

Spectral Methods for Learning Multivariate Latent Tree Structure

Authored by:

Tong Zhang
Sham M. Kakade
Le Song
Daniel J. Hsu
Kamalika Chaudhuri
Animashree Anandkumar

Abstract

This work considers the problem of learning the structure of multivariate linear tree models, which include a variety of directed tree graphical models with continuous, discrete, and mixed latent variables such as linear-Gaussian models, hidden Markov models, Gaussian mixture models, and Markov evolutionary trees. The setting is one where we only have samples from certain observed variables in the tree, and our goal is to estimate the tree structure (i.e., the graph of how the underlying hidden variables are connected to each other and to the observed variables). We propose the Spectral Recursive Grouping algorithm, an efficient and simple bottom-up procedure for recovering the tree structure from independent samples of the observed variables. Our finite sample size bounds for exact recovery of the tree structure reveal certain natural dependencies on underlying statistical and structural properties of the underlying joint distribution. Furthermore, our sample complexity guarantees have no explicit dependence on the dimensionality of the observed variables, making the algorithm applicable to many high-dimensional settings. At the heart of our algorithm is a spectral quartet test for determining the relative topology of a quartet of variables from second-order statistics.

1 Paper Body

Graphical models are a central tool in modern machine learning applications, as they provide a natural methodology for succinctly representing high-dimensional distributions. As such, they have enjoyed much success in various AI and machine learning applications such as natural language processing, speech recognition, robotics, computer vision, and bioinformatics. The main statistical challenges associated with graphical models include estimation and inference. While

the body of techniques for probabilistic inference in graphical models is rather rich [1], current methods for tackling the more challenging problems of parameter and structure estimation are less developed and understood, especially in the presence of latent (hidden) variables. The problem of parameter estimation involves determining the model parameters from samples of certain observed variables. Here, the predominant approach is the expectation maximization (EM) algorithm, and only rather recently is the understanding of this algorithm improving [2, 3]. The problem of structure learning is to estimate the underlying graph of the graphical model. In general, structure learning is NP-hard and becomes even more challenging when some variables are unobserved [4]. The main approaches for structure estimation are either greedy or local search approaches [5, 6] or, more recently, based on convex relaxation [7]. 1

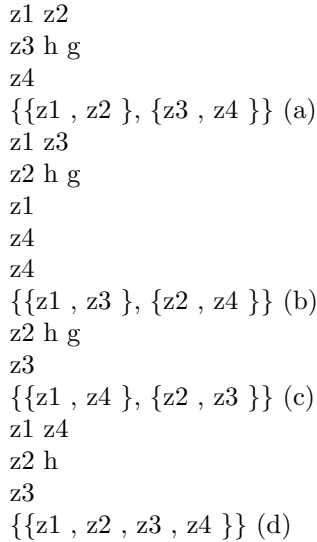


Figure 1: The four possible (undirected) tree topologies over leaves $\{z1, z2, z3, z4\}$. This work focuses on learning the structure of multivariate latent tree graphical models. Here, the underlying graph is a directed tree (e.g., hidden Markov model, binary evolutionary tree), and only samples from a set of (multivariate) observed variables (the leaves of the tree) are available for learning the structure. Latent tree graphical models are relevant in many applications, ranging from computer vision, where one may learn object/scene structure from the co-occurrences of objects to aid image understanding [8]; to phylogenetics, where the central task is to reconstruct the tree of life from the genetic material of surviving species [9]. Generally speaking, methods for learning latent tree structure exploit structural properties afforded by the tree that are revealed through certain statistical tests over every choice of four variables in the tree. These quartet tests, which have origins in structural equation modeling [10, 11], are hypothesis tests of the relative configuration of four (possibly non-adjacent) nodes/variables in the tree (see Figure 1); they are also related to the four point condition associated with a corresponding additive tree metric induced by the distribution [12]. Some early methods for learning tree structure are

