

---

# The Benefits of Learning with Strongly Convex Approximate Inference

## \*\*\* Supplemental \*\*\*

---

**Ben London**

University of Maryland, College Park, MD 20742 USA

BLONDON@CS.UMD.EDU

**Bert Huang**

Virginia Tech, Blacksburg, VA 24061 USA

BHUANG@VT.EDU

**Lise Getoor**

University of California, Santa Cruz, CA 95064 USA

GETOOR@SOE.UCSC.EDU

### A. Properties of Strong Convexity

Strong convexity can be characterized in a number of ways. The following facts provide some conditions that are equivalent to Definition 1.

**Fact 1.** *A differentiable function,  $\varphi : \mathcal{S} \rightarrow \mathbb{R}$ , of a convex set,  $\mathcal{S}$ , is  $\kappa$ -strongly convex w.r.t. a norm,  $\|\cdot\|$ , if and only if, for all  $s, s' \in \mathcal{S}$ ,*

$$\kappa \|s - s'\|^2 \leq \langle \nabla \varphi(s) - \nabla \varphi(s'), s - s' \rangle.$$

**Fact 2.** *A twice-differentiable function,  $\varphi : \mathcal{S} \rightarrow \mathbb{R}$ , of a convex set,  $\mathcal{S}$ , is  $\kappa$ -strongly convex w.r.t. a norm,  $\|\cdot\|$ , if and only if, for all  $s, s' \in \mathcal{S}$ ,*

$$\kappa \|s\|^2 \leq \langle s, \nabla^2 \varphi(s') s \rangle.$$

For the 2-norm, Fact 2 means that the minimum eigenvalue of the Hessian is lower-bounded by  $\kappa$ .

### B. Proofs from Section 3

This section contains all deferred proofs from Section 3.

#### B.1. Proof of Stability Lemma (Lemma 2)

Recall that  $\tilde{\mu}(\theta)$  and  $\tilde{\mu}(\theta')$  are the gradients of  $\tilde{\Phi}(\theta)$  and  $\tilde{\Phi}(\theta')$ , respectively. Since the conjugate function,  $\tilde{\Phi}^*$ , is assumed to be  $\kappa$ -strongly convex, we have via Lemma 1 and Definition 2 that

$$\begin{aligned} \|\tilde{\mu}(\theta) - \tilde{\mu}(\theta')\|_2 &= \left\| \nabla \tilde{\Phi}(\theta) - \nabla \tilde{\Phi}(\theta') \right\|_2 \\ &\leq \frac{1}{\kappa} \|\theta - \theta'\|_2. \end{aligned} \quad (14)$$

Dividing both sides by  $\sqrt{|G|}$  completes the proof.

#### B.2. The Marginals of the Expected NLL Minimizer are the True Marginals

Observe that  $\hat{\theta}_m$  effectively fits the empirical marginals of the dataset,  $\frac{1}{m} \sum_{j=1}^m \hat{\mathbf{y}}^{(j)}$ . Thus, as  $m \rightarrow \infty$ , the marginals induced by  $\hat{\theta}_m$  and  $\theta^*$  converge. This is formalized in the following lemma.

**Lemma 4.** *Let  $\mu(\theta^*)$  denote the true marginals of a distribution. Let  $\hat{\theta}$  denote the minimizer of the expected NLL, per Eq. 5. Then,*

$$\mu(\theta^*) = \tilde{\mu}(\hat{\theta}).$$

**Proof** Expanding the expected NLL, we have

$$\mathbb{E}[-\ln \tilde{p}(\mathbf{Y}; \theta)] = \tilde{\Phi}(\theta) - \mathbb{E}[\theta \cdot \hat{\mathbf{Y}}].$$

The gradient of this is

$$\nabla \mathbb{E}[-\ln \tilde{p}(\mathbf{Y}; \theta)] = \tilde{\mu}(\theta) - \mathbb{E}[\hat{\mathbf{Y}}] = \tilde{\mu}(\theta) - \mu(\theta^*).$$

Since the NLL is differentiable, the gradient at the minimum is zero. Thus, when  $\nabla \mathbb{E}[-\ln \tilde{p}(\mathbf{Y}; \theta)] = 0$ , we have  $\tilde{\mu}(\hat{\theta}) = \mu(\theta^*)$ . ■

#### B.3. Proof of Error Bound (Proposition 1)

By Lemma 4,  $\mu(\theta^*) = \tilde{\mu}(\hat{\theta})$ . Further, because  $\tilde{\Phi}^*$  is assumed to be  $\kappa$ -strongly convex, using Lemma 2, we have that

$$\begin{aligned} \frac{1}{\sqrt{|G|}} \left\| \tilde{\mu}(\hat{\theta}_m) - \mu(\theta^*) \right\|_2 &= \frac{1}{\sqrt{|G|}} \left\| \tilde{\mu}(\hat{\theta}_m) - \tilde{\mu}(\bar{\theta}) \right\|_2 \\ &\leq \frac{1}{\kappa \sqrt{|G|}} \left\| \hat{\theta}_m - \bar{\theta} \right\|_2. \end{aligned} \quad (15)$$

The rest of the proof involves upper-bounding  $\left\| \hat{\theta}_m - \bar{\theta} \right\|_2$ .

Assumption 1 states that, with probability at least  $1 - \delta$ , there exists a convex set,  $\mathcal{S}$ , encompassing  $\bar{\theta}$  and  $\hat{\theta}_m$ , such that the minimum eigenvalue of  $\nabla^2 \mathcal{L}(\cdot; \theta) : \theta \in \mathcal{S}$  is lower-bounded by  $\gamma(\delta, m, G)$ . By Fact 2, this event implies that the NLL is  $\gamma(\delta, m, G)$ -strongly convex in  $\mathcal{S}$ . Since  $\nabla^2 \mathcal{L}(\cdot; \theta) = \nabla^2 \mathcal{L}_m(\theta)$ , the same can be said for  $\mathcal{L}_m$ , so the regularized NLL,

$$\mathcal{L}_m^R(\theta) \triangleq \mathcal{L}_m(\theta) + \Lambda_m \|\theta\|_2^2,$$

is also  $\gamma(\delta, m, G)$ -strongly convex in  $\mathcal{S}$ . Therefore, with probability at least  $1 - \delta$  over draws of  $m$  examples,

$$\begin{aligned} \|\bar{\theta} - \hat{\theta}_m\|_2^2 &\leq \frac{\langle \nabla \mathcal{L}_m^R(\bar{\theta}) - \nabla \mathcal{L}_m^R(\hat{\theta}_m), \bar{\theta} - \hat{\theta}_m \rangle}{\gamma(\delta, m, G)} \\ &= \frac{\langle \nabla \mathcal{L}_m^R(\bar{\theta}), \bar{\theta} - \hat{\theta}_m \rangle}{\gamma(\delta, m, G)} \\ &\leq \frac{\|\nabla \mathcal{L}_m^R(\bar{\theta})\|_2 \|\bar{\theta} - \hat{\theta}_m\|_2}{\gamma(\delta, m, G)}. \end{aligned}$$

The second line follows from the fact that  $\hat{\theta}_m$  is the minimizer of  $\mathcal{L}_m^R$ , which is differentiable, so  $\nabla \mathcal{L}_m^R(\hat{\theta}_m) = \mathbf{0}$ . The last line uses Cauchy-Schwarz. Dividing both sides by  $\|\bar{\theta} - \hat{\theta}_m\|_2$ , and combining with Eq. 15, we have that, with probability at least  $1 - \delta$ ,

$$\frac{1}{\sqrt{|G|}} \|\tilde{\mu}(\hat{\theta}_m) - \mu(\theta^*)\|_2 \leq \frac{\|\nabla \mathcal{L}_m^R(\bar{\theta})\|_2}{\kappa \gamma(\delta, m, G) \sqrt{|G|}}. \quad (16)$$

Using the triangle inequality, the norm of the gradient decomposes as

$$\begin{aligned} \|\nabla \mathcal{L}_m^R(\bar{\theta})\|_2 &= \|\nabla \mathcal{L}_m(\bar{\theta}) + 2\Lambda_m \bar{\theta}\|_2 \\ &\leq \|\nabla \mathcal{L}_m(\bar{\theta})\|_2 + 2\Lambda_m \|\bar{\theta}\|_2. \end{aligned} \quad (17)$$

