distribution $\rho(S) = \frac{1}{2}\delta(S - 1) + \frac{1}{2}\delta(S + 1)$ so that

$$Z_0^{(i)} = \cosh \left(\sum_{j, j \neq i} J_{ij} \langle S_j \rangle - V_i \langle S_i \rangle + \theta_i \right)$$

which leads to

$$\langle S_i \rangle = \tanh \left(\sum_{j:j \neq i} J_{ij} \langle S_j \rangle - V_i \langle S_i \rangle + \theta_i \right) .$$

and $\chi_{ii} = 1 - \langle S_i \rangle^2$.

4.1.1 SK-model

For the SK-model [37] the statistics of the couplings is given by $\overline{J_{ij}} = 0$, $\overline{J_{ij}^2} = \beta/N$ with $J_{ij} = J_{ji}$. Then the G -function eq. (42) becomes $G(r) = \frac{(\beta r)^2}{4}$ and according to eq. (50)

$$V = 2G'(1 - q) = \beta^2(1 - q) , \quad (64)$$

where $q = \frac{1}{N} \sum_i \langle S_i \rangle^2$ is the *Edwards Anderson* parameter. The stability condition simplifies to $1 - \frac{\beta^2}{N} \sum_i \chi_{ii}^2 > 0$ and $\rho(\gamma) = \frac{1}{\pi} \sqrt{\frac{\gamma}{\frac{\beta^6}{N} \sum_i \chi_{ii}^3}}$ in agreement with [3].

4.1.2 Hopfield Model

The coupling matrix of the Hopfield model [9] is $J_{ij} = \frac{\beta}{N} \sum_k^m x_{ki} x_{kj}$, where we assume that the x_{ki} are iid random variables of zero mean and unit variance. The G -function⁵ eq. (42) is then found to be

$$G(r) = -\frac{m}{2N} [\ln(1 - \beta r) + \beta r] \quad (65)$$

leading to

$$V = 2G'(1 - q) = \frac{m}{N} \frac{\beta^2(1 - q)}{1 - \beta(1 - q)} \quad (66)$$

in agreement with [21].

4.2 Perceptron

Perceptrons are single layer neural networks which are parametrized by a vector of weights \mathbf{S} . We consider both the learning of regression and binary classification problems from a training set that is given by $\{(\mathbf{x}_k, y_k), k = 1, \dots, m\}$. $\mathbf{x} \in R^N$ denotes a vector of inputs and $y \in \mathcal{R}$ is a real valued output for regression and a binary label $y = \pm 1$ for classification. In the first case the output of the perceptron is given by $\mathbf{S} \cdot \mathbf{x}$ and in the latter case by $\text{sign}(\mathbf{S} \cdot \mathbf{x})$.

Perceptrons can be understood as probabilistic models by defining a probability for the observations y given inputs \mathbf{x} and weights \mathbf{S} which is defined by the likelihood $P(y|\mathbf{S} \cdot \mathbf{x})$. For classification we consider the so-called *probit* model which can be derived by assuming that labels are generated as $y = \text{sign}(\mathbf{S} \cdot \mathbf{x} + u)$

⁵This is easily shown for x_{ki} Gaussian. For binary x_{ki} , the average must be restricted to the condensed patterns and the relation eq. (42) will hold only for $N \rightarrow \infty$

where u is a Gaussian noise of variance σ^2 . Hence, $P(y|\mathbf{S} \cdot \mathbf{x}) = \phi(y \frac{\mathbf{S} \cdot \mathbf{x}}{\sigma})$, where $\phi(z) \equiv \int_{-\infty}^z Dt$. In the noise-free limit ϕ reduces to the unit stepfunction. For regression with additive Gaussian noise the likelihood is $P(y|\mathbf{S} \cdot \mathbf{x}) \propto e^{-(y - \mathbf{S} \cdot \mathbf{x})^2 / 2\sigma^2}$.

The model is clearly of the form given by eq. (2) with $\hat{J}_{ki} = x_{ki}$ and $\mathbf{J} = 0$. We identify the likelihood $P(y|\hat{h})$ with $F(\hat{h})$ in eq. (3) where $\hat{h} = \mathbf{S} \cdot \mathbf{x}$. The explicit appearance of the hatted variables in the algorithm will be especially useful when we want to discuss the important effects of removing a *datapoint* (rather than a weight S_j) from the set of training examples.

