

Message-passing for Graph-structured Linear Programs: Proximal Methods and Rounding Schemes

Pradeep Ravikumar

*Department of Statistics
University of California, Berkeley
Berkeley, CA 94720*

PRADEEPR@STAT.BERKELEY.EDU

Alekh Agarwal

*Department of Electrical Engineering and Computer Sciences
University of California, Berkeley
Berkeley, CA 94720*

ALEKH@CS.BERKELEY.EDU

Martin J. Wainwright

Department of Statistics
University of California, Berkeley
Berkeley, CA 94720*

WAINWRIG@STAT.BERKELEY.EDU

Editor: Yair Weiss

Abstract

The problem of computing a maximum a posteriori (MAP) configuration is a central computational challenge associated with Markov random fields. There has been some focus on “tree-based” linear programming (LP) relaxations for the MAP problem. This paper develops a family of super-linearly convergent algorithms for solving these LPs, based on proximal minimization schemes using Bregman divergences. As with standard message-passing on graphs, the algorithms are distributed and exploit the underlying graphical structure, and so scale well to large problems. Our algorithms have a double-loop character, with the outer loop corresponding to the proximal sequence, and an inner loop of cyclic Bregman projections used to compute each proximal update. We establish convergence guarantees for our algorithms, and illustrate their performance via some simulations. We also develop two classes of rounding schemes, deterministic and randomized, for obtaining integral configurations from the LP solutions. Our deterministic rounding schemes use a “re-parameterization” property of our algorithms so that when the LP solution is integral, the MAP solution can be obtained even before the LP-solver converges to the optimum. We also propose graph-structured randomized rounding schemes applicable to iterative LP-solving algorithms in general. We analyze the performance of and report simulations comparing these rounding schemes.

Keywords: graphical models, MAP Estimation, LP relaxation, proximal minimization, rounding schemes

1. Introduction

A key computational challenge that arises in applications of discrete graphical models is to compute the most probable configuration(s), often referred to as the *maximum a posteriori* (MAP) problem. Although the MAP problem can be solved exactly in polynomial time on trees (and more generally, graphs with bounded treewidth) using the max-product algorithm, it is computationally challenging

*. Also in the Department of Electrical Engineering and Computer Sciences.

for general graphs. Indeed, the MAP problem for general discrete graphical models includes a large number of classical NP-complete problems as special cases, including MAX-CUT, independent set, and various satisfiability problems.

This intractability motivates the development and analysis of methods for obtaining approximate solutions, and there is a long history of approaches to the problem. One class of methods is based on simulated annealing (Geman and Geman, 1984), but the cooling schedules required for theoretical guarantees are often prohibitively slow. Besag (1986) proposed the iterated conditional modes algorithm, which performs a sequence of greedy local maximizations to approximate the MAP solution, but may be trapped by local maxima. Greig et al. (1989) observed that for binary problems with attractive pairwise interactions (the ferromagnetic Ising model in statistical physics terminology), the MAP configuration can be computed in polynomial-time by reduction to a max-flow problem. The ordinary max-product algorithm, a form of non-serial dynamic-programming (Bertele and Brioschi, 1972), computes the MAP configuration exactly for trees, and is also frequently applied to graphs with cycles. Despite some local optimality results (Freeman and Weiss, 2001; Wainwright et al., 2004), it has no general correctness guarantees for graph with cycles, and even worse, it can converge rapidly to non-MAP configurations (Wainwright et al., 2005), even for problems that are easily solved in polynomial time (e.g., ferromagnetic Ising models). For certain graphical models arising in computer vision, Boykov et al. (2001) proposed graph-cut based search algorithms that compute a local maximum over two classes of moves. A broad class of methods are based on the principle of convex relaxation, in which the discrete MAP problem is relaxed to a convex optimization problem over continuous variables. Examples of this convex relaxation problem include linear programming relaxations (Koval and Schlesinger, 1976; Chekuri et al., 2005; Wainwright et al., 2005), as well as quadratic, semidefinite and other conic programming relaxations (for instance, (Ravikumar and Lafferty, 2006; Kumar et al., 2006; Wainwright and Jordan, 2004)).

Among the family of conic programming relaxations, linear programming (LP) relaxation is the least expensive computationally, and also the best understood. The primary focus of this paper is a well-known LP relaxation of the MAP estimation problem for pairwise Markov random fields, one which has been independently proposed by several groups (Koval and Schlesinger, 1976; Chekuri et al., 2005; Wainwright et al., 2005). This LP relaxation is based on optimizing over a set of locally consistent pseudomarginals on edges and vertices of the graph. It is an exact method for any tree-structured graph, so that it can be viewed naturally as a tree-based LP relaxation.¹ The first connection between max-product message-passing and LP relaxation was made by Wainwright et al. (2005), who connected the tree-based LP relaxation to the class of tree-reweighted max-product (TRW-MP) algorithms, showing that TRW-MP fixed points satisfying a strong “tree agreement” condition specify optimal solutions to the LP relaxation.

For general graphs, this first-order LP relaxation could be solved—at least in principle—by various standard algorithms for linear programming, including the simplex and interior-point methods (Bertsimas and Tsitsikilis, 1997; Boyd and Vandenberghe, 2004). However, such generic methods fail to exploit the graph-structured nature of the LP, and hence do not scale favorably to large-scale problems (Yanover et al., 2006). A body of work has extended the connection between the LP relaxation and message-passing algorithms in various ways. Kolmogorov (2005) developed a serial form of TRW-MP updates with certain convergence guarantees; he also showed that there exist fixed points of the TRW-MP algorithm, not satisfying strong tree agreement, that do not cor-

1. In fact, this LP relaxation is the first in a hierarchy of relaxations, based on the treewidth of the graph (Wainwright et al., 2005).

respond to optimal solutions of the LP. This issue has a geometric interpretation, related to the fact that coordinate ascent schemes (to which TRW-MP is closely related), need not converge to the global optima for convex programs that are not strictly convex, but can become trapped in corners. Kolmogorov and Wainwright (2005) showed that this trapping phenomena does not arise for graphical models with binary variables and pairwise interactions, so that TRW-MP fixed points are always LP optimal. Globerson and Jaakkola (2007b) developed a related but different dual-ascent algorithm, which is guaranteed to converge but is not guaranteed to solve the LP. Weiss et al. (2007) established connections between convex forms of the sum-product algorithm, and exactness of reweighted max-product algorithms; Johnson et al. (2007) also proposed algorithms related to convex forms of sum-product. Various authors have connected the ordinary max-product algorithm to the LP relaxation for special classes of combinatorial problems, including matching (Bayati et al., 2005; Huang and Jebara, 2007; Bayati et al., 2007) and independent set (Sanghavi et al., 2007). For general problems, max-product does *not* solve the LP; Wainwright et al. (2005) describe a instance of the MIN-CUT problem on which max-product fails, even though LP relaxation is exact. Other authors (Feldman et al., 2002a; Komodakis et al., 2007) have implemented subgradient methods which are guaranteed to solve the linear program, but such methods typically have sub-linear convergence rates (Bertsimas and Tsitsikilis, 1997).

This paper makes two contributions to this line of work. Our first contribution is to develop and analyze a class of message-passing algorithms with the following properties: their only fixed points are LP-optimal solutions, they are provably convergent with at least a geometric rate, and they have a distributed nature, respecting the graphical structure of the problem. All of the algorithms in this paper are based on the well-established idea of *proximal minimization*: instead of directly solving the original linear program itself, we solve a sequence of so-called proximal problems, with the property that the sequence of associated solutions is guaranteed to converge to the LP solution. We describe different classes of algorithms, based on different choices of the proximal function: quadratic, entropic, and tree-reweighted entropies. For all choices, we show how the intermediate proximal problems can be solved by forms of message-passing on the graph that are similar to but distinct from the ordinary max-product or sum-product updates. An additional desirable feature, given the wide variety of lifting methods for further constraining LP relaxations (Wainwright and Jordan, 2003), is that new constraints can be incorporated in a relatively seamless manner, by introducing new messages to enforce them.

Our second contribution is to develop various types of rounding schemes that allow for early termination of LP-solving algorithms. There is a substantial body of past work (e.g., Raghavan and Thompson, 1987) on rounding fractional LP solutions so as to obtain integral solutions with approximation guarantees. Our use of rounding is rather different: instead, we consider rounding schemes applied to problems for which the LP solution is integral, so that rounding would be unnecessary if the LP were solved to optimality. In this setting, the benefit of certain rounding procedures (in particular, those that we develop) is allowing an LP-solving algorithm to be terminated *before* it has solved the LP, while still returning the MAP configuration, either with a deterministic or high probability guarantee. Our deterministic rounding schemes apply to a class of algorithms which, like the proximal minimization algorithms that we propose, maintain a certain invariant of the original problem. We also propose and analyze a class of graph-structured randomized rounding procedures that apply to any algorithm that approaches the optimal LP solution from the interior of the relaxed polytope. We analyze these rounding schemes, and give finite bounds on the number of iterations required for the rounding schemes to obtain an integral MAP solution.

The remainder of this paper is organized as follows. We begin in Section 2 with background on Markov random fields, and the first-order LP relaxation. In Section 3, we introduce the notions of proximal minimization and Bregman divergences, then derive various of message-passing algorithms based on these notions, and finally discuss their convergence properties. Section 4 is devoted to the development and analysis of rounding schemes, both for our proximal schemes as well as other classes of LP-solving algorithms. We provide experimental results in Section 5, and conclude with a discussion in Section 6.

2. Background

We begin by introducing some background on Markov random fields, and the LP relaxations that are the focus of this paper. Given a discrete space $\mathcal{X} = \{0, 1, 2, \dots, m-1\}$, let $X = (X_1, \dots, X_N) \in \mathcal{X}^N$ denote a N -dimensional discrete random vector. (While we have assumed the variables take values in the same set \mathcal{X} , we note that our results easily generalize to the case where the variables take values in different sets with differing cardinalities.) We assume that the distribution \mathbb{P} of the random vector is a Markov random field, meaning that it factors according to the structure of an undirected graph $G = (V, E)$, with each variable X_s associated with one node $s \in V$, in the following way. Letting $\theta_s : \mathcal{X} \rightarrow \mathbb{R}$ and $\theta_{st} : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$ be singleton and edgewise potential functions respectively, we assume that the distribution takes the form

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

The problem of *maximum a posteriori* (MAP) estimation is to compute a configuration with maximum probability—that is, an element

$$x^* \in \arg \max_{x \in \mathcal{X}^N} \left\{ \sum_{s \in V} \theta_s(x_s) + \sum_{(s,t) \in E} \theta_{st}(x_s, x_t) \right\}, \quad (1)$$

where the $\arg \max$ operator extracts the configurations that achieve the maximal value. This problem is an integer program, since it involves optimizing over the discrete space \mathcal{X}^N . For future reference, we note that the functions $\theta_s(\cdot)$ and $\theta_{st}(\cdot)$ can always be represented in the form

$$\begin{aligned} \theta_s(x_s) &= \sum_{j \in \mathcal{X}} \theta_{s;j} \mathbb{I}[x_s = j], \\ \theta_{st}(x_s, x_t) &= \sum_{j,k \in \mathcal{X}} \theta_{st;jk} \mathbb{I}[x_s = j; x_t = k], \end{aligned}$$

where the m -vectors $\{\theta_{s;j}, j \in \mathcal{X}\}$ and $m \times m$ matrices $\{\theta_{st;jk}, (j,k) \in \mathcal{X} \times \mathcal{X}\}$ parameterize the problem.

The first-order linear programming (LP) relaxation (Koval and Schlesinger, 1976; Chekuri et al., 2005; Wainwright et al., 2005) of this problem is based on a set of pseudomarginals μ_s and μ_{st} , associated with the nodes and vertices of the graph. These pseudomarginals are constrained to be non-negative, as well to normalize and be locally consistent in the following sense:

$$\sum_{x_s \in \mathcal{X}} \mu_s(x_s) = 1, \quad \text{for all } s \in V, \text{ and} \quad (2)$$

$$\sum_{x_t \in \mathcal{X}} \mu_{st}(x_s, x_t) = \mu_s(x_s) \quad \text{for all } (s,t) \in E, x_s \in \mathcal{X}. \quad (3)$$

The polytope defined by the non-negativity constraints $\mu \geq 0$, the normalization constraints (2) and the marginalization constraints (3), is denoted by $\mathbb{L}(G)$. The LP relaxation is based on maximizing the linear function

$$\langle \theta, \mu \rangle := \sum_{s \in V} \sum_{x_s} \theta_s(x_s) \mu_s(x_s) + \sum_{(s,t) \in E} \sum_{x_s, x_t} \theta_{st}(x_s, x_t) \mu_{st}(x_s, x_t), \quad (4)$$

subject to the constraint $\mu \in \mathbb{L}(G)$.

In the sequel, we write the linear program (4) more compactly in the form $\max_{\mu \in \mathbb{L}(G)} \langle \theta, \mu \rangle$. By construction, this relaxation is guaranteed to be exact for any problem on a tree-structured graph (Wainwright et al., 2005), so that it can be viewed as a tree-based relaxation. The main goal of this paper is to develop efficient and distributed algorithms for solving this LP relaxation,² as well as strengthenings based on additional constraints. For instance, one natural strengthening is by “lifting”: view the pairwise MRF as a particular case of a more general MRF with higher order cliques, define higher-order pseudomarginals on these cliques, and use them to impose higher-order consistency constraints. This particular progression of tighter relaxations underlies the Bethe to Kikuchi (sum-product to generalized sum-product) hierarchy (Yedidia et al., 2005); see Wainwright and Jordan (2003) for further discussion of such LP hierarchies.

3. Proximal Minimization Schemes

We begin by defining the notion of a proximal minimization scheme, and various types of divergences (among them Bregman) that we use to define our proximal sequences. Instead of dealing with the maximization problem $\max_{\mu \in \mathbb{L}(G)} \langle \theta, \mu \rangle$, it is convenient to consider the equivalent minimization problem,

$$\min_{\mu \in \mathbb{L}(G)} -\langle \theta, \mu \rangle.$$

3.1 Proximal Minimization

The class of methods that we develop are based on the notion of proximal minimization (Bertsekas and Tsitsiklis, 1997). Instead of attempting to solve the LP directly, we solve a sequence of problems of the form

$$\mu^{n+1} = \arg \min_{\mu \in \mathbb{L}(G)} \left\{ -\langle \theta, \mu \rangle + \frac{1}{\omega^n} D_f(\mu \| \mu^n) \right\}, \quad (5)$$

where for iteration numbers $n = 0, 1, 2, \dots$, the vector μ^n denotes current iterate, the quantity ω^n is a positive weight, and D_f is a generalized distance function, known as the proximal function. (Note that we are using superscripts to represent the iteration number, *not* for the power operation.)

The purpose of introducing the proximal function is to convert the original LP—which is convex but not strictly so—into a strictly convex problem. The latter property is desirable for a number of reasons. First, for strictly convex programs, coordinate descent schemes are guaranteed to converge to the global optimum; note that they may become trapped for non-strictly convex problems, such as the piecewise linear surfaces that arise in linear programming. Moreover, the dual of a strictly convex problem is guaranteed to be differentiable (Bertsekas, 1995); a guarantee which need not hold

2. The relaxation could fail to be exact though, in which case the optimal solution to the relaxed problem will be suboptimal on the original MAP problem

for non-strictly convex problems. Note that differentiable dual functions can in general be solved more easily than non-differentiable dual functions. In the sequel, we show how for appropriately chosen generalized distances, the proximal sequence $\{\mu^n\}$ can be computed using message passing updates derived from cyclic projections.

We note that the proximal scheme (5) is similar to an annealing scheme, in that it involves perturbing the original cost function, with a choice of weights $\{\omega^n\}$. While the weights $\{\omega^n\}$ can be adjusted for faster convergence, they can also be set to a constant, unlike for standard annealing procedures in which the annealing weight is taken to 0. The reason is that $D_f(\mu \parallel \mu^{(n)})$, as a generalized distance, itself converges to zero as the algorithm approaches the optimum, thus providing an “adaptive” annealing. For appropriate choice of weights and proximal functions, these proximal minimization schemes converge to the LP optimum with at least geometric and possibly superlinear rates (Bertsekas and Tsitsiklis, 1997; Iusem and Teboulle, 1995).

