
A Spectral Algorithm for Latent Tree Graphical Models

Ankur P. Parikh

Le Song

Eric P. Xing

APPARIKH@CS.CMU.EDU

LESONG@CS.CMU.EDU

EPXING@CS.CMU.EDU

School of Computer Science, Carnegie Mellon University, Pittsburgh, PA 15213, USA

Abstract

Latent variable models are powerful tools for probabilistic modeling, and have been successfully applied to various domains, such as speech analysis and bioinformatics. However, parameter learning algorithms for latent variable models have predominantly relied on local search heuristics such as expectation maximization (EM). We propose a fast, local-minimum-free spectral algorithm for learning latent variable models with arbitrary tree topologies, and show that the joint distribution of the observed variables can be reconstructed from the marginals of triples of observed variables irrespective of the maximum degree of the tree. We demonstrate the performance of our spectral algorithm on synthetic and real datasets; for large training sizes, our algorithm performs comparable to or better than EM while being orders of magnitude faster.

1 Introduction

Latent variable models usually refer to probabilistic graphical models that relate a set of observed variables to an additional set of unobserved or hidden variables. Introducing latent variables can greatly improve the flexibility of probabilistic modeling, allowing it to address a diverse range of problems with hidden factors such as in document analysis (Blei et al., 2002), social network modeling (Hoff et al., 2002), speech recognition (Rabiner & Juang, 1986) and bioinformatics (Clark, 1990). Latent variables can also lead to significant savings in model parametrization. By defining a joint model over observed and latent variables, the marginal distribution of the observed variables is obtained by integrating out the latent ones. This allows complex distributions over observed variables (*e.g.*, clique models) to be expressed in terms of more tractable joint models (*e.g.*, tree models) over the augmented variable space.

Although latent variable models are very flexible and can

be represented in a compact way, learning the model parameters has predominantly relied on likelihood maximization and local search heuristics such as expectation maximization (EM) (Dempster et al., 1977). Besides the problem of local minima, EM can require many iterations to reach a prescribed training precision, and high dimensional problems can dramatically slow down EM.

While EM tries to recover the full set of parameters in latent variable models, in many applications it is the inference task that is most interesting. For instance, in speech classification, we are interested in estimating the likelihood of a test sequence under different models; in quantitative finance, we are interested in predicting the price of one stock given the prices of other stocks; or in biological analysis, we are interested in forecasting the expression of one gene given perturbations to other genes. In all these examples, the inference task involves estimating either the joint or conditional distribution of a set of observed variables. Ideally, we want to avoid explicitly recovering the parameters related to latent variables (which leads to non-convex problems), and proceed directly to the interested quantities.

Recently, Hsu et al. (2009) proposed a spectral algorithm for learning hidden Markov models (HMM) which directly estimates the joint distribution of the observed variables without recovering the HMM model parameters. The major computation of the algorithms involves a singular value decomposition (SVD) of small marginal probability matrices involving pairs of observed variables. Compared to EM, this spectral algorithm does not have the problem of local optima, and one can formally study its statistical properties. However, this spectral algorithm is specific to HMMs, and it is not clear whether their techniques can be extended to latent variable models with other topologies.

Mossel & Roch (2006) also proposed a spectral algorithm for latent variable models which applies to arbitrary tree topologies, but they made very restrictive assumptions: all variables (observed and latent) have exactly the same number of states, and all conditional probability tables (CPT) are invertible. Under these conditions, they derived a spectral algorithm that can explicitly recover all CPTs from marginals of triples of observed variables. In many applications, however, latent variables can represent factors sim-

pler than the noisy observations, and the number of hidden states can be smaller than that of the observed states. In these cases, the CPTs are no longer invertible, which renders this spectral algorithm no longer applicable.

In this paper, we propose a novel spectral algorithm for latent variable models with arbitrary tree topologies where the number of hidden states is smaller than or equal to that of the observed states. Instead of first explicitly learning the model parameters and then performing inference, we directly compute the joint distribution of the observed variables without explicitly recovering the model parameters.

We first express the joint distribution of the observed variables using 3rd order tensors, and then show that the components in this tensor representation can be reconstructed from the marginals of triples of observed variables. Given a finite number of samples, our spectral algorithm estimates the desired joint distributions by performing singular value decompositions on a collection of small marginal probability matrices, and hence is very efficient. In addition to estimating the joint distribution, our method can also recover the marginal of any set of observed variables. We conducted experiments on both synthetic and real world data, and demonstrated the competitive performance of our algorithm to EM. For large training sizes, our algorithm performs comparably or better than EM while being orders of magnitude faster.

2 Tensor Algebra

We first give a brief introduction to tensor algebra (for more details, see [Kolda & Bader \(2009\)](#)). A tensor is a multidimensional array, and its order is the number of dimensions, also known as modes. In this paper, vectors (tensors of order one) are denoted by boldface lowercase letters, *e.g.*, \mathbf{a} . Matrices (tensors of order two) are denoted by boldface capital letters, *e.g.*, \mathbf{A} . Higher-order tensors (order three or higher) are denoted by boldface caligraphic letters, *e.g.*, \mathcal{T} . Scalars are denoted by lowercase letters, *e.g.*, a .

Subarrays of a tensor are formed when a subset of the indices is fixed. Particularly, a fiber is defined by fixing every index but one. Fibers are the higher-order analogue of matrix rows and columns. A colon is used to indicate all elements of a mode. Thus, the j th column of a matrix \mathbf{A} is $\mathbf{A}(:, j)$, and the i th row of \mathbf{A} is $\mathbf{A}(i, :)$. Analogously, the mode- n fiber of a N th order tensor \mathcal{T} is then denoted as $\mathcal{T}(i_1, i_2, \dots, i_{n-1}, :, i_{n+1}, \dots, i_N)$. Fibers can be used to construct higher order tensors from lower order ones. For instance, a third order tensor \mathcal{A} which is diagonal in mode-2 and 3 can be constructed from a matrix \mathbf{B} by setting $\mathcal{T}(:, i, i) = \mathbf{B}(:, i)$.

