

---

# On the Use of Variational Inference for Learning Discrete Graphical Models

---

Eunho Yang  
Pradeep Ravikumar

EUNHO@CS.UTEXAS.EDU  
PRADEEPR@CS.UTEXAS.EDU

Department of Computer Science, The University of Texas, Austin, TX 78712, USA

## Abstract

We study the general class of estimators for graphical model structure based on optimizing  $\ell_1$ -regularized approximate log-likelihood, where the approximate likelihood uses tractable variational approximations of the partition function. We provide a message-passing algorithm that *directly* computes the  $\ell_1$  regularized approximate MLE. Further, in the case of certain reweighted entropy approximations to the partition function, we show that surprisingly the  $\ell_1$  regularized approximate MLE estimator has a *closed-form*, so that we would no longer need to run through many iterations of approximate inference and message-passing. Lastly, we analyze this general class of estimators for graph structure recovery, or its *sparsistency*, and show that it is indeed sparsistent under certain conditions.

## 1. Introduction

A Markov random field (MRF) over a  $p$ -dimensional discrete random vector  $X = (X_1, X_2, \dots, X_p)$  is specified by an undirected graph  $G = (V, E)$ , with vertex set  $V = \{1, 2, \dots, p\}$  – one for each variable – and edge set  $E \subset V \times V$ . The structure of this graph encodes conditional independence assumptions among subsets of the variables. In structure learning, the task is to estimate this underlying graph from  $n$  independent and identically distributed samples  $\{x^{(1)}, x^{(2)}, \dots, x^{(n)}\}$ .

Recent results on the discrete graphical model structure learning problem have built on a natural connection between such structure learning and *parameter*

*estimation*: indeed, learning the graph structure is equivalent to learning the parameters of a fully saturated graphical model under the assumption of a *sparse* set of parameters, where parameters underlying non-edges are equal to zero. Such parameter estimation in turn hinges crucially on graphical model *inference*: since it involves solving optimization problems that require computing quantities such as the partition function, or the marginals.

Thus, structure learning when reduced to sparse parameter estimation hinges on two components: sparsity recovering regularization methods, and methods for approximate inference in graphical models. While methods that have been proposed in this class of approaches use state of the art sparsity recovery methods ( $\ell_1$ -regularization and the like), their approximate inference components are far from the state of the art. For instance, Ravikumar et al. (2010) use node-wise logistic regressions to estimate node-neighborhoods, which can be thought of as employing a pseudo-likelihood approximation (with asymmetric edge parameters) of the partition function in the log-likelihood. Lee et al. (2007) compute approximate estimates of the gradient using Belief Propagation (Pearl, 1988), which can be thought of as using a Bethe entropy approximation of the partition function. The state of the art in approximate inference on the other hand involves convex variational approximations of the entropy of the graphical model, and dual decompositions involving tractable subcomponents of the graphical model.

We thus have a gap: between state of the art in inference and the use of inference in state of the art in structure learning methods. Towards this, we study a *general* class of estimators for graphical model structure that use tractable approximations of the partition function and  $\ell_1$ -regularization. The resulting optimization problem can be solved naturally using composite gradient descent, since the gradients of the log-likelihood take the form of marginals which can be

---

Appearing in *Proceedings of the 28<sup>th</sup> International Conference on Machine Learning*, Bellevue, WA, USA, 2011. Copyright 2011 by the author(s)/owner(s).

approximated using the approximate inference procedure corresponding to the partition function approximation. However this involves performing graphical model inference for each of the gradient steps, which could be expensive. To address this, we provide a message-passing algorithm that *directly* computes the solution that optimizes the  $\ell_1$  regularized approximate log-likelihood. Further, in the case of certain reweighted entropy approximations to the partition function such as the tree-reweighted approximation, we show that surprisingly the  $\ell_1$  regularized approximate MLE estimator has a *closed-form*, so that we would no longer need to run through many iterations of approximate inference and message-passing. Lastly, we analyze this general class of estimators for graph structure recovery, or its *sparsistency*. For such estimators, one might imagine that even though the approximate inference is tight with respect to the partition function or the marginals, the corresponding approximate MLE need not even be consistent; see (Wainwright, 2006) for instance for the case of a weakly regularized approximate MLE. However, note that this is also the case with high-dimensional sparse parameter estimation where typical MLE estimators are not consistent, *unless* one carefully chooses the magnitude of  $\ell_1$  regularization (Candes & Tao, 2007; Tropp, 2006). Indeed, even for our general class of  $\ell_1$  regularized *approximate* log-likelihood estimators, we show that under certain conditions on the edge weights, the methods do succeed in recovering the graph structure. Indeed, the development in this paper raises the research agenda of tuning approximate inference procedures to structure learning by developing partition function approximations that would impose the weakest conditions on the parameters.

## 2. Review, Setup and Notation

### 2.1. Markov Random Fields

Let  $X = (X_1, \dots, X_p)$  be a random vector, each variable  $X_i$  taking values in a discrete set  $\mathcal{X}$  of cardinality  $m$ . Let  $G = (V, E)$  denote a graph with  $p$  nodes, corresponding to the  $p$  variables  $\{X_1, \dots, X_p\}$ . A pairwise Markov random field over  $X = (X_1, \dots, X_p)$  is then specified by nodewise and pairwise functions  $\theta_s : \mathcal{X} \mapsto \mathbb{R}$  for all  $s \in V$ , and  $\theta_{st} : \mathcal{X} \times \mathcal{X} \mapsto \mathbb{R}$  for all  $(s, t) \in E$ , as

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

In this paper, we largely focus on the case where the variables are binary with  $\mathcal{X} = \{-1, +1\}$ , where we can rewrite (1) to the Ising model form (Ising, 1925)

$$\mathbb{P}(x) \propto \exp \left\{ \sum_{s \in V} \theta_s x_s + \sum_{(s,t) \in E} \theta_{st} x_s x_t \right\}, \quad (2)$$

for some set of parameters  $\{\theta_s\}$  and  $\{\theta_{st}\}$ . It will be useful to rewrite (2) as the member of an exponential family,  $\mathbb{P}(x) = \exp(\langle \theta, \phi \rangle - A(\theta))$ , where  $\phi(x) = \{\phi_s(x) = x_s; \phi_{st}(x) \equiv x_s x_t\}$  are the set of Ising potentials, and  $\theta = \{\theta_s; \theta_{st}\}$  are the corresponding set of parameters; and  $A : \Theta \mapsto \mathbb{R}$  is the log of the normalization constant, also called the log-partition function,  $A(\theta) = \log \sum_{x \in \mathcal{X}^p} \exp(\langle \theta, \phi \rangle)$ .

### 2.2. Variational Approximations

A complication for discrete undirected graphical models is that typical inference tasks, even calculation of the log-partition function  $A(\theta)$ , is computationally intractable. Here, we briefly review variational approximations to the partition function, following the development in Wainwright & Jordan (2008). By properties of exponential families, the moments  $\mu(\theta) = \mathbb{E}_\theta(\phi) = \nabla A(\theta)$ . Denote the conjugate of the log-partition function by  $A^*(\mu) = \sup_{\theta \in \Theta} \langle \theta, \mu \rangle - A(\theta)$ . It can be shown that  $A^*(\mu)$  is the negative entropy of the graphical model distribution with parameter  $\theta = (\nabla A)^{-1}(\mu)$ . Consider the set of all possible mean parameters,  $\mathcal{M} = \{\mu : \exists \text{ distribution } p \text{ s.t. } \mathbb{E}_p(\phi) = \mu\}$ , which is also called the *marginal polytope* of the graphical model. Then, by convex duality, we can write,

$$A(\theta) = \sup_{\mu \in \mathcal{M}} \langle \theta, \mu \rangle - A^*(\mu). \quad (3)$$

(3) thus provides a variational formulation of the log-partition function  $A(\theta)$ . Following the development in Wainwright & Jordan (2008), it is easier to describe approximations to this log-partition function using so called *overcomplete representations*.

Since  $\mathcal{X}$  is discrete, any potential function  $\theta_c$  can be parameterized as linear combinations of  $\{0, 1\}$ -valued indicator functions. For each  $s \in V$  and  $j \in \{1, \dots, m\}$ , we can define node-wise indicators,  $[x_s = j] = 1$  if  $x_s = j$  and equal to 0 otherwise. With this notation, any set of potential functions can then be written as  $\theta_s(x_s) = \sum_{j \in [m]} \theta_{s;j} [x_s = j]$  for  $s \in V$  and  $\theta_{st}(x_s, x_t) = \sum_{j,k \in [m]} \theta_{st;jk} [x_s = j, x_t = k]$  for  $(s, t) \in E$ . Thus, (1) can be rewritten as,  $\mathbb{P}(x) \propto \exp \left\{ \sum_{s \in V; j \in [m]} \theta_{s;j} [x_s = j] + \sum_{(s,t) \in E; j,k \in [m]} \theta_{st;jk} [x_s = j, x_t = k] \right\}$ , for a set of parameters  $\theta := \{\theta_{s;j}, \theta_{st;jk} : s, t \in V; (s, t) \in E; j, k \in [m]\}$ .

Given these sufficient statistics, the mean parameters  $\{\mu_{s;j}\}$  and  $\{\mu_{st;jk}\}$  are just the node and pairwise marginals. Wainwright & Jordan (2008) then describe variational approximations of the log-partition function  $A(\theta)$ , as involving approximating the two intractable components in its variational formulation (3) (a) the marginal polytope  $\mathcal{M}$ , and (b) the graphical model entropy  $A^*(\mu)$ .