In this paper, we focus primarily on proximal functions that are Bregman divergences (Censor and Zenios, 1997), a class that includes various well-known divergences, among them the squared ℓ_2 -distance and norm, and the Kullback-Leibler divergence. We say that a function $f : S \mapsto \mathbb{R}$, with domain $S \subseteq \mathbb{R}^p$, is a *Bregman function* if $\text{int } S \neq \emptyset$ and it is continuously differentiable and strictly convex on $\text{int } S$. Any such function induces a *Bregman divergence* $D_f : S \times \text{int } S \mapsto \mathbb{R}$ as follows:

$$D_f(\mu' \parallel \mathbf{v}) := f(\mu') - f(\mathbf{v}) - \langle \nabla f(\mathbf{v}), \mu' - \mathbf{v} \rangle. \tag{6}$$

Figure 1 illustrates the geometric interpretation of this definition in terms of the tangent approximation. A Bregman divergence satisfies $D_f(\mu' \parallel \mathbf{v}) \geq 0$ with equality if and only if $\mu' = \mathbf{v}$, but need not be symmetric or satisfy the triangle inequality, so it is only a generalized distance. Further restrictions on the inducing function f are thus required for the divergence to be “well-behaved,” for instance that it satisfy the property that for any sequence $\mathbf{v}^n \rightarrow \mathbf{v}^*$, where $\mathbf{v}^n \in \text{int } S$, $\mathbf{v}^* \in S$, then $D_f(\mathbf{v}^* \parallel \mathbf{v}^n) \rightarrow 0$. Censor and Zenios (1997) impose such technical conditions explicitly in their definition of a Bregman function; in this paper, we impose the stronger yet more easily stated condition that the Bregman function f (as defined above) be of Legendre type (Rockafellar, 1970). In this case, in addition to the Bregman function properties, it satisfies the following property: for any sequence $\mu^n \rightarrow \mu^*$ where $\mu^n \in \text{int } S$, $\mu^* \in \partial S$, it holds that $\|\nabla f(\mu^n)\| \rightarrow +\infty$. Further, we assume that the range $\nabla f(\text{int } S) = \mathbb{R}^p$.

Let us now look at some choices of divergences, proximal minimizations (5) based on which we will be studying in the sequel.

3.1.1 QUADRATIC DIVERGENCE

This choice is the simplest, and corresponds to setting the inducing Bregman function f in (6) to be the quadratic function

$$q(\mu) := \frac{1}{2} \left\{ \sum_{s \in V} \sum_{x_s \in \mathcal{X}} \mu_s^2(x_s) + \sum_{(s,t) \in E} \sum_{(x_s, x_t) \in \mathcal{X} \times \mathcal{X}} \mu_{st}^2(x_s, x_t) \right\},$$

defined over nodes and edges of the graph. The divergence is then simply the quadratic norm across nodes and edges

$$Q(\mu \parallel \mathbf{v}) := \frac{1}{2} \sum_{s \in V} \|\mu_s - \mathbf{v}_s\|^2 + \frac{1}{2} \sum_{(s,t) \in E} \|\mu_{st} - \mathbf{v}_{st}\|^2, \tag{7}$$

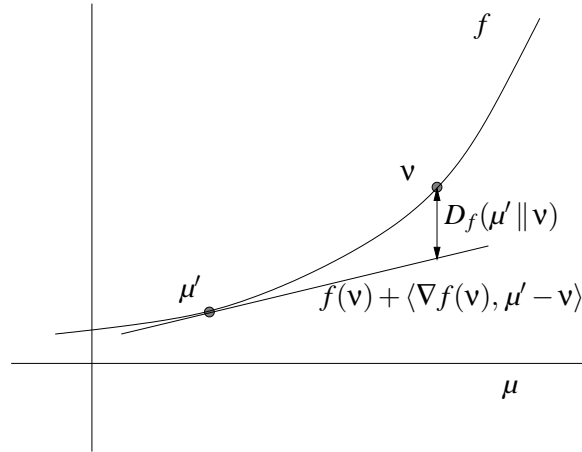


Figure 1: Graphical illustration of a Bregman divergence.

where we have used the shorthand $\|\mu_s - \nu_s\|^2 = \sum_{x_s \in \mathcal{X}} |\mu_s(x_s) - \nu_s(x_s)|^2$, with similar notation for the edges.

3.1.2 WEIGHTED ENTROPIC DIVERGENCE

Another choice for the inducing Bregman function is the weighted sum of negative entropies

$$\bar{H}_\alpha(\mu) = \sum_{s \in V} \alpha_s \bar{H}_s(\mu_s) + \sum_{(s,t) \in E} \alpha_{st} \bar{H}_{st}(\mu_{st}), \quad (8)$$

where \bar{H}_s and \bar{H}_{st} are defined by

$$\begin{aligned} \bar{H}_s(\mu_s) &:= \sum_{x_s \in \mathcal{X}} \mu_s(x_s) \log \mu_s(x_s), \text{ and} \\ \bar{H}_{st}(\mu_{st}) &:= \sum_{(x_s, x_t) \in \mathcal{X} \times \mathcal{X}} \mu_{st}(x_s, x_t) \log \mu_{st}(x_s, x_t), \end{aligned}$$

corresponding to the node-based and edge-based negative entropies, respectively. The corresponding Bregman divergence is a weighted sum of Kullback-Leibler (KL) divergences across the nodes and edges. In particular, letting $\alpha_s > 0$ and $\alpha_{st} > 0$ be positive weights associated with node s and edge (s, t) respectively, we define

$$D_\alpha(\mu \| \nu) = \sum_{s \in V} \alpha_s D(\mu_s \| \nu_s) + \sum_{(s,t) \in E} \alpha_{st} D(\mu_{st} \| \nu_{st}), \quad (9)$$

where $D(p \| q) := \sum_x (p(x) \log \frac{p(x)}{q(x)} - [p(x) - q(x)])$ is the KL divergence. An advantage of the KL divergence, relative to the quadratic norm, is that it automatically acts to enforce non-negativity constraints on the pseudomarginals in the proximal minimization problem. (See Section 3.4 for a more detailed discussion of this issue.)

3.1.3 TREE-REWEIGHTED ENTROPIC DIVERGENCE

Our last example is based on a *tree-reweighted* entropy. The notion of a tree-reweighted entropy was first proposed by Wainwright et al. (2002). Their entropy function however while a Bregman

function is not of the Legendre type. Nonetheless let us first describe their proposed function. Given a set \mathcal{T} of spanning trees $T = (V, E(T))$, and a probability distribution ρ over \mathcal{T} , we can obtain edge weights $\rho_{st} \in (0, 1]$ for each edge (s, t) of the graph G as $\rho_{st} = \sum_{T \in \mathcal{T}} \mathbb{I}((s, t) \in E)$. Given such edge weights, define

$$f_{\text{trw}}(\mu) := \sum_{s \in V} \bar{H}_s(\mu_s) + \sum_{(s,t) \in E} \rho_{st} I_{st}(\mu_{st}), \quad (10)$$

where \bar{H} is the negative entropy as defined earlier, while the quantity I_{st} defined as

$$I_{st}(\mu_{st}) := \sum_{(x_s, x_t) \in \mathcal{X} \times \mathcal{X}} \mu_{st}(x_s, x_t) \log \frac{\mu_{st}(x_s, x_t)}{[\sum_{x'_t} \mu_{st}(x_s, x'_t)][\sum_{x'_s} \mu_{st}(x'_s, x_t)]},$$

is the mutual information associated with edge (s, t) . It can be shown that the function f_{trw} is strictly convex and continuously differentiable when restricted to $\mu \in \mathbb{L}(G)$; and in particular that it is a Bregman function with domain $\mathbb{L}(G)$. Within its domain $\mathbb{L}(G)$, the function can be re-expressed as a weighted negative entropy family (8),

$$f_{\text{trw}}(\mu) = \sum_{s \in V} (1 - \sum_{t: (s,t) \in E} \rho_{st}) \bar{H}_t(\mu_t) + \sum_{(s,t) \in E} \rho_{st} \bar{H}_{st}(\mu_{st}),$$

but where the node entropy weights $\alpha_s := 1 - \sum_{t: (s,t) \in E} \rho_{st}$ are not always positive. The corresponding Bregman divergence belongs to the weighted entropic family (9), with node weights α_s defined above, and edge-weights $\alpha_{st} = \rho_{st}$. However as stated above, this tree-reweighted entropy function is not of Legendre type, and hence is not admissible for our proximal minimization procedure (5).

However, Globerson and Jaakkola (2007a) proposed an alternative tree reweighted entropy that while equal to $f_{\text{trw}}(\mu)$ for $\mu \in \mathbb{L}(G)$ is yet convex for all μ (not just when restricted to $\mathbb{L}(G)$). Their proposed function is described as follows. For each undirected edge in E , construct two oriented edges in both directions; denote the set of these oriented edges by \bar{E} . Then given node weights $\rho_{os} \in (0, 1]$ for each node $s \in V$, and edge weights $\rho_{s|t} \in (0, 1]$ for oriented edges $(t \rightarrow s) \in \bar{E}$, define

$$f_{\text{otw}}(\mu) := \sum_{s \in V} \rho_{os} \bar{H}_s(\mu_s) + \sum_{(t \rightarrow s) \in \bar{E}} \rho_{s|t} \bar{H}_{s|t}(\mu_{st}), \quad (11)$$

where the quantity $\bar{H}_{s|t}$ defined as

$$\bar{H}_{s|t}(\mu_{st}) := \sum_{(x_s, x_t) \in \mathcal{X} \times \mathcal{X}} \mu_{st}(x_s, x_t) \log \frac{\mu_{st}(x_s, x_t)}{\sum_{x'_s} \mu_{st}(x'_s, x_t)},$$

is the conditional entropy of X_s given X_t with respect to the joint distribution μ_{st} . It can be shown that this oriented tree-reweighted entropy is not only a Bregman function with domain the non-negative orthant \mathbb{R}_+^p , but is also of Legendre type, so that it is indeed admissible for our proximal minimization procedure. The corresponding divergence is given as,

$$D_\rho(\mu \| \nu) = \sum_{s \in V} \rho_{os} D(\mu_s \| \nu_s) + \sum_{t \rightarrow s \in \bar{E}} \rho_{s|t} (D(\mu_{st} \| \nu_{st}) + \tilde{D}(\mu_{st} \| \nu_{st})),$$

where $D(p \| q)$ is the KL divergence, and $\tilde{D}(\cdot \| \cdot)$ is a KL divergence like term, defined as

$$\begin{aligned} \tilde{D}(\mu_{st} \| \nu_{st}) := & \sum_{(x_s, x_t) \in \mathcal{X} \times \mathcal{X}} \mu_{st}(x_s, x_t) \log \frac{[\sum_{x'_s} \nu_{st}(x'_s, x_t)]}{[\sum_{x'_s} \mu_{st}(x'_s, x_t)]} \\ & + \frac{\nu_{st}(x_s, x_t)}{[\sum_{x'_s} \nu_{st}(x'_s, x_t)]} [\mu_{st}(x_s, x_t) - \nu_{st}(x_s, x_t)]. \end{aligned}$$

3.2 Proximal Sequences via Bregman Projection

The key in designing an efficient proximal minimization scheme is ensuring that the proximal sequence $\{\mu^n\}$ can be computed efficiently. In this section, we first describe how sequences of proximal minimizations (when the proximal function is a Bregman divergence) can be reformulated as a particular Bregman projection. We then describe how this Bregman projection can itself be computed iteratively, in terms of a sequence of cyclic Bregman projections (Censor and Zenios, 1997) based on a decomposition of the constraint set $\mathbb{L}(G)$. In the sequel, we then show how these cyclic Bregman projections reduce to very simple message-passing updates.

Given a Bregman divergence D , the *Bregman projection* of a vector \mathbf{v} onto a convex set C is given by

$$\hat{\mu} := \arg \min_{\mu \in C} D_f(\mu \| \mathbf{v}). \quad (12)$$

That this minimum is achieved and is unique follows from our assumption that the function f is of Legendre type and from Theorem 3.12 in Bauschke and Borwein (1997), so that the projection is well-defined. We define the projection operator

$$\Pi_C(\mathbf{v}) := \arg \min_{\mu \in C} D_f(\mu \| \mathbf{v}), \quad (13)$$

where we have suppressed the dependence on the Bregman function f in the notation. When the constraint set $C = \cap_{i=1}^M C_i$ is an intersection of simpler constraint sets, then a candidate algorithm for the Bregman projection is to compute it in a *cyclic manner*: by iteratively projecting onto the simple constraint sets $\{C_i\}$ (Censor and Zenios, 1997). Define the sequence

$$\mu^{t+1} = \Pi_{C_{i(t)}}(\mu^t),$$

for some control sequence parameter $i : \mathbb{N} \mapsto \{1, \dots, M\}$ that takes each output value an infinite number of times, for instance $i(t) = t \bmod M$. It can be shown that *when the constraints are affine* then such cyclic Bregman projections μ^t converge to the projection $\hat{\mu}$ onto the entire constraint set as defined in (12) so that $\mu^t \rightarrow \hat{\mu}$ (Censor and Zenios, 1997). But when a constraint C_i is non-affine, the individual projection would have to be followed by a correction (Dykstra, 1985; Han, 1988; Censor and Zenios, 1997) in order for such convergence to hold. In Appendix A we have outlined these corrections briefly for the case where the constraints are linear inequalities. For ease of notation, we will now subsume these corrections into the iterative projection notation, $\mu^{t+1} = \Pi_{C_{i(t)}}(\mu^t)$, so that the notation assumes that the Bregman projections are suitably corrected when the constraints $C_{i(t)}$ are non-affine. In this paper, other than positivity constraints, we will be concerned only with affine constraints, for which no corrections are required.

Let us now look at the stationary condition characterizing the optimum $\hat{\mu}$ of (12). As shown in for instance Bertsekas (1995), the optimum $\hat{\mu}$ of any constrained optimization problem $\min_{\mu \in C} g(\mu)$ is given by the stationary condition,

$$\langle \nabla g(\hat{\mu}), \mu - \hat{\mu} \rangle \geq 0, \quad (14)$$

for all $\mu \in C$. For the projection problem (12), the gradient of the objective $D_f(\mu \| \mathbf{v}) := f(\mu) - f(\mathbf{v}) - \langle \nabla f(\mathbf{v}), \mu - \mathbf{v} \rangle$ with respect to the first argument μ is given by $\nabla f(\mu) - \nabla f(\mathbf{v})$, which when substituted in (14) yields the stationary condition of the optimum $\hat{\mu}$ as

$$\langle \nabla f(\hat{\mu}) - \nabla f(\mathbf{v}), \mu - \hat{\mu} \rangle \geq 0, \quad (15)$$

for all $\mu \in C$. Now consider the proximal minimization problem to be solved at step n , namely the strictly convex problem

$$\min_{\mu \in \mathbb{L}(G)} \left\{ -\langle \theta, \mu \rangle + \frac{1}{\omega^n} D_f(\mu \| \mu^n) \right\}. \quad (16)$$

Solving for the derivative of this objective with respect to μ as $-\theta + \frac{1}{\omega^n}(\nabla f(\mu) - \nabla f(\mu^n))$, and substituting in (14), we obtain the conditions defining the optimum μ^{n+1} as

$$\langle \nabla f(\mu^{n+1}) - \nabla f(\mu^n) - \omega^n \theta, \mu - \mu^{n+1} \rangle \geq 0,$$

for all $\mu \in \mathbb{L}(G)$. Comparing these with the conditions for Bregman projection (15), we see that if there exists a vector v such that

$$\nabla f(v) = \nabla f(\mu^n) + \omega^n \theta, \quad (17)$$

then the proximal iterate μ^{n+1} is the Bregman projection of this vector v onto the set $\mathbb{L}(G)$. As shown in Bauschke and Borwein (1997), for any function f of Legendre type with domain S , the gradient ∇f is a one-to-one function with domain $\text{int } S$, so that its inverse $(\nabla f)^{-1}$ is a well-defined function on the range $\nabla f(\text{int } S)$ of ∇f . Since we have assumed that this range is \mathbb{R}^p , we can thus obtain the unique v which satisfies the condition in (17) as $v = (\nabla f)^{-1}(\nabla f(\mu^n) + \omega^n \theta)$ (Note that the range constraint could be relaxed to only require that the range of ∇f be a cone containing θ). Accordingly, we set up the following notation: for any Bregman function f , induced divergence D_f , and convex set C , we define the operator

$$J_f(\mu, v) := (\nabla f)^{-1}(\nabla f(\mu) + v).$$

We can then write the proximal update (16) in a compact manner as the compounded operation

$$\mu^{n+1} = \Pi_{\mathbb{L}(G)}(J_f(\mu^n, \omega^n \theta)).$$

Consequently, efficient algorithms for computing the Bregman projection (12) can be leveraged to compute the proximal update (16). In particular, we consider a decomposition of the constraint set as an intersection— $\mathbb{L}(G) = \cap_{k=1}^M \mathbb{L}_k(G)$ —and then apply the method of cyclic Bregman projections discussed above. Initializing $\mu^{n,0} = \mu^n$ and updating from $\mu^{n,\tau} \mapsto \mu^{n,\tau+1}$ by projecting $\mu^{n,\tau}$ onto constraint set $\mathbb{L}_{i(\tau)}(G)$, where $i(\tau) = \tau \bmod M$, for instance, we obtain the meta-algorithm summarized in Algorithm 1.