Tensors can be multiplied together. For matrices and vectors, we will use standard notation for their multiplications, *e.g.*, $\mathbf{B}\mathbf{a}$ and $\mathbf{A}\mathbf{B}$. For tensors of higher order, we

are particularly interested in multiplying a tensor by matrices and vectors. The n -mode matrix product is the multiplication of a tensor with a matrix in mode n of the tensor. Let $\mathcal{T} \in \mathbb{R}^{I_1 \times I_2 \times \dots \times I_N}$ be an N th order tensor and $\mathbf{A} \in \mathbb{R}^{J \times I_n}$ be a matrix. Then

$$\mathcal{T}' = \mathcal{T} \times_n \mathbf{A} \in \mathbb{R}^{I_1 \times \dots \times I_{n-1} \times J \times I_{n+1} \times \dots \times I_N}, \quad (1)$$

where the entries $\mathcal{T}'(i_1, \dots, i_{n-1}, j, i_{n+1}, \dots, i_N)$ are defined as $\sum_{i_n=1}^{I_n} \mathcal{T}(i_1, \dots, i_n, \dots, i_N) \mathbf{A}(j, i_n)$. We will further introduce two useful properties of n -mode matrix product. First, for distinct modes in a series of multiplications, the order of the multiplication can be exchanged

$$\mathcal{T} \times_n \mathbf{A} \times_m \mathbf{B} = \mathcal{T} \times_m \mathbf{B} \times_n \mathbf{A} \quad (m \neq n). \quad (2)$$

Second, the matrices can be combined first, if the modes in a series of multiplications are the same

$$\mathcal{T} \times_n \mathbf{A} \times_n \mathbf{B} = \mathcal{T} \times_n (\mathbf{B}\mathbf{A}). \quad (3)$$

We note that n -mode matrix product does not change the order of a tensor, but the size of the tensor may change.

Multiplication of a tensor with a vector in mode n of the tensor is called n -mode vector product. Let $\mathcal{T} \in \mathbb{R}^{I_1 \times I_2 \times \dots \times I_N}$ and $\mathbf{a} \in \mathbb{R}^{I_n}$. Then

$$\mathcal{T}' = \mathcal{T} \bar{\times}_n \mathbf{a} \in \mathbb{R}^{I_1 \times \dots \times I_{n-1} \times I_{n+1} \times \dots \times I_N} \quad (4)$$

where the entries $\mathcal{T}'(i_1, \dots, i_{n-1}, i_{n+1}, \dots, i_N)$ is defined as $\sum_{i_n=1}^{I_n} \mathcal{T}(i_1, i_2, \dots, i_n, \dots, i_N) \mathbf{a}(i_n)$. We note that n -mode vector product actually reduces the order of the tensor, *i.e.*, \mathcal{T}' is order $N - 1$ if \mathcal{T} is order N . Using n -mode vector product, we can turn a diagonal operation on vector-matrix product into tensor multiplications, *i.e.*,

$$\text{diag}(\mathbf{a}^\top \mathbf{B}) = \mathcal{T} \bar{\times}_1 \mathbf{a}, \quad \text{where } \mathcal{T}(:, i, i) = \mathbf{B}(:, i). \quad (5)$$

3 Latent Tree Graphical Models (LTMs)

In this paper, we will focus on discrete latent variable models where the conditional independence structures are specified by a tree. Furthermore, we follow the convention that uppercase letters denote random variables (*e.g.*, X_i) and lowercase letters their instantiations (*e.g.*, x_i). A latent tree model defines a joint probability distribution over a set of O observed variables $\mathcal{O} = \{X_1, \dots, X_O\}$ and a set of H hidden variables $\mathcal{H} = \{X_{O+1}, \dots, X_{O+H}\}$. For simplicity, we assume that all observed variables have S_O states and all hidden variables have S_H states, and $S_O \geq S_H$. The complete set of variables is denoted by $\mathcal{X} = \mathcal{O} \cup \mathcal{H}$.

The joint distribution of \mathcal{X} in a latent tree model is fully characterized by a set of conditional probability tables (CPTs). More specifically, we can select an arbitrary (observed or latent) node in the tree as the root, and sort the nodes in the tree in topological order. Then the set of CPTs between nodes and their parents $\mathbb{P}[X_i | X_{\pi_i}]$ are sufficient to characterize the joint distribution (the root node X_r has no parent, *i.e.*, $\mathbb{P}[X_r | X_{\pi_r}] = \mathbb{P}[X_r]$),

$$\mathbb{P}[x_1, \dots, x_{O+H}] = \prod_{i=1}^{O+H} \mathbb{P}[x_i | x_{\pi_i}]. \quad (6)$$

Compared to tree models which are defined solely on ob-

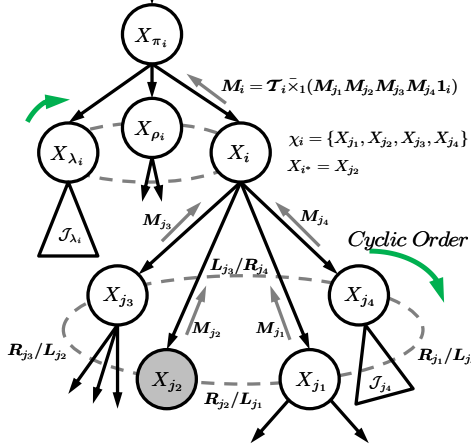


Figure 1. Notation for latent tree models. After rooting the tree and sorting the nodes in topological order, we denote the parent of a node X_i as X_{π_i} , and the set of children of X_i as χ_i . We order the sibling of X_i in a clockwise cyclic order, such that the left (next) sibling of X_i is denoted as X_{λ_i} , and the right (previous) sibling as X_{ρ_i} . We denote the subtree induced by X_i and its descendants as \mathcal{J}_i , and an observed variable in \mathcal{J}_i as X_{i^*} . We note that there may be multiple observed variables in \mathcal{J}_i , and we will use X_{i^*} to refer to either of them. Similarly, we will also use $X_{\lambda_i^*}$ and $X_{\rho_i^*}$ to denote observed variables at subtrees rooted at X_{λ_i} and X_{ρ_i} respectively. Last, we use shaded nodes to denote observed variables, and un-shaded ones for hidden variables.

served variables (e.g., models obtained from Chow & Liu (1968) algorithm), latent tree models encompass a much larger classes of models, allowing more flexibility in modeling observed variables. This is evident if we compute the marginal distribution of the observed variables by summing out the latent ones,

$$\mathbb{P}[x_1, \dots, x_O] = \sum_{x_{O+1}} \dots \sum_{x_{O+H}} \prod_{i=1}^{O+H} \mathbb{P}[x_i | x_{\pi_i}]. \quad (7)$$

This expression leads to complicated conditional independence structures between observed variables depending on the tree topology. In other words, latent tree models allow complex distributions over observed variables (e.g., clique models) to be expressed in terms of more tractable joint models over the augmented variable space. This is a significant saving in model parametrization. We also note that for latent tree models, observed variables can be internal nodes as well as leaf nodes, allowing diverse structures, such as cliques connected by trees, for observed variables. Other notation related to the topological ordering of the nodes in a latent tree model are illustrated in Figure 1.