based on the use of exact correlation statistics or distance measurements (e.g., [13, 14]). Unfortunately, these methods ignore the crucial aspect of estimation error, which ultimately governs their sample complexity. Indeed, this (lack of) robustness to estimation error has been quantified for various algorithms (notably, for the popular Neighbor Joining algorithm [15, 16]), and therefore serves as a basis for comparing different methods. Subsequent work in the area of mathematical phylogenetics has focused on the sample complexity of evolutionary tree reconstruction [17, 15, 18, 19]. The basic model there corresponds to a directed tree over discrete random variables, and much of the recent effort deals exclusively in the regime for a certain model parameter (the Kesten-Stigum regime [20]) that allows for a sample complexity that is polylogarithmic in the number of leaves, as opposed to polynomial [18, 19]. Finally, recent work in machine learning has developed structure learning methods for latent tree graphical models that extend beyond the discrete distributions of evolutionary trees [21], thereby widening their applicability to other problem domains. This work extends beyond previous studies, which have focused on latent tree models with either discrete or scalar Gaussian variables, by directly addressing the multivariate setting where hidden and observed nodes may be random vectors rather than scalars. The generality of our techniques allows us to handle a much wider class of distributions than before, both in terms of the conditional independence properties imposed by the models (i.e., the random vector associated with a node need not follow a distribution that corresponds to a tree model), as well as other characteristics of the node distributions (e.g., some nodes in the tree could have discrete state spaces and others continuous, as in a Gaussian mixture model). We propose the Spectral Recursive Grouping algorithm for learning multivariate latent tree structure. The algorithm has at its core a multivariate spectral quartet test, which extends the classical quartet tests for scalar variables by applying spectral techniques from multivariate statistics (specifically canonical correlation analysis [22, 23]). Spectral methods have enjoyed recent success in the context of parameter estimation [24, 25, 26, 27]; our work shows that they are also useful for structure learning. We use the spectral quartet test in a simple modification of the recursive grouping algorithm of [21] to perform the tree reconstruction. The algorithm is essentially a robust method for reasoning about the results of quartet tests (viewed simply as hypothesis tests); the tests either confirm or reject hypotheses about the relative topology over quartets of variables. By carefully choosing which tests to consider and properly interpreting their results, the algorithm is able to recover the correct latent tree structure (with high probability) in a provably efficient manner, in terms of both computational and sample complexity. The recursive grouping procedure is similar to the short quartet method from phylogenetics [15], which also guarantees efficient reconstruction in the context of evolutionary trees. However, our method and analysis applies to considerably more general high-dimensional settings; for instance, our sample complexity bound is given in terms of natural correlation con_2

ditions that generalize the more restrictive effective depth conditions of previous works [15, 21]. Finally, we note that while we do not directly address the

question of parameter estimation, provable parameter estimation methods may be derived using the spectral techniques from [24, 25].

2.2.1

Preliminaries Latent variable tree models

Let T be a connected, directed tree graphical model with leaves $V_{\text{obs}} := \{x_1, x_2, \dots, x_n\}$ and internal nodes $V_{\text{hid}} := \{h_1, h_2, \dots, h_m\}$ such that every node has at most one parent. The leaves are termed the observed variables and the internal nodes hidden variables. Note that all nodes in this work generally correspond to multivariate random vectors; we will abuse terminology and still refer to these random vectors as random variables. For any $h \in V_{\text{hid}}$, let $\text{Children}_T(h) \subseteq V_T$ denote the children of h in T . Each observed variable $x \in V_{\text{obs}}$ is modeled as random vector in \mathbb{R}^d , and each hidden variable $h \in V_{\text{hid}}$ as a random vector in \mathbb{R}^k . The joint distribution over all the variables $V_T := V_{\text{obs}} \cup V_{\text{hid}}$ is assumed satisfy conditional independence properties specified by the tree structure over the variables. Specifically, for any disjoint subsets $V_1, V_2, V_3 \subseteq V_T$ such that V_3 separates V_1 from V_2 in T , the variables in V_1 are conditionally independent of those in V_2 given V_3 .

Structural and distributional assumptions

The class of models considered are specified by the following structural and distributional assumptions. Condition 1 (Linear conditional means). Fix any hidden variable $h \in V_{\text{hid}}$. For each hidden child $g \in \text{Children}_T(h) \cap V_{\text{hid}}$, there exists a matrix $A(g|h) \in \mathbb{R}^{k \times k}$ such that $E[g|h] = A(g|h)h$; and for each observed child $x \in \text{Children}_T(h) \cap V_{\text{obs}}$, there exists a matrix $C(x|h) \in \mathbb{R}^{d \times k}$ such that $E[x|h] = C(x|h)h$. We refer to the class of tree graphical models satisfying Condition 1 as linear tree models. Such models include a variety of continuous and discrete tree distributions (as well as hybrid combinations of the two, such as Gaussian mixture models) which are widely used in practice. Continuous linear tree models include linear-Gaussian models and Kalman filters. In the discrete case, suppose that the observed variables take on d values, and hidden variables take k values. Then, each variable is represented by a binary vector in $\{0, 1\}^s$, where $s = d$ for the observed variables and $s = k$ for the hidden variables (in particular, if the variable takes value i , then the corresponding vector is the i -th coordinate vector), and any conditional distribution between the variables is represented by a linear relationship. Thus, discrete linear tree models include discrete hidden Markov models [25] and Markovian evolutionary trees [24]. In addition to the linearity, the following conditions are assumed in order to recover the hidden tree structure. For any matrix M , let $\sigma_t(M)$ denote its t -th largest singular value. Condition 2 (Rank condition). The variables in $V_T = V_{\text{hid}} \cup V_{\text{obs}}$ obey the following rank conditions. 1. For all $h \in V_{\text{hid}}$, $E[hh^T]$ has rank k (i.e., $\sigma_k(E[hh^T]) > 0$). 2. For all $h \in V_{\text{hid}}$ and hidden child $g \in \text{Children}_T(h) \cap V_{\text{hid}}$, $A(g|h)$ has rank k . 3. For all $h \in V_{\text{hid}}$ and observed child $x \in \text{Children}_T(h) \cap V_{\text{obs}}$, $C(x|h)$ has rank k . The rank condition is a generalization of parameter identifiability conditions in latent variable models [28, 24, 25] which rules out various (provably) hard instances in discrete variable settings [24].