As shown in the following sections, by using a decomposition of $\mathbb{L}(G)$ over the edges of the graph, the inner loop steps correspond to local message-passing updates, slightly different in nature depending on the choice of Bregman distance. Iterating the inner and outer loops yields a provably convergent message-passing algorithm for the LP. Convergence follows from the convergence properties of proximal minimization (Bertsekas and Tsitsiklis, 1997), combined with convergence guarantees for cyclic Bregman projections (Censor and Zenios, 1997). In the following section, we derive the message-passing updates corresponding to various Bregman functions of interest.

3.3 Quadratic Projections

Consider the proximal sequence with the quadratic distance Q from Equation (7); the Bregman function inducing this distance is the quadratic function $q(y) = \frac{1}{2}y^2$, with gradient $\nabla q(y) = y$. A little calculation shows that the operator J_q takes the form

$$J_q(\mu, \omega\theta) = \mu + \omega\theta,$$

Algorithm 1 Basic proximal-Bregman LP solver

Given a Bregman distance D , weight sequence $\{\omega^n\}$ and problem parameters θ :

- Initialize μ^0 to the uniform distribution: $\mu_s^{(0)}(x_s) = \frac{1}{m}$, $\mu_{st}^{(0)}(x_s, x_t) = \frac{1}{m^2}$.
 - **Outer Loop:** For iterations $n = 0, 1, 2, \dots$, update $\mu^{n+1} = \Pi_{\mathbb{L}(G)}\left(\mathbf{J}_f(\mu^n, \omega^n \theta)\right)$.
 - Solve Outer Loop via **Inner Loop:**
 - (a) Inner initialization $\mu^{n,0} = \mathbf{J}_f(\mu^n, \omega^n \theta)$.
 - (b) For $t = 0, 1, 2, \dots$, set $i(t) = t \bmod M$.
 - (c) Update $\mu^{n,t+1} = \Pi_{\mathbb{L}_{i(t)}(G)}(\mu^{n,t})$.
-

whence we obtain the initialization in Equation (19).

We now turn to the projections $\mu^{n,\tau+1} = \Pi_q(\mu^{n,\tau}, \mathbb{L}_i(G))$ onto the individual constraints $\mathbb{L}_i(G)$. For each such constraint, the local update is based on the solving the problem

$$\mu^{n,\tau+1} = \arg \min_{\mathbf{v} \in \mathbb{L}_i(G)} \left\{ q(\mathbf{v}) - \langle \mathbf{v}, \nabla q(\mu^{n,\tau}) \rangle \right\}. \quad (18)$$

In Appendix B.1, we show how the solution to these inner updates takes the form (20) given in Algorithm 2. The $\{Z_s, Z_{st}\}$ variables correspond to the dual variables used to correct the Bregman projections for positivity (and hence inequality) constraints, as outlined in (33) in Section 3.2.

3.4 Entropic Projections

Consider the proximal sequence with the Kullback-Leibler distance $D(\mu \parallel \mathbf{v})$ defined in Equation (9). The Bregman function h_α inducing the distance is a sum of negative entropy functions $f(\mu) = \mu \log \mu$, and its gradient is given by $\nabla f(\mu) = \log(\mu) + \vec{1}$. In this case, some calculation shows that the map $\mathbf{v} = \mathbf{J}_f(\mu, \omega \theta)$ is given by

$$\mathbf{v} = \mu \exp(\omega \theta / \alpha),$$

whence we obtain the initialization Equation (21). In Appendix B.2, we derive the message-passing updates summarized in Algorithm 3.

3.5 Tree-reweighted Entropy Proximal Sequences

In the previous sections, we saw how to solve the proximal sequences following the algorithmic template 1 and using message passing updates derived from cyclic Bregman projections. In this section, we show that for the tree-reweighted entropic divergences (11), in addition to the cyclic Bregman projection recipe of the earlier sections, we can also use tree-reweighted sum-product or related methods (Globerson and Jaakkola, 2007b; Hazan and Shashua, 2008) to compute the proximal sequence.

Algorithm 2 Quadratic Messages for μ^{n+1}

Initialization:

$$\begin{aligned}
 \mu_{st}^{(n,0)}(x_s, x_t) &= \mu_{st}^{(n)}(x_s, x_t) + w^n \theta_{st}(x_s, x_t), \\
 \mu_s^{(n,0)}(x_s) &= \mu_s^{(n)}(x_s) + w^n \theta_s(x_s). \\
 Z_s(x_s) &= \mu_s^{(n,0)}(x_s), \\
 Z_{st}(x_s, x_t) &= \mu_{st}^{(n)}(x_s, x_t).
 \end{aligned} \tag{19}$$

repeat
for each edge $(s, t) \in E$ **do**

$$\begin{aligned}
 \mu_{st}^{(n,\tau+1)}(x_s, x_t) &= \mu_{st}^{(n,\tau)}(x_s, x_t) + \left(\frac{1}{m+1} \right) \left(\mu_s^{(n,\tau)}(x_s) - \sum_{x_t} \mu_{st}^{(n,\tau)}(x_s, x_t) \right), \\
 \mu_s^{(n,\tau+1)}(x_s) &= \mu_s^{(n,\tau)}(x_s) + \left(\frac{1}{m+1} \right) \left(-\mu_s^{(n,\tau)}(x_s) + \sum_{x_t} \mu_{st}^{(n,\tau)}(x_s, x_t) \right), \\
 C_{st}(x_s, x_t) &= \min\{Z_{st}(x_s, x_t), \mu_{st}^{(n,\tau+1)}(x_s, x_t)\}, \\
 Z_{st}(x_s, x_t) &= Z_{st}(x_s, x_t) - C_{st}(x_s, x_t), \\
 \mu_{st}^{(n,\tau+1)}(x_s, x_t) &= \mu_{st}^{(n,\tau+1)}(x_s, x_t) - C_{st}(x_s, x_t).
 \end{aligned} \tag{20}$$

end for
for each node $s \in V$ **do**

$$\begin{aligned}
 \mu_s^{(n,\tau+1)}(x_s) &= \mu_s^{(n,\tau)}(x_s) + \frac{1}{m} \left(1 - \sum_{x_s} \mu_s^{(n,\tau)}(x_s) \right), \\
 C_s(x_s) &= \min\{Z_s(x_s), \mu_s^{(n,\tau+1)}(x_s)\}, \\
 Z_s(x_s) &= Z_s(x_s) - C_s(x_s), \\
 \mu_s^{(n,\tau+1)}(x_s) &= \mu_s^{(n,\tau+1)}(x_s) - C_s(x_s).
 \end{aligned}$$

end for
until convergence

Recall the proximal sequence optimization problem (5) written as

$$\begin{aligned}
 \mu^{n+1} &= \arg \min_{\mathbf{v} \in \mathbb{L}(G)} \left\{ -\langle \theta, \mathbf{v} \rangle + \frac{1}{\omega^n} D_f(\mathbf{v} \| \mu^n) \right\} \\
 &= \arg \min_{\mathbf{v} \in \mathbb{L}(G)} \left\{ -\langle \theta, \mathbf{v} \rangle + \frac{1}{\omega^n} (f(\mathbf{v}) - f(\mu^n) - \langle \nabla f(\mu^n), \mathbf{v} - \mu^n \rangle) \right\}.
 \end{aligned} \tag{24}$$

 Let us denote $\theta^n := \omega^n \theta + \nabla f(\mu^n)$, and set the Bregman function f to the tree-reweighted entropy f_{trw} defined in (10) (or equivalently the oriented tree-reweighted entropy f_{otw} (11) since both are

Algorithm 3 Entropic Messages for μ^{n+1}

Initialization:

$$\begin{aligned}\mu_{st}^{(n,0)}(x_s, x_t) &= \mu_{st}^{(n)}(x_s, x_t) \exp(\omega^n \theta_{st}(x_s, x_t) / \alpha_{st}), & \text{and} \\ \mu_s^{(n,0)}(x_s) &= \mu_s^{(n)}(x_s) \exp(\omega^n \theta_s(x_s) / \alpha_s).\end{aligned}\quad (21)$$

repeat
for each edge $(s, t) \in E$ **do**

$$\begin{aligned}\mu_{st}^{(n,\tau+1)}(x_s, x_t) &= \mu_{st}^{(n,\tau)}(x_s, x_t) \left(\frac{\mu_s^{(n,\tau)}(x_s)}{\sum_{x_t} \mu_{st}^{(n,\tau)}(x_s, x_t)} \right)^{\frac{\alpha_s}{\alpha_s + \alpha_{st}}}, & \text{and} \\ \mu_s^{(n,\tau+1)}(x_s) &= \mu_s^{(n,\tau)}(x_s)^{\frac{\alpha_s}{\alpha_s + \alpha_{st}}} \left(\sum_{x_t} \mu_{st}^{(n,\tau)}(x_s, x_t) \right)^{\frac{\alpha_{st}}{\alpha_s + \alpha_{st}}}.\end{aligned}\quad (22)$$

end for
for each node $s \in V$ **do**

$$\mu_s^{(n,\tau+1)}(x_s) = \frac{\mu_s^{(n,\tau)}(x_s)}{\sum_{x_s} \mu_s^{(n,\tau)}(x_s)}.\quad (23)$$

end for
until convergence

equivalent over the constraint set $\mathbb{L}(G)$). The proximal optimization problem as stated above (24) reduces to,

$$\mu^{n+1} = \arg \min_{\mathbf{v} \in \mathbb{L}(G)} \{ \langle \theta^n, \mathbf{v} \rangle + f_{\text{trw}}(\mathbf{v}) \}.$$

But this is precisely the optimization problem solved by the tree-reweighted sum-product (Wainwright and Jordan, 2003), as well as other related methods (Globerson and Jaakkola, 2007b; Hazan and Shashua, 2008), for a graphical model with parameters θ^n .

Computing the gradient of the function f_{trw} , and performing some algebra yields the algorithmic template of Algorithm 4.

3.6 Convergence

We now turn to the convergence of the message-passing algorithms that we have proposed. At a high-level, for any Bregman proximal function, convergence follows from two sets of known results: (a) convergence of proximal algorithms; and (b) convergence of cyclic Bregman projections.

For completeness, we re-state the consequences of these results here. For any positive sequence $\omega^n > 0$, we say that it satisfies the *infinite travel condition* if $\sum_{n=1}^{\infty} (1/\omega^n) = +\infty$. We let $\mu^* \in \mathbb{L}(G)$ denote an optimal solution (not necessarily unique) of the LP, and use $f^* = f(\mu^*) = \langle \theta, \mu^* \rangle$ to denote

Algorithm 4 TRW proximal solver

- For outer iterations $n = 0, 1, 2, \dots$,

(a) Update the parameters:

$$\begin{aligned}\theta_s^n(x_s) &= \omega^n \theta_s(x_s) + \log(\mu^n(x_s)) + 1, \\ \theta_{st}^n(x_s, x_t) &= \omega^n \theta_{st}(x_s, x_t) + \rho_{st} \left(\log \frac{\mu_{st}^n(x_s, x_t)}{\sum_{x'_s} \mu_{st}^n(x'_s, x_t) \sum_{x'_t} \mu_{st}^n(x_s, x'_t)} - 1 \right).\end{aligned}$$

(b) Run a convergent TRW-solver on a graphical model with parameters θ^n , so as to compute

$$\mu^{n+1} = \arg \min_{\mathbf{v} \in \mathbb{L}(G)} \left\{ -\langle \theta^n, \mathbf{v} \rangle + f_{\text{trw}}(\mathbf{v}) \right\}.$$

the LP optimal value. We say that the convergence rate is *superlinear* if

$$\lim_{n \rightarrow +\infty} \frac{|f(\mu^{n+1}) - f^*|}{|f(\mu^n) - f^*|} = 0,$$

and *linear* if

$$\lim_{n \rightarrow +\infty} \frac{|f(\mu^{n+1}) - f^*|}{|f(\mu^n) - f^*|} \leq \gamma,$$

for some $\gamma \in (0, 1)$. We say the convergence is *geometric* if there exists some constant $C > 0$ and $\gamma \in (0, 1)$ such that for all n ,

$$|f(\mu^n) - f^*| \leq C\gamma^n.$$

Proposition 1 (Rate of outer loop convergence) *Consider the sequence of iterates produced by a proximal algorithm (5) for LP-solving.*

- (a) *Using the quadratic proximal function and positive weight sequence $\omega^n \rightarrow +\infty$ satisfying infinite travel, the proximal sequence $\{\mu^n\}$ converges superlinearly.*
- (b) *Using the entropic proximal function and positive weight sequence ω^n satisfying infinite travel, the proximal sequence $\{\mu^n\}$ converges:*
 - (i) *superlinearly if $\omega^n \rightarrow 0$, and*
 - (ii) *at least linearly if $1/\omega^n \geq c$ for some constant $c > 0$.*

The quadratic case is covered in Bertsekas and Tsitsiklis (1997), whereas the entropic case was analyzed by Tseng and Bertsekas (1993), and Iusem and Teboulle (1995).

Our inner loop message updates use cyclic Bregman projections, for which there is also a substantial literature on convergence. Censor and Zenios (1997) show that with dual feasibility correction, cyclic projections onto general convex sets are convergent. For Euclidean projections with

linear constraints, Deutsch and Hundal (2006) establish a linear rate of convergence, with the rate dependent on angles between the half-spaces defining the constraints. The intuition is that the more orthogonal the half-spaces, the faster the convergence; for instance, a single iteration suffices for completely orthogonal constraints. Our inner updates thus converge linearly to the solution within each outer proximal step.

We note that the rate-of-convergence results for the outer proximal loops assume that the proximal update (computed within each inner loop) has been performed exactly. In practice, the inner loop iterations do not converge finitely (though they have a linear rate of convergence), so that an early stopping entails that the solution to each proximal update would be computed only approximately, up to some accuracy ε . That is, if the proximal optimization function at outer iteration n is $h^n(\mu)$ with minimum μ^{n+1} , then the computed proximal update $\underline{\mu}^{n+1}$ is sub-optimal, with $h^n(\underline{\mu}^{n+1}) - h^n(\mu^{n+1}) \leq \varepsilon$. Some recent theory has addressed whether superlinear convergence can still be obtained in such a setting; for instance, Solodov and Svaiter (2001) shows that that under mild conditions superlinear rates still hold for proximal iterates with inner-loop solutions that are ε -suboptimal. In practice, we cannot directly use ε -suboptimality as the stopping criterion for the inner loop iterations since we do not have the optimal solution μ^{n+1} . However, since we are trying to solve a feasibility problem, it is quite natural to check for violation in the constraints defining $\mathbb{L}(G)$. We terminate our inner iterations when the violation in all the constraints below a tolerance ε . As we show in Section 5, our experiments show that setting this termination threshold to $\underline{\varepsilon} = 10^{-4}$ is small enough for sub-optimality to be practically irrelevant and that superlinear convergence still occurs.