Under the assumptions of a Gaussian cavity field the Bayes predictors for regression and classification become $\langle \hat{h} \rangle = \langle \mathbf{S} \rangle \cdot \mathbf{x}$ and $\text{sign}(\hat{h})$ [28]. For the weight variables S_j we consider both Ising weights $\rho(S) = \frac{1}{2}\delta(S-1) + \frac{1}{2}\delta(S+1)$ and weights with a Gaussian prior distribution $\rho(S) = e^{-S^2/2}/\sqrt{2\pi}$. In the first case, we recover the Ising result

$$\langle S_i \rangle = \tanh \left(\sum_k x_{ki} \langle \hat{S}_k \rangle - V_i \langle S_i \rangle + \theta_i \right) \quad (67)$$

and in the Gaussian case we simply get

$$\langle S_i \rangle = \sum_k x_{ki} \langle \hat{S}_k \rangle + \frac{\theta_i}{1 - V_i}.$$

The TAP eqs. for the hatted variables $\langle \hat{S}_k \rangle = \frac{\partial \ln \hat{Z}_0^{(k)}}{\partial \hat{\theta}_k}$ are obtained from eq. (37)

$$\hat{Z}_0^{(k)} = \int Dz P(y_k | \langle \hat{h}_k \rangle_{\setminus k} + \hat{\theta}_k + \sqrt{\hat{V}_k} z)$$

with $\langle \hat{h}_k \rangle_{\setminus k} = \sum_i x_{ki} \langle S_i \rangle - \hat{V}_k \langle \hat{S}_k \rangle$. Explicit expressions are for classification:

$$\hat{Z}_0^{(k)} = \phi \left(y_k \frac{\langle \hat{h}_k \rangle_{\setminus k} + \hat{\theta}_k}{\sqrt{\sigma^2 + \hat{V}_k}} \right)$$

and for regression:

$$\hat{Z}_0^{(k)} = \frac{1}{\sqrt{2\pi(\sigma^2 + \hat{V}_k)}} \exp \left(-\frac{(y_k - \langle \hat{h}_k \rangle_{\setminus k} - \hat{\theta}_k)^2}{2(\sigma^2 + \hat{V}_k)} \right).$$

To connect with results known in the literature, we derive the selfaveraging properties for the case, where the x_{ki} are iid random variables with zero means and variance $1/N$. The G -function eq. (52) becomes

$$G_{\hat{\Lambda}}(r) = -\frac{1}{2N} \sum_k \ln(1 - r/\hat{\Lambda}_k).$$

The selfaveraging value for the variances of the original variables $V_i = V$ is given by eq. (50) and the variance for the hatted variables is given by eq. (39). Taken together they lead to the symmetric result in the two sets of variables

$$\hat{V} = \bar{\chi} = \frac{1}{N} \sum_i [\chi_{ii}]_{\mathbf{J}, \hat{\mathbf{J}}} \quad (68)$$

$$V = \bar{\hat{\chi}} \equiv \frac{1}{N} \sum_k [\hat{\chi}_{kk}]_{\mathbf{J}, \hat{\mathbf{J}}} . \quad (69)$$

The Onsager term eq. (40), stability condition eq. (61) and the the eigenvalue spectrum eq. (63) become

$$\Delta\Phi^{\text{self}} = -\frac{N}{2} V \hat{V} \quad (70)$$

$$1 - \frac{1}{N^2} \sum_k \hat{\chi}_{kk}^2 \sum_i \chi_{ii}^2 > 0 \quad (71)$$

$$\rho(\gamma) \simeq \frac{1}{\pi} \sqrt{\frac{\gamma}{\left[\frac{1}{N} \sum_k \hat{\chi}_{kk}^2\right]^3 \frac{1}{N} \sum_i \chi_{ii}^3 + \frac{1}{N} \sum_k \hat{\chi}_{kk}^3}} \quad \text{for } \gamma \rightarrow 0 . \quad (72)$$

Specializing to a Gaussian weight prior for which $\chi_{ii} = 1/(1 - V_i)$, we find in accordance with previous results [19, 44, 28] that $\hat{V} = \frac{1}{1-V}$ and the stability condition reduces to $1 - \hat{V}^2 \frac{1}{N} \sum_k \hat{\chi}_{kk}^2 > 0$.