h4 h2

x_3
 $h_1 \ x_1$
 T_1
 $h_3 \ x_4$
 T_2
 x_6
 T_3
 x_5
 x_2

Figure 2: Set of trees $F_{h_4} = \{T_1, T_2, T_3\}$ obtained if h_4 is removed. Condition 3 (Non-redundancy condition). Each hidden variable has at least three neighbors. Furthermore, there exists $\epsilon_{\max} \in (0, 1]$ such that for each pair of distinct hidden variables $h, g \in V_{\text{hid}}$, $\det(E[hg])^2 \geq \epsilon_{\max} \cdot \det(E[h]) \cdot \det(E[g])$. The requirement for each hidden node to have three neighbors is natural; otherwise, the hidden node can be eliminated. The quantity ϵ_{\max} is a natural multivariate generalization of correlation. First, note that $\epsilon_{\max} \leq 1$, and that if $\epsilon_{\max} = 1$ is achieved with some h and g , then h and g are completely correlated, implying the existence of a deterministic map between hidden nodes h and g ; hence simply merging the two nodes into a single node h (or g) resolves this issue. Therefore the non-redundancy condition simply means that any two hidden nodes h and g cannot be further reduced to a single node. Clearly, this condition is necessary for the goal of identifying the correct tree structure, and it is satisfied as soon as h and g have limited correlation in just a single direction. Previous works [13, 29] show that an analogous condition ensures identifiability for general latent tree models (and in fact, the conditions are identical in the Gaussian case). Condition 3 is therefore a generalization of this condition suitable for the multivariate setting. Our learning guarantees also require a correlation condition that generalize the explicit depth conditions considered in the phylogenetics literature [15, 24]. To state this condition, first define F_h to be the set of subtrees of T that remain after a hidden variable $h \in V_{\text{hid}}$ is removed from T (see Figure 2). Also, for any subtree $T' \in F_h$, let $V_{\text{obs}}[T']$ be the observed variables in T' . Condition 4 (Correlation condition). There exists $\epsilon_{\min} \in (0, 1]$ such that for all hidden variables $h \in V_{\text{hid}}$ and all triples of subtrees $\{T_1, T_2, T_3\} \in F_h$ in the forest obtained if h is removed from T , \max

$$\begin{aligned}
& \min \\
& \quad x_1 \in V_{\text{obs}}[T_1], x_2 \in V_{\text{obs}}[T_2], x_3 \in V_{\text{obs}}[T_3] \text{ } \{i, j\} \in \{1, 2, 3\} \\
& \quad \epsilon_k(E[x_i x_j]) \geq \epsilon_{\min}.
\end{aligned}$$

The quantity ϵ_{\min} is related to the effective depth of T , which is the maximum graph distance between a hidden variable and its closest observed variable [15, 21]. The effective depth is at most logarithmic in the number of variables (as achieved by a complete binary tree), though it can also be a constant if every hidden variable is close to an observed variable (e.g., in a hidden Markov model, the effective depth is 1, even though the true depth, or diameter, is $m + 1$). If the matrices giving the (conditionally) linear relationship between neighboring variables in T are all well-conditioned, then ϵ_{\min} is at worst exponentially small

in the effective depth, and therefore at worst polynomially small in the number of variables. Finally, also define $\gamma_{\max} :=$

$$\max_{\{x_1, x_2\} \in \mathcal{V}_{\text{obs}}} \left\{ \frac{1}{\sqrt{2}} \left(\mathbb{E}[x_1 x_2^2] \right) \right\}$$

to be the largest spectral norm of any second-moment matrix between observed variables. Note $\gamma_{\max} \leq 1$ in the discrete case, and, in the continuous case, $\gamma_{\max} \leq 1$ if each observed random vector is in isotropic position. In this work, the Euclidean norm of a vector x is denoted by $\|x\|$, and the (induced) spectral norm of a matrix A is denoted by $\|A\|$, i.e., $\|A\| := \sqrt{\lambda_1(A)} = \sup\{\|Ax\| : \|x\| = 1\}$. 4

