

---

# Structure learning of antiferromagnetic Ising models

---

Guy Bresler<sup>1</sup>   David Gamarnik<sup>2</sup>   Devavrat Shah<sup>1</sup>

Laboratory for Information and Decision Systems  
Department of Electrical Engineering and Computer Science<sup>1</sup>  
Operations Research Center and Sloan School of Management<sup>2</sup>  
Massachusetts Institute of Technology  
{gbresler,gamarnik,devavrat}@mit.edu

## Abstract

In this paper we investigate the computational complexity of learning the graph structure underlying a discrete undirected graphical model from i.i.d. samples. Our first result is an unconditional computational lower bound of  $\Omega(p^{d/2})$  for learning general graphical models on  $p$  nodes of maximum degree  $d$ , for the class of so-called statistical algorithms recently introduced by Feldman et al. [1]. The construction is equivalent to the notoriously difficult learning parities with noise problem in computational learning theory. Our lower bound suggests that the  $\tilde{O}(p^{d+2})$  runtime required by Bresler, Mossel, and Sly's [2] exhaustive-search algorithm cannot be significantly improved without restricting the class of models.

Aside from structural assumptions on the graph such as it being a tree, hypertree, tree-like, etc., many recent papers on structure learning assume that the model has the correlation decay property. Indeed, focusing on ferromagnetic Ising models, Bento and Montanari [3] showed that all known low-complexity algorithms fail to learn simple graphs when the interaction strength exceeds a number related to the correlation decay threshold. Our second set of results gives a class of repelling (antiferromagnetic) models that have the *opposite* behavior: very strong interaction allows efficient learning in time  $\tilde{O}(p^2)$ . We provide an algorithm whose performance interpolates between  $\tilde{O}(p^2)$  and  $\tilde{O}(p^{d+2})$  depending on the strength of the repulsion.

## 1 Introduction

Graphical models have had tremendous impact in a variety of application domains. For unstructured high-dimensional distributions, such as in social networks, biology, and finance, an important first step is to determine which graphical model to use. In this paper we focus on the problem of structure learning: Given access to  $n$  independent and identically distributed samples  $\sigma^{(1)}, \dots, \sigma^{(n)}$  from an undirected graphical model representing a discrete random vector  $\sigma = (\sigma_1, \dots, \sigma_p)$ , the goal is to find the graph  $G$  underlying the model. Two basic questions are 1) How many samples are required? and 2) What is the computational complexity?

In this paper we are mostly interested in the computational complexity of structure learning. We first consider the problem of learning a general discrete undirected graphical model of bounded degree.

## 1.1 Learning general graphical models

Several algorithms based on exhaustively searching over possible node neighborhoods have appeared in the last decade [4, 2, 5]. Abbeel, Koller, and Ng [4] gave algorithms for learning general graphical models close to the true distribution in Kullback-Leibler distance. Bresler, Mossel, and Sly [2] presented algorithms guaranteed to learn the true underlying graph. The algorithms in both [4] and [2] perform a search over candidate neighborhoods, and for a graph of maximum degree  $d$ , the computational complexity for recovering a graph on  $p$  nodes scales as  $\tilde{O}(p^{d+2})$  (where the  $\tilde{O}$  notation hides logarithmic factors).

While the algorithms in [2] are guaranteed to reconstruct general models under basic nondegeneracy conditions using an optimal number of samples  $n = O(d \log p)$  (sample complexity lower bounds were proved by Santhanam and Wainwright [6] as well as [2]), the exponent  $d$  in the  $\tilde{O}(p^{d+2})$  run-time is impractically high even for constant but large graph degrees. This has motivated a great deal of work on structure learning for special classes of graphical models. But before giving up on general models, we ask the following question:

**Question 1:** Is it possible to learn the structure of general graphical models on  $p$  nodes with maximum degree  $d$  using substantially less computation than  $p^d$ ?

Our first result suggests that the answer to Question 1 is negative. We show an unconditional computational lower bound of  $p^{\frac{d}{2}}$  for the class of *statistical algorithms* introduced by Feldman et al. [1]. This class of algorithms was introduced in order to understand the apparent difficulty of the Planted Clique problem, and is based on Kearns' statistical query model [7]. Kearns showed in his landmark paper that statistical query algorithms require exponential computation to learn parity functions subject to classification noise, and our hardness construction is related to this problem. Most known algorithmic approaches (including Markov chain Monte Carlo, semidefinite programming, and many others) can be implemented as statistical algorithms, so the lower bound is fairly convincing.

We give background and prove the following theorem in Section 4.

**Theorem 1.1.** *Statistical algorithms require at least  $\Omega(p^{\frac{d}{2}})$  computation steps in order to learn the structure of a general graphical models of degree  $d$ .*

If complexity  $p^d$  is to be considered intractable, what shall we consider as tractable? Writing algorithm complexity in the form  $c(d)p^{f(d)}$ , for high-dimensional (large  $p$ ) problems the exponent  $f(d)$  is of primary importance, and we will think of tractable algorithms as having an  $f(d)$  that is bounded by a constant independent of  $d$ . The factor  $c(d)$  is also important, and we will use it to compare algorithms with the same exponent  $f(d)$ .

In light of Theorem 1.1, reducing computation below  $p^{\Omega(d)}$  requires restricting the class of models. One can either restrict the graph structure or the nature of the interactions between variables. The seminal paper of Chow and Liu [8] makes a model restriction of the first type, assuming that the graph is a tree; generalizations include to polytrees [9], hypertrees [10], and others. Among the many possible assumptions of the second type, the correlation decay property is distinguished: to the best of our knowledge all existing low-complexity algorithms require the correlation decay property [3].

## 1.2 Correlation decay property

Informally, a graphical model is said to have the correlation decay property (CDP) if any two variables  $\sigma_s$  and  $\sigma_t$  are asymptotically independent as the graph distance between  $s$  and  $t$  increases. Exponential decay of correlations holds when the distance from independence decreases exponentially fast in graph distance, and we will mean this stronger form when referring to correlation decay. Correlation decay is known to hold for a number of pairwise graphical models in the so-called high-temperature regime, including Ising, hard-core lattice gas, Potts (multinomial) model, and others (see, e.g., [11, 12, 13, 14, 15, 16]).

Bresler, Mossel, and Sly [2] observed that it is possible to efficiently learn models with (exponential) decay of correlations, under the additional assumption that neighboring variables have correlation bounded away from zero (as is true, e.g., for the ferromagnetic Ising model in the high temperature regime). The algorithm they proposed for this setting pruned the candidate set of neighbors for each node to roughly size  $O(d)$  by retaining only those variables with sufficiently high correlations, and then within this set performed the exhaustive search over neighborhoods mentioned before, resulting in a computational cost of  $d^{O(d)}\tilde{O}(p^2)$ . The greedy algorithms of Netrapali et al. [17] and Ray et al. [18] also require the correlation decay property and perform a similar pruning step by retaining only nodes with high pairwise correlation; they then use a different method to select the true neighborhood.

A number of papers consider the problem of reconstructing Ising models on graphs with few short cycles, beginning with Anandkumar et al. [19]. Their results apply to the case of Ising models on sparsely connected graphs such as the Erdős-Renyi random graph  $\mathcal{G}(p, \frac{d}{p})$ . They additionally require the interaction parameters to be either generic or ferromagnetic. Ferromagnetic models have the benefit that neighbors always have a non-negligible correlation because the dependencies cannot cancel, but in either case the results still require the CDP to hold. Wu et al. [20] remove the assumption of generic parameters in [19], but again require the CDP.

