

# Convergence Analysis of Reweighted Sum-Product Algorithms

Tanya Roosta, *Student Member, IEEE*, Martin J. Wainwright, *Member, IEEE*, and  
Shankar Sastry, *Member, IEEE*,

## Abstract

Markov random fields are designed to represent structured dependencies among large collections of random variables, and are well-suited to capture the structure of real-world signals. Many fundamental tasks in signal processing (e.g., smoothing, denoising, segmentation etc.) require efficient methods for computing (approximate) marginal probabilities over subsets of nodes in the graph. The marginalization problem, though solvable in linear time for graphs without cycles, is computationally intractable for general graphs with cycles. This intractability motivates the use of approximate “message-passing” algorithms. This paper studies the convergence and stability properties of the family of *reweighted sum-product algorithms*, a generalization of the widely-used sum-product or belief propagation algorithm, in which messages are adjusted with graph-dependent weights. For pairwise Markov random fields, we derive various conditions that are sufficient to ensure convergence, and also provide bounds on the geometric convergence rates. When specialized to the ordinary sum-product algorithm, these results provide strengthening of previous analyses. We prove that some of our conditions are necessary and sufficient for subclasses of homogeneous models, but not for general models. The experimental simulations on various classes of graphs validate our theoretical results.

Manuscript received on July 29, 2007; revised on January 4, 2008. Martin J. Wainwright was supported by NSF grants DMS-0528488 and CAREER CCF-0545862. Tanya Roosta was supported by TRUST (The Team for Research in Ubiquitous Secure Technology), which receives support NSF grant CCF-0424422, and the following organizations: Cisco, ESCHER, HP, IBM, Intel, Microsoft, ORNL, Qualcomm, Pirelli, Sun and Symantec.

Tanya Roosta is with the Department of Electrical and Computer Engineering, University of California at Berkeley, Berkeley, CA, 94720 USA e-mail: (roosta@eecs.berkeley.edu).

Martin J. Wainwright is with the Department of Electrical and Computer Engineering and Department of Statistics, University of California at Berkeley, Berkeley, CA, 94720 USA e-mail: (wainwrig@eecs.berkeley.edu).

Shankar Sastry is with the Department of Electrical and Computer Engineering, University of California at Berkeley, Berkeley, CA, 94720 USA e-mail: (sastry@eecs.berkeley.edu).

## Index Terms

Markov random fields; graphical models; belief propagation; sum-product algorithm; convergence analysis; approximate marginalization.

## I. INTRODUCTION

Graphical models provide a powerful framework for capturing the complex statistical dependencies exhibited by real-world signals. Accordingly, they play a central role in many disciplines, including statistical signal processing [1], [2], image processing [3], statistical machine learning [4], and computational biology. A core problem common to applications in all of these domains is the *marginalization problem*—namely, to compute marginal distributions over local subsets of random variables. For graphical models without cycles, including Markov chains and trees (see Figure 1(a) and (b)), the marginalization problem is exactly solvable in linear-time via the sum-product algorithm, which operates in a distributed manner by passing “messages” between nodes in the graph. This sum-product framework includes many well-known algorithms as special cases, among them the  $\alpha$ - $\beta$  or forward-backward algorithm for Markov chains, the peeling algorithm in bioinformatics, and the Kalman filter; see the review articles [1], [2], [5] for further background on the sum-product algorithm and its uses in signal processing.

Although Markov chains/trees are tremendously useful, many classes of real-world signals are best captured by graphical models with cycles. (For instance, the lattice or grid-structured model in Figure 1(c) is widely used in computer vision and statistical image processing.) At least in principle, the nodes in any such graph with cycles can be clustered into “supernodes”, thereby converting the original graph into junction tree form [6], to which the sum-product algorithm can be applied to obtain exact results. However, the cluster sizes required by this junction tree formulation—and hence the computational complexity of the sum-product algorithm—grow *exponentially* in the treewidth of the graph. For many classes of graphs, among them the lattice model in Figure 1(c), the treewidth grows in an unbounded manner with graph size, so that the junction tree approach rapidly becomes infeasible. Indeed, the marginalization problem is known to be computationally intractable for general graphical models.

This difficulty motivates the use of efficient algorithms for computing *approximations* to the

marginal probabilities. In fact, one of the most successful approximate methods is based on applying the sum-product updates to the graphs with cycles. Convergence and correctness, though guaranteed for tree-structured graphs, are no longer ensured when the underlying graph has cycles. Nonetheless, this “loopy” form of the sum-product algorithm has proven very successful in many applications [1]–[3], [5]. However, there remain a variety of theoretical questions concerning the use of sum-product and related message-passing algorithms for approximate marginalization. It is well known that the standard form of sum-product message-passing is not guaranteed to converge, and in fact may have multiple fixed points in certain regimes. Recent work has shed some light on the fixed points and convergence properties of the ordinary sum-product algorithm. Yedidia et al. [7] showed that sum-product fixed points correspond to local minima of an optimization problem known as the Bethe variational principle. Tatikonda and Jordan [8] established an elegant connection between the convergence of the ordinary sum-product algorithm and the uniqueness of Gibbs measures on the associated computation tree, and provided several sufficient conditions for convergence. Wainwright et al. [9] showed that the sum-product algorithm can be understood as seeking an alternative reparameterization of the distribution, and used this to characterize the error in the approximation. Heskes [10] discussed convergence and its relation to stability properties of the Bethe variational problem. Other researchers [11], [12] have used contraction arguments to provide sharper sufficient conditions for convergence of the standard sum-product algorithm. Finally, several groups [13]–[15] have proposed modified algorithms for solving the Bethe variational problem with convergence guarantees, albeit at the price of increased complexity.

In this paper, we study the broader class of *reweighted sum-product* algorithms [16]–[19], including the ordinary sum-product algorithm as a special case, in which messages are adjusted by edge-based weights determined by the graph structure. For suitable choices of these weights, the reweighted sum-product algorithm is known to have a unique fixed point for any graph and any interaction potentials [16]. An additional desirable property of reweighted sum-product is that the message-passing updates tend to be more stable, as confirmed by experimental investigation [16], [18], [19]. This algorithmic stability should be contrasted with the ordinary sum-product algorithm, which can be highly unstable due to phase transitions in the Bethe variational problem [7], [8]. Despite these encouraging empirical results, current theoretical understanding of the stability and

convergence properties of reweighted message-passing remains incomplete.

The main contributions of this paper are a number of theoretical results characterizing the convergence properties of reweighted sum-product algorithms, including the ordinary sum-product updates as a special case. Beginning with the simple case of homogeneous binary models, we provide sharp guarantees for convergence, and prove that there always exists a choice of edge weights for which the associated reweighted sum-product algorithm converges. We then analyze more general inhomogeneous models, both for binary variables and the general multinomial model, and provide sufficient conditions for convergence of reweighted algorithms. Relative to the bulk of past work, a notable feature of our analysis is that it incorporates the benefits of making observations, whether partial or noisy, of the underlying random variables in the Markov random field to which message-passing is applied. Intuitively, the convergence of message-passing algorithms should be function of both the strength of the interactions between random variables, as well as the local observations, which tend to counteract the interaction terms. Indeed, when specialized to the ordinary sum-product algorithm, our results provide a strengthening of the best previously known convergence guarantees for sum-product [8], [10]–[12]. As pointed out after initial submission of this paper, independent work by Mooij and Kappen [20] yields a similar refinement for the case of the ordinary sum-product, and binary variables; the result given here applies more generally to reweighted sum-product, as well as higher-order state spaces. As we show empirically, the benefits of incorporating observations into convergence analysis can be substantial, particularly in the regimes most relevant to applications.

The remainder of this paper is organized as follows. In Section II, we provide basic background on graphical models (with cycles), and the class of reweighted sum-product algorithms that we study. Section III provides convergence analysis for binary models, which we then extend to general discrete models in Section IV. In Section V, we describe experimental results that illustrate our findings, and we conclude in Section VI.

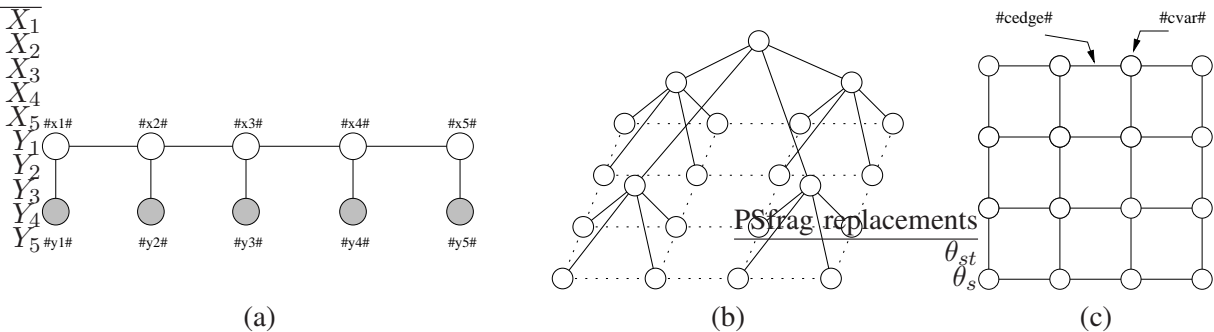
## II. BACKGROUND

In this section we provide some background on Markov random fields, and message-passing algorithms, including the reweighted sum-product that is the focus of this paper.

### A. Graphical models

Undirected graphical models, also known as Markov random fields, are based on associating a collection of random variables  $X = \{X_1, \dots, X_n\}$  with the vertices of a graph. More precisely, an undirected graph  $G = (V, E)$ , where  $V = \{1, \dots, n\}$  are vertices, and  $E \subset V \times V$  are edges joining pairs of vertices. Each random variable  $X_i$  is associated with node  $i \in V$ , and the edges in the graph (or more precisely, the absences of edges) encode Markov properties of the random vector  $X$ . These Markov properties are captured by a particular factorization of the probability distribution  $p$  of the random vector  $X$ , which is guaranteed to break into a product of local functions on the cliques of the graph. (A graph clique is a subset  $C$  of vertices that are all joined by edges.)

PSfrag replacements



**Fig. 1.** Examples of graphical models. (a) A hidden Markov chain model (with noisy observations  $Y_s$  of each hidden  $X_s$ ), on which the marginalization problem is solved by the forward-backward algorithm. (b) Marginalization can also be performed in linear time on a tree (graph without cycles), as widely used in multi-resolution signal processing [1]. (c) A lattice-based model frequently used in image processing [21], for which the marginalization problem is intractable in general.

In this paper, we focus on discrete (multinomial) random variables  $X_s \in \mathcal{X} := \{0, 1, \dots, m-1\}$  with distribution specified according to a pairwise Markov random field. Any such model has a probability distribution of the form

$$p(x; \theta) \propto \exp \left\{ \sum_{s \in V} \theta_s(x_s) + \sum_{(s,t) \in E} \theta_{st}(x_s, x_t) \right\}. \quad (1)$$

Here the quantities  $\theta_s$  and  $\theta_{st}$  are *potential functions* that depend only on the value  $X_s = x_s$ , and the pair values  $(X_s, X_t) = (x_s, x_t)$  respectively. Otherwise stated, each singleton potential  $\theta_s$  is a real-valued function of  $\mathcal{X} = \{0, 1, \dots, m\}$ , whose values can be represented as an  $m$ -vector, whereas each edge potential  $\theta_{st}$  is a real-valued mapping on the Cartesian product  $\mathcal{X} \times \mathcal{X}$ , whose

values can be represented as a  $m \times m$  matrix. With this set-up, the *marginalization problem* is to compute the singleton marginal distributions  $p(x_s; \theta) = \sum_{x_t, t \neq s} p(x; \theta)$ , and possibly higher-order marginal distributions (e.g.,  $p(x_s, x_t; \theta)$ ) as well. Note that if viewed naively, the summation defining  $p(x_s; \theta)$  involves an exponentially growing number of terms ( $m^{n-1}$  to be precise).