4 Tensor Representation for LTMs

The computation of the marginal distribution of the observed variables in (7) can be expressed in terms of tensor multiplications. Basically, the information contained in each tensor will correspond to the information in a conditional probability table (CPT) of the model and the tensor multiplications implement the summations. However,

there are multiple ways of rewriting (7) using tensor notation, and not all of them provide intuition or easy derivation to a spectral algorithm. In this section, we will derive a specific representation of latent tree models which requires only tensors up to 3rd order and provides us a basis for deriving a spectral algorithm. For simplicity, we assume that all internal nodes of the tree correspond to latent variables and leaf nodes correspond to observed variables. The general case where observed variables can appear as both internal and leaf nodes can be found in the supplementary.

Root. We associate the root node X_r with the marginal probability vector $\mathbf{r} = \mathbb{P}[X_r]$ of X_r . Here we use $\mathbb{P}[X_r]$ to denote a vector where its k th dimension is defined as $\mathbb{P}[X_r = k]$ and k ranges over all possible assignments of X_r . Similarly, we use $\mathbb{P}[X_i, X_j]$, $\mathbb{P}[X_i | X_j]$, and $\mathbb{P}[X_i, X_j, X_k]$ to denote the joint probability matrix, conditional probability matrix and joint probability tensor respectively, and we denote $\mathbb{P}[X_i, x_j, X_k]$ as a slice (or a fiber) of the tensor when the middle variable is fixed to x_j .

Internal nodes. We associate each internal node X_i with a 3rd order tensor \mathcal{T}_i related to the conditional probability matrix between X_i and its parent X_{π_i} . This tensor is diagonal in its 2nd and 3rd mode, and hence its nonzero entries can be accessed by two indices k and l . Furthermore, $\mathcal{T}_i(k, l, l) = \mathbb{P}[X_i = k | X_{\pi_i} = l]$. The reason for defining this tensor is to implement the marginalization operation over variable X_i using tensor vector multiplications, and return the result as a diagonal matrix. Let $\mathbf{v} = \mathbb{P}[x_j | X_i]$ be a likelihood vector. Then the mode-1 vector product, $\mathcal{T}_i \bar{\times}_1 \mathbf{v}$, results in a diagonal matrix with nonzero entries $M_i(l, l) = \sum_k \mathbb{P}[x_j | X_i = k] \mathbb{P}[X_i = k | X_{\pi_i} = l]$ (or $\mathbb{P}[x_j | X_{\pi_i} = l]$).

Leaf nodes. We associate each leaf node x_i , which is always observed, with a diagonal matrix \mathbf{M}_i related to the likelihood of x_i , i.e., $M_i(l, l) = \mathbb{P}[x_i | X_{\pi_i} = l]$. Furthermore, we let this \mathbf{M}_i be the messages passed from the leaf nodes to their parents. We can show that the marginal probability of the leaf nodes (equation (7)) can be computed recursively using a message passing algorithm (Pearl, 1988): each node in the tree sends a message to its parent according to the reverse topological order of the nodes, and the final messages are aggregated in the root to yield the desired quantity.

Message updates. The outgoing message from an internal node X_i to its parent can be computed as (also see Figure 1)

$$\mathbf{M}_i = \mathcal{T}_i \bar{\times}_1 (\mathbf{M}_{j_1} \mathbf{M}_{j_2} \dots \mathbf{M}_{j_J} \mathbf{1}_i) \quad (8)$$

where each \mathbf{M}_j (a diagonal matrix) is an incoming message from a child, and $j_1, j_2, \dots, j_J \in \chi_i$ range over all children of X_i ($J = |\chi_i|$). The $\mathbf{1}_i$ is a vector of all ones with suitable size, and it is used to reduce the incoming messages (all are diagonal matrices) to a single vector. The computation in (8) essentially implements the message up-

date we often see in ordinary message passing algorithm (Pearl, 1988),

$$\mathbf{m}_i(x_{\pi_i}) = \sum_{x_i} \mathbb{P}[x_i|x_{\pi_i}] \mathbf{m}_{j_1}(x_i) \dots \mathbf{m}_{j_J}(x_i), \quad (9)$$

where $\mathbf{m}_j(x_i)$ represents incoming messages to X_i . The $\mathbf{M}_{j_1} \mathbf{M}_{j_2} \dots \mathbf{M}_{j_J} \mathbf{1}_i$ corresponds to aggregating all incoming messages $\mathbf{m}_{j_1}(x_i) \dots \mathbf{m}_{j_J}(x_i)$, and the $\mathcal{T}_i \bar{\times}_1 *$ corresponds to the summation $\sum_{x_i} \mathbb{P}[x_i|x_{\pi_i}] *$. The characteristic feature of our update in (8) is that we use 3rd tensors to ensure that given incoming messages as diagonal matrices, the outgoing message is also a diagonal matrix, such that message aggregation can be carried on recursively.

Marginal probability. At root node, all incoming messages are combined to yield the final joint probability,

$$\mathbb{P}[x_1, \dots, x_O] = \mathbf{r}^\top (\mathbf{M}_{j_1} \mathbf{M}_{j_2} \dots \mathbf{M}_{j_J} \mathbf{1}_r), \quad (10)$$

where $\mathbf{r}^\top *$ operation basically marginalizes out the root variables, *i.e.*, $\sum_{x_r} \mathbb{P}[x_r] *$. Note that in both (8) and (10), we require that the message multiplications $\mathbf{M}_{j_1} \mathbf{M}_{j_2} \dots \mathbf{M}_{j_J}$ are ordered according to the cyclic order of the siblings illustrated in Figure 1.

5 Spectral Algorithm for LTM

The drawback of the representations in (8) and (10) is that they require the exact knowledge of the parameters (CPTs) associated with latent variables, but none of them are available in training. If we are not interested in recovering these model parameters but only in the marginal probability of the observed variables (*i.e.*, inference), we may not need to recover the transition tensors \mathcal{T} , the messages \mathbf{M} and the root marginal \mathbf{r} exactly. Our key observation from (8) and (10) is that as long as we can recover them up to some invertible transformations, we will still be able to compute the marginal probability correctly.