Other algorithms for structure learning are based on convex optimization, such as Ravikumar et al.'s [21] approach using regularized node-wise logistic regression. While this algorithm does not explicitly require the CDP, Bento and Montanari [3] found that the logistic regression algorithm of [21] provably fails to learn certain ferromagnetic Ising model on simple graphs not satisfying the CDP. Other convex optimization-based algorithms such as [22, 23, 24] require similar incoherence or restricted isometry-type conditions that are difficult to verify, but likely also require the CDP. Since all known algorithms for structure learning require the CDP, we ask the following question (paraphrasing Bento and Montanari):

**Question 2:** Is low-complexity structure learning possible for models which do not exhibit the CDP, on general bounded degree graphs?

Our second main result answers this question affirmatively by showing that a broad class of repelling models on general graphs can be learned using simple algorithms, even when the underlying model does not exhibit the CDP.

### 1.3 Repelling models

The antiferromagnetic Ising model has a negative interaction parameter, whereby neighboring nodes prefer to be in opposite states. Other popular antiferromagnetic models include the Potts or coloring model, and the hard-core model.

Antiferromagnetic models have the interesting property that correlations between neighbors can be zero due to cancellations. Thus algorithms based on pruning neighborhoods using pairwise correlations, such as the algorithm in [2] for models with correlation decay, does not work for anti-ferromagnetic models. To our knowledge there are no previous results that improve on the  $p^d$  computational complexity for structure learning of antiferromagnetic models on general graphs of maximum degree  $d$ .

Our first learning algorithm, described in Section 2, is for the hard-core model.

**Theorem 1.2 (Informal).** *It is possible to learn strongly repelling models, such as the hard-core model, with run-time  $\tilde{O}(p^2)$ .*

We extend this result to weakly repelling models (equivalent to the antiferromagnetic Ising model parameterized in a nonstandard way, see Section 3). Here  $\beta$  is a repelling strength and  $h$  is an external field.

**Theorem 1.3** (Informal). *Suppose  $\beta \geq (d - \alpha)(h + \ln 2)$  for a nonnegative integer  $\alpha$ . Then it is possible to learn an antiferromagnetic Ising model with interaction  $\beta$ , with run-time  $\tilde{O}(p^{2+\alpha})$ .*

The computational complexity of the algorithm interpolates between  $\tilde{O}(p^2)$ , achievable for strongly repelling models, and  $\tilde{O}(p^{d+2})$ , achievable for general models using exhaustive search. The complexity depends on the repelling strength of the model, rather than structural assumptions on the graph as in [19, 20].

We remark that the strongly repelling models exhibit long-range correlations, yet the algorithmic task of graph structure learning is possible using a certain local procedure.

The focus of this paper is on structure learning, but the problem of parameter estimation is equally important. It turns out that the structure learning problem is strictly more challenging for the models we consider: once the graph is known, it is not difficult to estimate the parameters with low computational complexity (see, e.g., [4]).

## 2 Learning the graph of a hard-core model

We warm up by considering the hard-core (independent set) model. The analysis in this section is straightforward, but serves as an example to highlight the fact that the CDP is not a necessary condition for structure learning.

Given a graph  $G = (V, E)$  on  $|V| = p$  nodes, denote by  $\mathcal{I}(G) \subseteq \{0, 1\}^p$  the set of independent set indicator vectors  $\sigma$ , for which at least one of  $\sigma_i$  or  $\sigma_j$  is zero for each edge  $\{i, j\} \in E(G)$ . The hardcore model with fugacity  $\lambda > 0$  assigns nonzero probability only to vectors in  $\mathcal{I}(G)$ , with

$$P(\sigma) = \frac{\lambda^{|\sigma|}}{Z}, \quad \sigma \in \mathcal{I}(G). \quad (2.1)$$

Here  $|\sigma|$  is the number of entries of  $\sigma$  equal to one and  $Z = \sum_{\sigma \in \mathcal{I}(G)} \lambda^{|\sigma|}$  is the normalizing constant called the partition function. If  $\lambda > 1$  then more mass is assigned to larger independent sets. (We use indicator vectors to define the model in order to be consistent with the antiferromagnetic Ising model in the next section.)

Our goal is to learn the graph  $G = (V, E)$  underlying the model (2.1) given access to independent samples  $\sigma^{(1)}, \dots, \sigma^{(n)}$ . The following simple algorithm reconstructs  $G$  efficiently.

---

### Algorithm 1 SIMPLEHC

---

Input:  $n$  samples  $(\sigma^{(1)}, \dots, \sigma^{(n)}) \in \{0, 1\}^p$ . Output: edge set  $\hat{E}$ .

- 1: Let  $S = \emptyset$
  - 2: For each  $i, j, k$ :
  - 3: If  $\sigma_i^{(k)} = \sigma_j^{(k)} = 1$ , then  $S \leftarrow S \cup \{i, j\}$
  - 4: Output  $\hat{E} = S^c$
- 

The idea behind the algorithm is very simple. If  $\{i, j\}$  belongs to the edge set  $E(G)$ , then for every sample  $\sigma^{(k)}$  either  $\sigma_i^{(k)} = 0$  or  $\sigma_j^{(k)} = 0$  (or both). Thus for every  $i, j$  and  $k$  such that  $\sigma_i^{(k)} = \sigma_j^{(k)} = 1$  we can safely declare  $\{i, j\}$  *not* to be an edge. To show correctness of the algorithm it is therefore sufficient to argue that for every non-edge  $\{i, j\}$  there is a high likelihood that such an independent set  $\sigma^{(k)}$  will be sampled.

Before doing this, we observe that SIMPLEHC actually computes the maximum-likelihood estimate for the graph  $G$ . To see this, note that an edge  $e = \{i, j\}$  for which  $\sigma_i^{(k)} = \sigma_j^{(k)} = 1$  for some  $k$  cannot be in  $\hat{G}$ , since  $P(\sigma^{(k)} | \hat{G} + e) = 0$  for any  $\hat{G}$ . Thus the ML estimate contains a subset of those edges  $e$  which have not been ruled out by  $\sigma^{(1)}, \dots, \sigma^{(n)}$ . But adding any such edge  $e$  to the graph *decreases* the value of the partition function in (2.1) (the sum is over fewer independent sets), thereby increasing the likelihood of each of the samples.

The sample complexity and computational complexity of SIMPLEHC is as follows:

**Theorem 2.1.** *Consider the hard-core model (2.1) on a graph  $G = (V, E)$  on  $|V| = p$  nodes and with maximum degree  $d$ . The sample complexity of SIMPLEHC is*

$$n = O(2^{2d} \max\{1, \lambda^{2d}\} \log p), \quad (2.2)$$

*i.e. with this many samples the algorithm SIMPLEHC correctly reconstructs the graph with probability  $1 - o(1)$ . The computational complexity is*

$$O(np^2) = O((2\lambda)^{2d-2} p^2 \log p). \quad (2.3)$$

*Proof.* Algorithm c correctly reconstructs the graph  $G$  if for every  $e = \{i, j\}$  not in  $E(G)$ , at least one observed independent set vector  $\sigma^{(k)}$  contains both  $i$  and  $j$ . Let  $A_{ij}^k = \{\sigma_i^{(k)} = 0 \text{ or } \sigma_j^{(k)} = 0\}$  be the event that at least one of  $i$  or  $j$  is missing from  $\sigma^{(k)}$ , and let  $A_{ij} = \bigcap_{k=1}^n A_{ij}^k$ . We have by the union bound and independence of  $A_{ij}^k$  for different  $k$ ,

$$\mathbb{P}(\text{error}) \leq \mathbb{P}\left(\bigcup_{(i,j) \in E^c} A_{ij}\right) \leq \binom{p}{2} \mathbb{P}(\bigcap_{k=1}^n A_{ij}^k) = \binom{p}{2} \mathbb{P}(A_{ij}^1)^n \leq \binom{p}{2} (1 - \gamma)^n.$$

