
Concavity of reweighted Kikuchi approximation

Po-Ling Loh

Department of Statistics
The Wharton School
University of Pennsylvania
loh@wharton.upenn.edu

Andre Wibisono

Computer Science Division
University of California, Berkeley
wibisono@berkeley.edu

Abstract

We analyze a reweighted version of the Kikuchi approximation for estimating the log partition function of a product distribution defined over a region graph. We establish sufficient conditions for the concavity of our reweighted objective function in terms of weight assignments in the Kikuchi expansion, and show that a reweighted version of the sum product algorithm applied to the Kikuchi region graph will produce global optima of the Kikuchi approximation whenever the algorithm converges. When the region graph has two layers, corresponding to a Bethe approximation, we show that our sufficient conditions for concavity are also necessary. Finally, we provide an explicit characterization of the polytope of concavity in terms of the cycle structure of the region graph. We conclude with simulations that demonstrate the advantages of the reweighted Kikuchi approach.

1 Introduction

Undirected graphical models are a familiar framework in diverse application domains such as computer vision, statistical physics, coding theory, social science, and epidemiology. In certain settings of interest, one is provided with potential functions defined over nodes and (hyper)edges of the graph. A crucial step in probabilistic inference is to compute the log partition function of the distribution based on these potential functions for a given graph structure. However, computing the log partition function either exactly or approximately is NP-hard in general [2, 17]. An active area of research involves finding accurate approximations of the log partition function and characterizing the graph structures for which such approximations may be computed efficiently [29, 22, 7, 19, 25, 18].

When the underlying graph is a tree, the log partition function may be computed exactly via the sum product algorithm in time linear in the number of nodes [15]. However, when the graph contains cycles, a generalized version of the sum product algorithm known as loopy belief propagation may either fail to converge or terminate in local optima of a nonconvex objective function [26, 20, 8, 13].

In this paper, we analyze the Kikuchi approximation method, which is constructed from a variational representation of the log partition function by replacing the entropy with an expression that decomposes with respect to a region graph. Kikuchi approximations were previously introduced in the physics literature [9] and reformalized by Yedidia et al. [28, 29] and others [1, 14] in the language of graphical models. The Bethe approximation, which is a special case of the Kikuchi approximation when the region graph has only two layers, has been studied by various authors [3, 28, 5, 25]. In addition, a reweighted version of the Bethe approximation was proposed by Wainwright et al. [22, 16]. As described in Vontobel [21], computing the global optimum of the Bethe variational problem may in turn be used to approximate the permanent of a nonnegative square matrix.

The particular objective function that we study generalizes the Kikuchi objective appearing in previous literature by assigning arbitrary weights to individual terms in the Kikuchi entropy expansion. We establish necessary and sufficient conditions under which this class of objective functions is concave, so a global optimum may be found efficiently. Our theoretical results synthesize known results on Kikuchi and Bethe approximations, and our main theorem concerning concavity conditions for the reweighted Kikuchi entropy recovers existing results when specialized to the unweighted

Kikuchi [14] or reweighted Bethe [22] case. Furthermore, we provide a valuable converse result in the reweighted Bethe case, showing that when our concavity conditions are violated, the entropy function cannot be concave over the whole feasible region. As demonstrated by our experiments, a message-passing algorithm designed to optimize the Kikuchi objective may terminate in local optima for weights outside the concave region. Watanabe and Fukumizu [24, 25] provide a similar converse in the unweighted Bethe case, but our proof is much simpler and our result is more general.

In the reweighted Bethe setting, we also present a useful characterization of the concave region of the Bethe entropy function in terms of the geometry of the graph. Specifically, we show that if the region graph consists of only singleton vertices and pairwise edges, then the region of concavity coincides with the convex hull of incidence vectors of single-cycle forest subgraphs of the original graph. When the region graph contains regions with cardinality greater than two, the latter region may be strictly contained in the former; however, our result provides a useful way to generate weight vectors within the region of concavity. Whereas Wainwright et al. [22] establish the concavity of the reweighted Bethe objective on the spanning forest polytope, that region is contained within the single-cycle forest polytope, and our simulations show that generating weight vectors in the latter polytope may yield closer approximations to the log partition function.

The remainder of the paper is organized as follows: In Section 2, we review background information about the Kikuchi and Bethe approximations. In Section 3, we provide our main results on concavity conditions for the reweighted Kikuchi approximation, including a geometric characterization of the region of concavity in the Bethe case. Section 4 outlines the reweighted sum product algorithm and proves that fixed points correspond to global optima of the Kikuchi approximation. Section 5 presents experiments showing the improved accuracy of the reweighted Kikuchi approximation over the region of concavity. Technical proofs and additional simulations are contained in the Appendix.

2 Background and problem setup

In this section, we review basic concepts of the Kikuchi approximation and establish some terminology to be used in the paper.

Let $G = (V, R)$ denote a *region graph* defined over the vertex set V , where each region $r \in R$ is a subset of V . Directed edges correspond to inclusion, so $r \rightarrow s$ is an edge of G if $s \subseteq r$. We use the following notation, for $r \in R$:

$$\begin{aligned}\mathcal{A}(r) &:= \{s \in R: r \subsetneq s\} && (\text{ancestors of } r) \\ \mathcal{F}(r) &:= \{s \in R: r \subseteq s\} && (\text{forebears of } r) \\ N(r) &:= \{s \in R: r \subseteq s \text{ or } s \subseteq r\} && (\text{neighbors of } r).\end{aligned}$$

For $R' \subseteq R$, we define $\mathcal{A}(R') = \bigcup_{r \in R'} \mathcal{A}(r)$, and we define $\mathcal{F}(R')$ and $N(R')$ similarly.

We consider joint distributions $x = (x_s)_{s \in V}$ that factorize over the region graph; i.e.,

$$p(x) = \frac{1}{Z(\alpha)} \prod_{r \in R} \alpha_r(x_r), \quad (1)$$

for potential functions $\alpha_r > 0$. Here, $Z(\alpha)$ is the normalization factor, or partition function, which is a function of the potential functions α_r , and each variable x_s takes values in a finite discrete set \mathcal{X} . One special case of the factorization (1) is the pairwise Ising model, defined over a graph $G = (V, E)$, where the distribution is given by

$$p_\gamma(x) = \exp \left(\sum_{s \in V} \gamma_s(x_s) + \sum_{(s,t) \in E} \gamma_{st}(x_s, x_t) - A(\gamma) \right), \quad (2)$$

and $\mathcal{X} = \{-1, +1\}$. Our goal is to analyze the log partition function

$$\log Z(\alpha) = \log \left\{ \sum_{x \in \mathcal{X}^{|V|}} \prod_{r \in R} \alpha_r(x_r) \right\}. \quad (3)$$

2.1 Variational representation

It is known from the theory of graphical models [14] that the log partition function (3) may be written in the variational form

$$\log Z(\alpha) = \sup_{\{\tau_r(x_r)\} \in \Delta_R} \left\{ \sum_{r \in R} \sum_{x_r} \tau_r(x_r) \log(\alpha_r(x_r)) + H(p_\tau) \right\}, \quad (4)$$

where p_τ is the maximum entropy distribution with marginals $\{\tau_r(x_r)\}$ and

$$H(p) := - \sum_x p(x) \log p(x)$$

is the usual entropy. Here, Δ_R denotes the R -marginal polytope; i.e., $\{\tau_r(x_r): r \in R\} \in \Delta_R$ if and only if there exists a distribution $\tau(x)$ such that $\tau_r(x_r) = \sum_{x_{\setminus r}} \tau(x_r, x_{\setminus r})$ for all r . For ease of notation, we also write $\tau \equiv \{\tau_r(x_r): r \in R\}$. Let $\theta \equiv \theta(x)$ denote the collection of log potential functions $\{\log(\alpha_r(x_r)): r \in R\}$. Then equation (4) may be rewritten as

$$\log Z(\theta) = \sup_{\tau \in \Delta_R} \{\langle \theta, \tau \rangle + H(p_\tau)\}. \quad (5)$$

Specializing to the Ising model (2), equation (5) gives the variational representation

$$A(\gamma) = \sup_{\mu \in \mathbb{M}} \{\langle \gamma, \mu \rangle + H(p_\mu)\}, \quad (6)$$

which appears in Wainwright and Jordan [23]. Here, $\mathbb{M} \equiv \mathbb{M}(G)$ denotes the marginal polytope, corresponding to the collection of mean parameter vectors of the sufficient statistics in the exponential family representation (2), ranging over different values of γ , and p_μ is the maximum entropy distribution with mean parameters μ .

2.2 Reweighted Kikuchi approximation

Although the set Δ_R appearing in the variational representation (5) is a convex polytope, it may have exponentially many facets [23]. Hence, we replace Δ_R with the set

$$\Delta_R^K = \left\{ \tau: \forall t, u \in R \text{ s.t. } t \subseteq u, \sum_{x_{u \setminus t}} \tau_u(x_t, x_{u \setminus t}) = \tau_t(x_t) \text{ and } \forall u \in R, \sum_{x_u} \tau_u(x_u) = 1 \right\}$$

of *locally consistent* R -pseudomarginals. Note that $\Delta_R \subseteq \Delta_R^K$ and the latter set has only polynomially many facets, making optimization more tractable.

In the case of the pairwise Ising model (2), we let $\mathbb{L} \equiv \mathbb{L}(G)$ denote the polytope Δ_R^K . Then \mathbb{L} is the collection of nonnegative functions $\tau = (\tau_s, \tau_{st})$ satisfying the marginalization constraints

$$\begin{aligned} \sum_{x_s} \tau_s(x_s) &= 1, & \forall s \in V, \\ \sum_{x_t} \tau_{st}(x_s, x_t) &= \tau_s(x_s) \text{ and } \sum_{x_s} \tau_{st}(x_s, x_t) = \tau_t(x_t), & \forall (s, t) \in E. \end{aligned}$$

Recall that $\mathbb{M}(G) \subseteq \mathbb{L}(G)$, with equality achieved if and only if the underlying graph G is a tree. In the general case, we have $\Delta_R = \Delta_R^K$ when the Hasse diagram of the region graph admits a minimal representation that is loop-free (cf. Theorem 2 of Pakzad and Anantharam [14]).

Given a collection of R -pseudomarginals τ , we also replace the entropy term $H(p_\tau)$, which is difficult to compute in general, by the approximation

$$H(p_\tau) \approx \sum_{r \in R} \rho_r H_r(\tau_r) := H(\tau; \rho), \quad (7)$$

where $H_r(\tau_r) := - \sum_{x_r} \tau_r(x_r) \log \tau_r(x_r)$ is the entropy computed over region r , and $\{\rho_r: r \in R\}$ are weights assigned to the regions. Note that in the pairwise Ising case (2), with $p := p_\gamma$, we have the equality

$$H(p) = \sum_{s \in V} H_s(p_s) - \sum_{(s, t) \in E} I_{st}(p_{st})$$

when G is a tree, where $I_{st}(p_{st}) = H_s(p_s) + H_t(p_t) - H_{st}(p_{st})$ denotes the mutual information and p_s and p_{st} denote the node and edge marginals. Hence, the approximation (7) is exact with

$$\rho_{st} = 1, \quad \forall (s, t) \in E, \quad \text{and} \quad \rho_s = 1 - \deg(s), \quad \forall s \in V.$$

Using the approximation (7), we arrive at the following *reweighted Kikuchi approximation*:

$$B(\theta; \rho) := \sup_{\tau \in \Delta_R^K} \underbrace{\{\langle \theta, \tau \rangle + H(\tau; \rho)\}}_{B_{\theta, \rho}(\tau)}. \quad (8)$$

Note that when $\{\rho_r\}$ are the *overcounting numbers* $\{c_r\}$, defined recursively by

$$c_r = 1 - \sum_{s \in \mathcal{A}(r)} c_s, \quad (9)$$

the expression (8) reduces to the usual (unweighted) Kikuchi approximation considered in Pakzad and Anantharam [14].

