Radboud Repository



PDF hosted at the Radboud Repository of the Radboud University Nijmegen

The following full text is an author's version which may differ from the publisher's version.

For additional information about this publication click this link. http://hdl.handle.net/2066/36525

Please be advised that this information was generated on 2017-12-06 and may be subject to change.

Sufficient Conditions for Convergence of the Sum-Product Algorithm

Joris M. Mooij and Hilbert J. Kappen

Abstract-Novel conditions are derived that guarantee convergence of the Sum-Product algorithm (also known as Loopy Belief Propagation or simply Belief Propagation) to a unique fixed point, irrespective of the initial messages, for parallel (synchronous) updates. The computational complexity of the conditions is polynomial in the number of variables. In contrast with previously existing conditions, our results are directly applicable to arbitrary factor graphs (with discrete variables) and are shown to be valid also in the case of factors containing zeros, under some additional conditions. The conditions are compared with existing ones, numerically and, if possible, analytically. For binary variables with pairwise interactions, sufficient conditions are derived that take into account local evidence (i.e., singlevariable factors) and the type of pair interactions (attractive or repulsive). It is shown empirically that this bound outperforms existing bounds.

Index Terms—Contraction, convergence, factor graphs, graphical models, marginalization, message passing, Loopy Belief Propagation, Sum-Product Algorithm

I. INTRODUCTION

THE Sum-Product Algorithm [2], also known as Loopy Belief Propagation, which we will henceforth abbreviate as BP, is a popular algorithm for approximate inference on graphical models. Applications can be found in diverse areas such as error correcting codes (iterative channel decoding algorithms for Turbo Codes and Low Density Parity Check Codes [3]), combinatorial optimization (satisfiability problems such as 3-SAT and graph coloring [4]) and computer vision (stereo matching [5] and image restoration [6]). BP can be regarded as the most elementary one in a family of related algorithms, consisting of double-loop algorithms [7], GBP [8], EP [9], EC [10], the Max-Product Algorithm [11], the Survey Propagation Algorithm [4], [12] and Fractional BP [13]. A good understanding of BP may therefore be beneficial to understanding these other algorithms as well.

In practice, there are two major obstacles in the application of BP to concrete problems: (i) if BP converges, it is not clear whether the results are a good approximation of the

J. M. Mooij is with the Department of Biophysics, Radboud University Nijmegen, PO Box 9101, 6500 HB Nijmegen, The Netherlands (e-mail: j.mooij@science.ru.nl).

H. J. Kappen is with the Department of Biophysics, Radboud University Nijmegen, PO Box 9101, 6500 HB Nijmegen, The Netherlands (e-mail: b.kappen@science.ru.nl).

Parts of this work have been presented at the 21st Conference on Uncertainty in Artificial Intelligence (UAI 2005) and published in the conference proceedings [1].

The research reported here is part of the Interactive Collaborative Information Systems (ICIS) project (supported by the Dutch Ministry of Economic Affairs, grant BSIK03024) and was also sponsored in part by the Dutch Technology Foundation (STW).

exact marginals; (ii) BP does not always converge, and in these cases gives no approximations at all. These two issues might actually be interrelated: the "folklore" is that failure of BP to converge often indicates low quality of the Bethe approximation on which it is based. This would mean that if one has to "force" BP to converge (e.g., by using damping or double-loop approaches), one may expect the results to be of low quality.

Although BP is an old algorithm that has been reinvented in many fields, a thorough theoretical understanding of the two aforementioned issues and their relation is still lacking. Significant progress has been made in recent years regarding the question under what conditions BP converges [14]–[16]¹, on the uniqueness of fixed points [18], and on the accuracy of the marginals [15], but the theoretical understanding is still incomplete. For the special case of a graphical model consisting of a single loop, it has been shown that convergence rate and accuracy are indeed related [19].

In this work, we study the question of convergence of BP and derive new sufficient conditions for BP to converge to a unique fixed point. Our results are more general and in some cases stronger than previously known sufficient conditions.

II. BACKGROUND

To introduce our notation, we give a short treatment of factorizing probability distributions, the corresponding visualizations called factor graphs, and the BP algorithm on factor graphs. For an excellent, extensive treatment of these topics we refer the reader to [2].

A. Graphical Models

Consider N discrete random variables x_i for $i \in \mathcal{V} := \{1, 2, \ldots, N\}$, with x_i taking values in \mathcal{X}_i for x_i . We write $x = (x_1, \ldots, x_N) \in \mathcal{X} := \prod_{i \in \mathcal{V}} \mathcal{X}_i$. We are interested in the class of probability measures on \mathcal{X} that can be written as a product of factors (also called potentials):

$$P(x_1, \dots, x_N) := \frac{1}{Z} \prod_{I \in \mathcal{F}} \psi^I(x_I). \tag{1}$$

The factors ψ^I are indexed by subsets of \mathcal{V} , i.e., $\mathcal{F} \subseteq \mathcal{P}(\mathcal{V})$. If $I \in \mathcal{F}$ is the subset $I = \{i_1, \ldots, i_m\} \subseteq \mathcal{V}$, we write $x_I := (x_{i_1}, \ldots, x_{i_m}) \in \prod_{i \in I} \mathcal{X}_i$. Each factor ψ^I is a positive

¹After submission of this work, we came to the attention of [17], which contains improved versions of results in [16], some of which are similar or identical to results presented here (c.f. Section V-B).

function $\psi^I:\prod_{i\in I}\mathcal{X}_i\to(0,\infty)$. Z is a normalizing constant ensuring that $\sum_{x\in\mathcal{X}}P(x)=1$. The class of probability measures described by (1) contains Markov Random Fields as well as Bayesian Networks. We will use uppercase letters for indices of factors $(I, J, K, \ldots \in \mathcal{F})$ and lowercase letters for indices of variables $(i, j, k, \ldots \in \mathcal{V})$.

The factor graph that corresponds to the probability distribution (1) is a bipartite graph with vertex set $\mathcal{V} \cup \mathcal{F}$. In the factor graph (see also Fig. 1), each variable node $i \in \mathcal{V}$ is connected with all the factors $I \in \mathcal{F}$ that contain the variable, i.e., the neighbors of i are the factor nodes $N_i := \{I \in \mathcal{F} : i \in \mathcal{F} :$ I). Similarly, each factor node $I \in \mathcal{F}$ is connected with all the variable nodes $i \in \mathcal{V}$ that it contains and we will simply denote the neighbors of I by $I = \{i \in \mathcal{V} : i \in I\}$. For each variable node $i \in \mathcal{V}$, we define the set of its neighboring variable nodes by $\partial i := (\bigcup N_i) \setminus \{i\}$, i.e., ∂i is the set of indices of those variables that interact directly with x_i .

B. Loopy Belief Propagation

Loopy Belief Propagation is an algorithm that calculates approximations to the marginals $\{P(x_I)\}_{I\in\mathcal{F}}$ and $\{P(x_i)\}_{i\in\mathcal{V}}$ of the probability measure (1). The calculation is done by message-passing on the factor graph: each node passes messages to its neighbors. One usually discriminates between two types of messages: messages $\mu^{I \to i}(x_i)$ from factors to variables and messages $\mu^{i\to I}(x_i)$ from variables to factors (where $i \in I \in \mathcal{F}$). Both messages are positive functions on \mathcal{X}_i , or, equivalently, vectors in $\mathbb{R}^{\mathcal{X}_i}$ (with positive components). The messages that are sent by a node depend on the incoming messages; the new messages, designated by $\tilde{\mu}$, are given in terms of the incoming messages by the following BP update rules³

$$\tilde{\mu}^{j \to I}(x_j) \propto \prod_{J \in N_j \setminus I} \mu^{J \to j}(x_j)$$
(2)

$$\tilde{\mu}^{j \to I}(x_j) \propto \prod_{J \in N_j \setminus I} \mu^{J \to j}(x_j)$$

$$\tilde{\mu}^{I \to i}(x_i) \propto \sum_{x_{I \setminus i}} \psi^I(x_I) \prod_{j \in I \setminus i} \mu^{j \to I}(x_j).$$
(2)

Usually, one normalizes the messages in the ℓ_1 -sense (i.e., such that $\sum_{x_i \in \mathcal{X}_i} \mu(x_i) = 1$). If all messages have converged to some fixed point μ_{∞} , one calculates the approximate marginals or beliefs

$$b_i(x_i) = C^i \prod_{I \in N_i} \mu_{\infty}^{I \to i}(x_i) \approx P(x_i)$$
$$b_I(x_I) = C^I \psi^I(x_I) \prod_{i \in I} \mu_{\infty}^{i \to I}(x_i) \approx P(x_I),$$

where the C^{i} 's and C^{I} 's are normalization constants, chosen such that the approximate marginals are normalized in ℓ_1 sense. A fixed point always exists if all factors are strictly positive [8]. However, the existence of a fixed point does not necessarily imply convergence towards the fixed point, and fixed points may be unstable.

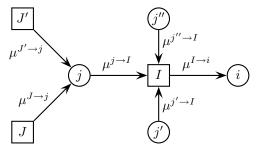


Fig. 1. Part of the factor graph illustrating the BP update rules (2) and (3). The factor nodes $I, J, J' \in \mathcal{F}$ are drawn as rectangles, the variable nodes $i, j, j', j'' \in \mathcal{V}$ as circles. Note that $N_i \setminus I = \{J, J'\}$ and $I \setminus i = \{j, j', j''\}$. Apart from the messages that have been drawn, each edge also carries a message flowing in the opposite direction.

Note that the beliefs are invariant under rescaling of the

$$\mu_{\infty}^{I \to i}(x_i) \mapsto \alpha^{I \to i} \mu_{\infty}^{I \to i}(x_i), \qquad \mu_{\infty}^{i \to I}(x_i) \mapsto \alpha^{i \to I} \mu_{\infty}^{i \to I}(x_i)$$

for positive constants α , which shows that the precise way of normalization in (2) and (3) is irrelevant. For numerical stability however, some way of normalization (not necessarily in ℓ_1 -sense) is desired to ensure that the messages stay in some compact domain.

In the following, we will formulate everything in terms of the messages $\mu^{I \to i}(x_i)$ from factors to variables; the update equations are then obtained by substituting (2) in (3):

$$\tilde{\mu}^{I \to i}(x_i) = C^{I \to i} \sum_{x_{I \setminus i}} \psi^I(x_I) \prod_{j \in I \setminus i} \prod_{J \in N_j \setminus I} \mu^{J \to j}(x_j). \quad (4)$$

with $C^{I \to i}$ such that $\sum_{x_i \in \mathcal{X}_i} \tilde{\mu}^{I \to i}(x_i) = 1$. We consider here BP with a *parallel* update scheme, which means that all message updates (4) are done in parallel.

III. SPECIAL CASE: BINARY VARIABLES WITH PAIRWISE INTERACTIONS

In this section we investigate the simple special case of binary variables (i.e., $|\mathcal{X}_i| = 2$ for all $i \in \mathcal{V}$), and in addition we assume that all potentials consist of at most two variables ("pairwise interactions"). Although this is a special case of the more general theory to be presented later on, we start with this simple case because it illustrates most of the underlying ideas without getting involved with the additional technicalities of the general case.

We will assume that all variables are ± 1 -valued, i.e., $\mathcal{X}_i =$ $\{-1,+1\}$ for all $i \in \mathcal{V}$. We take the factor index set as $\mathcal{F} = \mathcal{F}_1 \cup \mathcal{F}_2$ with $\mathcal{F}_1 = \mathcal{V}$ (the "local evidence") and $\mathcal{F}_2 \subseteq \{\{i,j\}: i,j \in \mathcal{V}, i \neq j\}$ (the "pair-potentials"). The probability measure (1) can then be written as

$$P(x) = \frac{1}{Z} \exp\left(\sum_{\{i,j\}\in\mathcal{F}_2} J_{ij} x_i x_j + \sum_{i\in\mathcal{F}_1} \theta_i x_i\right)$$
 (5)

for some choice of the parameters J_{ij} ("couplings") and θ_i ("local fields"), with $\psi^i(x_i) = \exp(\theta_i x_i)$ for $i \in \mathcal{F}_1$ and $\psi^{\{i,j\}}(x_i, x_j) = \exp(J_{ij}x_ix_j) \text{ for } \{i, j\} \in \mathcal{F}_2.$

²In subsection IV-E we will loosen this assumption and allow for factors containing zeros.