For example, we can introduce a pair of matrices, \mathbf{R} and \mathbf{L} ($\mathbf{R}\mathbf{L}^{-1} = \mathbf{I}$), between \mathbf{M}_{j_1} and \mathbf{M}_{j_2} in (10), *i.e.*,

$$\mathbb{P}[x_1, \dots, x_O] = \mathbf{r}^\top (\mathbf{M}_{j_1} \mathbf{R} \mathbf{L}^{-1} \mathbf{M}_{j_2} \dots \mathbf{M}_{j_J} \mathbf{1}_r),$$

without changing the final marginal probability. This is interesting because the transformed representation ($\mathbf{M}_{j_1} \mathbf{R}$ and $\mathbf{L}^{-1} \mathbf{M}_{j_2}$) provides us an additional degree of freedom for algorithm design: we want to choose \mathbf{R} and \mathbf{L} from the large class of invertible matrices, such that the transformed representation can be recovered from observed quantities without the need for accessing the latent variables.

We will show that such \mathbf{R} and \mathbf{L} can be constructed from singular vectors \mathbf{U} of the joint probability matrices of certain pairs of observed variables. Given a finite number of samples, this leads us to a very efficient algorithm for estimating the joint probability $\mathbb{P}[x_1, \dots, x_O]$: the main computation only involves a sequence of singular value decompositions of empirical pairwise joint probability matrices. Furthermore, our algorithm's sample complexity will de-

pend on the singular values of the pairwise joint probability matrices. The dependence of our method on the spectral properties of the model give the name ‘‘spectral algorithm’’.

5.1 Transformed Tensor Representation

More specifically, we transform each message \mathbf{M}_j by two invertible matrices \mathbf{L}_j and \mathbf{R}_j , one from its left and one from the right (see Figure 1 for illustration). Then the message update in (8) can be re-written as

$$\mathbf{M}_i = \mathcal{T}_i \bar{\times}_1 \quad (11)$$

$(\mathbf{L}_{j_1} \mathbf{L}_{j_1}^{-1} \mathbf{M}_{j_1} \mathbf{R}_{j_1} \mathbf{L}_{j_2}^{-1} \mathbf{M}_{j_2} \mathbf{R}_{j_2} \dots \mathbf{L}_{j_J}^{-1} \mathbf{M}_{j_J} \mathbf{R}_{j_J} \mathbf{R}_{j_J}^{-1} \mathbf{1}_i)$. We further require that $\mathbf{R}_{j_1} = \mathbf{L}_{j_2}$, $\mathbf{R}_{j_2} = \mathbf{L}_{j_3}$ *etc.* such that the transformations cancel out with each other, *e.g.*, $\mathbf{R}_{j_1} \mathbf{L}_{j_2}^{-1} = \mathbf{I}$. Since the message multiplications $\mathbf{M}_{j_1} \mathbf{M}_{j_2} \dots \mathbf{M}_{j_J}$ are ordered according to the cyclic order of the siblings, this is equivalent to requiring that \mathbf{R}_j of X_j be equal to matrix \mathbf{L}_{ρ_j} of its right sibling X_{ρ_j} , *i.e.*, $\mathbf{R}_j = \mathbf{L}_{\rho_j}$; similarly, we require $\mathbf{L}_j = \mathbf{R}_{\lambda_j}$.

The same can be done with \mathcal{T}_i . If we propagate message \mathbf{M}_i one step further to its parent X_{π_i} and the outgoing message at X_{π_i} can be written as

$$\mathbf{M}_{\pi_i} = \mathcal{T}_{\pi_i} \bar{\times}_1 (\dots \mathbf{L}_i^{-1} \mathbf{M}_i \mathbf{R}_i \dots \mathbf{1}_{\pi_i}) \quad (12)$$

We now show that we can re-define the tensors \mathcal{T} and messages \mathbf{M} by grouping them with these transformations such that the message recursion still works in this transformed representation.

From (11) and (12), we observe that the components in the tensor representation (\mathcal{T} , \mathbf{M} and $\mathbf{1}$) have been ‘‘sandwiched’’ by the invertible transformations (\mathbf{R} and \mathbf{L}). Therefore we can define a set of new quantities

$$\tilde{\mathcal{T}}_i = \mathcal{T}_i \times_1 \mathbf{L}_{j_1}^\top \times_2 \mathbf{L}_i^{-1} \times_3 \mathbf{R}_i^\top \quad (13)$$

$$\tilde{\mathbf{M}}_j = \mathbf{L}_j^{-1} \mathbf{M}_j \mathbf{R}_j \quad (14)$$

$$\tilde{\mathbf{1}}_i = \mathbf{R}_{j_J}^{-1} \mathbf{1}_i \quad (15)$$

where $\tilde{\mathcal{T}}_i$ in (13) is obtained by absorbing the leading transformation (\mathbf{L}_{j_1}) from (11), and the other two parent-level transformations from (12). Then the message update can be expressed in these new quantities as

$$\tilde{\mathbf{M}}_i = \tilde{\mathcal{T}}_i \bar{\times}_1 (\tilde{\mathbf{M}}_{j_1} \dots \tilde{\mathbf{M}}_{j_J} \tilde{\mathbf{1}}_i) \quad (16)$$

Similarly, we can transform the probability vector \mathbf{r} at root node by \mathbf{L}_{j_1} , which leads to $\tilde{\mathbf{r}}^\top = \mathbf{r}^\top \mathbf{L}_{j_1}$ and the final joint probability is

$$\mathbb{P}[x_1, \dots, x_O] = \tilde{\mathbf{r}}^\top (\tilde{\mathbf{M}}_{j_1} \dots \tilde{\mathbf{M}}_{j_J} \tilde{\mathbf{1}}_r). \quad (17)$$

Next we show that \mathbf{R} and \mathbf{L} can be chosen smartly such that all quantities in the transformed representation can be recovered from observed quantities.

5.2 Observable Representation

We now show that the transformed representation of latent tree models can be reconstructed from observed quantities. Each component requires at most 3 observed variables for the reconstruction, so the trees have to be tri-

connected, *i.e.*, each node has at least 3 neighbors. Our strategy is to relate latent quantities to observed quantities using the sum rule of probability; then based on these relations, we solve for the latent quantities. For notation, let O_{ij} be a conditional probability matrix: $O_{ij}(k, l) = \mathbb{P}[X_i = k | X_j = l]$. Recall that X_{λ_i} and X_{ρ_i} are the left/right siblings of X_i respectively, and X_{i^*} is an observed leaf in the subtree rooted at X_i (see Figure 1).