Finally, we test the adaptive and selfaveraging TAP equations in two learning scenarios. We first test the internal consistency of the theory by comparing the cavity field calculated from the solution of the TAP equations $\langle \hat{h}_k \rangle_{\setminus k} = \sum_i x_{ki} \langle S_i \rangle - \hat{V}_k \hat{S}_k$ with the ‘exact’ cavity field $\langle \hat{h}_k \rangle_{\setminus k}^{\text{exact}}$ computed by actually *removing* example k from the training set and solving the TAP equations for the remaining $m - 1$ examples and repeating this procedure for $k = 1, \dots, m$. A precise estimate of the cavity field is of practical relevance in machine learning since it can be used to define ‘leave-one-out’ estimators of the generalization error [28, 29, 30]. One such ‘leave-one-out’ estimate for classification is the fraction of negative terms $y_k \langle \hat{h}_k \rangle_{\setminus k}$ over the training set: $\epsilon_{\text{loo}} = \frac{1}{m} \sum_k \Theta(-y_k \langle \hat{h}_k \rangle_{\setminus k})$ since $\text{sign}(\langle \hat{h}_k \rangle_{\setminus k})$ is the leave-one-out prediction of y_k .

Figures 1 and 2 show the result of learning in the simple perceptron with Gaussian weight prior, the likelihood for noise-free classification and the iid distribution of inputs for which the selfaveraging theory is expected to become exact in the thermodynamic limit. The output labels are generated by a neural network teacher $y = \text{sign}(\mathbf{T} \cdot \mathbf{x})$. While for positive values of $y_k \langle \hat{h}_k \rangle_{\setminus k}$ (not shown in the figures), i.e. the examples for which the leave-one-out prediction is correct, the agreement between $\langle \hat{h}_k \rangle_{\setminus k}$ and $\langle \hat{h}_k \rangle_{\setminus k}^{\text{exact}}$ tends to be better, the negative values are more crucial for real applications because they give the desired error count. The results clearly show that the internal consistency of the adaptive TAP

is better than that of the selfaveraging theory. The results indicate that finite size effects are quite important even for reasonably large systems and that the adaptive theory is better at taking these into account.

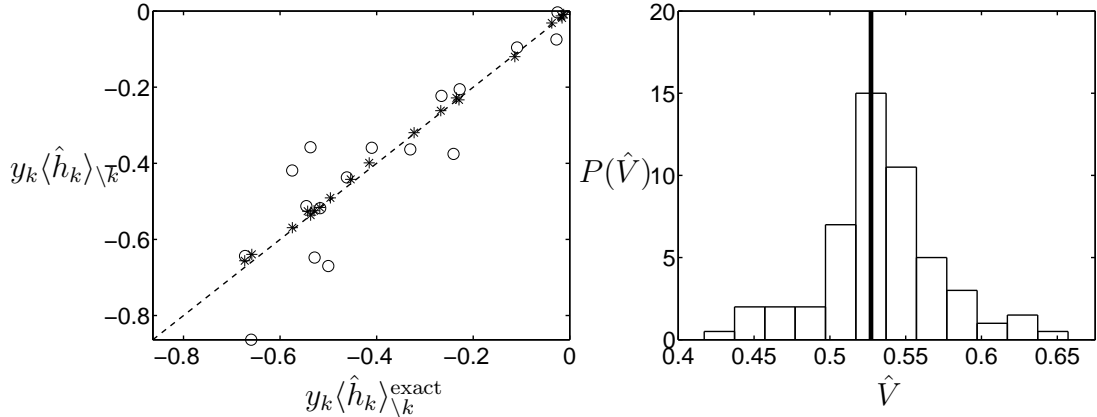


Figure 1: Test of self-consistency of TAP – $y_j \langle \hat{h}_j \rangle_{\setminus j}$ versus $y_j \langle \hat{h}_j \rangle_{\setminus j}^{\text{exact}}$. The stars/circles are for adaptive/conventional TAP. The right plot shows the distribution of the cavity variances \hat{V}_k . The line in the middle is the value found from the selfaveraging theory. The plot is for noise-free perceptron with $m = N = 100$.