Any variational approximation to the log-partition function (3) can then be written as,

$$B(\theta) = \sup_{\mu \in \mathcal{L}} \langle \theta, \mu \rangle - B^*(\mu), \quad (4)$$

where  $\mathcal{L}$  is a tractable bound on the marginal polytope  $\mathcal{M}$ , and  $B^*(\mu)$  is a tractable approximation to the graphical model entropy  $A^*(\mu)$ . A popular bound  $L_G$  of the

marginal polytope is given by

$$L_G = \left\{ \mu \mid \sum_j \mu_{s;j} = 1, \sum_k \mu_{st;jk} = \mu_{s;j}; \right. \\ \left. s, t \in V; j, k \in [m] \right\}. \quad (5)$$

Popular approximations to the negative entropy use weighted sums of node and edge-entropies. Let  $H_s := \sum_{x_s \in \mathcal{X}} \mu_s(x_s) \log \mu_s(x_s)$  and  $H_{st}(\mu_{st}) := \sum_{(x_s, x_t) \in \mathcal{X} \times \mathcal{X}} \mu_{st}(x_s, x_t) \log \mu_{st}(x_s, x_t)$  denote the node-based and edge-based negative entropies, respectively, and  $I_{st}(\mu_{st}) := H_{st}(\mu_{st}) - H_s(\mu_s) - H_t(\mu_t)$ .

The Bethe approximation (Yedidia et al., 2001) to the entropy  $A^*(\mu)$  is given by  $B_{\text{bethe}}^*(\mu) = \sum_s H_s(\mu_s) - \sum_{st} I_{st}(\mu_{st})$ . The tree-reweighted entropy (Wainwright et al., 2003) in turn is given by

$$B_{\text{trw}}^*(\mu) = \sum_s H_s(\mu_s) - \sum_{st} \rho_{st} I_{st}(\mu_{st}), \quad (6)$$

where  $\rho_{st}$  are edge-weights that lie in a so-called spanning tree polytope. Heskes (2006); Weiss et al. (2007) also discuss general convex entropic forms, the simplest of which are the weighted forms  $B_\alpha^*(\mu) = \sum_{s \in V} \alpha_s H_s(\mu_s) + \sum_{(s,t) \in E} \alpha_{st} H_{st}(\mu_{st})$ , for some weights  $\{\alpha_s, \alpha_{st} \geq 0\}$ .

The approximation given by  $B_{\text{bethe}}(\theta) = \sup_{\mu \in L_G} \langle \theta, \mu \rangle - B_{\text{bethe}}^*(\mu)$ , underlies belief propagation, while  $B_{\text{trw}}(\theta) = \sup_{\mu \in L_G} \langle \theta, \mu \rangle - B_{\text{trw}}^*(\mu)$  yields the tree-reweighted approximation to the log-partition function.

### 3. Graphical Model Selection

Suppose that we are given a collection  $D := \{x^{(1)}, \dots, x^{(n)}\}$  of  $n$  samples, where each  $p$ -dimensional vector  $x^{(i)} \in \{1, \dots, m\}^p$  is drawn i.i.d. from a distribution  $\mathbb{P}_{\theta^*}$  of the form (2), for parameters  $\theta^*$  and graph  $G = (V, E^*)$  over the  $p$  variables. The goal of *graphical model selection* is to infer the edge set  $E^*$  of the graphical model defining the probability distribution that generates the samples. Note that the true edge set  $E^*$  can also be expressed as a function of the parameters as

$$E^* = \{(s, t) \in V \times V : \theta_{st}^* \neq 0\}. \quad (7)$$

Given the data,  $D := \{x^{(1)}, \dots, x^{(n)}\}$ , the  $\ell_1$  regularized MLE can then be written as the solution of the optimization problem,

$$\hat{\theta} \in \arg \min_{\theta} -\langle \theta, \hat{\phi} \rangle + A(\theta) + \lambda \|\theta\|_{1,E}, \quad (8)$$

where  $\hat{\phi} = \frac{1}{n} \sum_{i=1}^n \phi(x^{(i)})$  is the average of the sufficient statistics, and where  $\|\cdot\|_{1,E}$  is the  $\ell_1$  norm of just the edge-parameters, so that  $\|\theta\|_{1,E} = \sum_{s \neq t} |\theta_{st}|$ .

The caveat with solving (8) is the intractable computation of the log-partition function  $A(\theta)$ . We thus consider the following class of  $M$ -estimators:

$$\hat{\theta} \in \arg \min_{\theta} -\langle \theta, \hat{\phi} \rangle + B(\theta) + \lambda \|\theta\|_{1,E}, \quad (9)$$

where  $B(\theta)$  is a variational approximation to the log-partition function of the form (3) outlined in the previous section. Given the solution  $\hat{\theta}$ , we can then estimate the graph structure as  $\hat{E} = \{(s, t) : \hat{\theta}_{st} \neq 0\}$ .

## 4. Optimization Methods

We now consider the task of solving the optimization problem in (9).

### 4.1. Gradient based methods

Let us first go back to the exact  $\ell_1$ -regularized MLE problem (8), and consider an approximate technique to solve this intractable optimization problem. In particular, as Lee et al. (2007) suggest we could solve (8) using approximate estimates of the gradient, computed using Belief Propagation. One could then perform the following *approximate* composite gradient descent (Nesterov, 2004):

---

**Algorithm 1** Solving (8) Using Approximate Marginals

---

**for**  $t = 1, 2, \dots$  **do**

$$\theta^{t+1} = \mathcal{S}_{\lambda \eta^t} \left( \theta^t - \eta^t (-\hat{\phi} + \mu^{\text{apx}}(\theta^t)) \right).$$

**end for**

where  $\mu^{\text{apx}}(\theta^t)$  are approximate estimates of “moments”  $\nabla A(\theta^t)$ .

---

Here  $\mathcal{S}_r^{(E)}$  denotes the soft-thresholding function applied to the edge elements, so that for  $\alpha \in E$ ,  $[\mathcal{S}_r^{(E)}(w)]_\alpha = \text{sign}(w_\alpha) \max\{|w_\alpha| - r, 0\}$ . Such composite gradient descent has been shown to have at least sublinear convergence provided the step-sizes are chosen appropriately (Nesterov, 2004). Indeed, one way to view these iterates given objective  $f(\theta)$  is as minimizing a composite quadratic approximation:  $\min_{\theta} \nabla f(\theta^t) \cdot (\theta - \theta^t) + 1/\eta^t \|\theta - \theta^t\|_2^2 + \lambda \|\theta\|_{1,E}$ . Vanilla gradient descent on the other hand solves  $\min_{\theta} \nabla f(\theta^t) \cdot (\theta - \theta^t) + 1/\eta^t \|\theta - \theta^t\|_2^2$ .

**Proposition 1.** *Suppose that in Algorithm 1 the approximate marginals  $\mu^{\text{aprox}}(\theta)$  satisfy  $\mu^{\text{aprox}}(\theta) =$*

$\nabla B(\theta)$ , for some approximation  $B(\theta)$  to the log-partition function  $A(\theta)$ , so that they are pseudo-moments under the approximation  $B(\theta)$ . Then the fixed point of Algorithm 1 if any is a local minimum of the optimization problem in (9) with  $B(\theta)$  as the log-partition function approximation.

Thus, estimating the marginals using belief propagation as in Lee et al. (2007) corresponds to a Bethe entropy approximation to the partition function  $B_{\text{bethe}}(\theta)$ .

## 4.2. Message Passing Updates

In the sequel, we derive a message-passing algorithm for solving the estimator in (9). It will again be useful to consider the overcomplete representation as outlined in Section 2.2. We overload notation and continue to use  $\hat{\phi}$  and  $\theta$  for the sufficient statistics and parameters. We thus solve the  $\ell_1$  regularized optimization problem

$$\hat{\theta} \in \arg \min_{\theta} -\langle \theta, \hat{\phi} \rangle + B(\theta) + \lambda \|\theta\|_{1,E}. \quad (10)$$

Note that even though the overcomplete representation is not identifiable, the added  $\ell_1$ -regularization makes the solution unique provided the smooth component is strictly convex. Indeed, for the binary Ising model case, we have the equivalence:

**Proposition 2.** *Suppose that  $\tilde{\theta}$  is the unique solution of the approximate MLE (9) for the Ising model with regularization penalty  $4\lambda$ . Then the overcomplete estimator in (10) with regularization penalty  $\lambda$  has a unique solution  $\hat{\theta}$  given by*

$$\hat{\theta}_{st}(x_s, x_t) = \tilde{\theta}_{st} x_s x_t.$$

Now, by duality, this overcomplete approximate MLE (10) can be rewritten as

$$\inf_{\theta} \sup_{Z \in C} -\langle \theta, \hat{\phi} \rangle + B(\theta) + \langle \theta, Z \rangle,$$

where  $C := \{Z \subseteq \Theta : \|Z\|_{\infty} \leq \lambda, Z_N = 0\}$ , where we use  $Z_N$  to denote the coordinates of  $Z$  corresponding to node potentials. By strong duality (Boyd & Vandenberghe, 2004), this in turn can be rewritten as

$$\sup_{Z \in C} \inf_{\theta} -\langle \theta, \hat{\phi} - Z \rangle + B(\theta) = - \inf_{\mu \in W} B^*(\mu), \quad (11)$$

where  $W := \{\mu : \mu \in L_G; \mu_N = \phi_N; \|\mu_E - \phi_E\|_{\infty} \leq \lambda\}$ , where  $L_G$  is the approximation to the marginal polytope  $\mathcal{M}$  underlying the variational approximation  $B(\theta)$ , as outlined in Section 2.2.

We note that this optimization problem is very similar to the typical variational optimization problems for approximation of the partition function, where a linear term and the entropy are optimized with  $\mu$  constrained to the outer polytope  $L$ . These latter optimization problems are typically solved using graph-structured message-passing algorithms, and here, we derive a message passing algorithm that solves the above objective instead, so that it would obtain the  $\ell_1$ -regularized approximate MLE in one shot in contrast to the iterative message passing in the previous section.

Towards this, we use an iterative projection method (Censor & Zenios, 1988), that iteratively projects the primal variables onto individual constraints, while maintaining dual feasibility.