### B. Sum-product algorithms

The sum-product algorithm is an iterative algorithm for computing either exact marginals (on trees), or approximate marginals (for graphs with cycles). It operates in a distributed manner, with nodes in the graph exchanging statistical information via a sequence of “message-passing” updates. For tree-structured graphical models, the updates can be derived as a form of non-serial dynamic programming, and are guaranteed to converge and compute the correct marginal distributions at each node. However, the updates are routinely applied to more general graphs with cycles, which is the application of interest in this paper. Here we describe the more general family of reweighted sum-product algorithms, which include the ordinary sum-product updates as a particular case.

In any sum-product algorithm, one message is passed in each direction of every edge  $(s, t)$  in the graph. The message from node  $t$  to node  $s$ , denoted by  $M_{ts}(x_s)$ , is a function of the possible states  $x_s \in \{0, 1, \dots, m-1\}$  at node  $s$ . Consequently, in the discrete case, the message can be represented by an  $m$ -vector of possible function values. The family of reweighted sum-product algorithms is parameterized by a set of *edge weights*, with  $\rho_{st} \in (0, 1]$  associated with edge  $(s, t)$ . Various choices of these edge weights have been proposed [16], [18], [19], and have different theoretical properties. The simplest case of all—namely, setting  $\rho_{st} = 1$  for all edges—recovers the ordinary sum-product algorithm. Given some fixed set of edge weights  $\rho_{st} \in (0, 1]$ , the reweighted sum-product updates are given by the recursion

$$M_{ts}(x_s) \leftarrow \sum_{x'_t} \exp \left\{ \frac{\theta_{st}(x_s, x'_t)}{\rho_{st}} + \theta_t(x'_t) \right\} \frac{\prod_{u \in N(t) \setminus s} [M_{ut}(x'_t)]^{\rho_{ut}}}{[M_{st}(x'_t)]^{1-\rho_{ts}}}, \quad (2)$$

where  $N(t) := \{s \in V \mid (s, t) \in E\}$  denotes the neighbors of node  $t$  in the graph. Typically, the message vector  $M_{ts}$  is normalized to unity after each iteration (i.e.,  $\sum_{x_s} M_{ts}(x_s) = 1$ ). Once the updates converge to some message fixed point  $M^*$ , then the fixed point can be used to compute (ap-

proximate) marginal probabilities  $\tau_s$  at each node via  $\tau_s(x_s) \propto \exp\{\theta_s(x_s)\} \prod_{t \in N(s)} [M_{ts}^*(x_s)]^{\rho_{st}}$ .

When the ordinary updates ( $\rho_{st} = 1$ ) are applied to a tree-structured graph, it can be shown by induction that the algorithm converges after a finite number of steps. Moreover, a calculation using Bayes' rule shows that  $\tau_s(x_s)$  is equal to the desired marginal probability  $p(x_s; \theta)$ . However, the sum-product algorithm is routinely applied to graphs with cycles, in which case the message updates (2) are not guaranteed to converge, and the quantities  $\tau_s(x_s)$  represent approximations to the true marginal distributions. Our focus in this paper is to determine conditions under which the reweighted sum-product message updates (2) are guaranteed to converge.

### III. CONVERGENCE ANALYSIS

In this section, we describe and provide proofs of our main results on the convergence properties of the reweighted sum-product updates (2) when the messages belong to a binary state space, which we represent as  $\mathcal{X} = \{-1, 1\}$ . In this special case, the general MRF distribution (1) can be simplified into the Ising model form

$$p(x; \theta) \propto \exp \left\{ \sum_{s \in V} \theta_s x_s + \sum_{(s,t) \in E} \theta_{st} x_s x_t \right\}, \quad (3)$$

so that the model is parameterized<sup>1</sup> by a single real number  $\theta_s$  for each node, and a single real number  $\theta_{st}$  for each edge.

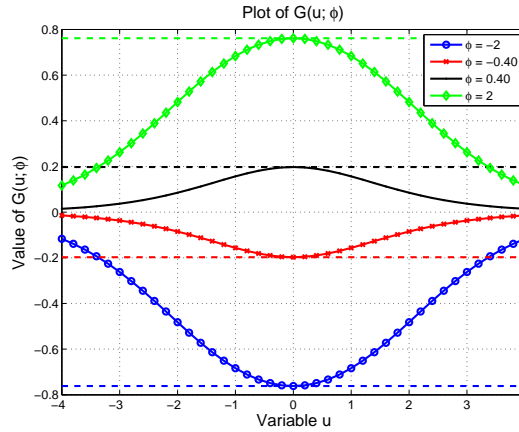
#### A. Convergence for binary homogeneous models

We begin by stating and proving some convergence conditions for a particularly simple class of models: homogeneous models on  $d$ -regular graphs. A graph is  $d$ -regular if each vertex has exactly  $d$  neighbors. Examples include single cycles ( $d = 2$ ), and lattice models with toroidal boundary conditions ( $d = 4$ ). In a homogeneous model, the edge weights  $\theta_{st}$  are all equal to a common value  $\theta_{ed}$ , and similarly the node parameters  $\theta_s$  are all equal to a common value  $\theta_{vx}$ .

In order to state our convergence result, we first define, for any real numbers  $u$  and  $\phi$ , the function

$$G(u; \phi) = \frac{\exp(\phi + u)}{1 + \exp(\phi + u)} - \frac{\exp(u)}{\exp(\phi) + \exp(u)} = \frac{\sinh \phi}{\cosh \phi + \cosh u}. \quad (4)$$

<sup>1</sup>This assumption is valid, because the distribution (1) does not change if we replace  $\theta_s(x_s)$  with  $\tilde{\theta}_s(x_s) := \theta_s(x_s) - \theta_s(-1)$ , with a similar calculation for the edges. See the paper [4] for details.



**Fig. 2.** Plots of the function  $G(\cdot; \phi)$  for  $\phi = 0.2$  and  $\phi = 3$ . For each fixed  $\phi \in \mathbb{R}$ , it is symmetric about 0 with  $\lim_{u \rightarrow \pm\infty} G(u; \phi) = 0$ . Moreover, the absolute value  $|G(u; \phi)|$  achieves its unconstrained maximum  $|G(0; \phi)| = \tanh(|\phi|/2) < 1$  at  $u^* = 0$ .

As illustrated in Figure 2, for any fixed  $\phi \in \mathbb{R}$ , the mapping  $u \mapsto G(u; \phi)$  is symmetric about zero, and  $\lim_{u \rightarrow \pm\infty} G(u; \phi) = 0$ . Moreover, the function  $G(\cdot; \phi)$  is bounded in absolute value by  $|G(0; \phi)| = \tanh(|\phi|/2)$

**Proposition 1.** Consider the reweighted sum-product algorithm with uniform weights  $\rho_{st} = \text{const}$  applied to a homogeneous binary model on a  $d$ -regular graph with arbitrary choice of  $(\theta_{vx}, \theta_{ed})$ .

(a) The reweighted updates (2) have a unique fixed point and converge as long as  $R < 1$ , where

$$R_d(\theta_{vx}, \theta_{ed}; \rho) := \begin{cases} |\rho d - 1| G\left(0; \frac{2|\theta_{ed}|}{\rho}\right) & \text{if } \delta \leq 0 \\ |\rho d - 1| G\left(\delta; \frac{2|\theta_{ed}|}{\rho}\right) & \text{otherwise.} \end{cases} \quad (5)$$

with  $\delta := 2|\theta_{vx}| - 2|\rho d - 1| \frac{|\theta_{ed}|}{\rho}$ .

(b) As a consequence, if  $\rho \leq 2/d$ , the reweighted updates (2) converge for all finite  $\theta_{ed}$ .

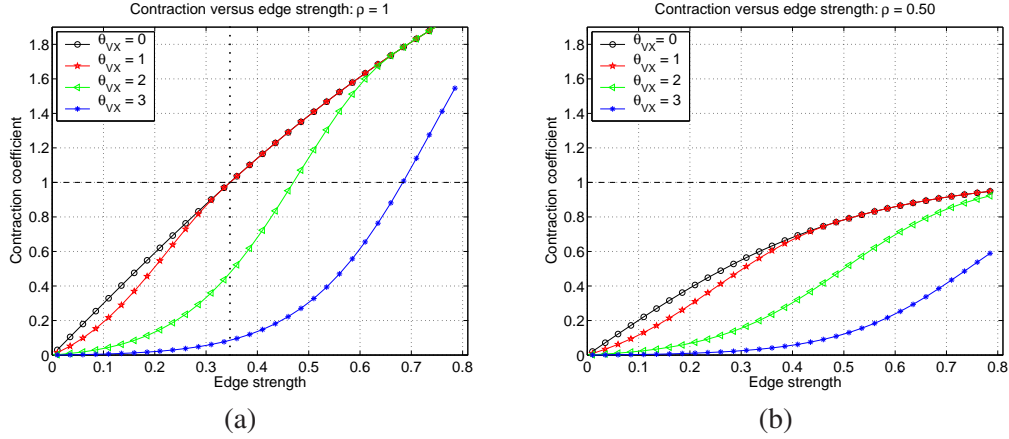
(c) Conversely, if  $\rho > 2/d$ , then there exist finite edge potentials  $\theta_{ed}$  for which the reweighted updates (2) have multiple fixed points.

**Remarks:** Consider the choice of edge weight  $\rho = 1$ , corresponding to the standard sum-product algorithm. If the graph is a single cycle ( $d = 2$ ), Proposition 1(b) implies that the standard sum-product algorithm always converges, consistent with previous work on the single cycle case [8], [22].

For more general graphs with  $d > 2$ , convergence depends on the relation between the observation



strength  $\theta_{vx}$  and the edge strengths  $\theta_{ed}$ . For the case  $d = 4$ , corresponding for instance to a lattice model with toroidal boundary as in Figure 1(c), Figure 3(a) provides a plot of the coefficient  $R_4(\theta_{vx}, \theta_{ed}; 1)$  as a function of the edge strength  $\theta_{ed}$ , for different choices of the observation potential  $\theta_{vx}$ . The curve marked with circles corresponds to  $\theta_{vx} = 0$ . Observe that it crosses the threshold  $R_4 = 1$  from convergence to non-convergence at the critical value  $\text{arctanh}(\frac{1}{3}) \approx 0.3466$ ,



**Fig. 3.** Plots of the contraction coefficient  $R_4(\theta_{vx}, \theta_{ed}; \rho)$  versus the edge strength  $\theta_{ed}$ . Each curve corresponds to a different choice of the observation potential  $\theta_{vx}$ . (a) For  $\rho = 1$ , the updates reduces to the standard sum-product algorithm; note that the transition from convergence to non-convergence occurs at  $\theta_{ed}^* \approx 0.3466$  in the case of no observations ( $\theta_{vx} = 0$ ). (b) Corresponding plots for reweighted sum-product with  $\rho = 0.50$ . Since  $\rho d = (0.50)4 = 2$ , the contraction coefficient is always less than one in this case, as predicted by Proposition 1.

corresponding with the classical result due to Bethe [23], and also confirmed in other analyses of standard sum-product [8], [11], [12]. The other curves correspond to non-zero observation potentials ( $\theta_{vx} \in \{1, 2, 3\}$ ) respectively. Here it is interesting to note with  $\theta_{vx} > 0$ , Proposition 1 reveals that the standard sum-product algorithm continues to converge well beyond the classical breakdown point without observations ( $\theta_{ed}^* \approx 0.3466$ ).