³We abuse notation slightly by writing $X \setminus x$ instead of $X \setminus \{x\}$ for sets

Note from (4) that the messages sent from single-variable factors \mathcal{F}_1 to variables are constant. Thus the question whether messages converge can be decided by studying only the messages sent from pair-potentials \mathcal{F}_2 to variables. It turns out to be advantageous to use the following "natural" parameterization of the messages

$$\tanh \nu^{i \to j} := \mu^{\{i,j\} \to j} (x_i = 1) - \mu^{\{i,j\} \to j} (x_i = -1), \quad (6)$$

where $\nu^{i \to j} \in \mathbb{R}$ is now interpreted as a message sent from variable i to variable j (instead of a message sent from the factor $\{i,j\}$ to variable j). Note that in the pairwise case, the product over $j \in I \setminus i$ in (4) becomes trivial. Some additional elementary algebraic manipulations show that the BP update equations (4) become particularly simple in this parameterization and can be written as:

$$\tanh \tilde{\nu}^{i \to j} = \tanh(J_{ij}) \tanh \left(\theta_i + \sum_{t \in \partial i \setminus j} \nu^{t \to i}\right)$$
 (7)

where $\partial i = \{t \in \mathcal{V} : \{i, t\} \in \mathcal{F}_2\}$ are the variables that interact with i via a pair-potential.

Defining the set of ordered pairs $D:=\{i\to j:\{i,j\}\in\mathcal{F}_2\}$, we see that the parallel BP update is a mapping $f:\mathbb{R}^D\to\mathbb{R}^D$; (7) specifies the component $(f(\nu))^{i\to j}:=\tilde{\nu}^{i\to j}$ in terms of the components of ν . Our goal is now to derive sufficient conditions under which the mapping f is a contraction. For this we need some elementary but powerful mathematical theorems.

A. Normed Spaces, Contractions And Bounds

In this subsection we introduce some (standard) notation and remind the reader of some elementary but important properties of vector norms, matrix norms, contractions and the Mean Value Theorem in arbitrary normed vector spaces, which are the main mathematical ingredients for our basic tool, Lemma 2. The reader familiar with these topics can skip this subsection and proceed directly to Lemma 2 in section III-B.

Let $(V, \|\cdot\|)$ be a normed finite-dimensional real vector space. Examples of norms that will be important later on are the ℓ_1 -norm on \mathbb{R}^N , defined by

$$||x||_1 := \sum_{i=1}^N |x_i|$$

and the ℓ_{∞} -norm on \mathbb{R}^N , defined by

$$||x||_{\infty} := \max_{i \in \{1, \dots, N\}} |x_i|.$$

A norm on a vector space V induces a metric on V by the definition $d(v,w):=\|v-w\|$. The resulting metric space is complete.⁴

Let (X,d) be a metric space. A mapping $f:X\to X$ is called a *contraction with respect to* d if there exists $0\le K<1$ such that

$$d(f(x), f(y)) \le Kd(x, y)$$
 for all $x, y \in X$. (8)

⁴Completeness is a topological property which we will not further discuss, but we need this to apply Theorem 1.

In case d is induced by a norm $\|\cdot\|$, we will call a contraction with respect to d a $\|\cdot\|$ -contraction. If (X,d) is complete, we can apply the following theorem, due to Banach:

Theorem 1 (Contracting Mapping Principle): Let $f: X \to X$ be a contraction of a complete metric space (X,d). Then f has a unique fixed point $x_\infty \in X$ and for any $x \in X$, the sequence $x, f(x), f^2(x), \ldots$ obtained by iterating f converges to x_∞ . The rate of convergence is at least linear, since $d(f(x), x_\infty) \leq K d(x, x_\infty)$ for all $x \in X$.

Proof: Can be found in many textbooks on analysis. \square Note that linear convergence means that the error decreases exponentially, indeed $d(x_n, x_\infty) \leq CK^n$ for some C.

Let $(V, \|\cdot\|)$ be a normed space. The norm induces a *matrix norm* (also called *operator norm*) on linear mappings $A: V \to V$, defined as follows:

$$||A|| := \sup_{\substack{v \in V, \\ ||v|| \le 1}} ||Av||.$$

The ℓ_1 -norm on \mathbb{R}^N induces the following matrix norm:

$$||A||_1 = \max_{j \in \{1,\dots,N\}} \sum_{i=1}^N |A_{ij}|$$
 (9)

where $A_{ij} := (Ae_j)_i$ with e_j the j^{th} canonical basis vector. The ℓ_{∞} -norm on \mathbb{R}^N induces the following matrix norm:

$$||A||_{\infty} = \max_{i \in \{1, \dots, N\}} \sum_{j=1}^{N} |A_{ij}|.$$
 (10)

In the following consequence of the well-known Mean Value Theorem, the matrix norm of the derivative ("Jacobian") f'(v) at $v \in V$ of a differentiable mapping $f: V \to V$ is used to bound the distance of the f-images of two vectors:

Lemma 1: Let $(V, \|\cdot\|)$ be a normed space and $f: V \to V$ a differentiable mapping. Then, for $x,y \in V$:

$$||f(y) - f(x)|| \le ||y - x|| \cdot \sup_{z \in [x,y]} ||f'(z)||$$

where we wrote [x, y] for the segment $\{\lambda x + (1 - \lambda)y : \lambda \in [0, 1]\}$ joining x and y.

Proof: See [20, Thm.
$$8.5.4$$
].

B. The Basic Tool

Combining Theorem 1 and Lemma 1 immediately yields our basic tool:

Lemma 2: Let $(V, \|\cdot\|)$ be a normed space, $f: V \to V$ differentiable and suppose that

$$\sup_{v \in V} \|f'(v)\| < 1.$$

Then f is a $\|\cdot\|$ -contraction by Lemma 1. Hence, for any $v \in V$, the sequence $v, f(v), f^2(v), \ldots$ converges to a unique fixed point $v_{\infty} \in V$ with a convergence rate that is at least linear by Theorem 1.

C. Sufficient Conditions For BP To Be A Contraction

We apply Lemma 2 to the case at hand: the parallel BP update mapping $f: \mathbb{R}^D \to \mathbb{R}^D$, written out in components in (7). Different choices of the vector norm on \mathbb{R}^D will yield different sufficient conditions for whether iterating f will converge to a unique fixed point. We will study two examples: the ℓ_1 norm and the ℓ_∞ norm.

The derivative of f is easily calculated from (7) and is given by

$$\left(f'(\nu)\right)_{i\to j,k\to l} = \frac{\partial \tilde{\nu}^{i\to j}}{\partial \nu^{k\to l}} = A_{i\to j,k\to l} B_{i\to j}(\nu) \tag{11}$$

where⁵

$$B_{i\to j}(\nu) := \frac{1 - \tanh^2(\theta_i + \sum_{t \in \partial i \setminus j} \nu^{t\to i})}{1 - \tanh^2(\tilde{\nu}^{i\to j}(\nu))} \operatorname{sgn} J_{ij} \quad (12)$$

$$A_{i \to j, k \to l} := \tanh |J_{ij}| \, \delta_{i,l} \mathbf{1}_{\partial i \setminus j}(k). \tag{13}$$

Note that we have absorbed all ν -dependence in the factor $B_{i \to j}(\nu)$; the reason for this will become apparent later on. The factor $A_{i \to j, k \to l}$ is nonnegative and independent of ν and captures the structure of the graphical model. Note that $\sup_{\nu \in V} |B_{i \to j}(\nu)| = 1$, implying that

$$\left| \frac{\partial \tilde{\nu}^{i \to j}}{\partial \nu^{k \to l}} \right| \le A_{i \to j, k \to l} \tag{14}$$

everywhere on V.

1) Example: the ℓ_{∞} -norm: The ℓ_{∞} -norm on \mathbb{R}^D yields the following condition:

Corollary 1: For binary variables with pairwise interactions: if

$$\max_{i \in \mathcal{V}} \left((|\partial i| - 1) \max_{j \in \partial i} \tanh |J_{ij}| \right) < 1, \tag{15}$$

BP is an ℓ_{∞} -contraction and converges to a unique fixed point, irrespective of the initial messages.

Proof: Using (10), (13) and (14):

$$||f'(\nu)||_{\infty} = \max_{i \to j} \sum_{k \to l} \left| \frac{\partial \tilde{\nu}^{i \to j}}{\partial \nu^{k \to l}} \right|$$

$$\leq \max_{i \to j} \sum_{k \to l} \tanh |J_{ij}| \, \delta_{il} \mathbf{1}_{\partial i \setminus j}(k)$$

$$= \max_{i \in \mathcal{V}} \max_{j \in \partial i} \sum_{k \in \partial i \setminus j} \tanh |J_{ij}|$$

$$= \max_{i \in \mathcal{V}} \left((|\partial i| - 1) \max_{j \in \partial i} \tanh |J_{ij}| \right),$$

and now simply apply Lemma 2.

2) Another Example: the ℓ_1 -norm: Using the ℓ_1 -norm instead, we find:

Corollary 2: For binary variables with pairwise interactions:

$$\max_{i \in \mathcal{V}} \max_{k \in \partial i} \sum_{j \in \partial i \setminus k} \tanh |J_{ij}| < 1, \tag{16}$$

BP is an ℓ_1 -contraction and converges to a unique fixed point, irrespective of the initial messages.

⁵For a set X, we define the indicator function $\mathbf{1}_X$ of X by $\mathbf{1}_X(x) = 1$ if $x \in X$ and $\mathbf{1}_X(x) = 0$ if $x \notin X$.

Proof: Similar to the proof of Corollary 1, now using (9) instead of (10):

$$||f'(\nu)||_1 \le \max_{k \to l} \sum_{i \to j} \tanh |J_{ij}| \, \delta_{il} \mathbf{1}_{\partial i \setminus j}(k)$$
$$= \max_{i \in \mathcal{V}} \max_{k \in \partial i} \sum_{j \in \partial i \setminus k} \tanh |J_{ij}|.$$

It is easy to see that condition (16) is implied by (15), but not conversely; thus in this case the ℓ_1 -norm yields a tighter bound than the ℓ_∞ -norm.

D. Beyond Norms: The Spectral Radius

Instead of pursuing a search for the optimal norm, we will derive a criterion for convergence based on the spectral radius of the matrix (13). The key idea is to look at several iterations of BP at once. This will yield a significantly stronger condition for convergence of BP to a unique fixed point.

For a square matrix A, we denote by $\sigma(A)$ its *spectrum*, i.e., the set of eigenvalues of A. By $\rho(A)$ we denote its *spectral radius*, which is defined as $\rho(A) := \sup |\sigma(A)|$, i.e., the largest modulus of eigenvalues of A.

Lemma 3: Let $f: X \to X$ be a mapping, d a metric on X and suppose that f^N is a d-contraction for some $N \in \mathbb{N}$. Then f has a unique fixed point x_∞ and for any $x \in X$, the sequence $x, f(x), f^2(x), \ldots$ obtained by iterating f converges to x_∞ .

Proof: Take any $x \in X$. Consider the N sequences obtained by iterating f^N , starting respectively in x, f(x), ..., $f^{N-1}(x)$:

$$x, f^{N}(x), f^{2N}(x), \dots$$

 $f(x), f^{N+1}(x), f^{2N+1}(x), \dots$
 \vdots
 $f^{N-1}(x), f^{2N-1}(x), f^{3N-1}(x), \dots$

Each sequence converges to x_{∞} since f^N is a d-contraction with fixed point x_{∞} . But then the sequence $x, f(x), f^2(x), \ldots$ must converge to x_{∞} .

Theorem 2: Let $f: \mathbb{R}^m \to \mathbb{R}^m$ be differentiable and suppose that f'(x) = B(x)A, where A has nonnegative entries and B is diagonal with bounded entries $|B_{ii}(x)| \leq 1$. If $\rho(A) < 1$ then for any $x \in \mathbb{R}^m$, the sequence $x, f(x), f^2(x), \ldots$ obtained by iterating f converges to a fixed point x_{∞} , which does not depend on x.

Proof: For a matrix B, we will denote by |B| the matrix with entries $|B|_{ij} = |B_{ij}|$. For two matrices B, C we will write $B \leq C$ if $B_{ij} \leq C_{ij}$ for all entries (i,j). Note that if $|B| \leq |C|$, then $\|B\|_1 \leq \|C\|_1$. Also note that $|BC| \leq |B| |C|$. Finally, if $0 \leq A$ and $B \leq C$, then $AB \leq AC$ and $BA \leq CA$.

⁶One should not confuse the spectral *radius* $\rho(A)$ with the spectral *norm* $||A||_2 = \sqrt{\rho(A^T A)}$ of A, the matrix norm induced by the ℓ_2 -norm.

Using these observations and the chain rule, we have for any $n = 1, 2, \ldots$ and any $x \in \mathbb{R}^m$:

$$|(f^n)'(x)| = \left| \prod_{i=1}^n f'(f^{i-1}(x)) \right|$$

$$\leq \prod_{i=1}^n \left(\left| B(f^{i-1}(x)) \right| A \right) \leq A^n,$$

hence $||(f^n)'(x)||_1 \le ||A^n||_1$.