The last inequality is from Lemma 2.3, proved at the end of this section, with the quantity  $\gamma$  defined in the statement of the Lemma. To make  $\mathbb{P}(\text{error})$  approach zero at the rate  $1/p$  it suffices to take

$$n = 3\gamma^{-1} \log p.$$

This proves the theorem.  $\square$

We next show that the sample complexity bound in Theorem 2.1 is basically tight:

**Theorem 2.2** (Sample complexity lower bound). *Consider the hard-core model (2.1). There is a family of graphs on  $p$  nodes with maximum degree  $d$  such that if the probability of successful reconstruction is above  $1/2$ , then the number of samples must be*

$$n = \Omega((1 + \lambda)^{d-2} \log p).$$

*Proof.* We give a set of graphs  $\mathcal{G}$  on  $p$  nodes with maximum degree at most  $d$  so that given samples generated from a graph selected uniformly at random from  $\mathcal{G}$ , the (optimal) maximum a posteriori (MAP) rule requires the number of samples stated in the theorem. It is possible to prove the theorem using Fano's inequality, but since we know the ML rule is equivalent to algorithm SIMPLEHC, we can give a direct proof.

We define a set of graphs  $\mathcal{G}_m$  as follows. Let  $G_0$  consist of  $m$  stars of degree  $d-1$ , i.e. for each  $1 \leq v \leq m$  add  $d-1$  nodes  $u_{v,1}, \dots, u_{v,d-1}$  with edges  $\{v, u_{v,i}\}$ . There are  $p = m \cdot (d-1)$  nodes in total. Now we let  $\mathcal{G}_m$  be the set of  $\binom{m}{2}$  graphs  $G_{ij}$  obtained by adding the edge  $\{i, j\}$  between a pair of star centers  $i$  and  $j$ . The graph  $G$  is selected uniformly at random from  $\mathcal{G}_m$  and samples are generated from the model (2.1).

The samples  $\sigma^{(1)}, \dots, \sigma^{(n)}$  do not rule out edge  $e = \{i, j\}$  if there is no  $\sigma^{(k)}$  with  $\sigma_i^{(k)} = \sigma_j^{(k)} = 1$ . Suppose that none of edges  $e_1, e_2, \dots, e_r$  have been ruled out. In this case the observations have the same likelihood under  $G_{e_t}$  for each  $1 \leq t \leq r$ , and it follows that the probability of error is at least  $1 - 1/(r-1)$  since the prior on the models is uniform.

From now onward we assume without loss of generality (by symmetry of the construction) that samples are generated from the model on  $G_{ab}$ . Call  $\sigma^{(k)}$  a *witness* for non-edge  $\{i, j\} \neq \{a, b\}$  if  $\sigma_i^{(k)} = \sigma_j^{(k)} = 1$ . We proceed by upper bounding the probability of observing a witness for each of the  $\binom{m}{2} - 1$  missing edges. Each star center  $i$  is included in a particular random independent set  $\sigma^{(k)}$  with probability at most

$$\frac{\lambda}{\lambda + \sum_{r=0}^{d-1} \binom{d-1}{r} \lambda^r} \leq \frac{1}{\lambda^{-1}(1 + \lambda)^{d-1}} \leq \frac{1}{(1 + \lambda)^{d-2}} := q,$$

even conditional on any assignment to other star centers. It follows that  $\sigma^{(k)}$  is a witness for non-edge  $\{i, j\}$  with probability at most  $q^2$ .

Take an arbitrary cardinality  $m/3$  matching  $\mathcal{M}$  of non-edges (i.e. no two of the non-edges share an endpoint) with each edge also disjoint from  $a$  and  $b$  (recall that we are focusing on graph  $G_{ab}$ ). For each  $e \in \mathcal{M}$  let  $X_e$  be the indicator variable for the event that in  $n$  samples, non-edge  $e$  has no witness. Note that the variables  $X_e$  are mutually independent. If we define  $Z = \sum_{e \in \mathcal{M}} X_e$ , then we have  $\mathbb{E}Z \geq (m/3)(1 - q^2)^n$ , and moreover,  $\mathbb{E}Z^2 \leq \mathbb{E}Z + (\mathbb{E}Z)^2$ .

By the Paley-Zygmund inequality,

$$\mathbb{P}\left(Z \geq \frac{\mathbb{E}Z}{10}\right) \geq \frac{4(\mathbb{E}Z)^2}{5\mathbb{E}Z^2} \geq \frac{4}{5(1 + \mathbb{E}Z/(\mathbb{E}Z)^2)}.$$

If  $\mathbb{E}Z \geq 40$ , then  $\mathbb{P}(Z \geq 3) \geq 2/3$ . If  $Z \geq 4$ , then by the above discussion, the probability of error is at least  $3/4$ , hence  $\mathbb{E}Z \geq 40$  implies  $\mathbb{P}(\text{error}) \geq \frac{2}{3} \cdot \frac{3}{4} = 1/2$ . Hence if the probability of successful reconstruction is above  $1/2$ , then  $\mathbb{E}Z < 40$ , which requires

$$n \geq (1 + o(1)) \frac{\log m/3}{-\log(1 - q^2)} = \Omega((1 + \lambda)^{d-2} \log p),$$

where we used the fact that  $-\log(1 - q^2) = q^2 + o(q^4)$  and  $q^{-1} = (1 + \lambda)^{d-2}$ .  $\square$

**Lemma 2.3.** *Suppose edge  $e = (i, j) \notin G$ , and let  $I$  be an independent set chosen according to the Gibbs distribution (2.1). Then  $\mathbb{P}(\{i, j\} \subseteq I) \geq (2^{2d+1} \max\{1, \lambda^{2d}\})^{-1} \triangleq \gamma$ .*

*Proof.* We can decompose the partition function as

$$\begin{aligned} Z &= \sum_I \lambda^{|I|} = \sum_{I \in S_{\emptyset, \emptyset}} \lambda^{|I|} + \sum_{I \in S_{\emptyset, j}} \lambda^{|I|} + \sum_{I \in S_{i, \emptyset}} \lambda^{|I|} + \sum_{I \in S_{i, j}} \lambda^{|I|} \\ &:= Z_{\emptyset, \emptyset} + Z_{\emptyset, j} + Z_{i, \emptyset} + Z_{i, j}, \end{aligned} \quad (2.4)$$

where  $S_{ij} = \{I : i, j \in I\}$ ,  $S_{i, \emptyset} = \{I : i \in I, j \notin I\}$ , etc. Our goal is to lower bound  $Z_{i, j}$ , since

$$\mathbb{P}(\{i, j\} \subseteq I) = \frac{\sum_{I: \{i, j\} \subseteq I} \lambda^{|I|}}{\sum_I \lambda^{|I|}} = \frac{Z_{i, j}}{Z}. \quad (2.5)$$

We begin by observing that

$$|S_{i, j}| \cdot 2^d \geq |S_{\emptyset, j}|, \quad (2.6)$$

because for each independent set  $I$  with  $i \in I$ , there are at most  $2^d$  distinct independent sets  $I'$  with  $i \notin I'$  with some subset of (at most  $d$ ) neighbors of  $i$  included. One way of observing this is defining the map  $f : S_{\emptyset, j} \rightarrow S_{i, j}$  by  $I \mapsto \{i\} \cup I \setminus \mathcal{N}(i)$ . The map  $f$  takes at most  $2^d$  sets  $I' \in S_{\emptyset, j}$  to each  $I \in S_{i, j}$ , which implies (2.6).