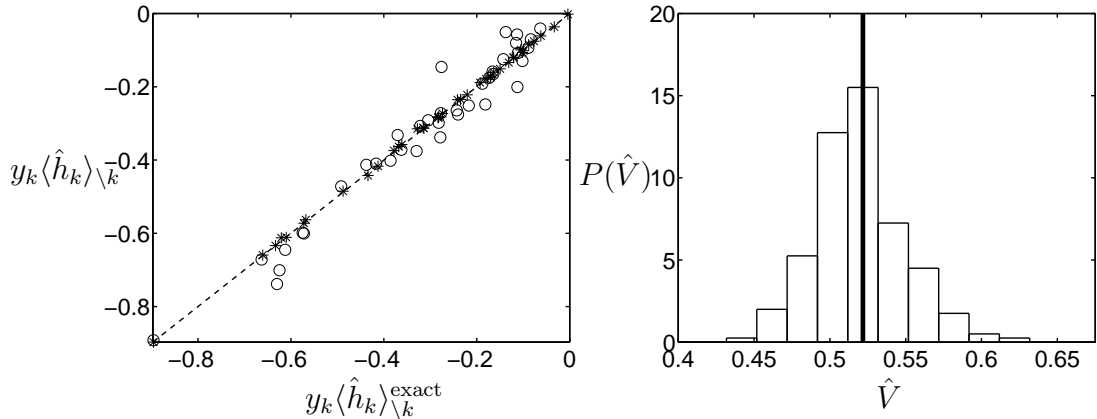


Figure 2: The same as above with $m = N = 200$.

Performing the same analysis for real data gives even more striking results. Here we consider the data set, ‘Sonar – Mines versus Rocks’ [7] of size $m = 104$ with binary class labels $y_k = \pm 1$ and a $N = 60$ dimensional input space. We use the Gaussian prior for the weights and $\sigma^2 = 0.5$ in the likelihood. In figure 3 we again plot $y_k \langle \hat{h}_k \rangle_{\setminus k}$ versus $y_k \langle \hat{h}_k \rangle_{\setminus k}^{\text{exact}}$. For the adaptive theory, we find a perfect agreement between the two computations of the leave-one-out estimate: $\epsilon_{\text{loo}} = \epsilon_{\text{loo}}^{\text{exact}} = 33/104$. For comparison, the selfaveraging TAP approach gives $\epsilon_{\text{loo}} = 41/104$ and $\epsilon_{\text{loo}}^{\text{exact}} = 33/104$.

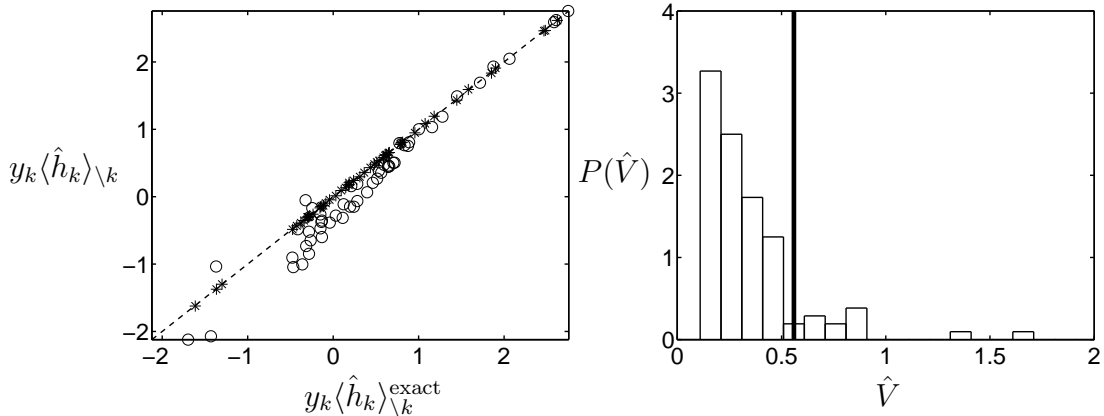


Figure 3: The same as above for the Sonar data set.

In the second set of simulations we have tested the influence of using a wrong cavity variance V in the mean field equations (and wrong Onsager term $\Delta\Phi$ in the free energy). Since the free energy is the negative log Likelihood of the observed data, i.e. *at equilibrium*, $\Phi = -\ln P(\mathbf{y})$, it can be used for deciding which model gives the best fit to data. It is therefore also of practical interest to get a reliable estimate of Φ .

In these simulations we consider a nontrivial regression model with binary weights. See Ref. [40] for a discussion of this model in the context of demodulation in communications systems. As it can be seen directly from the Likelihood, the regression problem can alternatively be regarded as an N -dimensional model of the type eq. (1), i.e. $\prod_k P(y_k | \mathbf{S} \cdot \mathbf{x}_k) \propto e^{\sum_{i>j} S_i J_{ij} S_j + \sum_i S_i \theta_i}$. The couplings and external fields are given by $J_{ij} = -\sum_k x_{ki} x_{kj} / \sigma^2$ and $\theta_i = \sum_k x_{ki} y_k / \sigma^2$. We can now directly compare the use of the correct selfaveraging V eq. (69) for this model with the one provided from the adaptive TAP approach and that of other Ising models with different matrix ensembles, namely the SK- and Hopfield-model eqs. (64) and (66).