By the Gelfand spectral radius theorem,

$$\lim_{n \to \infty} \|A^n\|_1^{-1/n} = \rho(A).$$

Choose $\epsilon>0$ such that $\rho(A)+\epsilon<1$. For some $N, \|A^N\|_1\leq (\rho(A)+\epsilon)^N<1$. Hence for all $x\in\mathbb{R}^m, \|(f^N)'(x)\|_1<1$. Applying Lemma 2, we conclude that f^N is a ℓ_1 -contraction. Now apply Lemma 3.

Using (11), (12) and (13), this immediately yields:

Corollary 3: For binary variables with pairwise interactions, BP converges to a unique fixed point, irrespective of the initial messages, if the spectral radius of the $|D| \times |D|$ -matrix

$$A_{i \to j, k \to l} := \tanh |J_{ij}| \, \delta_{i,l} \mathbf{1}_{\partial i \setminus j}(k)$$

is strictly smaller than 1.

The calculation of the spectral norm of the (sparse) matrix A can be done using standard numerical techniques in linear algebra.

Any matrix norm of A is actually an upper bound on the spectral radius $\rho(A)$, since for any eigenvalue λ of A with eigenvector x we have $|\lambda| \ \|x\| = \|\lambda x\| = \|Ax\| \le \|A\| \ \|x\|$, hence $\rho(A) \le \|A\|$. This implies that no norm in Lemma 2 will result in a sharper condition than Corollary 3, hence the title of this section.

Further, for a given matrix A and some $\epsilon>0$, there exists a vector norm $\|\cdot\|$ such that the induced matrix norm of A satisfies $\rho(A)\leq \|A\|\leq \rho(A)+\epsilon$; see [21] for a constructive proof. Thus for given A one can approximate $\rho(A)$ arbitrarily close by induced matrix norms. This immediately gives a result on the convergence rate of BP (in case $\rho(A)<1$): for any $\epsilon>0$, there exists a norm-induced metric such that the linear rate of contraction of BP with respect to that metric is bounded from above by $\rho(A)+\epsilon$.

One might think that there is a shorter proof of Corollary 3: it seems quite plausible intuitively that in general, for a continuously differentiable $f:\mathbb{R}^m\to\mathbb{R}^m$, iterating f will converge to a unique fixed point if $\sup_{x\in\mathbb{R}^m}\rho(f'(x))<1$. However, this conjecture (which has been open for a long time) has been shown to be true in two dimensions but false in higher dimensions [22].

E. Improved Bound For Strong Local Evidence

Empirically, it is known that the presence of strong local fields (i.e., single-variable factors which are far from uniform) often improves the convergence of BP. However, our results so far are completely independent of the parameters $(\theta_i)_{i\in\mathcal{V}}$ that measure the strength of the local evidence. By proceeding more carefully than we have done above, the results can easily

be improved in such a way that local evidence is taken into account.

Consider the quantity $B_{i\to j}$ defined in (12). We have bounded this quantity by noting that $\sup_{\nu\in V}|B_{i\to j}(\nu)|=1$. Note that for all BP updates (except for the first one), the argument ν (the incoming messages) is in f(V), which can be considerably smaller than the complete vector space V. Thus, after the first BP update, we can use

$$\sup_{\nu \in f(V)} |B_{i \to j}(\nu)| = \sup_{\nu \in f(V)} \frac{1 - \tanh^2(\theta_i + \sum_{k \in \partial i \setminus j} \nu^{k \to i})}{1 - \tanh^2(\tilde{\nu}^{i \to j}(\nu))}$$
$$= \sup_{\nu \in f(V)} \frac{1 - \tanh^2(h^{i \setminus j})}{1 - \tanh^2(J_{ij}) \tanh^2(h^{i \setminus j})}$$

where we used (7) and defined the cavity field

$$h^{i \setminus j}(\nu) := \theta_i + \sum_{k \in \partial i \setminus j} \nu^{k \to i}. \tag{17}$$

The function $x\mapsto \frac{1-\tanh^2 x}{1-\tanh^2(J_{ij})\tanh^2 x}$ is strictly decreasing for $x\geq 0$ and symmetric around x=0, thus, defining

$$h_*^{i\backslash j} := \inf_{\nu \in f(V)} \left| h^{i\backslash j}(\nu) \right|,\tag{18}$$

we obtain

$$\sup_{\nu \in f(V)} |B_{i \to j}(\nu)| = \frac{1 - \tanh^2(h_*^{i \setminus j})}{1 - \tanh^2(J_{ij}) \tanh^2(h_*^{i \setminus j})}.$$

Now, from (7) we derive that

$$\{\nu^{k\to i}: \nu \in f(V)\} = (-|J_{ki}|, |J_{ki}|),$$

hence

$$\{h^{i\backslash j}(\nu):\nu\in f(V)\}=(h_-^{i\backslash j},h_+^{i\backslash j})$$

where we defined

$$h_{\pm}^{i\setminus j} := \theta_i \pm \sum_{k\in\partial i\setminus j} |J_{ki}|.$$

We conclude that $h_*^{i\backslash j}$ is simply the distance between 0 and the interval $(h_-^{i\backslash j},h_+^{i\backslash j})$, i.e.,

$$h_*^{i\backslash j} = \begin{cases} \left|h_+^{i\backslash j}\right| & \text{if } h_+^{i\backslash j} < 0\\ h_-^{i\backslash j} & \text{if } h_-^{i\backslash j} > 0\\ 0 & \text{otherwise.} \end{cases}$$

Thus the element $A_{i \to j, k \to i}$ (for $i \in \partial j, k \in \partial i \setminus j$) of the matrix A defined in Corollary 3 can be replaced by

$$\tanh |J_{ij}| \frac{1 - \tanh^2(h_*^{i \setminus j})}{1 - \tanh^2(J_{ij}) \tanh^2(h_*^{i \setminus j})}$$
$$= \frac{\tanh(|J_{ij}| - h_*^{i \setminus j}) + \tanh(|J_{ij}| + h_*^{i \setminus j})}{2},$$

which is generally smaller than $\tanh |J_{ij}|$ and thus gives a tighter bound.

This trick can be repeated arbitrarily often: assume that $m \ge 0$ BP updates have been done already, which means that it

suffices to take the supremum of $|B_{i\to j}(\nu)|$ over $\nu \in f^m(V)$. Define for all $i\to j\in D$ and all $t=0,1,\ldots,m$:

$$\underline{h}_t^{i \setminus j} := \inf\{h^{i \setminus j}(\nu) : \nu \in f^t(V)\},\tag{19}$$

$$\overline{h}_t^{i\backslash j} := \sup\{h^{i\backslash j}(\nu) : \nu \in f^t(V)\},\tag{20}$$

and define the intervals

$$\mathcal{H}_t^{i\backslash j} := [\underline{h}_t^{i\backslash j}, \overline{h}_t^{i\backslash j}]. \tag{21}$$

Specifically, for t=0 we have $\underline{h}_0^{i\backslash j}=-\infty$ and $\overline{h}_0^{i\backslash j}=\infty$, which means that

$$\mathcal{H}_0^{i\backslash j} = \mathbb{R}.\tag{22}$$

Using (7) and (17), we obtain the following recursive relations for the intervals (where we use interval arithmetic defined in the obvious way):

$$\mathcal{H}_{t+1}^{i\setminus j} = \theta_i + \sum_{k\in\partial i\setminus j} \tanh^{-1} \left(\tanh J_{ki} \tanh \mathcal{H}_t^{k\setminus i}\right). \tag{23}$$

Using this recursion relation, one can calculate $\mathcal{H}_m^{i\setminus j}$ and define $h_*^{i\setminus j}$ as the distance (in absolute value) of the interval $\mathcal{H}_m^{i\setminus j}$ to 0:

$$h_*^{i\backslash j} = \begin{cases} \left| \overline{h}_m^{i\backslash j} \right| & \text{if } \overline{h}_m^{i\backslash j} < 0\\ \underline{h}_m^{i\backslash j} & \text{if } \underline{h}_m^{i\backslash j} > 0\\ 0 & \text{otherwise.} \end{cases}$$
 (24)

Thus by replacing the matrix A in Corollary 3 by

$$A_{i \to j, k \to l} = \frac{\tanh(|J_{ij}| - h_*^{i \setminus j}) + \tanh(|J_{ij}| + h_*^{i \setminus j})}{2} \, \delta_{i,l} \mathbf{1}_{\partial i \setminus j}(k),$$
(25)

we obtain stronger results that improve as m increases:

Corollary 4: Let $m \geq 0$. For binary variables with pairwise interactions, BP converges to a unique fixed point, irrespective of the initial messages, if the spectral radius of the $|D| \times |D|$ -matrix defined in (25) (with $h_*^{i \setminus j}$ defined in equations (21)–(24)) is strictly smaller than 1.

IV. GENERAL CASE

In the general case, when the domains \mathcal{X}_i are arbitrarily large (but finite), we do not know of a natural parameterization of the messages that automatically takes care of the invariance of the messages $\mu^{I \to j}$ under scaling (like (6) does in the binary case). Instead of handling the scale invariance by the parameterization and using standard norms and metrics, it seems easier to take a simple parameterization and to change the norms and metrics in such a way that they are insensitive to the (irrelevant) extra degrees of freedom arising from the scale invariance. This is actually the key insight in extending the previous results beyond the binary case: once one sees how to do this, the rest follows in a (more or less) straightforward way.

Another important point is to reparameterize the messages: a natural parameterization for our analysis is now in terms of logarithms of messages $\lambda^{I \to i} := \log \mu^{I \to i}$. The BP update equations (4) can be written in terms of the log-messages as:

$$\tilde{\lambda}^{I \to i}(x_i) = \log \sum_{x_{I \setminus i}} \psi^I(x_I) h^{I \setminus i}(x_{I \setminus i})$$
 (26)

where we dropped the normalization and defined

$$h^{I\setminus i}(x_{I\setminus i}) := \exp\left(\sum_{j\in I\setminus i} \sum_{J\in N_j\setminus I} \lambda^{J\to j}(x_j)\right). \tag{27}$$

Each log-message $\lambda^{I \to i}$ is a vector in the vector space $\mathcal{V}^{I \to i} := \mathbb{R}^{\mathcal{X}_i}$; we will use Greek letters as indices for the components, e.g., $\lambda_{\alpha}^{I \to i} := \lambda^{I \to i}(\alpha)$ with $\alpha \in \mathcal{X}_i$. We will call everything that concerns individual vector spaces $\mathcal{V}^{I \to i}$ local and define the global vector space \mathcal{V} as the direct sum of the local vector spaces:

$$\mathcal{V} := \bigoplus_{i \in I \in \mathcal{F}} \mathcal{V}^{I \to i}.$$

The parallel BP update is the mapping $f: \mathcal{V} \to \mathcal{V}$, written out in components in (26) and (27).

Note that the invariance of the messages $\mu^{I \to i}$ under scaling amounts to invariance of the log-messages $\lambda^{I \to i}$ under translation. More formally, defining linear subspaces

$$\mathcal{W}^{I \to i} := \{ \lambda \in \mathcal{V}^{I \to i} : \lambda_{\alpha} = \lambda_{\alpha'} \text{ for all } \alpha, \alpha' \in \mathcal{X}_i \} \quad (28)$$

and their direct sum

$$\mathcal{W} := \bigoplus_{i \in I \in \mathcal{F}} \mathcal{W}^{I \to i} \subseteq \mathcal{V},$$

the invariance amounts to the observation that

$$f(\lambda + w) - f(\lambda) \in \mathcal{W}$$
 for all $\lambda \in \mathcal{V}, w \in \mathcal{W}$.

Since $\lambda + w$ and λ are equivalent for our purposes, we want our measures of distance in $\mathcal V$ to reflect this equivalence. Therefore we will "divide out" the equivalence relation and work in the quotient space $\mathcal V/\mathcal W$, which is the topic of the next subsection.

A. Quotient Spaces

Let V be a finite-dimensional vector space. Let W be a linear subspace of V. We can consider the *quotient space* $V/W:=\{v+W:v\in V\}$, where $v+W:=\{v+w:w\in W\}$. Defining addition and scalar multiplication on the quotient space in the natural way, the quotient space is again a vector space. We will denote its elements as $\overline{v}:=v+W$. Note that the projection $\pi:V\to V/W:v\mapsto \overline{v}$ is linear.