Now, each such set  $I'$  mapping to  $I$  has weight at most a factor  $\max\{1, \lambda^{d-1}\}$  larger than  $I$ , so

$$2^d \max\{1, \lambda^{d-1}\} Z_{i, j} \geq Z_{\emptyset, j}. \quad (2.7)$$

Similar reasoning gives

$$2^d \max\{1, \lambda^{d-1}\} Z_{i, j} \geq Z_{i, \emptyset}, \quad \text{and} \quad 2^{2d} \max\{1, \lambda^{2d-2}\} Z_{i, j} \geq Z_{\emptyset, \emptyset}. \quad (2.8)$$

Using these estimates, we obtain

$$Z \leq Z_{i, j} (1 + 4 \cdot 2^{d-1} \max\{1, \lambda^{d-1}\} + 4 \cdot 2^{2d-2} \max\{1, \lambda^{2d-2}\}) \leq Z_{i, j} \cdot 2^{2d+1} \max\{1, \lambda^{2d}\},$$

and plugging into (2.5) proves the lemma.  $\square$

### 3 Learning anti-ferromagnetic Ising models

In this section we consider the anti-ferromagnetic Ising model on a graph  $G = (V, E)$ . We parametrize the model in such a way that each configuration has probability

$$P(\sigma) = \frac{1}{Z} \exp \{H(\sigma)\}, \quad \sigma \in \{0, 1\}^p, \quad (3.1)$$

where

$$H(\sigma) = -\beta \sum_{(i,j) \in E} \sigma_i \sigma_j + \sum_{i \in V} h_i \sigma_i. \quad (3.2)$$

Here  $\beta > 0$  and  $\{h_i\}_{i \in V}$  are real-valued parameters, and we assume that  $|h_i| \leq h$  for all  $i$ . Working with configurations in  $\{0, 1\}^p$  rather than the more typical  $\{-1, +1\}^p$  amounts to a reparametrization (which is without loss of generality as shown for example in Appendix 1 of [25]). Setting  $h_i = h = \ln \lambda$  for all  $i$ , we recover the hard-core model with fugacity  $\lambda$  in the limit  $\beta \rightarrow \infty$ , so we think of (3.2) as a “soft” independent set model.

#### 3.1 Strongly antiferromagnetic models

We start by considering the situation in which the repelling strength  $\beta$  is sufficiently large that we can modify the approach used for the hard-core model.

Define the empirical conditional probability

$$\widehat{P}(\sigma_a = 1 | \sigma_b = 1) := \frac{\widehat{P}(\sigma_a = 1, \sigma_b = 1)}{\widehat{P}(\sigma_b = 1)},$$

where for any set  $S \subset V$  and  $x_S \in \{0, 1\}^{|S|}$ ,

$$\widehat{P}(\sigma_V = x_V) = \frac{1}{n} \sum_{k=1}^n \mathbb{1}_{\{\sigma_V^{(k)} = x_V\}}.$$

The following lemma shows that we can obtain good estimates for  $P(\sigma_a = 1 | \sigma_b = 1)$ .

**Lemma 3.1.** *Suppose that  $P(\sigma_b = 1) \geq q$  for all  $b \in V$ . If the number of samples is  $n \geq (2/q^2\epsilon^2) \log(8p^2/\zeta)$ , then with probability at least  $1 - \zeta$  we have for all  $a, b \in V$*

$$|P(\sigma_a = 1 | \sigma_b = 1) - \widehat{P}(\sigma_a = 1 | \sigma_b = 1)| \leq \epsilon.$$

The proof is given in the Supplementary Material.

The structure estimation algorithm STRONGREPELLING, described next, determines whether each edge  $\{a, b\}$  is present based on comparing  $\widehat{P}$  to a threshold.

---

#### Algorithm 2 STRONGREPELLING

---

**Input:**  $\beta, h, d$ , and  $n$  samples  $\sigma^{(1)}, \dots, \sigma^{(n)} \in \{0, 1\}^p$ . **Output:** edge set  $\widehat{E}$ .

1: Let  $\delta = (1 + 2^d e^{h(d-1)})^{-2}$  and  $\widehat{E} = \emptyset$

2: For each possible edge  $\{a, b\} \in \binom{V}{2}$ :

3: If  $\widehat{P}(\sigma_a = 1 | \sigma_b = 1) \leq (1 + e^{\beta-h})^{-1} + \delta$ , then add edge  $(a, b)$  to  $\widehat{E}$

4: Output  $\widehat{E}$

---

The performance of algorithm STRONGREPELLING is stated next in Proposition 3.2. The proof is similar to that of Theorem 2.1, replacing Lemma 2.3 by Lemma 3.3 below. Theorem 3.7, given in the next subsection, subsumes Proposition 3.2, so we prove only the stronger Theorem 3.7.

**Proposition 3.2.** *Consider the antiferromagnetic Ising model (3.2) on a graph  $G = (V, E)$  on  $p$  nodes and with maximum degree  $d$ . If*

$$\beta \geq (d + 2)(h + \ln 2),$$



then algorithm STRONGREPELLING has sample complexity

$$n = O\left(2^{2d} e^{2h(d+1)} \log p\right),$$

i.e. this many samples are sufficient to reconstruct the graph with probability  $1 - o(1)$ . The computational complexity of STRONGREPELLING is

$$O(np^2) = O\left(2^{2d} e^{2h(d+1)} p^2 \log p\right).$$

When the interaction parameter  $\beta \geq (d+2)(h + \ln 2)$  it is possible to identify edges using pairwise statistics. The next lemma shows the necessary separation.

**Lemma 3.3.** *We have the following estimates:*

- (i) If  $(a, b) \notin E(G)$ , then  $\mathbb{P}(\sigma_a = 1 | \sigma_b = 1) \geq \frac{1}{1 + 2^{\deg(a)} e^{h(\deg(a)+1)}}$ .
- (ii) Conversely, if  $(a, b) \in E(G)$ , then  $\mathbb{P}(\sigma_a = 1 | \sigma_b = 1) \leq \frac{1}{1 + e^{\beta - h}}$ .
- (iii) For any  $b \in V$ ,  $\mathbb{P}(\sigma_b = 1) \geq \frac{1}{1 + 2^{\deg(b)} e^{h(\deg(b)+1)}}$ .

*Proof.* We start by defining restricted partition function summations: Let

$$\begin{aligned} S_{ab} &= \{\sigma \in \{0, 1\}^p : \sigma_a = \sigma_b = 1\}, \\ S_{a\emptyset} &= \{\sigma \in \{0, 1\}^p : \sigma_a = 1, \sigma_b = 0\}, \end{aligned}$$

and analogously for  $S_{\emptyset b}$  and  $S_{\emptyset\emptyset}$ . We then define  $Z_{ab} = \sum_{\sigma \in S_{ab}} \exp(H(\sigma))$  and again analogously for  $Z_{a\emptyset}, Z_{\emptyset b}, Z_{\emptyset\emptyset}$ .

We first prove part (i) of the lemma, in which we assume that  $(a, b) \notin E(G)$  and lower bound the probability

$$\mathbb{P}(\sigma_a = 1 | \sigma_b = 1) = \frac{Z_{ab}}{Z_{ab} + Z_{\emptyset b}}.$$

To this end, consider the map  $f : S_{\emptyset b} \rightarrow S_{ab}$  defined by taking a configuration  $\sigma$ , setting  $\sigma_i = 0$  for neighbors  $i \in N(a)$ , and then setting  $\sigma_a = 1$ . Since the assumption  $(a, b) \notin E(G)$  implies that  $\sigma_a = \sigma_b = 1$  is a valid assignment to these variables, the definition of  $f$  implies in particular that  $(f(\sigma))_b = 1$  if  $\sigma_b = 1$ , so indeed  $f(\sigma) \in S_{ab}$  for  $\sigma \in S_{\emptyset b}$ .