3.6.1 REMARKS

The quadratic proximal updates turn out to be equivalent to solving the primal form of the LP by the projected subgradient method (Bertsekas, 1995) for constrained optimization. (This use of the subgradient method should be contrasted with other work Feldman et al. (2002b); Komodakis et al. (2007) which performed subgradient descent to the dual of the LP.) For any constrained optimization problem:

$$\begin{aligned} \min_{\mu} \quad & f_0(\mu) \\ \text{s.t.} \quad & f_j(\mu) \leq 0, \quad j = 1, \dots, m, \end{aligned} \tag{25}$$

the projected subgradient method performs subgradient descent iteratively on (i) the objective function f_0 , as well as on (ii) the constraint functions $\{f_j\}_{j=1}^m$ till the constraints are satisfied. Casting it in the notation of Algorithm 1; over outer loop iterations $n = 1, \dots$, it sets

$$\mu^{n,0} = \mu^n - \alpha_n \nabla f_0(\mu^n),$$

and computes, over inner loop iterations $t = 1, \dots$,

$$\begin{aligned} j(t) &= t \bmod m, \\ \mu^{n,t+1} &= \mu^{n,t} - \alpha_{n,t} \nabla f_{j(t)}(\mu^{n,t}), \end{aligned}$$

and sets $\mu^{n+1} = \mu^{n,\infty}$, the converged estimate of the inner loops of outer iteration n . The constants $\{\alpha_n, \alpha_{n,t}\}$ are step-sizes for the corresponding subgradient descent steps.

The constraint set in our LP problem, $\mathbb{L}(G)$, has equality constraints so that it is not directly in the form of Equation (25). However any equality constraint $h(\mu) = 0$ can be rewritten equivalently as two inequality constraints $h(\mu) \leq 0$, and $-h(\mu) \leq 0$; so that one could cast our constrained LP in the form of (25) and solve it using the constrained subgradient descent method. As regards the step-sizes, suppose we set $\alpha_n = \omega^n$, and $\alpha_{n,t}$ according to Polyak’s step-size (Bertsekas, 1995) so that $\alpha_{n,t} = \frac{f_j(\mu^{n,t}) - f_j(\mu^*)}{\|\nabla f_j(\mu^{n,t})\|_2^2}$, where μ^* is the constrained optimum. Since μ^* is feasible by definition, $f_j(\mu^*) = 0$. Further, for the normalization constraints $C_{ss}(\mu) \leq 1$ where $C_{ss}(\mu) := \sum_{x_s \in \mathcal{X}} \mu_s(x_s) - 1$, we have $\|\nabla C_{ss}(\mu)\|^2 = m$, while for the marginalization constraints $C_{st}(\mu) \leq 0$, where $C_{st}(\mu) := \sum_{x_t \in \mathcal{X}} \mu_{st}(x_s, x_t) - \mu_s(x_s)$, we have $\|\nabla C_{st}(\mu)\|^2 = (m + 1)$. It can then be seen that the subgradient method for constrained optimization applied to our constrained LP with the above step-sizes yields the same updates as our quadratic proximal scheme.

4. Rounding Schemes with Optimality Guarantees

The graph-structured LP in (4) was a relaxation of the MAP integer program (1), so that there are two possible outcomes to solving the LP: either an integral vertex is obtained, which is then guaranteed to be a MAP configuration, or a fractional vertex is obtained, in which case the relaxation is loose. In the latter case, a natural strategy is to “round” the fractional solution, so as to obtain an integral solution (Raghavan and Thompson, 1987). Such rounding schemes may either be randomized or deterministic. A natural measure of the quality of the rounded solution is in terms of its value relative to the optimal (MAP) value. There is now a substantial literature on performance guarantees of various rounding schemes, when applied to particular sub-classes of MAP problems (e.g., Raghavan and Thompson, 1987; Kleinberg and Tardos, 1999; Chekuri et al., 2005).

In this section, we show that rounding schemes can be useful even when the LP optimum is integral, since they may permit an LP-solving algorithm to be *finitely terminated*—that is, before it has actually solved the LP—while retaining the same optimality guarantees about the final output. An attractive feature of our proximal Bregman procedures is the existence of precisely such rounding schemes—namely, that under certain conditions, rounding pseudomarginals at intermediate iterations yields the integral LP optimum. We describe these rounding schemes in the following sections, and provide two kinds of results. We provide certificates under which the rounded solution is guaranteed to be MAP optimal; moreover, we provide upper bounds on the number of outer iterations required for the rounding scheme to obtain the LP optimum.

In the next Section 4.1, we describe and analyze deterministic rounding schemes that are specifically tailored to the proximal Bregman procedures that we have described. Then in the following Section 4.2, we propose and analyze a graph-structured randomized rounding scheme, which applies not only to our proximal Bregman procedures, but more broadly to any algorithm that generates a sequence of iterates contained within the local polytope $\mathbb{L}(G)$.

4.1 Deterministic Rounding Schemes

We begin by describing three deterministic rounding schemes that exploit the particular structure of the Bregman proximal updates.

4.1.1 NODE-BASED ROUNDING

This method is the simplest of the deterministic rounding procedures, and applies to the quadratic and entropic updates. It operates as follows: given the vector μ^n of pseudomarginals at iteration n , obtain an integral configuration $x^n(\mu^n) \in \mathcal{X}^N$ by choosing

$$x_s^n \in \arg \max_{x'_s \in \mathcal{X}} \mu^n(x'_s), \quad \text{for each } s \in V.$$

We say that the node-rounded solution x^n is *edgewise-consistent* if

$$(x_s^n, x_t^n) \in \arg \max_{(x'_s, x'_t) \in \mathcal{X} \times \mathcal{X}} \mu_{st}^n(x'_s, x'_t) \quad \text{for all edges } (s, t) \in E. \quad (26)$$

4.1.2 NEIGHBORHOOD-BASED ROUNDING

This rounding scheme applies to all three proximal schemes. For each node $s \in V$, denote its star-shaped neighborhood graph by $N_s = \{(s, t) | t \in N(s)\}$, consisting of edges between node s and its neighbors. Let {QUA, ENT, TRW} refer to the quadratic, entropic, and tree-reweighted schemes respectively.

(a) Define the neighborhood-based energy function

$$F_s(x; \mu^n) := \begin{cases} 2\mu^n(x_s) + \sum_{t \in N(s)} \mu^n(x_s, x_t) & \text{for QUA} \\ 2\alpha_s \log \mu_s^n(x_s) + \sum_{t \in N(s)} \alpha_{st} \log \mu_{st}^n(x_s, x_t) & \text{for ENT} \\ 2 \log \mu^n(x_s) + \sum_{t \in N(s)} \rho_{st} \log \frac{\mu_{st}^n(x_s, x_t)}{\mu_s^n(x_s) \mu_t^n(x_t)} & \text{for TRW.} \end{cases} \quad (27)$$

(b) Compute a configuration $x^n(N_s)$ maximizing the function $F_s(x; \mu^n)$ by running two rounds of ordinary max-product on the star graph.

Say that such a rounding is *neighborhood-consistent* if the neighborhood MAP solutions $\{x^n(N_s), s \in V\}$ agree on their overlaps.

4.1.3 TREE-BASED ROUNDING

This method applies to all three proximal schemes, but most naturally to the TRW proximal method. Let T_1, \dots, T_K be a set of spanning trees that cover the graph (meaning that each edge appears in at least one tree), and let $\{\rho(T_i), i = 1, \dots, K\}$ be a probability distribution over the trees. For each edge (s, t) , define the *edge appearance probability* $\rho_{st} = \sum_{i=1}^K \rho(T_i) \mathbb{I}[(s, t) \in T_i]$. Then for each tree $i = 1, \dots, K$:

(a) Define the tree-structured energy function

$$F_i(x; \mu^n) := \begin{cases} \sum_{s \in V} \log \mu^n(x_s) + \sum_{(s, t) \in E(T_i)} \frac{1}{\rho_{st}} \log \mu_{st}^n(x_s, x_t) & \text{for QUA} \\ \sum_{s \in V} \alpha_s \log \mu^n(x_s) + \sum_{(s, t) \in E(T_i)} \frac{\alpha_{st}}{\rho_{st}} \log \mu_{st}^n(x_s, x_t) & \text{for ENT} \\ \sum_{s \in V} \log \mu^n(x_s) + \sum_{(s, t) \in E(T_i)} \log \frac{\mu_{st}^n(x_s, x_t)}{\mu_s^n(x_s) \mu_t^n(x_t)} & \text{for TRW.} \end{cases} \quad (28)$$

- (b) Run the ordinary max-product problem on energy $F_i(x; \mu^n)$ to find a MAP-optimal configuration $x^n(T_i)$.

Say that such a rounding is *tree-consistent* if the tree MAP solutions $\{x^n(T_i), i = 1, \dots, M\}$ are all equal. This notion of tree-consistency is similar to the underlying motivation of the tree-reweighted max-product algorithm (Wainwright et al., 2005).

4.1.4 OPTIMALITY CERTIFICATES FOR DETERMINISTIC ROUNDING

The following result characterizes the optimality guarantees associated with these rounding schemes, when they are consistent respectively in the *edge-consistency*, *neighborhood-consistency* and *tree-consistency* senses defined earlier.

Theorem 2 (Deterministic rounding with MAP certificate) *Consider a sequence of iterates $\{\mu^n\}$ generated by the quadratic or entropic proximal schemes. For any $n = 1, 2, 3, \dots$, any consistent rounded solution x^n obtained from μ^n via any of the node, neighborhood or tree-rounding schemes (when applicable) is guaranteed to be a MAP-optimal solution. For the iterates of TRW proximal scheme, the guarantee holds for both neighborhood and tree-rounding methods.*

We prove this claim in Section 4.1.6. It is important to note that such deterministic rounding guarantees do *not* apply to an arbitrary algorithm for solving the linear program. At a high-level, there are two key properties required to ensure guarantees in the rounding. First, the algorithm must maintain some representation of the cost function that (up to possible constant offsets) is equal to the cost function of the original problem, so that the set of maximizers of the invariance would be equivalent to the set of maximizers of the original cost function, and hence the MAP problem. Second, given a rounding scheme that maximizes tractable sub-parts of the reparameterized cost function, the rounding is said to be admissible if these partial solutions agree with one another. Our deterministic rounding schemes and optimality guarantees follow this approach, as we detail in the proof of Theorem 2.

We note that the invariances maintained by the proximal updates in this paper are closely related to the reparameterization condition satisfied by the sum-product and max-product algorithms (Wainwright et al., 2003). Indeed, each sum-product (or max-product) update can be shown to compute a new set of parameters for the Markov random field that preserves the probability distribution. A similar but slightly different notion of reparameterization underlies the tree-reweighted sum-product and max-product algorithms (Wainwright et al., 2005); for these algorithms, the invariance is preserved in terms of convex combinations over tree-structured graphs. The tree-reweighted max-product algorithm attempts to produce MAP optimality certificates that are based on verifying consistency of MAP solutions on certain tree-structured components whose convex combination is equal to the LP cost. The sequential TRW-S max-product algorithm of Kolmogorov (2006) is a version of tree-reweighted max-product using a clever scheduling of the messages to guarantee monotonic changes in a dual LP cost function. Finally, the elegant work of Weiss et al. (2007) exploits similar reparameterization arguments to derive conditions under which their convex free-energy based sum-product algorithms yield the optimal MAP solution.

An attractive feature of all the rounding schemes that we consider is their relatively low computational cost. The node-rounding scheme is trivial to implement. The neighborhood-based scheme requires running two iterations of max-product for each neighborhood of the graph. Finally, the tree-rounding scheme requires $O(KN)$ iterations of max-product, where K is the number of trees

that cover the graph, and N is the number of nodes. Many graphs with cycles can be covered with a small number K of trees; for instance, the lattice graph in 2-dimensions can be covered with two spanning trees, in which case the rounding cost is linear in the number of nodes.

4.1.5 BOUNDS ON ITERATIONS FOR DETERMINISTIC ROUNDING

Of course, the natural question is how many iterations are sufficient for a given rounding scheme to succeed. The following result provides a way of deriving such upper bounds:

Corollary 3 *Suppose that the LP optimum is uniquely attained at an integral vertex μ^* , and consider algorithms generating sequence $\{\mu^n\}$ converging to μ^* . Then we have the following guarantees:*

- (a) *for quadratic and entropic schemes, all three types of rounding recover the MAP solution once $\|\mu^n - \mu\|_\infty \leq 1/2$.*
- (b) *for the TRW-based proximal method, tree-based rounding recovers the MAP solution once $\|\mu^n - \mu\|_\infty \leq \frac{1}{4N}$.*

Proof We first claim that if the ℓ_∞ -bound $\|\mu^n - \mu^*\|_\infty < \frac{1}{2}$ is satisfied, then the node-based rounding returns the (unique) MAP configuration, and moreover this MAP configuration x^* is edge-consistent with respect to μ^n . To see these facts, note that the ℓ_∞ bound implies, in particular, that at every node $s \in V$, we have

$$|\mu_s^n(x_s^*) - \mu_s^*(x_s^*)| = |\mu_s^n(x_s^*) - 1| < \frac{1}{2},$$

which implies that $\mu_s^n(x_s^*) > 1/2$ as $\mu_s^*(x_s^*) = 1$. Due to the non-negativity constraints and marginalization constraint $\sum_{x_s \in \mathcal{X}} \mu^n(x_s) = 1$, at most one configuration can have mass above $1/2$. Thus, node-based rounding returns x_s^* at each node s , and hence overall, it returns the MAP configuration x^* . The same argument also shows that the inequality $\mu_{st}^n(x_s^*, x_t^*) > \frac{1}{2}$ holds, which implies that $(x_s^*, x_t^*) = \arg \max_{x_s, x_t} \mu^n(x_s, x_t)$ for all $(s, t) \in E$. Thus, we have shown x^* is edge-consistent for μ_{st}^n , according to the definition (26).

Next we turn to the performance of neighborhood and tree-rounding for the quadratic and entropic updates. For $n \geq n^*$, we know that x^* achieves the unique maximum of $\mu_s^n(x_s)$ at each node, and $\mu_{st}^n(x_s, x_t)$ on each edge. From the form of the neighborhood and tree energies (27),(28), this node- and edge-wise optimality implies that $x^*(N(s)) := \{x_t^*, t \in s \cup N(s)\}$ maximizes the neighborhood-based and tree-based cost functions as well, which implies success of neighborhood and tree-rounding. (Note that the positivity of the weights α_s and α_{st} is required to make this assertion.)

For the TRW algorithm in part (b), we note that when $\|\mu^n - \mu\|_\infty \leq 1/(4N)$, then we must have $\mu_s^n(x_s^*) \geq 1 - 1/(4N)$ for every node. We conclude that these inequalities imply that $x^* = (x_1^*, \dots, x_N^*)$ must be the unique MAP on every tree. Indeed, consider the set $S = \{x \in \mathcal{X}^N \mid x \neq x^*\}$. By union bound, we have

$$\begin{aligned} \mathbb{P}(S) &= \mathbb{P}[\exists s \in V \mid x_s \neq x_s^*] \\ &\leq \sum_{s=1}^N \mathbb{P}(x_s \neq x_s^*) \\ &= \sum_{s=1}^N (1 - \mu_s(x_s^*)) \leq \frac{1}{4}, \end{aligned}$$

showing that we have $\mathbb{P}(x^*) \geq 3/4$, so that x^* must be the MAP configuration.

To conclude the proof, note that the tree-rounding scheme computes the MAP configuration on each tree T_i , under a distribution with marginals μ_s and μ_{st} . Consequently, under the stated conditions, the configuration x^* must be the unique MAP configuration on each tree, so that tree rounding is guaranteed to find it. \blacksquare

Using this result, it is possible to bound the number of iterations required to achieve the ℓ_∞ -bounds. In particular, suppose that the algorithm has a linear rate of convergence—say that $|f(\mu^n) - f(\mu^*)| \leq |f(\mu^0) - f(\mu^*)|\gamma^n$ for some $\gamma \in (0, 1)$. For the quadratic or entropic methods, it suffices to show that $\|\mu^n - \mu^*\|_2 < 1/2$. For the entropic method, there exists some constant $C > 0$ such that $\|\mu^n - \mu^*\|_2 \leq \frac{1}{2C}|f(\mu^n) - f(\mu^*)|$ (cf. Prop. 8, Iusem and Teboulle, 1995). Consequently, we have

$$\|\mu^n - \mu^*\|_2 \leq \frac{|f(\mu^0) - f(\mu^*)|}{2C} \gamma^n.$$

Consequently, after $n^* := \frac{\log C |f(\mu^0) - f(\mu^*)|}{\log(1/\gamma)}$ iterations, the rounding scheme would be guaranteed to configuration for the entropic proximal method. Similar finite iteration bounds can also be obtained for the other proximal methods, showing finite convergence through use of our rounding schemes.