Let $\|\cdot\|$ be any vector norm on V. It induces a *quotient norm* on V/W, defined by

$$\|\overline{v}\| := \inf_{w \in W} \|v + w\|,$$
 (29)

which is indeed a norm, as one easily checks. The quotient norm in turn induces the *quotient metric* $d(\overline{v_1}, \overline{v_2}) := \|\overline{v_2} - \overline{v_1}\|$ on V/W. The metric space (V/W, d) is complete

⁷Indeed, we have a null vector 0+W, addition $(v_1+W)+(v_2+W):=(v_1+v_2)+W$ for $v_1,v_2\in V$ and scalar multiplication $\lambda(v+W):=(\lambda v)+W$ for $\lambda\in\mathbb{R},v\in V$.

(since any finite-dimensional normed vector space is complete).

Let $f: V \to V$ be a (possibly nonlinear) mapping with the following symmetry:

$$f(v+w) - f(v) \in W$$
 for all $v \in V$, $w \in W$. (30)

We can then unambiguously define the quotient mapping

$$\overline{f}: V/W \to V/W : \overline{v} \mapsto \overline{f(v)},$$

which yields the following commutative diagram:

$$V \xrightarrow{f} V$$

$$\downarrow \pi \qquad \qquad \downarrow \pi \qquad \qquad \pi \circ f = \overline{f} \circ \pi$$

$$V/W \xrightarrow{\overline{f}} V/W$$

For a linear mapping $A:V\to V$, condition (30) amounts to $AW\subseteq W$, i.e., A should leave W invariant; we can then unambiguously define the quotient mapping $\overline{A}:V/W\to V/W:\overline{v}\mapsto \overline{Av}$.

If $f:V\to V$ is differentiable and satisfies (30), the symmetry property (30) implies that $f'(x)W\subseteq W$, hence we can define $\overline{f'(x)}:V/W\to V/W$. The operation of taking derivatives is compatible with projecting onto the quotient space. Indeed, by using the chain rule and the identity $\pi\circ f=\overline{f}\circ\pi$, one finds that the derivative of the induced mapping $\overline{f}:V/W\to V/W$ at \overline{x} equals the induced derivative of f at x:

$$\overline{f}'(\overline{x}) = \overline{f'(x)}$$
 for all $x \in V$. (31)

By Lemma 2, \overline{f} is a contraction with respect to the quotient norm if

$$\sup_{\overline{x} \in V/W} \left\| \overline{f}'(\overline{x}) \right\| < 1.$$

Using (29) and (31), this condition can be written more explicitly as:

$$\sup_{x \in V} \sup_{\substack{v \in V, \\ \|v\| \le 1}} \inf_{w \in W} \|f'(x) \cdot v + w\| < 1.$$

B. Constructing A Norm On V

Whereas in the binary case, each message $\nu^{i \to j}$ was parameterized by a single real number, the messages are now $|\mathcal{X}_i|$ -dimensional vectors $\lambda^{I \to i}$ (with components $\lambda^{I \to i}_{\alpha}$ indexed by $\alpha \in \mathcal{X}_i$). In extending the ℓ_1 -norm that proved to be useful in the binary case to the more general case, we have the freedom to choose the "local" part of the generalized ℓ_1 -norm. Here we show how to construct such a generalization of the ℓ_1 -norm and its properties; for a more detailed account of the construction, see Appendix A.

The "global" vector space $\mathcal V$ is the direct sum of the "local" subspaces $\mathcal V^{I \to i}$. Suppose that for each subspace $\mathcal V^{I \to i}$, we have a local norm $\|\cdot\|_{I \to i}$. A natural generalization of the ℓ_1 -norm in the binary case is the following global norm on $\mathcal V$:

$$\|\lambda\| := \sum_{I \to i} \|\lambda^{I \to i}\|_{I \to i} . \tag{32}$$

It is easy to check that this is indeed a norm on V.

Each subspace $\mathcal{V}^{I \to i}$ has a 1-dimensional subspace $\mathcal{W}^{I \to i}$ defined in (28) and the local norm on $\mathcal{V}^{I \to i}$ induces a local quotient norm on the quotient space $\mathcal{V}^{I \to i}/\mathcal{W}^{I \to i}$. The global norm (32) on \mathcal{V} induces a global quotient norm on \mathcal{V}/\mathcal{W} , which is simply the sum of the local quotient norms (c.f. (A.57)):

$$\|\overline{\lambda}\| = \sum_{I \to i} \|\overline{\lambda^{I \to i}}\|_{I \to i}. \tag{33}$$

Let $\lambda \in \mathcal{V}$. The derivative $f'(\lambda)$ of $f: \mathcal{V} \to \mathcal{V}$ at λ is a linear mapping $f'(\lambda): \mathcal{V} \to \mathcal{V}$ satisfying $f'(\lambda)\mathcal{W} \subseteq \mathcal{W}$. It projects down to a linear mapping $\overline{f'(\lambda)}: \mathcal{V}/\mathcal{W} \to \mathcal{V}/\mathcal{W}$. The matrix norm of $\overline{f'(\lambda)}$ induced by the quotient norm (33) is given by (c.f. (A.58)):

$$\left\| \overline{f'(\lambda)} \right\| = \max_{J \to j} \sum_{I \to i} \left\| \overline{\left(f'(\lambda) \right)_{I \to i, J \to j}} \right\|_{I \to i}^{J \to j} \tag{34}$$

where the local quotient matrix norm of the "block" $(f'(\lambda))_{I \to i, J \to j}$ is given by (c.f. (A.59)):

$$\left\| \overline{(f'(\lambda))}_{I \to i, J \to j} \right\|_{I \to i}^{J \to j}$$

$$= \sup_{\substack{v \in \mathcal{V}^{J \to j}, \\ \|v\|_{I \to i} \le 1}} \left\| \overline{(f'(\lambda))}_{I \to i, J \to j} v \right\|_{I \to i}. \tag{35}$$

The derivative of the (unnormalized) parallel BP update (26) is easily calculated:

$$\frac{\partial \tilde{\lambda}^{I \to i}(x_i)}{\partial \lambda^{J \to j}(y_j)} = \mathbf{1}_{N_j \setminus I}(J) \mathbf{1}_{I \setminus i}(j) \\
\times \frac{\sum_{x_{I \setminus i}} \psi^I(x_i, x_j, x_{I \setminus \{i,j\}}) \delta_{x_j, y_j} h^{I \setminus i}(x_{I \setminus i})}{\sum_{x_{I \setminus i}} \psi^I(x_i, x_{I \setminus i}) h^{I \setminus i}(x_{I \setminus i})}.$$
(36)

To lighten the notation, we will use Greek subscripts instead of arguments: let α correspond to x_i , β to x_j , β' to y_j and γ to $x_{I\setminus\{i,j\}}$; for example, we write $h^{I\setminus i}(x_{I\setminus i})$ as $h^{I\setminus i}_{\beta\gamma}$. Taking the global quotient norm (34) of (36) yields:

$$\left\| \overline{f'(\lambda)} \right\| = \max_{J \to j} \sum_{I \to i} \mathbf{1}_{N_j \setminus I}(J) \mathbf{1}_{I \setminus i}(j) B_{I \to i, J \to j} \left(h^{I \setminus i}(\lambda) \right)$$
(37)

where

$$B_{I \to i, J \to j} \left(h^{I \setminus i}(\lambda) \right) := \left\| \frac{\sum_{\gamma} \psi_{\alpha\beta'\gamma}^{I} h_{\beta'\gamma}^{I \setminus i}(\lambda)}{\sum_{\beta} \sum_{\gamma} \psi_{\alpha\beta\gamma}^{I} h_{\beta\gamma}^{I \setminus i}(\lambda)} \right\|_{I \to i}^{J \to j} . \tag{38}$$

Note that $B_{I \to i, J \to j}$ depends on λ via the dependence of $h^{I \setminus i}$ on λ (c.f. (27)). We will for the moment simplify matters by assuming that λ can be any vector in \mathcal{V} , and later discuss the more careful estimate (where $\lambda \in f^m(\mathcal{V})$):

$$\sup_{\lambda \in \mathcal{V}} B_{I \to i, J \to j} (h^{I \setminus i}(\lambda)) \le \sup_{h^{I \setminus i} > 0} B_{I \to i, J \to j} (h^{I \setminus i}).$$
 (39)

Defining the matrix A by the expression on the r.h.s. and using (35) and (29), we obtain:

$$A_{I \to i, J \to j} := \sup_{h^{I \setminus i} > 0} B_{I \to i, J \to j}(h^{I \setminus i}) =$$

$$\sup_{h^{I\backslash i}>0} \sup_{\substack{v\in\mathcal{V}^{J\to j}\\\|v\|_{J\to j}\leq 1}} \inf_{w\in\mathcal{W}^{I\to i}} \left\| \frac{\sum_{\beta'} \sum_{\gamma} \psi_{\alpha\beta'\gamma}^{I} h_{\beta'\gamma}^{I\backslash i} v_{\beta'}}{\sum_{\beta} \sum_{\gamma} \psi_{\alpha\beta\gamma}^{I} h_{\beta\gamma}^{I\backslash i}} - w \right\|_{I\to i}$$

$$\tag{40}$$

for $I \to i$ and $J \to j$ such that $j \in I \setminus i$ and $J \in N_j \setminus I$. Surprisingly, it turns out that we can calculate (40) analytically if we take all local norms to be ℓ_∞ norms. We have also tried the ℓ_2 norm and the ℓ_1 norm as local norms, but were unable to calculate expression (40) analytically in these cases. Numerical calculations turned out to be difficult because of the nested suprema.

C. Local ℓ_{∞} Norms

Take for each local norm $\|\cdot\|_{I \to i}$ the ℓ_{∞} norm on $\mathcal{V}^{I \to i} = \mathbb{R}^{\mathcal{X}_i}$. The local subspace $\mathcal{W}^{I \to i}$ is spanned by the vector $\mathbf{1} := (1,1,\ldots,1) \in \mathbb{R}^{\mathcal{X}_i}$. The local quotient norm of a vector $v \in \mathcal{V}^{I \to i}$ is thus given by

$$\|\overline{v}\|_{I \to i} = \|\overline{v}\|_{\infty} = \inf_{w \in \mathbb{R}} \|v + w\mathbf{1}\|_{\infty}$$

$$= \frac{1}{2} \sup_{\alpha, \alpha' \in \mathcal{X}_i} |v_{\alpha} - v_{\alpha'}|.$$
(41)

For a linear mapping $A: \mathcal{V}^{J \to j} \to \mathcal{V}^{I \to i}$ that satisfies $A\mathcal{W}^{J \to j} \subseteq \mathcal{W}^{I \to i}$, the induced quotient matrix norm (35) is given by

$$\|\overline{A}\|_{I \to i}^{J \to j} = \sup_{\substack{v \in \mathcal{V}^{J \to j}, \\ \|v\|_{\infty} \le 1}} \|\overline{Av}\|_{\infty}$$

$$= \sup_{\substack{v \in \mathcal{V}^{J \to j}, \\ \|v\|_{\infty} \le 1}} \frac{1}{2} \sup_{\alpha, \alpha' \in \mathcal{X}_{i}} \left| \sum_{\beta} (A_{\alpha\beta} - A_{\alpha'\beta}) v_{\beta} \right|$$

$$= \frac{1}{2} \sup_{\alpha, \alpha' \in \mathcal{X}_{i}} \sum_{\beta} |A_{\alpha\beta} - A_{\alpha'\beta}|$$

$$(42)$$

Fixing for the moment $I \to i$ and $J \to j$ (such that $j \in I \setminus i$ and $J \in N_j \setminus I$) and dropping the superscripts from the notation, using (42), we can write (40) as

$$\sup_{h>0} \frac{1}{2} \sup_{\alpha,\alpha' \in \mathcal{X}_i} \sum_{\beta} \left| \frac{\sum_{\gamma} \psi_{\alpha\beta\gamma} h_{\beta\gamma}}{\sum_{\beta} \sum_{\gamma} \psi_{\alpha\beta\gamma} h_{\beta\gamma}} - \frac{\sum_{\gamma} \psi_{\alpha'\beta\gamma} h_{\beta\gamma}}{\sum_{\beta} \sum_{\gamma} \psi_{\alpha'\beta\gamma} h_{\beta\gamma}} \right|.$$

Interchanging the two suprema, fixing (for the moment) α and α' , defining $\tilde{\psi}_{\beta\gamma}:=\psi_{\alpha\beta\gamma}/\psi_{\alpha'\beta\gamma}$ and $\tilde{h}_{\beta\gamma}:=h_{\beta\gamma}\psi_{\alpha'\beta\gamma}$, noting that we can assume (without loss of generality) that \tilde{h} is normalized in ℓ_1 sense, the previous expression (apart from the $\frac{1}{2}\sup_{\alpha,\alpha'}$) simplifies to

$$\sup_{\substack{\tilde{h} > 0, \\ \|\tilde{h}\|_{1} = 1}} \sum_{\beta} \left| \sum_{\gamma} \tilde{h}_{\beta\gamma} \left(\frac{\tilde{\psi}_{\beta\gamma}}{\sum_{\beta} \sum_{\gamma} \tilde{\psi}_{\beta\gamma} \tilde{h}_{\beta\gamma}} - 1 \right) \right|. \tag{43}$$

In Appendix B we show that this equals

$$2 \sup_{\beta \neq \beta'} \sup_{\gamma, \gamma'} \tanh \left(\frac{1}{4} \log \frac{\tilde{\psi}_{\beta \gamma}}{\tilde{\psi}_{\beta' \gamma'}} \right). \tag{44}$$

We conclude that if we take all local norms to be the ℓ_∞ norms, then $A_{I \to i, J \to j}$ equals

$$N(\psi^{I}, i, j) = \sup_{\alpha \neq \alpha'} \sup_{\beta \neq \beta'} \sup_{\gamma, \gamma'} \tanh \left(\frac{1}{4} \log \frac{\psi^{I}_{\alpha\beta\gamma}}{\psi^{I}_{\alpha'\beta\gamma}} \frac{\psi^{I}_{\alpha'\beta'\gamma'}}{\psi^{I}_{\alpha\beta'\gamma'}} \right), \tag{45}$$

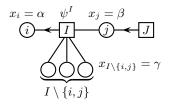


Fig. 2. Part of the factor graph relevant in expressions (45), (46) and (47). Here $i, j \in I$ with $i \neq j$, and $J \in N_j \setminus I$.

which is defined for $i,j \in I$ with $i \neq j$ and where $\psi^I_{\alpha\beta\gamma}$ is shorthand for $\psi^I(x_i = \alpha, x_j = \beta, x_{I\setminus\{i,j\}} = \gamma)$; see Fig. 2 for an illustration.