3 Main results and consequences

In this section, we analyze the concavity of the Kikuchi variational problem (8). We derive a sufficient condition under which the function $B_{\theta,\rho}(\tau)$ is concave over the set Δ_R^K , so global optima of the reweighted Kikuchi approximation may be found efficiently. In the Bethe case, we also show that the condition is necessary for $B_{\theta,\rho}(\tau)$ to be concave over the entire region Δ_R^K , and we provide a geometric characterization of Δ_R^K in terms of the edge and cycle structure of the graph.

3.1 Sufficient conditions for concavity

We begin by establishing sufficient conditions for the concavity of $B_{\theta,\rho}(\tau)$. Clearly, this is equivalent to establishing conditions under which $H(\tau; \rho)$ is concave. Our main result is the following:

Theorem 1. *If $\rho \in \mathbb{R}^{|R|}$ satisfies*

$$\sum_{s \in \mathcal{F}(S)} \rho_s \geq 0, \quad \forall S \subseteq R, \quad (10)$$

then the Kikuchi entropy $H(\tau; \rho)$ is strictly concave on Δ_R^K .

The proof of Theorem 1 is contained in Appendix A.1, and makes use of a generalization of Hall's marriage lemma for weighted graphs (cf. Lemma 1 in Appendix A.2).

The condition (10) depends heavily on the structure of the region graph. For the sake of interpretability, we now specialize to the case where the region graph has only two layers, with the first layer corresponding to vertices and the second layer corresponding to hyperedges. In other words, for $r, s \in R$, we have $r \subseteq s$ only if $|r| = 1$, and $R = V \cup F$, where F is the set of hyperedges and V denotes the set of singleton vertices. This is the *Bethe case*, and the entropy

$$H(\tau; \rho) = \sum_{s \in V} \rho_s H_s(\tau_s) + \sum_{\alpha \in F} \rho_\alpha H_\alpha(\tau_\alpha) \quad (11)$$

is consequently known as the Bethe entropy.

The following result is proved in Appendix A.3:

Corollary 1. *Suppose $\rho_\alpha \geq 0$ for all $\alpha \in F$, and the following condition also holds:*

$$\sum_{s \in U} \rho_s + \sum_{\alpha \in F: \alpha \cap U \neq \emptyset} \rho_\alpha \geq 0, \quad \forall U \subseteq V. \quad (12)$$

Then the Bethe entropy $H(\tau; \rho)$ is strictly concave over Δ_R^K .

3.2 Necessary conditions for concavity

We now establish a converse to Corollary 1 in the Bethe case, showing that condition (12) is also necessary for the concavity of the Bethe entropy. When $\rho_\alpha = 1$ for $\alpha \in F$ and $\rho_s = 1 - |N(s)|$ for $s \in V$, we recover the result of Watanabe and Fukumizu [25] for the unweighted Bethe case. However, our proof technique is significantly simpler and avoids the complex machinery of graph zeta functions. Our approach proceeds by considering the Bethe entropy $H(\tau; \rho)$ on appropriate slices of the domain Δ_R^K so as to extract condition (12) for each $U \subseteq V$. The full proof is provided in Appendix B.1.

Theorem 2. *If the Bethe entropy $H(\tau; \rho)$ is concave over Δ_R^K , then $\rho_\alpha \geq 0$ for all $\alpha \in F$, and condition (12) holds.*

Indeed, as demonstrated in the simulations of Section 5, the Bethe objective function $B_{\theta,\rho}(\tau)$ may have multiple local optima if ρ does not satisfy condition (12).

3.3 Polytope of concavity

We now characterize the polytope defined by the inequalities (12). We show that in the pairwise Bethe case, the polytope may be expressed geometrically as the convex hull of single-cycle forests

formed by the edges of the graph. In the more general (non-pairwise) Bethe case, however, the polytope of concavity may strictly contain the latter set.

Note that the Bethe entropy (11) may be written in the alternative form

$$H(\tau; \rho) = \sum_{s \in V} \rho'_s H_s(\tau_s) - \sum_{\alpha \in F} \rho_\alpha \tilde{I}_\alpha(\tau_\alpha), \quad (13)$$

where $\tilde{I}_\alpha(\tau_\alpha) := \{\sum_{s \in \alpha} H_s(\tau_s)\} - H_\alpha(\tau_\alpha)$ is the KL divergence between the joint distribution τ_α and the product distribution $\prod_{s \in \alpha} \tau_s$, and the weights ρ'_s are defined appropriately.

We show that the polytope of concavity has a nice geometric characterization when $\rho'_s = 1$ for all $s \in V$, and $\rho_\alpha \in [0, 1]$ for all $\alpha \in F$. Note that this assignment produces the expression for the reweighted Bethe entropy analyzed in Wainwright et al. [22] (when all elements of F have cardinality two). Equation (13) then becomes

$$H(\tau; \rho) = \sum_{s \in V} \left(1 - \sum_{\alpha \in N(s)} \rho_\alpha\right) H_s(\tau_s) + \sum_{\alpha \in F} \rho_\alpha H_\alpha(\tau_\alpha), \quad (14)$$

and the inequalities (12) defining the polytope of concavity are

$$\sum_{\alpha \in F: \alpha \cap U \neq \emptyset} (|\alpha \cap U| - 1) \rho_\alpha \leq |U|, \quad \forall U \subseteq V. \quad (15)$$

Consequently, we define

$$\mathbb{C} := \left\{ \rho \in [0, 1]^{|F|} : \sum_{\alpha \in F: \alpha \cap U \neq \emptyset} (|\alpha \cap U| - 1) \rho_\alpha \leq |U|, \quad \forall U \subseteq V \right\}.$$

By Theorem 2, the set \mathbb{C} is the region of concavity for the Bethe entropy (14) within $[0, 1]^{|F|}$.

We also define the set

$$\mathbb{F} := \{1_{F'} : F' \subseteq F \text{ and } F' \cup N(F') \text{ is a single-cycle forest in } G\} \subseteq \{0, 1\}^{|F|},$$

where a *single-cycle forest* is defined to be a subset of edges of a graph such that each connected component contains at most one cycle. (We disregard the directions of edges in G .) The following theorem gives our main result. The proof is contained in Appendix C.1.

Theorem 3. *In the Bethe case (i.e., the region graph G has two layers), we have the containment $\text{conv}(\mathbb{F}) \subseteq \mathbb{C}$. If in addition $|\alpha| = 2$ for all $\alpha \in F$, then $\text{conv}(\mathbb{F}) = \mathbb{C}$.*

The significance of Theorem 3 is that it provides us with a convenient graph-based method for constructing vectors $\rho \in \mathbb{C}$. From the inequalities (15), it is not even clear how to efficiently verify whether a given $\rho \in [0, 1]^{|F|}$ lies in \mathbb{C} , since it involves testing $2^{|V|}$ inequalities.

Comparing Theorem 3 with known results, note that in the pairwise case ($|\alpha| = 2$ for all $\alpha \in F$), Theorem 1 of Wainwright et al. [22] states that the Bethe entropy is concave over $\text{conv}(\mathbb{T})$, where $\mathbb{T} \subseteq \{0, 1\}^{|E|}$ is the set of edge indicator vectors for spanning forests of the graph. It is trivial to check that $\mathbb{T} \subseteq \mathbb{F}$, since every spanning forest is also a single-cycle forest. Hence, Theorems 2 and 3 together imply a stronger result than in Wainwright et al. [22], characterizing the precise region of concavity for the Bethe entropy as a superset of the polytope $\text{conv}(\mathbb{T})$ analyzed there. In the unweighted Kikuchi case, it is also known [1, 14] that the Kikuchi entropy is concave for the assignment $\rho = 1_F$ when the region graph G is connected and has at most one cycle. Clearly, $1_F \in \mathbb{C}$ in that case, so this result is a consequence of Theorems 2 and 3, as well. However, our theorems show that a much more general statement is true.

It is tempting to posit that $\text{conv}(\mathbb{F}) = \mathbb{C}$ holds more generally in the Bethe case. However, as the following example shows, settings arise where $\text{conv}(\mathbb{F}) \subsetneq \mathbb{C}$. Details are contained in Appendix C.2.

Example 1. Consider a two-layer region graph with vertices $V = \{1, 2, 3, 4, 5\}$ and factors $\alpha_1 = \{1, 2, 3\}$, $\alpha_2 = \{2, 3, 4\}$, and $\alpha_3 = \{3, 4, 5\}$. Then $(1, \frac{1}{2}, 1) \in \mathbb{C} \setminus \text{conv}(\mathbb{F})$.

In fact, Example 1 is a special case of a more general statement, which we state in the following proposition. Here, $\mathfrak{F} := \{F' \subseteq F : 1_{F'} \in \mathbb{F}\}$, and an element $F^* \in \mathfrak{F}$ is *maximal* if it is not contained in another element of \mathfrak{F} .

Proposition 1. Suppose (i) G is not a single-cycle forest, and (ii) there exists a maximal element $F^* \in \mathfrak{F}$ such that the induced subgraph $F^* \cup N(F^*)$ is a forest. Then $\text{conv}(\mathbb{F}) \subsetneq \mathbb{C}$.

The proof of Proposition 1 is contained in Appendix C.3. Note that if $|\alpha| = 2$ for all $\alpha \in F$, then condition (ii) is violated whenever condition (i) holds, so Proposition 1 provides a partial converse to Theorem 3.

4 Reweighted sum product algorithm

In this section, we provide an iterative message passing algorithm to optimize the Kikuchi variational problem (8). As in the case of the generalized belief propagation algorithm for the unweighted Kikuchi approximation [28, 29, 11, 14, 12, 27] and the reweighted sum product algorithm for the Bethe approximation [22], our message passing algorithm searches for stationary points of the Lagrangian version of the problem (8). When ρ satisfies condition (10), Theorem 1 implies that the problem (8) is strictly concave, so the unique fixed point of the message passing algorithm globally maximizes the Kikuchi approximation.

Let $G = (V, R)$ be a region graph defining our Kikuchi approximation. Following Pakzad and Anantharam [14], for $r, s \in R$, we write $r \prec s$ if $r \subsetneq s$ and there does not exist $t \in R$ such that $r \subsetneq t \subsetneq s$. For $r \in R$, we define the parent set of r to be $\mathcal{P}(r) = \{s \in R: r \prec s\}$ and the child set of r to be $\mathcal{C}(r) = \{s \in R: s \prec r\}$. With this notation, $\tau = \{\tau_r(x_r): r \in R\}$ belongs to the set Δ_R^K if and only if $\sum_{x_{s \setminus r}} \tau_s(x_r, x_{s \setminus r}) = \tau_r(x_r)$ for all $r \in R, s \in \mathcal{P}(r)$.