Now, at most  $2^{\deg(a)}$  sets are mapped by  $f$  to any one set (since the neighbors of  $a$  can be in any configuration), and for any  $\sigma \in S_{\emptyset b}$ ,  $\exp(H(f(\sigma))) \geq \exp(H(\sigma) - h(\deg(a) + 1))$ . This shows that  $2^{\deg(a)} \exp[h(\deg(a) + 1)] Z_{ab} \geq Z_{\emptyset b}$ , and proves part (i) of the lemma. The proof of part (iii) is omitted as it is almost identical to part (i).

We now turn to part (ii), assuming that  $(a, b) \in E(G)$ . Consider the map  $g : S_{ab} \rightarrow S_{\emptyset b}$  taking  $\sigma \in S_{ab}$  and setting  $\sigma_a = 0$  (removing node  $a$  from the independent set). The map  $g$  is one-to-one, and  $H$  increases by  $\beta$  due to resolving the conflict on edge  $(a, b)$ , but decreases by  $h_a \leq h$  due to omitting node  $a$ :  $\exp(H(g(\sigma))) \geq \exp(H(\sigma) + \beta - h)$ . This shows that  $Z_{ab} \geq e^{-\beta+h} Z_{\emptyset b}$ , and completes the proof.  $\square$

### 3.2 Weakly antiferromagnetic models

In this section we focus on learning weakly repelling models and show a trade-off between computational complexity and strength of the repulsion. Recall that for strongly repelling models (with  $\beta \geq d(h + \ln 2)$ ) our algorithm has run-time  $O(p^2 \log p)$ , the same as for the hard-core model (infinite repulsion).

For a subset of nodes  $U \subseteq V$ , let  $G \setminus U$  denote the graph obtained from  $G$  by removing nodes in  $U$  (as well as any edges incident to nodes in  $U$ ).

We can effectively remove nodes from the graph by conditioning: The family of models (3.2) has the property that conditioning on  $\sigma_i = 0$  amounts to removing node  $i$  from the graph.



**Fact 3.4** (Self-reducibility). *Let  $G = (V, E)$ , and consider the model (3.2). Then for any subset of nodes  $U \subseteq V$ , the probability law  $\mathbb{P}_G(\sigma \in \cdot | \sigma_U = \mathbf{0})$  is equal to  $\mathbb{P}_{G \setminus U}(\sigma_{V \setminus U} \in \cdot)$  with the same  $\beta$  and the natural restriction of  $(h_i)_{i \in V}$  to  $(h_i)_{i \in V \setminus U}$ .*

The following corollary is immediate from Lemma 3.3.

**Corollary 3.5.** *We have the conditional probability estimates for deleting subsets of nodes:*

(i) *If  $(a, b) \notin E(G)$ , then for any subset of nodes  $U \subset V \setminus \{a, b\}$ ,*

$$\mathbb{P}_{G \setminus U}(\sigma_a = 1 | \sigma_b = 1) \geq \frac{1}{1 + 2^{\deg_{G \setminus U}(a)} e^{h(\deg_{G \setminus U}(a)+1)}}.$$

(ii) *Conversely, if  $(a, b) \in E(G)$ , then for any subset of nodes  $U \subseteq V \setminus \{a, b\}$*

$$\mathbb{P}_{G \setminus U}(\sigma_a = 1 | \sigma_b = 1) \leq \frac{1}{1 + e^{\beta - h}}.$$

The final ingredient is to show that we can condition by restricting attention to a subset of the observed data,  $\sigma^{(1)}, \dots, \sigma^{(n)}$ , without throwing away too many samples.

**Lemma 3.6.** *Let  $U \subseteq V$  be a subset of nodes and denote the subset of samples with variables  $\sigma_U$  equal to zero by  $A_U = \{\sigma^{(k)} : \sigma_u^{(k)} = 0 \text{ for all } u \in U\}$ . Then with probability at least  $1 - \exp(-n/8(1 + e^h)^{2|U|})$  the number  $|A_U|$  of such samples is at least  $\frac{n}{2} \cdot (1 + e^h)^{-|U|}$ .*

*Proof.* We start by computing the probability that a particular sample  $\sigma^{(k)}$  is in  $A_U$ , or equivalently that  $\sigma_U^{(k)} = \mathbf{0}$ . Let  $W \subseteq V$  be any subset of nodes, and denote by  $x_W$  an assignment to the corresponding variables. Due to the antiferromagnetic nature of the interaction, the distribution (3.2) satisfies the monotonicity property

$$\mathbb{P}(\sigma_a = 1 | \sigma_W = x_W) \leq \mathbb{P}(\sigma_a = 1 | \sigma_W = x_W, \sigma_b = 0)$$

for any neighbor  $b \in N(a) \setminus W$ . This monotonicity together with Bayes' rule gives

$$\begin{aligned} \mathbb{P}(\sigma_U = \mathbf{0}) &= \prod_{i=1}^{|U|} \mathbb{P}(\sigma_{u_i} = 0 | \sigma_{u_1} = \dots = \sigma_{u_{i-1}} = 0) \geq \prod_{i=1}^{|U|} \mathbb{P}(\sigma_{u_i} = 0 | \sigma_{N(u_i)} = \mathbf{0}) \\ &= \prod_{i=1}^{|U|} [1 + e^{h_i}]^{-1} \geq (1 + e^h)^{-|U|}. \end{aligned}$$

Denoting the last displayed quantity by  $q$ , we see that the number of samples obtained,  $|A_U|$ , stochastically dominates a  $\text{Binom}(n, q)$  random variable. We apply Azuma's inequality, which states that

$$\mathbb{P}(\text{Bin}(n, q) - nq \leq -nt) \leq \exp(-nt^2/2),$$

with  $t = q/2$  and this proves the lemma.  $\square$

We now present the algorithm. Effectively, it reduces node degree by removing nodes (which can be done by conditioning on value zero as discussed above), and then applies the strong repelling algorithm to the residual graph.

---

**Algorithm 3** WEAKREPELLING

---

Input:  $\beta, h, d$ , and  $n$  samples  $\sigma^{(1)}, \dots, \sigma^{(n)} \in \{0, 1\}^p$ . Output: edge set  $\widehat{E}$ .

- 1: Let  $\delta = (4 + 4 \cdot 2^{d-\alpha} e^{h(d-\alpha+1)})^{-1}$ ,  $\widehat{E} = \emptyset$ , and  $\alpha = \lceil d - \beta/(h + \ln 2) \rceil$
- 2: For each  $\{a, b\} \in \binom{V}{2}$ :
- 3: For each  $U \subseteq V \setminus \{a, b\}$  of size  $|U| \leq \alpha$
- 4: Compute  $\widehat{P}_{G \setminus U}(\sigma_a = 1 | \sigma_b = 1)$
- 5: If  $\max_{U: |U| \leq \alpha} \widehat{P}_{G \setminus U}(\sigma_a = 1 | \sigma_b = 1) \leq (1 + e^{\beta-h})^{-1} + \delta$ , then add  $\{a, b\}$  to  $\widehat{E}$
- 6: Output  $\widehat{E}$

---

**Theorem 3.7.** Let  $\alpha \leq d$  be a nonnegative integer, and consider the antiferromagnetic Ising model 3.2 with

$$\beta \geq (d + 2 - \alpha)(h + \ln 2)$$

on a graph  $G$ . Algorithm WEAKREPELLING reconstructs the graph with probability  $1 - o(1)$  as  $p \rightarrow \infty$  using

$$n = O\left((1 + e^h)^{2\alpha}(d + 2)2^{4d}e^{4h(d+1)} \log p\right)$$

i.i.d. samples, with run-time

$$O(np^{2+\alpha}) = \tilde{O}_{\beta, h, d}(p^{2+\alpha}).$$

*Proof.* We first argue that all of the empirical conditional probabilities  $\widehat{P}_{G \setminus U}(\sigma_a = 1 | \sigma_b = 1)$  computed in Step 4 of algorithm WEAKREPELLING are accurate up to tolerance  $\delta$  when considering subsets  $U$  of cardinality up to  $\alpha$ , i.e.,

$$|\widehat{P}_{G \setminus U}(\sigma_a = 1 | \sigma_b = 1) - P_{G \setminus U}(\sigma_a = 1 | \sigma_b = 1)| \leq \delta. \quad (3.3)$$