Now combining (37), (39) and (45), we finally obtain:

$$\left\| \overline{f'(\lambda)} \right\| = \left\| \overline{f'(\lambda)} \right\| \le \max_{J \to j} \sum_{I \in N_j \setminus J} \sum_{i \in I \setminus j} N(\psi^I, i, j).$$

Applying Lemma 2 now yields that f is a contraction with respect to the quotient norm on V/W if the right-hand side is strictly smaller than 1.

Consider the mapping $\eta: \mathcal{V}/\mathcal{W} \to \mathcal{V}$ that maps $\overline{\lambda}$ to the normalized $\lambda \in \mathcal{V}$, i.e., such that $\left\|\exp \lambda^{I \to i}\right\|_1 = 1$ for all components $I \to i$. If we take for f the ℓ_1 -normalized BP update (in the log-domain), the following diagram commutes:

$$\begin{array}{ccc}
\mathcal{V} & \stackrel{f}{\longrightarrow} & \mathcal{V} \\
\downarrow^{\pi} & & \uparrow^{\eta} & f = \eta \circ \overline{f} \circ \pi. \\
\mathcal{V}/\mathcal{W} & \stackrel{\overline{f}}{\longrightarrow} & \mathcal{V}/\mathcal{W}
\end{array}$$

Since both π and $\underline{\eta}$ are continuous, we can translate convergence results for \overline{f} back to similar results for f. We have proved:

Theorem 3: If

$$\max_{J \to j} \sum_{I \in N_i \setminus J} \sum_{i \in I \setminus j} N(\psi^I, i, j) < 1, \tag{46}$$

BP converges to a unique fixed point irrespective of the initial messages.

Now we can also generalize Corollary 3:

Theorem 4: If the spectral radius of the matrix

$$A_{I \to i, J \to j} = \mathbf{1}_{N_i \setminus I}(J) \mathbf{1}_{I \setminus i}(j) N(\psi^I, i, j), \tag{47}$$

is strictly smaller than 1, BP converges to a unique fixed point irrespective of the initial messages.

Proof: Similar to the binary pairwise case; see Theorem 10 in Appendix A for details.

Note that Theorem 3 is a trivial consequence of Theorem 4, since the ℓ_1 -norm is an upper bound on the spectral radius. However, to prove the latter, it seems that we have to go through all the work (and some more) needed to prove the former.

D. Special Cases

In this subsection we study the implications for two special cases, namely factor graphs that contain no cycles and the case of pairwise interactions.

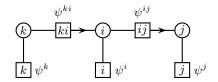


Fig. 3. Part of the factor graph in the pairwise case relevant in (48) and (49). Here $k \in \partial i$ and $j \in \partial i \setminus k$.

1) Trees: Theorem 4 gives us a proof of the well-known fact that BP converges on trees (whereas Theorem 3 is not strong enough to prove that result):

Corollary 5: If the factor graph is a tree, BP converges to a unique fixed point irrespective of the initial messages.

Proof: The spectral radius of (47) is easily shown to be zero in this special case, for any choice of the potentials. \Box

2) Pairwise Interactions: We formulate Theorems 3 and 4 for the special case of pairwise interactions (which corresponds to γ taking on only one value), i.e., all factors consists of either one or two variables. For a pair-potential $\psi^{ij} = \psi^{ij}_{\alpha\beta}$, expression (45) simplifies to (see also Fig. 3)

$$N(\psi^{ij}) := \sup_{\alpha \neq \alpha'} \sup_{\beta \neq \beta'} \tanh \left(\frac{1}{4} \left(\log \frac{\psi^{ij}_{\alpha\beta}}{\psi^{ij}_{\alpha'\beta}} \frac{\psi^{ij}_{\alpha'\beta'}}{\psi^{ij}_{\alpha\beta'}} \right) \right). \quad (48)$$

Note that this quantity is invariant to "reallocation" of single-variable factors ψ^i or ψ^j to the pairwise factor ψ^{ij} (i.e., $N(\psi^{ij}) = N(\psi^{ij}\psi^i\psi^j)$). $N(\psi^{ij})$ can be regarded as a measure of the strength of the potential ψ^{ij} .

The ℓ_1 -norm based condition (46) can be written in the pairwise case as:

$$\max_{i \in \mathcal{V}} \max_{k \in \partial i} \sum_{j \in \partial i \setminus k} N(\psi^{ij}) < 1. \tag{49}$$

The matrix defined in (47), relevant for the spectral radius condition, can be replaced by the following $|D| \times |D|$ -matrix in the pairwise case:

$$A_{i \to j, k \to l} := N(\psi^{ij}) \delta_{i,l} \mathbf{1}_{\partial i \setminus j}(k). \tag{50}$$

For the binary case, we reobtain our earlier results, since $N(\exp(J_{ij}x_ix_j)) = \tanh |J_{ij}|$.

E. Factors Containing Zeros

Until now, we have assumed that all factors are strictly positive. In many interesting applications of the Sum-Product Algorithm, this assumption is violated: the factors may contain zeros. It is thus interesting to see if and how our results can be extended towards this more general case.

The easiest way to extend the results is by assuming that—although the factors may contain zeros—the messages are guaranteed to remain strictly positive (i.e., the log-messages remain finite) after each update.⁸ Even more general extensions with milder conditions may exist, but we believe that considerably more work would be required to overcome the technical problems that arise due to messages containing zeros.

Assume that each factor ψ^I is a nonnegative function $\psi^I:\prod_{i\in I}\mathcal{X}_i\to [0,\infty)$. In addition, assume that all factors involving only a single variable are strictly positive. This can be assumed without loss of generality, since the single-variable factors that contain one or more zeros can simply be absorbed into multi-variable factors involving the same variable. Additionally, for each $I\in\mathcal{F}$ consisting of more than one variable, assume that

$$\forall_{i \in I} \, \forall_{x_i \in \mathcal{X}_i} \exists_{x_{I \setminus i} \in \mathcal{X}_{I \setminus i}} : \psi^I(x_i, x_{I \setminus i}) > 0.$$
 (51)

These conditions guarantee that strictly positive messages remain strictly positive under the update equations (4), as one easily checks, implying that we can still use the logarithmic parameterization of the messages and that the derivative (36) is still well-defined.

The expression for the potential strength (45) can be written in a way that is also well-defined if the potential ψ^I contains zeros:

$$N(\psi^{I}, i, j) = \sup_{\alpha \neq \alpha'} \sup_{\beta \neq \beta'} \sup_{\gamma, \gamma'} \frac{\sqrt{\psi^{I}_{\alpha\beta\gamma}\psi^{I}_{\alpha'\beta'\gamma'}} - \sqrt{\psi^{I}_{\alpha'\beta\gamma}\psi^{I}_{\alpha\beta'\gamma'}}}{\sqrt{\psi^{I}_{\alpha\beta\gamma}\psi^{I}_{\alpha'\beta'\gamma'}} + \sqrt{\psi^{I}_{\alpha'\beta\gamma}\psi^{I}_{\alpha\beta'\gamma'}}}$$
(52)

which is defined for $i, j \in I$ with $i \neq j$ and where $\psi^I_{\alpha\beta\gamma}$ is shorthand for $\psi^I(x_i = \alpha, x_j = \beta, x_{I\setminus\{i,j\}} = \gamma)$.

The immediate generalization of Corollary 4 is then as follows:

Theorem 5: Under the assumptions on the potentials described above (strict positivity of single-variable factors and (51) for the other factors): if the spectral radius of the matrix

$$A_{I \to i, J \to j} = \mathbf{1}_{N_j \setminus I}(J) \mathbf{1}_{I \setminus i}(j) N(\psi^I, i, j), \tag{53}$$

(with $N(\psi^I, i, j)$ defined in (52)) is strictly smaller than 1, BP converges to a unique fixed point irrespective of the initial messages.

Proof: Similar to the strictly positive case. The only slight subtlety occurs in Appendix B where one has to take a limit of strictly positive factors converging to the desired nonnegative factor and use the continuity of the relevant expressions with respect to the factor entries to prove that the bound also holds in this limit.

1) Example: Define, for $\epsilon \geq 0$, the ("ferromagnetic") pairwise factor $\psi(\epsilon)$ by the following matrix:

$$\psi(\epsilon) := \begin{pmatrix} 1 & \epsilon \\ \epsilon & 1 \end{pmatrix}.$$

Now consider a binary pairwise factor graph, consisting of a single loop of N binary variables, i.e., the network topology is that of a circle. Take for the N-1 pair interactions $\psi^{\{i,i+1\}}$ (for $i=1,2,\ldots,N-1$) the identity matrices (i.e., the above pair factors for $\epsilon=0$) and take for the remaining one $\psi^{\{1,N\}}=\psi(\epsilon)$ for some $\epsilon\geq0$. Note that the potential strength $N(\psi(\epsilon))=\frac{1-\epsilon}{1+\epsilon}$ converges to 1 as $\epsilon\downarrow0$. The spectral radius of the corresponding matrix $A_{I\to i,J\to j}$ can be shown to be equal to

$$\rho(A) = \left(\frac{1 - \epsilon}{1 + \epsilon}\right)^{1/N}$$

⁸Additionally, the initial messages are required to be strictly positive, but this requirement is easily met and is necessary for obtaining good BP results.

which is strictly smaller than 1 if and only if $\epsilon > 0$. Hence BP converges to a unique fixed point if $\epsilon > 0$. This result is sharp, since for $\epsilon = 0$, BP simply "rotates" the messages around without changing them and hence no convergence occurs (except, obviously, if the initial messages already correspond to the fixed point of uniform messages).

V. COMPARISON WITH OTHER WORK

In this section we explore the relations of our results with previously existing work.

A. Comparison With Work Of Tatikonda And Jordan

In [14], [15], a connection is made between two seemingly different topics, namely the Sum-Product Algorithm on the one hand and the theory of Gibbs measures [23] on the other hand. The main result of [14] states that BP converges uniformly (to a unique fixed point) if the Gibbs measure on the corresponding computation tree⁹ is unique.

This is a remarkable and beautiful result; however, the question of convergence of BP is replaced by the question of uniqueness of the Gibbs measure, which is not trivial. Fortunately, sufficient conditions for the uniqueness of the Gibbs measure exist; the most well-known are Dobrushin's condition and a weaker (but more easily verifiable) condition known as Simon's condition. In combination with the main result of [14], they yield directly testable sufficient conditions for convergence of BP to a unique fixed point. For reference, we will state both results in our notation below. For details, see [14], [15] and [23]. Note that the results are valid for the case of positive factors consisting of at most two variables.

1) BP Convergence via Dobrushin's Condition: Define Dobrushin's interdependence matrix as the $N \times N$ matrix C with entries

$$C_{ij} := \sup_{x_{\partial i \setminus j}} \sup_{x_j, x_j'} \frac{1}{2} \sum_{x_i} \left| P(x_i \mid x_{\partial i \setminus j}, x_j) - P(x_i \mid x_{\partial i \setminus j}, x_j') \right|$$
(54)

for $j \in \partial i$ and 0 otherwise.