Let  $N = |\bar{\theta}|$ , and note that  $N = \ell |\mathcal{V}| + \ell^2 |\mathcal{E}| \leq \ell^2 |G|$ . Therefore, using the definition of  $\Lambda_m$ , and leveraging the assumption that  $\|\bar{\theta}\|_\infty \leq 1$ , we have that

$$2\Lambda_m \|\bar{\theta}\|_2 \leq 2\sqrt{\frac{N}{m}} \|\bar{\theta}\|_\infty \leq 2\ell \sqrt{\frac{|G|}{m}}. \quad (18)$$

Turning now to the gradient of  $\mathcal{L}_m$ , we can expand Eq. 4 as

$$\mathcal{L}_m(\theta) = \frac{1}{m} \sum_{j=1}^m \tilde{\Phi}(\theta) - \theta \cdot \hat{y}^{(j)}.$$

Since  $\tilde{\mu}(\bar{\theta})$  is the gradient of  $\tilde{\Phi}(\theta)$ , and is in fact equal to the true marginals,  $\mu(\theta^*)$ , we have that the gradient of  $\mathcal{L}_m$  is

$$\begin{aligned} \nabla \mathcal{L}_m(\bar{\theta}) &= \frac{1}{m} \sum_{j=1}^m \tilde{\mu}(\bar{\theta}) - \hat{y}^{(j)} \\ &= \mu(\theta^*) - \frac{1}{m} \sum_{j=1}^m \hat{y}^{(j)}. \end{aligned}$$

Note that the gradient is a zero-mean random vector; random because it depends on the draw of the training set. We will bound this quantity with high probability, using a technique borrowed from London et al. (2014).

It helps to denote the gradient by a vector,  $\nabla \mathcal{L}_m(\bar{\theta}) \triangleq \mathbf{g} \in \mathbb{R}^N$ . Fix some value  $\epsilon > 0$ . For  $\mathbf{g}$  to be greater than  $\epsilon$ , at least one of its coordinates must have magnitude at least  $\epsilon/\sqrt{N}$ ; otherwise, we would have

$$\|\mathbf{g}\|_2 = \sqrt{\sum_{i=1}^N |g_i|^2} < \sqrt{\sum_{i=1}^N \frac{\epsilon^2}{N}} = \epsilon.$$

Thus, using the union bound, we have that

$$\begin{aligned} \Pr \{\|\mathbf{g}\|_2 \geq \epsilon\} &\leq \Pr \left\{ \exists i : |g_i| \geq \frac{\epsilon}{\sqrt{N}} \right\} \\ &\leq \sum_{i=1}^N \Pr \left\{ |g_i| \geq \frac{\epsilon}{\sqrt{N}} \right\}. \end{aligned}$$

Each  $g_i$  is the difference of the mean and sample average of a sufficient statistic for some node variable  $Y_v$  (or edge variable  $Y_e$ ) having label  $y_v$  (or  $y_e$ ). The sufficient statistics are bounded in the interval  $[0, 1]$ , so  $|g_i| \leq 1$ . Moreover, the sample average is taken from  $m$  i.i.d. draws from the target distribution. Therefore, applying Hoeffding's inequality to each  $i$ , we have that

$$\Pr \left\{ |g_i| \geq \frac{\epsilon}{\sqrt{N}} \right\} \leq 2 \exp \left( -\frac{2m\epsilon^2}{N} \right).$$

Summing over  $i = 1, \dots, N$ , we have

$$\Pr \{\|\mathbf{g}\|_2 \geq \epsilon\} \leq 2N \exp \left( -\frac{2m\epsilon^2}{N} \right).$$

Thus, with probability at least  $1 - \delta$ ,

$$\|\nabla \mathcal{L}_m(\bar{\theta})\|_2 \leq \sqrt{\frac{N \ln \frac{2N}{\delta}}{2m}} \leq \ell \sqrt{\frac{|G| \ln \frac{2\ell^2 |G|}{\delta}}{2m}}. \quad (19)$$

The last inequality uses the fact that  $N \leq \ell^2 |G|$ .

Substituting Eqs. 18 and 19 into Eq. 17, and rearranging the terms, we have that with probability at least  $1 - \delta$ ,

$$\|\nabla \mathcal{L}_m^R(\bar{\theta})\|_2 \leq \ell \sqrt{\frac{|G|}{m}} \left( \sqrt{\frac{1}{2} \ln \frac{2\ell^2 |G|}{\delta}} + 2 \right).$$

Then, combining the above with Eq. 16, we have that with probability at least  $1 - 2\delta$  over draws of the training set,

$$\frac{1}{\sqrt{|G|}} \|\tilde{\mu}(\hat{\theta}_m) - \mu(\theta^*)\|_2 \leq \frac{\ell \left( \sqrt{\frac{1}{2} \ln \frac{2\ell^2 |G|}{\delta}} + 2 \right)}{\kappa \gamma(\delta, m, G) \sqrt{m}},$$

which completes the proof.

## C. Tree-Structured Models

In this section, we analyze tree-structured models. We show that the negative entropy of a tree-structured model is strongly convex, with a modulus that depends on the contraction coefficients induced by the model. This result is used in the proof of Proposition 2. We also show how the contraction coefficients of a tree-structured model can be measured efficiently.

### C.1. Strong Convexity of the Tree Negative Entropy

When the model is structured according to a tree,  $T$ , the marginal polytope,  $\mathcal{M}$ , is exactly equivalent to the local marginal polytope,  $\tilde{\mathcal{M}}$ . Further, its entropy function,  $H_T$ , can be expressed succinctly as a function of the marginals, using the Bethe entropy formula (see Eq. 8). Wainwright (2006) showed that  $-H_T$  is  $\Omega(1/|G|)$ -strongly convex. This is a pessimistic lower bound, since it considers all models in the exponential family. Indeed, we can show that tree-structured models with good contraction (see Definition 3) and bounded degree induce a negative entropy that is  $\Omega(1)$ -strongly convex.

**Proposition 4.** *Fix a tree,  $T$ , with maximum degree  $\Delta_T = O(1)$ , independent of  $|\mathcal{V}|$ . Let  $\Theta \subseteq \mathbb{R}^{|\theta|}$  denote the set of potentials with maximum contraction coefficient  $\vartheta_{\theta}^* \leq 1/\Delta_T$ , and let  $\mathcal{M}(\Theta) \triangleq \{\mu(\theta) : \theta \in \Theta\}$  denote the corresponding set of realizable marginals. Then, the negative entropy,  $-H_T$ , is  $\Omega(1)$ -strongly convex in  $\mathcal{M}(\Theta)$ .*

**Proof** The Hessian of the log-partition,  $\Phi(\theta)$ , is the covariance matrix,

$$\Sigma(\mathbf{Y}; \theta) \triangleq \mathbb{E}[\hat{y}\hat{y}^\top; \theta] - \mathbb{E}[\hat{y}; \theta] \mathbb{E}[\hat{y}^\top; \theta],$$

where  $\mathbb{E}[\cdot; \theta]$  denotes an expectation over the distribution parameterized by  $\theta$ . (For a derivation of this fact, see Wainwright & Jordan (2008).) Let  $\Sigma^{-1}(\mathbf{Y}; \theta)$  denote the inverse covariance (i.e., precision) matrix. Since  $\Phi$  is the convex conjugate of the negative entropy,  $-H$ , the Hessian of one is the inverse Hessian of the other. This insight yields the following lemma.

**Lemma 5.** *The negative entropy,  $-H$ , is  $(1/\lambda_{\max})$ -strongly convex in  $\mathcal{M}(\Theta)$ , where  $\lambda_{\max} \triangleq \max_{\theta \in \Theta} \|\Sigma(\mathbf{Y}; \theta)\|_2$  is the maximum eigenvalue of the covariance matrix, over all potentials in  $\Theta$ .*

**Proof** Via Fact 2,  $-H$  is  $\kappa$ -strongly convex in  $\mathcal{M}(\Theta)$  if the eigenvalues of  $\nabla^2(-H(\mu(\theta)))$ , for every  $\mu(\theta) \in \mathcal{M}(\Theta)$  (i.e., every  $\theta \in \Theta$ ), are bounded away from zero by  $\kappa$ . Via convex conjugacy,

$$\nabla^2(-H(\mu(\theta))) = (\nabla^2\Phi(\theta))^{-1} = \Sigma^{-1}(\mathbf{Y}; \theta).$$

Therefore, the minimum eigenvalue of  $-H(\mu(\theta))$  is equal to the maximum eigenvalue of  $\Sigma(\mathbf{Y}; \theta)$ . ■

Thus, to lower-bound the convexity of  $-H$ , it suffices to uniformly upper-bound the spectral norm of  $\Sigma(\mathbf{Y}; \theta)$ , over all  $\theta \in \Theta$ . A simple way to do this (used by Wainwright, 2006) is to analyze the trace norm (i.e., sum of the diagonal), which upper-bounds the spectral norm. The diagonal elements of the covariance matrix are uniformly upper-bounded by  $1/4$ , since the sufficient statistics are bounded in  $[0, 1]$ . This yields a (loose) upper bound of  $O(|G|)$ . For our purposes, this bound is too loose, since it grows with the size of the graph.