Now, note that the dual in (11) of the  $\ell_1$  regularized approximate MLE has the following form:

$$\begin{aligned} \inf_{\mu} B^*(\mu) \\ \text{s.t. } \langle a_i, \mu \rangle = b_i, \quad i = 1, \dots, m. \\ l_j \leq \langle c_j, \mu \rangle \leq u_j, \quad j = 1, \dots, r, \end{aligned}$$

which has a mix of linear equality and some interval constraints.

In the Supplementary Material, we outline a row-action algorithm for this class of optimization problems that use iterative projections (with corrections to maintain dual-feasibility) onto these interval constraints (Censor & Zenios, 1988). We briefly outline this algorithm below for completeness, but the reader could skip to the next section, where we describe these updates for convex entropic approximations of the partition function.

It will be useful to define the following notation: Given a convex function  $f$ , an iterate,  $x$ , and a linear constraint  $h \equiv \langle x, a \rangle = b$ , suppose we project  $x$  onto the hyperplane defining the equality constraint, under the Bregman divergence induced by  $f$ . This can be rewritten quite simply as computing  $y$  such that:

$$\begin{aligned} \nabla f(y) &= \nabla f(x) - \theta a, \\ \langle y, a \rangle &= b. \end{aligned}$$

We are interested in the value of  $\theta$  above; let us denote this by  $\Pi_D(f; x, h)$ . We can now detail the row-action algorithm:

Let  $\Phi$  denote the  $(m+r) \times p$  matrix with rows as  $\{a_j\}_{j=1}^m$  stacked above  $\{c_j\}_{j=1}^r$ .

Initialization:  $(\mu^0, z^0)$  such that  $\nabla B^*(\mu^0) = -\Phi^T z^0$ .

Iterative Step: Given  $\mu^t$  and  $z^t$ , and the current constraint  $h_j$  corresponding to the  $j$ -th row of  $\Phi$ , calculate the next primal and dual iterates  $\mu^{t+1}$  and  $z^{t+1}$  as

$$\begin{aligned}\nabla B^*(\mu^{t+1}) &= \nabla B^*(\mu^t) + d^t \Phi_j, \\ z^{t+1} &= z^t - d^t e_j,\end{aligned}$$

where if the constraint  $h_j \equiv \langle a_j, \mu \rangle = b_j$  is a linear equality, then  $d^t = \Pi_D(B^*; \mu^t, h_j)$ ; and if the constraint  $h_j \equiv l_j \leq \langle c_j, \mu \rangle \leq u_j$  is an interval constraint, then denoting  $h_j^- \equiv \langle c_j, \mu \rangle = l_j$  and  $h_j^+ \equiv \langle c_j, \mu \rangle = u_j$ , then  $d^t = \text{median}(u^t, A_t, B_t)$ , where  $A_t = \Pi_D(B^*; \mu^t, h_j^-)$  and  $B_t = \Pi_D(B^*; \mu^t, h_j^+)$ .

#### 4.2.1. WEIGHTED ENTROPIC APPROXIMATION

Recall from Section 2.2 the general weighted free-energy approximation consisting of a weighted sum of negative entropies

$$B_\alpha^*(\mu) = \sum_{s \in V} \alpha_s H_s(\mu_s) + \sum_{(s,t) \in E} \alpha_{st} H_{st}(\mu_{st}), \quad (12)$$

where  $H_s := \sum_{x_s \in \mathcal{X}} \mu_s(x_s) \log \mu_s(x_s)$  and  $H_{st}(\mu_{st}) := \sum_{(x_s, x_t) \in \mathcal{X} \times \mathcal{X}} \mu_{st}(x_s, x_t) \log \mu_{st}(x_s, x_t)$  denote the node-based and edge-based negative entropies respectively. Now, consider the variational approximation given by,

$$B_\alpha(\theta) = \sup_{\mu \in L_G} \langle \theta, \mu \rangle - B_\alpha^*(\mu),$$

where  $L_G$  is the marginal polytope outer bound in (5). The optimization problem in (11) then becomes:

$$\begin{aligned}\inf_{\mu} \left\{ \sum_s \alpha_s \sum_{x_s} \mu_s(x_s) \log \mu_s(x_s) \right. \\ \left. + \sum_{st} \alpha_{st} \sum_{x_s, x_t} \mu_{st}(x_s, x_t) \log \mu_{st}(x_s, x_t) \right\} \\ \text{s.t. } \mu_s(x_s) = \widehat{\phi}_s(x_s), \mu_{st}(x_s, x_t) \geq 0, \\ \sum_{x_t} \mu_{st}(x_s, x_t) = \widehat{\phi}_s(x_s). \\ \widehat{\phi}_{st}(x_s, x_t) - \lambda \leq \mu_{st}(x_s, x_t) \leq \widehat{\phi}_{st}(x_s, x_t) + \lambda.\end{aligned}$$

Thus, we have a set of equality marginalization constraints, and some box constraints. We can then apply the algorithm template above in the previous section to get the set of iterative updates in Algorithm 2.

#### 4.3. Closed Form Updates

The stationary condition characterizing the solution of (9) is given by

$$-\widehat{\phi} + \widehat{\mu}(\widehat{\theta}) + \lambda \widehat{Z}(\widehat{\theta}) = 0, \quad (14)$$

where  $\widehat{\mu} = \nabla B(\widehat{\theta})$ , and  $\widehat{Z}$  is a subgradient vector, with  $\widehat{Z}_N = 0$ , and  $\widehat{Z}_E \in \partial \|\widehat{\theta}_E\|_1$ , so that

---

#### Algorithm 2 Weighted Entropic Message Passing

---

Initialization:

$$\mu_s(x_s) = \widehat{\phi}(x_s), \quad (13a)$$

$$\mu_{st}^{(0)}(x_s, x_t) = \widehat{\phi}(x_s, x_t), \quad (13b)$$

$$Z_{st}(x_s, x_t) = -\alpha_{st}(\log \widehat{\phi}(x_s, x_t) + 1). \quad (13c)$$

**repeat**

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

$$\mu_{st}^{(t+1)}(x_s, x_t) = \mu_{st}^{(t)}(x_s, x_t) \left( \frac{\widehat{\phi}_s(x_s)}{\sum_{x_t} \mu_{st}^{(t)}(x_s, x_t)} \right),$$

$$\mu_{st}^{(t+1)}(x_s, x_t) = \mu_{st}^{(t+1)}(x_s, x_t) \left( \frac{\widehat{\phi}_t(x_t)}{\sum_{x_s} \mu_{st}^{(t+1)}(x_s, x_t)} \right).$$

$$\Delta_+ = \alpha_{st} \log \frac{(\widehat{\phi}_{st}(x_s, x_t) + \lambda)}{\mu_{st}^{(t+1)}(x_s, x_t)},$$

$$\Delta_- = \alpha_{st} \log \frac{(\widehat{\phi}_{st}(x_s, x_t) - \lambda)}{\mu_{st}^{(t+1)}(x_s, x_t)}.$$

$$C_{st}(x_s, x_t) = \text{median}(Z_{st}(x_s, x_t), \Delta_+, \Delta_-).$$

$$Z_{st}(x_s, x_t) = Z_{st}(x_s, x_t) - C_{st}(x_s, x_t),$$

$$\mu_{st}^{(t+1)}(x_s, x_t) = \mu_{st}^{(t+1)}(x_s, x_t) \exp(C_{st}(x_s, x_t)/\alpha_{st}).$$

**end for**

**until** convergence

---

(a) if  $\widehat{\theta}_{st;jk} \neq 0$ , then  $\widehat{Z}_{st;jk} = \text{sign}(\widehat{\theta}_{st;jk})$ ,

(b) if  $\widehat{\theta}_{st;jk} = 0$ , then  $|\widehat{Z}_{st;jk}| \leq 1$ ,

(c)  $\widehat{Z}_{s;j} = 0$ . That is,  $\widehat{\mu}_{s;j} = \widehat{\phi}_{s;j}$ .

Now suppose we use the tree-reweighted entropy approximation (6) as the variational approximation to the log-partition function. An interesting property satisfied by the pseudo-moments  $\widehat{\mu}$  of this approximation is a reparameterization condition (Wainwright et al., 2003), so that tuple  $(\widehat{\theta}, \widehat{\mu})$  is a valid primal dual pair iff they satisfy

$$\widehat{\theta}_{s;j} = \log \widehat{\mu}_{s;j} + C_{sj} + C_s \quad (15)$$

$$\widehat{\theta}_{st;jk} = \rho_{st} \log \frac{\widehat{\mu}_{st;jk}}{\widehat{\mu}_{s;j} \cdot \widehat{\mu}_{t;k}} - C_{sj} - C_{tk}, \quad (16)$$

for some constants  $\{C_s, C_{sj}\}$ , and further that  $\widehat{\mu}$  lies in the pseudomarginal polytope  $L_G$  (5), so that

$$\sum_k \widehat{\mu}_{st;jk} = \widehat{\mu}_{s;j}, \widehat{\mu}_{st;jk} \geq 0. \quad (17)$$

Now consider the three cases of the sign of any edge parameter  $\widehat{\theta}_{st;jk}$ :

- (a)  $\widehat{\theta}_{st;jk} > 0$ : Then from (14), we get  $\widehat{\mu}_{st;jk} = \widehat{\phi}_{st;jk} - \lambda$ , and substituting this in (16), we get  $\widehat{\phi}_{st;jk} > \widehat{\phi}_{sj} \exp(C_{sj}) \widehat{\phi}_{tk} \exp(C_{tk}) + \lambda$ .
- (b) Similarly, if  $\widehat{\theta}_{st;jk} < 0$ , then  $\widehat{\mu}_{st;jk} = \widehat{\phi}_{st;jk} + \lambda$ , which entails that:  $\widehat{\phi}_{st;jk} < \widehat{\phi}_{sj} \exp(C_{sj}) \widehat{\phi}_{tk} \exp(C_{tk}) - \lambda$ .
- (c) Finally, if  $\widehat{\theta}_{st;jk} = 0$ , then  $|\widehat{\phi}_{st;jk} - \widehat{\phi}_{sj} \exp(C_{sj}) \widehat{\phi}_{tk} \exp(C_{tk})| < \lambda$ , and where  $\widehat{\mu}_{st;jk} = \widehat{\phi}_{st;jk} \exp(C_{sj}) \widehat{\phi}_{tk} \exp(C_{tk})$ .