Theorem 6: For pairwise (positive) factors, BP converges to a unique fixed point if

$$\max_{i \in \mathcal{V}} \sum_{j \in \partial i} C_{ij} < 1.$$

 $\max_{i\in\mathcal{V}}\sum_{j\in\partial i}C_{ij}<1.$ Proof: For a proof sketch, see [15]. For the proof of Dobrushin's condition see chapter 8 in [23].

We can rewrite the conditional probabilities in terms of factors:

$$P(x_i \mid x_{\partial i \setminus j}, x_j) = \frac{\psi^i(x_i)\psi^{ij}(x_{ij}) \prod_{k \in \partial i \setminus j} \psi^{ik}(x_{ik})}{\sum_{x_i} \psi^i(x_i)\psi^{ij}(x_{ij}) \prod_{k \in \partial i \setminus j} \psi^{ik}(x_{ik})}.$$

Note that the complexity of the calculation of this quantity is generally exponential in the size of the neighborhood ∂j , which may prohibit practical application of Dobrushin's condition.

For the case of binary ± 1 -valued variables, some elementary algebraic manipulations yield

$$C_{ij} = \sup_{x_{\partial i \setminus j}} \frac{\sinh 2|J_{ij}|}{\cosh 2J_{ij} + \cosh 2(\theta_i + \sum_{k \in \partial i \setminus j} x_k J_{ik})}$$
$$= \frac{\tanh(|J_{ij}| - H_{ij}) + \tanh(|J_{ij}| + H_{ij})}{2}$$

with

$$H_{ij} := \inf_{x_{\partial i \setminus j}} \left| \theta_i + \sum_{k \in \partial i \setminus j} x_k J_{ik} \right|.$$

2) BP Convergence via Simon's Condition: Simon's condition is a sufficient condition for Dobrushin's condition (see proposition 8.8 in [23]). This leads to a looser, but more easily verifiable, bound:

Theorem 7: For pairwise (positive) factors, BP converges to a unique fixed point if

$$\max_{i \in \mathcal{V}} \sum_{j \in \partial i} \left(\frac{1}{2} \sup_{\alpha, \alpha'} \sup_{\beta, \beta'} \log \frac{\psi_{\alpha\beta}^{ij}}{\psi_{\alpha'\beta'}^{ij}} \right) < 1.$$

It is not difficult to show that this bound is weaker than (49). Furthermore, unlike Dobrushin's condition and Corollary 4, it does not take into account single-variable factors.

B. Comparison With Work Of Ihler et al.

In the recent and independent work [16] of Ihler et al., a methodology was used which is very similar to the one used in this work. In particular, the same local ℓ_{∞} quotient metric is used to derive sufficient conditions for BP to be a contraction. In the work presented here, the Mean Value Theorem (in the form of Lemma 1) is used in combination with a bound on the derivative in order to obtain a bound on the convergence rate K in (8). In contrast, in [16] a direct bound on the distance of two outgoing messages is derived in terms of the distance of two different products of incoming messages (equation (13) in [16]). This bound becomes relatively stronger as the distance of the products of incoming messages increases. This has the advantage that it can lead to stronger conclusions about the effect of finite message perturbations than would be possible with our bound, based on the Mean Value Theorem. However, for the question of *convergence*, the relevant limit turns out to be that of infinitesimal message perturbations, i.e., it suffices to study the derivative of the BP updates as we have done here.

In the limit of infinitesimal message perturbations, the fundamental bound (13) in [16] leads to the following measure of potential strength:

$$D(\psi^{ij}) := \tanh \left(\frac{1}{2} \left(\sup_{\alpha,\beta} \sup_{\alpha',\beta'} \log \frac{\psi^{ij}_{\alpha\beta}}{\psi^{ij}_{\alpha'\beta'}} \right) \right).$$

Using this measure, Ihler et. al derive two different conditions for convergence of BP. The first one is similar to our (49) and the second condition is equivalent to our spectral radius result (50), except that in both conditions, $N(\psi^{ij})$ is used instead of $D(\psi^{ij})$. The latter condition is formulated in [16] in terms of

⁹The computation tree is an "unwrapping" of the factor graph with respect to the Sum-Product Algorithm; specifically, the computation tree starting at variable $i \in \mathcal{V}$ consists of all paths starting at i that never backtrack.

the convergence properties of an iterative BP-like algorithm. The equivalence of this formulation with a formulation in terms of the spectral radius of a matrix can be seen from the fact that for any square matrix A, $\rho(A) < 1$ if and only if $\lim_{n\to\infty} A^n = 0$. However, our result also gives a contraction rate, unlike the iterative formulation in [16].

Thus, the results in [16] are similar to ours in the pairwise case, except for the occurrence of $D(\psi^{ij})$ instead of $N(\psi^{ij})$. It is not difficult to see that $N(\psi^{ij}) \leq D(\psi^{ij})$ for any pair factor ψ^{ij} ; indeed, for any choice of $\alpha, \beta, \gamma, \delta$:

$$\sqrt{\psi_{\alpha\gamma}\psi_{\beta\delta}}\Big/\sqrt{\psi_{\beta\gamma}\psi_{\alpha\delta}}\quad \leq\quad \Big(\sup_{\sigma\tau}\psi_{\sigma\tau}\Big)\Big/\Big(\inf_{\sigma\tau}\psi_{\sigma\tau}\Big).$$

Thus the convergence results in [16] are similar to, but weaker than the results derived in the present work.

After initial submission of this work, [17] was published, which improves upon [16] by exploiting the freedom of choice of the single-variable factors (which can be "absorbed" to an arbitrary amount by corresponding pair factors). This leads to an improved measure of potential strength, which turns out to be identical to our measure $N(\psi^{ij})$. Thus, for pairwise, strictly positive potentials, the results in [17] are equivalent to the results (49) and (50) presented here. Our most general results, Theorems 3, 4 and 5 and Corollary 4, are not present in [17].

C. Comparison With Work Of Heskes

A completely different methodology to obtain sufficient conditions for the uniqueness of the BP fixed point is used in [18]. By studying the Bethe free energy and exploiting the relationship between properties of the Bethe free energy and the BP algorithm, conclusions are drawn about the uniqueness of the BP fixed point; however, whether uniqueness of the fixed point also implies convergence of BP seems to be an open question. We state the main result of [18] in our notation below.

The following measure of potential strength is used in [18]. For $I \in \mathcal{F}$, let

$$\omega_I := \sup_{x_I} \sup_{x_I'} \left(\log \psi^I(x_I) + (|I| - 1) \log \psi^I(x_I') - \sum_{i \in I} \log \psi^I(x_{I \setminus i}', x_i) \right).$$

The potential strength is then defined as $\sigma_I := 1 - e^{-\omega_I}$.

Theorem 8: BP has a unique fixed point if there exists an "allocation matrix" X_{Ii} between factors $I \in \mathcal{F}$ and variables $i \in \mathcal{V}$ such that

- 1) $X_{Ii} \geq 0$ for all $I \in \mathcal{F}, i \in I$;
- 2) $(1 \sigma_I) \max_{i \in I} X_{Ii} + \sigma_I \sum_{i \in I} X_{Ii} \le 1$ for all $I \in \mathcal{F}$; 3) $\sum_{I \in N_i} X_{Ii} \ge |N_i| 1$ for all $i \in \mathcal{V}$.

Proof: See Theorem 8.1 in [18].

The (non)existence of such a matrix can be determined using standard linear programming techniques.

VI. NUMERICAL COMPARISON OF VARIOUS BOUNDS

In this subsection, we compare various bounds on binary pairwise graphical models, defined in (5), for various choices

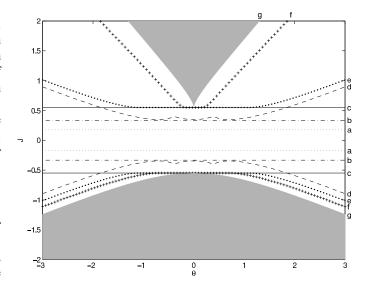


Fig. 4. Comparison of various BP convergence bounds for the fully connected N=4 binary Ising model with uniform coupling J and uniform local field θ . (a) Heskes' condition (b) Simon's condition (c) spectral radius condition (d) Dobrushin's condition (e) improved spectral radius condition for m=1(f) improved spectral radius condition for m = 5 (g) uniqueness of Gibbs' measure condition. See the main text (section VI-A) for more explanation.

of the parameters. First we study the case of a completely uniform model (i.e., full connectivity, uniform couplings and uniform local fields). Then we study nonuniform couplings J_{ij} , in the absence of local fields. Finally, we take fully random models in various parameter regimes (weak/strong local fields, strong/weak ferromagnetic/spin-glass/anti-ferromagnetic couplings).

A. Uniform Couplings, Uniform Local Field

The fully connected Ising model consisting of N binary ± 1 -valued variables with uniform couplings J and uniform local field θ is special in the sense that an exact description of the parameter region for which the Gibbs measure on the computation tree is unique, is available. Using the results of Tatikonda and Jordan, this yields a strong bound on the parameter region for which BP converges to a unique fixed point. Indeed, the corresponding computation tree is a uniform Ising model on a Cayley tree of degree N-2, for which (semi-)analytical expressions for the paramagnetic-ferromagnetic and paramagnetic-antiferromagnetic phase transition boundaries are known (see section 12.2 in [23]). Since the Gibbs measure is known to be unique in the paramagnetic phase, this gives an exact description of the (J, θ) region for which the Gibbs measure on the computation tree is unique, and hence a bound on BP convergence on the original model.

In Fig. 4 we have plotted various bounds on BP convergence in the (J,θ) plane for N=4 (other values of N yield qualitatively similar results). The gray area (g) marks regions where the Gibbs measure on the computation tree is not unique; in the white area, the Gibbs measure is unique and hence BP is guaranteed to converge. Note that this bound is only available due to the high symmetry of the model. In [24] it is shown that parallel BP does not converge in the lower

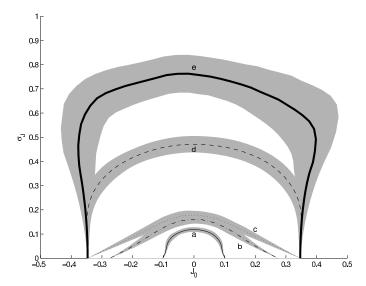


Fig. 5. Comparison of various bounds for BP convergence for toroidal Ising model of size 10×10 with normally distributed couplings $J_{ij} \sim \mathcal{N}(J_0, \sigma_J^2)$ and zero local fields. (a) Heskes' condition (b) Dobrushin's condition (c) ℓ_1 -norm condition (d) spectral radius condition (e) empirical convergence boundary. See the main text (section VI-B) for more explanation.

(anti-ferromagnetic) gray region. In the upper (ferromagnetic) region on the other hand, parallel BP does converge, but it may be that the fixed point is no longer unique.

The various lines correspond to different sufficient conditions for BP convergence; the regions enclosed by two lines of the same type (i.e., the inner regions for which J is small) mark the regions of guaranteed convergence. The lightly dotted lines (a) correspond with Heskes' condition, Theorem 8. The dash-dotted lines (b) correspond with Simon's condition, Theorem 7. The dashed lines (d) correspond with Dobrushin's condition (Theorem 6), which is seen to improve upon Simon's condition for $\theta \neq 0$, but is nowhere sharp. The solid lines (c) correspond with the spectral radius condition Corollary 3 (which coincides with the ℓ_1 -norm based condition Corollary 2 in this case and is also equivalent to the result of [16]), which is independent of θ but is actually sharp for $\theta = 0$. The heavily dotted lines (e) correspond to Corollary 4 with m=1, the +-shaped lines (f) to the same Corollary with m=5. Both (e) and (f) are seen to coincide with (c) for small θ , but improve for large θ .

We conclude that the presence of local fields makes it more difficult to obtain sharp bounds on BP convergence; only Dobrushin's condition (Theorem 6) and Corollary 4 take into account local fields. Furthermore, in this case, our result Corollary 4 is stronger than the other bounds. Note that the calculation of Dobrushin's condition is exponential in the number of variables N, whereas the time complexity of our bound is polynomial in N. Similar results are obtained for higher values of N.