Figure 3(b) shows the corresponding curves of  $R_4(\theta_{vx}, \theta_{ed}; 0.50)$ , corresponding to the reweighted sum-product algorithm with  $\rho = 0.50$ . Note that  $\rho d = 0.50(4) = 2$ , so that as predicted by Proposition 1(b), the contraction coefficient  $R_4$  remains below one for all values of  $\theta_{vx}$  and  $\theta_{ed}$ , meaning that the reweighted sum-product algorithm with  $\rho = 0.50$  always converges for these graphs.

**Proof of Proposition 1:** Given the edge and node homogeneity of the model and the  $d$ -regularity

of the graph, the message-passing updates can be completely characterized by a single log message  $z = \log M(1)/M(-1) \in \mathbb{R}$ , and the update

$$F(z; \theta_{\text{vx}}, \theta_{\text{ed}}, \rho) = \log \left[ \frac{\exp[\frac{2\theta_{\text{ed}}}{\rho} + (\rho d - 1)z + 2\theta_{\text{vx}}] + 1}{\exp[(\rho d - 1)z + 2\theta_{\text{vx}}] + \exp(\frac{2\theta_{\text{ed}}}{\rho})} \right]. \quad (6)$$

(a) In order to prove the sufficiency of condition (5), we begin by observing that for any choice of  $z \in \mathbb{R}$ , we have  $|F(z; \theta_{\text{vx}}, \theta_{\text{ed}}, \rho)| \leq 2\frac{|\theta_{\text{ed}}|}{\rho}$ , so that the message  $z$  must belong to the *admissible interval*  $[-2\frac{|\theta_{\text{ed}}|}{\rho}, 2\frac{|\theta_{\text{ed}}|}{\rho}]$ . Next we compute and bound the derivative of  $F$  over this set of admissible messages. A straightforward calculation yields that  $F'(z) = (\rho d - 1)G\left(2\theta_{\text{vx}} + (\rho d - 1)z; 2\frac{\theta_{\text{ed}}}{\rho}\right)$ , where the function  $G$  was defined previously in (4). Note that for any fixed  $\phi \in \mathbb{R}$ , the function  $|G(u; \phi)|$  achieves its maximum at  $u^* = 0$ . Consequently, the unconstrained maximum of  $|F'(z)|$  is achieved at the point  $z^* = \frac{-2\theta_{\text{vx}}}{\rho d - 1}$  satisfying  $2\theta_{\text{vx}} + (\rho d - 1)z^* = 0$ , with  $\frac{F'(z^*)}{\rho d - 1} = G(0; 2\frac{|\theta_{\text{ed}}|}{\rho})$ . Otherwise, if  $|z^*| > 2\frac{|\theta_{\text{ed}}|}{\rho}$ , then the constrained maximum is obtained at the boundary point of the admissible region closest to 0—namely, at the point  $2\theta_{\text{vx}} - 2\text{sign}(\theta_{\text{vx}})\frac{|\theta_{\text{ed}}|(\rho d - 1)}{\rho}$ . Overall, we conclude that for all admissible messages  $z$ , we have

$$\frac{|F'(z)|}{|\rho d - 1|} \leq \begin{cases} G(0; 2\frac{|\theta_{\text{ed}}|}{\rho}) & \text{if } |\theta_{\text{vx}}| \leq \frac{|\rho d - 1||\theta_{\text{ed}}|}{\rho} \\ G(2|\theta_{\text{vx}}| - 2\frac{|\theta_{\text{ed}}||\rho d - 1|}{\rho}; 2\frac{|\theta_{\text{ed}}|}{\rho}) & \text{otherwise,} \end{cases} \quad (7)$$

so that  $|F'(z)| < R$  as defined in the statement. Note that if  $R < 1$ , the update is an iterated contraction, and hence converges [24].

(b) Since for all finite  $\theta_{\text{ed}}$ , we have  $\sup_{u \in \mathbb{R}} |G(u; \frac{2|\theta_{\text{ed}}|}{\rho})| = \tanh(\frac{|\theta_{\text{ed}}|}{\rho}) < 1$ , the condition  $|d\rho - 1| \leq 1$ —or equivalently  $d\rho \leq 2$ —implies that  $R < 1$ , showing that the updates are strictly contractive, which implies convergence and uniqueness of the fixed point [24].

(c) In our discussion below (following statement of Theorem 2), we establish that the condition (5) is actually necessary for the special case of zero observations ( $\theta_{\text{vx}} = 0$ ). Given  $\theta_{\text{vx}} = 0$  and  $\rho > 2/d$ , we can always find a finite  $\theta_{\text{ed}}$  such that the condition (5) is violated, so that the algorithm does not converge.

### B. Extension to binary inhomogeneous models

We now turn to the generalization of the previous result to the case of inhomogeneous models, in which the node parameters  $\theta_s$  and edge parameters  $\theta_{st}$  may differ across nodes and edges, respectively. For each *directed* edge ( $t \rightarrow s$ ), define the quantity

$$D_{t \rightarrow s}(\theta; \rho) = 2 \left\{ |\theta_t| - \sum_{u \in N(t) \setminus s} |\theta_{ut}| + (1 - \rho_{st}) \frac{|\theta_{st}|}{\rho_{st}} \right\} \quad (8)$$

and the weight

$$L_{t \rightarrow s} = \begin{cases} G(0; 2 \frac{|\theta_{st}|}{\rho_{st}}) & \text{if } D_{t \rightarrow s}(\theta; \rho) \leq 0. \\ G(D_{t \rightarrow s}(\theta; \rho); 2 \frac{|\theta_{st}|}{\rho_{st}}) & \text{otherwise.} \end{cases} \quad (9)$$

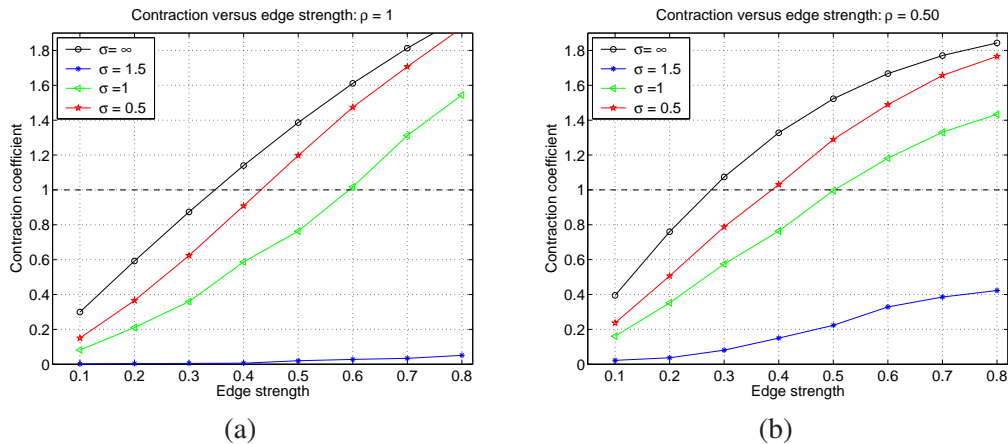
where the function  $G$  was previously defined (4). Finally, define a  $2|E| \times 2|E|$  matrix  $M = M(\theta; \rho)$ , with entries indexed by directed edges ( $t \rightarrow s$ ), and of the form

$$M_{(t \rightarrow s), (u \rightarrow v)} = \begin{cases} \rho_{ut} L_{t \rightarrow s} & \text{if } v = t \text{ and } u \neq s \\ (1 - \rho_{st}) L_{t \rightarrow s} & \text{if } v = t \text{ and } u = s \\ 0 & \text{otherwise.} \end{cases} \quad (10)$$

**Theorem 2.** *For an arbitrary pairwise Markov random field over binary variables, if the spectral radius of  $M(\theta; \rho)$  is less than one, then the reweighted sum-product algorithm converges, and the associated fixed point is unique.*

When specialized to the case of uniform edge weights  $\rho_{st} = 1$ , then Theorem 2 strengthens previous results, due independently to Ihler et al. [11] and Mooij and Kappen [12], on the ordinary sum-product algorithm. This earlier work provided conditions based on matrices that involved only terms of the form  $G(0; 2|\theta_{st}|)$ , as opposed to the smaller and observation-dependent weights  $G(D_{t \rightarrow s}(\theta; \rho); 2|\theta_{st}|)$  that our analysis yields once  $D_{t \rightarrow s}(\theta; \rho) > 0$ . As a consequence, Theorem 2 can yield sharper estimates of convergence by incorporating the benefits of having observations. For the ordinary sum-product algorithm and binary Markov random fields, independent work by Mooij and Kappen [20] has also made similar refinements of earlier results. In addition to these consequences for the ordinary sum-product algorithm, our Theorem 2 also provides sufficient

conditions for convergence of reweighted sum-product algorithms.



**Fig. 4.** Illustration of the benefits of observations. Plots of the contraction coefficient versus the edge strength. Each curve corresponds to a different setting of the noise variance  $\sigma^2$  as indicated. (a) Ordinary sum-product algorithm  $\rho = 1$ . Upper-most curve labeled  $\sigma^2 = +\infty$  corresponds to bounds taken from previous work [8], [11], [12]. (b) Corresponding curves for the reweighted sum-product algorithm  $\rho = 0.50$ .

In order to illustrate the benefits of including observations in the convergence analysis, we conducted experiments on grid-structured graphical models in which a binary random vector, with a prior distribution of the form (3), is observed in Gaussian noise (see Section V-A for the complete details of the experimental set-up). Figure 4 provides summary illustrations of our findings, for the ordinary sum-product ( $\rho = 1$ ) in panel (a), and reweighted sum-product ( $\rho = 0.50$ ) in panel (b). Each plot shows the contraction coefficient predicted by Theorem 2 as a function of an edge strength parameter. Different curves show the effect of varying the noise variance  $\sigma^2$  specifying the signal-to-noise ratio in the observation model (see Section V for the complete details). The extreme case  $\sigma^2 = +\infty$  corresponds to the case of no observations. Notice how the contraction coefficient steadily decreases as the observations become more informative, both for the ordinary and reweighted sum-product algorithms.

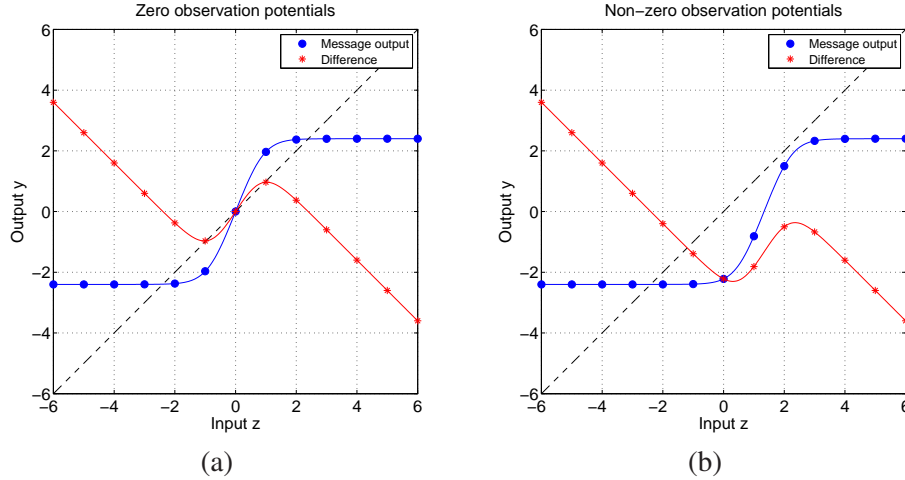
From past work, one sufficient (but not necessary) condition for uniqueness of the reweighted sum-product fixed point is convexity of the associated free energy [10], [16]. The contraction condition of Theorem 2 is also sufficient (but not necessary) to guarantee uniqueness of the fixed point. In general, these two sufficient conditions are incomparable, in that neither implies (nor is implied by) the other. For instance, Figure 4(b) shows that the condition from Theorem 2 is in

some cases weaker than the convexity argument. In this example, with the edge weights  $\rho = 0.50$ , the associated free energy is known to be convex [10], [16], so that the fixed point is always unique. However, the contraction bound from Theorem 2 guarantees uniqueness only up to some edge strength. On the other, the advantage of the contraction coefficient approach is illustrated by Figures 3(a) and 4(a): for the grid graph, the Bethe free energy is not convex for any edge strength, yet the contraction coefficient still yields uniqueness for suitably small couplings. An intriguing possibility raised by Heskes [10] is whether uniqueness of the fixed point implies convergence of the updates. It turns out that this implication does not hold for the reweighted sum-product algorithm: in particular, there are weights  $\{\rho_{st}\}$  for which the optimization is strictly convex (and so has a unique global optimum), but the reweighted sum-product updates do not converge.