The message passing algorithm we propose is as follows: For each $r \in R$ and $s \in \mathcal{P}(r)$, let $M_{sr}(x_r)$ denote the message passed from s to r at assignment x_r . Starting with an arbitrary positive initialization of the messages, we repeatedly perform the following updates for all $r \in R, s \in \mathcal{P}(r)$:

$$M_{sr}(x_r) \leftarrow C \left[\frac{\sum_{x_{s \setminus r}} \exp(\theta_s(x_s)/\rho_s) \prod_{v \in \mathcal{P}(s)} M_{vs}(x_s)^{\rho_v/\rho_s} \prod_{w \in \mathcal{C}(s) \setminus r} M_{sw}(x_w)^{-1}}{\exp(\theta_r(x_r)/\rho_r) \prod_{u \in \mathcal{P}(r) \setminus s} M_{ur}(x_r)^{\rho_u/\rho_r} \prod_{t \in \mathcal{C}(r)} M_{rt}(x_t)^{-1}} \right]^{\frac{\rho_r}{\rho_r + \rho_s}}. \quad (16)$$

Here, $C > 0$ may be chosen to ensure a convenient normalization condition; e.g., $\sum_{x_r} M_{sr}(x_r) = 1$. Upon convergence of the updates (16), we compute the pseudomarginals according to

$$\tau_r(x_r) \propto \exp\left(\frac{\theta_r(x_r)}{\rho_r}\right) \prod_{s \in \mathcal{P}(r)} M_{sr}(x_r)^{\rho_s/\rho_r} \prod_{t \in \mathcal{C}(r)} M_{rt}(x_t)^{-1}, \quad (17)$$

and we obtain the corresponding Kikuchi approximation by computing the objective function (8) with these pseudomarginals. We have the following result, which is proved in Appendix D:

Theorem 4. The pseudomarginals τ specified by the fixed points of the messages $\{M_{sr}(x_r)\}$ via the updates (16) and (17) correspond to the stationary points of the Lagrangian associated with the Kikuchi approximation problem (8).

As with the standard belief propagation and reweighted sum product algorithms, we have several options for implementing the above message passing algorithm in practice. For example, we may perform the updates (16) using serial or parallel schedules. To improve the convergence of the algorithm, we may damp the updates by taking a convex combination of new and previous messages using an appropriately chosen step size. As noted by Pakzad and Anantharam [14], we may also use a minimal graphical representation of the Hasse diagram to lower the complexity of the algorithm.

Finally, we remark that although our message passing algorithm proceeds in the same spirit as classical belief propagation algorithms by operating on the Lagrangian of the objective function, our algorithm as presented above does not immediately reduce to the generalized belief propagation algorithm for unweighted Kikuchi approximations or the reweighted sum product algorithm for tree-reweighted pairwise Bethe approximations. Previous authors use algebraic relations between the overcounting numbers (9) in the Kikuchi case [28, 29, 11, 14] and the two-layer structure of the Hasse diagram in the Bethe case [22] to obtain a simplified form of the updates. Since the coefficients ρ in our problem lack the same algebraic relations, following the message-passing protocol used in previous work [11, 28] leads to more complicated updates, so we present a slightly different algorithm that still optimizes the general reweighted Kikuchi objective.

5 Experiments

In this section, we present empirical results to demonstrate the advantages of the reweighted Kikuchi approximation that support our theoretical results. For simplicity, we focus on the binary pairwise Ising model given in equation (2). Without loss of generality, we may take the potentials to be $\gamma_s(x_s) = \gamma_s x_s$ and $\gamma_{st}(x_s, x_t) = \gamma_{st} x_s x_t$ for some $\gamma = (\gamma_s, \gamma_{st}) \in \mathbb{R}^{|V|+|E|}$. We run our experiments on two types of graphs: (1) K_n , the complete graph on n vertices, and (2) T_n , the $\sqrt{n} \times \sqrt{n}$ toroidal grid graph where every vertex has degree four.

Bethe approximation. We consider the pairwise Bethe approximation of the log partition function $A(\gamma)$ with weights $\rho_{st} \geq 0$ and $\rho_s = 1 - \sum_{t \in N(s)} \rho_{st}$. Because of the regularity structure of K_n and T_n , we take $\rho_{st} = \rho \geq 0$ for all $(s, t) \in E$ and study the behavior of the Bethe approximation as ρ varies. For this particular choice of weight vector $\vec{\rho} = \rho \mathbf{1}_E$, we define

$$\rho_{\text{tree}} = \max\{\rho \geq 0 : \vec{\rho} \in \text{conv}(\mathbb{T})\}, \quad \text{and} \quad \rho_{\text{cycle}} = \max\{\rho \geq 0 : \vec{\rho} \in \text{conv}(\mathbb{F})\}.$$

It is easily verified that for K_n , we have $\rho_{\text{tree}} = \frac{2}{n}$ and $\rho_{\text{cycle}} = \frac{2}{n-1}$; while for T_n , we have $\rho_{\text{tree}} = \frac{n-1}{2n}$ and $\rho_{\text{cycle}} = \frac{1}{2}$.

Our results in Section 3 imply that the Bethe objective function $B_{\gamma, \rho}(\tau)$ in equation (8) is concave if and only if $\rho \leq \rho_{\text{cycle}}$, and Wainwright et al. [22] show that we have the bound $A(\gamma) \leq B(\gamma; \rho)$ for $\rho \leq \rho_{\text{tree}}$. Moreover, since the Bethe entropy may be written in terms of the edge mutual information (13), the function $B(\gamma; \rho)$ is decreasing in ρ . In our results below, we observe that we may obtain a tighter approximation to $A(\gamma)$ by moving from the upper bound region $\rho \leq \rho_{\text{tree}}$ to the concavity region $\rho \leq \rho_{\text{cycle}}$. In addition, for $\rho > \rho_{\text{cycle}}$, we observe multiple local optima of $B_{\gamma, \rho}(\tau)$.

Procedure. We generate a random potential $\gamma = (\gamma_s, \gamma_{st}) \in \mathbb{R}^{|V|+|E|}$ for the Ising model (2) by sampling each potential $\{\gamma_s\}_{s \in V}$ and $\{\gamma_{st}\}_{(s,t) \in E}$ independently. We consider two types of models:

$$\text{Attractive: } \gamma_{st} \sim \text{Uniform}[0, \omega_{st}], \quad \text{and} \quad \text{Mixed: } \gamma_{st} \sim \text{Uniform}[-\omega_{st}, \omega_{st}].$$

In each case, $\gamma_s \sim \text{Uniform}[0, \omega_s]$. We set $\omega_s = 0.1$ and $\omega_{st} = 2$. Intuitively, the attractive model encourages variables in adjacent nodes to assume the same value, and it has been shown [18, 19] that the ordinary Bethe approximation ($\rho_{st} = 1$) in an attractive model lower-bounds the log partition function. For $\rho \in [0, 2]$, we compute stationary points of $B_{\gamma, \rho}(\tau)$ by running the reweighted sum product algorithm of Wainwright et al. [22]. We use a damping factor of $\lambda = 0.5$, convergence threshold of 10^{-10} for the average change of messages, and at most 2500 iterations. We repeat this process with at least 8 random initializations for each value of ρ . Figure 1 shows the scatter plots of ρ and the Bethe approximation $B_{\gamma, \rho}(\tau)$. In each plot, the two vertical lines are the boundaries $\rho = \rho_{\text{tree}}$ and $\rho = \rho_{\text{cycle}}$, and the horizontal line is the value of the true log partition function $A(\gamma)$.

Results. Figures 1(a)–1(d) show the results of our experiments on small graphs (K_5 and T_9) for both attractive and mixed models. We see that the Bethe approximation with $\rho \leq \rho_{\text{cycle}}$ generally provides a better approximation to $A(\gamma)$ than the Bethe approximation computed over $\rho \leq \rho_{\text{tree}}$. However, in general we cannot guarantee whether $B(\gamma; \rho)$ will give an upper or lower bound for $A(\gamma)$ when $\rho \leq \rho_{\text{cycle}}$. As noted above, we have $B(\gamma; 1) \leq A(\gamma)$ for attractive models.

We also observe from Figures 1(a)–1(d) that shortly after ρ leaves the concavity region $\{\rho \leq \rho_{\text{cycle}}\}$, multiple local optima emerge for the Bethe objective function. The presence of the point clouds near $\rho = 1$ in Figures 1(a) and 1(c) arises because the sum product algorithm has not converged after 2500 iterations. Indeed, the same phenomenon is true for all our results: in the region where multiple local optima begin to appear, it is more difficult for the algorithm to converge. See Figure 2 and the accompanying text in Appendix E for a plot of the points $(\rho, \log_{10}(\Delta))$, where Δ is the final average change in the messages at termination of the algorithm. From Figure 2, we see that the values of Δ are significantly higher for the values of ρ near where multiple local optima emerge. We suspect that for these values of ρ , the sum product algorithm fails to converge since distinct local optima are close together, so messages oscillate between the optima. For larger values of ρ , the local optima become sufficiently separated and the algorithm converges to one of them. However, it is interesting to note that this point cloud phenomenon does not appear for attractive models, despite the presence of distinct local optima.

Simulations for larger graphs are shown in Figures 1(e)–1(h). If we zoom into the region near $\rho \leq \rho_{\text{cycle}}$, we still observe the same behavior that $\rho \leq \rho_{\text{cycle}}$ generally provides a better Bethe

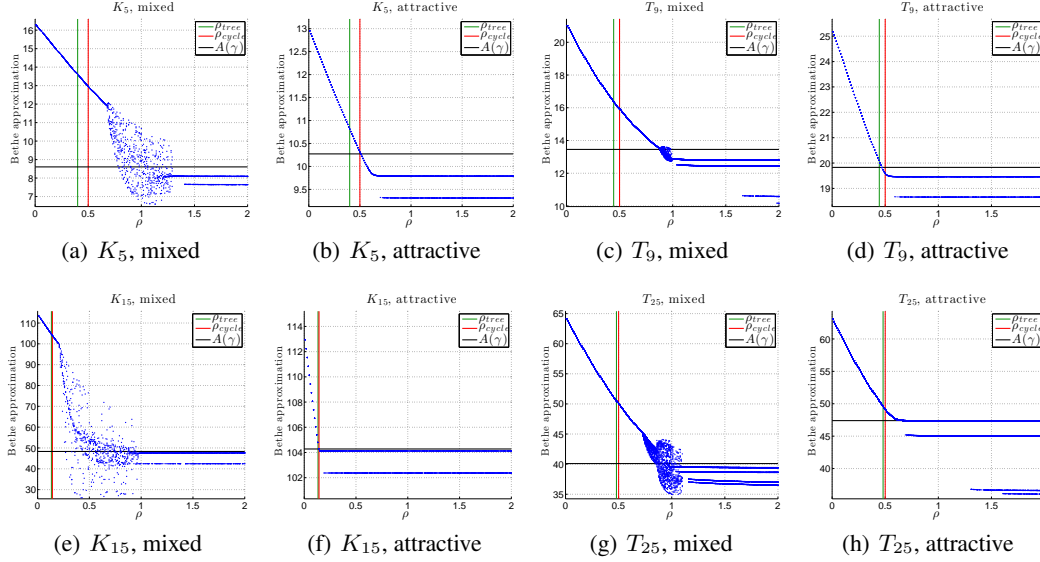


Figure 1: Values of the reweighted Bethe approximation as a function of ρ . See text for details.

approximation than $\rho \leq \rho_{\text{tree}}$. Moreover, the presence of the point clouds and multiple local optima are more pronounced, and we see from Figures 1(c), 1(g), and 1(h) that new local optima with even worse Bethe values arise for larger values of ρ . Finally, we note that the same qualitative behavior also occurs in all the other graphs that we have tried (K_n for $n \in \{5, 10, 15, 20, 25\}$ and T_n for $n \in \{9, 16, 25, 36, 49, 64\}$), with multiple random instances of the Ising model p_γ .

6 Discussion

In this paper, we have analyzed the reweighted Kikuchi approximation method for estimating the log partition function of a distribution that factorizes over a region graph. We have characterized necessary and sufficient conditions for the concavity of the variational objective function, generalizing existing results in literature. Our simulations demonstrate the advantages of using the reweighted Kikuchi approximation and show that multiple local optima may appear outside the region of concavity.

An interesting future research direction is to obtain a better understanding of the approximation guarantees of the reweighted Bethe and Kikuchi methods. In the Bethe case with attractive potentials θ , several recent results [22, 19, 18] establish that the Bethe approximation $B(\theta; \rho)$ is an upper bound to the log partition function $A(\theta)$ when ρ lies in the spanning tree polytope, whereas $B(\theta; \rho) \leq A(\theta)$ when $\rho = 1_F$. By continuity, we must have $B(\theta; \rho^*) = A(\theta)$ for some values of ρ^* , and it would be interesting to characterize such values where the reweighted Bethe approximation is exact.