B. Nonuniform Couplings, Zero Local Fields

We have investigated in more detail the influence of the distribution of the couplings J_{ij} , in the absence of local fields, and have also compared with the empirical convergence behavior of BP. We have taken a binary Ising model on a rectangular

toroidal grid (i.e., with periodic boundary conditions) of size 10×10 . The couplings were random independent normally distributed nearest-neighbor couplings $J_{ij} \sim \mathcal{N}(J_0, \sigma_J^2)$, the local fields were $\theta_i=0$. Let (r_J,ϕ_J) be the polar coordinates corresponding to the Cartesian coordinates (J_0, σ_J) . For various angles $\phi_J \in [0, \pi]$, we have determined the critical radius r_J for each bound. The results have been averaged over 40 instances of the model and can be found in Fig. 5. The lines correspond to the mean bounds, the gray areas are "error bars" of one standard deviation. The inner area (for which the couplings are small) bounded by each line means "convergence", either guaranteed or empirical (thus the larger the enclosed area, the tighter the bound). From bottom to top: the thin solid line (a) corresponds with Heskes' result (Theorem 8), the dash-dotted line (b) with Dobrushin's condition (Theorem 6), the dotted line (c) corresponds with the ℓ_1 -norm based condition Corollary 2, the dashed line (d) with the spectral radius condition Corollary 3 and the thick solid line (e) with the empirical convergence behavior of BP.

We conclude from Fig. 5 that the spectral radius condition improves upon the ℓ_1 -norm based condition for nonuniform couplings and that the improvement can be quite substantial. For uniform couplings (and zero local fields), both conditions coincide and it can be proved that they are sharp [25].

C. Fully Random Models

Finally, we have considered fully connected binary pairwise graphical models with completely random couplings and local fields (in various parameter regimes). We drew random couplings and local fields as follows: first, we drew i.i.d. random parameters $J_0, \sigma_J, \theta_0, \sigma_\theta$ from a normal distribution with mean 0 and variance 1. Then, for each variable i we independently drew a local field parameter $\theta_i \sim \mathcal{N}(\theta_0, \sigma_\theta^2)$, and for each pair $\{i, j\}$ we independently drew a coupling parameter $J_{ij} \sim \mathcal{N}(J_0, \sigma_J^2)$.

For the resulting graphical model, we have verified whether various sufficient conditions for BP convergence hold. If condition A holds whereas condition B does not hold, we say that A wins from B. We have counted for each ordered pair (A,B) of conditions how often A wins from B. The results (for 50000 random models consisting of N=4,8 variables) can be found in Table I: the number at row A, column B is the number of trials for which bound A wins from bound B. On the diagonal A=B is the total number of trials for which bound A predicts convergence. Theorem 6 is due to [15], Corollary 3 was first published (for the binary case) in [16] and Theorem 8 is due to [18].

Our result Corollary 4 (for m=1) outperforms the other bounds in each trial. For other values of N, we obtain similar results.

VII. DISCUSSION

In this paper we have derived sufficient conditions for convergence of BP to a unique fixed point. Our conditions are directly applicable to arbitrary graphical models with discrete variables and nonnegative factors. This is in contrast with the sufficient conditions of Tatikonda and Jordan and with

Comparison of bounds (50000 trials, for N=4 and N=8)

N = 4	Th. 6	Cor. 3	Th. 8	Cor. 4
Th. 6, [15]	(5779)	170	3564	0
Cor. 3, [16]	10849	(16458)	13905	0
Th. 8, [18]	338	0	(2553)	0
Cor. 4, $m = 1$, this work	13820	3141	17046	(19599)

N = 8	Th. 6	Cor. 3	Th. 8	Cor. 4
Th. 6, [15]	(668)	39	597	0
Cor. 3, [16]	507	(1136)	1065	0
Th. 8, [18]	0	0	(71)	0
Cor. 4, $m=1$, this work	972	504	1569	(1640)

the results of Ihler, Fisher and Willsky, which were only formulated for pairwise, positive factors. We have shown cases where our results are stronger than previously known sufficient

Our numerical experiments lead us to conjecture that Corollary 4 is stronger than the other bounds. We have no proof for this conjecture at the moment, apart from the obvious fact that Corollary 3 is weaker than Corollary 4. To prove that Corollary 4 is stronger than Theorem 6 seems subtle, since it is generally not the case that $\rho(A) \leq \|C\|_{\infty}$, although it seems that the weaker relation $\|C\|_{\infty} < 1 \implies \rho(A) < 1$ does hold in general. The relation with the condition in Theorem 8 is not evident as well.

In the binary pairwise case, it turned out to be possible to derive sufficient conditions that take into account local evidence (Corollary 4). In the general case, such an improvement is possible in principle but seems to be more involved. The resulting optimization problem (essentially (43) with additional assumptions on h) looks difficult in general. If the variables' cardinalities and connectivities are small, the resulting optimization problem can be solved, but writing down a general solution does not appear to be trivial. The question of finding an efficient solution in the general case is left for future investigation.

The work reported here raises new questions, some of which have been (partially) answered elsewhere after the initial submission of this paper. The influence of damping the BP update equations has been considered for the binary pairwise case in [25], where it was shown that damping has the most effect for anti-ferromagnetic interactions. Furthermore, it has been proved in [25] that the bounds for BP convergence derived in the present work are sharp in the case of binary variables with (anti-)ferromagnetic pairwise interactions and zero local fields, as suggested by Fig. 5. An extension of the results towards sequential update schemes has been given in [26]. Likewise, in [24] it is shown that Dobrushin's condition is also valid for sequential BP.

APPENDIX A Generalizing the ℓ_1 -norm

Let $(V_i, \|\cdot\|_i)$ be a finite collection of normed vector spaces and let $V = \bigoplus_i V_i$ be the direct sum of the V_i 's. The function $\|\cdot\|:V\to\mathbb{R}$ defined by

$$||v|| := \sum_{i} ||v_{i}||_{i}$$
 (A.55)

is a norm on V, as one easily checks. Let $A: V \to V$ be a linear mapping with "blocks" $A_{ij}: V_j \to V_i$ defined by

$$\forall v_j \in V_j : Av_j = \sum_i A_{ij} v_j, A_{ij} v_j \in V_i$$

for all j.

Theorem 9: The matrix norm of A induced by the vector norm $\|\cdot\|$ is given by:

$$||A|| = \max_{j} \sum_{i} ||A_{ij}||_{i}^{j}$$
 (A.56)

where

$$\begin{split} \|A_{ij}\|_i^j := \sup_{\substack{x \in V_j, \\ \|x\|_j \leq 1}} \|A_{ij}x\|_i \;. \\ \textit{Proof:} \;\; \text{Let} \; v_k \in V_k \; \text{such that} \;\; \|v_k\|_k = 1. \; \text{Then} \end{split}$$

$$||Av_k|| = \left\| \sum_{i} A_{ik} v_k \right\| = \sum_{i} ||A_{ik} v_k||_{i}$$

$$\leq \sum_{i} ||A_{ik}||_{i}^{k} \leq \max_{j} \sum_{i} ||A_{ij}||_{i}^{j}.$$

Now let $v \in V$ such that ||v|| = 1. Then v can be written as the convex combination $v = \sum_k \|v_k\|_k \tilde{v}_k$, where

$$\tilde{v}_k := \begin{cases} \frac{v_k}{\|v_k\|_k} & \text{if } v_k \neq 0\\ 0 & \text{if } v_k = 0. \end{cases}$$

Hence:

$$||Av|| = \left\| \sum_{k} ||v_{k}||_{k} A\tilde{v}_{k} \right\| \leq \sum_{k} ||v_{k}||_{k} ||A\tilde{v}_{k}||$$

$$\leq \max_{j} \sum_{i} ||A_{ij}||_{i}^{j}.$$

It is evident that this value is also achieved for some $v \in V$ with ||v|| = 1.

An illustrative example is obtained by considering $V = \mathbb{R}^N$ to be the direct sum of N copies of \mathbb{R} with the absolute value as norm; then the norm (A.55) on \mathbb{R}^N is simply the ℓ_1 -norm and the induced matrix norm (A.56) reduces to (9).

Suppose that each V_i has a linear subspace W_i . We can consider the quotient spaces V_i/W_i with quotient norms $\|\bar{\cdot}\|_i$. The direct sum $W := \bigoplus_i W_i$ is itself a subspace of V, yielding a quotient space V/W. For $v \in V$ we have $\overline{v} = \sum_i \overline{v_i}$ and hence $V/W = \bigoplus_i (V_i/W_i)$. The quotient norm on V/Wis simply the sum of the quotient norms on the V_i/W_i :

$$\|\overline{v}\| := \inf_{w \in W} \|v + w\| = \inf_{w \in W} \sum_{i} \|v_{i} + w_{i}\|_{i}$$

$$= \sum_{i} \inf_{w_{i} \in W_{i}} \|v_{i} + w_{i}\|_{i} = \sum_{i} \|\overline{v_{i}}\|_{i}.$$
(A.57)

Let $A: V \to V$ be a linear mapping such that $AW \subseteq W$. Then A induces a linear $\overline{A}: V/W \to V/W$; since $A_{ij}W_i \subseteq$ W_i , each block $A_{ij}: V_j \to V_i$ induces a linear $\overline{A_{ij}}: V_j/W_j \to V_i$ V_i/W_i , and \overline{A} can be regarded as consisting of the blocks $\overline{A_{ij}}$.

Corollary 6: The matrix norm of $\overline{A}:V/W\to V/W$ induced by the quotient norm $\|\overline{\cdot}\|$ on V/W is:

$$\|\overline{A}\| = \max_{j} \sum_{i} \|\overline{A_{ij}}\|_{i}^{j}$$
 (A.58)

where

$$\|\overline{A_{ij}}\|_{i}^{j} = \sup_{\substack{x \in V_{j}, \\ \|x\|_{i} \le 1}} \|\overline{A_{ij}x}\|_{i}.$$
 (A.59)

Proof: We can directly apply the previous Theorem to the quotient spaces to obtain (A.58); because

$$\{\overline{x} \in V_j/W_j: \|\overline{x}\|_j \le 1\} = \overline{\{x \in V_j: \|x\|_j \le 1\}},$$

we have:

$$\left\| \overline{A_{ij}} \right\|_{i}^{j} := \sup_{\substack{\overline{x} \in V_{j}/W_{j} \\ \|\overline{x}\|_{i} \leq 1}} \left\| \overline{A_{ij}} \overline{x} \right\|_{i} = \sup_{\substack{x \in V_{j} \\ \|x\|_{i} \leq 1}} \left\| \overline{A_{ij}} x \right\|_{i}.$$

For a linear $A:V\to V$ such that $AW\subseteq W$, we define the matrix |A| with entries $|A|_{ij}:=\left\|\overline{A_{ij}}\right\|_i^j$. Let A,B be two such linear mappings; then

$$|AB|_{ij} = \left\| \overline{(AB)_{ij}} \right\|_{i}^{j} = \left\| \sum_{k} \overline{A_{ik} B_{kj}} \right\|_{i}^{j}$$

$$\leq \sum_{k} \left\| \overline{A_{ik} \overline{B_{kj}}} \right\|_{i}^{j} \leq \sum_{k} \left\| \overline{A_{ik}} \right\|_{i}^{k} \left\| \overline{B_{kj}} \right\|_{k}^{j}$$

$$= \sum_{k} |A|_{ik} |B|_{kj}$$

hence $|AB| \le |A| |B|$. Note that $|||A|||_1 = ||\overline{A}||$. We can generalize Theorem 2:

Theorem 10: Let $f:V\to V$ be differentiable and suppose that it satisfies (30). Suppose further that $|f'(v)|\leq A$ for some matrix A_{ij} (which does not depend on v) with $\rho(A)<1$. Then for any $\overline{v}\in V/W$, the sequence $\overline{v},\overline{f}(\overline{v}),\overline{f}^2(\overline{v}),\ldots$ obtained by iterating \overline{f} converges to a unique fixed point \overline{v}_{∞} .

Proof: Using the chain rule, we have for any $n=1,2,\ldots$ and any $v\in V$:

$$\left\| (\overline{f}^n)'(\overline{v}) \right\| = \left\| \overline{(f^n)'(v)} \right\| = \left\| \overline{\prod_{i=1}^n f'(f^{i-1}(v))} \right\|$$

$$= \left\| \left| \prod_{i=1}^n f'(f^{i-1}(v)) \right| \right\|_1 \le \left\| \prod_{i=1}^n \left| f'(f^{i-1}(v)) \right| \right\|_1$$

$$\le \left\| \prod_{i=1}^n A \right\|_1 = \left\| A^n \right\|_1.$$

By the Gelfand Spectral Radius Theorem, $\|A^n\|_1^{-1/n} \to \rho(A)$ for $n \to \infty$. Choose $\epsilon > 0$ such that $\rho(A) + \epsilon < 1$. For some N, $\|A^N\|_1 \le (\rho(A) + \epsilon)^N < 1$. Hence $\|(\overline{f}^N)'(\overline{v})\| < 1$ for all $\overline{v} \in V/W$. By Lemma 2, \overline{f}^N is a contraction with respect to the quotient norm on V/W. Now apply Lemma 3.