Algorithm 1 SpectralQuartetTest on observed variables $\{z_1, z_2, z_3, z_4\}$. $\gamma_{i,j}$ of the second-moment matrix Input: For each pair $\{i, j\} \in \{1, 2, 3, 4\}$, an empirical estimate $\hat{\gamma}_{i,j} = \mathbb{E}[z_i z_j]$ and a corresponding confidence parameter $\gamma_{i,j} \in [0, 1]$. Output: Either a pairing $\{\{z_i, z_j\}, \{z_{i'}, z_{j'}\}\}$ or \perp : 1: if there exists a partition of $\{z_1, z_2, z_3, z_4\} = \{z_i, z_j\} \cup \{z_{i'}, z_{j'}\}$ such that $k \leq$

$s=1$
 $\gamma_{i,j} \leq \gamma_{i',j'} + \gamma_{s(i',j')} \leq \gamma_{i',j'} + \gamma_{s(i',j')} + \gamma_{s(i',j')}$
then return the pairing $\{\{z_i, z_j\}, \{z_{i'}, z_{j'}\}\}$.

$k \leq$
 $\gamma_{i',j'} + \gamma_{i',j'} \leq \gamma_{s(i',j')} + \gamma_{i',j'} \leq \gamma_{s(i',j')} + \gamma_{i',j'}$
 $s=1$

2: else return \perp .

3

Spectral quartet tests

This section describes the core of our learning algorithm, a spectral quartet test that determines topology of the subtree induced by four observed variables $\{z_1, z_2, z_3, z_4\}$. There are four possibilities for the induced subtree, as shown in Figure 1. Our quartet test either returns the correct induced subtree among possibilities in Figure 1(a)-(c); or it outputs \perp to indicate abstinence. If the test returns \perp , then no guarantees are provided on the induced subtree topology. If it does return a subtree, then the output is guaranteed to be the correct induced subtree (with high probability). The quartet test proposed is described in Algorithm 1 (SpectralQuartetTest). The notation $[a]_+$ denotes $\max\{0, a\}$ and $[t]$ (for an integer t) denotes the set $\{1, 2, \dots, t\}$.

The quartet test is defined with respect to four observed variables $Z := \{z_1, z_2, z_3, z_4\}$. For each $\gamma_{i,j}$ of the second-moment matrix pair of variables z_i and z_j , it takes as input an empirical estimate $\hat{\gamma}_{i,j} = \mathbb{E}[z_i z_j]$, and confidence bound parameters $\gamma_{i,j}$ which are functions of N , the number of samples $\gamma_{i,j}$, a confidence parameter γ , and of properties of the distributions of z_i and used to compute the $\gamma_{i,j}$. In practice, one uses a single threshold γ for all pairs, which is tuned by the algorithm. Our theoretical analysis also applies to this case. The output of the test is either \perp or a pairing of the variables $\{\{z_i, z_j\}, \{z_{i'}, z_{j'}\}\}$. For example, if the output is the pairing is $\{\{z_1, z_2\}, \{z_3, z_4\}\}$, then Figure 1(a) is the output topology. Even though the configuration in Figure 1(d) is a possibility, the spectral quartet test never returns $\{\{z_1, z_2, z_3\}, \{z_4\}\}$.

, z_4 }, as there is no correct pairing of Z . The topology $\{\{z_1, z_2, z_3, z_4\}\}$ can be viewed as a degenerate case of $\{\{z_1, z_2\}, \{z_3, z_4\}\}$ (say) where the hidden variables h and g are deterministically identical, and Condition 3 fails to hold with respect to h and g . 3.1

Properties of the spectral quartet test

With exact second moments: The spectral quartet test is motivated by the following lemma, which shows the relationship between the singular values of second-moment matrices of the z_i 's and the τ_k induced topology among them in the latent tree. Let $\text{det}_k(M) := \sum_{s=1}^k \sigma_s(M)$ denote the product of the k largest singular values of a matrix M . Lemma 1 (Perfect quartet test). Suppose that the observed variables $Z = \{z_1, z_2, z_3, z_4\}$ have the true induced tree topology shown in Figure 1(a), and the tree model satisfies Condition 1 and Condition 2. Then $\text{det}_k(E[z_1 z_3^T]) \text{det}_k(E[z_2 z_4^T]) \leq \text{det}_k(E[z_1 z_4^T]) \text{det}_k(E[z_2 z_3^T])$, $\det(E[hg^T])^2 = \tau_1 \det(E[hh^T]) \det(E[gg^T])$, $\text{det}_k(E[z_1 z_2^T]) \text{det}_k(E[z_3 z_4^T]) \leq \text{det}_k(E[z_1 z_3^T]) \text{det}_k(E[z_2 z_4^T]) = \text{det}_k(E[z_1 z_4^T]) \text{det}_k(E[z_2 z_3^T])$. This lemma shows that given the true second-moment matrices and assuming Condition 3, the inequality in (1) becomes strict and thus can be used to deduce the correct topology: the correct pairing is $\{\{z_i, z_j\}, \{z_i^?, z_j^?\}\}$ if and only if $\text{det}_k(E[z_i z_j^T]) \text{det}_k(E[z_i^? z_j^{??T}]) \leq \text{det}_k(E[z_i^? z_j^T]) \text{det}_k(E[z_i z_j^{??T}])$. 5