A related question raised by a referee concerns the tightness of the sufficient conditions in Proposition 1 and Theorem 2: that is, to what extent are they also necessary conditions? Focusing on the conditions of Proposition 1, we can analyze this issue by studying the function  $F$  defined in the message update equation (6), and determining for what pairs  $(\theta_{vx}, \theta_{ed})$  it has multiple fixed points. For clarity, we state the following

**Corollary 3.** *Regarding uniqueness of fixed points, the conditions of Proposition 1 are necessary and sufficient for  $\theta_{vx} = 0$ , but only sufficient for general  $\theta_{vx} \neq 0$ .*

*Proof:* For  $\theta_{ed} > 0$  (or respectively  $\theta_{ed} < 0$ ), the function  $F(z)$  has (for any  $\theta_{vx}$ ) the following general properties: it is monotonically increasing (decreasing) in  $z$ , and converges to  $\pm 2\theta_{ed}/\rho$  as  $z \rightarrow \pm\infty$ , as shown in the curves marked with circles in Figure 5. Its fixed point structure is most easily visualized by plotting the difference function  $F(z) - z$ , as shown by the curves marked with \*-symbols. Note that the \*-curve in panel (a) crosses the horizontal  $y = 0$  three times, corresponding to three fixed points, whereas the curve in panel (b) crosses only once, showing that the fixed point is unique. In contrast, the sufficient conditions of Proposition 1 and Theorem 2 are based on whether the magnitude of the derivative  $|F'(z)|$  exceeds one, or equivalently whether the quantity  $\frac{d}{dz}[F(z) - z]$  ever becomes zero. Since the \*-curves (corresponding to  $F(z) - z$ ) in both panels (a) and (b) have flat portions, we see that the sufficient conditions do not hold in either case. In particular, panel (b) is an instance where the fixed point is unique, yet the conditions of



**Fig. 5.** Illustration of the necessity of the conditions in Proposition 1, which fail to be satisfied in both cases. (a) Zero observation potentials  $\theta_{vx} = 0$ : here there are three fixed points, and the conditions are necessary in general. (b) Non-zero observation potentials  $\theta_{vx} > 0$ : in this case, the fixed point is still unique, yet the conditions of Proposition 1 are not satisfied.

Proposition 1 fail to hold, showing that they are not necessary in general. Whereas panel (b) involves  $\theta_{vx} = -2 \neq 0$ , panel (a) has zero observation potentials  $\theta_{vx} = 0$ . In this case, the conditions of Proposition 1 are actually necessary and sufficient. To establish this claim rigorously, note that for  $\theta_{vx} = 0$ , the point  $z^* = 0$  is always a fixed point, and the maximum of  $F'(z)$  is achieved at  $z^* = 0$ . If  $|F'(z)| > 1$  (so that the condition in Proposition 1) is violated, then continuity implies that  $F(z) - z > 0$  for all  $z \in (0, \epsilon)$ , for some suitably small  $\epsilon > 0$ . But since  $\lim_{z \rightarrow +\infty} F(z) = 2\theta_{ed}/\rho$ , the curve must eventually cross the identity line again, implying that there is a second fixed point  $\hat{z} > 0$ . By symmetry, there must also exist a third fixed point  $\tilde{z} < 0$ , as illustrated in Figure 5(a). Therefore, the condition of Proposition 1 is actually necessary and sufficient for the case  $\theta_{vx} = 0$ . ■

### C. Proof of Theorem 2

We begin by establishing a useful auxiliary result that plays a key role in this proof, as well as other proofs in the sequel:

**Lemma 4.** For real numbers  $\phi$  and  $u$ , define the function

$$H(u; \phi) = \log \frac{\exp(\phi + u) + 1}{\exp(u) + \exp(\phi)}. \quad (11)$$

For each fixed  $\phi$ , we have  $\sup_{u \in \mathbb{R}} |H(u; \phi)| \leq |\phi|$ .

*Proof:* Computing the derivative of  $H$  with respect to  $u$ , we have

$$H'(u; \phi) = \frac{\exp(\phi + u)}{1 + \exp(\phi + u)} - \frac{\exp(u)}{\exp(u) + \exp(\phi)} = G(u; \phi),$$

where the function  $G$  was previously defined (4). Therefore, the function  $H$  is strictly increasing if  $\phi > 0$  and strictly decreasing if  $\phi < 0$  (see Figure 2 for a plot of  $G = H'$ ). Consequently, the supremum is obtained by taking  $u \rightarrow \pm\infty$ , and is equal to  $|\phi|$  as claimed. ■

With this lemma in hand, we begin by re-writing the message update (2) in a form more amenable to analysis. For each directed edge ( $t \rightarrow s$ ), define the log message ratio  $z_{ts} = \log \frac{M_{t \rightarrow s}(1)}{M_{t \rightarrow s}(-1)}$ . From the standard form of the updates, a few lines of algebra show that it is equivalent to update these log ratios via

$$F_{t \rightarrow s}(z) := \log \frac{\exp \left[ \frac{2\theta_{st}}{\rho_{st}} + 2\theta_t + \sum_{v \in N(t) \setminus s} \rho_{vt} z_{v \rightarrow t} + (1 - \rho_{st}) z_{st} \right] + 1}{\exp \left[ 2\theta_t + \sum_{v \in N(t) \setminus s} \rho_{vt} z_{vt} + (1 - \rho_{st}) z_{st} \right] + \exp \left[ \frac{2\theta_{st}}{\rho_{st}} \right]}. \quad (12)$$

A key property of the message update function  $F_{t \rightarrow s}$  is that it can be written as a function  $H$  of the form (11), with  $\phi = 2\frac{\theta_{st}}{\rho_{st}}$  and  $u = 2\theta_t + \sum_{v \in N(t) \setminus s} \rho_{vt} z_{vt} + (1 - \rho_{st}) z_{st}$ . Consequently, if we apply Lemma 4, we may conclude that  $|F_{t \rightarrow s}(z)| \leq 2\frac{|\theta_{st}|}{\rho_{st}}$  for all  $z \in \mathbb{R}$ , and consequently that  $|z_{ts}^n| \leq 2\frac{|\theta_{st}|}{\rho_{st}}$  for all iterations  $n \geq 1$ . Consequently, we may assume that message vector  $z^n$  for all iterations  $n \geq 1$  belongs to the box of admissible messages defined by

$$\mathbb{B}(\theta; \rho) := \left\{ z \in \mathbb{R}^{2|E|} \mid |z_{ts}| \leq 2\frac{|\theta_{st}|}{\rho_{st}} \quad \text{for all edges } (t \rightarrow s) \right\}. \quad (13)$$

We now bound the derivative of the message-update equation over this set of admissible messages:

**Lemma 5.** For all  $z \in \mathbb{B}(\theta; \rho)$ , the elements of  $\nabla F_{t \rightarrow s}(z)$  are bounded as

$$\left| \frac{\partial F_{t \rightarrow s}}{\partial z_{ut}}(z) \right| \leq \rho_{ut} L_{t \rightarrow s} \quad \forall \quad u \in N(t) \setminus s, \quad \text{and} \quad \left| \frac{\partial F_{t \rightarrow s}}{\partial z_{st}}(z) \right| \leq (1 - \rho_{st}) L_{t \rightarrow s}, \quad (14)$$

where the directed weights  $L_{t \rightarrow s}$  were defined previously in (9). All other gradient elements are zero.

See Appendix A for the proof. In order to exploit Lemma 5, for any iteration  $n \geq 2$ , let us use the mean-value theorem to write

$$z_{st}^{n+1} - z_{st}^n = F_{t \rightarrow s}(z^n) - F_{t \rightarrow s}(z^{n-1}) = \nabla F_{t \rightarrow s}(z^\lambda)^T (z^n - z^{n-1}), \quad (15)$$

where  $z^\lambda = \lambda z^n + (1 - \lambda)z^{n-1}$  for some  $\lambda \in (0, 1)$ . Since  $z^n$  and  $z^{n-1}$  both belong to the convex set  $\mathbb{B}(\theta; \rho)$ , so does the convex combination  $z^\lambda$ , and we can apply Lemma 5. Starting from equation (15), we have

$$\begin{aligned} |z_{t \rightarrow s}^{n+1} - z_{t \rightarrow s}^n| &\leq \left| \nabla F_{t \rightarrow s}(z^\lambda) \right|^T |z^n - z^{n-1}| \\ &\leq \left( \sum_{u \in N(t) \setminus s} \rho_{ut} L_{t \rightarrow s} |z_{u \rightarrow t}^n - z_{u \rightarrow t}^{n-1}| \right) + (1 - \rho_{st}) L_{t \rightarrow s} |z_{s \rightarrow t}^n - z_{s \rightarrow t}^{n-1}|. \end{aligned} \quad (16)$$

Since this bound holds for each directed edge, we have established that the vector of message differences obeys  $|z^{n+1} - z^n| \leq M(\theta, \rho) |z^n - z^{n-1}|$ , where the non-negative matrix  $M = M(\theta, \rho)$  was defined previously in 10. By standard results on non-negative matrix recursions [25], if the spectral radius of  $M$  is less than 1, then the sequence  $|z^n - z^{n-1}|$  converges to zero. Thus, the sequence  $\{z^n\}$  is a Cauchy sequence, and so must converge.

#### D. Explicit conditions for convergence

A drawback of Theorem 2 is that it requires computing the spectral radius of the  $2|E| \times 2|E|$  matrix  $M$ , which can be a non-trivial computation for large problems. Accordingly, we now specify some corollaries that are sufficient to ensure convergence of the reweighted sum-product algorithm. As in the work of Mooij and Kappen [12], the first two conditions follow by upper bounding the spectral norm by standard matrix norms. Conditions (c) and (d) are refinements that require further work.

**Corollary 6.** *Convergence of reweighted sum-product is guaranteed by any of the following conditions:*

(a) *Row sum condition:*  $\max_{(t \rightarrow s)} \left( \sum_{u \in N(t) \setminus s} \rho_{ut} + (1 - \rho_{st}) \right) L_{t \rightarrow s} < 1$



(b) *Column sum condition:*

$$\max_{(t \rightarrow s)} C_{t \rightarrow s} = \max_{(t \rightarrow s)} \left\{ \rho_{ts} \left( \sum_{u \in N(t) \setminus s} L_{u \rightarrow t} \right) + (1 - \rho_{ts}) L_{s \rightarrow t} \right\} < 1. \quad (17)$$

(c) *Reweighted norm condition:*

$$K(\theta) := \max_{(t \rightarrow s)} \left\{ \left( \sum_{u \in N(t) \setminus s} \rho_{ut} L_{u \rightarrow t} \right) + (1 - \rho_{ts}) L_{s \rightarrow t} \right\} < 1. \quad (18)$$

(d) *Pairwise neighborhood condition: the quantity*

$$\min_{\lambda \in [0,1]} \max_{(t \rightarrow s)} \left\{ \rho_{ts} \left( \sum_{w \in N(t) \setminus s} L_{w \rightarrow t} \right) + (1 - \rho_{ts}) L_{s \rightarrow t} \right\}^\lambda \max_{u \in N(t)} \left\{ \rho_{ts} \left( \sum_{v \in N(u) \setminus t} L_{v \rightarrow u} \right) + (1 - \rho_{tu}) L_{t \rightarrow u} \right\}^{1-\lambda}$$

*is less than one.*

**Remarks:** To put these results in perspective, if we specialize to  $\rho_{st} = 1$  for all edges and use the *weaker version* of the weights  $L_{t \rightarrow s}$  that ignore the effects of observations, then the  $\ell_\infty$ -norm condition (17) is equivalent to earlier results on the ordinary sum-product algorithm [11], [12]. In addition, one may observe that for the ordinary sum-product algorithm (where  $\rho_{st} = 1$  for all edges), condition (18) is equivalent to the  $\ell_\infty$ -condition (17). However, for the general reweighted algorithm, these two conditions are distinct.