#### 4.3.1. EXPLICIT CONSTRUCTION

Given any constants  $\{C_s, C_{sj}\}$ , suppose we write out a tuple  $(\widehat{\mu}, \widehat{\theta}, \widehat{Z})$  as follows:

- (a) If  $\widehat{\phi}_{st;jk} > \widehat{\phi}_{sj} \exp(C_{sj}) \widehat{\phi}_{tk} \exp(C_{tk}) + \lambda$ , then  $\widehat{\mu}_{st;jk} = \widehat{\phi}_{st;jk} - \lambda$ ,  $\widehat{Z}_{st;jk} = 1$ ;
- (b) If  $\widehat{\phi}_{st;jk} < \widehat{\phi}_{sj} \exp(C_{sj}) \widehat{\phi}_{tk} \exp(C_{tk}) - \lambda$ , then  $\widehat{\mu}_{st;jk} = \widehat{\phi}_{st;jk} + \lambda$ ,  $\widehat{Z}_{st;jk} = -1$ ;
- (c) Otherwise,  $\widehat{\mu}_{st;jk} = \widehat{\phi}_{sj} \exp(C_{sj}) \widehat{\phi}_{tk} \exp(C_{tk})$ , and  $\widehat{Z}_{st;jk} = \widehat{\phi}_{st;jk} - \widehat{\phi}_{sj} \exp(C_{sj}) \widehat{\phi}_{tk} \exp(C_{tk})$ .
- (d) Set  $\widehat{\theta}$  from (15), (16).

Note that the tuple  $(\widehat{\theta}, \widehat{\mu})$  satisfies the stationary condition (14) by construction. The subgradient condition  $\widehat{Z} \in \partial \|\widehat{\theta}\|_1$  also holds by construction. Thus, this is a valid tuple and  $\widehat{\theta}$  is the solution of (14) provided the resulting  $\widehat{\mu} \in L$ . The goal then is to derive constants  $\{C_s, C_{sj}\}$  so that the  $\widehat{\mu}$  constructed as above lies in the polytope  $L_G$ .

While this is difficult to do in general, we show that when the variables are binary, the appropriate constants  $\{C_s, C_{sj}\}$  can be written out explicitly.

**Proposition 3.** *When the variables are binary, the solution  $\widehat{\theta}$  of (9) can be computed in closed form. Specifically, setting  $\{C_s = 0, C_{sj} = 0\}$  in the construction above yields solution  $\widehat{\theta}$ , with pseudomarginal  $\widehat{\mu}$ .*

The proof consists of showing that with the constants set to zero, the pseudomarginal  $\widehat{\mu}$  resulting from the construction satisfies  $\widehat{\mu} \in L$ ; and is provided in detail in the Supplementary Material.

As noted in Proposition 2, even though we have an overcomplete parameterization, the resulting solution has the Ising form, and in particular is equivalent to solving a corresponding  $\ell_1$ -regularized approximate MLE in (9) for the Ising model. Further, the solution

$\widehat{\theta}$  is available in closed-form, which thus leads to this very simple estimator for the structure of graph:

$$\widehat{E} = \{(s, t) : \exists j, k; |\widehat{\phi}_{st;jk} - \widehat{\phi}_{s;j} \cdot \widehat{\phi}_{t;k}| > \lambda\}. \quad (18)$$

Similar algorithms based on correlations have been proposed elsewhere, see Montanari & Pereira (2009) for instance for the case of a homogeneous Ising model, where the edge parameters are all equal. Here, we obtain the very interesting connection between the specific correlation based edge-detection algorithm above and the tree-reweighted entropy based approximate MLE. This thus opens up new avenues for analyzing such methods, see for instance our sparsistency analysis in Section 5.

*Discussion.* Note that the closed form estimator in (18) requires time scaling as  $O(p^2)$ , whereas other state of the art methods such as (Ravikumar et al., 2010) require time that could scale as  $O(p^5)$ . As we will see in the experiments section in spite of this, their graph structure recovery performance is nonetheless comparable.

## 5. Sparsistency

In this section, we show that any structure learning estimator obtained as a solution of (9) given any partition function approximation  $B(\theta)$  is *sparsistent* under certain conditions. While the analysis builds on standard tools such as the dual-witness technique from (Wainwright, 2009), here we face multiple subtleties: the objective function in (9) does not arise from the likelihood of the data, which has the consequence that the gradient of the objective need not be small at the true parameter. Moreover, the number of non-zero elements here equal the number of edges which scale at least linearly with the number of nodes; so that we needed additional tools such as Brouwers fixed point theorem (Ortega & Rheinboldt, 2000). Lastly, the sparsistency theorem holds not just for one, but a whole family of estimators obtained as solutions of (9).

Let us first study the stationary condition characterizing the solution of (9):

$$-\widehat{\phi} + \widehat{\mu} + \lambda \widehat{Z} = 0. \quad (19)$$

We will now introduce some notation to simplify this condition. Let  $\mu^* = \nabla A(\theta^*)$  be the true marginals and let  $\bar{\mu} = \nabla B(\theta^*)$  be the “true” variational pseudomarginals. Denote  $W_1 = \widehat{\phi} - \mu^*$  and  $W_2 = \mu^* - \bar{\mu}$ , and further that  $W = W_1 + W_2$ . Then, the stationary

condition can be rewritten as,

$$\hat{\mu} - \bar{\mu} - W + \lambda Z = 0. \quad (20)$$

We define the second-order Taylor's expansion remainder of  $\nabla B$  around  $\theta^*$  as,

$$R(\Delta; \theta^*) = \nabla B(\theta^* + \Delta) - \nabla B(\theta^*) - \nabla^2 B(\theta^*)\Delta.$$

Denoting  $\hat{\Delta} = \hat{\theta} - \theta^*$ , and noting that  $\hat{\mu} = \nabla B(\hat{\theta})$  and  $\bar{\mu} = \nabla B(\theta^*)$  we can rewrite (20) as

$$\nabla^2 B(\theta^*)\hat{\Delta} + R(\hat{\Delta}) - W + \lambda Z = 0. \quad (21)$$

We can now state our assumptions. We will assume that the difference between the true marginals  $\mu^* = \nabla A(\theta^*)$  and the ‘‘true’’ variational pseudomarginals  $\bar{\mu} = \nabla B(\theta^*)$  is bounded so that

**Assumption 1.**  $\kappa_B := \|\nabla B(\theta^*) - \nabla A(\theta^*)\|_\infty < 1$ .

Thus,  $\|W_2\|_\infty = \kappa_B$  is controlled.

We will also assume a mild regularity condition on the remainder:

**Assumption 2.** For all sparse and bounded  $\Delta$  where  $\|\Delta\|_\infty \leq \lambda$ , and  $\Delta_{S^c} = 0$ , the second-order remainder is bounded by,

$$\|R(\Delta; \theta^*)\|_\infty \leq \kappa_R \|\Delta\|_\infty^2, \quad (22)$$

for some  $\kappa_R > 0$ .

We note that such conditions have been considered in other analyses of  $\ell_1$ -regularized non-linear objectives. For instance, in the case where the objective is the log-likelihood of a Gaussian distribution with true covariance matrix  $\Sigma^*$ , [Ravikumar et al. \(2008\)](#) showed that the second-order remainder is bounded as in the assumption for  $\kappa_R = 3d/2 \|\Sigma^*\|_\infty^3$ . Note that since both  $\theta^*$  and  $\Delta$  have support restricted to  $S$ , and  $\kappa_R$  only depends on the support size and is independent of the ambient dimension.

Lastly, let  $Q = \nabla^2 B(\theta^*)$  denote the Hessian of  $B(\theta)$  at the true parameters  $\theta^*$ . Let  $S$  denote the support of the true parameters  $\theta^*$  and let  $S^c$  denote its complement. Thus,  $S = V \cup \{(s, t) : \theta_{st}^* \neq 0\}$ , and  $S^c = \{(s, t) : \theta_{st}^* = 0\}$ . Let  $Q_{AB}$  denote the submatrix of  $Q$  with rows indexed by  $A$  and columns indexed by  $B$ . We then assume the following incoherence assumption:

**Assumption 3.**  $\|Q_{SS}^{-1} Q_{S^c S}\|_\infty \leq 1 - \alpha$ , for some constant  $\alpha > 0$ .

**Theorem 1.** Consider a graphical model distribution with parameters  $\theta^*$  satisfying assumptions 1,2,3. Suppose we solve (9) by setting the regularization parameter  $\lambda$  as  $\lambda \geq \frac{4}{\alpha} \left( \sqrt{\frac{\log p}{n}} + \kappa_B \right)$ , and where the sample

size  $n$  scales as  $n \geq (\alpha^2 / (32\kappa_Q^2 \kappa_R) - \kappa_B)^2 \log p$ . Then with probability greater than  $1 - \exp(-c \log p)$  for some constant  $c > 0$ , we have:

- (a) the estimate  $\hat{\theta}$  from (9) satisfies the elementwise  $\ell_\infty$  bound:  $\|\hat{\theta} - \theta^*\| \leq \min\{1/(2\kappa_Q \kappa_R), 4\kappa_Q \lambda\}$ ,
- (b) it specifies an edgeset  $E(\hat{\theta})$  that has no false inclusions (i.e.  $E(\hat{\theta}) \subseteq E(\theta^*)$ ), and moreover includes all edges  $(s, t)$  such that,  $|\theta_{st}^*| > \min\{1/(2\kappa_Q \kappa_R), 4\kappa_Q \lambda\}$ .