Reliability: The next lemma shows that even if the singular values of $E[z_i z_j^T]$ are not known exactly, then with valid confidence intervals (that contain these singular values) a robust test can be constructed which is reliable in the following sense: if it does not output τ , then the output topology is indeed the correct topology. Lemma 2 (Reliability). Consider the setup of Lemma 1, and suppose that Figure 1(a) is the $\tau_{i,j} = \tau_{i,j}^*$ correct topology. If for all pairs $\{z_i, z_j\} \in Z$ and all $s \in [k]$, $\tau_s(\tau_{i,j}^*) \leq \tau_s(E[z_i z_j^T]) \leq \tau_s(\tau_{i,j}^*) + \tau_{i,j}^*$, and if $\text{SpectralQuartetTest}$ returns a pairing $\{\{z_i, z_j\}, \{z_i^?, z_j^?\}\}$, then $\{\{z_i, z_j\}, \{z_i^?, z_j^?\}\} = \{\{z_1, z_2\}, \{z_3, z_4\}\}$. In other words, the spectral quartet test never returns an incorrect pairing as long as the singular $\tau_{i,j}^*$. The lemma values of $E[z_i z_j^T]$ lie in an interval of length $2\tau_{i,j}^*$ around the singular values of τ below shows how to set the $\tau_{i,j}^*$ s as a function of N , τ and properties of the distributions of z_i and z_j so that this required event holds with probability at least $1 - \epsilon$. We remark that any valid confidence intervals may be used; the one described below is particularly suitable when the observed variables are high-dimensional random vectors. Lemma 3 (Confidence intervals). Let $Z = \{z_1, z_2, z_3, z_4\}$ be four random vectors. Let $\tau_{i,j}^* = \tau_{i,j}^*$ is computed using almost surely, and let $\epsilon \in (0, 1/6)$. If each empirical second-moment matrix τ is based on N iid copies of z_i and z_j , and if $E[\tau_{i,j}^* \tau_{i,j}^*] \leq \text{tr}(E[z_i z_j^T] E[z_i z_j^T])$

, $\tau_{i,j} := 1.55 \ln(24d\tau_{i,j}^*/\epsilon)$, $\max\{\tau_{i,j}^* \tau_{i,j}^*, \tau_{i,j}^* \tau_{i,j}^*\} \leq \tau_{i,j}^* \tau_{i,j}^* + \tau_{i,j}^* \tau_{i,j}^* \leq 2 \max\{\tau_{i,j}^* \tau_{i,j}^*, \tau_{i,j}^* \tau_{i,j}^*\} \leq \tau_{i,j}^* \tau_{i,j}^* + \tau_{i,j}^* \tau_{i,j}^* \leq 2 \tau_{i,j}^* \tau_{i,j}^* + \tau_{i,j}^* \tau_{i,j}^* \leq 3 \tau_{i,j}^* \tau_{i,j}^* \leq 3N$ then with probability $1 - \epsilon$, for all pairs $\{z_i, z_j\} \in Z$ and all $s \in [k]$, $\tau_s(\tau_{i,j}^*) \leq \tau_s(E[z_i z_j^T]) \leq \tau_s(\tau_{i,j}^*) + \tau_{i,j}^* \leq \tau_s(\tau_{i,j}^*) + d\tau_{i,j}^* :=$

(2)

Conditions for returning a correct pairing: The conditions under which Spec-

trivialQuartetTest returns an induced topology (as opposed to τ) are now provided.

An important quantity in this analysis is the level of non-redundancy between the hidden variables h and g . Let $\det(E[hg \mid \tau])^2 := \frac{\det(E[hg \mid \tau])^2}{\det(E[gg \mid \tau]) \det(E[hh \mid \tau])}$. (3) If Figure 1(a) is the correct induced topology among $\{z_1, z_2, z_3, z_4\}$, then the smaller ρ is, the greater the gap between $\det_k(E[z_1 z_2 \mid \tau]) \det_k(E[z_3 z_4 \mid \tau])$ and either of $\det_k(E[z_1 z_3 \mid \tau]) \det_k(E[z_2 z_4 \mid \tau])$ and $\det_k(E[z_1 z_4 \mid \tau]) \det_k(E[z_2 z_3 \mid \tau])$. Therefore, ρ also governs how small the $\rho_{i,j}$ need to be for the quartet test to return a correct pairing; this is quantified in Lemma 4. Note that Condition 3 implies $\rho \geq \rho_{\max} \geq 1$. Lemma 4 (Correct pairing). Suppose that (i) the observed variables $Z = \{z_1, z_2, z_3, z_4\}$ have the true induced tree topology shown in Figure 1(a); (ii) the tree model satisfies Condition 1, Condition 2, and $\rho \geq 1$ (where ρ is defined in (3)), and (iii) the confidence bounds in (2) hold for all $\{i, j\}$ and all $s \in [k]$. If $\rho \geq 1 - \rho_{i,j} \geq \min(1, \frac{1}{k} \min\{\frac{1}{k} \det_k(E[z_i z_j \mid \tau])\} \geq \frac{1}{8k} \rho_{i,j}$ for each pair $\{i, j\}$, then SpectralQuartetTest returns the correct pairing $\{\{z_1, z_2\}, \{z_3, z_4\}\}$.