In figure 4, we compare the TAP mean field free energy Φ found in simulations using the different expressions for the Onsager term with the prediction of replica theory [36]. In the simulations $N = 60$, $\sigma^2 = 0.2$ and the training set is generated by a noise-free binary teacher: $y = \mathbf{T} \cdot \mathbf{x}$ with $T_i = \pm 1$. The simulation are averaged over 100 runs and the errorbars are of the size of the symbols. We set $\beta = 1$ for the SK-model and $\beta = 0.99$ for the Hopfield-model. The latter to avoid numerical problems when $q \approx 1$. As it can be seen from the figure both adaptive and selfaveraging TAP using the Onsager term eq. (70) are in excellent agreement with replica theory. Using the SK and Hopfield Onsager terms tend to produce saturated solutions, i.e. $\langle S_i \rangle \approx \pm 1$ even for training set sizes where this is not expected theoretically and lead to a completely wrong estimate of the free energy.

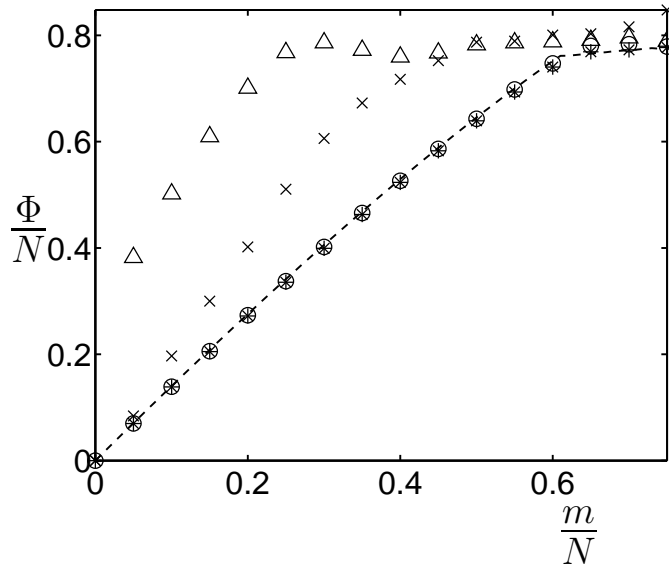


Figure 4: The free energy Φ as a function of the training set size m . The dashed line is the prediction of replica theory. Stars/circles (almost coinciding) are the results for adaptive TAP/correct selfaveraging TAP. Crosses/triangles are the results for selfaveraging TAP with the Hopfield/SK Onsager term.

5 Summary and Outlook

We have presented a generalization of the TAP approach for disordered systems which is able to cope with the lack of knowledge of the disorder distribution. Such a generalization is necessary for the recent applications of mean field methods to probabilistic data models.

We have demonstrated the significance of our approach in two ways: We have shown that it reproduces the correct thermodynamic limit results for a class of disorder distributions compatible with fully connected models in replica symmetry. Second, the application of our approach to toy models as well as real data models, has shown the importance of using the correct TAP approximation when good approximations for *leave-one-out estimators* of errors and *free energies* are required. Such quantities serve as a practical yardsticks for comparing different data models and assessing the validity of model predictions.

While the present framework may be sufficient for a variety of practical applications, it may nevertheless break down when the probabilistic model is multimodal having many modes of almost equal weights (free energies). For systems in such ‘glassy’ phases our TAP solutions are expected to violate the stability condition and an extension of our framework to a theory which includes Parisi’s hierarchical organization of states would be highly important. We expect that this is possible by generalizing the ideas presented in chapter V of [21]. However, one may speculate that in such cases solving the TAP equations may be highly

nontrivial.

We conclude with two other problem areas which have high priority for our future research. These are the limitations of our method to models with extensive connectivities and the algorithmic aspects of our approach, i.e. the development of efficient algorithms for solving the TAP equations. Recent studies on other advanced mean field techniques indicate that both problems have interesting relationships and also give promising directions for their solutions.