The detailed proof is provided in the Supplementary Material.

## 6. Experiments

We now briefly illustrate our results on 25 node Ising models (2). Further details on our experimental settings, as well as more exhaustive simulations are provided in the appendix. Figure 1 (a) compares our message-passing updates (Algorithm 2) to gradient descent for tree-reweighted entropic approximation, on Ising models with four-nearest neighbor lattice graph-structure. Here, we plot the  $\ell_2$  deviation of the pseudomarginals to the optimum against iterations: it can be seen that our message-passing updates (Algorithm 2) converge very fast. Figure 1(b) plots the convergence of our message-passing updates for a more general weighted entropic approximation. Figure 1 shows the edge-recovery rate of the TRW-approx estimator that we have available in closed form. It has comparable performance to the state of the art method ([Ravikumar et al., 2010](#)) that uses nodewise  $\ell_1$  regularized logistic regressions, which is impressive considering that our estimator in this case has a simple closed-form solution. In Figure 1(d) we follow ([Wainwright, 2006](#)). We compare two parameter estimators: our TRW-approx estimator, and an oracle estimator that knows the true graph structure and estimates the parameters using the tree-reweighted entropy approximation to the log-partition function (called pseudo-moment matching in [Wainwright et al. \(2003\)](#)). We then plot the  $\ell_2$  error in moment estimates after perturbing this parameter estimate. This is a surrogate metric for gauging the use of the estimated model for prediction; we see that our estimator is very close to the oracle estimator.

**Summary.** We investigate a whole class of estimators that recover graphical model structure by minimizing  $\ell_1$ -regularized surrogate log-likelihoods based on variational approximations to the partition function. As we note in the introduction, many state of the art methods fall into

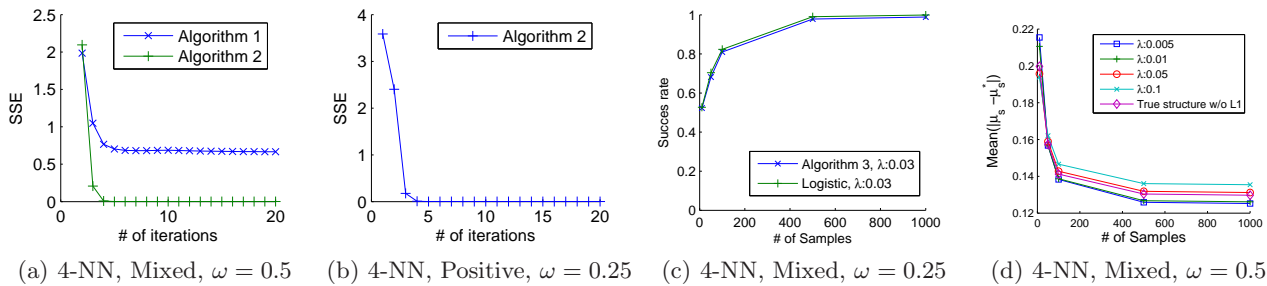


Figure 1. The four panels illustrate the following four experiments:

- (a) Convergence Rate: Gradient Descent (Alg. 1) vs Our Mesg. Passing updates (Alg. 2) for tree-reweighted entropic approx.  
 (b) Our Mesg. Passing updates (Alg. 2) for weighted Entropic approximation with  $\alpha_s = \alpha_{st} = 1$ .  
 (c) Structure Recovery: TRW-approx estimator (in closed form, Alg. 3) vs Nodewise logistic regressions of (Ravikumar et al., 2010).  
 (d) Prediction (Error in moments after perturbing parameter estimate): TRW-approx estimator vs Oracle parameter estimator that knows the true graph structure.

this category. For this general setting, we provide (a) a general message passing algorithm for *directly* solving the resulting  $\ell_1$  regularized optimization problems (in contrast to iterative calls to a separate approximate inference procedure), and (b) sparsistency results for this entire class of estimators. Our study also revealed that in special cases, the resulting estimator is available in closed form.

## References

- Boyd, S. and Vandenberghe, L. *Convex optimization*. Cambridge University Press, Cambridge, UK, 2004.
- Candes, E. and Tao, T. The Dantzig selector: Statistical estimation when  $p$  is much larger than  $n$ . *Annals of Statistics*, 35(6):2313–2351, 2007.
- Censor, Y. and Zenios, S. A. *Parallel Optimization: Theory, Algorithms, and Applications*. Numerical Mathematics and Scientific Computation. Oxford University Press, 1988.
- Heskes, T. Convexity arguments for efficient minimization of the bethe and kikuchi free energies. *J. Artif. Intell. Res. (JAIR)*, 26:153–190, 2006.
- Ising, E. Beitrag zur theorie der ferromagnetismus. *Zeitschrift für Physik*, 31:253–258, 1925.
- Lee, S.-I., Ganapathi, V., and Koller, D. Efficient structure learning of markov networks using  $\ell_1$ -regularization. In *NIPS 19*, 2007.
- Montanari, A. and Pereira, J. A. Which graphical models are difficult to learn. In *Neural Information Processing Systems (NIPS)*, 2009.
- Nesterov, Y. Introductory lectures on convex optimization: a basic course. 2004.
- Ortega, J. M. and Rheinboldt, W. C. *Iterative solution of nonlinear equations in several variables*. Classics in applied mathematics. SIAM, New York, 2000.
- Pearl, J. *Probabilistic reasoning in intelligent systems*. Morgan Kaufman, San Mateo, 1988.
- Ravikumar, P., Wainwright, M. J., Raskutti, G., and Yu, B. High-dimensional covariance estimation by minimizing  $\ell_1$ -penalized log-determinant divergence. Technical Report 767, Dept. of Stat., UC Berkeley, 2008.
- Ravikumar, P., Wainwright, M. J., and Lafferty, J. High-dimensional ising model selection using  $\ell_1$ -regularized logistic regression. *Ann. Stat.*, 38(3):1287–1319, 2010.
- Tropp, J. A. Just relax: Convex programming methods for identifying sparse signals. *IEEE Trans. Info. Theory*, 51(3):1030–1051, March 2006.
- Wainwright, M. J. Estimating the “wrong” graphical model: Benefits in the computation-limited regime. *JMLR*, 7:1829–1859, September 2006.
- Wainwright, M. J. Sharp thresholds for noisy and high-dimensional recovery of sparsity using  $\ell_1$ -constrained quadratic programming (lasso). *IEEE Transactions on Information Theory*, 55:2183–2202, 2009.
- Wainwright, M. J. and Jordan, M. I. Graphical models, exponential families, and variational inference. *Foundations and Trends in Machine Learning*, 1(1–2):1–305, 2008.
- Wainwright, M. J., Jaakkola, T. S., and Willsky, A. S. Tree-reweighted belief propagation algorithms and approximate ML estimation by pseudomoment matching. In *AISTATS*, 2003.
- Wainwright, M. J., Jaakkola, T. S., and Willsky, A. S. A new class of upper bounds on the log partition function. *IEEE Trans. Info. Theory*, 51(7):2313–2335, July 2005.
- Weiss, Y., Yanover, C., and Meltzer, T. MAP estimation, linear programming, and belief propagation with convex free energies. In *UAI*, 2007.
- Yedidia, J. S., Freeman, W. T., and Weiss, Y. Generalized belief propagation. In *NIPS 13*, pp. 689–695, 2001.



## 7. Proofs of Lemmas and Propositions

**Proposition 4.** *When the variables are binary, the solution  $\hat{\theta}$  of (9) can be computed in closed form. Specifically, setting  $\{C_s = 0, C_{sj} = 0\}$  in the construction above yields the solution  $\hat{\theta}$  and its pseudomarginal  $\hat{\mu}$ .*

The proof consists of showing that with the constants set to zero, the pseudomarginal  $\hat{\mu}$  resulting from the construction satisfies  $\hat{\mu} \in L$ ; and is provided in detail in the Supplementary Material.

*Proof.* As noted in the discussion above, we only need guarantee that with the constants set to zero, the pseudomarginal  $\hat{\mu}$  resulting from the construction satisfies  $\hat{\mu} \in L$ .

*Node Marginals.* The node marginals  $\{\hat{\mu}_{s;j}\}$  naturally satisfy the normalization constraint in  $L$  since they are set as  $\hat{\mu}_{s;j} = \hat{\phi}_{s;j}$ .

*Edge Marginals.* From the construction above, there are three possible cases for  $\hat{\mu}_{st;jk}$ . Suppose  $\hat{\mu}_{st;jk}$  falls under the first case, which means  $\hat{\mu}_{st;jk} = \hat{\phi}_{st;jk} - \lambda$  so that

$$\hat{\phi}_{st;jk} > \hat{\phi}_{s;j} \hat{\phi}_{t;k} + \lambda.$$

Then noting that  $\hat{\phi}_{t;\bar{k}} = \hat{\phi}_{s;j} - \hat{\phi}_{st;jk}$ , and  $\hat{\phi}_{t;\bar{k}} = 1 - \hat{\phi}_{t;k}$ , we get,

$$\hat{\phi}_{st;j\bar{k}} < \hat{\phi}_{s;j} \hat{\phi}_{t;\bar{k}} - \lambda,$$

and thus  $\hat{\mu}_{st;j\bar{k}}$  in turn satisfies  $\hat{\mu}_{st;jk} = \hat{\phi}_{st;jk} + \lambda$ . Thus, from the construction,

$$\begin{aligned} \sum_k \hat{\mu}_{st;jk} &= \sum_k \hat{\phi}_{st;jk} + \lambda - \lambda \\ &= \hat{\phi}_{s;j} = \hat{\mu}_{s;j}, \end{aligned}$$

so that the marginalization condition is satisfied. Moreover,

$$\begin{aligned} \hat{\mu}_{st;jk} &= \hat{\phi}_{st;jk} - \lambda \\ &> \hat{\phi}_{s;j} \hat{\phi}_{t;k} \geq 0, \end{aligned}$$

so that the non-negativity condition is satisfied as well. The other two cases are similar.  $\square$

## 8. Proof of Sparsistency

We will follow the dual-witness technique from (Wainwright, 2009). Here we face the subtlety that

the objective is not the likelihood, and moreover that the number of non-zero elements equal the number of edges which scale atleast linearly with the number of nodes. Some additional tools we use include using Brouwers fixed point theorem (Ortega & Rheinboldt, 2000) to obtain an  $\ell_\infty$  bound on the deviation ( $\hat{\theta} - \theta^*$ ). The detailed proof is provided in the Supplementary Material.