Transition tensor in (13). If we choose $L_{j_1} = O_{j_1^* i}^\top U_{j_1^*}$, $L_i = O_{i^* \pi_i}^\top U_{i^*}$ and $R_i = O_{\rho_i^* \pi_i} U_{\rho_i^*}$, then the tensor in (13) becomes

$$\tilde{\mathcal{T}}_i = \mathcal{T}_i \times_1 (U_{j_1^*}^\top O_{j_1^* i}) \times_2 (O_{i^* \pi_i}^\top U_{i^*})^{-1} \times_3 (U_{\rho_i^*}^\top O_{\rho_i^* \pi_i}), \quad (18)$$

where U_{j^*} is a matrix specific to $O_{j^* i}$ such that $O_{j^* i}^\top U_{j^*}$ is invertible. One choice for U_{j^*} that meets this requirement is to perform a “thin” singular value decompositions (SVDs) of the pair marginal $\mathbb{P}[X_{\lambda_j^*}, X_{j^*}]$, and then take the first S_H right principal singular vectors to form U_{j^*} . In the rest of the paper, we use this SVD approach to obtain U_{j^*} .

Let \mathbf{a}_i be a marginal probability vector with entries $\mathbf{a}_i(k) = \mathbb{P}[X_i = k]$, then $\mathbb{P}[X_{\lambda_i^*}, X_{i^*}] = O_{\lambda_i^* \pi_i} \text{diag}(\mathbf{a}_{\pi_i}) O_{i^* \pi_i}^\top$ by summing out X_{π_i} . Next, we multiply tensor $\tilde{\mathcal{T}}_i$ in its mode-2 by $\mathbb{P}[X_{\lambda_i^*}, X_{i^*}] U_{i^*}$, resulting in

$$\tilde{\mathcal{T}}_i \times_2 (\mathbb{P}[X_{\lambda_i^*}, X_{i^*}] U_{i^*}) \quad (19)$$

$$= \tilde{\mathcal{T}}_i \times_2 O_{\lambda_i^* \pi_i} \text{diag}(\mathbf{a}_{\pi_i}) O_{i^* \pi_i}^\top U_{i^*} \\ = \mathbb{P}[X_{j_1^*}, X_{\lambda_i^*}, X_{\rho_i^*}] \times_1 U_{j_1^*}^\top \times_3 U_{\rho_i^*}^\top, \quad (20)$$

where in (20) we cancel out the $O_{i^* \pi_i}^\top U_{i^*}$ in mode-2 of $\tilde{\mathcal{T}}_i$ in (18), and use the fact that the tensor multiplication effectively marginalizes out variable X_i and X_{π_i} .

Based on (19) and (20), we can recover the hidden $\tilde{\mathcal{T}}_i$ as

$$\tilde{\mathcal{T}}_i = \mathbb{P}[X_{j_1^*}, X_{\lambda_i^*}, X_{\rho_i^*}] \times_1 U_{j_1^*}^\top \\ \times_2 (\mathbb{P}[X_{\lambda_i^*}, X_{i^*}] U_{i^*})^\dagger \times_3 U_{\rho_i^*}^\top, \quad (21)$$

where we multiple the pseudo-inverse of $\mathbb{P}[X_{\lambda_i^*}, X_{i^*}] U_{i^*}$ to the mode-2 of the tensor to recover $\tilde{\mathcal{T}}_i$ from (19). We also note that since X_{j_1} is a child of X_i , we can set $i^* = j_1^*$.

Message in (14). We only need to derive the case for leaf nodes, and the messages from internal nodes will be automatically taken care of by the message recursion. Choose $L_i = O_{i \pi_i}^\top U_i$ and $R_i = O_{\rho_i^* \pi_i}^\top U_{\rho_i^*}$ respectively, then

$$\tilde{M}_i = (O_{i \pi_i}^\top U_i)^{-1} M_i O_{\rho_i^* \pi_i}^\top U_{\rho_i^*}. \quad (22)$$

Next we relate \tilde{M}_i to the marginal probability of triples of observed variables in the following way

$$O_{\lambda_i^* \pi_i} \text{diag}(\mathbf{a}_{\pi_i}) O_{i \pi_i}^\top U_i \tilde{M}_i \\ = \mathbb{P}[X_{\lambda_i^*}, X_i] U_i \tilde{M}_i \quad (23)$$

$$= O_{\lambda_i^* \pi_i} \text{diag}(\mathbf{a}_{\pi_i}) M_i O_{\rho_i^* \pi_i}^\top U_{\rho_i^*} \\ = \mathbb{P}[X_{\lambda_i^*}, x_i, X_{\rho_i^*}] U_{\rho_i^*}. \quad (24)$$

where we use the fact that M_i is a diagonal matrix

with nonzero entries coming from the likelihood vector $\mathbb{P}[x_i | X_{\pi_i}]$. Note that $\mathbb{P}[X_{\lambda_i^*}, x_i, X_{\rho_i^*}]$ is a slice of the joint probability tensor where the middle variable X_i is fixed at x_i . Based on (23) and (24), we can recover \tilde{M}_i as

$$\tilde{M}_i = (\mathbb{P}[X_{\lambda_i^*}, X_i] U_i)^\dagger \mathbb{P}[X_{\lambda_i^*}, x_i, X_{\rho_i^*}] U_{\rho_i^*}, \quad (25)$$

where the tensor slice $\mathbb{P}[X_{\lambda_i^*}, x_i, X_{\rho_i^*}]$ behaves as a matrix.

One in (15) and root marginal in (17). We let $R_{j_j} = O_{\rho_{j_j}^* i}^\top U_{\rho_{j_j}^*}$ (the parent of X_{j_j} is X_i), then the transformed $\tilde{\mathbf{1}}$ becomes $(O_{\rho_{j_j}^* i}^\top U_{\rho_{j_j}^*})^{-1} \mathbf{1}_i$. Next we multiply it by $\mathbb{P}[X_{j_j^*}, X_{\rho_{j_j}^*}] = O_{j_j^* i} \text{diag}(\mathbf{a}_i) (O_{\rho_{j_j}^* i}^\top U_{\rho_{j_j}^*})$, resulting in

$$O_{j_j^* i} \text{diag}(\mathbf{a}_i) (O_{\rho_{j_j}^* i}^\top U_{\rho_{j_j}^*}) (O_{\rho_{j_j}^* i}^\top U_{\rho_{j_j}^*})^{-1} \mathbf{1}_i \\ = \mathbb{P}[X_{j_j^*}, X_{\rho_{j_j}^*}] U_{\rho_{j_j}^*} \mathbf{1}_i = \mathbb{P}[X_{j_j^*}] \quad (26)$$

Based on (26), we can recover $\tilde{\mathbf{1}}_i$ as

$$\tilde{\mathbf{1}}_i = (\mathbb{P}[X_{j_j^*}, X_{\rho_{j_j}^*}] U_{\rho_{j_j}^*})^\dagger \mathbb{P}[X_{j_j^*}]. \quad (27)$$

Similarly, we can also recover the transformed $\tilde{\mathbf{r}}$ as

$$\tilde{\mathbf{r}}^\top = \mathbf{r}^\top L_{j_1^*} = \mathbf{r}^\top O_{j_1^* i}^\top U_{j_1^*} = \mathbb{P}[X_{j_1^*}]^\top U_{j_1^*}, \quad (28)$$

where here X_{j_1} refers to a child of X_r .