There are at most  $\alpha \binom{p}{\alpha} \leq \alpha p^\alpha$  subsets  $|U|$  of size at most  $\alpha$ . By Lemma 3.6, for each such  $U$ , with probability at least  $1 - \exp(-n/8(1 + e^h)^{2|U|}) \geq 1 - \exp(-n/8(1 + e^h)^{2\alpha})$  the number  $|A_U|$  of samples with  $\sigma_U = 0$  is at least  $\frac{n}{2} \cdot (1 + e^h)^{-\alpha}$ . It follows from the union bound that with probability at least

$$1 - \alpha p^\alpha \exp(-n/8(1 + e^h)^{2\alpha})$$

we have  $|A_U| \geq \frac{n}{2} \cdot (1 + e^h)^{-\alpha}$  for all  $U$  with  $|U| \leq \alpha$ . By the assumed  $n$  in the theorem statement, this holds with probability  $1 - o(1)$ . Denote the *effective sample size* by  $n' = \frac{n}{2} \cdot (1 + e^h)^{-\alpha}$ .

We now apply Lemma 3.1 with

$$\epsilon = \delta := \frac{1}{4(1 + 2^{d-\alpha}e^{h(d-\alpha+1)})}.$$

This requires  $n' \geq (2/q^2\epsilon^2) \log(8p^2/\zeta)$ , where  $q = (1 + 2^{\deg(b)}e^{h(\deg(b)+1)})^{-1}$  and we can take  $\zeta = 1/p$ . The value of  $n$  given in the theorem statement suffices in order that (3.3) holds for all  $a, b \in V \setminus U$ .

We first argue that  $E \subseteq \widehat{E}$ , that is, all true edges are added to  $\widehat{E}$ . Consider an arbitrary edge  $e = (a, b) \in E$ . By Corollary 3.5 and (3.3),

$$\max_{U: |U| \leq \alpha} \widehat{P}_{G \setminus U}(\sigma_a = 1 | \sigma_b = 1) \leq (1 + e^{\beta-h})^{-1} + \delta := A^{-1} + \delta,$$

so in Line 5 of algorithm WEAKREPELLING the edge  $e$  is added to  $\widehat{E}$ .

We next show that  $\widehat{E} \subseteq E$ , so only true edges are included. Suppose  $e = (a, b) \notin E$ . By choosing  $U \subseteq \partial a \setminus \{b\}$ , Corollary 3.5 and (3.3) imply that

$$\widehat{P}_{G \setminus U}(\sigma_a = 1 | \sigma_b = 1) \geq (1 + 2^{d-\alpha}e^{h(d-\alpha+1)})^{-1} - \delta := B^{-1} - \delta,$$

hence the same inequality applies to the maximum computed in Line 5 of the algorithm. Now, under the assumption  $\beta \geq (d + 2 - \alpha)(h + \ln 2)$ , we have

$$A - 1 = e^{\beta-h} \geq 4 \cdot e^{(d-\alpha)h} e^{h2^{d-\alpha}} = 4(B - 1).$$

Hence

$$B^{-1} - A^{-1} \geq \frac{1}{B} - \frac{1}{4B - 3} = \frac{3B - 3}{B(4B - 3)} > \frac{1}{2B} = 2\delta,$$

where the last inequality used the fact that  $B \geq 2$ . This shows that  $e \notin E$  is not added to  $\widehat{E}$  and completes the proof.  $\square$

## 4 Statistical algorithms and proof of Theorem 1.1

We start by describing the statistical algorithm framework introduced by [1]. In this section it is convenient to work with variables taking values in  $\{-1, +1\}$  rather than  $\{0, 1\}$ .

### 4.1 Background on statistical algorithms

Let  $\mathcal{X} = \{-1, +1\}^p$  denote the space of configurations and let  $\mathcal{D}$  be a set of distributions over  $\mathcal{X}$ . Let  $\mathcal{F}$  be a set of solutions (in our case, graphs) and  $\mathcal{Z} : \mathcal{D} \rightarrow 2^{\mathcal{F}}$  be a map taking each distribution  $D \in \mathcal{D}$  to a subset of solutions  $\mathcal{Z}(D) \subseteq \mathcal{F}$  that are defined to be valid solutions for  $D$ . In our setting, since each graphical model under our consideration will be identifiable, there is a single graph  $\mathcal{Z}(D)$  corresponding to each distribution  $D$ . For  $n > 0$ , the *distributional search problem*  $\mathcal{Z}$  over  $\mathcal{D}$  and  $\mathcal{F}$  using  $n$  samples is to find a valid solution  $f \in \mathcal{Z}(D)$  given access to  $n$  random samples from an unknown  $D \in \mathcal{D}$ .

The class of algorithms we are interested in are called *unbiased statistical algorithms*, defined by access to an unbiased oracle. Other related classes of algorithms are defined in [1], and similar lower bounds can be derived for those as well.

**Definition 4.1** (Unbiased Oracle). Let  $D$  be the true distribution. The algorithm is given access to an oracle, which when given any function  $h : \mathcal{X} \rightarrow \{0, 1\}$ , takes an independent random sample  $x$  from  $D$  and returns  $h(x)$ .

These algorithms access the sampled data only through the oracle: unbiased statistical algorithms outsource the computation. Because the data is accessed through the oracle, it is possible to prove *unconditional* lower bounds using information-theoretic methods. As noted in the introduction, many algorithmic approaches can be implemented as statistical algorithms.

We now define a key quantity called average correlation. The *average correlation* of a subset of distributions  $\mathcal{D}' \subseteq \mathcal{D}$  relative to a distribution  $D$  is denoted  $\rho(\mathcal{D}', D)$ ,

$$\rho(\mathcal{D}', D) := \frac{1}{|\mathcal{D}'|^2} \sum_{D_1, D_2 \in \mathcal{D}'} \left| \left\langle \frac{D_1}{D} - 1, \frac{D_2}{D} - 1 \right\rangle_D \right|, \quad (4.1)$$

where  $\langle f, g \rangle_D := \mathbb{E}_{x \sim D}[f(x)g(x)]$  and the ratio  $D_1/D$  represents the ratio of probability mass functions, so  $(D_1/D)(x) = D_1(x)/D(x)$ .

We quote the definition of statistical dimension with average correlation from [1], and then state a lower bound on the number of queries needed by any statistical algorithm.

**Definition 4.2** (Statistical dimension). Fix  $\gamma > 0, \eta > 0$ , and search problem  $\mathcal{Z}$  over set of solutions  $\mathcal{F}$  and class of distributions  $\mathcal{D}$  over  $X$ . We consider pairs  $(D, \mathcal{D}_D)$  consisting of a “reference distribution”  $D$  over  $\mathcal{X}$  and a finite set of distributions  $\mathcal{D}_D \subseteq \mathcal{D}$  with the following property: for any solution  $f \in \mathcal{F}$ , the set  $\mathcal{D}_f = \mathcal{D}_D \setminus \mathcal{Z}^{-1}(f)$  has size at least  $(1 - \eta) \cdot |\mathcal{D}_D|$ . Let  $\ell(D, \mathcal{D}_D)$  be the largest integer  $\ell$  so that for any subset  $\mathcal{D}' \subseteq \mathcal{D}_f$  with  $|\mathcal{D}'| \geq |\mathcal{D}_f|/\ell$ , the average correlation is  $|\rho(\mathcal{D}', D)| < \gamma$  (if there is no such  $\ell$  one can take  $\ell = 0$ ). The *statistical dimension* with average correlation  $\gamma$  and solution set bound  $\eta$  is defined to be the largest  $\ell(D, \mathcal{D}_D)$  for valid pairs  $(D, \mathcal{D}_D)$  as described, and is denoted by  $\text{SDA}(\mathcal{Z}, \gamma, \eta)$ .