We begin by unpacking the stationary condition in (19) via the following lemma:

**Lemma 1** (Optimality Conditions). *Any optimal primal-dual pair  $(\hat{\theta}, \hat{Z})$  of (9) satisfies*

1. (Stationary Condition).  $-\hat{\phi} + \hat{\mu} + \lambda \hat{Z} = 0$ .
2. (Dual Feasibility).  $\hat{Z}$  is equal to the subgradient  $\partial \|\hat{\theta}\|_{1,E}$  so that
  - (a) if  $\hat{\theta}_u \neq 0$  then  $\hat{Z}_u = \text{sign}(\hat{\theta}_u)$ ,
  - (b) if  $\hat{\theta}_u = 0$  then  $|\hat{Z}_u| \leq 1$ .

The next lemma states that structure recovery is guaranteed if the dual is *strictly* feasible.

**Lemma 2** (Strict Dual Feasibility). *Suppose that there exists an optimal primal-dual pair  $(\hat{\theta}, \hat{Z})$  for (9) such that  $\|\hat{Z}_{S^c}\|_\infty < 1$ . Then, any optimal primal solution  $\hat{\theta}$  satisfies  $\hat{\theta}_S^c = 0$ . Moreover, if the Hessian sub-matrix  $Q_{SS} \succ 0$  then  $\hat{\theta}$  is the unique optimal solution.*

We are now ready to sketch the proof of Theorem 1.

*Part (a).* The proof proceeds by a primal-dual witness technique, and consists of the construction of a feasible primal-dual pair in the following two steps:

- (i) Primal Candidate using an oracle subproblem: Let  $\hat{\theta}$  be the optimal solution of the restricted problem:

$$\hat{\theta} = \arg \min_{\{\theta: \theta_{S^c} = 0\}} -\langle \theta, \hat{\phi} \rangle + B(\theta) + \lambda \|\theta\|_{1,E}. \quad (23)$$

- (ii) Dual Candidate from Stationary Optimality Condition: Set  $\hat{Z}_S = \text{sign}(\hat{\theta}_S)$ . Set  $\hat{Z}_{S^c}$  from the stationary condition (19).

It only remains to show strict dual feasibility. By construction, the  $(\hat{\theta}, \hat{Z})$  pair satisfies the stationary condition (19). It remains to show that the the dual  $\hat{Z}$  is strictly feasible. We show that this holds, and also that the solution is unique in the sequel.

We have:

$$Q_{SS}\Delta_S + R_S - W_S + \lambda Z_S = 0,$$

from which we have:

$$\Delta_S = -(Q_{SS})^{-1}(R_S - W_S + \lambda Z_S).$$

From (19) we have:

$$\begin{aligned} \lambda Z_{S^c} &= Q_{S^cS}\Delta_S - R_{S^c} + W_{S^c} \\ &= Q_{S^cS}(Q_{SS})^{-1}(R_S - W_S + \lambda Z_S) - R_{S^c} + W_{S^c}, \end{aligned}$$

so that

$$\begin{aligned} \lambda \|Z_{S^c}\|_\infty &\leq \|Q_{S^cS}(Q_{SS})^{-1}\|_\infty \|R_S - W_S + \lambda Z_S\|_\infty \\ &\quad + \|R_{S^c} + W_{S^c}\|_\infty \\ &\leq (2 - \alpha)(\|R\|_\infty + \|W\|_\infty) + (1 - \alpha)\lambda. \end{aligned} \quad (24)$$

Now from Lemma 3 and Assumption 1,

$$\|W\|_\infty \leq 2\sqrt{\frac{\log p}{n}} + \kappa_B \leq \frac{\alpha\lambda}{2}, \quad (25)$$

from the setting of  $\lambda$  specified in the theorem. Further, from Assumption 2,

$$\begin{aligned} \|R\|_\infty &\leq \kappa_R \|\delta\|_\infty^2 \\ &\leq \frac{8\kappa_Q^2 \kappa_R}{\alpha} \lambda \cdot \frac{\alpha\lambda}{2} \\ &\leq \frac{\alpha\lambda}{2}, \end{aligned} \quad (26)$$

where the second inequality uses the elementwise  $\ell_\infty$  bound on  $\Delta$  from Lemma 28, and the last inequality uses the scaling of the sample size from the theorem. Plugging (25) and (26) into (24), we get that

$$\begin{aligned} \lambda \|Z_{S^c}\|_\infty &< \frac{\alpha\lambda}{2} + \frac{\alpha\lambda}{2} + (1 - \alpha)\lambda \\ &= \lambda, \end{aligned}$$

as required.

*Part (b).* From Lemma 4,

$$\|\tilde{\theta}_S - \theta^*\|_\infty \leq r,$$

so that  $\tilde{\theta}_{st}$  and  $\theta_{st}^*$  would share the same sign provided that  $|\theta_{st}^*| > r$ , as stated in the theorem.

This proves the result.

**Lemma 3** (Sample Noise).

$$\|\hat{\phi} - \mu^*\|_\infty \leq 2\sqrt{\frac{\log p}{n}}, \quad (27)$$

with probability at least  $1 - 2\exp(-2\log p)$ .

*Proof.* By Hoeffding's inequality, for any node or edge potential with index denoted by  $\alpha$

$$\mathbb{P}\{|\hat{\phi}_\alpha - \mu_\alpha^*| > t\} \leq 2\exp(-2nt^2).$$

Since there are at most  $\binom{p}{2} + p < p^2$  node and edge potentials, by a union bound,

$$\mathbb{P}\{\|\hat{\phi} - \mu^*\|_\infty > t\} \leq 2\exp(-2nt^2 + 2\log p),$$

from which the statement of the theorem follows.  $\square$

**Lemma 4.** Let  $r$  denote the value

$$r = \min\{1/(2\kappa_Q\kappa_R), 4\kappa_Q\lambda\}. \quad (28)$$

Then, we have the elementwise  $\ell_\infty$  bound:

$$\|\Delta\|_\infty \leq \|\tilde{\theta} - \theta^*\| \leq r. \quad (29)$$

*Proof.* Note that the restricted problem (23) has a unique optimum  $\tilde{\theta}$ . Let  $\tilde{Z}$  be any member of the sub-differential of  $\|\theta\|_{1,E}$ , evaluated at  $\tilde{\theta}$ . By Lagrangian theory, the witness  $\tilde{\theta}$  must be an optimum of the associated Lagrangian problem,

$$\min_{\theta: \theta_{S^c} = 0} \left\{ -\langle \theta, \hat{\phi} \rangle + B(\theta) + \lambda \langle \theta, \tilde{Z} \rangle \right\}$$

In fact, since this Lagrangian is strictly convex,  $\tilde{\theta}$  is the only optimum of this problem. We have the zero-gradient condition,

$$G(\theta_S) = -\hat{\phi}_S + \nabla_S B(\theta) + \lambda \tilde{Z}_S = 0.$$

Our goal is to bound the deviation of this solution from  $\theta_S^*$ , or equivalently to bound the deviation  $\Delta = \tilde{\theta} - \theta^*$ . Our strategy is to show the existence of a solution  $\Delta$  to the zero-gradient condition above, that is contained inside a ball  $\mathbb{B}(r)$ , defined below in equation (30). By uniqueness of the optimal solution, we can conclude that  $\tilde{\theta} - \theta^*$  lies in this ball.

Let us define a map,

$$F(\Delta_S) = -Q_{SS}^{-1}[G(\theta_S^* + \Delta_S)] + \Delta_S.$$

Consider the ball of radius  $r$  from (28),

$$\mathbb{B}(r) := \{\theta_S \mid \|\theta_S\|_\infty \leq r\}. \quad (30)$$

Suppose we show that

$$F(\mathbb{B}(r)) \subseteq \mathbb{B}(r). \quad (31)$$

Since  $F$  is continuous and  $\mathbb{B}(r)$  is convex and compact, this inclusion implies, by Brouwer's fixed point

theorem (Ortega & Rheinboldt, 2000), that there exists some fixed point  $\Delta \in \mathbb{B}(r)$ . By uniqueness of the zero gradient condition (and hence fixed points of  $F$ ), we can thereby conclude that  $\|\tilde{\theta}_S - \theta_S^*\|_\infty \leq r$ .

It thus remains to show (31).

$$\begin{aligned} F(\Delta_S) &= -Q_{SS}^{-1}[G(\theta_S^* + \Delta_S)] + \Delta_S \\ &= -Q_{SS}^{-1}(R(\Delta_S) + W_S + \lambda \tilde{Z}_S). \end{aligned}$$

Thus,

$$\begin{aligned} \|F(\Delta_S)\|_\infty &\leq \|Q_{SS}^{-1}\|_\infty (\|R(\Delta_S)\| + \|W_S\| + \lambda) \\ &\leq \kappa_Q \kappa_R r^2 + 2\kappa_Q \lambda, \end{aligned}$$

where we used  $\|R(\Delta_S)\| \leq \kappa_R r^2$  from Assumption 2, and  $\|W_S\| \leq \lambda$  from the assumption on  $\lambda$ . Thus, from the conditions in the theorem,

$$\|F(\Delta_S)\|_\infty \leq r/2 + r/2 = r,$$

as required.  $\square$

## 9. Message Passing Derivation

Given a function  $f$ , an iterate,  $x$ , and a linear constraint  $h \equiv \langle x, a \rangle = b$ , we define  $\Pi_D(f; x, h)$  as the value of  $\theta$  such that the following are satisfied:

$$\begin{aligned} \nabla f(y) &= \nabla f(x) - \theta a, \\ \langle y, a \rangle &= b. \end{aligned}$$

It can be shown that  $\Pi_D$  corresponds to the deviation in a corresponding dual coordinate given a projection of  $x$  onto the hyperplane defining the equality constraint.