We note that our derivation of the observable representation requires a variable to have at least two different siblings. This basically requires that each internal node have at least 3 children. However, with slight modifications, our reasoning in this section also applies to the cases where an internal node X_i has only 2 children. In this case, we only need to conceptually “adopt” a child from the sibling X_{λ_i} of X_i , and order this “adopted” child with the other 2 children of X_i in cyclic order.

5.3 Spectral Algorithm

We now present a spectral algorithm for latent tree models based on the observable representation in the previous section. We assume that the topologies of the trees are given. Then, given N *i.i.d.* samples of the observed variables $\{x_1, \dots, x_O\}$ from a latent tree model (latent variables are not observed at either training or test time), our algorithm first selects a root and then sorts the other nodes in topological order. Then we estimate the empirical marginal distributions of up to triples of variables needed to recover the observable representations in (21), (25), (27) and (28).

Next for each U_{j^*} required in (21), (25), (27) and (28), we perform a “thin” singular value decompositions (SVDs) of the empirical pair marginal $\mathbb{P}[X_{\lambda_j^*}, X_{j^*}]$, and taking the first S_H right principal singular vectors to form an estimate \hat{U}_{j^*} . Since X_{j^*} is an arbitrary observed variable in the subtree induced by X_j and its descendants, there can be multiple choices for X_{j^*} . In practice, we choose X_{j^*} such that the S_H^{th} largest singular value of $\mathbb{P}[X_{\lambda_j^*}, X_{j^*}]$ is large (justified by Theorem 1). Finally, we compute the U_j for each observed variable X_j only once at initialization, and store it for later use. We summarize the algorithm in Algorithm 1.

Algorithm 1 Spectral Algorithm for Latent Tree Models

In: Tree topology and N i.i.d. samples $\{x_1^s, \dots, x_O^s\}_{s=1}^N$

Out: Estimated marginal $\hat{\mathbb{P}}[x_1, \dots, x_O]$

- 1: Select a root node and sort other nodes in topological order.
- 2: For each $\tilde{\mathcal{T}}_i, \tilde{\mathbf{1}}_i$ in an internal node X_i , each $\tilde{\mathcal{M}}_i$ in a leaf node, and $\tilde{\mathbf{r}}$ in the root node, estimate the following empirical marginals from training samples:

	Triple	Pair	Singleton
$\tilde{\mathcal{T}}_i$	$\hat{\mathbb{P}}[X_{j_1}^*, X_{\lambda_i}^*, X_{\rho_i}^*]$	$\hat{\mathbb{P}}[X_{\lambda_i}^*, X_{i^*}]$	-
$\tilde{\mathbf{1}}_i$	-	$\hat{\mathbb{P}}[X_{j_J}^*, X_{\rho_{j_J}}^*]$	$\hat{\mathbb{P}}[X_{j_J}^*]$
$\tilde{\mathcal{M}}_i$	$\hat{\mathbb{P}}[X_{\lambda_i}^*, x_i, X_{\rho_i}^*]$	$\hat{\mathbb{P}}[X_{\lambda_i}^*, X_i]$	-
$\tilde{\mathbf{r}}$	-	-	$\hat{\mathbb{P}}[X_{j_1}^*]$

($j_1, j_J \in \chi_i$ for internal nodes and $j_1 \in \chi_r$ for the root.)

- 3: For each leaf node X_i , perform a “thin” singular value decomposition of $\hat{\mathbb{P}}[X_{\lambda_i}^*, X_i] = \mathbf{V}\Sigma\mathbf{U}^\top$; let $\hat{\mathbf{U}}_i = \mathbf{U}(:, 1 : S_H)$ be the first S_H principal right singular vectors.
- 4: Estimate each $\tilde{\mathcal{T}}_i, \tilde{\mathbf{1}}_i$ in internal nodes, each $\tilde{\mathcal{M}}_i$ in leaf nodes, and $\tilde{\mathbf{r}}$ via

$$\tilde{\mathcal{T}}_i = \hat{\mathbb{P}}[X_{j_1}^*, X_{\lambda_i}^*, X_{\rho_i}^*] \times_1 \hat{\mathbf{U}}_{j_1}^\top \times_2 (\hat{\mathbb{P}}[X_{\lambda_i}^*, X_{i^*}] \hat{\mathbf{U}}_{i^*})^\dagger \times_3 \hat{\mathbf{U}}_{\rho_i}^\top \quad (29)$$

$$\tilde{\mathbf{1}}_i = (\hat{\mathbb{P}}[X_{j_J}^*, X_{\rho_{j_J}}^*] \hat{\mathbf{U}}_{\rho_{j_J}})^\dagger \hat{\mathbb{P}}[X_{j_J}^*] \quad (30)$$

$$\tilde{\mathcal{M}}_i = (\hat{\mathbb{P}}[X_{\lambda_i}^*, X_i] \hat{\mathbf{U}}_i)^\dagger \hat{\mathbb{P}}[X_{\lambda_i}^*, x_i, X_{\rho_i}^*] \hat{\mathbf{U}}_{\rho_i}^\top \quad (31)$$

$$\tilde{\mathbf{r}}^\top = \hat{\mathbb{P}}[X_{j_1}^*]^\top \hat{\mathbf{U}}_{j_1} \quad (32)$$

- 5: In reverse topological order, internal nodes send messages

$$\tilde{\mathcal{M}}_i = \tilde{\mathcal{T}}_i \bar{\times}_1 (\tilde{\mathcal{M}}_{j_1} \dots \tilde{\mathcal{M}}_{j_J} \tilde{\mathbf{1}}_i), \quad (33)$$

and at root node, all incoming messages are combined

$$\hat{\mathbb{P}}[x_1, \dots, x_O] = \tilde{\mathbf{r}}^\top (\tilde{\mathcal{M}}_{j_1} \dots \tilde{\mathcal{M}}_{j_J} \tilde{\mathbf{1}}_r) \quad (34)$$

5.4 Sample Complexity

We analyze the sample complexity of Algorithm 1 and find that it depends on the tree topology and the spectral properties of the true model. See the supplementary for a proof¹.