Note that we proved correctness of the neighborhood and tree-based rounding schemes by leveraging the correctness of the node-based rounding scheme. In practice, it is possible for neighborhood- or tree-based rounding to succeed even if node-based rounding fails; however, we currently do not have any sharper sufficient conditions for these rounding schemes.

4.1.6 PROOF OF THEOREM 2

We now turn to the proof of Theorem 2. At a high level, the proof consists of two main steps. First, we show that each proximal algorithm maintains a certain invariant of the original MAP cost function $F(x; \theta)$; in particular, the iterate μ^n induces a reparameterization $F(x; \mu^n)$ of the cost function such that the set of maximizers is preserved—viz.:

$$\arg \max_{x \in \mathcal{X}^N} F(x; \theta) := \arg \max_{x \in \mathcal{X}^N} \sum_{s \in V, x_s \in \mathcal{X}} \theta_s(x_s) + \sum_{(s,t) \in E, x_s, x_t \in \mathcal{X}} \theta_{st}(x_s, x_t) = \arg \max_{x \in \mathcal{X}^N} F(x; \mu^n). \quad (29)$$

Second, we show that the consistency conditions (edge, neighborhood or tree, respectively) guarantee that the rounded solution belongs to $\arg \max_{x \in \mathcal{X}^N} F(x; \mu^n)$

We begin with a lemma on the invariance property:

Lemma 4 (Invariance of maximizers) *Define the function*

$$F(x; \mu) := \begin{cases} \sum_{s \in V} \mu_s(x_s) + \sum_{(s,t) \in E} \mu_{st}(x_s, x_t) & \text{for QUA} \\ \sum_{s \in V} \alpha_s \log \mu_s(x_s) + \sum_{(s,t) \in E} \alpha_{st} \log \mu_{st}(x_s, x_t) & \text{for ENT} \\ \sum_{s \in V} \log \mu_s(x_s) + \sum_{(s,t) \in E} \rho_{st} \log \frac{\mu_{st}(x_s, x_t)}{\mu_s(x_s) \mu_t(x_t)} & \text{for TRW.} \end{cases} \quad (30)$$

At each iteration $n = 1, 2, 3, \dots$ for which $\mu^n > 0$, the function $F(x; \mu^n)$ preserves the set of maximizers (29).

The proof of this claim, provided in Appendix C, is based on exploiting the necessary (Lagrangian) conditions defined by the optimization problems characterizing the sequence of iterations $\{\mu^n\}$.

For the second part of the proof, we show how a solution x^* , obtained by a rounding procedure, is guaranteed to maximize the function $F(x; \mu^n)$, and hence (by Lemma 4) the original cost function $F(x; \theta)$. In particular, we state the following simple lemma:

Lemma 5 *The rounding procedures have the following guarantees:*

- (a) *Any edge-consistent configuration from node rounding maximizes $F(x; \mu^n)$ for the quadratic and entropic schemes.*
- (b) *Any neighborhood-consistent configuration from neighborhood rounding maximizes $F(x; \mu^n)$ for the quadratic and entropic schemes.*
- (c) *Any tree-consistent configuration from tree rounding maximizes $F(x; \mu^n)$ for all three schemes.*

Proof We begin by proving statement (a). Consider an edge-consistent integral configuration x^* obtained from node rounding. By definition, it maximizes $\mu^n(x_s)$ for all $s \in V$, and $\mu_{st}^n(x_s, x_t)$ for all $(s, t) \in E$, and so by inspection, also maximizes $F(x; \mu^n)$ for the quadratic and proximal cases.

We next prove statement (b) on neighborhood rounding. Suppose that neighborhood rounding outputs a single neighborhood-consistent integral configuration x^* . Since $x_{N(s)}^*$ maximizes the neighborhood energy (27) at each node $s \in V$, it must also maximize the sum $\sum_{s \in V} F_s(x; \mu^n)$. A little calculation shows that this sum is equal to $2F(x; \mu^n)$, the factor of two arising since the term on edge (s, t) arises twice, one for neighborhood rooted at s , and once for t .

Turning to claim (c), let x^* be a tree-consistent configuration obtained from tree rounding. Then for each $i = 1, \dots, K$, the configuration x^* maximizes the tree-structured function $F_i(x; \mu^n)$, and hence also maximizes the convex combination $\sum_{i=1}^K \rho(T_i) F_i(x; \mu^n)$. By definition of the edge appearance probabilities ρ_{st} , this convex combination is equal to the function $F(x; \mu^n)$. ■

4.2 Randomized Rounding Schemes

The schemes considered in the previous section were all deterministic, since (disregarding any possible ties), the output of the rounding procedure was a deterministic function of the given pseudomarginals $\{\mu_s^n, \mu_{st}^n\}$. In this section, we consider randomized rounding procedures, in which the output is a random variable.

Perhaps the most naive randomized rounding scheme is the following: for each node $r \in V$, assign it value $x_r \in \mathcal{X}$ with probability $\mu_r^n(x_r)$. We propose a graph-structured generalization of this naive randomized rounding scheme, in which we perform the rounding in a dependent way across sub-groups of nodes, and establish guarantees for its success. In particular, we show that when the LP relaxation has a unique integral optimum that is well-separated from the second best configuration, then the rounding scheme succeeds with high probability after a pre-specified number of iterations.

4.2.1 THE RANDOMIZED ROUNDING SCHEME

Our randomized rounding scheme is based on any given subset E' of the edge set E . Consider the subgraph $G(E \setminus E')$, with vertex set V , and edge set $E \setminus E'$. We assume that E' is chosen such that

the subgraph $G(E \setminus E')$ is a forest. That is, we can decompose $G(E \setminus E')$ into a union of disjoint trees, $\{T_1, \dots, T_K\}$, where $T_i = (V_i, E_i)$, such that the vertex subsets V_i are all disjoint and $V = V_1 \cup V_2 \cup \dots \cup V_K$. We refer to the edge subset as *forest-inducing* when it has this property. Note that such a subset always exists, since $E' = E$ is trivially forest-inducing. In this case, the “trees” simply correspond to individual nodes, without any edges; $V_i = \{i\}$, $E_i = \emptyset$, $i = 1, \dots, N$.

For any forest-inducing subset $E' \subseteq E$, Algorithm 5 defines our randomized rounding scheme.

Algorithm 5 RANDOMIZED ROUNDING SCHEME

for subtree indices $i = 1, \dots, K$ **do**

Sample a sub-configuration X_{V_i} from the probability distribution

$$p(x_{V_i}; \mu(T_i)) = \prod_{s \in V_i} \mu^n(x_s) \prod_{(s,t) \in E_i} \frac{\mu^n(x_s, x_t)}{\mu^n(x_s) \mu^n(x_t)}.$$

end for

Form the global configuration $X \in \mathcal{X}^N$ by concatenating all the local random samples:

$$X := \left(X_{V_1}, \dots, X_{V_K} \right).$$

To be clear, the randomized solution X is a function of both the pseudomarginals μ^n , and the choice of forest-inducing subset E' , so that we occasionally use the notation $X(\mu^n; E')$ to reflect explicitly this dependence. Note that the simplest rounding scheme of this type is obtained by setting $E' = E$. Then the “trees” simply correspond to individual nodes without any edges, and the rounding scheme is the trivial node-based scheme.

The randomized rounding scheme can be “derandomized” so that we obtain a deterministic solution $x^d(\mu^n; E')$ that does at least well as the randomized scheme does in expectation. This derandomization scheme is shown in Algorithm 6, and its correctness is guaranteed in the following theorem, proved in Appendix D.

Theorem 6 *Let $(G = (V, E), \theta)$ be the given MAP problem instance, and let $\mu^n \in \mathbb{L}(G)$ be any set of pseudomarginals in the local polytope $\mathbb{L}(G)$. Then, for any subset $E' \subseteq E$ of the graph G , the (E', μ^n) -randomized rounding scheme in Algorithm 5, when derandomized as in Algorithm 6 satisfies,*

$$F(x^d(\mu^n; E'); \theta) \geq \mathbb{E} \left(F(X(\mu^n; E'); \theta) \right),$$

where $X(\mu^n; E')$ and $x^d(\mu^n; E')$ denote the outputs of the randomized and derandomized schemes respectively.

4.2.2 OSCILLATION AND GAPS

In order to state some theoretical guarantees on our randomized rounding schemes, we require some notation. For any edge $(s, t) \in E$, we define the *edge-based oscillation*

$$\delta_{st}(\theta) := \max_{x_s, x_t} [\theta_{st}(x_s, x_t)] - \min_{x_s, x_t} [\theta_{st}(x_s, x_t)].$$

Algorithm 6 DERANDOMIZED ROUNDING SCHEME

 Initialize: $\bar{\mu} = \mu^l$.

for subtree indices $i = 1, \dots, K$ **do**

Solve

$$x_{V_i}^d = \arg \max_{x_{V_i}} \sum_{s \in V_i} \left\{ \theta_s(x_s) + \sum_{t: (s,t) \in E'} \sum_{x_t} \bar{\mu}_t(x_t) \theta_{st}(x_s, x_t) \right\} + \sum_{(s,t) \in E_i} \theta_{st}(x_s, x_t).$$

 Update $\bar{\mu}$:

$$\bar{\mu}_s(x_s) = \begin{cases} \bar{\mu}_s(x_s) & \text{if } s \notin V_i \\ 0 & \text{if } s \in V_i, x_s^d \neq x_s \\ 1 & \text{if } s \in V_i, x_s^d = x_s. \end{cases}$$

$$\bar{\mu}_{st}(x_s, x_t) = \begin{cases} \bar{\mu}_{st}(x_s, x_t) & \text{if } (s,t) \notin E_i \\ \bar{\mu}_s(x_s) \bar{\mu}_t(x_t) & \text{if } (s,t) \in E_i. \end{cases}$$

end for

 Form the global configuration $x^d \in \mathcal{X}^N$ by concatenating all the subtree configurations:

$$x^d := \left(x_{V_1}^d, \dots, x_{V_K}^d \right).$$

We define the *node-based oscillation* $\delta_s(\theta)$ in the analogous manner. The quantities $\delta_s(\theta)$ and $\delta_{st}(\theta)$ are measures of the strength of the potential functions.

We extend these measures of interaction strength to the full graph in the natural way

$$\delta_G(\theta) := \max \left\{ \max_{(s,t) \in E} \delta_{st}(\theta), \max_{s \in V} \delta_s(\theta) \right\}.$$

Using this oscillation function, we now define a measure of the quality of a unique MAP optimum, based on its separation from the second most probable configuration. In particular, letting $x^* \in \mathcal{X}^N$ denote a MAP configuration, and recalling the notation $F(x; \theta)$ for the LP objective, we define the *graph-based gap*

$$\Delta(\theta; G) := \frac{\min_{x \neq x^*} \left[F(x^*; \theta) - F(x; \theta) \right]}{\delta_G(\theta)}.$$

This gap function is a measure of how well-separated the MAP optimum x^* is from the remaining integral configurations. By definition, the gap $\Delta(\theta; G)$ is always non-negative, and it is strictly positive whenever the MAP configuration x^* is unique. Finally, note that the gap is invariant to the translations ($\theta \mapsto \theta' = \theta + C$) and rescalings ($\theta \mapsto \theta' = c\theta$) of the parameter vector θ . These invariances are appropriate for the MAP problem since the optima of the energy function $F(x; \theta)$ are not affected by either transformation (i.e., $\arg \max_x F(x; \theta) = \arg \max_x F(x; \theta')$ for both $\theta' = \theta + C$ and $\theta' = c\theta$).

Finally, for any forest-inducing subset, we let $d(E')$ be the maximum degree of any node with respect to edges in E' —namely,

$$d(E') := \max_{s \in V} |\{t \in V \mid (s, t) \in E'\}|.$$

4.2.3 OPTIMALITY GUARANTEES FOR RANDOMIZED ROUNDING

We show, in this section, that when the pseudomarginals μ^n are within a specified ℓ_1 norm ball around the unique MAP optimum μ^* , the randomized rounding scheme outputs the MAP configuration with high probability.

Theorem 7 *Consider a problem instance (G, θ) for which the MAP optimum x^* is unique, and let μ^* be the associated vertex of the polytope $\mathbb{L}(G)$. For any $\varepsilon \in (0, 1)$, if at some iteration n , we have $\mu^n \in \mathbb{L}(G)$, and*

$$\|\mu^n - \mu^*\|_1 \leq \frac{\varepsilon \Delta(\theta; G)}{1 + d(E')}, \quad (31)$$

then (E', μ^n) -randomized rounding succeeds with probability greater than $1 - \varepsilon$,

$$\mathbb{P}[X(\mu^n; E') = x^*] \geq 1 - \varepsilon.$$

We provide the proof of this claim in Appendix E. It is worthwhile observing that the theorem applies to any algorithm that generates a sequence $\{\mu^n\}$ of iterates contained within the local polytope $\mathbb{L}(G)$. In addition to the proximal Bregman updates discussed in this paper, it also applies to interior-point methods (Boyd and Vandenberghe, 2004) for solving LPs. For the naive rounding based on $E' = E$, the sequence $\{\mu^n\}$ need not belong to $\mathbb{L}(G)$, but instead need only satisfy the milder conditions $\mu_s^n(x_s) \geq 0$ for all $s \in V$ and $x_s \in \mathcal{X}$, and $\sum_{x_s} \mu_s^n(x_s) = 1$ for all $s \in V$.

The derandomized rounding scheme enjoys a similar guarantee, as shown in the following theorem, proved in Appendix F.

Theorem 8 *Consider a problem instance (G, θ) for which the MAP optimum x^* is unique, and let μ^* be the associated vertex of the polytope $\mathbb{L}(G)$. If at some iteration n , we have $\mu^n \in \mathbb{L}(G)$, and*

$$\|\mu^n - \mu^*\|_1 \leq \frac{\Delta(\theta; G)}{1 + d(E')},$$

then the (E', μ^n) -derandomized rounding scheme in Algorithm 6 outputs the MAP solution,

$$x^d(\mu^n; E') = x^*.$$

4.2.4 BOUNDS ON ITERATIONS FOR RANDOMIZED ROUNDING

Although Theorems 7 and 8 apply even for sequences $\{\mu^n\}$ that need not converge to μ^* , it is most interesting when the LP relaxation is tight, so that the sequence $\{\mu^n\}$ generated by any LP-solver satisfies the condition $\mu^n \rightarrow \mu^*$. In this case, we are guaranteed that for any fixed $\varepsilon \in (0, 1)$, the bound (31) will hold for an iteration number n that is “large enough”. Of course, making this intuition precise requires control of convergence rates. Recall that N is the number of nodes in the graph, and m is cardinality of the set \mathcal{X} from which all variables takes their values. With this notation, we have the following.

Corollary 9 *Under the conditions of Theorem 7, suppose that the sequence of iterates $\{\mu^n\}$ converge to the LP (and MAP) optimum at a linear rate: $\|\mu^n - \mu^*\|_2 \leq \gamma^n \|\mu^0 - \mu^*\|_2$. Then:*

- (a) *The randomized rounding in Algorithm 5 succeeds with probability at least $1 - \varepsilon$ for all iterations greater than*

$$n^* := \frac{\frac{1}{2} \log(Nm + N^2m^2) + \log(\|\mu^0 - \mu^*\|_2) + \log\left(\frac{1+d(E')}{\Delta(\theta; G)}\right) + \log(1/\varepsilon)}{\log(1/\gamma)}.$$

- (b) *The derandomized rounding in Algorithm 6 yields the MAP solution for all iterations greater than*

$$n^* := \frac{\frac{1}{2} \log(Nm + N^2m^2) + \log(\|\mu^0 - \mu^*\|_2) + \log\left(\frac{1+d(E')}{\Delta(\theta; G)}\right)}{\log(1/\gamma)}.$$

This corollary follows by observing that the vector $(\mu^n - \mu^*)$ has less than $Nm + N^2m^2$ elements, so that $\|\mu^n - \mu^*\|_1 \leq \sqrt{Nm + N^2m^2} \|\mu^n - \mu^*\|_2$. Moreover, Theorems 7 and 8 provide an ℓ_1 -ball radius such that the rounding schemes succeed (either with probability greater than $1 - \varepsilon$, or deterministically) for all pseudomarginal vectors within these balls.