Now, consider the optimization problem:

$$\begin{aligned} \inf_{\mu} B^*(\mu) \\ \text{s.t. } \langle a_i, \mu \rangle &= b_i, \quad i = 1, \dots, m. \\ l_j &\leq \langle c_j, \mu \rangle \leq u_j, \quad j = 1, \dots, r, \end{aligned}$$

which has a mix of linear equality and some interval constraints.

Then the following row-action algorithm can be shown to converge to its solution. Let  $\Phi$  denote the  $(m+r) \times p$  matrix whose  $j$ -th row corresponds to  $a_j$  for  $j \in [m]$ , and to  $c_{j-m}$  for  $j \in \{m+1, \dots, p\}$ .

Initialization:  $(\mu^0, z^0)$  such that

$$\nabla B^*(\mu^0) = -\Phi^T z^0.$$

Iterative Step: Given  $\mu^t$  and  $z^t$ , and the current constraint  $h_j$  corresponding to the  $j$ -th row of  $\Phi$ , calculate the next primal and dual iterates  $\mu^{t+1}$  and  $z^{t+1}$  as

$$\begin{aligned} \nabla B^*(\mu^{t+1}) &= \nabla B^*(\mu^t) + d^t \Phi_j, \\ z^{t+1} &= z^t - d^t e_j, \end{aligned}$$

where if the constraint  $h_j \equiv \langle a_j, \mu \rangle = b_j$  is a linear equality, then  $d^t = \Pi_D(B^*; \mu^t, h_j)$ ; and if the constraint  $h_j \equiv l_j \leq \langle c_j, \mu \rangle \leq u_j$  is an interval constraint, then denoting  $h_j^- \equiv \langle c_j, \mu \rangle = l_j$  and  $h_j^+ \equiv \langle c_j, \mu \rangle = u_j$ , then  $d^t = \text{median}(u^t, A_t, B_t)$ , where  $A_t = \Pi_D(B^*; \mu^t, h_j^-)$  and  $B_t = \Pi_D(B^*; \mu^t, h_j^+)$ .

Using these updates for convex entropic approximations of the partition function yields Algorithm 2.

## 10. Additional Experiments

For the true graph structure, three different topologies of graphs were selected: (a) four-nearest neighbor lattices, (b) eight-nearest neighbor lattices, and (c) star-shaped graph. For a given graph topology and edge strength  $\omega > 0$ , we examined two special classes of Ising model: (a) positive couplings, meaning that  $\theta_{st}^* = \omega$  for all true edges  $(s, t)$  and (b) mixed couplings, meaning that  $\theta_{st}^* = \pm\omega$  with equal probabilities.

To compute true marginal probabilities  $\mu^*$  as the ground truth, we restricted the number of nodes to 25. And, for all three algorithms, we use the tree-reweighted approximation (Wainwright et al., 2005) of true log partition function. All the results presented here are the average from 100 trials.

In the first experiment, we evaluate the efficiency of Algorithm 2 against Algorithm 1. Solutions of both algorithms will approach to the optimal point as the number of iterations increases. Figure 2 shows how fast each algorithm converges to the closed-form solution under the some environment: the number of samples=1000,  $\omega=0.25$  and  $\lambda=0.05$ . To evaluate the distance between the pseudo-moments approximation after the  $i$ -th iteration,  $\mu_{st}^{(i)}$  and the optimal pseudo-moments  $\hat{\mu}_{st}$  by the closed-form, we used the sum of squared errors over all edges. Note that learning rate  $\eta^t$  for Algorithm 1 is set to  $1/\sqrt{t}$ . If we increase it, convergence rates for some topologies will be faster while it will fail to converge for some topologies. Since two algorithms have different initial points -  $\theta^{(0)}$  and  $\mu^{(0)}$  respectively, accuracies after one iteration have large differences. However, for all the cases, Algorithm 2 using proposed message passing converges much faster to the optimal points as iterations increases. Moreover, Algorithm 2 can also be

applied to the reweighted entropic approximation of  $\alpha_s = \alpha_{st} = 1$  where conventional message passing algorithms do not work well. Even in this regime, proposed message passing algorithm converges within 5 iterations as shown in the Figure 4 .

In the second experiments, we compare our closed-form estimator in 18, henceforth referred to as Algorithm 3 with  $\ell_1$ -regularized logistic approach Ravikumar et al. (2010). Figure 5 shows the edge recovery rate in terms of the number of samples where the edge recovery rate is defined as  $\frac{1}{2} \left( \frac{\# \text{ of correct inclusion}}{\# \text{ of true edges}} + \frac{\# \text{ of correct exclusion}}{\# \text{ of true non-edge}} \right)$ .  $\lambda$  for the weight of  $\ell_1$  is selected by the oracle fashion from  $\{0.01, 0.03, 0.05, 0.07, 0.1, 0.12\}$ . We can see that algorithm 3 is comparable on the most cases even with its simple computation complexity.

In the third experiment, we consider joint the estimation/prediction procedure (Wainwright, 2006): Initially, we are given a set of i.i.d data samples of  $X$  and compute the estimator  $\hat{\theta}$  according to Section 4.3.1. Then, we model new noisy observation  $Y$  whose distribution is perturbed by random vector  $\delta$  from the original problem:  $\mathbb{P}_{\theta+\delta}(y)$ . Here, we considers the perturbation only at the single nodes. That is,  $\delta_{st} = 0$  for all edge  $(s, t)$  And we assumes that the perturbation at each node is independent:  $\delta_{s;j} \sim \mathcal{N}(0, 0.5^2)$ . Finally, we compute approximate marginals  $\hat{\mu}(\hat{\theta} + \delta)$  using tree reweighted approximation algorithm.

With same samples and parameters, two approaches are compared: i) Algorithm 3 without the knowledge of structure for different parameter  $\lambda$  and ii) Oracle approach with the assumption that it knows the structure a priori. To estimate  $\hat{\theta}$  given the structure, we use the pseudo-moment matching method (Wainwright et al., 2003). To evaluate how accurate two approaches are, we compute the true marginal probabilities by the junction tree algorithm. And then we get the mean value of  $|\hat{\mu}_s - \mu_s^*|$  over the nodes. Figure 6 represents the results as the functions of the number of sample size for different settings. Even though structure recovering by the entropy approximations with  $\ell_1$ -regularization do not recover perfectly every kind of MRFs, it is still good approximation robust to perturbations on the single nodes.

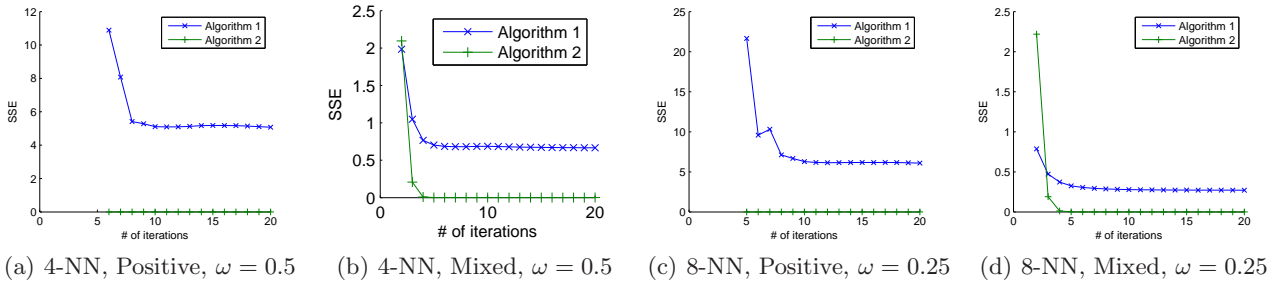


Figure 2. Converges rates using TRW approximation where algorithm 2 performs better

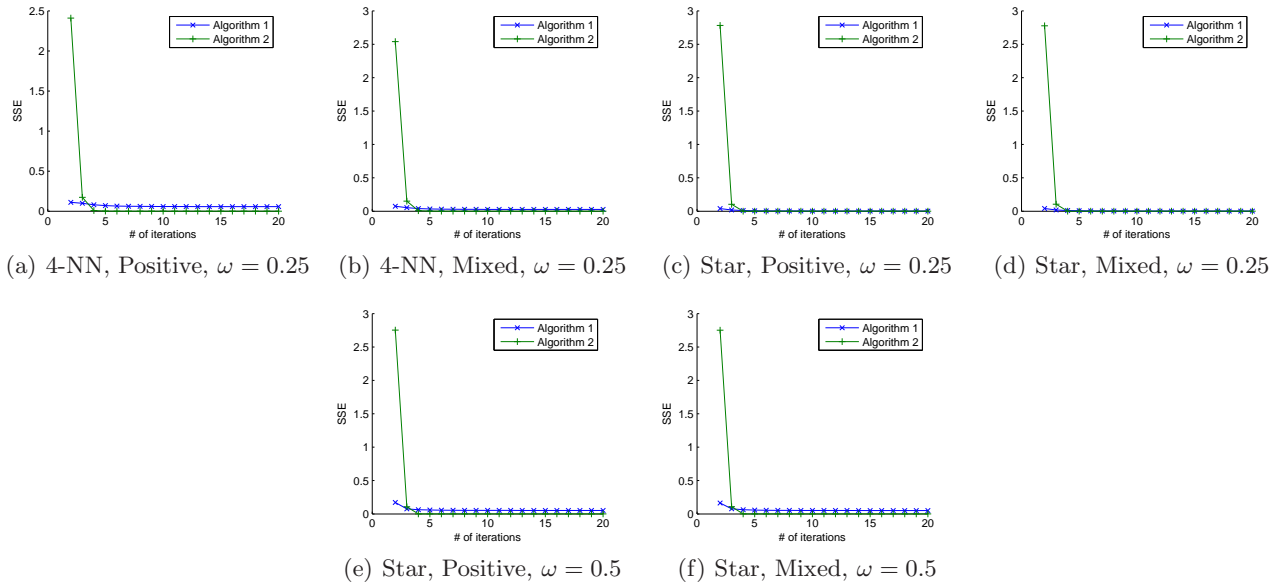


Figure 3. Converges rates using TRW approximation where both algorithms work well