A better approach is to analyze the induced 1-norm (i.e., maximum column sum) or  $\infty$ -norm (i.e., maximum row sum), which, for symmetric matrices, are equivalent, and conveniently upper-bound the spectral norm. (This is because  $\|\mathbf{A}\|_2 \leq \sqrt{\|\mathbf{A}\|_1 \|\mathbf{A}\|_\infty} = \sqrt{\|\mathbf{A}\|_1 \|\mathbf{A}\|_1} = \|\mathbf{A}\|_1$ .) Intuitively, the 1-norm of the covariance matrix captures the maximum dependence as a function of graph distance. To bound the 1-norm, we will relate each covariance coefficient to a product of contraction coefficients. For contraction less than 1—i.e., without determinism—this product will decrease geometrically with graph distance. This geometric series converges, provided the structure has bounded degree and sufficiently small contraction.

Our proof uses a technical lemma that is often credited to Dobrushin. We use a version of this given by Kontorovich (2012).

**Lemma 6** (Kontorovich, 2012, Lemma 2.1). *Let  $\nu : \Omega \rightarrow \mathbb{R}$  be a signed, balanced measure, such that  $\sum_{\omega \in \Omega} \nu(\omega) = 0$ . Let  $K : \Omega \times \Omega \rightarrow \mathbb{R}$  be a Markov kernel, where  $K(\omega | \omega') \geq 0$ ,  $\sum_{\omega} K(\omega | \omega') = 1$ , and*

$$(K\nu)(\omega) \triangleq \sum_{\omega' \in \Omega} K(\omega | \omega') \nu(\omega').$$

Then

$$\begin{aligned} \|K\nu\|_{\text{TV}} &= \sum_{\omega} \left| \sum_{\omega'} K(\omega | \omega') \nu(\omega') \right| \\ &\leq \vartheta \sum_{\omega'} |\nu(\omega')| \\ &= \vartheta \|\nu\|_{\text{TV}}, \end{aligned}$$

where

$$\vartheta \triangleq \sup_{\omega, \omega' \in \Omega} \|K(\cdot | \omega) - K(\cdot | \omega')\|_{\text{TV}}.$$

is the contraction coefficient of  $K$ .

Fix any  $\theta \in \Theta$ . For the following, we use the shorthand  $p_{\theta}(y)$  to denote  $p(Y = y; \theta)$ , and similar probabilities. We also let  $\sigma_{\theta}(y_u, y_v)$  denote the entry of the covariance matrix corresponding to  $Y_u = y_u$  and  $Y_v = y_v$ .

Let  $\pi(1), \dots, \pi(l)$  denote the sequence of nodes along a path. Note that  $\pi$  is the unique path connecting its end points, since the model is tree-structured. The covariance entries corresponding to  $Y_{\pi(1)} = y_{\pi(1)}$  and  $Y_{\pi(l)} = y_{\pi(l)}$  can be written recursively as

$$\begin{aligned}
 & \sigma_{\theta}(y_{\pi(1)}, y_{\pi(l)}) \\
 &= p_{\theta}(y_{\pi(1)}, y_{\pi(l)}) - p_{\theta}(y_{\pi(1)})p_{\theta}(y_{\pi(l)}) \\
 &= \sum_{y_{\pi(l-1)}} p_{\theta}(y_{\pi(1)}, y_{\pi(l-1)}, y_{\pi(l)}) \\
 &\quad - p_{\theta}(y_{\pi(1)})p_{\theta}(y_{\pi(l-1)}, y_{\pi(l)}) \\
 &= \sum_{y_{\pi(l-1)}} p_{\theta}(y_{\pi(1)}, y_{\pi(l-1)})p_{\theta}(y_{\pi(l)} | y_{\pi(l-1)}) \\
 &\quad - p_{\theta}(y_{\pi(1)})p_{\theta}(y_{\pi(l-1)})p_{\theta}(y_{\pi(l)} | y_{\pi(l-1)}) \\
 &= \sum_{y_{\pi(l-1)}} p_{\theta}(y_{\pi(l)} | y_{\pi(l-1)}) \\
 &\quad \times (p_{\theta}(y_{\pi(1)}, y_{\pi(l-1)}) - p_{\theta}(y_{\pi(1)})p_{\theta}(y_{\pi(l-1)})) \\
 &= \sum_{y_{\pi(l-1)}} p_{\theta}(y_{\pi(l)} | y_{\pi(l-1)}) \sigma_{\theta}(y_{\pi(1)}, y_{\pi(l-1)}).
 \end{aligned}$$

Note that the second equality follows from the Markov property; since  $Y_{\pi(l)}$  is conditionally independent of  $Y_{\pi(1)}$  given  $Y_{\pi(l-1)}$ , we have that  $p_{\theta}(y_{\pi(l)} | y_{\pi(l-1)}, y_{\pi(1)}) = p_{\theta}(y_{\pi(l)} | y_{\pi(l-1)})$ .

In the righthand expression, the conditional probability under  $p_{\theta}$  defines a Markov kernel. Moreover, the covariance with  $y_{\pi(1)}$  defines a signed measure,

$$\nu(Y; y_{\pi(1)}) \triangleq \sigma_{\theta}(y_{\pi(1)}, Y),$$

which is balanced, since

$$\begin{aligned}
 \sum_y \nu(y; y_{\pi(1)}) &= \sum_y \sigma_{\theta}(y_{\pi(1)}, y) \\
 &= \sum_y p_{\theta}(y_{\pi(1)}, y) - p_{\theta}(y_{\pi(1)})p_{\theta}(y) \\
 &= p_{\theta}(y_{\pi(1)}) - p_{\theta}(y_{\pi(1)}) = 0.
 \end{aligned}$$

Therefore, via Lemma 6, we have that

$$\begin{aligned}
 & \sum_{y_{\pi(l)}} |\sigma_{\theta}(y_{\pi(1)}, y_{\pi(l)})| \\
 &= \sum_{y_{\pi(l)}} \left| \sum_{y_{\pi(l-1)}} p_{\theta}(y_{\pi(l)} | y_{\pi(l-1)}) \sigma_{\theta}(y_{\pi(1)}, y_{\pi(l-1)}) \right| \\
 &\leq \vartheta_{\theta}^* \sum_{y_{\pi(l-1)}} |\sigma_{\theta}(y_{\pi(1)}, y_{\pi(l-1)})|
 \end{aligned}$$

Applying this identity recursively, we have that

$$\begin{aligned}
 & \sum_{y_{\pi(l)}} |\sigma_{\theta}(y_{\pi(1)}, y_{\pi(l)})| \\
 &\leq \vartheta_{\theta}^* \sum_{y_{\pi(l-1)}} |\sigma_{\theta}(y_{\pi(1)}, y_{\pi(l-1)})| \\
 &\vdots \\
 &\leq (\vartheta_{\theta}^*)^{l-2} \sum_{y_{\pi(2)}} |\sigma_{\theta}(y_{\pi(1)}, y_{\pi(2)})| \\
 &\leq (\vartheta_{\theta}^*)^{l-1} \sum_{y'_{\pi(1)}} |\sigma_{\theta}(y_{\pi(1)}, y'_{\pi(1)})| \\
 &\leq \frac{\ell}{4} (\vartheta_{\theta}^*)^{l-1}.
 \end{aligned}$$

The last inequality follows from the fact that the covariance of any variable assignment is at most  $1/4$  in magnitude, and the covariance between any two assignments to the same variable is also at most  $1/4$ .