*Proof:* Conditions (a) and (b) follows immediately from the fact that the spectral norm of  $M$  is upper bounded by any other matrix norm [25]. It remains to prove conditions (c) and (d) in the corollary statement.

(c) Defining the vector  $\Delta^n = |z^n - z^{n-1}|$  of successive changes, from equation (16), we have

$$\Delta_{t \rightarrow s}^{n+1} \leq L_{t \rightarrow s} \left\{ \sum_{u \in N(t) \setminus s} \rho_{ut} \Delta_{u \rightarrow t}^n + (1 - \rho_{st}) \Delta_{s \rightarrow t}^n \right\} \quad (19)$$

The previous step of the updates yields a similar equation—namely

$$\Delta_{u \rightarrow t}^n \leq L_{u \rightarrow t} \left\{ \sum_{v \in N(u) \setminus t} \rho_{vu} \Delta_{v \rightarrow u}^{n-1} + (1 - \rho_{ut}) \Delta_{t \rightarrow u}^{n-1} \right\}. \quad (20)$$

Now let us define a norm on  $\Delta$  by  $\|\Delta\|_* = \max_{(t \rightarrow s) \in E} \left\{ \sum_{u \in N(t) \setminus s} \rho_{ut} |\Delta_{u \rightarrow t}| + (1 - \rho_{st}) |\Delta_{s \rightarrow t}| \right\}$ . With this notation, the bound (20) implies that  $\Delta_{u \rightarrow t}^n \leq L_{u \rightarrow t} \|\Delta^{n-1}\|_*$ . Substituting this bound into equation (19) yields that

$$\Delta_{t \rightarrow s}^{n+1} \leq L_{t \rightarrow s} \left\{ \sum_{u \in N(t) \setminus s} \rho_{ut} L_{u \rightarrow t} + (1 - \rho_{st}) L_{s \rightarrow t} \right\} \|\Delta^{n-1}\|_* \leq K(\theta) \|\Delta^{n-1}\|_*.$$

For any edge  $(s \rightarrow u)$ , summing weighted versions of this equation over all neighbors of  $s$  yields that

$$\begin{aligned} \sum_{v \in N(s) \setminus u} \rho_{vs} \Delta_{v \rightarrow s}^{n+1} + (1 - \rho_{us}) \Delta_{u \rightarrow s}^{n+1} &\leq \left\{ \sum_{v \in N(s) \setminus u} \rho_{vs} L_{v \rightarrow s} + (1 - \rho_{us}) L_{u \rightarrow s} \right\} K(\theta) \|\Delta^{n-1}\|_* \\ &\leq K^2(\theta) \|\Delta^{n-1}\|_*. \end{aligned}$$

Finally, since the edge  $(s \rightarrow u)$  was arbitrary, we can maximize over it, which proves that  $\|\Delta^{n+1}\|_* \leq K^2(\theta) \|\Delta^{n-1}\|_*$ . Therefore, if  $K(\theta) < 1$ , the updates are an iterated contraction in the  $\|\cdot\|_*$  norm, and hence converge by standard contraction results [24].

(d) Given a non-negative matrix  $A$ , let  $C_\alpha(A)$  denote the column sum indexed by some element  $\alpha$ . In general, it is known [25] that for any  $\lambda \in [0, 1]$ , the spectral radius of  $A$  is upper bounded by the quantity  $\max_{\alpha, \beta} [C_\alpha(A)]^\lambda [C_\beta(A)]^{1-\lambda}$ , where  $\alpha$  and  $\beta$  range over all column indices. A more refined result due to Kolotilina [26] asserts that if  $A$  is a sparse matrix, then one need only optimize over column index pairs  $\alpha, \beta$  such that  $A_{\alpha\beta} \neq 0$ . For our problem, the matrix  $M$  is indexed by directed edges  $(s \rightarrow t)$ , and  $M_{(t \rightarrow s), (u \rightarrow v)}$  is non-zero only if  $v = t$ . Consequently, we can reduce the maximization over column sums to maximizing over directed edge pairs  $(t \rightarrow s)$  with  $(u \rightarrow t)$ , which yields the stated claim. ■

#### IV. CONVERGENCE FOR GENERAL DISCRETE MODELS

In this section, we describe how our results generalize to multinomial random variables, with the variable  $X_s$  at each node  $s$  taking a total of  $m \geq 2$  states in the space  $\mathcal{X} = \{0, 1, \dots, m-1\}$ .

Given our Markov assumptions, the distribution takes the factorized form

$$p(x; \theta) \propto \exp \left\{ \sum_{s \in V} \theta_s(x_s) + \sum_{(s,t) \in E} \theta_{st}(x_s, x_t) \right\}, \quad (21)$$

where each  $\theta_s(\cdot)$  is a vector of  $m$  numbers, and each  $\theta_{st}(\cdot, \cdot)$  is an  $m \times m$  matrix of numbers.

In our analysis, it will be convenient to work with an alternative parameter vector  $\tilde{\theta}$  that represents the same Markov random field as  $p(x; \theta)$ , given by

$$\tilde{\theta}_{st}(x_s, x_t) := \theta_{st}(x_s, x_t) - \theta_{st}(x_s, 0) - \theta_{st}(0, x_t) + \theta_{st}(0, 0), \text{ and} \quad (22a)$$

$$\tilde{\theta}_s(x_s) = [\theta_s(x_s) - \theta_s(0)] + \sum_{t \in N(s)} [\theta_{st}(x_s, 0) - \theta_{st}(0, 0)]. \quad (22b)$$

This set of functions  $\tilde{\theta}$  is a different parameterization of the distribution  $p(x; \theta)$  because

$$\sum_{s \in V} \tilde{\theta}_s(x_s) + \sum_{(s,t) \in E} \tilde{\theta}_{st}(x_s, x_t) = \sum_{s \in V} \theta_s(x_s) + \sum_{(s,t) \in E} \theta_{st}(x_s, x_t) + C,$$

where  $C$  is a constant independent of  $x$ . Moreover, note that  $\tilde{\theta}_s(0) = 0$  for all nodes  $s \in V$ , and  $\tilde{\theta}_{st}(x_s, 0) = \tilde{\theta}_{st}(0, x_t) = 0$  for all  $x_s, x_t \in \{0, 1, \dots, m-1\}$ .

#### A. Convergence Theorem and Some Consequences

In order to state a result about convergence for multinomial Markov random fields, we require some preliminary definitions. For each directed edge ( $t \rightarrow s$ ) and states  $i, k \in \{1, \dots, m-1\}$ , define functions of the vector  $\vec{v} = (v(1), \dots, v(m-1))$  as follows

$$\psi_{t \rightarrow s}(\vec{v}; k, i) := \frac{1}{2} \left| -\beta_{t \rightarrow s}(\vec{v}; i, k) - \alpha_{t \rightarrow s}(\vec{v}; i, k) + v(k) + \tilde{\theta}_t(k) + \frac{\tilde{\theta}_{ts}(k, i)}{\rho_{st}} \right|, \text{ and} \quad (23a)$$

$$\phi_{t \rightarrow s}(\vec{v}; k, i) := \frac{1}{2} \left| \beta_{t \rightarrow s}(\vec{v}; k, i) - \alpha_{t \rightarrow s}(\vec{v}; k, i) + \frac{\tilde{\theta}_{ts}(k, i)}{\rho_{st}} \right|, \quad (23b)$$

where

$$\alpha_{t \rightarrow s}(\vec{v}; k, i) := \log \left( 1 + \sum_{x_t \neq 0, k} \exp \left\{ \frac{\tilde{\theta}_{ts}(i, x_t)}{\rho_{st}} + v(x_t) + \tilde{\theta}_t(k) \right\} \right) \quad (24a)$$

$$\beta_{t \rightarrow s}(\vec{v}; k, i) := \log \left( 1 + \sum_{x_t \neq 0, k} \exp(v(x_t) + \tilde{\theta}_t(k)) \right). \quad (24b)$$

With these definitions, define for each directed edge ( $t \rightarrow s$ ) the non-negative weight

$$L_{t \rightarrow s} := \max_{i,k \in \{1, \dots, m-1\}} \max_{\vec{v} \in \mathbb{B}_{ts}(\tilde{\theta}; \rho)} |G(\psi_{t \rightarrow s}(\vec{v}; i, k); \phi_{t \rightarrow s}(\vec{v}; i, k))|, \quad (25)$$

where the function  $G$  was defined previously (4), and the box of admissible vectors is given by

$$\mathbb{B}_{ts}(\tilde{\theta}; \rho) := \left\{ \vec{v} \in \mathbb{R}^{m-1} \mid |v(k)| \leq \sum_{u \in N(t) \setminus s} \max_j \frac{|\tilde{\theta}_{ut}(j, k)|}{\rho_{ut}} + (1 - \rho_{st}) \max_j \frac{|\tilde{\theta}_{st}(j, k)|}{\rho_{st}} \right\}. \quad (26)$$

Finally, using the choice of weights  $L_{t \rightarrow s}$  in equation (25), we define the  $2|E| \times 2|E|$  matrix  $M = M(L)$  as before (see equation (10)).

**Theorem 7.** *If the spectral radius of  $M$  is less than one, then the reweighted sum-product algorithm converges, and the associated fixed point is unique.*

We provide the proof of Theorem 7 in the appendix. Despite its notational complexity, Theorem 7 is simply a natural generalization of our earlier results for binary variables. When  $m = 2$ , note that the functions  $\beta_{t \rightarrow s}$  and  $\alpha_{t \rightarrow s}$  are identically zero (since there are no states other than  $k = 1$  and 0), so that the form of  $\phi$  and  $\psi$  simplifies substantially. Moreover, as in our earlier development on the binary case, when specialized to  $\rho_{st} = 1$ , Theorem 7 provides a strengthening of previous results [11], [12]. In particular, we now show how these previous results can be recovered from Theorem 7 by ignoring the box constraints (26):

**Corollary 8.** *The reweighted sum-product algorithm converges if*

$$\max_{t \rightarrow s} \sum_{u \in N(t) \setminus s} \rho_{ut} W_{u \rightarrow t} + (1 - \rho_{st}) W_{s \rightarrow t} < 1, \quad (27)$$

where  $W_{u \rightarrow t} := \tanh \left( \frac{1}{4\rho_{ut}} \max_{i \neq j} \max_{\ell \neq k} |\theta_{ts}(\ell, i) - \theta_{ts}(\ell, j) - \theta_{ts}(k, i) + \theta_{ts}(k, j)| \right)$ .