The so-called belief propagation algorithm [32], used in the field of artificial intelligence for approximate probabilistic computations on graphical models with sparse connectivities was recently identified [14, 42, 43] as an efficient method to solve the Bethe approximation (a cavity type of approximation) of statistical physics. This observation has already led to principled ways of combining the improved accuracy of higher order (Kikuchi) Bethe approximations [43] with the efficiency of the belief propagation method.

Another interesting approach to an approximate propagation of probability distributions when data arrive sequentially is the Bayesian on-line method introduced in [25, 26] and further developed in [38, 4, 5]. This technique can be formulated for fairly general model classes but was sofar limited to a single sweep through the data, thereby making the approximation dependent on the ordering of the data sequence. In a recent study by Minka [22] it was shown that by a proper recycling of the data, a convergence to the solutions of the TAP equations for the case of a Gaussian process classifier [29] was achieved. We expect that by a consequent and principled combination of the cavity idea with algorithms that are similar to the on-line or the belief propagation technique, we will not only get efficient methods for the solution of TAP equations but also be able to significantly extend the range of applications of the TAP approach.

Acknowledgments

This research is supported by the Swedish Foundation for Strategic Research as well as the Danish Research Councils through the THOR Center for Neuroinformatics and Center for Biological Sequence Analysis.

A Correctness of TAP equations for the Gaussian model

We will show that the relation eq. (24)

$$\chi = (\mathbf{A} - \mathbf{J})^{-1} \tag{73}$$

is correct for a Gaussian model. Gaussian models are defined by $\rho(S_i) \propto e^{-\frac{1}{2}A_i x_i^2}$ and we have always $\chi = (\mathbf{A} - \mathbf{J})^{-1}$ where $\mathbf{A} = \text{diag}(A_1, \dots, A_N)$. Hence, we

only have to show that $A_j = \Lambda_j = V_j + 1/\chi_{jj}$. To see this we look at the single variable partition function eq. (18) derived from the Gaussian cavity field assumption which is exact for the Gaussian model

$$Z_0^{(i)} = \int dS \rho_i(S) e^{S(\langle h_i \rangle_{\lambda_i + \theta_i} + \frac{1}{2} V_i S^2)} \quad (74)$$

which gives in fact $\chi_{ii} = \frac{\partial^2 \ln Z_0^{(i)}}{\partial \theta_i^2} = \frac{1}{A_i - V_i}$.

B Adaptive TAP Free Energy II

Parisi and Potters [31] in their analysis of a spin glass model with random orthogonal couplings made the important observation (based on a high temperature expansion) that (within the TAP approximation) two models having the same interactions $\sum_{i < j} S_i J_{ij} S_j$ but differ only in their single spin constraints $\rho_i(S)$, should have free energies Φ which differ also only in the ‘single variable’ (or entropic) contribution Φ_0 .

Hence, it is possible to compute the TAP approximation for the free energy Φ from the free energy for an exactly solvable model Φ^s , the entropic term for the solvable model Φ_0^s and the single variable term for our model Φ_0 , i.e.,

$$\Phi = \Phi^s - \Phi_0^s + \Phi_0 . \quad (75)$$

Since the TAP equations for a Gaussian model are exact (see appendix A) we choose $\rho(S) = e^{-S^2/2}/\sqrt{2\pi}$ for the solvable model. We easily get the following exact result for its Gibbs free energy (after eliminating the Lagrange multipliers γ_i and with hindsight redefining $\Lambda_i \equiv 1 - \lambda_i$)

$$\Phi^s(\mathbf{m}, \mathbf{M}) = \frac{1}{2} \ln \det(\mathbf{\Lambda} - \mathbf{J}) + \frac{1}{2} \sum_{ij} m_i J_{ij} m_j - \frac{1}{2} \sum_i \chi_{ii} \Lambda_i + \sum_i \frac{M_i}{2} , \quad (76)$$

where we have to insert the value for Λ_i which solves

$$\chi_{ii} = \langle S_i^2 \rangle - \langle S_i \rangle^2 = [(\mathbf{\Lambda} - \mathbf{J})^{-1}]_{ii} . \quad (77)$$

The single variable term for the Gaussian model Φ_0^s is found by setting $J_{ij} = 0$ in eq. (76). Eliminating $\mathbf{\Lambda}$ using $\partial_{\mathbf{\Lambda}} \Phi_0^s = 0$:

$$\Phi_0^s(\mathbf{m}, \mathbf{M}) = -\frac{1}{2} \sum_i \ln \chi_{ii} - \frac{N}{2} + \sum_i \frac{M_i}{2} . \quad (78)$$