Given an upper bound on the covariances of node assignments, we can bound the covariance of edge assignments. Consider edges  $\{a, b\}, \{c, d\} \in \mathcal{E}$ . Due to the tree structure, the edges lie at opposite ends of a unique path connecting their constituent nodes. Without loss of generality, assume that this path has the order  $a, b, \dots, c, d$ , and that the length of the path from  $b$  to  $c$  is  $l$ . By the Markov property,  $Y_a$  and  $Y_d$  are conditionally independent given  $Y_b$  and  $Y_c$ . Thus, for any configuration  $(Y_a, Y_b) = (y_a, y_b)$  and  $(Y_c, Y_d) = (y_c, y_d)$ , we have that

$$\begin{aligned}
 & \sum_{y_c, y_d} |\sigma_{\theta}((y_a, y_b), (y_c, y_d))| \\
 &= \sum_{y_c, y_d} |p_{\theta}(y_a, y_b, y_c, y_d) - p_{\theta}(y_a, y_b)p_{\theta}(y_c, y_d)| \\
 &= \sum_{y_c, y_d} |p_{\theta}(y_a, y_d | y_b, y_c)p_{\theta}(y_b, y_c) \\
 &\quad - p_{\theta}(y_a | y_b)p_{\theta}(y_b)p_{\theta}(y_d | y_c)p_{\theta}(y_c)| \\
 &= \sum_{y_c, y_d} |p_{\theta}(y_a | y_b)p_{\theta}(y_d | y_c)p_{\theta}(y_b, y_c) \\
 &\quad - p_{\theta}(y_a | y_b)p_{\theta}(y_d | y_c)p_{\theta}(y_b)p_{\theta}(y_c)| \\
 &= \sum_{y_c, y_d} p_{\theta}(y_a | y_b)p_{\theta}(y_d | y_c) |\sigma_{\theta}(y_b, y_c)| \\
 &= p_{\theta}(y_a | y_b) \sum_{y_c} |\sigma_{\theta}(y_b, y_c)| \sum_{y_d} p_{\theta}(y_d | y_c) \\
 &= p_{\theta}(y_a | y_b) \sum_{y_c} |\sigma_{\theta}(y_b, y_c)| \\
 &\leq \frac{\ell}{4} (\vartheta_{\theta}^*)^{l-1}.
 \end{aligned}$$

The same argument can be used to bound the covariance between node and edge variables, where the relevant path

length  $l$  becomes the length from the node to the closest endpoint of the edge. The base case of covariance between a node or edge state indicator and another state is also at most  $1/4$ .

Thus far, we have derived upper bounds on the entries of the covariance matrix, which correspond to covariances between three types of pairs: node variables and node variables; node variables and edge variables; and edge variables and edge variables. For a distribution induced by a tree-structured model, with maximum degree  $\Delta_T$ , the 1-norm of a column corresponding to a node assignment  $Y_u = y_u$  is

$$\begin{aligned}
 \sigma_{\theta}(Y_u = y_u) &= \sum_{y'_u} |\sigma_{\theta}(y_u, y'_u)| + \sum_{v \in \mathcal{V} \setminus u} \sum_{y_v} |\sigma_{\theta}(y_u, y_v)| \\
 &\quad + \sum_{\{v, v'\} \in \mathcal{E}} \sum_{y_v, y_{v'}} |\sigma_{\theta}(y_u, (y_v, y_{v'}))| \\
 &\leq \frac{\ell}{4} + \frac{\ell}{4} \sum_{v \in \mathcal{V} \setminus u} (\vartheta_{\theta}^*)^{l(u, v) - 1} \\
 &\quad + \frac{\ell}{4} \sum_{\{v, v'\} \in \mathcal{E}} (\vartheta_{\theta}^*)^{\max\{0, \min\{l(u, v), l(u, v')\} - 1\}} \\
 &\leq \frac{\ell}{4} + \frac{\ell}{4} \sum_{d=1}^{\infty} \Delta_T^d (\vartheta_{\theta}^*)^{d-1} \\
 &\quad + \frac{\ell \Delta_T}{4} + \frac{\ell}{4} \sum_{d=1}^{\infty} \Delta_T^{d+1} (\vartheta_{\theta}^*)^{d-1} \\
 &= \frac{\ell}{4} + \frac{\ell \Delta_T}{4} \sum_{d=1}^{\infty} (\Delta_T \vartheta_{\theta}^*)^{d-1} \\
 &\quad + \frac{\ell \Delta_T}{4} + \frac{\ell \Delta_T^2}{4} \sum_{d=1}^{\infty} (\Delta_T \vartheta_{\theta}^*)^{d-1} \\
 &= \frac{\ell}{4} + \frac{\ell \Delta_T}{4(1 - \Delta_T \vartheta_{\theta}^*)} + \frac{\ell \Delta_T}{4} + \frac{\ell \Delta_T^2}{4(1 - \Delta_T \vartheta_{\theta}^*)}.
 \end{aligned}$$

where  $l(u, v)$  is the length of the path from node  $u$  to  $v$ . The second inequality holds because the number of nodes at distance  $d$  is at most  $\Delta_T^d$ , and the maximum number of edges with endpoints at distance  $d$  is at most  $\Delta_T^{d+1}$ , where we adjust for node and edge variables at distance zero. The last line applies the geometric series identity, since  $\Delta_T \vartheta_{\theta}^* < \Delta_T / \Delta_T = 1$ . An analogous argument bounds the 1-norm of any column corresponding to an edge assignment.

Since the 1-norm of every column of the covariance matrix is upper-bounded independently of  $|G|$ , it follows that the induced 1-norm of  $\Sigma(\mathbf{Y}; \theta)$  is bounded independently of  $|G|$ ; that is,

$$\|\Sigma(\mathbf{Y}; \theta)\|_1 = O(1).$$

This holds for every  $\theta \in \Theta$ , though the constant may differ, depending on  $\vartheta_{\theta}^*$ . Recall that the 1-norm of the covariance matrix upper-bounds the spectral norm, since the covariance matrix is symmetric. Thus, the minimum eigenvalue of  $\nabla^2(-H(\mu(\theta)))$ , for every  $\mu(\theta) \in \mathcal{M}(\Theta)$ , is lower-bounded by a constant, which means that the negative entropy is  $\Omega(1)$ -strongly convex in  $\mathcal{M}(\Theta)$ . ■

## C.2. Measuring Contraction

In the previous section, we relate the convexity of  $-H_T$  to the model's maximum contraction coefficient. For general graphical models, measuring the contraction coefficients may be intractable. However, when the model is tree-structured, there is an efficient algorithm.

For a tree-structured model, exact inference can be computed efficiently using message passing. Given the node and edge marginals, one can compute the conditional probabilities via

$$p(Y_u = y_u | Y_v = y_v; \theta) = \frac{p(Y_u = y_u, Y_v = y_v; \theta)}{p(Y_v = y_v; \theta)}.$$

One can then compute the total variation distance; hence, the contraction coefficient. For variables with small domains (e.g., binary), this is efficient. Given the contraction coefficient for each  $(u, v) : \{u, v\} \in \mathcal{E}$ , computing the maximum contraction coefficient is trivial.

Note that marginal inference only needs to be computed once in this procedure. The time complexity of inference in a tree-structured model, with  $\ell$  labels and  $|\mathcal{E}|$  edges is  $O(\ell^2 |\mathcal{E}|)$ . For each undirected edge, there are two contraction coefficients (one per direction), each of which involves  $\ell^2$  operations ( $\ell$  additions to compute the total variation distance conditioned on  $Y_v$ ; and  $\ell$  values of  $Y_v$  to condition on to compute the supremum). Since there are  $|\mathcal{E}|$  edges, the overall time complexity of computing the contraction coefficients is  $O(\ell^2 |\mathcal{E}|)$ .

## D. Tree-Reweighting

In this section, we prove Proposition 2, which gives a model-dependent lower bound on the modulus of convexity for the tree-reweighted negative entropy. We also explore the ramifications of Proposition 2 for a grid-structured model.

### D.1. Proof of $-H^{\text{TR}}$ Strong Convexity (Proposition 2)

The following lemma relates the convexity of  $-H^{\text{TR}}$  to the convexity of its constituent tree entropies, as well as the tree distribution.

**Lemma 7.** (Wainwright, 2006, Appendix C) Fix a graph,  $G \triangleq (\mathcal{V}, \mathcal{E})$ , and a distribution,  $\rho$ , over the spanning trees,  $\mathcal{T}(G)$ , such that  $\rho(e) > 0$  for all  $e \in \mathcal{E}$ . Let  $\rho_e^* \triangleq \min_{e \in \mathcal{E}} \rho(e)$  denote the minimum edge probability. Let  $\kappa_T^*$  denote the minimum convexity of  $-H_T$  for any tree  $T \in \mathcal{T}(G)$  with positive probability under  $\rho$ . Then the tree-reweighted negative entropy,  $-H^{\text{TR}}$ , is  $(\rho_e^* \kappa_T^*)$ -strongly convex.