*Proof:* We begin by proving that  $L_{u \rightarrow t} \leq W_{u \rightarrow t}$ . First of all, ignoring the box constraints (26), then certainly

$$\begin{aligned} L_{t \rightarrow s} &\leq \max_{i,k \in \{1, \dots, m-1\}} \max_{\vec{v} \in \mathbb{R}^{m-1}} |G(\psi_{t \rightarrow s}(\vec{v}; i, k); \phi_{t \rightarrow s}(\vec{v}; i, k))| \\ &\leq \max_{i,k \in \{1, \dots, m-1\}} \max_{\vec{v} \in \mathbb{R}^{m-1}} |G(0; \phi_{t \rightarrow s}(\vec{v}; i, k))| = \max_{i,k \in \{1, \dots, m-1\}} \max_{\vec{v} \in \mathbb{R}^{m-1}} \tanh \left( \frac{1}{2} |\phi_{t \rightarrow s}(\vec{v}; i, k)| \right), \end{aligned}$$

since for any fixed  $\phi$ , the function  $G(u^*; \phi)$  is maximized at  $u^* = 0$ , and  $G(0; |\phi|) = \tanh(|\phi|/2)$ . Due to the monotonicity of  $G(0; |\phi|)$  in  $\phi$ , it now suffices to maximize the absolute value of  $|\phi_{t \rightarrow s}(\vec{v}; i, k)|$ . Since  $\phi$  is defined in terms of  $\beta$  and  $\alpha$ , we first bound their difference. In one direction, we have

$$\alpha_{t \rightarrow s}(\vec{v}; i, k) - \beta_{t \rightarrow s}(\vec{v}; i, k) = \log \frac{1 + \sum_{x_t \neq 0, k} \exp \left\{ \frac{\tilde{\theta}_{ts}(x_t, i)}{\rho_{st}} + v(x_t) + \tilde{\theta}_t(x_t) \right\}}{1 + \sum_{x_t \neq 0, k} \exp \left\{ v(x_t) + \tilde{\theta}_t(x_t) \right\}} \geq \min_{x_t \neq k, 0} \frac{\tilde{\theta}_{ts}(x_t, i)}{\rho_{st}},$$

and hence

$$\phi_{t \rightarrow s}(\vec{v}; k, i) \leq \frac{1}{2\rho_{st}} \max_{x_t \neq k, 0} \left\{ \tilde{\theta}_{ts}(k, i) - \tilde{\theta}_{ts}(x_t, i) \right\}. \quad (28)$$

In the other direction, we have  $\beta_{t \rightarrow s}(\vec{v}; i, k) - \alpha_{t \rightarrow s}(\vec{v}; i, k) \geq -\max_{x_t \neq k, 0} \frac{\tilde{\theta}_{ts}(x_t, i)}{\rho_{st}}$ , and hence

$$\phi_{t \rightarrow s}(\vec{v}; i, k) \geq -\frac{1}{2\rho_{st}} \max_{x_t \neq k, 0} \left\{ \tilde{\theta}_{ts}(x_t, i) - \tilde{\theta}_{ts}(i, k) \right\}. \quad (29)$$

Combining equations (28) and (29), we conclude that  $\max_{i, k \neq 0} \max_{\vec{v} \in \mathbb{R}^{m-1}} |\phi_{t \rightarrow s}(\vec{v}; i, k)|$  is upper bounded by

$$\frac{1}{2\rho_{st}} \max_{i, k} \max_{\ell \neq k, 0} \left| \tilde{\theta}_{ts}(\ell, i) - \tilde{\theta}_{ts}(k, i) \right| = \frac{1}{2\rho_{st}} \max_{i \neq 0} \max_{\ell \neq k} \left| \theta_{ts}(\ell, i) - \theta_{ts}(\ell, 0) - \theta_{ts}(k, i) + \theta_{ts}(k, 0) \right|$$

Therefore, we have proved that  $L_{t \rightarrow s} \leq W_{t \rightarrow s}$ , where  $W_{t \rightarrow s}$  was defined in the corollary statement. Consequently, if we define a matrix  $M(W)$  using the weights  $W$ , we have  $M(L) \leq M(W)$  in an elementwise sense, and therefore, the spectral radius of  $M(W)$  is an upper bound on the spectral radius of  $M(L)$  (see Bertsekas and Tsitsiklis [25]). ■

A challenge associated with verifying the sufficient conditions in Theorem 7 is that in principle, it requires solving a set of  $(m-1)^2$  optimization problems, each involving the  $m-1$  dimensional vector  $\vec{v}$  restricted to a box. As pointed out by one of the referees, in the absence of additional structure (e.g., convexity), one might think of obtaining the global maximum by discretizing the  $(m-1)$ -cube, albeit with prohibitive complexity (in particular, growing as  $(1/\epsilon)^m$  where  $\epsilon \rightarrow 0$  is the discretization accuracy). However, there are many natural weakenings of the bound in Theorem 7, still stronger than Corollary 8, that can be computed efficiently. Below we state one such result:

**Corollary 9.** *Define the edge weights*

$$K_{t \rightarrow s} := \max_{i,k \in \{1, \dots, m-1\}} \max_{\vec{v} \in \mathbb{B}_{ts}(\tilde{\theta}; \rho)} \left| G\left(\psi_{t \rightarrow s}(\vec{v}; i, k); \tanh^{-1}(2W_{t \rightarrow s})\right) \right|, \quad (30)$$

where the weights  $W_{t \rightarrow s}$  are defined in Corollary 8. Then the reweighted sum-product algorithm converges if  $\max_{t \rightarrow s} \sum_{u \in N(t) \setminus s} \rho_{ut} K_{u \rightarrow t} + (1 - \rho_{st}) K_{s \rightarrow t} < 1$ .

*Proof:* The proof of this result is analogous to Corollary 8: we weaken the bound from Theorem 7 by optimizing  $G$  separately over its arguments:

$$\begin{aligned} \max_{\vec{v} \in \mathbb{B}_{ts}(\tilde{\theta}; \rho)} |G(\psi_{t \rightarrow s}(\vec{v}; i, k); \phi_{t \rightarrow s}(\vec{v}; i, k))| &\leq \max_{\vec{v} \in \mathbb{B}_{ts}(\tilde{\theta}; \rho)} \max_{\vec{v}' \in \mathbb{B}_{ts}(\tilde{\theta}; \rho)} |G(\psi_{t \rightarrow s}(\vec{v}; i, k); \phi_{t \rightarrow s}(\vec{v}'; i, k))| \\ &\leq \max_{\vec{v} \in \mathbb{B}_{ts}(\tilde{\theta}; \rho)} |G(\psi_{t \rightarrow s}(\vec{v}; i, k); \tanh^{-1}(2W_{t \rightarrow s}))|, \end{aligned}$$

where we have used the fact that  $G(u; \phi)$  is increasing in  $\phi$  (for any fixed  $u$ ), and the bound (29), which can be restated as  $\max_{\vec{v}'} |\phi_{t \rightarrow s}(\vec{v}'; i, k)| \leq 2 \tanh^{-1}(2W_{t \rightarrow s})$ .  $\blacksquare$

The advantage of Corollary 9 is that it is relatively straightforward to compute the weights  $K_{t \rightarrow s}$ : in particular, for a fixed  $i, k$ , we need to solve the maximization problem over  $\vec{v} \in \mathbb{B}_{ts}(\tilde{\theta}; \rho)$  in equation (30). Since  $|G(u; \phi)|$  achieves its unconstrained maximum at  $u^* = 0$ , and decreases as  $u$  moves away from zero, it is equivalent to minimize the function  $\psi_{t \rightarrow s}(\vec{v}; i, k)$  over the admissible box. In order to so, it suffices by the definition of  $\psi$  to minimize and maximize the function

$$f(\vec{v}; i, k) := \frac{1}{2} \left\{ -\beta_{t \rightarrow s}(\vec{v}; i, k) - \alpha_{t \rightarrow s}(\vec{v}; i, k) + v(k) + \tilde{\theta}_t(k) + \frac{\tilde{\theta}_{ts}(k, i)}{\rho_{st}} \right\}.$$

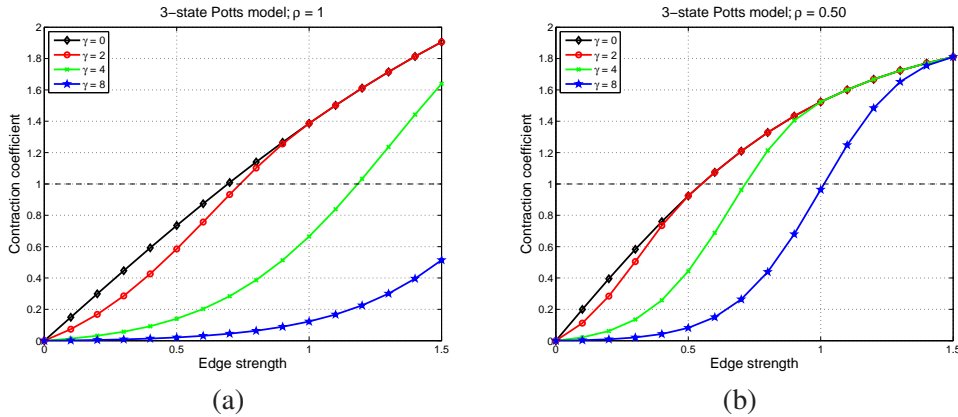
over the admissible box, and then take the minimum of the absolute values of these two quantities. Since both  $\beta$  and  $\alpha$  are convex and differentiable functions of  $\vec{v}$ , we see that  $f$  is a concave and differentiable function of  $\vec{v}$ , so that the maximum can be computed with standard gradient methods. On the other hand, the minimum of a concave function over a convex set is always achieved at an extreme point [27], which in this case, are simply the vertices of the admissible cube. Therefore, the weight  $K_{t \rightarrow s}$  can be computed in time that is at worst  $\mathcal{O}(2^{m-1})$ , corresponding to the number of vertices of the  $(m-1)$ -cube.

If the graphical model has very weak observations (i.e.,  $\theta_s(k) \approx 0$  uniformly for all states  $k \in \{0, 1, \dots, m-1\}$ ), then the observation-dependent conditions provided in Theorem 7 and

Corollary 9 would be expected to provide little to no benefit over Corollary 8. However, as with the earlier results on binary models (see Figure 4), the benefits can be substantial when the model has stronger observations, as would be common in applications. To be concrete, let us consider particular type of multi-state discrete model, known as the Potts model, in which the edge potentials are of the form

$$\theta_{st}(x_s, x_t) = \begin{cases} 0 & \text{if } x_s = x_t \\ -\beta & \text{otherwise,} \end{cases}$$

for some edge strength parameter  $\beta \in \mathbb{R}$ . We then suppose that the node potential vector  $\theta_s(x_s)$  at each node takes the form  $\gamma \begin{bmatrix} 0 & -1 & +1 \end{bmatrix}$ , for a signal-to-noise parameter  $\gamma$  to be chosen. Figure 6 illustrates the resulting contraction coefficients predicted by Corollary 9, for both the



**Fig. 6.** Illustration of the benefits of observations. Plots of the contraction coefficient from Corollary 9 versus the edge strength for a 3-state Potts model. Each curve corresponds to a different setting of the SNR parameter  $\gamma$ . (a) Ordinary sum-product algorithm  $\rho = 1$ . Upper-most curve labeled  $\gamma = 0$  corresponds to the best bounds from previous work [11], [12], [20]. (b) Reweighted sum-product algorithm  $\rho = 0.50$ .

ordinary sum-product updates ( $\rho = 1$ ) in panel (a), and the reweighted sum-product algorithm with  $\rho = 0.50$  in panel (b). As would be expected, the results are qualitatively similar to those from the binary state models in Figure 4.

## V. EXPERIMENTAL RESULTS

In this section, we present the results of some additional simulations to illustrate and support our theoretical findings.