We can now write down the general result for the TAP mean field Gibbs free energy for a model of the type eq. (1). Collecting the term eqs. (76) and (78),

we arrive at the free energy

$$\Phi = \Phi^s - \Phi_0^s + \Phi_0 = \Phi_0 - \frac{1}{2} \sum_{ij} m_i J_{ij} m_j + \Delta\Phi \quad (79)$$

$$\Delta\Phi = \frac{1}{2} \ln \det(\Lambda - \mathbf{J}) - \frac{1}{2} \sum_i \Lambda_i \chi_{ii} + \frac{1}{2} \sum_i \ln \chi_{ii} + \frac{N}{2}. \quad (80)$$

This result should be compared to eq. (31). Using the saddlepoint condition $\partial_{M_i} \Phi = 0$ which implies $\Lambda_i = \frac{1}{\chi_{ii}} + \lambda_i$, we can rewrite the Onsager term in the form of eq. (32) where it should be noted that $V_i = \lambda_i$. We have thus rederived the result obtained in section 2.4.

C Free Energy for the spherical model

For the spherical model defined by the constraint $\sum_i S_i^2 = N$ we obtain

$$\Phi = \Phi_0 - \frac{1}{2} \sum_{ij} \langle S_i \rangle J_{ij} \langle S_j \rangle + \frac{1}{2} \ln \det(\Lambda \mathbf{I} - \mathbf{J}) - \frac{N\Lambda\chi}{2} + \frac{N}{2} \ln \chi - \frac{N}{2}, \quad (81)$$

where $\chi \equiv \frac{1}{N} \sum_i \chi_{ii}$ and there is also only a single λ^0 in Φ_0 . Λ is determined by

$$\chi = \frac{1}{N} \text{Tr}(\Lambda \mathbf{I} - \mathbf{J})^{-1}. \quad (82)$$

Secondly, we have $\Lambda = 1/\chi + V$. Repeating the same averaging step as before yields again eq. (50).

D Eigenvalue spectrum of \mathbf{J}

In this appendix we will show how the G -function eq. (42) can be expressed in terms of the eigenvalue spectrum of the matrix \mathbf{J} . We define

$$r \equiv \int d\mu \frac{p(\mu)}{\Lambda - \mu} = \frac{1}{N} \text{Tr}(\Lambda \mathbf{I} - \mathbf{J})^{-1} \quad (83)$$

where $p(\mu)$ is the density of eigenvalues of \mathbf{J} . By adding a small imaginary part to Λ we get apart from a factor, directly the density. Using again the Gaussian representation of the determinant yields the equations $r = \frac{1}{\Lambda - \hat{r}}$, where \hat{r} is the orderparameter conjugate to $\sum_i z_i^2$. It obeys $\hat{r} = 2G'(r)$ and $r = \frac{1}{\Lambda - \hat{r}}$. Hence

$$\frac{1}{2r} - \frac{\Lambda}{2} + G'(r) = 0 \quad (84)$$

Solving eq. (84) enables us to compute $r(\Lambda)$ when the function $G(r)$ is given. We may also get G as a function of Λ by integrating eq. (84).

References

- [1] S. Amari, S. Ikeda H. Shimokawa in [27].
- [2] C. Bhattacharyya and S. S. Keerthi, *J. Phys. A (Math. Gen.)* **10**, 1307 (2000).
- [3] A. J. Bray and M. A. Moore, *J. Phys. C* **13**, L469 (1980).
- [4] L. Csató, E. Fokoué, M. Opper, B. Schottky, and O. Winther, in *Advances in Neural Information Processing Systems 12*, S. A. Solla, T. K. Leen, and K.-R. Müller (Eds.), 251–257. The MIT Press (2000).
- [5] L. Csató, and M. Opper, in *Advances in Neural Information Processing Systems 13*, T. K. Leen, T. G. Diettrich, and V. Tresp (Eds.), The MIT Press (2001).
- [6] C. De Dominicis, *Phys. Rep.* **67**, 37 (1980).
- [7] R. P. Gorman and T. J. Sejnowski, *Neural Networks* **1**, 75 (1988).
- [8] G. E. Hinton and T. J. Sejnowski, in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (Washington, 1983)*, 448, New York, IEEE Press (1983).
- [9] J. J. Hopfield, *Proc. Nat. Acad. Sci. USA*, **79** 2554 (1982).
- [10] M. Jordan, ed., *Learning in Graphical Models*, MIT Press (1999).
- [11] H. J. Kappen and F. B. Rodríguez, *Neural Computation* **10**, 1137 (1998).
- [12] H. J. Kappen and F. B. Rodríguez, in *Advances in Neural Information Processing Systems 11*, eds. M. S. Kearns, S. A. Solla and D. A. Cohn, 280 MIT Press (1999).
- [13] H. J. Kappen and W. Wiegerinck, in [27].
- [14] Y. Kabashima and D. Saad, *Euro. Phys. Lett.* **44**, 668 (1998).
- [15] Y. Kabashima and D. Saad, in [27].
- [16] T.-W. Lee: *Independent Component Analysis*, Kluwer Academic Publishers, Boston (1998).
- [17] D. J. C. MacKay, *Maximum Likelihood and Covariant Algorithms for Independent Component Analysis*, University of Cambridge, Cavendish Laboratory, Draft 3.7, (1996).