Thus, to prove  $\Omega(1)$ -strong convexity, one must show that the minimum edge probability,  $\rho_e^*$ , and the minimum tree convexity,  $\kappa_T^*$  are both lower-bounded by values that are independent of  $|G|$ .

In Proposition 2, we assume that  $\rho_e^*$  is lower-bounded by a positive constant,  $C > 0$ . Since  $H^{\text{TR}}$  can be defined using any distribution over spanning trees, it is usually possible to construct an edge distribution for which this holds. (An example for a grid is given in Appendix D.2.) Therefore, the real challenge is to show that  $\kappa_T^* = \Omega(1)$ . For each  $T \in \mathcal{T}(G)$ , denote the set of admissible potentials by  $\Theta_T \subseteq \mathbb{R}^{|\theta|}$ , where dimensions corresponding to edges that don't exist in  $T$  have unbounded range. Note that

$$\Theta = \bigcap_{T \in \mathcal{T}(G): \rho(T) > 0} \Theta_T,$$

so

$$\tilde{\mathcal{M}}(\Theta) = \bigcap_{T \in \mathcal{T}(G): \rho(T) > 0} \tilde{\mathcal{M}}(\Theta_T).$$

Let  $\tilde{\mathcal{M}}_T(\Theta)$  denote the projection of  $\tilde{\mathcal{M}}(\Theta)$  onto the subspace defined by the nodes and edges in  $T$ , and note that  $\tilde{\mathcal{M}}_T(\Theta) \subseteq \tilde{\mathcal{M}}_T(\Theta_T)$ . In Proposition 4, we showed that, under suitable structural and contraction conditions,  $-H_T$  is  $\Omega(1)$ -strongly convex in  $\tilde{\mathcal{M}}_T(\Theta_T)$ ; hence, in  $\tilde{\mathcal{M}}_T(\Theta)$  as well. When combined with Lemma 7, with  $\rho_e^* > C$ , this proves that  $-H^{\text{TR}}$  is  $\Omega(1)$ -strongly convex in  $\tilde{\mathcal{M}}(\Theta)$ .

## D.2. Example Tree-Rewighting for a Grid Graph

Suppose the model is structured according to an  $m \times n$  grid. This graph can be covered using a set of 4 chains, using the “snake-like” pattern illustrated in Figure 1. Observe that each internal edge is covered by 2 chains, and each boundary edge is covered by 3 chains. Therefore, using a uniform distribution over the chains, we have that each internal edge,  $e$ , has probability  $\rho(e) = 1/2$ , and each boundary edge,  $e'$ , has probability  $\rho(e') = 3/4$ .

To apply Proposition 2 to this spanning tree distribution, we take  $C = 1/2$  as the minimum edge probability. The maximum degree of a chain is 2, so the maximum contraction coefficient,  $\vartheta_{\theta, T}^*$ , must be at most  $1/2$ . It may be possible to upper-bound  $\vartheta_{\theta, T}^*$  analytically for all  $\theta$  in some space. Alternately, one could map out the space of feasi-

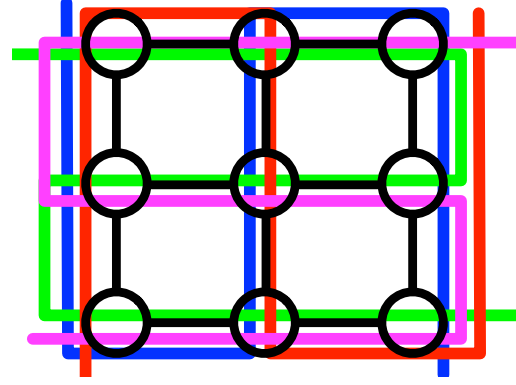


Figure 1. Covering the edges of a grid graph with 4 chains.

ble potentials by measuring  $\vartheta_{\theta, T}^*$ , using the procedure from Appendix C.2.

## E. Counting Numbers

In this section, we prove Proposition 3, which characterizes the modulus of convexity for counting number entropies. We also present a slackened version of the counting number QP, which can be used when the variable validity constraints are not satisfied.

### E.1. Proof of $-H^c$ Strong Convexity (Proposition 3)

The proof of Proposition 3 requires two technical lemmas.

**Lemma 8** (Shalev-Schwartz, 2007, Lemma 16). *The function  $\varphi(\mathbf{z}) \triangleq \sum_i^d z_i \log z_i$  is 1-strongly convex in the probability simplex,  $\{\mathbf{z} \in [0, 1]^d : \|\mathbf{z}\|_1 = 1\}$ , w.r.t. the 1-norm.*

**Lemma 9** (Heskes, 2006, Lemma A.1). *The difference of entropies, equivalent to the negative conditional entropy,  $H_v(\tilde{\mu}_v) - H_e(\tilde{\mu}_e) = -H_{e|v}(\tilde{\mu}_e)$ , for  $v \in e$ , is a convex function of  $\tilde{\mu}_e$ .*

We now prove Proposition 3.

**Proof** [Proposition 3] Every edge,  $e$ , is composed of exactly two nodes,  $\{u, v\}$ . By assumption, we have that  $\alpha_e \geq \kappa > 0$ . Therefore, we can shift  $(2\kappa/3)$  weight from  $\alpha_e$  to  $\alpha_u$  and  $\alpha_v$  without affecting the counting numbers or Heskes’s convexity conditions. Let:

$$\begin{aligned} \forall e \in \mathcal{E}, \quad \tilde{\alpha}_e &\triangleq \alpha_e - \frac{2\kappa}{3}; \\ \forall (v, e) : v \in e, \quad \tilde{\alpha}_{v,e} &\triangleq \alpha_{v,e} + \frac{\kappa}{3}; \\ \forall v \in \mathcal{V}, \quad \tilde{\alpha}_v &\triangleq \alpha_v + \sum_{e:v \in e} \frac{\kappa}{3}. \end{aligned}$$

Observe that the new auxiliary counts satisfy Eqs. 11

and 12:

$$\begin{aligned} \forall v \in \mathcal{V}, c_v &= \alpha_v - \sum_{e:v \in e} \left( \alpha_{v,e} + \frac{\kappa}{3} - \frac{\kappa}{3} \right) \\ &= \tilde{\alpha}_v - \sum_{e:v \in e} \tilde{\alpha}_{v,e}; \end{aligned} \quad (20)$$

$$\begin{aligned} \forall e \in \mathcal{E}, c_e &= \alpha_e + \sum_{v:v \in e} \left( \alpha_{v,e} + \frac{\kappa}{3} - \frac{\kappa}{3} \right) \\ &= \tilde{\alpha}_e + \sum_{v:v \in e} \tilde{\alpha}_{v,e}. \end{aligned} \quad (21)$$

Now, every  $e$  has  $\tilde{\alpha}_e \geq \kappa/3$ . Further, because we assume that every node is involved in at least one edge, every  $v$  has  $\tilde{\alpha}_v \geq \kappa/3$ . (We could extend Proposition 3 to arbitrary graphs by assuming that every isolated node has  $c_v \geq \kappa/3$ .)

Substituting Eqs. 20 and 21 into Eq. 10 and rearranging the terms, we obtain

$$\begin{aligned} -H^c(\tilde{\boldsymbol{\mu}}) &= - \sum_{v \in \mathcal{V}} \tilde{\alpha}_v H_v(\tilde{\boldsymbol{\mu}}_v) - \sum_{e \in \mathcal{E}} \tilde{\alpha}_e H_e(\tilde{\boldsymbol{\mu}}_e) \\ &\quad + \sum_{e \in \mathcal{E}} \sum_{v:v \in e} \tilde{\alpha}_{v,e} (H_v(\tilde{\boldsymbol{\mu}}_v) - H_e(\tilde{\boldsymbol{\mu}}_e)) \\ &= - \sum_{v \in \mathcal{V}} \tilde{\alpha}_v H_v(\tilde{\boldsymbol{\mu}}_v) - \sum_{e \in \mathcal{E}} \tilde{\alpha}_e H_e(\tilde{\boldsymbol{\mu}}_e) \\ &\quad - \sum_{e \in \mathcal{E}} \sum_{v:v \in e} \tilde{\alpha}_{v,e} H_{e|v}(\tilde{\boldsymbol{\mu}}_e). \end{aligned} \quad (22)$$

We will analyze the entropy terms individually, using the gradient definition of (strong) convexity.