4

The Spectral Recursive Grouping algorithm

The Spectral Recursive Grouping algorithm, presented as Algorithm 2, uses the spectral quartet test discussed in the previous section to estimate the structure of a multivariate latent tree distribution from iid samples of the observed leaf variables.¹ The algorithm is a modification of the recursive $\tau_{x,y}$ and threshold parameters $\tau_{x,y} \geq 0$ for all pairs $\{x, y\} \in \text{Vobs}$ are globally defined. In particular, we assume the spectral quartet tests use these quantities.

6

Algorithm 2 Spectral Recursive Grouping. $\tau_{x,y}$ for all pairs $\{x, y\} \in \text{Vobs}$ computed from N iid Input: Empirical second-moment matrices Σ samples from the distribution over Vobs ; threshold parameters $\tau_{x,y}$ for all pairs $\{x, y\} \in \text{Vobs}$. τ or τ_{failure} . Output: Tree structure T 1: let $R := \text{Vobs}$, and for all $x \in R$, $T[x] :=$ rooted single-node tree x and $L[x] := \{x\}$. 2: while $|R| > 1$ do 3: let pair $\{u, v\} \in \{\{x, y\} \in R : \text{Mergeable}(R, L[x], \tau_{x,y}, u, v) = \text{true}\}$ be such that $\tau_{x,y}(\{x, y\}) = (x, y) \in L[u] \cap L[v]$ is maximized. If no such pair exists, then halt $\max\{\frac{1}{k} \det_k(E[z_i z_j \mid \tau])\} \geq \frac{1}{8k} \rho_{i,j}$ and return τ_{failure} . 4: let $\text{result} := \text{Relationship}(R, L[u], T[u], u, v)$. 5: if $\text{result} = \text{siblings}$ then 6: Create a new variable h , create subtree $T[h]$ rooted at h by joining $T[u]$ and $T[v]$ to h with edges $\{h, u\}$ and $\{h, v\}$, and set $L[h] := L[u] \cup L[v]$. 7: Add h to R , and remove u and v from R . 8: else if $\text{result} = u$ is parent of v then 9: Modify subtree $T[u]$ by joining $T[v]$ to u with an edge $\{u, v\}$, and modify $L[u] := L[u] \cup L[v]$. 10: Remove v from R . 11: else if $\text{result} = v$ is parent of u then 12: {Analogous to above case.} 13: end if 14: end while $\tau := T[h]$ where $R = \{h\}$. 15: Return T grouping (RG) procedure proposed in [21]. RG builds the tree in a bottom-up fashion, where the initial working set of variables are the observed variables. The variables in the working set always correspond to roots of disjoint subtrees of T discovered by the algorithm. (Note that because these subtrees are rooted, they naturally induce parent/child relationships, but these may differ from those

implied by the edge directions in T .) In each iteration, the algorithm determines which variables in the working set to combine. If the variables are combined as siblings, then a new hidden variable is introduced as their parent and is added to the working set, and its children are removed. If the variables are combined as neighbors (parent/child), then the child is removed from the working set. The process repeats until the entire tree is constructed. Our modification of RG uses the spectral quartet tests from Section 3 to decide which subtree roots in the current working set to combine. Note that because the test may return \perp (a null result), our algorithm uses the tests to rule out possible siblings or neighbors among variables in the working set; this is encapsulated in the subroutine Mergeable (Algorithm 3), which tests quartets of observed variables (leaves) in the subtrees rooted at working set variables. For any pair $\{u, v\} \subseteq R$ submitted to the subroutine (along with the current working set R and leaf sets $L[\cdot]$): Mergeable returns false if there is evidence (provided by a quartet test) that u and v should first be joined with different variables (u' and v' , respectively) before joining with each other; and Mergeable returns true if no quartet test provides such evidence. The subroutine is also used by the subroutine Relationship (Algorithm 4) which determines whether a candidate pair of variables should be merged as neighbors (parent/child) or as siblings: essentially, to check if u is a parent of v , it checks if v is a sibling of each child of u . The use of unreliable estimates of long-range correlations is avoided by only considering highly-correlated variables as candidate pairs to merge (where correlation is measured using observed variables in their corresponding subtrees as proxies). This leads to a sample-efficient algorithm for recovering the hidden tree structure. The Spectral Recursive Grouping algorithm enjoys the following guarantee. Theorem 1. Let $\beta \in (0, 1)$. Assume the directed tree graphical model T over variables (random vectors) $VT = V_{\text{obs}} \cup V_{\text{hid}}$ satisfies Conditions 1, 2, 3, and 4. Suppose the Spectral Recursive