**Theorem 4.3** ([1]). Let  $\mathcal{X}$  be a domain and  $\mathcal{Z}$  a search problem over a set of solutions  $\mathcal{F}$  and a class of distributions  $\mathcal{D}$  over  $\mathcal{X}$ . For  $\gamma > 0$  and  $\eta \in (0, 1)$ , let  $\ell = \text{SDA}(\mathcal{Z}, \gamma, \eta)$ . Any (possibly randomized) unbiased statistical algorithm that solves  $\mathcal{Z}$  with probability  $\delta$  requires at least  $m$  calls to the Unbiased Oracle for

$$m = \min \left\{ \frac{\ell(\delta - \eta)}{2(1 - \eta)}, \frac{(\delta - \eta)^2}{12\gamma} \right\}.$$

In particular, if  $\eta \leq 1/6$ , then any algorithm with success probability at least  $2/3$  requires at least  $\min\{\ell/4, 1/48\gamma\}$  samples from the Unbiased Oracle.

In order to show that a graphical model on  $p$  nodes of maximum degree  $d$  requires computation  $p^{\Omega(d)}$  in this computational model, we therefore would like to show that  $\text{SDA}(\mathcal{Z}, \gamma, \eta) = p^{\Omega(d)}$  with  $\gamma = p^{-\Omega(d)}$ .

## 4.2 Soft parities

For any subset  $S \subset [p]$  of cardinality  $|S| = d$ , let  $\chi_S(x) = \prod_{i \in S} x_i$  be the parity of variables in  $S$ . Define a probability distribution by assigning mass to  $x \in \{-1, +1\}^p$  according to

$$p_S(x) = \frac{1}{Z} \exp(c \cdot \chi_S(x)). \quad (4.2)$$

Here  $c$  is a constant, and the partition function is

$$Z = \sum_x \exp(c \cdot \chi_S(x)) = 2^{p-1} (e^c + e^{-c}). \quad (4.3)$$

Our family of distributions  $\mathcal{D}$  is given by these soft parities over subsets  $S \subset [p]$ , and  $|\mathcal{D}| = \binom{p}{d}$ .

**Lemma 4.4.** *Let  $U$  denote the uniform distribution on  $\{-1, +1\}^p$ . For  $S \neq T$ , the correlation  $\langle \frac{p_S}{U} - 1, \frac{p_T}{U} - 1 \rangle$  is exactly equal to zero for any value of  $c$ . If  $S = T$ , the correlation  $\langle \frac{p_S}{U} - 1, \frac{p_S}{U} - 1 \rangle = 1 - \frac{4}{(e^c + e^{-c})^2} \leq 1$ .*

**Lemma 4.5.** *For any set  $\mathcal{D}' \subseteq \mathcal{D}$  of size at least  $|\mathcal{D}|/p^{d/2}$ , the average correlation satisfies  $\rho(\mathcal{D}', U) \leq d^d p^{-d/2}$ .*

*Proof.* By the preceding lemma, the only contributions to the sum (4.1) comes from choosing the same set  $S$  in the sum, of which there are a fraction  $1/|\mathcal{D}'|$ . Each such correlation is at most one by Lemma 4.4, so  $\rho \leq 1/|\mathcal{D}'| \leq p^{d/2}/|\mathcal{D}| = p^{d/2}/\binom{p}{d} \leq d^d/p^{d/2}$ . Here we used the estimate  $\binom{n}{k} \geq (\frac{n}{k})^k$ .  $\square$

*Proof of Theorem 1.1.* Let  $\eta = 1/6$  and  $\gamma = d^d p^{-d/2}$ , and consider the set of distributions  $\mathcal{D}$  given by soft parities as defined above. With reference distribution  $D = U$ , the uniform distribution, Lemma 4.5 implies that SDA( $\mathcal{Z}, \gamma, \eta$ ) of the structure learning problem over distribution (4.2) is at least  $\ell = p^{d/2}/d^d$ . The result follows from Theorem 4.3.  $\square$

## Acknowledgments

This work was supported in part by NSF grants CMMI-1335155 and CNS-1161964, and by Army Research Office MURI Award W911NF-11-1-0036.

## References

- [1] V. Feldman, E. Grigorescu, L. Reyzin, S. Vempala, and Y. Xiao, “Statistical algorithms and a lower bound for detecting planted cliques,” in *STOC*, pp. 655–664, ACM, 2013.
- [2] G. Bresler, E. Mossel, and A. Sly, “Reconstruction of Markov random fields from samples: Some observations and algorithms,” *Approximation, Randomization and Combinatorial Optimization*, pp. 343–356, 2008.
- [3] J. Bento and A. Montanari, “Which graphical models are difficult to learn?,” in *NIPS*, 2009.
- [4] P. Abbeel, D. Koller, and A. Y. Ng, “Learning factor graphs in polynomial time and sample complexity,” *The Journal of Machine Learning Research*, vol. 7, pp. 1743–1788, 2006.
- [5] I. Csiszár and Z. Talata, “Consistent estimation of the basic neighborhood of markov random fields,” *The Annals of Statistics*, pp. 123–145, 2006.
- [6] N. P. Santhanam and M. J. Wainwright, “Information-theoretic limits of selecting binary graphical models in high dimensions,” *Info. Theory, IEEE Trans. on*, vol. 58, no. 7, pp. 4117–4134, 2012.
- [7] M. Kearns, “Efficient noise-tolerant learning from statistical queries,” *Journal of the ACM (JACM)*, vol. 45, no. 6, pp. 983–1006, 1998.
- [8] C. Chow and C. Liu, “Approximating discrete probability distributions with dependence trees,” *Information Theory, IEEE Transactions on*, vol. 14, no. 3, pp. 462–467, 1968.
- [9] S. Dasgupta, “Learning polytrees,” in *Proceedings of the Fifteenth conference on Uncertainty in artificial intelligence*, pp. 134–141, Morgan Kaufmann Publishers Inc., 1999.
- [10] N. Srebro, “Maximum likelihood bounded tree-width markov networks,” in *Proceedings of the Seventeenth conference on Uncertainty in artificial intelligence*, pp. 504–511, Morgan Kaufmann Publishers Inc., 2001.
- [11] R. L. Dobrushin, “Prescribing a system of random variables by conditional distributions,” *Theory of Probability & Its Applications*, vol. 15, no. 3, pp. 458–486, 1970.
- [12] R. L. Dobrushin and S. B. Shlosman, “Constructive criterion for the uniqueness of gibbs field,” in *Statistical physics and dynamical systems*, pp. 347–370, Springer, 1985.
- [13] J. Salas and A. D. Sokal, “Absence of phase transition for antiferromagnetic potts models via the dobrushin uniqueness theorem,” *Journal of Statistical Physics*, vol. 86, no. 3-4, pp. 551–579, 1997.
- [14] D. Gamarnik, D. A. Goldberg, and T. Weber, “Correlation decay in random decision networks,” *Mathematics of Operations Research*, vol. 39, no. 2, pp. 229–261, 2013.
- [15] D. Gamarnik and D. Katz, “Correlation decay and deterministic fpts for counting list-colorings of a graph,” in *Proceedings of the eighteenth annual ACM-SIAM symposium on Discrete algorithms*, pp. 1245–1254, Society for Industrial and Applied Mathematics, 2007.
- [16] D. Weitz, “Counting independent sets up to the tree threshold,” in *Proceedings of the thirty-eighth annual ACM symposium on Theory of computing*, pp. 140–149, ACM, 2006.
- [17] P. Netrapalli, S. Banerjee, S. Sanghavi, and S. Shakkottai, “Greedy learning of markov network structure,” in *48th Allerton Conference*, pp. 1295–1302, 2010.
- [18] A. Ray, S. Sanghavi, and S. Shakkottai, “Greedy learning of graphical models with small girth,” in *50th Allerton Conference*, 2012.
- [19] A. Anandkumar, V. Tan, F. Huang, and A. Willsky, “High-dimensional structure estimation in Ising models: Local separation criterion,” *Ann. of Stat.*, vol. 40, no. 3, pp. 1346–1375, 2012.
- [20] R. Wu, R. Srikant, and J. Ni, “Learning loosely connected Markov random fields,” *Stochastic Systems*, vol. 3, no. 2, pp. 362–404, 2013.
- [21] P. Ravikumar, M. Wainwright, and J. Lafferty, “High-dimensional Ising model selection using  $\ell_1$ -regularized logistic regression,” *The Annals of Statistics*, vol. 38, no. 3, pp. 1287–1319, 2010.
- [22] S.-I. Lee, V. Ganapathi, and D. Koller, “Efficient structure learning of markov networks using  $l_1$ -regularization,” in *Advances in neural information processing systems*, pp. 817–824, 2006.
- [23] A. Jalali, C. C. Johnson, and P. D. Ravikumar, “On learning discrete graphical models using greedy methods,” in *NIPS*, pp. 1935–1943, 2011.
- [24] A. Jalali, P. Ravikumar, V. Vasuki, S. Sanghavi, and U. ECE, “On learning discrete graphical models using group-sparse regularization,” in *Inter. Conf. on AI and Statistics (AISTATS)*, vol. 14, 2011.
- [25] A. Sinclair, P. Srivastava, and M. Thurley, “Approximation algorithms for two-state anti-ferromagnetic spin systems on bounded degree graphs,” *Journal of Statistical Physics*, vol. 155, no. 4, pp. 666–686, 2014.