APPENDIX B

PROOF THAT (43) EQUALS (44)

Let $\tilde{\psi}_{\beta\gamma}$ be a matrix of positive numbers. Let

$$\mathcal{H} := \{ h : h_{\beta\gamma} \ge 0, \sum_{\beta,\gamma} h_{\beta\gamma} = 1 \}.$$

Define the function $g: \mathcal{H} \to \mathbb{R}$ by

$$g(h) = \sum_{\beta} \left| \sum_{\gamma} h_{\beta\gamma} \left(\frac{\tilde{\psi}_{\beta\gamma}}{\sum_{\beta} \sum_{\gamma} \tilde{\psi}_{\beta\gamma} h_{\beta\gamma}} - 1 \right) \right|.$$

Theorem 11:

$$\sup_{h \in \mathcal{H}} g(h) = 2 \sup_{\beta \neq \beta'} \sup_{\gamma, \gamma'} \tanh \left(\frac{1}{4} \log \frac{\tilde{\psi}_{\beta \gamma}}{\tilde{\psi}_{\beta' \gamma'}} \right).$$

Proof: First note that we can assume without loss of generality that all $\tilde{\psi}_{\beta\gamma}$ are different, because of continuity. Define

$$\begin{split} \tilde{\psi}_{-} &:= \inf_{\beta \gamma} \tilde{\psi}_{\beta \gamma}, \qquad \tilde{\psi}_{+} := \sup_{\beta \gamma} \tilde{\psi}_{\beta \gamma}, \\ X &:= [\tilde{\psi}_{-}, \tilde{\psi}_{+}], \qquad X' := X \setminus \{\tilde{\psi}_{\beta \gamma} : \beta, \gamma\}. \end{split}$$

For $\Psi \in X$, define

$$\mathcal{H}_{\Psi} := \{ h \in \mathcal{H} : \sum_{\beta,\gamma} \tilde{\psi}_{\beta\gamma} h_{\beta\gamma} = \Psi \},$$

which is evidently a closed convex set. The function

$$g_{\Psi}: \mathcal{H}_{\Psi}
ightarrow \mathbb{R}: h \mapsto \sum_{eta} \left| \sum_{\gamma} h_{eta \gamma} \left(rac{ ilde{\psi}_{eta \gamma}}{\Psi} - 1
ight)
ight|$$

obtained by restricting g to \mathcal{H}_{Ψ} is convex. Hence it achieves its maximum on an extremal point of its domain.

Define

$$\mathcal{H}_2 := \{ h \in \mathcal{H} : \#\{(\beta, \gamma) : h_{\beta\gamma} > 0 \} = 2 \}$$

as those $h \in \mathcal{H}$ with exactly two nonzero components. For $h \in \mathcal{H}_2$, define $\tilde{\psi}_-(h) := \inf\{\tilde{\psi}_{\beta\gamma} : h_{\beta\gamma} \neq 0\}$ and $\tilde{\psi}_+(h) := \sup\{\tilde{\psi}_{\beta\gamma} : h_{\beta\gamma} \neq 0\}$. Because of continuity, we can restrict ourselves to the $\Psi \in X'$, in which case the extremal points of \mathcal{H}_Ψ are precisely $\mathcal{H}_\Psi^* = \mathcal{H}_\Psi \cap \mathcal{H}_2$ (i.e., the extremal points have exactly two nonzero components).

Now

$$\sup_{h \in \mathcal{H}} g(h) = \sup_{\Psi \in X} \sup_{h \in \mathcal{H}_{\Psi}} g_{\Psi}(h)$$

$$= \sup_{\Psi \in X'} \sup_{h \in \mathcal{H}_{\Psi}^*} g_{\Psi}(h)$$

$$= \sup_{h \in \mathcal{H}_2} \sup_{\tilde{\psi}_{-}(h) \le \Psi \le \tilde{\psi}_{+}(h)} g_{\Psi}(h)$$

$$= \sup_{h \in \mathcal{H}_2} g(h).$$

For those $h \in \mathcal{H}_2$ with components with different β , we can use the Lemma below. The $h \in \mathcal{H}_2$ with components with equal β are suboptimal, since the two contributions in the sum over γ in g(h) have opposite sign. Hence

$$\sup_{h \in \mathcal{H}_2} g(h) = 2 \sup_{\beta \neq \beta'} \sup_{\gamma, \gamma'} \tanh \left(\frac{1}{4} \log \frac{\tilde{\psi}_{\beta \gamma}}{\tilde{\psi}_{\beta' \gamma'}} \right).$$

Lemma 4: Let 0 < a < b. Then

$$\sup_{\substack{\eta \in (0,1)^2 \\ \eta_1 + \eta_2 = 1}} \eta_1 \left| \frac{a}{\eta_1 a + \eta_2 b} - 1 \right| + \eta_2 \left| \frac{b}{\eta_1 a + \eta_2 b} - 1 \right|$$

$$= 2 \tanh \left(\frac{1}{4} \log \frac{b}{a}\right) = 2 \frac{\sqrt{b} - \sqrt{a}}{\sqrt{b} + \sqrt{a}}.$$
Proof: Elementary. The easiest way to see this is to

reparameterize $\eta = (\frac{e^{\nu}}{2\cosh \nu}, \frac{e^{-\nu}}{2\cosh \nu})$ with $\nu \in (-\infty, \infty)$. \square

ACKNOWLEDGMENT

We thank Martijn Leisink for stimulating discussions and the reviewers for their critique, which has led to a considerable improvement of the initial version of the manuscript.

REFERENCES

- [1] J. M. Mooij and H. J. Kappen, "Sufficient conditions for convergence of loopy belief propagation," in Proceedings of the 21st Annual Conference on Uncertainty in Artificial Intelligence (UAI-05), F. Bacchus and T. Jaakkola, Eds. Corvallis, Oregon: AUAI Press, 2005, pp. 396–403.
- [2] F. R. Kschischang, B. J. Frey, and H.-A. Loeliger, "Factor graphs and the sum-product algorithm," IEEE Trans. Inform. Theory, vol. 47, no. 2, pp. 498-519, Feb. 2001.
- [3] R. J. McEliece, D. J. C. MacKay, and J.-F. Cheng, "Turbo decoding as an instance of Pearl's 'belief propagation' algorithm," IEEE J. Select. Areas Commun., vol. 16, pp. 140-152, Feb. 1998.
- [4] A. Braunstein and R. Zecchina, "Survey propagation as local equilibrium equations," Journal of Statistical Mechanics: Theory and Experiment, vol. 2004, no. 06, p. P06007, 2004. [Online]. Available: http://stacks.iop.org/1742-5468/2004/P06007
- [5] J. Sun, N.-N. Zheng, and H.-Y. Shum, "Stereo matching using belief propagation," IEEE Transactions on Pattern Analysis and Machine Intelligence, vol. 25, no. 7, pp. 787-800, 2003.
- [6] K. Tanaka, "Statistical-mechanical approach to image processing," Journal of Physics A: Mathematical and General, vol. 35, no. 37, pp. R81-R150, 2002. [Online]. Available: http://stacks.iop.org/0305-4470/ 35/R81
- [7] T. Heskes, C. A. Albers, and H. J. Kappen, "Approximate inference and constrained optimization," in Proceedings of the 19th Annual Conference on Uncertainty in Artificial Intelligence (UAI-03). San Francisco, CA: Morgan Kaufmann Publishers, 2003, pp. 313–320.
- [8] J. S. Yedidia, W. T. Freeman, and Y. Weiss, "Constructing free-energy approximations and generalized belief propagation algorithms," IEEE Transactions on Information Theory, vol. 51, no. 7, pp. 2282-2312,
- [9] T. Minka, "Expectation propagation for approximate Bayesian inference," in Proceedings of the 17th Annual Conference on Uncertainty in Artificial Intelligence (UAI-01). San Francisco, CA: Morgan Kaufmann Publishers, 2001, pp. 362–369.
- [10] M. Opper and O. Winter, "Expectation consistent approximate inference," Journal of Machine Learning Research, vol. 6, pp. 2177-2204, Dec. 2005.
- [11] Y. Weiss and W. T. Freeman, "On the optimality of solutions of the max-product belief propagation algorithm in arbitrary graphs," IEEE Transactions on Information Theory, vol. 47, no. 2, pp. 736-744, Feb.
- [12] A. Braunstein, M. Mézard, and R. Zecchina, "Survey propagation: an algorithm for satisfiability," Random Structures and Algorithms, vol. 27, no. 2, pp. 201-226, 2005. [Online]. Available: http: //dx.doi.org/10.1002/rsa.20057
- [13] W. Wiegerinck and T. Heskes, "Fractional belief propagation," in Advances in Neural Information Processing Systems 15, S. T. S. Becker and K. Obermayer, Eds. Cambridge, MA: MIT Press, 2003, pp. 438-445.

- [14] S. C. Tatikonda and M. I. Jordan, "Loopy belief propogation and Gibbs measures," in Proc. of the 18th Annual Conf. on Uncertainty in Artificial Intelligence (UAI-02). San Francisco, CA: Morgan Kaufmann Publishers, 2002, pp. 493-500.
- [15] S. C. Tatikonda, "Convergence of the sum-product algorithm," in Proceedings 2003 IEEE Information Theory Workshop, 2003.
- A. T. Ihler, J. W. Fisher, and A. S. Willsky, "Message errors in belief propagation," in Advances in Neural Information Processing Systems 17 (NIPS*2004), L. K. Saul, Y. Weiss, and L. Bottou, Eds. Cambridge, MA: MIT Press, 2005, pp. 609-616.
- -, "Loopy belief propagation: Convergence and effects of message errors," Journal of Machine Learning Research, vol. 6, pp. 905-936, 2005
- [18] T. Heskes, "On the uniqueness of loopy belief propagation fixed points," Neural Computation, vol. 16, no. 11, pp. 2379-2413, Nov. 2004.
- [19] Y. Weiss, "Correctness of local probability propagation in graphical models with loops," Neur. Comp., vol. 12, pp. 1-41, 2000.
- [20] J. Dieudonné, Foundations of Modern Analysis, ser. Pure and Applied Mathematics. New York: Academic Press, 1969, vol. 10-I.
- [21] E. Deutsch, "On matrix norms and logarithmic norms," Numerische Mathematik, vol. 24, no. 1, pp. 49-51, Feb. 1975.
- [22] A. Cima, A. van den Essen, A. Gasull, E. Hubbers, and F. Manosas, "A polynomial counterexample to the Markus-Yamabe conjecture," Advances in Mathematics, vol. 131, no. 2, pp. 453-457, Nov. 1997.
- [23] H.-O. Georgii, Gibbs Measures and Phase Transitions. Berlin: Walter de Gruyter, 1988.
- N. Taga and S. Mase, "On the convergence of loopy belief propagation algorithm for different update rules," IEICE Trans. Fundamentals, vol. E89-A, no. 2, pp. 575-582, Feb. 2006.
- [25] J. M. Mooij and H. J. Kappen, "On the properties of the Bethe approximation and loopy belief propagation on binary networks, Journal of Statistical Mechanics: Theory and Experiment, vol. 2005, no. 11, p. P11012, 2005. [Online]. Available: http://stacks.iop.org/ 1742-5468/2005/P11012
- [26] G. Elidan, I. McGraw, and D. Koller, "Residual belief propagation: Informed scheduling for asynchronous message passing," in *Proceedings* of the 22nd Annual Conference on Uncertainty in Artificial Intelligence (UAI-06), Boston, Massachussetts, July 2006.

BIOGRAPHIES OF THE AUTHORS

Joris M. Mooij has been a Ph.D. student at the Radboud University Nijmegen, The Netherlands since 2003 and he expects to obtain his degree in 2007. He received his M.Sc. in physics (cum laude) at the same university in 2002 and one year later he received his M.Sc. in mathematics (*cum laude*). His research interests include algorithms for approximate probabilistic inference, their applications in machine learning and artificial intelligence, and their relationship with methods used in statistical physics. He is author of 10 publications.

Hilbert J. Kappen studied particle physics in Groningen, The Netherlands and completed his Ph.D. in this field in 1987 at the Rockefeller University in New York. From 1987 until 1989 he worked as a scientist at the Philips Research Laboratories in Eindhoven, The Netherlands. Since 1989, he is conducting research on neural networks at the laboratory for biophysics of the University of Nijmegen (now known as the Radboud University Nijmegen), The Netherlands. Since 1997 he is associate professor and since 2004 full professor at this university. He is author of approximately 120 publications.