Another interesting direction is to extend our theoretical results on properties of the reweighted Kikuchi approximation, which currently depend solely on the structure of the region graph and the weights ρ , to incorporate the effect of the model potentials θ . For example, several authors [20, 6] present conditions under which loopy belief propagation applied to the unweighted Bethe approximation has a unique fixed point. The conditions for uniqueness of fixed points slightly generalize the conditions for convexity, and they involve both the graph structure and the strength of the potentials. We suspect that similar results would hold for the reweighted Kikuchi approximation.

Acknowledgments. The authors thank Martin Wainwright for introducing the problem to them and providing helpful guidance. The authors also thank Varun Jog for discussions regarding the generalization of Hall’s lemma. The authors thank the anonymous reviewers for feedback that improved the clarity of the paper. PL was partly supported from a Hertz Foundation Fellowship and an NSF Graduate Research Fellowship while at Berkeley.

References

- [1] S. M. Aji and R. J. McEliece. The generalized distributive law and free energy minimization. In *Proceedings of the 39th Allerton Conference*, 2001.
- [2] F. Barahona. On the computational complexity of Ising spin glass models. *Journal of Physics A: Mathematical and General*, 15(10):3241, 1982.
- [3] H. A. Bethe. Statistical theory of superlattices. *Proceedings of the Royal Society of London. Series A, Mathematical and Physical Sciences*, 150(871):552–575, 1935.
- [4] P. Hall. On representatives of subsets. *Journal of the London Mathematical Society*, 10:26–30, 1935.
- [5] T. Heskes. Stable fixed points of loopy belief propagation are minima of the Bethe free energy. In *Advances in Neural Information Processing Systems 15*, 2002.
- [6] T. Heskes. On the uniqueness of loopy belief propagation fixed points. *Neural Computation*, 16(11):2379–2413, 2004.
- [7] T. Heskes. Convexity arguments for efficient minimization of the Bethe and Kikuchi free energies. *Journal of Artificial Intelligence Research*, 26:153–190, 2006.
- [8] A. T. Ihler, J. W. Fischer III, and A. S. Willsky. Loopy belief propagation: Convergence and effects of message errors. *Journal of Machine Learning Research*, 6:905–936, December 2005.
- [9] R. Kikuchi. A theory of cooperative phenomena. *Phys. Rev.*, 81:988–1003, March 1951.
- [10] B. Korte and J. Vygen. *Combinatorial Optimization: Theory and Algorithms*. Springer, 4th edition, 2007.
- [11] R. J. McEliece and M. Yildirim. Belief propagation on partially ordered sets. In *Mathematical Systems Theory in Biology, Communications, Computation, and Finance*, pages 275–300, 2002.
- [12] T. Meltzer, A. Globerson, and Y. Weiss. Convergent message passing algorithms: a unifying view. In *Proceedings of the Twenty-Fifth Conference on Uncertainty in Artificial Intelligence*, UAI ’09, 2009.
- [13] J. M. Mooij and H. J. Kappen. Sufficient conditions for convergence of the sum-product algorithm. *IEEE Transactions on Information Theory*, 53(12):4422–4437, December 2007.
- [14] P. Pakzad and V. Anantharam. Estimation and marginalization using Kikuchi approximation methods. *Neural Computation*, 17:1836–1873, 2003.
- [15] J. Pearl. *Probabilistic Reasoning in Intelligent Systems: Networks of Plausible Inference*. Morgan Kaufmann Publishers Inc., San Francisco, CA, USA, 1988.
- [16] T. Roosta, M. J. Wainwright, and S. S. Sastry. Convergence analysis of reweighted sum-product algorithms. *IEEE Transactions on Signal Processing*, 56(9):4293–4305, 2008.
- [17] D. Roth. On the hardness of approximate reasoning. *Artificial Intelligence*, 82(12):273 – 302, 1996.
- [18] N. Ruozzi. The Bethe partition function of log-supermodular graphical models. In *Advances in Neural Information Processing Systems 25*, 2012.
- [19] E. B. Sudderth, M. J. Wainwright, and A. S. Willsky. Loop series and Bethe variational bounds in attractive graphical models. In *Advances in Neural Information Processing Systems 20*, 2007.
- [20] S. C. Tatikonda and M. I. Jordan. Loopy belief propagation and Gibbs measures. In *Proceedings of the Eighteenth Conference on Uncertainty in Artificial Intelligence*, UAI ’02, 2002.
- [21] P. O. Vontobel. The Bethe permanent of a nonnegative matrix. *IEEE Transactions on Information Theory*, 59(3):1866–1901, 2013.
- [22] M. J. Wainwright, T. S. Jaakkola, and A. S. Willsky. A new class of upper bounds on the log partition function. *IEEE Transactions on Information Theory*, 51(7):2313–2335, 2005.
- [23] M. J. Wainwright and M. I. Jordan. Graphical models, exponential families, and variational inference. *Foundations and Trends in Machine Learning*, 1(1–2):1–305, January 2008.
- [24] Y. Watanabe and K. Fukumizu. Graph zeta function in the Bethe free energy and loopy belief propagation. In *Advances in Neural Information Processing Systems 22*, 2009.
- [25] Y. Watanabe and K. Fukumizu. Loopy belief propagation, Bethe free energy and graph zeta function. *arXiv preprint arXiv:1103.0605*, 2011.
- [26] Y. Weiss. Correctness of local probability propagation in graphical models with loops. *Neural Computation*, 12(1):1–41, 2000.
- [27] T. Werner. Primal view on belief propagation. In *UAI 2010: Proceedings of the Conference of Uncertainty in Artificial Intelligence*, pages 651–657, Corvallis, Oregon, July 2010. AUAI Press.
- [28] J. S. Yedidia, W. T. Freeman, and Y. Weiss. Generalized belief propagation. In *Advances in Neural Information Processing Systems 13*, 2000.
- [29] J. S. Yedidia, W. T. Freeman, and Y. Weiss. Constructing free energy approximations and generalized belief propagation algorithms. *IEEE Transactions on Information Theory*, 51:2282–2312, 2005.

A Proofs for Section 3.1

A.1 Proof of Theorem 1

We use the proof technique of Theorem 1 in Pakzad and Anantharam [14] for the unweighted Bethe entropy, together with Lemma 1 in Appendix A.2, which provides a generalization of Hall's marriage lemma for weighted bipartite graphs.

We construct a bipartite graph according to

$$V_1 := \{r \in R: \rho_r < 0\}, \quad \text{and} \quad V_2 := \{r \in R: \rho_r > 0\},$$

where $(s, t) \in E$ for $s \in V_1$ and $t \in V_2$ when $s \subset t$. Let weights w be defined such that $w(s) = -\rho_s$ for $s \in V_1$ and $w(s) = \rho_s$ for $s \in V_2$. We claim that condition (19) of Lemma 1 is satisfied. Indeed, for $U \subseteq V_1$, we have

$$w(U) = -\sum_{s \in U} \rho_s \leq \sum_{s \in A(U)} \rho_s = \sum_{s \in A(U): \rho_s > 0} \rho_s + \sum_{s \in A(U): \rho_s < 0} \rho_s \leq \sum_{s \in A(U): \rho_s > 0} \rho_s = w(N(U)),$$

where the first inequality is a direct application of the assumption (10). Hence, by Lemma 1, we have a saturating edge labeling γ .

For each $t \in V_2$, define

$$\rho'_t := \rho_t - \sum_{s \in N(t)} \gamma_{st} \geq 0.$$

We may then write

$$\begin{aligned} H(\tau; \rho) &= \sum_{s \in V_1} \rho_s H_s(\tau_s) + \sum_{t \in V_2} \rho_t H_t(\tau_t) \\ &= \sum_{(s,t) \in E} \gamma_{st} \{-H_s(\tau_s) + H_t(\tau_t)\} + \sum_{t \in V_2} \rho'_t H_t(\tau_t) \\ &= \sum_{(s,t) \in E} \gamma_{st} \left\{ \sum_{x_s} \tau_s(x_s) \log \tau_s(x_s) - \sum_{x_t} \tau_t(x_t) \log \tau_t(x_t) \right\} + \sum_{t \in V_2} \rho'_t H_t(\tau_t) \\ &= \sum_{(s,t) \in E} \gamma_{st} \sum_{x_t} \tau_t(x_t) \log \left(\frac{\tau_s(x_s)}{\tau_t(x_t)} \right) + \sum_{t \in V_2} \rho'_t H_t(\tau_t), \end{aligned} \tag{18}$$

where we have used the fact that $\sum_{x_t \in s} \tau_t(x_t) \log \tau_t(x_t) = \tau_s(x_s) \log \tau_s(x_s)$, since $\tau \in \Delta_R^K$, to obtain the last equality.

Note that for each pair (s, t) , we have

$$\sum_{x_t} \tau_t(x_t) \log \left(\frac{\tau_s(x_s)}{\tau_t(x_t)} \right) = -D_{\text{KL}}(\tau_t \| \tau_s),$$

which is strictly concave in the pair (τ_t, τ_s) . Furthermore, each term $H_t(\tau_t)$ is concave in τ_t . It follows by the expansion (18) that $H(\tau; \rho)$ is strictly concave, as wanted.

A.2 Generalization of Hall's marriage lemma

In this section, we prove a generalization of Hall's marriage lemma, which is useful in proving concavity of the Bethe entropy function $H(\tau; \rho)$.

Let $G = (V_1 \cup V_2, E)$ be a bipartite graph, where each vertex $v \in V := V_1 \cup V_2$ is assigned a weight $w(v) > 0$. For a set $U \subseteq V$, define

$$w(U) := \sum_{s \in U} w(s).$$

Also define the neighborhood set

$$N(U) := \bigcup_{s \in U} N(s),$$

where $N(s) := \{t : (s, t) \in E\}$ is the usual neighborhood set of a single node.

We say that an edge labeling $\gamma = (\gamma_{st} : (s, t) \in E) \in \mathbb{R}_{\geq 0}^{|E|}$ saturates V_1 if the following conditions hold:

1. For all $s \in V_1$, we have $\sum_{t \in N(s)} \gamma_{st} = w(s)$.
2. For all $t \in V_2$, we have $\sum_{s \in N(t)} \gamma_{st} \leq w(t)$.

Lemma 1. *Suppose*

$$w(U) \leq w(N(U)), \quad \forall U \subseteq V_1. \quad (19)$$

Then there exists an edge labeling γ that saturates V_1 .

Proof. We prove the lemma in stages. First, assume $w(v) \in \mathbb{Q}$ for all $v \in V$ and condition (19) holds. With an appropriate rescaling, we may assume that all weights are integers. Call the new weights w' . We then construct a graph G' such that each node $v \in V$ is expanded into a set U_v of $w'(v)$ nodes, and edges of G' are constructed by connecting all nodes in U_s to all nodes in U_t , for each $(s, t) \in E$. By the usual version of Hall's marriage lemma [4], there exists a matching of G' that saturates $V'_1 := \bigcup_{v \in V_1} U_v$. Indeed, it follows immediately from condition (19) that

$$w'(U) \leq w'(N(U)), \quad \forall U \subseteq V_1.$$

Suppose $T' \subseteq V'_1$, and let $T := \{s \in V_1 : U_s \cap T' \neq \emptyset\}$. Then

$$|T'| \leq \left| \bigcup_{s \in T} U_s \right| = w'(T) \leq w'(N(T)) = |N(T')|,$$

so the sufficient condition of Hall's marriage lemma is met, implying the existence of a matching. The edge labeling γ is obtained by setting

$$\gamma_{st} = \{\# \text{ of edges between } U_s \text{ and } U_t \text{ in matching}\}$$

and rescaling.