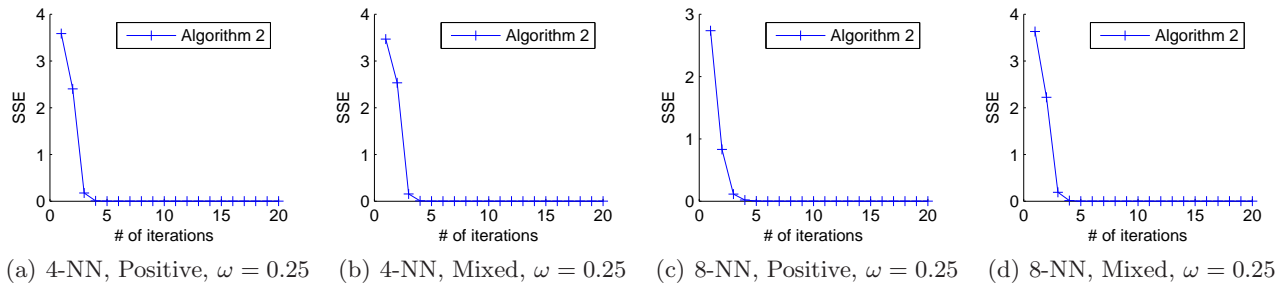


Figure 4. Converges rates for Weighted Entropic approximation where  $\alpha_s = \alpha_{st} = 1$

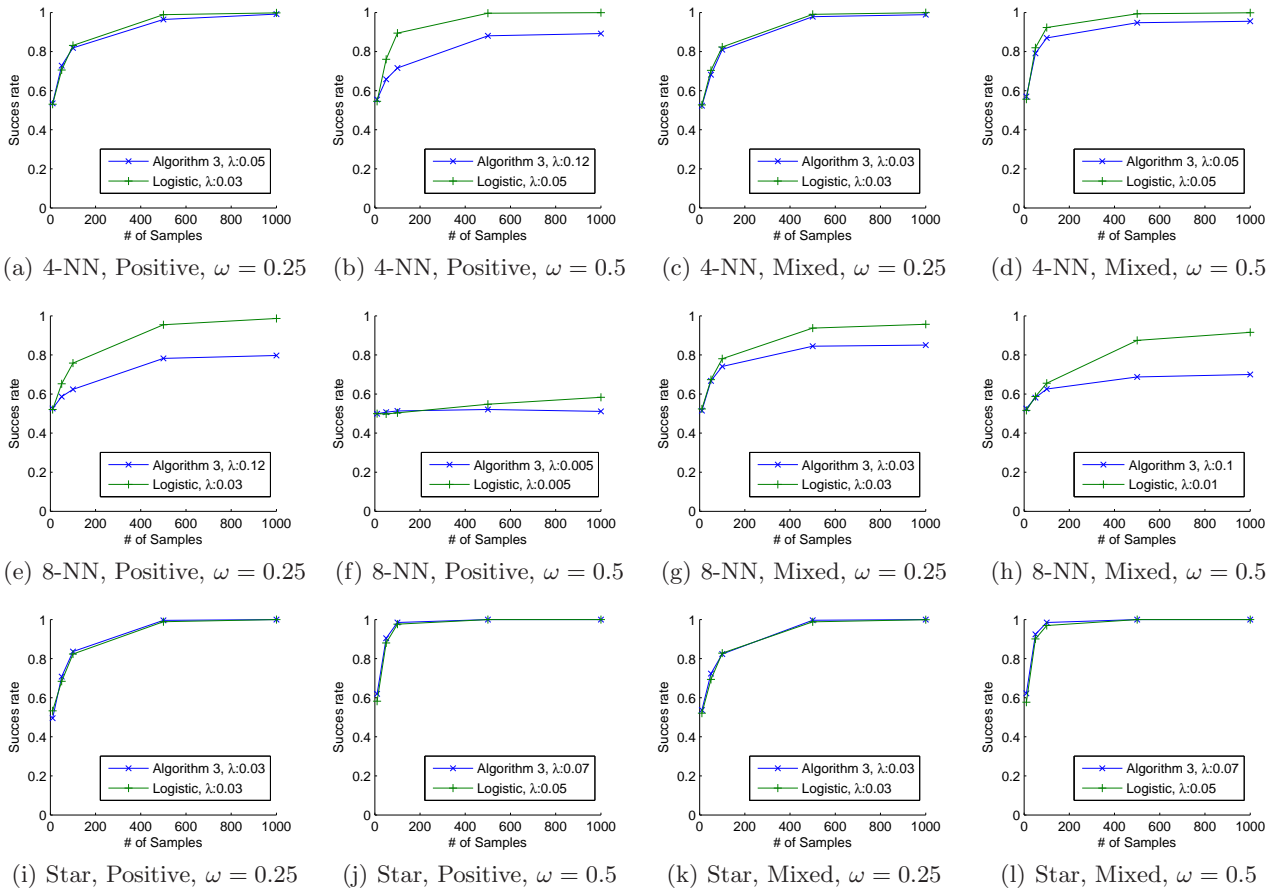


Figure 5. Edge recovery rate

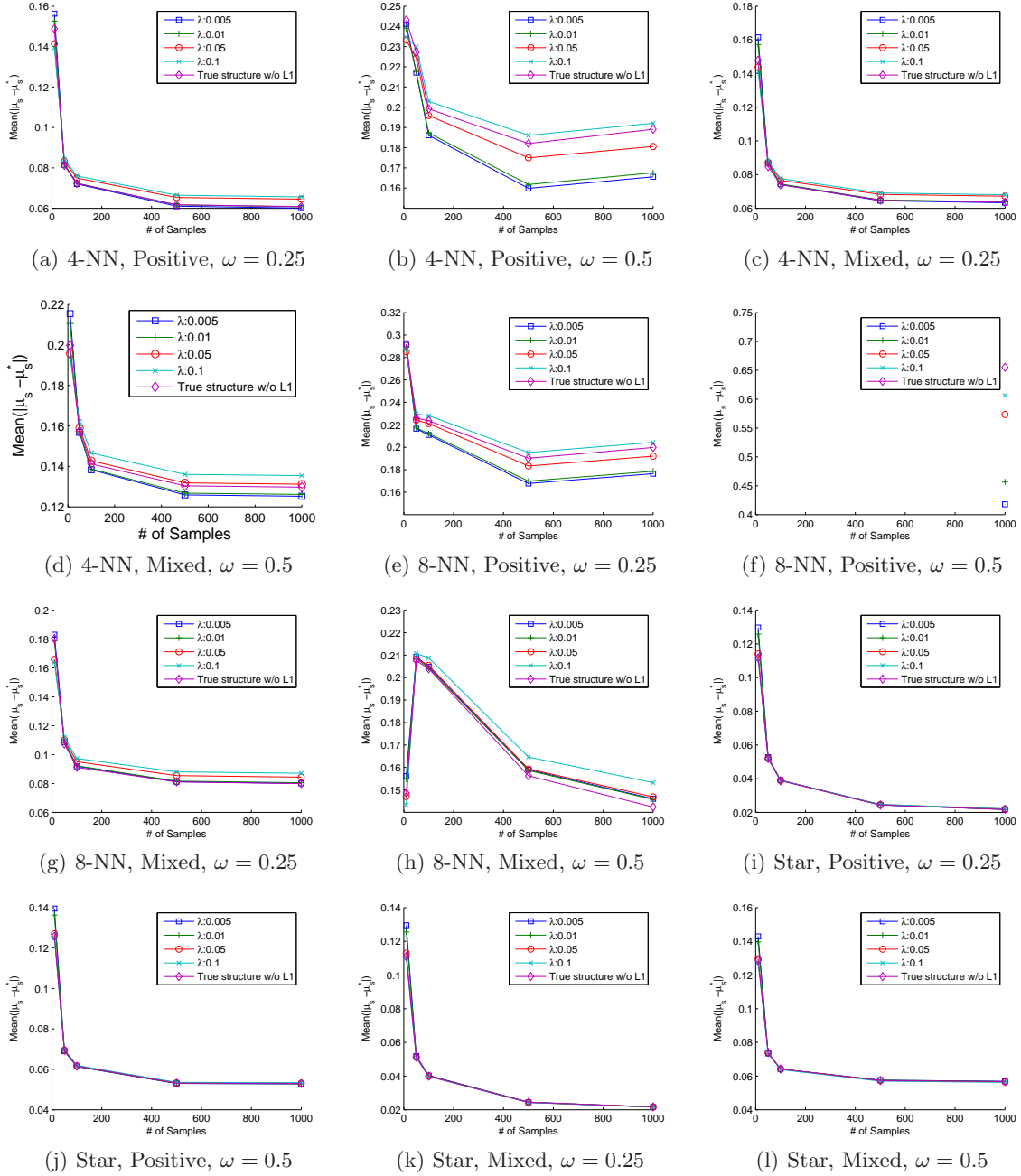


Figure 6. Robustness to perturbation