### A. Dependence on Signal-to-Noise Ratio

We begin by describing the experimental set-up used to generate the plots in Figure 4, which illustrate the effect of increased signal-to-noise ratio (SNR) on convergence bounds. In these simulations, the random vector  $X \in \{-1, +1\}^n$  is posited to have a prior distribution  $p(x; \theta)$  of the form (3), with the edge parameters  $\theta_{st}$  set uniformly to some fixed number  $\theta_{ed}$ , and symmetric node potentials  $\theta_s = 0$ . Now suppose that we make a noisy observation of the random vector  $X$ , say of the form  $Y_s = X_s + W_s$ , where  $W_s \sim N(0, \sigma^2)$ , so that we have a conditional distribution of the form  $p(y_s | x_s) \propto \exp(-\frac{1}{2\sigma^2}(y_s - x_s)^2)$ . We then examined the convergence behavior of both ordinary and reweighted sum-product algorithms for the posterior distribution  $p(x | y) \propto p(x; \theta) \prod_{s=1}^n p(y_s | x_s)$ .

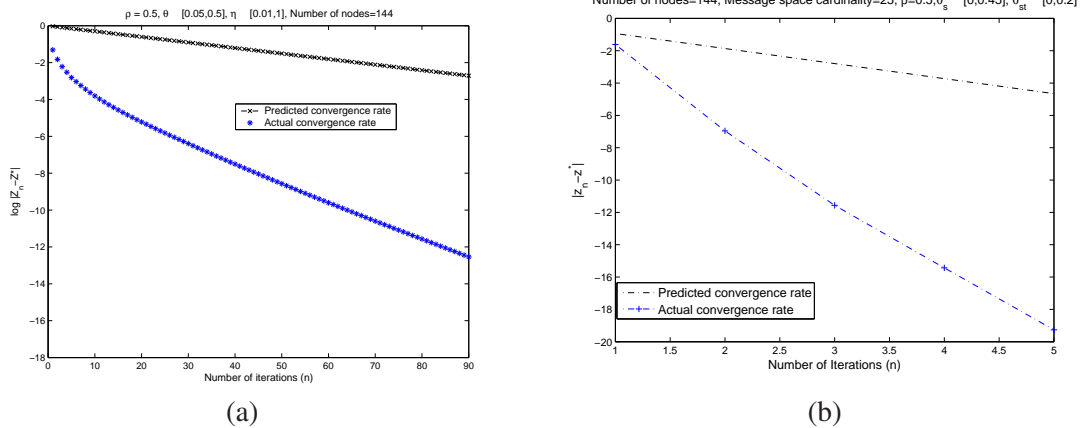
The results in Figure 4 were obtained from a grid with  $n = 100$  nodes, and by varying the observation noise  $\sigma^2$  from  $\sigma^2 = +\infty$  corresponding to  $SNR = 0$ , down to  $\sigma^2 = 0.5$ . For any fixed setting of  $\sigma^2$ , each curve plots the average of the spectral radius bound from Theorem 2 over 20 trials versus the edge strength parameter  $\theta_{ed}$ . Note how the convergence guarantees are substantially strengthened, relative to the case of zero SNR, as the benefits of observations are incorporated.

### B. Convergence rates

We now turn to a comparison of the empirical convergence rates of the reweighted sum-product algorithm to the theoretical upper bounds provided by the inductive norm (18) in Corollary 8. We have performed a large number of simulations for different values of number of nodes, edge weights  $\rho$ , node potentials, edge potentials, and message space sizes. Figure 7 shows a few plots that are representative of our findings, for binary state spaces (panel (a)) and higher order state spaces (panel (b)). The numbers parameterizing the node potentials,  $\theta_s$ , and the edge potentials,  $\theta_{st}$ , are shown on the corresponding plots. As shown in these plots, the convergence rates predicted by Corollary 9 are consistent with the empirical performance, but tend to be overly conservative.

Figure 8 compares the convergence rates predicted by Theorem 2 to the empirical rates in both the *symmetric* and *asymmetric* settings. The symmetric case corresponds to computing the weights  $L_{t \rightarrow s}$  while ignoring the observation potentials, so that the overall matrix  $M$  is symmetric in the edges (i.e.,  $L_{t \rightarrow s} = L_{s \rightarrow t}$ ). The asymmetric case explicitly incorporates the observation potentials,





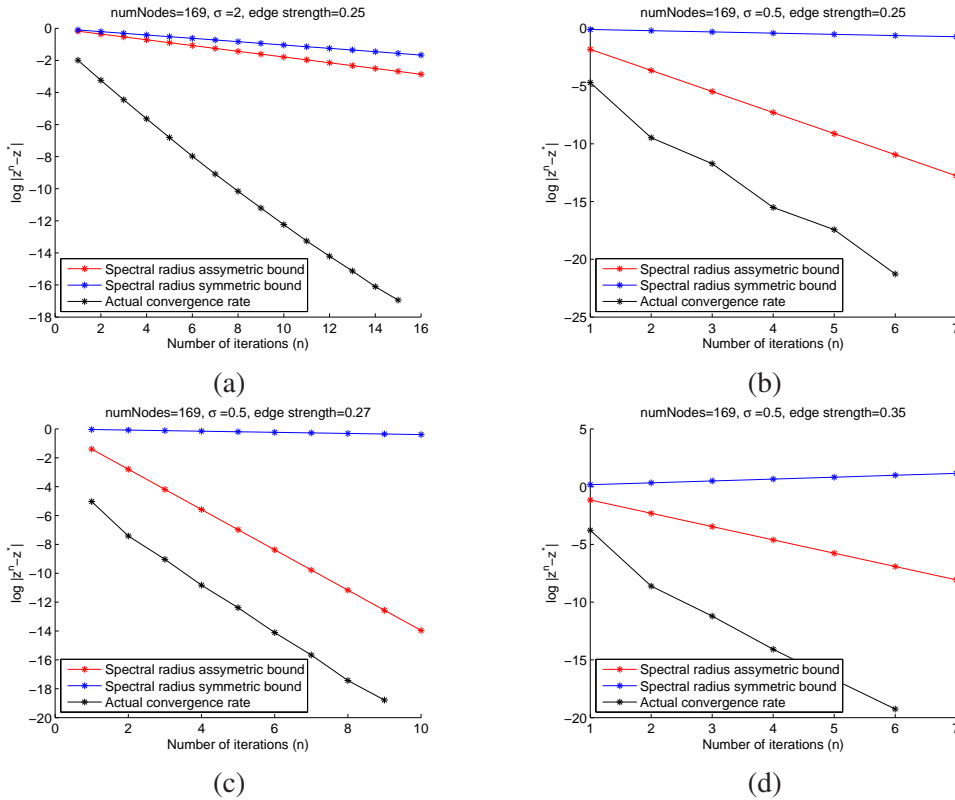
**Fig. 7.** Convergence rates of the reweighted sum-product algorithm as compared to the rate predicted by reweighted norm condition 18. (a) Binary state spaces ( $m = 2$ ). (b) Higher-order spaces ( $m = 25$ ).

and leads to bounds that are as good or better than the symmetric case. Figure 8 illustrates the benefits of including observations in the convergence analysis. Perhaps most strikingly, panel (d) both shows a case where the symmetric bound predicts divergence of the algorithm, whereas the asymmetric bound predicts convergence.

### VI. CONCLUSION AND FUTURE WORK

In this paper, we have studied the convergence and stability properties of the family of reweighted sum-product algorithms in general pairwise Markov random fields. For homogeneous models as well as more general inhomogeneous models, we derived a set of sufficient conditions that ensure convergence, and provide upper bounds on the (geometric) convergence rates. We demonstrated that these sufficient conditions are necessary for homogeneous models without observations, but not in general. We also provided simulation results to complement the theoretical results presented.

There are a number of open questions associated with the results presented here. Even though we have established the benefits of including observation potentials, the conditions provided are still somewhat conservative, since they require that the message updates be contractive at every update, as opposed to requiring that they be attractive in some suitably averaged sense—e.g., when averaged over multiple iterations, or over different updating sequences. An interesting direction would be to derive sharper “average-case” conditions for message-passing convergence. Another open direction concerns analysis of the effect of damped updates, which (from empirical results) can improve



**Fig. 8.** Empirical rates of convergence of the reweighted sum-product algorithm as compared to the rate predicted by the symmetric and asymmetric bounds from Theorem 2.

convergence behavior, especially for reweighted sum-product algorithms. Finally, our current results apply only to the reweighted sum-product algorithm; an interesting open question concerns to what extent similar techniques might be extended to reweighted max-product algorithms.

*Acknowledgments*

We thank the anonymous referees for helpful comments that helped to improve the manuscript.

APPENDIX

**Proof of Lemma 5:**

Setting  $\Delta_{ts} = \sum_{u \in N(t) \setminus s} \rho_{ut} z_{ut} + (1 - \rho_{st}) z_{st}$ , we compute via chain rule

$$\frac{\partial F_{t \rightarrow s}}{\partial z_{ut}}(z) = \begin{cases} \rho_{ut} \frac{\partial F_{t \rightarrow s}}{\partial \Delta_{ts}} & \text{for } u \in N(t) \setminus s, \\ (1 - \rho_{ut}) \frac{\partial F_{t \rightarrow s}}{\partial \Delta_{ts}} & \text{for } u = s. \end{cases}$$

so that it suffices to upper bound  $|\frac{\partial F_{t \rightarrow s}}{\partial \Delta_{ts}}|$ . Computing this partial derivative from the message update (12) yields

$$\frac{\partial F_{t \rightarrow s}}{\partial \Delta_{ts}} = \frac{\exp\left[\frac{2\theta_{st}}{\rho_{st}} + 2\theta_t + \Delta_{ts}\right]}{1 + \exp\left[\frac{2\theta_{st}}{\rho_{st}} + 2\theta_t + \Delta_{ts}\right]} - \frac{\exp[2\theta_t + \Delta_{ts}]}{\exp[2\theta_t + \Delta_{ts}] + \exp\left[\frac{2\theta_{st}}{\rho_{st}}\right]} = G\left(2\theta_t + \Delta_{ts}; \frac{2\theta_{st}}{\rho_{st}}\right)$$

where the function  $G$  was previously defined in (4). Since the message vector  $z^n$  must belong to the box (13) of admissible messages, the vector  $\Delta_{ts}$  must satisfy the bound

$$|\Delta_{ts}| \leq \sum_{u \in N(t) \setminus s} 2|\theta_{ut}| + 2(1 - \rho_{st}) \frac{|\theta_{st}|}{\rho_{st}} := U_{ts}.$$

For any fixed  $\phi$ , the function  $|G(u; \phi)|$  achieves its maximal value  $|G(0; \phi)| = G(0; |\phi|)$  at  $u^* = 0$ .

Noting that by its definition (8), we have  $D_{t \rightarrow s}(\theta; \rho) = 2|\theta_t| - U_{ts}$ , we conclude that

$$\left| \frac{\partial F_{t \rightarrow s}}{\partial \Delta_{ts}} \right| \leq \max_{|\Delta_{ts}| \leq U_{ts}} \left| G\left(2\theta_t + \Delta_{ts}; \frac{2\theta_{st}}{\rho_{st}}\right) \right| = \begin{cases} |G(0; \frac{2|\theta_{st}|}{\rho_{st}})| & \text{if } D_{t \rightarrow s}(\theta; \rho) \leq 0. \\ G\left(D_{t \rightarrow s}(\theta; \rho); 2\frac{|\theta_{st}|}{\rho_{st}}\right) & \text{otherwise.} \end{cases}$$

### Proof of Theorem 7:

We begin by parameterizing the reweighted sum-product messages in terms of the log ratios  $z_{st}(i) := \log \frac{M_{st}(i)}{M_{st}(0)}$ . For each  $i \in \{1, \dots, m-1\}$ , the message updates (2) can be re-written, following some straightforward algebra, in terms of these log messages and the modified potentials  $\tilde{\theta}$  as

$$F_{t \rightarrow s}(z) = \log \frac{1 + \sum_{x_t \neq 0} \exp\left\{\frac{\tilde{\theta}_{ts}(i, x_t)}{\rho_{st}} + \tilde{\theta}_t(x_t) + \Delta_{ts}(x_t)\right\}}{1 + \sum_{x_t \neq 0} \exp\left\{\tilde{\theta}_t(x_t) + \Delta_{ts}(x_t)\right\}}, \quad (31)$$

where  $\Delta_{ts}(x_t) := (1 - \rho_{st}) z_{st}(x_t) + \sum_{v \in N(t) \setminus s} \rho_{vt} z_{vt}(x_t)$ . Analogously to the proof of Lemma 4, we have  $|z_{ts}(i)| \leq \max_{x_t} \left| \frac{\tilde{\theta}_{ts}(i, x_t)}{\rho_{st}} \right|$ . Consequently, each vector  $\Delta_{ts}(x_t)$  must belong the admissible box (26).