- [18] E. Marinari, G. Parisi and F. Ritort, *J. Phys. A (Math. Gen.)* **27**, 7647 (1994).
- [19] M. Mézard, *J. Phys. A (Math. Gen.)* **22**, 2181 (1989).
- [20] M. Mézard and G. Parisi, *J. Phys. Lett.* **46**, L771 (1995).
- [21] M. Mézard, G. Parisi and M. A. Virasoro, *Spin Glass Theory and Beyond*, Lecture Notes in Physics, 9, World Scientific (1987).
- [22] T. P. Minka, Expectation Propagation for Approximate Bayesian Inference, Preprint MIT Media Lab (2000).
- [23] R. Neal, *Artificial Intelligence* **56**, 71(1992).
- [24] M. Ohlsson, C. Peterson and B. Söderberg, *Neural Computation* **5**, 331 (1993).
- [25] M. Opper, *Phys. Rev. Lett.* **77**, 4671 (1996).
- [26] M. Opper, in *On-Line Learning in Neural Networks*, ed. D. Saad, Cambridge Univ. Press, 363(1998).
- [27] M. Opper and D. Saad, eds. *Advanced Mean Field Methods, Theory and Practice*, MIT Press (2001).
- [28] M. Opper and O. Winther, *Phys. Rev. Lett.* **76**, 1964 (1996).
- [29] M. Opper and O. Winther, *Neural Computation* **12**, 2655 (2000).
- [30] M. Opper and O. Winther, *Phys. Rev. Lett.* **86**, 3695 (2001).
- [31] G. Parisi and M. Potters, *J. Phys. A (Math. Gen.)* **28**, 5267 (1995).
- [32] J. Pearl, *Probabilistic Reasoning in Intelligent Systems: Networks of Plausible Inference*, Morgan Kaufmann, San Francisco (1988).
- [33] T. Plefka, *J. Phys. A* **15**, 1971 (1982).
- [34] W. H. Press, S. A. Teukolsky, W. T. Vetterling and B. P. Flannery, *Numerical Recipes in C*, 77, Cambridge University Press (1992).
- [35] L. K. Saul, T. Jaakkola, M. I. Jordan, *J. Artificial Intelligence Research* **4**, 61(1996).
- [36] H. S. Seung, H. Sompolinsky and N. Tishby, *Phys. Rev. A.* **45**, 6056 (1992).
- [37] D. Sherrington, and K. Kirkpatrick, *Phys. Rev. Lett.* **35**, 1792 (1975).

- [38] S. A. Solla and O. Winther, in *On-Line Learning in Neural Networks*, ed. D. Saad, Cambridge Univ. Press, 379(1998).
- [39] T. Tanaka, Phys. Rev. E **58**, 2302(1998).
- [40] T. Tanaka, in *Advances in Neural Information Processing Systems 13*, T. K. Leen, T. G. Diettrich, and V. Tresp (Eds.), The MIT Press (2001).
- [41] D. J. Thouless, P. W. Anderson and R. G. Palmer, Phil. Mag. **35**, 593 (1977).
- [42] J. Yedidia, in [27].
- [43] J. S. Yedidia, W. T. Freeman and Y. Weiss, in *Advances in Neural Information Processing Systems 13*, T. K. Leen, T. G. Diettrich, and V. Tresp (Eds.), The MIT Press (2001).
- [44] K. Y. M. Wong, Europhys. Lett. **30**, 245 (1995).