Theorem 1 For any $\epsilon > 0, 0 < \delta < 1$, let

$$N \geq O\left(\frac{(d_{\max} S_H)^{2\ell+1}}{\alpha\beta\epsilon^2}\right) \log \frac{|\mathcal{O}|}{\delta}$$

where $\sigma_{S_H}^{\text{th}}(*)$ returns the S_H^{th} largest singular value and

$$\alpha = \min_{i \neq j, i, j \in \mathcal{O}} \sigma_{S_H}(\mathbb{P}[X_i, X_j])^4$$

$$\beta = \min_{i \in \mathcal{O}} \sigma_{S_H}(\mathbf{O}_{i\pi_i})^2$$

Then $\sum_{x_1, \dots, x_O} \left| \hat{\mathbb{P}}[x_1, \dots, x_O] - \mathbb{P}[x_1, \dots, x_O] \right| \leq \epsilon$ with probability $1 - \delta$.

This result implies that the estimation problem gets harder as the maximum degree d_{\max} of the hidden nodes, the number S_H of the hidden states, and the length ℓ of the chain of hidden variables increase. Furthermore, the sample complexity depends exponentially in ℓ , which suggests that we should choose the root of the tree to make ℓ small. However, we believe that such adverse dependence on ℓ is due to the artifact of our analysis.

¹can be found at <http://www.sailing.cs.cmu.edu>

A special case of latent tree models is hidden Markov models (HMMs). Recently, Hsu et al. (2009) derived a spectral algorithm specific to HMMs. Their reasoning relies heavily on a single connected chain of hidden variables and each hidden variable has an observed variable attached. Although this excludes many interesting tree topologies, they obtained a tighter sample complexity bound which is $O(\frac{S_H \ell^2}{\alpha\beta\epsilon^2})$ (polynomial in ℓ). This also suggests that our analysis can be further improved.

6 Discussion

There can be other representations of latent tree models. However, we are not aware of other representations for the marginal computation in (7) which lead to an efficient estimation procedure that requires only marginals of triples of observed variables. For instance, we can represent each hidden node X_i by a $(J+1)$ th order transition tensor \mathcal{T}_i ($J = |\chi_i|$) where \mathcal{T}_i is diagonal in the first J modes ($\mathcal{T}_i(k, \dots, k, l) = \mathbb{P}[X_i = k | X_{\pi_i}]$). The root node can be represented as a J th order diagonal tensor \mathcal{R} ($J = |\chi_r|$, and $\mathcal{R}(k, \dots, k) = \mathbb{P}[X_r = k]$). If we represent message from leaf node x_i as vectors $\mathbf{m}_i = \mathbb{P}[x_i | X_{\pi_i}]$, then the message update can be written as

$$\mathbf{m}_i = \mathcal{T}_i \bar{\times}_1 \mathbf{m}_{j_1} \dots \bar{\times}_J \mathbf{m}_{j_J}. \quad (35)$$

and the final marginal probability is computed as

$$\mathbb{P}[x_1, \dots, x_O] = \mathcal{R} \bar{\times}_1 \mathbf{m}_{j_1} \dots \bar{\times}_J \mathbf{m}_{j_J}. \quad (36)$$

Although this representation is more elegant, it does not suggest a simple spectral algorithm. We conjecture that in order to recover an observable representation for \mathcal{R} , we will need the marginal of J observed variables, which will be impractical for large J .

Another representation is to use elementwise vector products. Let $\mathbf{T}_i = \mathbb{P}[X_i | X_{\pi_i}]$ be the transition matrix at internal node X_i , and represent message from leaf node x_i as vectors $\mathbf{m}_i = \mathbb{P}[x_i | X_{\pi_i}]$, then message update becomes

$$\mathbf{m}_i = \mathbf{T}_i (\mathbf{m}_{j_1} \circ \dots \circ \mathbf{m}_{j_J}), \quad (37)$$

and the final marginal probability is computed as

$$\mathbb{P}[x_1, \dots, x_O] = \mathbf{r}^\top (\mathbf{m}_{j_1} \circ \dots \circ \mathbf{m}_{j_J}). \quad (38)$$

However, in this form, it is not clear how one can derive an observable representation and a spectral algorithm.

7 Experiments

We evaluate our spectral algorithm using 3 sets of experiments. Overall, our spectral algorithm is orders of magnitude faster than EM (main competitor) while being more accurate at sufficiently large sample sizes.

7.1 Comparisons with EM and Chow-Liu Tree

We first generate synthetic data from latent tree models with four different topologies (broad4, broad9, deep4 and deep5 shown in Figure 2). These tree topologies are designed so that the tree either grows broader (broad4 to broad9), or grows deeper (deep4 to deep5). We set $S_O = 6$,