Algorithm 3 Subroutine Mergeable($R, L[\cdot], u, v$). Input: Set of nodes R ; leaf sets $L[v]$ for all $v \in R$; distinct $u, v \in R$. Output: true or false. 1: if there exists distinct $u', v' \in R \setminus \{u, v\}$ and $(x, y, x', y') \in L[u] \times L[v] \times L[u'] \times L[v']$ s.t. SpectralQuartetTest($\{x, y, x', y'\}$) returns $\{\{x, x'\}, \{y, y'\}\}$ or $\{\{x, y'\}, \{x', y\}\}$ then return false. 2: else return true. Algorithm 4 Subroutine Relationship($R, L[\cdot], T[\cdot], u, v$). Input: Set of nodes R ; leaf sets $L[v]$ for all $v \in R$; rooted subtrees $T[v]$ for all $v \in R$; distinct $u, v \in R$. Output: ?siblings? , $\text{?}u$ is parent of v ? ($\text{?}u \text{? } v$?), or $\text{?}v$ is parent of u ? ($\text{?}v \text{? } u$?). 1: if u is a leaf then assert $u \neq v$. 2: if v is a leaf then assert $v \neq u$. 3: let $R[w] := (R \setminus \{w\}) \cup \{w' : w' \text{ is a child of } w \text{ in } T[w]\}$ for each $w \in \{u, v\}$. 4: if there exists child u_1 of u in $T[u]$ s.t. Mergeable($R[u], L[\cdot], u_1, v$) = false then assert $\text{?}u \text{? } v$?. 5: if there exists child v_1 of v in $T[v]$ s.t. Mergeable($R[v], L[\cdot], u, v_1$) = false then assert $\text{?}v \text{? } u$?. 6: if both $\text{?}u \text{? } v$? and $\text{?}v \text{? } u$? were asserted then return ?siblings? . 7: else if $\text{?}u \text{? } v$? was asserted then return $\text{?}v$ is parent of u ? ($\text{?}v \text{? } u$?). 8: else return $\text{?}u$ is parent of v ? ($\text{?}u \text{? } v$?). Grouping algorithm (Algorithm 2) is provided N independent samples from the distribution over V_{obs} , and uses parameters given by $\beta, 2B_{\text{xi}}, x_j, t_{\text{xi}}, x_j, M_{\text{xi}}, M_{\text{xj}}, t_{\text{xi}}, x_j, \text{?xi}, x_j := + (4) N^{3N}$ where

$$B_{xi, xj} := \max_{x \in \mathcal{X}} \{E[\sum_{i=1}^n x_i x_j], E[\sum_{i=1}^n x_i^2]\}$$

$$d_{xi, xj} := \frac{E[\sum_{i=1}^n x_i^2] + E[\sum_{j=1}^n x_j^2] - \text{tr}(E[\sum_{i=1}^n x_i x_j])}{2}$$

$$M_{xi} = \max_{x \in \mathcal{X}} \{B_{xi, xj}\}, M = \max_{xi} M_{xi}, t = \max_{xi} \{d_{xi, xj}\}$$

$$N_{\epsilon} = \frac{200 \epsilon^2 k^2}{\min\{1, \epsilon\} \max\{M^2, t\}}$$

ϵ with then with probability at least $1 - \epsilon$, the Spectral Recursive Grouping algorithm returns a tree T the same undirected graph structure as T .

Consistency is implied by the above theorem with an appropriate scaling of ϵ with N . The theorem reveals that the sample complexity of the algorithm depends solely on intrinsic spectral properties of the distribution. Note that there is no explicit dependence on the dimensions of the observable variables, which makes the result applicable to high-dimensional settings. Acknowledgements Part of this work was completed while DH was at the Wharton School of the University of Pennsylvania and at Rutgers University. AA was supported by in part by the setup funds at UCI and the AFOSR Award FA9550-10-1-0310.