Fix any two vectors  $\tilde{\boldsymbol{\mu}}, \tilde{\boldsymbol{\mu}}' \in \tilde{\mathcal{M}}$ , and let  $\boldsymbol{\delta} \triangleq \tilde{\boldsymbol{\mu}} - \tilde{\boldsymbol{\mu}}'$ . Recall that  $\forall v, \|\tilde{\boldsymbol{\mu}}_v\|_1 = \|\tilde{\boldsymbol{\mu}}'_v\|_1 = 1$  and  $\forall e, \|\tilde{\boldsymbol{\mu}}_e\|_1 = \|\tilde{\boldsymbol{\mu}}'_e\|_1 = 1$ . Via Lemma 8,  $-H_v$  and  $-H_e$  are 1-strongly convex in the probability simplex with respect to the 1-norm. By Fact 1, this means that every node  $v$  satisfies,

$$\langle \nabla(-H_v(\tilde{\boldsymbol{\mu}}_v)) - \nabla(-H_v(\tilde{\boldsymbol{\mu}}'_v)), \boldsymbol{\delta}_v \rangle \geq \|\boldsymbol{\delta}_v\|_1^2.$$

Therefore,

$$\begin{aligned} \tilde{\alpha}_v \langle \nabla(-H_v(\tilde{\boldsymbol{\mu}}_v)) - \nabla(-H_v(\tilde{\boldsymbol{\mu}}'_v)), \boldsymbol{\delta}_v \rangle &\geq \tilde{\alpha}_v \|\boldsymbol{\delta}_v\|_1^2 \\ &\geq \tilde{\alpha}_v \|\boldsymbol{\delta}_v\|_2^2 \\ &\geq \frac{\kappa}{3} \|\boldsymbol{\delta}_v\|_2^2. \end{aligned}$$

The same holds for every edge  $e$ . Further, by Lemma 9,  $H_{e|v}(\tilde{\boldsymbol{\mu}}_e) = H_v(\tilde{\boldsymbol{\mu}}_v) - H_e(\tilde{\boldsymbol{\mu}}_e)$  is convex, meaning

$$\langle \nabla(-H_{e|v}(\tilde{\boldsymbol{\mu}}_e)) - \nabla(-H_{e|v}(\tilde{\boldsymbol{\mu}}'_e)), \boldsymbol{\delta}_e \rangle \geq 0.$$

Thus, taking the gradient of Eq. 22, we have that

$$\begin{aligned} &\langle \nabla(-H^c(\tilde{\boldsymbol{\mu}})) - \nabla(-H^c(\tilde{\boldsymbol{\mu}}')), \boldsymbol{\delta} \rangle \\ &= \sum_{v \in \mathcal{V}} \tilde{\alpha}_v \langle \nabla(-H_v(\tilde{\boldsymbol{\mu}}_v)) - \nabla(-H_v(\tilde{\boldsymbol{\mu}}'_v)), \boldsymbol{\delta}_v \rangle \\ &\quad + \sum_{e \in \mathcal{E}} \tilde{\alpha}_e \langle \nabla(-H_e(\tilde{\boldsymbol{\mu}}_e)) - \nabla(-H_e(\tilde{\boldsymbol{\mu}}'_e)), \boldsymbol{\delta}_e \rangle \\ &\quad + \sum_{e \in \mathcal{E}} \sum_{v:v \in e} \tilde{\alpha}_{v,e} \langle \nabla(-H_{e|v}(\tilde{\boldsymbol{\mu}}_e)) - \nabla(-H_{e|v}(\tilde{\boldsymbol{\mu}}'_e)), \boldsymbol{\delta}_e \rangle \\ &\geq \frac{\kappa}{3} \sum_{v \in \mathcal{V}} \|\boldsymbol{\delta}_v\|_2^2 + \frac{\kappa}{3} \sum_{e \in \mathcal{E}} \|\boldsymbol{\delta}_e\|_2^2 + 0 \\ &= \frac{\kappa}{3} \|\boldsymbol{\delta}\|_2^2, \end{aligned}$$

which completes the proof, via Fact 1.  $\blacksquare$

## E.2. Slackened Variable-Valid Counting Number Optimization

For certain values of  $\kappa$ , the variable validity constraints in Eq. 13 create an infeasible optimization problem. When this happens, we propose switching to a slackened QP. This introduces a free parameter,  $C \geq 0$ , that adjusts the trade-off between fitting the target counts (in the equation below, the Bethe counts) and variable validity. The slackened QP is then

$$\begin{aligned} &\min_{\mathbf{c}, \boldsymbol{\alpha} \geq 0, \boldsymbol{\xi}} \|\mathbf{c} - \mathbf{c}^B\|_2^2 + C \|\boldsymbol{\xi}\|_2^2 \quad (23) \\ &\text{s.t. } \forall v \in \mathcal{V}, c_v + \sum_{e:v \in e} \alpha_{v,e} \geq 0; \\ &\quad \forall e \in \mathcal{E}, c_e - \sum_{v:v \in e} \alpha_{v,e} \geq 3\kappa; \\ &\quad \forall v \in \mathcal{V}, c_v + \sum_{e:v \in e} c_e = 1 + \xi_v. \end{aligned}$$

## F. Figures for Experimental Results

In all plots, results are averaged over 20 trials and the  $y$ -axis has been rescaled to fit the data. See Section 5.3 for discussion.

## References

- Heskes, T. Convexity arguments for efficient minimization of the Bethe and Kikuchi free energies. *Journal of Artificial Intelligence Research*, 26:153–190, 2006.
- Kontorovich, A. Obtaining measure concentration from Markov contraction. *Markov Processes and Related Fields*, 18:613–638, 2012.
- London, B., Huang, B., Taskar, B., and Getoor, L. PAC-Bayesian collective stability. In *Artificial Intelligence and Statistics*, 2014.

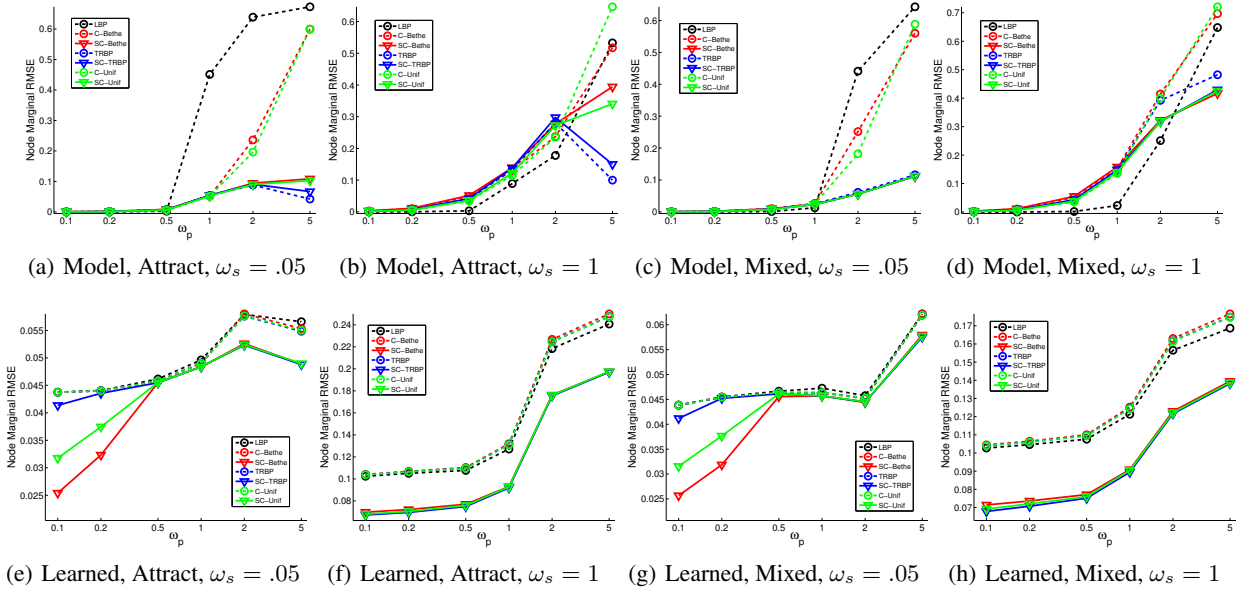


Figure 2. Plots of RMSE of the node marginals as a function of the interaction parameter,  $\omega_p$ . Inference is performed using the true model in (a)-(d), and the learned model in (e)-(h). The first two columns correspond to a model with attractive potentials; the third and fourth to a model with mixed potentials. The black dotted line is LBP; color dotted lines are the convex baselines, and solid lines are their SC counterparts. The SC methods use the post hoc optimal value of  $\kappa$  (and  $C$ ) in the counting number optimization. For learned marginals, SC offers statistically significant error reduction—sometimes over 40%—for all data models and baselines, except C-Bethe at  $\omega_p = .5$  in (g).