5. Experiments

In this section, we provide the results of several experiments to illustrate the behavior of our methods on different problems. We performed experiments on 4-nearest neighbor grid graphs with sizes varying from $N = 100$ to $N = 900$, using models with either $m = 3$ or $m = 5$ labels. The edge potentials were set to Potts functions, of the form

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

for a parameter $\beta_{st} \in \mathbb{R}$. These potential functions penalize disagreement of labels if $\beta_{st} > 0$, and penalize agreement if $\beta_{st} < 0$. The Potts weights on edges β_{st} were chosen randomly as $\text{Uniform}(-1, +1)$. We set the node potentials as $\theta_s(x_s) \sim \text{Uniform}(-\text{SNR}, \text{SNR})$, for some signal-to-noise parameter $\text{SNR} \geq 0$ that controls the ratio of node to edge strengths. In applying all of the proximal procedures, we set the proximal weights as $\omega^n = n$.

5.1 Rates of Convergence

We begin by reporting some results on the convergence rates of proximal updates. Figure 2(a) plots the logarithmic distance $\log \|\mu^n - \mu^*\|_2$ versus the number of iterations for grids of different sizes (node numbers $N \in \{100, 400, 900\}$). Here μ^n is the iterate at step n entropic proximal method and μ^* is the LP optimum. In all cases, note how the curves have an inverted quadratic shape, corresponding to a superlinear rate of convergence, which is consistent with Proposition 1. On other hand, Figure 2(b) provides plots of the logarithmic distance versus iteration number for problem sizes $N = 900$, and over a range of signal-to-noise ratios SNR (in particular, $\text{SNR} \in \{0.05, 0.25, 0.50, 1.0, 2.0\}$). Notice how the plots still show the same inverted quadratic shape, but that the rate of convergence slows down as the SNR decreases, as is to be expected.

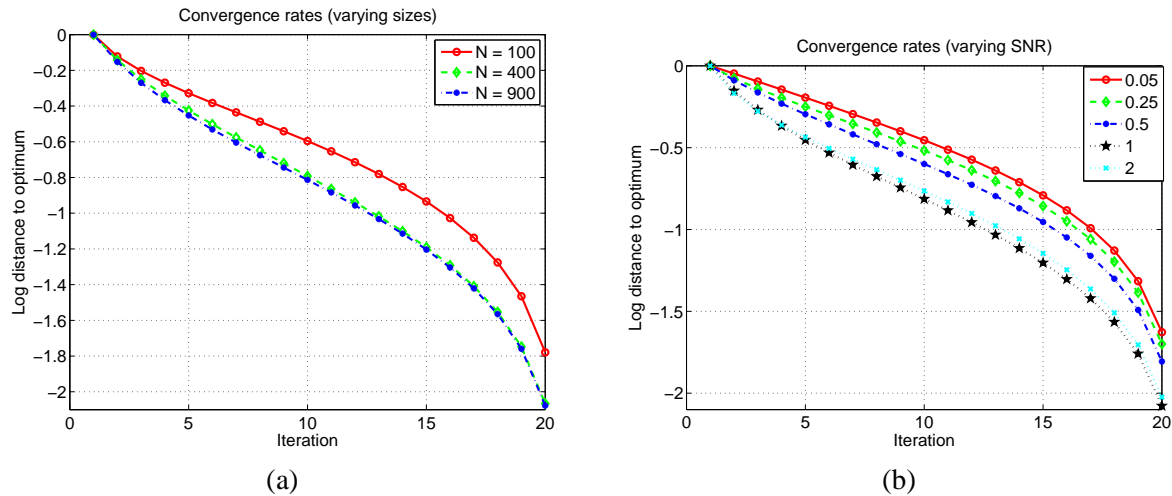


Figure 2: (a) Plot of distance $\log_{10} \|\mu^n - \mu^*\|_2$ between the current entropic proximal iterate μ^n and the LP optimum μ^* versus iteration number for Potts models on grids with $N \in \{100, 400, 900\}$ vertices, $m = 5$ labels and $\text{SNR} = 1$. Note the superlinear rate of convergence, consistent with Proposition 1. (b) Plot of distance $\log_{10} \|\mu^n - \mu^*\|_2$ between the current entropic proximal iterate μ^n and the LP optimum μ^* versus iteration number for Potts models on grids with $m = 5$ labels, $N = 900$ vertices, and a range of signal-to-noise ratios $\text{SNR} \in \{0.05, 0.25, 0.50, 1.0, 2.0\}$. The rate of convergence remains superlinear but slows down as the SNR is decreased.

In Figure 3, we compare two of our proximal schemes—the entropic and the quadratic schemes—with a subgradient descent method, as previously proposed (Feldman et al., 2002a; Komodakis et al., 2007). For the comparison, we used a Potts model on a grid of 400 nodes, with each node taking three labels. The Potts weights were set as earlier, with $\text{SNR} = 2$. Plotted in Figure 3(a) are the log probabilities of the solutions from the TRW-proximal and entropic proximal methods, compared to the dual upper bound that is provided by the sub-gradient method. Each step on the horizontal axis is a single outer iteration for the proximal methods, and five steps of the subgradient method. (We note that it is slower to perform five subgradient steps than a single proximal outer iteration.) Both the primal proximal methods and the dual subgradient method converge to the same point. The TRW-based proximal scheme converges the fastest, essentially within four outer iterations, whereas the entropic scheme requires a few more iterations. The convergence rate of the subgradient ascent method is slower than both of these proximal schemes, even though we allowed it to take more steps per “iteration”. In Figure 3(b), we plot a number of traces showing the number of inner iterations (vertical axis) required as a function of outer iteration (horizontal axis). The average number of inner iterations is around 20, and only rarely does the algorithm require substantially more.

5.2 Comparison of Rounding Schemes

In Figure 4, we compare five of our rounding schemes on a Potts model on grid graphs with $N = 400$, $m = 3$ labels and $\text{SNR} = 2$. For the graph-structured randomized rounding schemes, we used the node-based rounding scheme (so that $E \setminus E' = \emptyset$), and the chain-based rounding scheme (so that

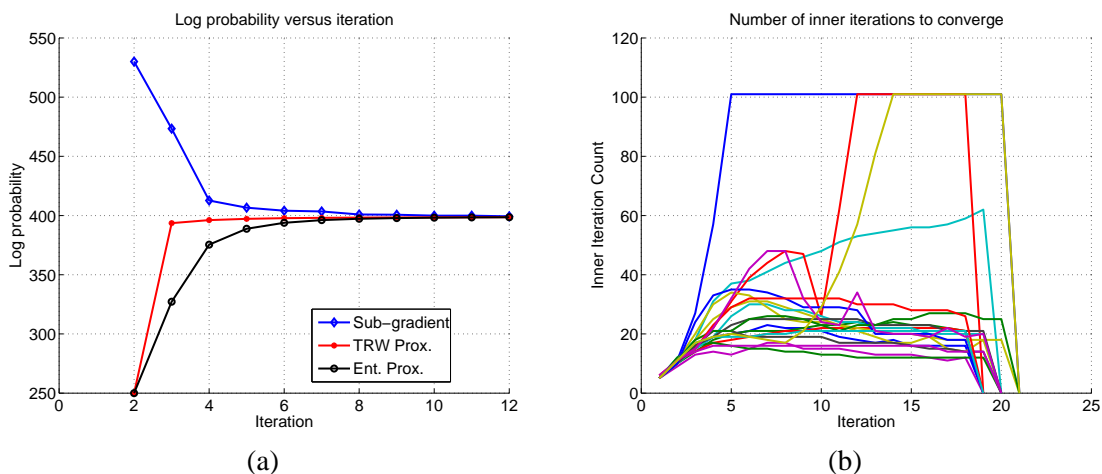


Figure 3: (a) Plots of the function value (for fractional iterates μ^i) versus number of iterations for a Potts model with $N = 400$ vertices, $m = 3$ labels and $\text{SNR} = 2$. Three methods are compared: a subgradient method (Feldman et al., 2002b; Komodakis et al., 2007), the entropic proximal method (Ent. Prox.), and the TRW-based proximal method (TRW Prox.). (b) Traces of different algorithm runs showing the number of inner iterations (vertical axis) versus the outer iteration number (horizontal axis). Typically around 20 inner iterations are required.

$E \setminus E'$ is the set of horizontal chains in the grid). For the deterministic rounding schemes, we used the node-based, neighborhood-based and the tree-based rounding schemes. Panel (a) of Figure 4 shows rounding schemes as applied to the entropic proximal algorithm, whereas panel (b) shows rounding schemes applied to the TRW proximal scheme. In both plots, the tree-based and star-based deterministic schemes are the first to return an optimal solution, whereas the node-based randomized scheme is the slowest in both plots. Of course, this type of ordering is to be expected, since the tree and star-based schemes look over larger neighborhoods of the graph, but incur larger computational cost.

6. Discussion

In this paper, we have developed distributed algorithms, based on the notion of proximal sequences, for solving graph-structured linear programming (LP) relaxations. Our methods respect the graph structure, and so can be scaled to large problems, and they exhibit a superlinear rate of convergence. We have also developed a series of graph-structured rounding schemes that can be used to generate integral solutions along with a certificate of optimality. These optimality certificates allow the algorithm to be terminated in a finite number of iterations.

The structure of our algorithms naturally lends itself to incorporating additional constraints, both linear and other types of conic constraints. It would be interesting to develop an adaptive version of our algorithm, which selectively incorporated new constraints as necessary, and then used the same proximal schemes to minimize the new conic program. Our algorithms for solving the LP are primal-based, so that the updates are in terms of the pseudo-marginals μ that are the primal parameters of the LP. This is contrast to typical message-passing algorithms such as tree-reweighted

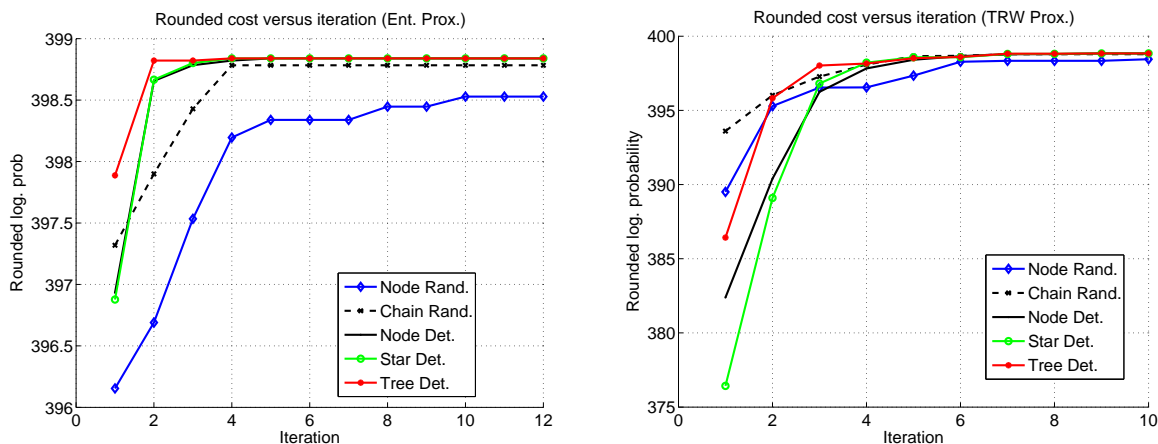


Figure 4: Plots of the log probability of rounded solutions versus the number of iterations for the entropic proximal scheme (panel (a)), and the TRW proximal scheme (panel (b)). In both cases, five different rounding schemes are compared: node-based randomized rounding (Node Rand.), chain-based randomized rounding (Chain Rand.), node-based deterministic rounding (Node Det.), star-based deterministic rounding (Star Det.), and tree-based deterministic rounding (Tree Det.).

max-product, which are dual-based and where the updates are entirely in terms of *message* parameters that are the dual parameters of the LP. However, the dual of the LP is non-differentiable, so that these dual-based updates could either get trapped in local minima (dual co-ordinate ascent) or have sub-linear convergence rates (dual sub-gradient ascent). On the one hand, our primal-based algorithm converges to the LP minimum, and has at least linear convergence rates. On the other, it is more memory-intensive because of the need to maintain $O(|E|)$ edge pseudo-marginal parameters. It would be interesting to modify our algorithms so that maintaining these explicitly could be avoided; note that our derandomized rounding scheme (Algorithm 4.2.1) does not make use of the edge pseudo-marginal parameters.

Acknowledgments

This work was partially supported by NSF grants CCF-0545862 and DMS-0528488 to MJW. AA is partially supported by NSF award DMS-0830410, DARPA award HR0011-08-2-0002 and MSR PhD Fellowship.

Appendix A. Corrections to Bregman Projections

We briefly outline the corrections needed to cyclic Bregman projections for the case where the constraints are linear inequalities. It is useful, in order to characterize these needed corrections, to first note that these cyclic projections are equivalent to co-ordinate ascent steps on the dual of the Bregman projection problem (13). Let the linear constraint set for the Bregman projection

problem (13) be $C \equiv \cap_i \{\langle a_i, \mu \rangle \leq b_i\}$. Its Lagrangian can be written as

$$\mathcal{L}(\mu, z) = D_f(\mu \| \mathbf{v}) + \sum_i z_i (\langle a_i, \mu \rangle - b_i),$$

where $z \geq 0$ are the Lagrangian or dual parameters. The dual function is given as $g(z) = \min_{\mu} \mathcal{L}(\mu, z)$, so that the dual problem can be written as

$$\min_{z \geq 0} g(z).$$

If the constraints were linear *equalities*, the dual variables $\{z\}$ would be unconstrained, and iterative co-ordinate ascent—which can be verified to be equivalent to cyclic projections of the primal variables onto individual constraints—would suffice to solve the dual problem. However, when the constraints have inequalities, the dual problem is no longer unconstrained: the dual variables are constrained to be positive. We would thus need to constrain the co-ordinate ascent steps. This can also be understood as the following primal-dual algorithmic scheme. Note that a necessary KKT condition for optimality of a primal-dual pair (μ, z) for (13) is

$$\nabla f(\mu) = \nabla f(\mathbf{v}) - \sum_i z_i a_i. \tag{32}$$

The primal-dual algorithmic scheme then consists of maintaining primal-dual iterates (μ^t, z^t) which satisfy the equality (32), are dual feasible with $z^t \geq 0$, and which entail co-ordinate ascent on the dual problem, so that $g(z^{t+1}) \geq g(z^t)$ with at most one co-ordinate of μ^t updated in μ^{t+1} . We can now write down the corrected-projection update of μ^t given the single constraint $C_i \equiv \{\langle a_i, \mu \rangle \leq b_i\}$. According to the primal-dual algorithmic scheme this corresponds to co-ordinate ascent on the i -th co-ordinate of z^t so that (32) is maintained, whereby

$$\begin{aligned} \nabla f(\mu^{t+1}) &= \nabla f(\mu^t) + C a_i, \\ z^{t+1} &= z^t - C e_i, \\ C &:= \min\{z_i^t, \beta\}, \end{aligned} \tag{33}$$

where e_i is the co-ordinate vector with one in the i -th co-ordinate and zero elsewhere, and β is the i -th dual parameter setting corresponding to an unconstrained co-ordinate ascent update,

$$\begin{aligned} \nabla f(\mu) &= \nabla f(\mu^n) + \beta a_i, \\ \langle \mu, a_i \rangle &= b_i. \end{aligned} \tag{34}$$

One could derive such corrections corresponding to constrained dual ascent for general convex constraints (Dykstra, 1985; Han, 1988).

Appendix B. Detailed Derivation of Message-passing Updates

In this appendix, we provided detailed derivation of the message-passing updates for the inner loops of the algorithms.