## Supplementary material

### Proof of Lemma 4.4

Calculating correlation relative to the uniform distribution  $U$  (see Equation (4.1)), we have for  $S \neq T$  with  $|S \cap T| = \lambda$

$$\begin{aligned} \left\langle \frac{p_S}{U} - 1, \frac{p_T}{U} - 1 \right\rangle_U &= \sum_{x \in \{-1, +1\}^p} 2^{-p} (2^p p_S(x) - 1) (2^p p_T(x) - 1) \\ &= \sum_{x \in \{-1, +1\}^p} 2^p p_S(x) p_T(x) - 1. \end{aligned} \quad (4.4)$$

Now

$$\begin{aligned} \sum_{x \in \{-1, +1\}^p} 2^p p_S(x) p_T(x) &= \frac{2^p}{Z^2} \sum_x \exp(c \cdot (\chi_S(x) + \chi_T(x))) \\ &= \frac{2^p \cdot 2^{p-2N+\lambda}}{Z^2} \sum_{x_{S \cap T}} \sum_{x_{S \Delta T}} \exp(c \cdot (\chi_S(x) + \chi_T(x))) \\ &\stackrel{(a)}{=} \frac{2^p \cdot 2^{p-2N+\lambda}}{Z^2} \sum_{x_{S \cap T}} 2^{2N-2\lambda} \cdot \frac{1}{4} \cdot (e^{2c} + e^{-2c} + 2) \\ &= \frac{2^{2p-2}}{Z^2} (e^c + e^{-c})^2 \stackrel{(b)}{=} 1. \end{aligned}$$

Step (a) follows from the fact that for any fixed  $x_{S \cap T}$ , half the assignments to  $x_{S \setminus T}$  result in  $\chi_S = 1$  and half  $\chi_S = -1$ , and similarly for  $x_{T \setminus S}$ ; step (b) is from the formula (4.3) for  $Z$ .

For the case  $S = T$ , we have

$$\begin{aligned} \sum_{x \in \{-1, +1\}^p} 2^p p_S(x) p_T(x) &= \frac{2^p}{Z^2} \sum_x \exp(c \cdot (\chi_S(x) + \chi_T(x))) \\ &= \frac{2^p \cdot 2^{p-1}}{Z^2} (e^{2c} + e^{-2c}) \\ &= \frac{2^{2p-2}}{Z^2} 2(e^c + e^{-c})^2 - \frac{4}{(e^c + e^{-c})^2} = 2 - \frac{4}{(e^c + e^{-c})^2}. \end{aligned}$$

Plugging this into (4.4) completes the proof.  $\square$

### Proof of Lemma 3.1

Azuma's inequality states that if  $Y \sim \text{Bin}(n, \mu)$ , then

$$P(|Y - n\mu| > \gamma n) \leq 2 \exp(-2\gamma^2 n),$$

so for any subset of nodes  $W \subseteq V$  and configuration  $x_W \in \{0, 1\}^{|W|}$  we have

$$P\left(\left|\widehat{P}(\sigma_W = x_W) - P(\sigma_W = x_W)\right| \geq \gamma\right) \leq 2 \exp(-2\gamma^2 n). \quad (4.5)$$

There are  $2^{|W|} \binom{p}{|W|} \leq (2p)^{|W|}$  such choices of  $W$  and  $x_W$  of a given cardinality, and hence at most  $2(2p)^2 = 8p^2$  choices of  $W$  and  $x_W$  with  $|W| \leq 2$ .

Suppose  $n \geq (2\gamma^2)^{-1} \log(8p^2/\zeta)$ . An application of the union bound implies that with probability at least

$$1 - 8p^2 \cdot 2 \exp(-2\gamma^2 n) \geq 1 - \zeta$$

it holds that

$$\left|\widehat{P}(X_W = x_W) - P(X_W = x_W)\right| \leq \gamma \quad (4.6)$$

for all  $\mathcal{W}$  and  $x_{\mathcal{W}}$  with  $|\mathcal{W}| \leq 2$ . For the remainder of the proof assume (4.6) holds. An application of the triangle inequality leads to

$$\begin{aligned}
& \left| \mathbb{P}(\sigma_a = 1 | \sigma_b = 1) - \widehat{\mathbb{P}}(\sigma_a = 1 | \sigma_b = 1) \right| \\
&= \left| \frac{\mathbb{P}(\sigma_a = 1, \sigma_b = 1)}{\mathbb{P}(\sigma_b = 1)} - \frac{\widehat{\mathbb{P}}(\sigma_a = 1, \sigma_b = 1)}{\widehat{\mathbb{P}}(\sigma_b = 1)} \right| \\
&\leq \left| \frac{\mathbb{P}(\sigma_a = 1, \sigma_b = 1)}{\mathbb{P}(\sigma_b = 1)} - \frac{\widehat{\mathbb{P}}(\sigma_a = 1, \sigma_b = 1)}{\mathbb{P}(\sigma_b = 1)} \right| + \left| \frac{\widehat{\mathbb{P}}(\sigma_a = 1, \sigma_b = 1)}{\mathbb{P}(\sigma_b = 1)} - \frac{\widehat{\mathbb{P}}(\sigma_a = 1, \sigma_b = 1)}{\widehat{\mathbb{P}}(\sigma_b = 1)} \right| \\
&\leq \frac{2\gamma}{q},
\end{aligned}$$

Taking  $\gamma = q\epsilon/2$  and plugging into the above value for  $n$  proves the lemma.  $\square$