Next, suppose $w(v) \in \mathbb{R}$ for all $v \in V$ and condition (19) holds with *strict* inequality; i.e.,

$$w(U) < w(N(U)), \quad \forall U \subseteq V_1. \quad (20)$$

We claim that there exists an edge labeling γ that saturates V_1 . Indeed, let

$$\epsilon := \min_{U \subseteq V_1} \{w(N(U)) - w(U)\} > 0.$$

Define a new weighting w' with only rational values, such that

$$\begin{aligned} w'(s) &\in \left[w(s), w(s) + \frac{\epsilon}{2 \cdot \deg(G)} \right), \quad \forall s \in V_1, \\ w'(t) &\in \left(w(t) - \frac{\epsilon}{2 \cdot \deg(G)}, w(t) \right], \quad \forall t \in V_2, \end{aligned}$$

where $\deg(G) = |E|$ is the number of edges in G . It is clear that Hall's condition (19) still holds for w' . Hence, by the result of the last paragraph, there exists an edge labeling γ' that saturates V_1 with respect to w' . Observe that by decreasing the weights of γ' slightly, we easily obtain an edge labeling γ that saturates V_1 with respect to the original weighting w .

Finally, consider the most general case: condition (19) holds and $w(v) \in \mathbb{R}$ for all $v \in V$. Note that the problem of finding an edge labeling that saturates V_1 may be rephrased as follows. Let $b_1 \in \mathbb{R}^{|V_1|}$ be the vector of weights $(w(s) : s \in V_1)$. Then for an appropriate choice of the matrix $A_1 \in \{0, 1\}^{|V_1| \times |E|}$, the conditions

$$\sum_{t \in N(s)} \gamma_{st} = w(s), \quad \forall s \in V_1,$$

may be expressed as a system of linear equations,

$$A_1 \gamma = b_1. \quad (21)$$

Similarly, letting $b_2 = (w(t) : t \in V_2) \in \mathbb{R}^{|V_2|}$, the conditions

$$\sum_{s \in N(t)} \gamma_{st} \leq w(t), \quad \forall t \in V_2,$$

may be expressed in the form

$$A_2 \gamma \leq b_2, \tag{22}$$

where $A_2 \in \{0, 1\}^{|V_2| \times |E|}$. A saturating edge labeling exists if and only if there exists $\gamma \in \mathbb{R}_{\geq 0}^{|E|}$ that simultaneously satisfies conditions (21) and (22). Now consider a sequence of weight vectors $\{b_1^n\}_{n \geq 1}$, such that $b_1^n \rightarrow b_1$ and the convergence is from below and strictly monotone for each component. Let $w^n = (b_1^n, b_2)$ denote the full sequence of weights. Then

$$w^n(U) < w(U) \leq w(N(U)) = w^n(N(U)), \quad \forall U \subseteq V.$$

It follows by the result of the previous paragraph that there exists an edge labeling $\gamma^n \in \mathbb{R}_{\geq 0}^{|E|}$ such that

$$A_1 \gamma^n = b_1^n, \quad \text{and} \quad \gamma^n \in D := \left\{ \gamma \in \mathbb{R}_{\geq 0}^{|E|} : A_2 \gamma \leq b_2 \right\}.$$

Clearly, D is a closed set; furthermore, it is easy to see that the constraint $A_2 \gamma \leq b_2$ implies that each component of γ is bounded from above, since A_2 contains only nonnegative entries. It follows that the sequence $\{\gamma^n\}_{n \geq 1}$ has a limit point $\gamma^* \in D$. By continuity of the linear map A_1 , we must have

$$A_1 \gamma^* = \lim_{n \rightarrow \infty} A_1 \gamma^n = \lim_{n \rightarrow \infty} b_1^n = b_1.$$

Hence, γ^* is a valid edge labeling that saturates V_1 . □

A.3 Proof of Corollary 1

By Theorem 1, $H(\tau; \rho)$ is strictly concave provided condition (10) holds. Note that

$$\mathcal{F}(\alpha) = \{\alpha\}, \quad \forall \alpha \in F,$$

whereas

$$\mathcal{F}(s) = \{s\} \cup N(s), \quad \forall s \in V.$$

Condition (10) applied to the set $S = \{\alpha\}$ gives the inequality

$$\rho_\alpha \geq 0, \quad \forall \alpha \in F. \tag{23}$$

For a subset $U \subseteq V$, we can write

$$\mathcal{F}(U) = \bigcup_{s \in U} \mathcal{F}(s) = U \cup N(U) = U \cup \{\alpha \in F : \alpha \cap U \neq \emptyset\},$$

so (10) translates into

$$\sum_{s \in U} \rho_s + \sum_{\alpha \in F : \alpha \cap U \neq \emptyset} \rho_\alpha \geq 0, \quad \forall U \subseteq V, \tag{24}$$

which is condition (12). It is easy to see that conditions (23) and (24) together also imply the validity of condition (10) for any other set of regions $S \subseteq R$.

B Proofs for Section 3.2

B.1 Proof of Theorem 2

Our result relies on the property that if the Bethe entropy $H(\tau; \rho)$ is concave over Δ_R^K , then $H(\tau; \rho)$ is also concave over any subset $\Delta' \subseteq \Delta_R^K$. In particular, it is sufficient to assume that \mathcal{X} is binary, say $\mathcal{X} = \{-1, +1\}$; the general multinomial case $|\mathcal{X}| > 2$ follows by restricting the distribution of X_s to be supported on only two points.

The first lemma shows that $\rho_\alpha \geq 0$ for all $\alpha \in F$. The proof is contained in Appendix B.2.

Lemma 2. *If the Bethe entropy $H(\tau; \rho)$ is concave over Δ_R^K , then $\rho_\alpha \geq 0$ for all $\alpha \in F$.*

To establish the necessity of condition (12), consider a nonempty subset $U \subseteq V$ and the corresponding sub-region graph $R_U = U \cup F_U$, where $F_U = \{\alpha \cap U : \alpha \in F, \alpha \cap U \neq \emptyset\}$. From the original weights $\rho \in \mathbb{R}^{|V|+|F|}$, construct the sub-region weights $\rho^U \in \mathbb{R}^{|U|+|F_U|}$ given by

$$\rho_s^U = \rho_s, \quad \forall s \in U, \quad \text{and} \quad \rho_{\alpha \cap U}^U = \rho_\alpha, \quad \forall \alpha \cap U \in F_U.$$

For simplicity, we consider R_U to be a multiset by remembering which factor $\alpha \in F$ each $\beta = \alpha \cap U \in F_U$ comes from; we can equivalently work with R_U as a set by defining the weights ρ^U to be the sum over the pre-images of the factors in R_U . Consider the set of locally consistent R_U -pseudomarginals $\Delta_{R_U}^K$. Define a map that sends $\tilde{\tau} \in \Delta_{R_U}^K$ to $\tau \in \Delta_R^K$ defined by

$$\begin{aligned} \tau_s(x_s) &= \begin{cases} \tilde{\tau}_s(x_s) & \text{if } s \in U, \\ \frac{1}{2} & \text{otherwise,} \end{cases} \\ \tau_\alpha(x_\alpha) &= \begin{cases} \tilde{\tau}_{\alpha \cap U}(x_{\alpha \cap U}) \cdot \prod_{s \in \alpha \setminus U} \tau_s(x_s) & \text{if } \alpha \cap U \neq \emptyset \text{ (so } \alpha \cap U \in F_U), \\ \prod_{s \in \alpha} \tau_s(x_s) & \text{otherwise.} \end{cases} \end{aligned}$$

Let Δ_U denote the image of $\Delta_{R_U}^K$ under the mapping above, and note that $\Delta_U \subseteq \Delta_R^K$. Therefore, $H(\tau; \rho)$ is concave over Δ_U . Now let $\tau \in \Delta_U$ and let $\tilde{\tau} \in \Delta_{R_U}^K$ be a pre-image of τ . With this construction, we have the following lemma, proved in Appendix B.3:

Lemma 3. *The entropy $H(\tau; \rho)$ differs from $H_U(\tilde{\tau}; \rho^U)$ by a constant, where $H_U(\tilde{\tau}; \rho^U)$ is the Bethe entropy defined over the sub-region graph R_U .*

Finally, we have a lemma showing that we can extract condition (12) for $U = V$. The proof is provided in Appendix B.4.

Lemma 4. *If the Bethe entropy $H(\tau; \rho)$ is concave over Δ_R^K , then $\sum_{s \in V} \rho_s + \sum_{\alpha \in F} \rho_\alpha \geq 0$.*

By Lemma 3, the concavity of $H(\tau; \rho)$ over Δ_U implies the concavity of $H_U(\tilde{\tau}; \rho^U)$ over $\Delta_{R_U}^K$. Then by Lemma 4 applied to R_U , we have

$$\sum_{s \in U} \rho_s + \sum_{\alpha \in F : \alpha \cap U \neq \emptyset} \rho_\alpha = \sum_{s \in U} \rho_s^U + \sum_{\beta \in F_U} \rho_\beta^U \geq 0,$$

finishing the proof.

B.2 Proof of Lemma 2

Fix $\alpha \in F$, and let Δ_α be the set of pseudomarginals $\tau \in \Delta_R^K$ with the property that for all $s \in V$ and $\beta \in F \setminus \{\alpha\}$, τ_s and τ_β are uniform distributions over X_s and X_β , respectively, while τ_α is an arbitrary distribution on X_α with uniform single-node marginals. Then $H(\tau; \rho)$ is concave over Δ_α . On the other hand, note that for $\tau \in \Delta_\alpha$, $H_s(\tau_s) = \log 2$ and $H_\beta(\tau_\beta) = |\beta| \log 2$ are constants for $s \in V$ and $\beta \in F \setminus \{\alpha\}$, so we can write

$$H(\tau; \rho) = \rho_\alpha H_\alpha(\tau_\alpha) + \text{constant}.$$

Since $H_\alpha(\tau_\alpha)$ is concave in τ_α , this implies $\rho_\alpha \geq 0$, as claimed.

B.3 Proof of Lemma 3

By construction, for $s \in V \setminus U$, we have $H_s(\tau_s) = \log 2$; and for $\alpha \in F$ with $\alpha \cap U = \emptyset$, we have $H_\alpha(\tau_\alpha) = |\alpha| \log 2$. Moreover, for $\alpha \in F$ with $\alpha \cap U \neq \emptyset$, we have

$$H_\alpha(\tau_\alpha) = H_{\alpha \cap U}(\tilde{\tau}_{\alpha \cap U}) + \sum_{s \in \alpha \setminus U} H_s(\tau_s) = H_{\alpha \cap U}(\tilde{\tau}_{\alpha \cap U}) + |\alpha \setminus U| \log 2.$$

Therefore, for $\tau \in \Delta_U$, we can write

$$\begin{aligned}
H(\tau; \rho) &= \sum_{s \in V} \rho_s H_s(\tau_s) + \sum_{\alpha \in F} \rho_\alpha H_\alpha(\tau_\alpha) \\
&= \sum_{s \in U} \rho_s H_s(\tilde{\tau}_s) + \left(\sum_{s \in V \setminus U} \rho_s \right) \log 2 \\
&\quad + \sum_{\alpha \in F: \alpha \cap U \neq \emptyset} \rho_\alpha \left(H_{\alpha \cap U}(\tilde{\tau}_{\alpha \cap U}) + |\alpha \setminus U| \log 2 \right) + \sum_{\alpha \in F: \alpha \cap U = \emptyset} \rho_\alpha |\alpha| \log 2 \\
&= \sum_{s \in U} \rho_s^U H_s(\tilde{\tau}_s) + \sum_{\beta \in F_U} \rho_\beta^U H_\beta(\tilde{\tau}_\beta) + \text{constant} \\
&= H_U(\tilde{\tau}; \rho^U) + \text{constant},
\end{aligned}$$

as wanted.