Shalev-Schwartz, S. *Online Learning: Theory, Algorithms, and Applications*. PhD thesis, The Hebrew University of Jerusalem, 2007.

Wainwright, M. Estimating the “wrong” graphical model: Benefits in the computation-limited setting. *Journal of Machine Learning Research*, 7:1829–1859, 2006.

Wainwright, M. and Jordan, M. *Graphical Models, Exponential Families, and Variational Inference*. Now Publishers Inc., 2008.



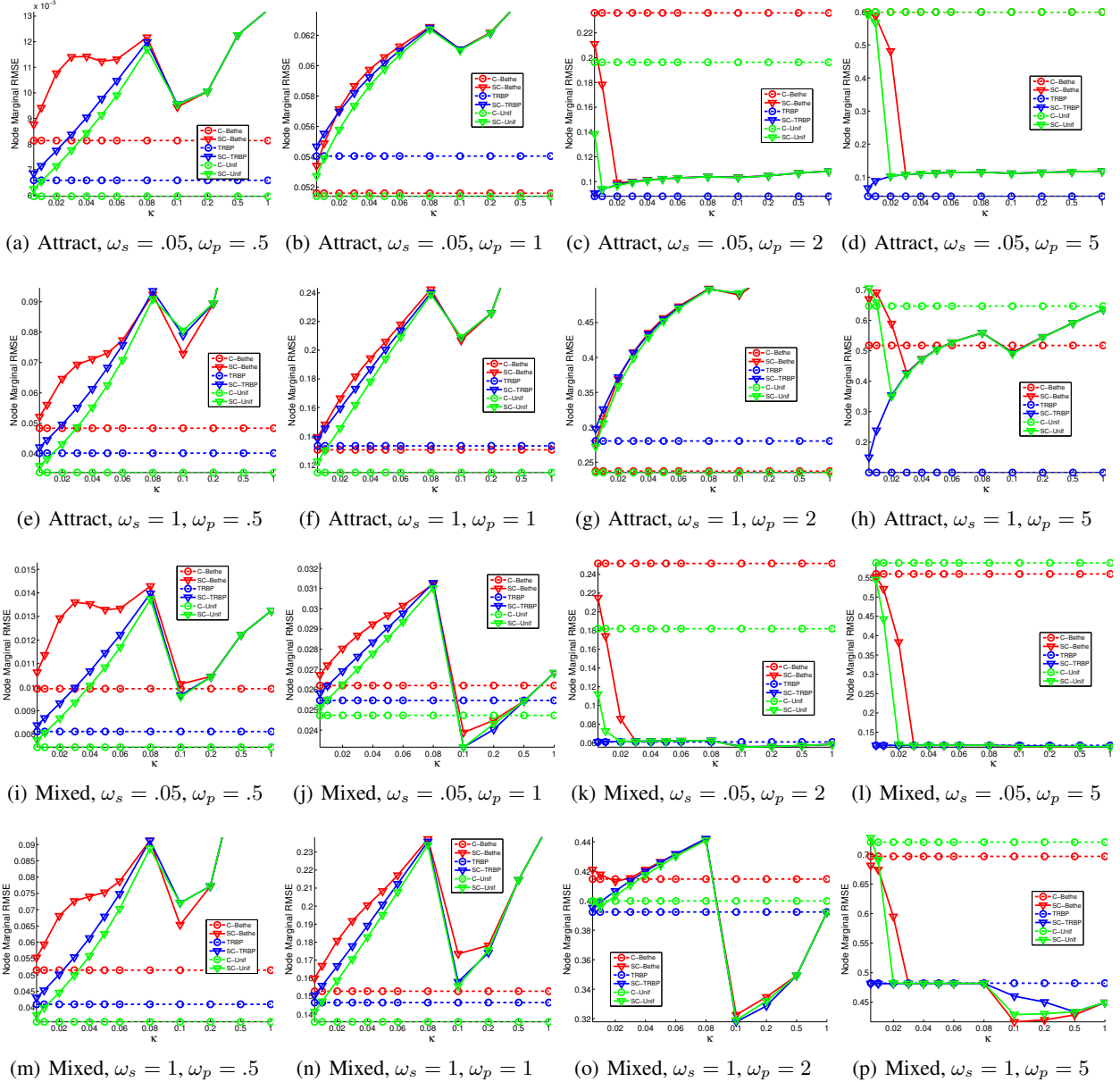


Figure 3. Plots of RMSE of the node marginals as a function of the convexity parameter,  $\kappa$ , which determines the minimum modulus of convexity used in the counting number QP. For  $\kappa < .1$ , we use Eq. 13; for  $\kappa \geq .1$ , we use Eq. 23 and report the score for the post hoc optimal  $C$ . SC algorithms are plotted as solid lines, and their respective counterparts are overlaid as dashed lines. Inference is performed using the true model. The first two rows correspond to a model with attractive potentials; the third and fourth to a model with mixed potentials. In all plots, the  $x$ -axis scales logarithmically for  $\kappa > .1$ . Certain plots have been truncated vertically to better fit the data.

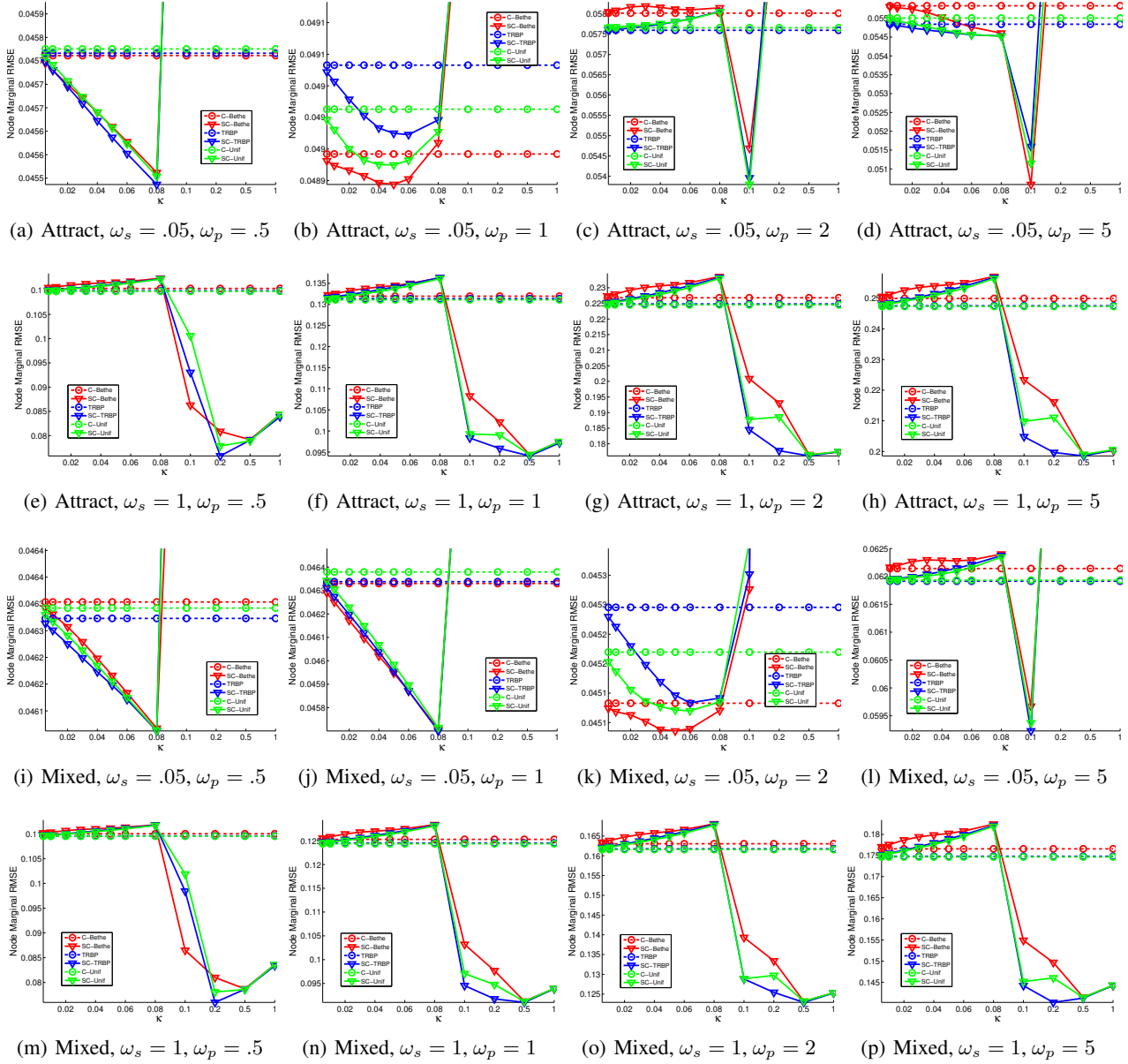


Figure 4. Plots of RMSE of the node marginals as a function of the convexity parameter,  $\kappa$ , when using the learned model for inference.