B.1 Derivation of Algorithm 2

Consider the edge marginalization constraint for edge (s, t) , $\mathbb{L}_i(G) \equiv \sum_{x_t} \mu_{st}(x_s, x_t) = \mu_s(x_s)$. Denoting the dual (Lagrange) parameter corresponding to the constraint by $\lambda_{st}(x_s)$, the Karush-Kuhn-Tucker conditions for the quadratic update (18) are given by

$$\begin{aligned} \nabla q(\mu_{st}^{n,\tau+1}(x_s, x_t)) &= \nabla q(\mu_{st}^{n,\tau}(x_s, x_t)) + \lambda_{st}(x_s), \\ \nabla q(\mu_s^{n,\tau+1}(x_s)) &= \nabla q(\mu_s^{n,\tau}(x_s)) - \lambda_{st}(x_s), \\ \mu_{st}^{n,\tau+1}(x_s, x_t) &= \mu_{st}^{n,\tau}(x_s, x_t) + \lambda_{st}(x_s), \\ \mu_s^{n,\tau+1}(x_s) &= \mu_s^{n,\tau}(x_s) - \lambda_{st}(x_s), \end{aligned}$$

while the constraint itself gives

$$\sum_{x_t} \mu_{st}^{n,\tau+1}(x_s, x_t) = \mu_s^{n,\tau}(x_s).$$

Solving for $\lambda_{st}(x_s)$ yields Equation (20). The node marginalization follows similarly.

The only inequalities are the positivity constraints, requiring that the node and edge pseudo-marginals be non-negative. Following the correction procedure for Bregman projections in (33), we maintain Lagrange dual variables corresponding to these constraints. We use $Z_s(x_s)$ as the Lagrange variables for the node positivity constraints $\mu_s(x_s) \geq 0$, and $Z_{st}(x_s, x_t)$ for the edge-positivity constraints $\mu_{st}(x_s, x_t) \geq 0$.

Consider the projection of $\{\mu^{n,\tau+1}\}$ onto the constraint $\mu_s(x_s) \geq 0$. Following (34), we first solve for $\beta_s(x_s)$ that satisfies

$$\begin{aligned} \mu_s(x_s) &= \mu_s^{n,\tau+1}(x_s) - \beta_s(x_s), \\ \mu_s(x_s) &= 0, \end{aligned}$$

so that $\beta_s(x_s) = \mu_s^{n,\tau+1}(x_s)$. Substituting in (33), we obtain the update

$$\begin{aligned} C_s(x_s) &= \min\{Z_s(x_s), \mu_s^{(n,\tau+1)}(x_s)\}, \\ Z_s(x_s) &= Z_s(x_s) - C_s(x_s), \\ \mu_s^{(n,\tau+1)}(x_s) &= \mu_s^{(n,\tau+1)}(x_s) - C_s(x_s). \end{aligned}$$

The edge positivity constraint updates follow similarly.

Thus overall, we obtain message-passing Algorithm 2 for the inner loop.

B.2 Derivation of Algorithm 3

Note that we do not need to explicitly impose positivity constraints in this case. Because the domain of the entropic Bregman function is the positive orthant, if we start from a positive point, any further Bregman projections would also result in a point in the positive orthant.

The projection $\mu^{n,\tau+1} = \Pi_h(\mu^{n,\tau}, \mathbb{L}_i(G))$ onto the individual constraint $\mathbb{L}_i(G)$ is defined by the optimization problem:

$$\mu^{n,\tau+1} = \min_{\mathbb{L}_i(G)} \{h(\mu) - \mu^\top \nabla h(\mu^{n,\tau})\}.$$

Consider the subset $\mathbb{L}_i(G)$ defined by the marginalization constraint along edge (s, t) , namely $\sum_{x'_t \in \mathcal{X}} \mu_{st}(x_s, x'_t) = \mu_s(x_s)$ for each $x_s \in \mathcal{X}$. Denoting the dual (Lagrange) parameters corresponding to these constraint by $\lambda_{st}(x_s)$, the KKT conditions are given by

$$\begin{aligned} \nabla h(\mu_{st}^{n, \tau+1}(x_s, x_t)) &= \nabla h(\mu_{st}^{n, \tau}(x_s, x_t)) + \lambda_{st}(x_s), \quad \text{and} \\ \nabla h(\mu_s^{n, \tau+1}(x_s)) &= \nabla h(\mu_s^{n, \tau}(x_s)) - \lambda_{st}(x_s). \end{aligned}$$

Computing the gradient ∇h and performing some algebra yields the relations

$$\begin{aligned} \mu_{st}^{(n, \tau+1)}(x_s, x_t) &= \mu_{st}^{(n, \tau)}(x_s, x_t) \exp(\lambda_{st}^{(n, \tau+1)}(x_s)), \\ \mu_s^{(n, \tau+1)}(x_s) &= \mu_s^{(n, \tau)}(x_s) \exp(-\lambda_{st}^{(n, \tau+1)}(x_s)), \quad \text{and} \\ \exp(2\lambda_{st}^{(n, \tau+1)}(x_s)) &= \frac{\mu_s^{(n, \tau)}(x_s)}{\sum_{x_t} \mu_{st}^{(n, \tau)}(x_s, x_t)}, \end{aligned}$$

from which the updates (22) follow.

Similarly, for the constraint set defined by the node marginalization constraint $\sum_{x_s \in \mathcal{X}} \mu_s(x_s) = 1$, we have $\nabla h(\mu_s^{(n, \tau+1)}(x_s)) = \nabla h(\mu_s^{(n, \tau)}(x_s)) + \lambda_s^{(n, \tau+1)}$, from which

$$\begin{aligned} \mu_s^{(n, \tau+1)}(x_s) &= \mu_s^{(n, \tau)}(x_s) \exp(\lambda_s^{(n, \tau+1)}), \quad \text{and} \\ \exp(\lambda_s^{(n, \tau+1)}) &= 1 / \sum_{x_s \in \mathcal{X}} \mu_s^{(n, \tau)}(x_s). \end{aligned}$$

The updates in Equation (23) follow.

Appendix C. Proof of Lemma 4

We provide a detailed proof for the entropic scheme; the arguments for other proximal algorithms are analogous. The key point is the following: regardless of how the proximal updates are computed, they must satisfy the necessary Lagrangian conditions for optimal points over the set $\mathbb{L}(G)$. Accordingly, we define the following sets of Lagrange multipliers:

$$\begin{aligned} \lambda_{ss} & \quad \text{for the normalization constraint } C_{ss}(\mu_s) = \sum_{x'_s} \mu_s(x'_s) - 1 = 0, \\ \lambda_{st}(x_s) & \quad \text{for the marginalization constraint } C_{ts}(x_s) = \sum_{x'_t} \mu_{st}(x_s, x'_t) - \mu_s(x_s) = 0, \\ \gamma_{st}(x_s, x_t) & \quad \text{for the non-negativity constraint } \mu_{st}(x_s, x_t) \geq 0. \end{aligned}$$

(There is no need to enforce the non-negativity constraint $\mu_s(x_s) \geq 0$ directly, since it is implied by the non-negativity of the joint pseudo-marginals and the marginalization constraints.)

With this notation, consider the Lagrangian associated with the entropic proximal update at step n :

$$L(x; \lambda, \gamma) = C(\mu; \theta, \mu^n) + \langle \gamma, \mu \rangle + \sum_{s \in \mathcal{V}} \lambda_{ss} C_{ss}(x_s) + \sum_{(s,t) \in E} [\lambda_{ts}(x_s) C_{ts}(x_s) + \lambda_{st}(x_t) C_{st}(x_t)],$$

where $C(\mu; \theta, \mu^n)$ is shorthand for the cost component $-\langle \theta, \mu \rangle + \frac{1}{\omega^n} D_\alpha(\mu \| \mu^n)$. Using C, C' to denote constants (whose value can change from line to line), we now take derivatives to find the necessary

Lagrangian conditions:

$$\begin{aligned}\frac{\partial L}{\partial \mu_s(x_s)} &= -\theta_s(x_s) + \frac{2\alpha_s}{\omega^n} \log \frac{\mu_s(x_s)}{\mu_s^n(x_s)} + C + \lambda_{ss} + \sum_{t \in N(s)} \lambda_{ts}(x_s), \quad \text{and} \\ \frac{\partial L}{\partial \mu_{st}(x_s, x_t)} &= -\theta_{st}(x_s, x_t) + \frac{2\alpha_{st}}{\omega^n} \log \frac{\mu_{st}(x_s, x_t)}{\mu_{st}^n(x_s, x_t)} + C' + \gamma_{st}(x_s, x_t) - \lambda_{ts}(x_s) - \lambda_{st}(x_t).\end{aligned}$$

Solving for the optimum $\mu = \mu^{n+1}$ yields

$$\begin{aligned}\frac{2\alpha_s}{\omega^n} \log \mu_s^{n+1}(x_s) &= \theta_s(x_s) + \frac{2\alpha_s}{\omega^n} \log \mu_s^n(x_s) - \sum_{t \in N(s)} \lambda_{ts}(x_s) + C, \\ \frac{2\alpha_{st}}{\omega^n} \log \mu_{st}^{n+1}(x_s, x_t) &= \theta_{st}(x_s, x_t) + \frac{2\alpha_{st}}{\omega^n} \log \mu_{st}^n(x_s, x_t) - \gamma_{st}(x_s, x_t) \\ &\quad + \lambda_{ts}(x_s) + \lambda_{st}(x_t) + C'.\end{aligned}$$

From these conditions, we can compute the energy invariant (30):

$$\begin{aligned}\frac{2}{\omega^n} F(x; \mu^{n+1}) &= \sum_{s \in V} \frac{2\alpha_s}{\omega^n} \log \mu_s^{n+1}(x_s) + \sum_{(s,t) \in E} \frac{2\alpha_{st}}{\omega^n} \log \mu_{st}^{n+1}(x_s, x_t) + C \\ &= F(x; \theta) + \frac{2}{\omega^n} \left\{ \sum_{s \in V} \alpha_s \log \mu_s^n(x_s) + \sum_{(s,t) \in E} \alpha_{st} \log \mu_{st}^n(x_s, x_t) \right\} \\ &\quad - \sum_{(s,t) \in E} \gamma_{st}(x_s, x_t) + C \\ &= F(x; \theta) + \frac{2}{\omega^n} F(x; \mu^n) - \sum_{(s,t) \in E} \gamma_{st}(x_s, x_t) + C.\end{aligned}$$

Now since $\mu^n > 0$, by complementary slackness, we must have $\gamma_{st}(x_s, x_t) = 0$, which implies that

$$\frac{2}{\omega^n} F(x; \mu^{n+1}) = F(x; \theta) + \frac{2}{\omega^n} F(x; \mu^n) + C. \quad (35)$$

From this equation, it is a simple induction to show for some constants $\gamma_n > 0$ and $C_n \in \mathbb{R}$, we have $F(x; \mu^n) = \gamma_n F(x; \theta) + C_n$ for all iterations $n = 1, 2, 3, \dots$, which implies preservation of the maximizers. If at iteration $n = 0$, we initialize $\mu^0 = 0$ to the all-uniform distribution, then we have $\frac{2}{\omega^0} F(x; \mu^1) = F(x; \theta) + C'$, so the statement follows for $n = 1$. Suppose that it holds at step n ; then $\frac{2}{\omega^n} F(x; \mu^n) = \frac{2}{\omega^n} \gamma_n F(x; \theta) + \frac{2C_n}{\omega^n}$, and hence from the induction step (35), we have $F(x; \mu^{n+1}) = \gamma_{n+1} F(x; \theta) + C_{n+1}$, where $\gamma_{n+1} = \frac{\omega^n}{2} \gamma_n$.

Appendix D. Proof of Theorem 6

Consider the expected cost of the configuration $X(\mu^n; E')$ obtained from the randomized rounding procedure of Algorithm 5. A simple computation shows that

$$\mathbb{E}[F(X(\mu^n; E'); \theta)] = G(\bar{\mu}) := \sum_{i=1}^K H(\mu^n; T_i) + H(\mu^n; E'),$$

where

$$\begin{aligned} H(\mu^n; T_i) &:= \sum_{s \in V_i} \sum_{x_s} \mu_s^n(x_s) \theta_s(x_s) + \sum_{(s,t) \in E_i} \sum_{x_s, x_t} \mu_{st}^n(x_s, x_t) \theta_{st}(x_s, x_t), \\ H(\mu^n; E') &:= \sum_{(u,v) \in E'} \sum_{x_u, x_v} \mu_u^n(x_u) \mu_v^n(x_v) \theta_{st}(x_u, x_v). \end{aligned} \quad (36)$$

We now show by induction that the de-randomized rounding scheme achieves cost at least as large as this expected value. Let $\bar{\mu}^{(i)}$ denote the updated pseudomarginals at the end of the i -th iteration. Since we initialize with $\bar{\mu}^{(0)} = \mu^n$, we have $G(\bar{\mu}^{(0)}) = \mathbb{E}[F(X(\mu^n; E'); \theta)]$. Consider the i -th step of the algorithm; the algorithm computes the portion of the de-randomized solution $x_{V_i}^d$ over the i -th tree. It will be convenient to use the decomposition $G = G_i + G_{\setminus i}$, where

$$\begin{aligned} G_i(\bar{\mu}) &:= \sum_{s \in V_i} \sum_{x_s} \bar{\mu}_s(x_s) \left\{ \theta_s(x_s) + \sum_{\{t \mid (s,t) \in E'\}} \sum_{x_t} \bar{\mu}_t(x_t) \theta_{st}(x_s, x_t) \right\} + \\ &\quad \sum_{(s,t) \in E_i} \sum_{x_s, x_t} \bar{\mu}_{st}(x_s, x_t) \theta_{st}(x_s, x_t), \end{aligned}$$

and $G_{\setminus i} = G - G_i$. If we define

$$\bar{F}_i(x_{V_i}) := \sum_{s \in V_i} \left\{ \theta_s(x_s) + \sum_{t: (s,t) \in E'} \sum_{x_t} \bar{\mu}_t^{(i-1)}(x_t) \theta_{st}(x_s, x_t) \right\} + \sum_{(s,t) \in E_i} \theta_{st}(x_s, x_t),$$

it can be seen that $G_i(\bar{\mu}^{(i-1)}) = \mathbb{E}[\bar{F}_i(x_{V_i})]$ where the expectation is under the tree-structured distribution over X_{V_i} given by

$$p(x_{V_i}; \bar{\mu}^{(i-1)}(T_i)) = \prod_{s \in V_i} \bar{\mu}^{(i-1)}(x_s) \prod_{(s,t) \in E_i} \frac{\bar{\mu}^{(i-1)}(x_s, x_t)}{\bar{\mu}^{(i-1)}(x_s) \bar{\mu}^{(i-1)}(x_t)}.$$

Thus when the algorithm makes the choice $x_{V_i}^d = \arg \max_{x_{V_i}} \bar{F}_i(x_{V_i})$, it holds that

$$G_i(\bar{\mu}^{(i-1)}) = \mathbb{E}[\bar{F}_i(x_{V_i})] \leq \bar{F}_i(x_{V_i}^d).$$

The updated pseudomarginals $\bar{\mu}^{(i)}$ at the end the i -th step of the algorithm are given by,

$$\begin{aligned} \bar{\mu}_s^{(i)}(x_s) &= \begin{cases} \bar{\mu}_s^{(i-1)}(x_s) & \text{if } s \notin V_i \\ 0 & \text{if } s \in V_i, X_{d,s} \neq x_s \\ 1 & \text{if } s \in V_i, X_{d,s} = x_s. \end{cases} \\ \bar{\mu}_{st}^{(i)}(x_s, x_t) &= \begin{cases} \bar{\mu}_{st}^{(i-1)}(x_s, x_t) & \text{if } (s,t) \notin E_i \\ \bar{\mu}_s^{(i)}(x_s) \bar{\mu}_t^{(i)}(x_t) & \text{if } (s,t) \in E_i. \end{cases} \end{aligned}$$

In other words, $\bar{\mu}^{(i)}(T_i)$ is the indicator vector of the maximum energy sub-configuration $x_{V_i}^d$. Consequently, we have

$$G_i(\bar{\mu}^{(i)}) = \bar{F}_i(x_{V_i}^d) \geq G_i(\bar{\mu}^{(i-1)}),$$

and $G_{\setminus i}(\bar{\mu}^{(i)}) = G_{\setminus i}(\bar{\mu}^{(i-1)})$, so that at the end of the i -th step, $G(\bar{\mu}^{(i)}) \geq G(\bar{\mu}^{(i-1)})$. By induction, we conclude that $G(\bar{\mu}^{(K)}) \geq G(\bar{\mu}^{(0)})$, where K is the total number of trees in the rounding scheme.