B.4 Proof of Lemma 4

Given $m_o, m_e \in \mathbb{R}$, we define a pseudomarginal $\tau = (\tau_s, \tau_\alpha)$ by¹

$$\tau_s(x_s) = \frac{1 + x_s m_o}{2}, \quad \forall s \in V, x_s \in X = \{-1, +1\},$$

and for $\alpha \in F$ with $|\alpha| = k$,

$$\tau_\alpha(x_\alpha) = \begin{cases} 2^{-k} (1 + 2^{k-1} m_o + (2^{k-1} - 1) m_e) & \text{if } x_\alpha = (1, \dots, 1), \\ 2^{-k} (1 - 2^{k-1} m_o + (2^{k-1} - 1) m_e) & \text{if } x_\alpha = (-1, \dots, -1), \\ 2^{-k} (1 - m_e) & \text{otherwise.} \end{cases}$$

It is easy to see that $\sum_{x_s} \tau_s(x_s) = \sum_{x_\alpha} \tau_\alpha(x_\alpha) = 1$, and that τ_s is the single-node marginal of τ_α . Thus, for τ to lie in Δ_R^K , we only need to ensure that $\tau_s(x_s) \geq 0$ and $\tau_\alpha(x_\alpha) \geq 0$, or equivalently,

$$-1 \leq m_o \leq 1, \quad \frac{1 + 2^{k-1} |m_o|}{2^{k-1} - 1} \leq m_e \leq 1, \quad \forall 2 \leq k \leq K,$$

where $K = \max\{|\alpha| : \alpha \in F\}$. Let M denote the set of (m_o, m_e) satisfying the constraints above, and let Δ_M denote the set of pseudomarginals $\tau[m_o, m_e] \in \Delta_R^K$ given by the construction above for each $(m_o, m_e) \in M$.

Observe that the function $(m_o, m_e) \mapsto \tau[m_o, m_e]$ is additive for convex combinations; i.e., for $(m_o^{(1)}, m_e^{(1)}), \dots, (m_o^{(j)}, m_e^{(j)}) \in M$ and $\lambda_1, \dots, \lambda_j \geq 0$ with $\lambda_1 + \dots + \lambda_j = 1$, we have

$$\sum_{i=1}^j \lambda_i \tau[m_o^{(i)}, m_e^{(i)}] = \tau \left[\sum_{i=1}^j \lambda_i m_o^{(i)}, \sum_{i=1}^j \lambda_i m_e^{(i)} \right].$$

Since M is convex, this shows that Δ_M is a convex subset of Δ_R^K . Therefore, $H(\tau; \rho)$ is concave over Δ_M , and the additivity property above implies that the function

$$\zeta(m_o, m_e) := H(\tau[m_o, m_e]; \rho)$$

is concave over M . We now compute the Hessian of ζ and show how it relates to the required quantity that we want to prove is nonnegative.

Fix $(m_o, m_e) \in M$, and note that $\tau \equiv \tau[m_o, m_e]$ has the property that $\tau_\alpha = \tau_\beta$ whenever $|\alpha| = |\beta|$. Therefore, we can collect the terms in $H(\tau; \rho)$ based on the cardinality of $\alpha \in V \cup F$. The single-node entropy is, as a function of m_o ,

$$\zeta_1(m_o) := H_s(\tau_s) = -\eta \left(\frac{1 + m_o}{2} \right) - \eta \left(\frac{1 - m_o}{2} \right),$$

¹The definition of $\tau[m_o, m_e]$ above is equivalent to imposing the conditions

$$\mathbb{E}_{\tau_\alpha} \left[\prod_{s \in \beta} X_s \right] = m_o \text{ if } |\beta| \text{ is odd} \quad \text{and} \quad \mathbb{E}_{\tau_\alpha} \left[\prod_{s \in \beta} X_s \right] = m_e \text{ if } |\beta| \text{ is even,}$$

for all $\alpha \in V \cup F$ and $\emptyset \neq \beta \subseteq \alpha$.

where $\eta(t) := t \log t$. For $\alpha \in F$ with $|\alpha| = k \geq 2$, the entropy corresponding to τ_α is

$$\begin{aligned} \zeta_k(m_o, m_e) := H_\alpha(\tau_\alpha) = & -\eta\left(\frac{1 + 2^{k-1}m_o + (2^{k-1} - 1)m_e}{2^k}\right) - \eta\left(\frac{1 - 2^{k-1}m_o + (2^{k-1} - 1)m_e}{2^k}\right) \\ & - (2^k - 2)\eta\left(\frac{1 - m_e}{2^k}\right). \end{aligned}$$

The Bethe entropy can then be written as

$$\zeta(m_o, m_e) = H(\tau; \rho) = c_1 \zeta_1(m_o) + \sum_{k=2}^K c_k \zeta_k(m_o, m_e),$$

where $c_1 = \sum_{s \in V} \rho_s$ and $c_k = \sum_{\alpha \in F: |\alpha|=k} \rho_\alpha$ for $k \geq 2$.

Let us compute the Hessian matrix of $\zeta(m_o, m_e)$ along the axis $m_o = 0$. The function ζ_1 has second derivative $\zeta_1''(m_o) = -1/(1 - m_o^2)$, so at $m_o = 0$, the contribution of ζ_1 to the Hessian of ζ is

$$\nabla^2 \zeta_1(0, m_e) = \begin{pmatrix} -1 & 0 \\ 0 & 0 \end{pmatrix}.$$

For $k \geq 2$, the first partial derivatives of ζ_k are

$$\begin{aligned} \frac{\partial \zeta_k}{\partial m_o} &= -\frac{1}{2} \left\{ \log(1 + 2^{k-1}m_o + (2^{k-1} - 1)m_e) - \log(1 - 2^{k-1}m_o + (2^{k-1} - 1)m_e) \right\}, \\ \frac{\partial \zeta_k}{\partial m_e} &= -\frac{(2^{k-1} - 1)}{2^k} \left\{ \log(1 + 2^{k-1}m_o + (2^{k-1} - 1)m_e) + \log(1 - 2^{k-1}m_o + (2^{k-1} - 1)m_e) \right. \\ &\quad \left. - 2 \log(1 - m_e) \right\}. \end{aligned}$$

The Hessian $\nabla^2 \zeta_k$ at $m_o = 0$ is then given by

$$\nabla^2 \zeta_k(0, m_e) = \begin{pmatrix} -\frac{2^{k-1}}{1 + (2^{k-1} - 1)m_e} & 0 \\ 0 & -\frac{2^{k-1} - 1}{(1 + (2^{k-1} - 1)m_e)(1 - m_e)} \end{pmatrix}.$$

Therefore, the Hessian of ζ at $m_o = 0$ is the diagonal matrix

$$\nabla^2 \zeta(0, m_e) = \begin{pmatrix} -c_1 - \sum_{k=2}^K \frac{2^{k-1}c_k}{1 + (2^{k-1} - 1)m_e} & 0 \\ 0 & -\sum_{k=2}^K \frac{(2^{k-1} - 1)c_k}{(1 + (2^{k-1} - 1)m_e)(1 - m_e)} \end{pmatrix}.$$

In particular, the eigenvalues of $\nabla^2 \zeta(0, m_e)$ are its diagonal entries. Taking $m_e \rightarrow 1$, we see that the eigenvalue corresponding to the first diagonal entry satisfies

$$\lim_{m_e \rightarrow 1} \lambda_1(m_e) = \lim_{m_e \rightarrow 1} \left\{ -c_1 - \sum_{k=2}^K \frac{2^{k-1}c_k}{1 + (2^{k-1} - 1)m_e} \right\} = -\sum_{k=1}^K c_k.$$

Since $(0, m_e) \in M$ as $m_e \rightarrow 1$ and $\zeta(m_o, m_e)$ is concave over M , we see that the eigenvalue above is nonpositive, which implies

$$\sum_{s \in V} \rho_s + \sum_{\alpha \in F} \rho_\alpha = \sum_{k=1}^K c_k \geq 0,$$

as desired.

C Proofs for Section 3.3

C.1 Proof of Theorem 3

We first show that $\text{conv}(\mathbb{F}) \subseteq \mathbb{C}$ in the general Bethe case. Since \mathbb{C} is convex, it suffices to show that $\mathbb{F} \subseteq \mathbb{C}$, so consider $1_{F'} \in \mathbb{F}$. We need to show that inequality (15) holds for $\rho = 1_{F'}$.

Let W_1, \dots, W_m denote the connected components of $F' \cup N(F')$ in G . Consider an arbitrary $U \subseteq V$, and define $U_i := W_i \cap U$ for $1 \leq i \leq m$, and $U_0 := U \setminus \{U_1, \dots, U_m\}$. Then each W_i has at most one cycle. Furthermore, we may write

$$\sum_{\substack{\alpha \in F: \\ \alpha \cap U \neq \emptyset}} (|\alpha \cap U| - 1) \rho_\alpha = \sum_{\substack{\alpha \in F': \\ \alpha \cap U \neq \emptyset}} (|\alpha \cap U| - 1) = \sum_{i=1}^m \left\{ \sum_{\substack{\alpha \in W_i: \\ \alpha \cap U_i \neq \emptyset}} (|\alpha \cap U_i| - 1) \right\}. \quad (25)$$

We claim that

$$\sum_{\alpha \in W_i: \alpha \cap U_i \neq \emptyset} (|\alpha \cap U_i| - 1) \leq |U_i|, \quad \forall 1 \leq i \leq m. \quad (26)$$

Indeed, consider the induced subgraph W'_i of W_i with vertex set $V_i := U_i \cup \{\alpha \in W_i: \alpha \cap U_i \neq \emptyset\}$. Since W_i has at most one cycle, W'_i has at most one cycle, as well. Furthermore, the number of edges of W'_i is given by

$$|E(W'_i)| = \sum_{\alpha \in W_i: \alpha \cap U_i \neq \emptyset} |\alpha \cap U_i|,$$

and the number of vertices is $|V_i| = |U_i| + |\{\alpha \in W_i: \alpha \cap U_i \neq \emptyset\}|$.

We have the following simple lemma:

Lemma 5. *A connected graph G has at most one cycle if and only if*

$$|E(U)| \leq |U|, \quad \forall U \subseteq V.$$

Proof. First suppose G has at most one cycle. For any subset $U \subseteq V$, the induced subgraph H clearly also contains at most one cycle. Hence, we may remove at most one edge to obtain a graph H' which is a forest. Then

$$|E(H')| \leq |V(H')| - 1 = |U| - 1. \quad (27)$$

Furthermore, $|E(U)| \leq |E(H')| + 1$. It follows that $|E(U)| \leq |U|$.

Conversely, if G is a connected graph with more than one cycle, we may pick U to be the union of vertices in the two cycles, along with a path connecting the two cycles (in case the cycles are disconnected). It is easy to check that condition (27) is violated in this case. \square

Applying Lemma 5 to the graph W'_i and rearranging then yields inequality (26). Combining with equation (25) then yields

$$\sum_{\alpha \in F: \alpha \cap U \neq \emptyset} (|\alpha \cap U| - 1) \rho_\alpha \leq \sum_{i=1}^m |U_i| = |U| - |U_0| \leq |U|,$$

proving the condition (15).

We now specialize to the case where $|\alpha| = 2$ for all $\alpha \in F$. Note that in this case, we may identify the region graph with an ordinary graph $\bar{G} = (V, E)$, where the edge set E is given by F . It is easy to check that $1_{F'} \in \mathbb{F}$ if and only if the subgraph of \bar{G} with edge set F' is a single-cycle forest. In the following argument, we abuse notation and refer to \bar{G} as G .