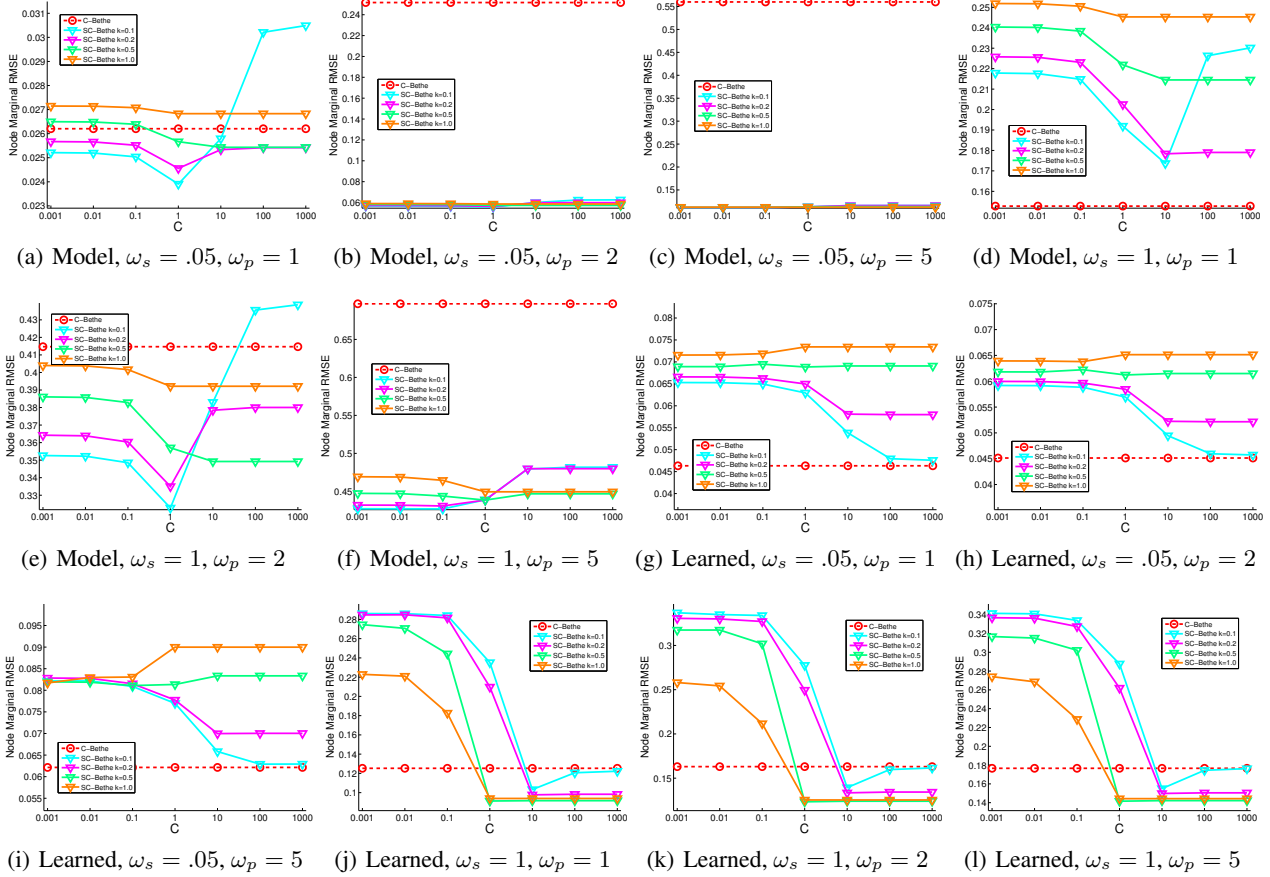


Figure 5. Select plots of RMSE as a function of the slack parameter,  $C$ , in the slackened counting number QP (Eq. 23), at higher values of  $\kappa$ . The slack parameter trades off between fitting the target counting numbers and satisfying variable validity. Data is generated using mixed potentials in all plots. These plots focus on the Bethe approximation. SC versions are solid color lines; C-Bethe is overlaid as a dashed red line.

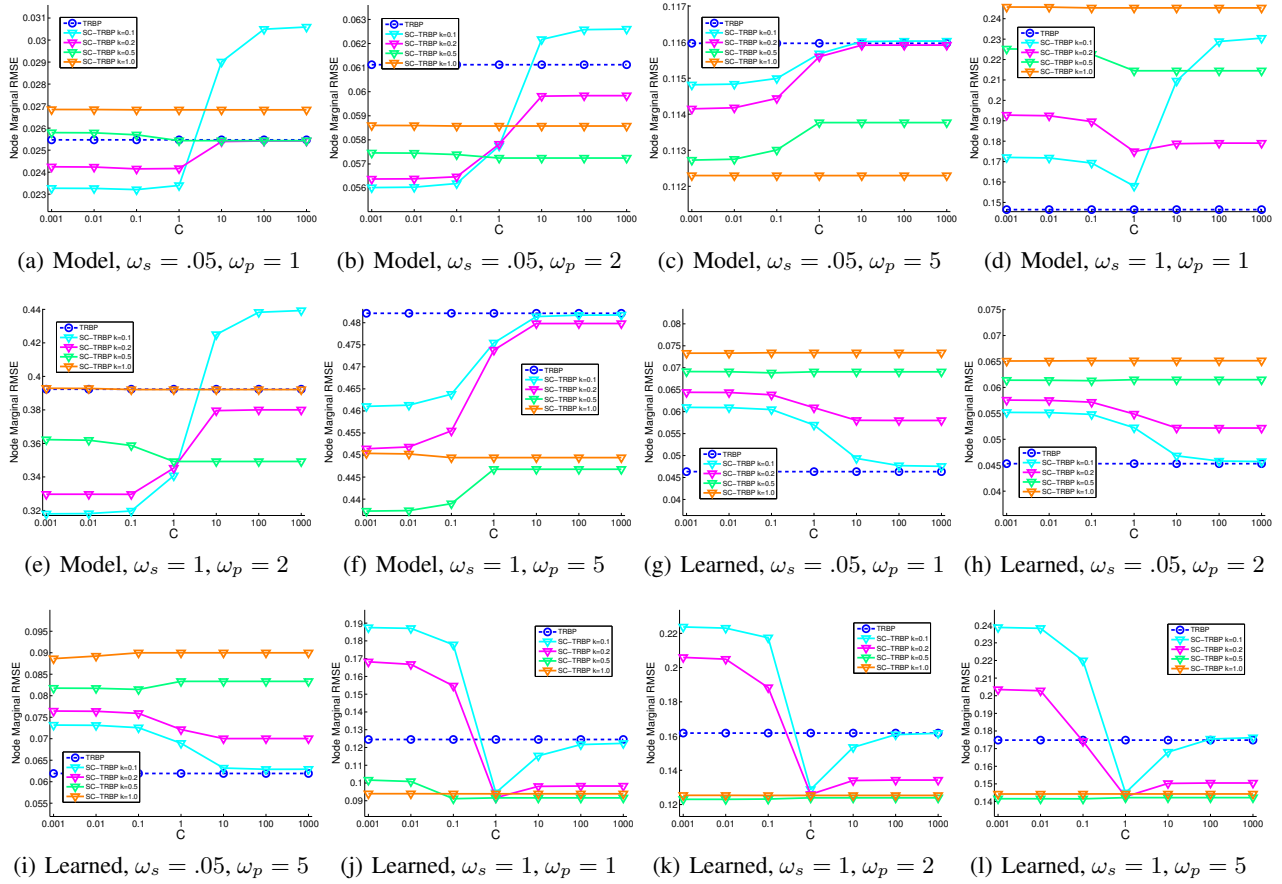


Figure 6. Select plots of RMSE as a function of the slack parameter,  $C$ , for the tree-reweighting approximation. SC versions are solid color lines; C-TRBP is overlaid as a dashed blue line.