At the end of K steps, the quantity $\bar{\mu}^{(K)}$ is the indicator vector for $x^d(\mu^n; E')$ so that $G(\bar{\mu}^{(K)}) = F(X_d(\mu^n; E'); \theta)$. We have also shown that $G(\bar{\mu}^{(0)}) = \mathbb{E}[F(X(\mu^n; E'); \theta)]$. Combining these pieces, we conclude that $F(x^d(\mu^n; E'); \theta) \geq \mathbb{E}[F(X(\mu^n; E'); \theta)]$, thereby completing the proof.

Appendix E. Proof of Theorem 7

Let $p_{\text{succ}} = \mathbb{P}[X(\mu^n; E') = x^*]$, and let $R(\mu^n; E')$ denote the (random) integral vertex of $\mathbb{L}(G)$ that is specified by the random integral solution $X(\mu^n; E')$. (Since E' is some fixed forest-inducing subset, we frequently shorten this notation to $R(\mu^n)$.) We begin by computing the expected cost of the random solution, where the expectation is taken over the rounding procedure. A simple computation shows that $\mathbb{E}[\langle \theta, R(\mu^n) \rangle] := \sum_{i=1}^K H(\mu^n; T_i) + H(\mu^n; E')$, where $H(\mu^n; T_i)$ and $H(\mu^n; E')$ were defined previously (36).

We now upper bound the difference $\langle \theta, \mu^* \rangle - \mathbb{E}[\langle \theta, R(\mu^n) \rangle]$. For each subtree $i = 1, \dots, K$, the quantity $D_i := H(\mu^*; T_i) - H(\mu^n; T_i)$ is upper bounded as

$$\begin{aligned} D_i &= \sum_{s \in V_i} \sum_{x_s} \left[\mu_s^*(x_s) - \mu_s^n(x_s) \right] \theta_s(x_s) + \sum_{(s,t) \in E_i} \sum_{x_s, x_t} \left[\mu_s^*(x_s) \mu_t^*(x_t) - \mu_{st}^n(x_s, x_t) \right] \theta_{st}(x_s, x_t) \\ &\leq \sum_{s \in V_i} \delta_s(\theta) \sum_{x_s} |\mu_s^*(x_s) - \mu_s^n(x_s)| + \sum_{(s,t) \in E_i} \delta_{st}(\theta) \sum_{x_s, x_t} |\mu_{st}^*(x_s, x_t) - \mu_{st}^n(x_s, x_t)|. \end{aligned}$$

In asserting this inequality, we have used the fact that the matrix with entries given by $\mu_s^*(x_s) \mu_t^*(x_t) - \mu_{st}^n(x_s, x_t)$ is a difference of probability distributions, meaning that all its entries are between -1 and 1 , and their sum is zero.

Similarly, we can upper bound the difference $D(E') = H(\mu^*; E') - H(\mu^n; E')$ associated with E' :

$$\begin{aligned} D(E') &= \sum_{(u,v) \in E'} \sum_{x_u, x_v} \left[\mu_u^*(x_u) \mu_v^*(x_v) - \mu_{uv}^n(x_u, x_v) \right] \theta_{uv}(x_u, x_v) \\ &\leq \sum_{(u,v) \in E'} \delta_{uv}(\theta) \sum_{x_u, x_v} \left| \mu_u^*(x_u) \mu_v^*(x_v) - \mu_{uv}^n(x_u, x_v) \right| \\ &\leq \sum_{(u,v) \in E'} \delta_{uv}(\theta) \sum_{x_u, x_v} \left\{ \left| \mu_u^*(x_u) [\mu_v^*(x_v) - \mu_v^n(x_v)] \right| + \left| \mu_v^n(x_v) [\mu_u^*(x_u) - \mu_u^n(x_u)] \right| \right\} \\ &\leq \sum_{(u,v) \in E'} \delta_{uv}(\theta) \left\{ \sum_{x_u} |\mu_u^n(x_u) - \mu_u^*(x_u)| + \sum_{x_v} |\mu_v^n(x_v) - \mu_v^*(x_v)| \right\}. \end{aligned}$$

Combining the pieces, we obtain

$$\begin{aligned} \langle \theta, \mu^* \rangle - \mathbb{E}[\langle \theta, R(\mu^n) \rangle] &\leq \delta_G(\theta) \left\{ \|\mu^n - \mu^*\|_1 + \sum_{s \in V} d(s; E') \sum_{x_s} |\mu_s^n(x_s) - \mu_s^*(x_s)| \right\} \\ &\leq (1 + d(E')) \delta_G(\theta) \|\mu^n - \mu^*\|_1. \end{aligned} \tag{37}$$

In the other direction, we note that when the rounding fails, then we have

$$\langle \theta, \mu^* \rangle - \langle \theta, R(\mu^n) \rangle \geq \max_{x \neq x^*} [F(x^*; \theta) - F(x; \theta)].$$

Consequently, conditioning on whether the rounding succeeds or fails, we have

$$\begin{aligned} \langle \theta, \mu^* \rangle - \mathbb{E}[\langle \theta, R(\mu^n) \rangle] &\geq p_{\text{succ}} [\langle \theta, \mu^* \rangle - \langle \theta, \mu^* \rangle] + (1 - p_{\text{succ}}) \max_{x \neq x^*} [F(x^*; \theta) - F(x; \theta)] \\ &= (1 - p_{\text{succ}}) \max_{x \neq x^*} [F(x^*; \theta) - F(x; \theta)]. \end{aligned}$$

Combining this lower bound with the upper bound (37), performing some algebra, and using the definition of the gap $\Delta(\theta; G)$ yields that the probability of successful rounding is at least

$$p_{\text{succ}} \geq 1 - \frac{(1 + d(E'))}{\Delta(\theta; G)} \|\mu^n - \mu^*\|_1.$$

If the condition (31) holds, then this probability is at least $1 - \varepsilon$, as claimed.

Appendix F. Proof of Theorem 8

The proof follows that of Theorem 7 until Equation (37), which gives

$$\langle \theta, \mu^* \rangle - \mathbb{E}[\langle \theta, R(\mu^n) \rangle] \leq (1 + d(E')) \delta_G(\theta) \|\mu^n - \mu^*\|_1.$$

Let $v^d(\mu^n; E')$ denote the integral vertex of $\mathbb{L}(G)$ that is specified by the de-randomized integral solution $x^d(\mu^n; E')$. Since E' is some fixed forest-inducing subset, we frequently shorten this notation to $v^d(\mu^n)$. Theorem 6 shows that

$$\mathbb{E}[\langle \theta, R(\mu^n) \rangle] \leq \langle \theta, v^d(\mu^n) \rangle.$$

Suppose the de-randomized solution is not optimal so that $v^d(\mu^n) \neq \mu^*$. Then, from the definition of the graph-based gap $\Delta(\theta; G)$, we obtain

$$\langle \theta, \mu^* \rangle - \langle \theta, v^d(\mu^n) \rangle \geq \delta_G(\theta) \Delta(\theta; G).$$

Combining the pieces, we obtain

$$\begin{aligned} \delta_G(\theta) \Delta(\theta; G) &\leq \langle \theta, \mu^* \rangle - \langle \theta, v^d(\mu^n) \rangle \\ &\leq \langle \theta, \mu^* \rangle - \mathbb{E}[\langle \theta, R(\mu^n) \rangle] \\ &\leq (1 + d(E')) \delta_G(\theta) \|\mu^n - \mu^*\|_1, \end{aligned}$$

which implies $\|\mu^n - \mu^*\|_1 \geq \frac{\Delta(\theta; G)}{1 + d(E')}$. However, this conclusion is a contradiction under the given assumption on $\|\mu^n - \mu^*\|_1$ in the theorem. It thus holds that the de-randomized solution $v^d(\mu^n)$ is equal to the MAP optimum μ^* , thereby completing the proof.

References

- H. H. Bauschke and J. M. Borwein. Legendre functions and the method of random bregman projections. *Journal of Convex Analysis*, 4(1):27–67, 1997.
- M. Bayati, D. Shah, and M. Sharma. Maximum weight matching for max-product belief propagation. In *Int. Symp. Info. Theory*, Adelaide, Australia, September 2005.
- M. Bayati, C. Borgs, J. Chayes, and R. Zecchina. Belief-propagation for weighted b-matchings on arbitrary graphs and its relation to linear programs with integer solutions. Technical Report arxiv:0709.1190, Microsoft Research, September 2007.
- U. Bertele and F. Brioschi. *Nonserial Dynamic Programming*. Academic Press, New York, 1972.

- D. P. Bertsekas and J. N. Tsitsiklis. *Parallel and Distributed Computation: Numerical Methods*. Athena Scientific, Boston, MA, 1997.
- D.P. Bertsekas. *Nonlinear Programming*. Athena Scientific, Belmont, MA, 1995.
- D. Bertsimas and J. Tsitsiklis. *Introduction to Linear Optimization*. Athena Scientific, Belmont, MA, 1997.
- J. Besag. On the statistical analysis of dirty pictures. *Journal of the Royal Statistical Society, Series B*, 48(3):259–279, 1986.
- S. Boyd and L. Vandenberghe. *Convex Optimization*. Cambridge University Press, Cambridge, UK, 2004.
- Y. Boykov, O. Veksler, and R. Zabih. Fast approximate energy minimization via graph cuts. *IEEE Trans. Pattern Anal. Mach. Intell.*, 23(11):1222–1239, 2001.
- Y. Censor and S. A. Zenios. *Parallel Optimization - Theory, Algorithms and Applications*. Oxford University Press, 1997.
- C. Chekuri, S. Khanna, J. Naor, and L. Zosin. A linear programming formulation and approximation algorithms for the metric labeling problem. *SIAM Journal on Discrete Mathematics*, 18(3):608–625, 2005.
- F. Deutsch and H. Hundal. The rate of convergence for the cyclic projection algorithm I: Angles between convex sets. *Journal of Approximation Theory*, 142:36–55, 2006.
- R. L. Dykstra. An iterative procedure for obtaining i-projections onto the intersection of convex sets. *Annals of Probability*, 13(3):975–984, 1985.
- J. Feldman, D. R. Karger, and M. J. Wainwright. Linear programming-based decoding of turbo-like codes and its relation to iterative approaches. In *Proc. 40th Annual Allerton Conf. on Communication, Control, and Computing*, October 2002a.
- J. Feldman, M. J. Wainwright, and D. R. Karger. Linear programming-based decoding of turbo-like codes and its relation to iterative approaches. In *Proc. Allerton Conf. Comm., Control and Computing*, October 2002b.
- W. T. Freeman and Y. Weiss. On the optimality of solutions of the max-product belief propagation algorithm in arbitrary graphs. *IEEE Trans. Info. Theory*, 47:736–744, 2001.
- S. Geman and D. Geman. Stochastic relaxation, Gibbs distributions, and the Bayesian restoration of images. *IEEE Trans. PAMI*, 6:721–741, 1984.
- A. Globerson and T. Jaakkola. Convergent propagation algorithms via oriented trees. In *Proc. Uncertainty in Artificial Intelligence*, Vancouver, Canada, July 2007a.
- A. Globerson and T. Jaakkola. Fixing max-product: Convergent message passing algorithms for MAP LP-relaxations. In *Neural Information Processing Systems*, Vancouver, Canada, December 2007b.

- D.M. Greig, B.T. Porteous, and A.H. Seheult. Exact maximum a posteriori estimation for binary images. *Journal of the Royal Statistical Society, Series B*, 51, 1989.
- S.-P. Han. A successive projection method. *Math. Programming*, 40:114, 1988.
- T. Hazan and A. Shashua. Convergent message-passing algorithms for inference over general graphs with convex free energy. In *The 24th Conference on Uncertainty in Artificial Intelligence (UAI)*, Helsinki, 2008.
- B. Huang and T. Jebara. Loopy belief propagation for bipartite maximum weight b-matching. In *AISTATS*, San Juan, Puerto Rico, March 2007.
- A. N. Iusem and M. Teboulle. Convergence rate analysis of nonquadratic proximal methods for convex and linear programming. *Mathematics of Operations Research*, 20(3):657–677, 1995.
- J. K. Johnson, D. M. Malioutov, and A. S. Willsky. Lagrangian relaxation for MAP estimation in graphical models. In *Proc. 45th Allerton Conf. Comm. Cont. Comp.*, Urbana-Champaign, IL, September 2007.
- J. Kleinberg and E. Tardos. Approximation algorithms for classification problems with pairwise relationships: Metric labeling and Markov random fields. In *IEEE Symp. Found. Comp. Science*, pages 14–24, 1999.
- V. Kolmogorov. Convergent tree-reweighted message-passing for energy minimization. In *International Workshop on Artificial Intelligence and Statistics*, January 2005.
- V. Kolmogorov. Convergent tree-reweighted message-passing for energy minimization. *IEEE Trans. PAMI*, 28(10):1568–1583, October 2006.
- V. Kolmogorov and M. J. Wainwright. On optimality properties of tree-reweighted message-passing. In *Uncertainty in Artificial Intelligence*, July 2005.
- N. Komodakis, N. Paragios, and G. Tziritas. MRF optimization via dual decomposition: Message-passing revisited. In *ICCV*, Rio di Janeiro, Brazil, 2007.
- V. K. Koval and M. I. Schlesinger. Two-dimensional programming in image analysis problems. *USSR Academy of Science, Automatics and Telemechanics*, 8:149–168, 1976. In Russian.
- P. Kumar, P.H.S. Torr, and A. Zisserman. Solving markov random fields using second order cone programming. *IEEE Conference of Computer Vision and Pattern Recognition*, 2006.
- P. Raghavan and C. D. Thompson. Randomized rounding: A technique for provably good algorithms and algorithmic proofs. *Combinatorica*, 7(4):365–374, 1987.
- P. Ravikumar and J. Lafferty. Quadratic programming relaxations for metric labeling and markov random field map estimation. In *ICML '06: Proceedings of the 23rd international conference on Machine learning*, pages 737–744, 2006.
- G. Rockafellar. *Convex Analysis*. Princeton University Press, Princeton, 1970.

- S. Sanghavi, D. Shah, and A. Willsky. Message-passing for max-weight independent set. In *Neural Information Processing Systems*, Vancouver, Canada, December 2007.
- M.V. Solodov and B.F. Svaiter. A unified framework for some inexact proximal point algorithms. *Numerical Functional Analysis and Optimization*, 22:1013–1035, 2001.
- P. Tseng and D. P. Bertsekas. On the convergence of the exponential multiplier method for convex programming. *Math. Programming*, 60:1—19, 1993.
- M. J. Wainwright and M. I. Jordan. Graphical models, exponential families, and variational inference. Technical report, UC Berkeley, Department of Statistics, No. 649, September 2003.
- M. J. Wainwright and M. I. Jordan. Treewidth-based conditions for exactness of the Sherali-Adams and Lasserre relaxations. Technical report, UC Berkeley, Department of Statistics, No. 671, September 2004.
- M. J. Wainwright, T. S. Jaakkola, and A. S. Willsky. A new class of upper bounds on the log partition function. In *Uncertainty in Artificial Intelligence*, volume 18, August 2002.
- 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*, 49(5):1120–1146, May 2003.
- M. J. Wainwright, T. S. Jaakkola, and A. S. Willsky. Tree consistency and bounds on the max-product algorithm and its generalizations. *Statistics and Computing*, 14:143–166, April 2004.
- M. J. Wainwright, T. S. Jaakkola, and A. S. Willsky. Exact MAP estimates via agreement on (hyper)trees: Linear programming and message-passing. *IEEE Trans. Information Theory*, 51(11):3697–3717, November 2005.
- Y. Weiss, C. Yanover, and T. Meltzer. Map estimation, linear programming and belief propagation with convex free energies. In *Uncertainty in Artificial Intelligence*, 2007.
- C. Yanover, T. Meltzer, and Y. Weiss. Linear programming relaxations and belief propagation: An empirical study. *Journal of Machine Learning Research*, 7:1887–1907, September 2006.
- J.S. Yedidia, W. T. Freeman, and Y. Weiss. Constructing free energy approximations and generalized belief propagation algorithms. *IEEE Trans. Info. Theory*, 51(7):2282–2312, July 2005.