Recall that a *rational polyhedron* is a set of the form $\{x \in \mathbb{R}^p: Ax \leq b\}$, such that A and b have rational entries. Clearly, \mathbb{C} is a rational polyhedron. Furthermore, a polyhedron is *integral* if all vertices are elements of the integer lattice \mathbb{Z}^p . The following result is standard in integer programming:

Lemma 6. [Theorem 5.12, [10]] *Let P be a rational polyhedron. Then P is integral if and only if $\max\{c^T x: x \in P\}$ is attained by an integral vector for each c for which the maximum is finite.*

We have already established that $1_{F'} \in \mathbb{C}$ for all $1_{F'} \in \mathbb{F}$. Furthermore, any lattice point in \mathbb{C} is of the form 1_H , where $H \subseteq E$. By Lemma 5, each connected component of H must contain at most one cycle, implying that H is a single-cycle forest. Hence, $1_H \in \mathbb{F}$, as well. We then combine Lemma 6 with the following proposition to obtain the desired result.

Proposition 2. *Let $G = (V, E)$ be a graph. For any set of weights $c = (c_{st}) \in \mathbb{R}^{|E|}$, the LP*

$$\max \sum_{(s,t) \in E} c_{st} x_{st} \quad (28)$$

$$\begin{aligned} \text{s.t. } \sum_{(s,t) \in E(U)} x_{st} &\leq |U|, \quad \forall U \subseteq V, \\ 0 &\leq x_{st} \leq 1, \quad \forall (s,t) \in E, \end{aligned} \quad (29)$$

attains its maximum value at an integral vector x^ .*

Proof. We first argue that it suffices to consider rational weights $c \in \mathbb{Q}^{|E|}$. Let X denote the feasible set of the LP, and let $F(c) = \max_{x \in X} c^\top x$ denote the maximum value of the LP. Note that $F(c)$ is continuous in c .

Suppose the claim in the proposition holds for $c \in \mathbb{Q}^{|E|}$. Given $c \in \mathbb{R}^{|E|}$, let $x^* \in \arg \max_{x \in X} c^\top x$. Let $(c^{(n)})_{n \geq 1}$ be a sequence of weights in $\mathbb{Q}^{|E|}$ converging to c elementwise as $n \rightarrow \infty$. Given $\epsilon > 0$, choose n sufficiently large such that $\|c^{(n)} - c\|_1 < \epsilon$ and $|F(c) - F(c^{(n)})| < \epsilon$. Applying our hypothesis, we know there exists an integral vector $z^* \in X$ such that $F(c^{(n)}) = (c^{(n)})^\top z^*$. Then

$$|F(c) - c^\top z^*| \leq |F(c) - F(c^{(n)})| + |(c^{(n)} - c)^\top z^*| \leq \epsilon + \|c^{(n)} - c\|_1 \|z^*\|_\infty \leq 2\epsilon.$$

Thus, we can find an integral vector $z^* \in X$ that achieves the objective function that is within 2ϵ from the optimal value. Since $\epsilon > 0$ is arbitrary, we conclude by continuity that we may find an integral vector in X arbitrarily close to x^* . This implies that x^* is an integral vector.

It now remains to prove the claim in the proposition for $c \in \mathbb{Q}^{|E|}$. If $c_{st} < 0$ for some $(s, t) \in E$, then any optimal solution x^* will have $x_{st}^* = 0$. If $c_{st} = 0$, then we can set $x_{st}^* = 0$ without changing the objective value. Thus, we can assume $c_{st} > 0$ for all $(s, t) \in E$. By scaling the weights, we can further assume that $c_{st} \in \{1, \dots, K\}$ for all $(s, t) \in E$, for some $K \in \mathbb{N}$.

We first upper-bound the objective function. For $1 \leq i \leq K$, let $E_i = \{(s, t) \in E : c_{st} \geq i\}$ denote the set of edges with weights at least i , and let V_i denote the set of vertices in E_i . By construction, we have

$$V = V_1 \supset \dots \supset V_K, \quad \text{and} \quad E = E_1 \supset \dots \supset E_K.$$

Suppose the subgraph $G_i = (V_i, E_i)$ is decomposed into connected components

$$G_i = T_{i1} \cup \dots \cup T_{i\alpha_i} \cup H_{i1} \cup \dots \cup H_{i\beta_i}, \quad (30)$$

where each $T_{ij} = (V(T_{ij}), E(T_{ij}))$ is a tree and each $H_{i\ell} = (V(H_{i\ell}), E(H_{i\ell}))$ is a connected graph with at least one loop. Thus, we have the disjoint partitions

$$V_i = \bigcup_{j=1}^{\alpha_i} V(T_{ij}) \cup \bigcup_{\ell=1}^{\beta_i} V(H_{i\ell}), \quad \text{and} \quad E_i = \bigcup_{j=1}^{\alpha_i} E(T_{ij}) \cup \bigcup_{\ell=1}^{\beta_i} E(H_{i\ell}).$$

Then we can write the objective function of the LP as

$$\sum_{(s,t) \in E} c_{st} x_{st} = \sum_{i=1}^K \sum_{(s,t) \in E_i} x_{st} = \sum_{i=1}^K \left(\sum_{j=1}^{\alpha_i} \sum_{(s,t) \in E(T_{ij})} x_{st} + \sum_{\ell=1}^{\beta_i} \sum_{(s,t) \in E(H_{i\ell})} x_{st} \right). \quad (31)$$

For $i = 1, \dots, K$ and $j = 1, \dots, \alpha_i$, since T_{ij} is a tree, we have

$$\sum_{(s,t) \in E(T_{ij})} x_{st} \leq |E(T_{ij})| = |V(T_{ij})| - 1, \quad \forall x \in X. \quad (32)$$

For $\ell = 1, \dots, \beta_i$, note that the set $E(H_{i\ell})$ of edges in $H_{i\ell}$ is contained within the set $E(V(H_{i\ell}))$ of edges in the subgraph of G induced by $V(H_{i\ell})$. Thus, by inequality (29), we have

$$\sum_{(s,t) \in E(H_{i\ell})} x_{st} \leq \sum_{(s,t) \in E(V(H_{i\ell}))} x_{st} \leq |V(H_{i\ell})|. \quad (33)$$

Plugging in the bounds (32) and (33) to inequality (31), we arrive at the upper bound

$$\sum_{(s,t) \in E} c_{st} x_{st} \leq \sum_{i=1}^K \left(\sum_{j=1}^{\alpha_i} \{|V(T_{ij})| - 1\} + \sum_{\ell=1}^{\beta_i} |V(H_{i\ell})| \right) = \sum_{i=1}^K (|V_i| - \alpha_i). \quad (34)$$

We now prove the claim in the proposition by explicitly constructing an integral vector x^* that achieves the upper bound (34). Since $x^* \in \{0, 1\}^{|E|}$, it is the indicator vector of a subset $E^* \subseteq E$.

Our approach is to construct, for each $1 \leq i \leq K$, a spanning single-cycle forest $F_i = (V_i, C_i)$ of $G_i = (V_i, E_i)$ with the following properties:

1. The restriction of F_i to $V_{i+1} \subseteq V_i$ is equal to $F_{i+1} = (V_{i+1}, C_{i+1})$, or equivalently, $C_i \cap E_{i+1} = C_{i+1}$. By induction, this implies $C_1 \cap E_i = C_i$, for $1 \leq i \leq K$.
2. For $1 \leq i \leq K$, we have $|C_i| = |V_i| - \alpha_i$.

Suppose we can construct such F_i 's. Setting $E^* = C_1$, we see that this construction yields a vector $x^* = 1_{E^*}$ satisfying

$$\begin{aligned} \sum_{(s,t) \in E} c_{st} x_{st}^* &= \sum_{i=1}^K \sum_{(s,t) \in E_i} x_{st}^* = \sum_{i=1}^K \sum_{(s,t) \in E_i} \mathbf{1}\{(s,t) \in C_1\} \\ &= \sum_{i=1}^K |C_1 \cap E_i| = \sum_{i=1}^K |C_i| = \sum_{i=1}^K (|V_i| - \alpha_i), \end{aligned}$$

so x^* achieves the bound (34), as desired.

It now remains to construct the F_i 's. We start by taking F_K to be a spanning single-cycle forest of G_K . Specifically, for each connected component H of G_K , we do the following: If H is a tree, we take H to be in F_K . If H contains at least one loop, then we take an arbitrary spanning single-cycle subgraph (i.e., a spanning tree with an additional edge to form one cycle) of H to be in F_K . Then $F_K = (V_K, C_K)$ satisfies $|C_K| = |V_K| - \alpha_K$, since there are α_K trees among the connected components of G_K .

Suppose that for some $1 \leq i \leq K - 1$, we have constructed a spanning single-cycle forest F_{i+1} satisfying the desired properties. Now consider $G_i = (V_i, E_i)$, and construct $F_i = (V_i, C_i)$ as follows: Consider each connected component of G_i in the decomposition (30).

- (a) For each tree $T_{ij} = (V(T_{ij}), E(T_{ij}))$, for all $1 \leq j \leq \alpha_i$, take T_{ij} to be in F_i . This component of F_i is clearly consistent with F_{i+1} , and the contribution to the total number of edges $|C_i|$ is

$$\sum_{j=1}^{\alpha_i} |E(T_{ij})| = \sum_{j=1}^{\alpha_i} (|V(T_{ij})| - 1) = \sum_{j=1}^{\alpha_i} |V(T_{ij})| - \alpha_i.$$

- (b) Consider $H_{i\ell} = (V(H_{i\ell}), E(H_{i\ell}))$, for some $1 \leq \ell \leq \beta_i$, so $H_{i\ell}$ has at least one loop. There may be several connected components of F_{i+1} in $H_{i\ell}$; suppose there are $\gamma_{i\ell}$ trees and $\delta_{i\ell}$ single-cycle graphs from F_{i+1} in $H_{i\ell}$. From each of the $\delta_{i\ell}$ single-cycle graphs, remove one edge to reduce it to a tree, and complete the $\gamma_{i\ell} + \delta_{i\ell}$ trees into a spanning tree of $H_{i\ell}$. Add the $\delta_{i\ell}$ edges back, so the spanning tree now has $\delta_{i\ell}$ cycles. Remove $\delta_{i\ell} - 1$ edges to break this graph into $\delta_{i\ell}$ connected components, such that each of the original $\delta_{i\ell}$ single-cycle graphs is in a separate connected components, and the last connected component is a tree. Set this new graph to be in F_i . It is clear by construction that this component of F_i is consistent with F_{i+1} since we keep all the edges from F_{i+1} . Moreover, its contribution to the total number of edges C_i is precisely

$$\sum_{\ell=1}^{\beta_i} (|V(H_{i\ell})| - 1) + \delta_{i\ell} - \{\delta_{i\ell} - 1\} = \sum_{\ell=1}^{\beta_i} |V(H_{i\ell})|.$$

Combining the two cases above, for each $1 \leq i \leq K$ we have constructed a spanning single-cycle forest F_i that is consistent with F_{i+1} and satisfies $|C_i| = \sum_{j=1}^{\alpha_i} |V(T_{ij})| - \alpha_i + \sum_{\ell=1}^{\beta_i} |V(H_{i\ell})| = |V_i| - \alpha_i$, as desired. This completes the proof of the proposition. \square

C.2 Details for Example 1

It is easy to check that $\mathbb{F} = \{0, 1\}^3 \setminus (1, 1, 1)$. Hence, $(1, \frac{1}{2}, 1) \notin \text{conv}(\mathbb{F})$. By enumerating the inequalities defining the boundary of \mathbb{C} for different values of $U \subseteq V$, one may check that the only inequalities that are not trivially satisfied by $\rho \in [0, 1]^3$ are

$$\begin{aligned} \rho_1 + 2\rho_2 + \rho_3 &\leq 3, \\ 2\rho_1 + 2\rho_2 + \rho_3 &\leq 4, \\ \rho_1 + 2\rho_2 + 2\rho_3 &\leq 4, \\ 2\rho_1 + 2\rho_2 + 2\rho_3 &\leq 5. \end{aligned}$$

The first inequality together with the condition $\rho \in [0, 1]^3$ implies the remaining three inequalities, so

$$\mathbb{C} = \{\rho \in [0, 1]^3 : \rho_1 + 2\rho_2 + \rho_3 \leq 3\}.$$

Clearly, $(1, \frac{1}{2}, 1) \in \mathbb{C}$.