2 References

- [1] M. J. Wainwright and M. I. Jordan. Graphical models, exponential families, and variational inference. *Foundations and Trends in Machine Learning*, 1(1-2):1-305, 2008.
- [2] S. Dasgupta and L. Schulman. A probabilistic analysis of EM for mixtures of separated, spherical Gaussians. *Journal of Machine Learning Research*, 8(Feb):203-226, 2007.
- [3] K. Chaudhuri, S. Dasgupta, and A. Vattani. Learning mixtures of Gaussians using the k-means algorithm, 2009. arXiv:0912.0086.
- [4] D. M. Chickering, D. Heckerman, and C. Meek. Large-sample learning of Bayesian networks is NP-hard. *Journal of Machine Learning Research*, 5:1287-1330, 2004.
- [5] C. Chow and C. Liu. Approximating discrete probability distributions with dependence trees. *IEEE Transactions on Information Theory*, 14(3):462-467, 1968.
- [6] N. Friedman, I. Nachman, and D. Pe'er. Learning Bayesian network structure from massive datasets: the "sparse candidate" algorithm. In *Fifteenth Conference on Uncertainty in Artificial Intelligence*, 1999.
- [7] P. Ravikumar, M. J. Wainwright, and J. Lafferty. High-dimensional Ising

model selection using ℓ_1 regularized logistic regression. *Annals of Statistics*, 38(3):1287–1319, 2010. [8] M. J. Choi, J. J. Lim, A. Torralba, and A. S. Willsky. Exploiting hierarchical context on a large database of object categories. In *IEEE Conference on Computer Vision and Pattern Recognition*, 2010. [9] R. Durbin, S. R. Eddy, A. Krogh, and G. Mitchison. *Biological Sequence Analysis: Probabilistic Models of Proteins and Nucleic Acids*. Cambridge University Press, 1999. [10] J. Wishart. Sampling errors in the theory of two factors. *British Journal of Psychology*, 19:180–187, 1928. [11] K. Bollen. *Structural Equation Models with Latent Variables*. John Wiley & Sons, 1989. [12] P. Buneman. The recovery of trees from measurements of dissimilarity. In F. R. Hodson, D. G. Kendall, and P. Tautu, editors, *Mathematics in the Archaeological and Historical Sciences*, pages 387–395. 1971. [13] J. Pearl and M. Tarsi. Structuring causal trees. *Journal of Complexity*, 2(1):60–77, 1986. [14] N. Saitou and M. Nei. The neighbor-joining method: A new method for reconstructing phylogenetic trees. *Molecular Biology and Evolution*, 4:406–425, 1987. [15] P. L. Erdős, L. A. Székely, M. A. Steel, and T. J. Warnow. A few logs suffice to build (almost) all trees: Part II. *Theoretical Computer Science*, 221:77–118, 1999. [16] M. R. Lacey and J. T. Chang. A signal-to-noise analysis of phylogeny estimation by neighbor-joining: insufficiency of polynomial length sequences. *Mathematical Biosciences*, 199(2):188–215, 2006. [17] P. L. Erdős, L. A. Székely, M. A. Steel, and T. J. Warnow. A few logs suffice to build (almost) all trees (I). *Random Structures and Algorithms*, 14:153–184, 1999. [18] E. Mossel. Phase transitions in phylogeny. *Transactions of the American Mathematical Society*, 356(6):2379–2404, 2004. [19] C. Daskalakis, E. Mossel, and S. Roch. Evolutionary trees and the Ising model on the Bethe lattice: A proof of Steel’s conjecture. *Probability Theory and Related Fields*, 149(1–2):149–189, 2011. [20] H. Kesten and B. P. Stigum. Additional limit theorems for indecomposable multidimensional galtonwatson processes. *Annals of Mathematical Statistics*, 37:1463–1481, 1966. [21] M. J. Choi, V. Tan, A. Anandkumar, and A. Willsky. Learning latent tree graphical models. *Journal of Machine Learning Research*, 12:1771–1812, 2011. [22] M. S. Bartlett. Further aspects of the theory of multiple regression. *Mathematical Proceedings of the Cambridge Philosophical Society*, 34:33–40, 1938. [23] R. J. Muirhead and C. M. Waternaux. Asymptotic distributions in canonical correlation analysis and other multivariate procedures for nonnormal populations. *Biometrika*, 67(1):31–43, 1980. [24] E. Mossel and S. Roch. Learning nonsingular phylogenies and hidden Markov models. *Annals of Applied Probability*, 16(2):583–614, 2006. [25] D. Hsu, S. M. Kakade, and T. Zhang. A spectral algorithm for learning hidden Markov models. In *Twenty-Second Annual Conference on Learning Theory*, 2009. [26] S. M. Siddiqi, B. Boots, and G. J. Gordon. Reduced-rank hidden Markov models. In *Thirteenth International Conference on Artificial Intelligence and Statistics*, 2010. [27] L. Song, S. M. Siddiqi, G. J. Gordon, and A. J. Smola. Hilbert space embeddings of hidden Markov models. In *International Conference on Machine Learning*, 2010. [28] E. S. Allman, C. Matias, and J. A. Rhodes. Identifiability of parameters in latent structure models with many observed variables. *The Annals of Statistics*, 37(6A):3099–3132, 2009. [29] J. Pearl. Probabilis-

tic Reasoning in Intelligent Systems? Networks of Plausible Inference. Morgan Kaufmann, 1988. [30] D. Hsu, S. M. Kakade, and T. Zhang. Dimension-free tail inequalities for sums of random matrices, 2011. arXiv:1104.1672.

9