As in our earlier proof, we now seek to bound the partial derivatives of the message update  $z_{ts} \leftarrow F_{t \rightarrow s}(z)$ . By chain rule, we have  $\frac{\partial z_{ts}}{\partial z_{vt}} = \rho_{vt} \frac{\partial z_{ts}}{\partial \Delta_{ts}}$  if  $v \in N(t) \setminus s$ , and  $\frac{\partial z_{ts}}{\partial z_{vt}} = (1 - \rho_{st}) \frac{\partial z_{ts}}{\partial \Delta_{ts}}$  if  $v = s$ . Consequently, it suffices to bound  $\frac{\partial z_{ts}}{\partial \Delta_{ts}}$ . Computing the partial derivative of component

$x_s = i$  with respect to message index  $x_t = k$  yields

$$\frac{\partial z_{ts}(i)}{\partial \Delta_{st}(k)} = \frac{\exp\left\{\frac{\tilde{\theta}_{ts}(i,k)}{\rho_{st}} + \tilde{\theta}_t(x_t) + \Delta_{ts}(k)\right\}}{1 + \sum_{x_t \neq 0} \exp\left\{\frac{\tilde{\theta}_{ts}(i,x_t)}{\rho_{st}} + \tilde{\theta}_t(x_t) + \Delta_{ts}(x_t)\right\}} - \frac{\exp\left\{\tilde{\theta}_t(x_t) + \Delta_{ts}(k)\right\}}{1 + \sum_{x_t \neq 0} \exp\left\{\tilde{\theta}_t(x_t) + \Delta_{ts}(x_t)\right\}}$$

Isolating the term involving  $x_t = k$ , we have

$$\begin{aligned} \frac{\partial z_{st}(i)}{\partial \Delta_{st}(k)} &= \frac{\exp\left\{\frac{\tilde{\theta}_{ts}(i,k)}{\rho_{st}} + \tilde{\theta}_t(k) + \Delta_{ts}(k)\right\}}{1 + \sum_{x_t \neq 0, k} \exp\left\{\frac{\tilde{\theta}_{ts}(i,x_t)}{\rho_{st}} + \tilde{\theta}_t(x_t) + \Delta_{ts}(x_t)\right\} + \exp\left\{\frac{\tilde{\theta}_{ts}(i,k)}{\rho_{st}} + \tilde{\theta}_t(k) + \Delta_{ts}(k)\right\}} \\ &\quad - \frac{\exp\left\{\tilde{\theta}_t(k) + \Delta_{ts}(k)\right\}}{1 + \sum_{x_t \neq 0, k} \exp\left\{\tilde{\theta}_t(x_t) + \Delta_{ts}(x_t)\right\} + \exp\left\{\tilde{\theta}_t(k) + \Delta_{ts}(k)\right\}}. \end{aligned}$$

Further simplifying

$$\frac{\partial z_{st}(i)}{\partial \Delta_{st}(k)} = \frac{\exp\left\{\frac{\tilde{\theta}_{ts}(i,k)}{\rho_{st}} - \alpha_{t \rightarrow s}(i, k) + \tilde{\theta}_t(k) + \Delta_{ts}(k)\right\}}{1 + \exp\left\{\frac{\tilde{\theta}_{ts}(i,k)}{\rho_{st}} - \alpha_{t \rightarrow s}(i, k) + \tilde{\theta}_t(k) + \Delta_{ts}(k)\right\}} - \frac{\exp\left\{-\beta_{t \rightarrow s}(i, k) + \tilde{\theta}_t(k) + \Delta_{ts}(k)\right\}}{1 + \exp\left\{-\beta_{t \rightarrow s}(i, k) + \tilde{\theta}_t(k) + \Delta_{ts}(k)\right\}}$$

where  $\alpha_{t \rightarrow s}(\Delta; i, k)$  and  $\beta_{t \rightarrow s}(\Delta; i, k)$  were previously defined.

Setting  $v = -\beta_{t \rightarrow s}(i, k) + \tilde{\theta}_t(k) + \Delta_{ts}(k)$  and  $\varphi = \frac{\tilde{\theta}_{ts}(i,k)}{\rho_{st}} - \alpha_{t \rightarrow s}(i, k) + \beta_{t \rightarrow s}(i, k)$ , we have

$$\begin{aligned} \frac{\partial z_{st}(i)}{\partial \Delta_{st}(k)} &= \frac{\exp(\varphi + v)}{1 + \exp(\varphi + v)} - \frac{\exp(v)}{1 + \exp(v)} \\ &= \frac{\exp(\varphi + v)}{1 + \exp(\varphi + v)} - \frac{\exp\left(v + \frac{\varphi}{2}\right)}{\exp\left(\frac{\varphi}{2}\right) + \exp\left(v + \frac{\varphi}{2}\right)} = G\left(v + \frac{\varphi}{2}; \frac{\varphi}{2}\right), \end{aligned}$$

where  $G$  was previously defined (4). Using the monotonicity properties of  $G$ , we have

$$\left| \frac{\partial z_{st}(i)}{\partial \Delta_{st}(k)} \right| \leq G\left(\left|v + \frac{\varphi}{2}\right|; \left|\frac{\varphi}{2}\right|\right). \quad (32)$$

The claim follows by noting that as defined, we have

$$\psi_{t \rightarrow s}(\vec{v}; k, i) = \left|v + \frac{\varphi}{2}\right| = \frac{1}{2} \left| -\beta_{t \rightarrow s}(\vec{v}; i, k) - \alpha_{t \rightarrow s}(\vec{v}; i, k) + v(k) + \tilde{\theta}_t(k) + \frac{\tilde{\theta}_{ts}(k, i)}{\rho_{st}} \right|, \text{ and}$$

$$\phi_{t \rightarrow s}(\vec{v}; k, i) = \left|\frac{\varphi}{2}\right| = \frac{1}{2} \left| \beta_{t \rightarrow s}(\vec{v}; k, i) - \alpha_{t \rightarrow s}(\vec{v}; k, i) + \frac{\tilde{\theta}_{ts}(k, i)}{\rho_{st}} \right|.$$

## REFERENCES

- [1] A. S. Willsky, "Multiresolution Markov models for signal and image processing," *Proceedings of the IEEE*, vol. 90, no. 8, pp. 1396–1458, 2002.
- [2] H. A. Loeliger, "An introduction to factor graphs," *IEEE Signal Processing Magazine*, vol. 21, pp. 28–41, 2004.
- [3] W. T. Freeman, E. C. Pasztor, and O. T. Carmichael, "Learning low-level vision," in *International Journal of Computer Vision*, October 2000.
- [4] M. J. Wainwright and M. I. Jordan, "Graphical models, exponential families, and variational inference," Tech. Rep., UC Berkeley, Department of Statistics, No. 649, September 2003.
- [5] F. Kschischang, "Codes defined on graphs," *IEEE Signal Processing Magazine*, vol. 41, pp. 118–125, August 2003.
- [6] S. L. Lauritzen and D. J. Spiegelhalter, "Local computations with probabilities on graphical structures and their application to expert systems (with discussion)," *Journal of the Royal Statistical Society B*, vol. 50, pp. 155–224, January 1988.
- [7] J.S. Yedidia, W. T. Freeman, and Y. Weiss, "Constructing free energy approximations and generalized belief propagation algorithms," *IEEE Trans. Info. Theory*, vol. 51, no. 7, pp. 2282–2312, July 2005.
- [8] S. Tatikonda and M. I. Jordan, "Loopy belief propagation and Gibbs measures," in *Proc. Uncertainty in Artificial Intelligence*, August 2002, vol. 18, pp. 493–500.
- [9] M. J. Wainwright, T. S. Jaakkola, and A. S. Willsky, "Tree-based reparameterization framework for analysis of sum-product and related algorithms," *IEEE Trans. Info. Theory*, vol. 49, no. 5, pp. 1120–1146, May 2003.
- [10] Tom Heskes, "On the uniqueness of loopy belief propagation fixed points," *Neural Computation*, vol. 16, no. 11, 2004.
- [11] A. T. Ihler, J. W. Fisher III, and A. S. Willsky, "Loopy belief propagation: Convergence and effects of message errors," *Journal of Machine Learning Research*, vol. 6, pp. 905–936, 2005.
- [12] J. M. Mooij and H. J. Kappen, "Sufficient conditions for convergence of loopy belief propagation," in *Proc. Uncertainty in Artificial Intelligence*, July 2005, vol. 21.
- [13] A. Yuille, "CCCP algorithms to minimize the Bethe and Kikuchi free energies: Convergent alternatives to belief propagation," *Neural Computation*, vol. 14, pp. 1691–1722, 2002.
- [14] M. Welling and Y. Teh, "Belief optimization: A stable alternative to loopy belief propagation," in *Uncertainty in Artificial Intelligence*, July 2001.
- [15] T. Heskes, K. Albers, and B. Kappen, "Approximate inference and constrained optimization," in *Uncertainty in Artificial Intelligence*, July 2003, vol. 13, pp. 313–320.
- [16] M. J. Wainwright, T. S. Jaakkola, and A. S. Willsky, "A new class of upper bounds on the log partition function," *IEEE Trans. Info. Theory*, vol. 51, no. 7, pp. 2313–2335, July 2005.
- [17] W. Wiegand and T. Heskes, "Fractional belief propagation," in *NIPS*, 2002, vol. 12, pp. 438–445.
- [18] W. Wiegand, "Approximations with reweighted generalized belief propagation," in *Workshop on Artificial Intelligence and Statistics*, January 2005.

- [19] A. Levin and Y. Weiss, "Learning to combine bottom-up and top-down segmentation," in *European Conference on Computer Vision (ECCV)*, June 2006.
- [20] J. M. Mooij and H. J. Kappen, "Sufficient conditions for convergence of loopy belief propagation," Tech. Rep. [arxiv.org/abs/cs/0504030v2](https://arxiv.org/abs/cs/0504030v2), University of Nijmegen, Revised in May 2007.
- [21] S. Geman and D. Geman, "Stochastic relaxation, Gibbs distributions, and the Bayesian restoration of images," *IEEE Trans. PAMI*, vol. 6, pp. 721–741, 1984.
- [22] Y. Weiss, "Correctness of local probability propagation in graphical models with loops," *Neural Computation*, vol. 12, pp. 1–41, 2000.
- [23] H. A. Bethe, "Statistics theory of superlattices," *Proc. Royal Soc. London, Series A*, vol. 150, no. 871, pp. 552–575, 1935.
- [24] J. M. Ortega and W. C. Rheinboldt, *Iterative solution of nonlinear equations in several variables*, Classics in applied mathematics. SIAM, New York, 2000.
- [25] D. P. Bertsekas and J. N. Tsitsiklis, *Parallel and Distributed Computation: Numerical Methods*, Athena Scientific, Boston, MA, 1997.
- [26] L. Yu. Kolotilina, "Bounds for the singular values of a matrix involving its sparsity pattern," *Journal of Mathematical Sciences*, vol. 137, 2006.
- [27] G. Rockafellar, *Convex Analysis*, Princeton University Press, Princeton, 1970.