C.3 Proof of Proposition 1

The first condition implies $F \notin \mathfrak{F}$. In particular, $F^* \neq F$ and we can find $\alpha^* \in F \setminus F^*$. Since F^* is maximal, $\tilde{F} = F^* \cup \{\alpha^*\} \notin \mathfrak{F}$. This means $1_{F^*} \in \mathbb{F}$ but $1_{\tilde{F}} = 1_{F^*} + 1_{\{\alpha^*\}} \notin \mathbb{F}$. Define

$$\rho = 1_{F^*} + \epsilon 1_{\{\alpha^*\}}, \quad \text{with} \quad \epsilon = \frac{1}{|\alpha^*| - 1} \in (0, 1).$$

We claim that $\rho \in \mathbb{C}$, which will give us the desired conclusion since $\rho \notin \text{conv}(\mathbb{F})$.

To show $\rho \in \mathbb{C}$, since we already know that $1_{F^*} \in \mathbb{F} \subseteq \mathbb{C}$, we only need to verify inequality (15) for $U \subseteq V$ with $U \cap \alpha^* \neq \emptyset$. Given such a subset U , note that since $F^* \cup N(F^*)$ is a forest, the subgraph induced by the nodes $U \cup \{\alpha \in F^* : \alpha \cap U \neq \emptyset\}$ is also a forest, so

$$\sum_{\alpha \in F^* : \alpha \cap U \neq \emptyset} (|\alpha \cap U| - 1) \leq |U| - 1.$$

Therefore,

$$\sum_{\alpha \in F : \alpha \cap U \neq \emptyset} (|\alpha \cap U| - 1) \rho_\alpha = \sum_{\alpha \in F^* : \alpha \cap U \neq \emptyset} (|\alpha \cap U| - 1) + \frac{|\alpha^* \cap U| - 1}{|\alpha^*| - 1} \leq |U| - 1 + 1 = |U|,$$

verifying condition (15), as desired.

D Proof of Theorem 4

For $r \in R$ and $s \in \mathcal{P}(r)$, let $\lambda_{sr}(x_r)$ be a Lagrange multiplier associated with the consistency constraint $\sum_{x_{s \setminus r}} \tau_s(x_r, x_{s \setminus r}) = \tau_r(x_r)$. We enforce the nonnegativity constraint $\tau_r(x_r) \geq 0$ and normalization constraint $\sum_{x_r} \tau_r(x_r) = 1$ explicitly. Then the Lagrangian associated with the optimization problem (8) is

$$\begin{aligned} \mathcal{L}_{\theta, \rho}(\tau; \lambda) &= \sum_{r \in R} \sum_{x_r} \tau_r(x_r) \theta_r(x_r) - \sum_{r \in R} \rho_r \sum_{x_r} \tau_r(x_r) \log \tau_r(x_r) \\ &\quad + \sum_{r \in R} \sum_{t \in \mathcal{C}(r)} \sum_{x_t} \lambda_{rt}(x_t) \left(\tau_t(x_t) - \sum_{x_{r \setminus t}} \tau_r(x_t, x_{r \setminus t}) \right). \end{aligned} \quad (35)$$

Setting the partial derivatives of $\mathcal{L}_{\theta,\rho}$ with respect to the Lagrange multipliers equal to zero recovers the consistency constraints. Taking the derivative of $\mathcal{L}_{\theta,\rho}$ with respect to $\tau_r(x_r)$ and setting it equal to zero yields

$$\log \tau_r(x_r) = C + \frac{\theta_r(x_r)}{\rho_r} + \sum_{s \in \mathcal{P}(r)} \frac{\lambda_{sr}(x_r)}{\rho_r} - \sum_{t \in \mathcal{C}(r)} \frac{\lambda_{rt}(x_t)}{\rho_r},$$

where C is a constant that enforces the normalization condition $\sum_{x_r} \tau_r(x_r) = 1$. Defining the messages by

$$\log M_{sr}(x_r) = \frac{\lambda_{sr}(x_r)}{\rho_s},$$

we can write the equation above as

$$\tau_r(x_r) \propto \exp\left(\frac{\theta_r(x_r)}{\rho_r}\right) \frac{\prod_{s \in \mathcal{P}(r)} M_{sr}(x_r)^{\rho_s/\rho_r}}{\prod_{t \in \mathcal{C}(r)} M_{rt}(x_t)},$$

recovering equation (17).

For $s \in R$ and $r \in \mathcal{C}(s)$, enforcing the consistency condition $\sum_{x_{s \setminus r}} \tau_s(x_r, x_{s \setminus r}) = \tau_r(x_r)$ gives us

$$\begin{aligned} \exp\left(\frac{\theta_r(x_r)}{\rho_r}\right) \frac{M_{sr}(x_r)^{\rho_s/\rho_r} \prod_{u \in \mathcal{P}(r) \setminus s} M_{ur}(x_r)^{\rho_u/\rho_r}}{\prod_{t \in \mathcal{C}(r)} M_{rt}(x_t)} \\ \propto \sum_{x_{s \setminus r}} \exp\left(\frac{\theta_s(x_s)}{\rho_s}\right) \frac{\prod_{v \in \mathcal{P}(s)} M_{vs}(x_s)^{\rho_v/\rho_s}}{M_{sr}(x_r) \prod_{w \in \mathcal{C}(s) \setminus r} M_{sw}(x_w)}. \end{aligned}$$

Rearranging the equation to collect $M_{sr}(x_r)$ on the left hand side and taking the $(1 + \rho_s/\rho_r)$ -th root on both sides gives us the update equation (16).

From the derivation above, it is clear that if $\{M_{sr}(x_r)\}$ is a fixed point of the update equation (16), then the collection τ of pseudomarginals defined by (17) is a stationary point of the Lagrangian (35), since it sets the derivatives of $\mathcal{L}_{\theta,\rho}$ equal to zero.

E Additional Simulation Results

In this section, we provide additional plots to better illustrate the observations that we make in Section 5. For convenience, Figures 2(a)–2(d) and Figures 2(i)–2(l) show the same plots as in Figure 1. Figures 2(e)–2(h) show the plots of $(\rho, \log_{10}(\Delta))$ for the Ising models in Figures 2(a)–2(d), and similarly for Figures 2(m)–2(p). Here, Δ is the final average change of the messages in the sum product algorithm at termination; i.e., either when $\Delta \leq 10^{-10}$ or after 2500 iterations of the algorithm with parallel updates.

For $\rho \leq \rho_{\text{cycle}}$, in which the Bethe variational problem (8) is concave, there is a unique optimal value for the Bethe approximation. The values of Δ in this region are slightly higher than the convergence threshold, which means sum product has not converged after 2500 iterations, but the final value of Δ is sufficiently small that the messages have stabilized.

Shortly after ρ becomes larger than ρ_{cycle} , the curve of the Bethe values splits into multiple lines, which indicates that the Bethe objective function has multiple local optima. These lines are evidently distinct local optima since the values of Δ are at the convergence threshold, which means sum product converges and yields stationary points of the Lagrangian.

In the models with mixed potentials, we observe that for the values of ρ where the multiple local optima begin to emerge, the values of Δ are significantly higher and sum product does not converge. This behavior is reflected in the presence of the point cloud in the plots of the Bethe values. As noted in Section 5, we suspect that this behavior arises because distinct local optima are initially close together, so messages oscillate between them. For larger values of ρ , however, the local optima are sufficiently separated, so sum product converges and there are multiple lines in the graphs of the Bethe values.

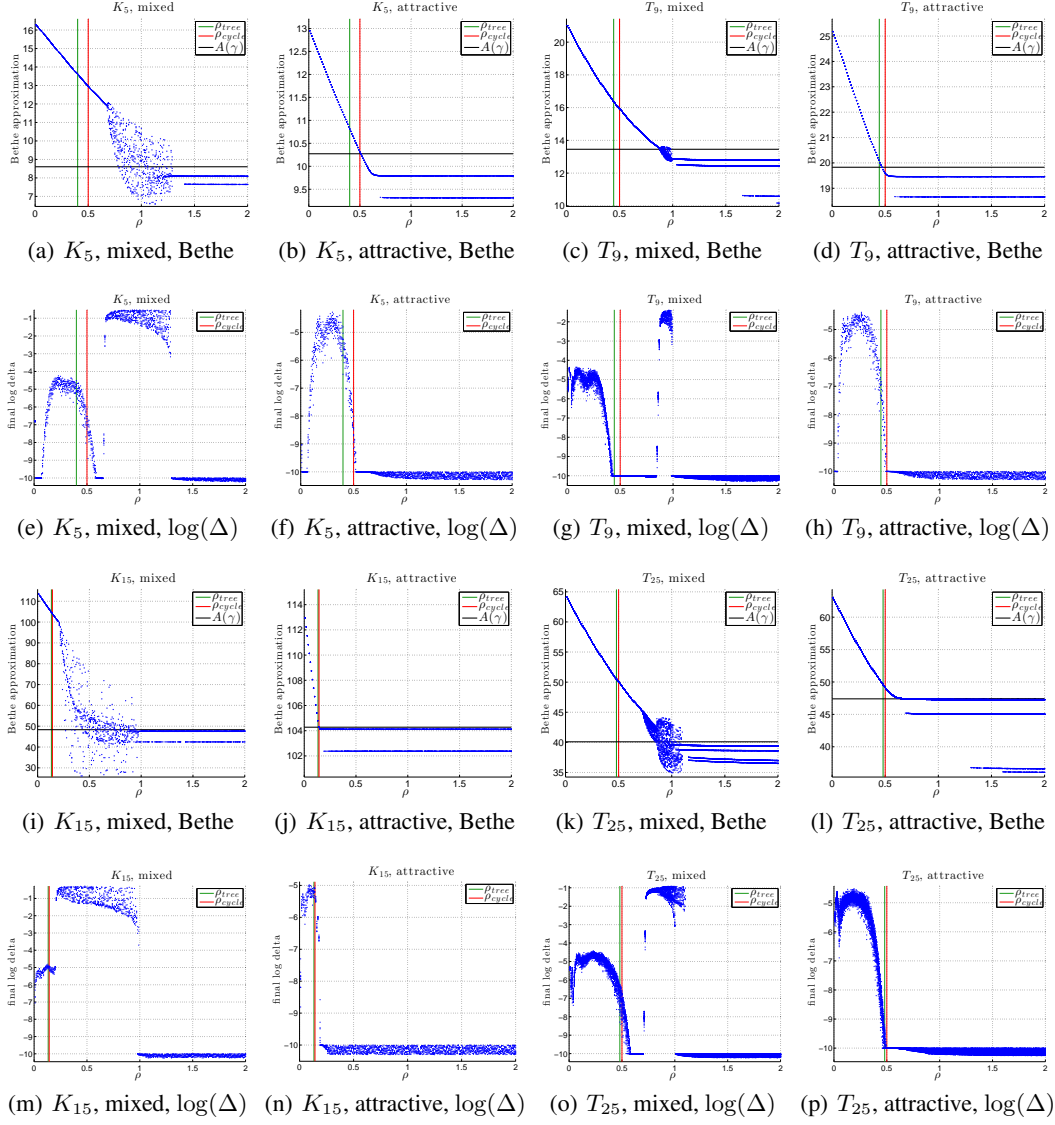


Figure 2: Values of the reweighted Bethe approximation and the final $\log_{10}(\Delta)$ as a function of ρ .