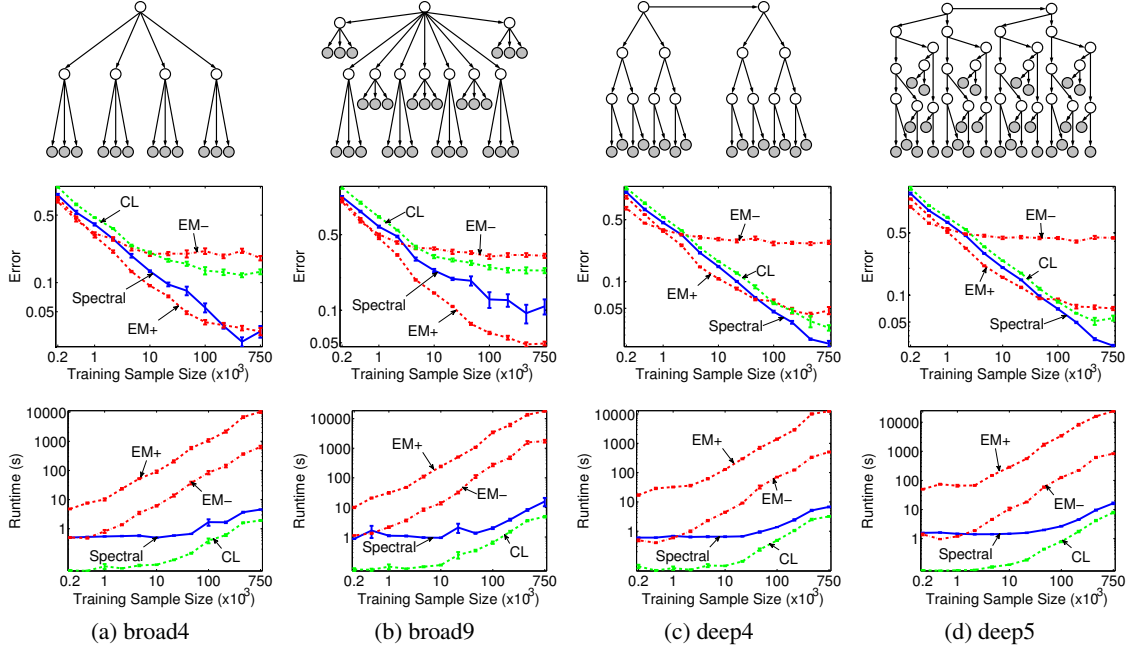


Figure 2. Comparison of our spectral algorithm (Spectral) to EM algorithm with high precision (EM+) and low precision (EM-), and Chow-Liu tree learning (CL) for 4 different tree topologies shown in the first row. Both errors and runtimes are plotted in log scale.

and $S_H = 2$, i.e., $S_O > S_H$. In this case, the Mossel & Roch (2006) algorithm no longer applies, and therefore we compare with EM and inference on the Chow-Liu tree.

For an experiment on a given tree type with N training points, we randomly generate 10 sets of model parameters and sample N training points and 1000 test points for each parameter set. For EM, we learn the CPTs (with 5 restarts) based on the training points and the true latent tree topology, and then perform inference on test points using message passing. We experiment with a low precision EM (0.01, denoted as EM-) and a high precision EM (0.0005, denoted as EM+) to give a better perspective on the time versus accuracy tradeoff as compared to our approach. For the Chow-Liu tree, we first learn the topology of a fully observable tree model using the Chow-Liu algorithm and the true pairwise marginals; then we learn the CPTs in the Chow-Liu tree using training points and perform inference on test points using message passing. We measure the performance of joint estimation using $\epsilon = \frac{|\mathbb{P}[x_1, \dots, x_O] - \mathbb{P}[x_1, \dots, x_O]|}{\mathbb{P}[x_1, \dots, x_O]}$, and we vary the training sample size N from 200 to 750,000, and report both the runtime for training and the test error for inference in Figure 2. (Test runtimes are about the same for all methods).

Figure 2(a)(b)(c)(d), show that our spectral algorithm can be orders of magnitude faster than EM for large sample sizes and that Chow-Liu tree learning is fastest (we do not count the time for tree topology learning), since it only needs to estimate a collection of CPTs for pairs of random variables. In terms of estimation errors, EM+ and EM- perform the best for small training sizes. However,

when the sample sizes go beyond 5,000, the performance of EM- levels off and our spectral algorithm overtakes EM-. This is because EM- learns the models with a low precision (0.01), and as we increase sample size beyond certain point, this fixed precision simply dominates the estimation error. Similarly, we also see that for large sample sizes our algorithm overtakes EM+ on 2 of the 4 experiments (deep4 and deep5) and performs equally on broad4 (EM+ does better on broad9). Furthermore, our spectral algorithm is significantly better than the Chow-Liu tree learning over the range of sample sizes. This is expected since both our spectral algorithm and EM use the correct tree topology while the fully observable tree learned by Chow-Liu has introduced large bias into the model.

Finally, our method’s performance does degrade as the topologies become more complex. However, the performance does not seem to degrade exponentially with the length of the chain of hidden variables, which suggests that our sample complexity analysis can be further improved.

7.2 Comparison with Mossel and Roch Algorithm

We now make a separate comparison with the spectral algorithm by (Mossel & Roch, 2006), since it only applies to case where the number of observed states S_O is the same as the number of hidden states S_H . We use the same experimental settings as in the previous section, but set $S_O = S_H = 2$. Although this method is theoretically interesting, it can perform poorly in practice.

The results are shown in Figure 3(a)-(d) (the runtime of both methods are similar, and thus not reported). Our spec-

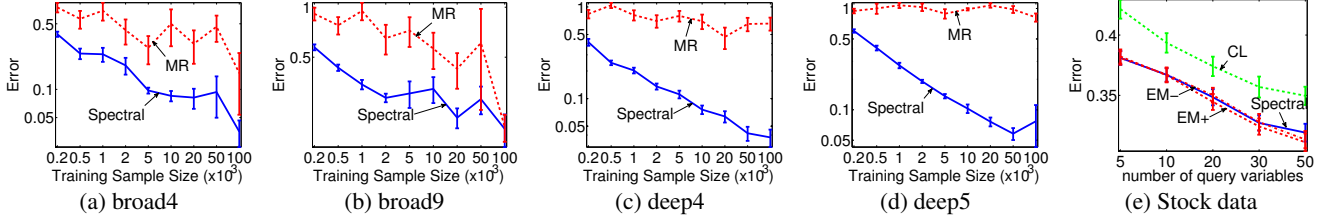


Figure 3. (a)-(d) Comparison of our spectral algorithm (Spectral) with the Mossel and Roch algorithm (MR) for 4 different latent tree topologies. The errors are plotted in log scale. (e) Comparison of our spectral algorithm (Spectral, blue line) with EMs (EM+ and EM-, red lines) and Chow-Liu based algorithm (CL, green line) on stock dataset.

tral algorithm significantly outperforms the MR algorithm on all trees for practically all sample sizes. This is because our method does not explicitly recover the CPTs, and is thus more robust. We also note that our approach is more general: it can allow for the observation state space to be larger than the hidden state space, which may be preferable in many applications where the observation space can be large (e.g., quantization of a continuous variable), but the hidden factors are simple and have lower dimensions.

7.3 Stock Trend Prediction

Finally, we evaluate our algorithm on a stock trend prediction problem. Our goal is to predict whether a stock X_i will go up or down on a particular day given the trends of a set \mathcal{E} of other stocks. We acquired closing prices of 59 stocks from 1984 to 2011, which provides us 6800 samples.² We randomly partition these samples to 6300 training points and 500 test points. Since we are only predicting whether a stock goes up or down, the data are binarized. From the training data, we learn the latent tree topology using an algorithm by Choi et al. (2010), and a fully observable Chow-Liu tree (Chow & Liu, 1968). A visualization of the learned tree topologies are in the supplementary.

We compare our spectral algorithm to EM+ and EM- using the latent tree, and with inference over the Chow-Liu tree. For the prediction task, we need to estimate the conditional, i.e., $\mathbb{P}[X_i | x_{j_1}, \dots, x_{j_{|\mathcal{E}|}}]$ and $j_1, \dots, j_{|\mathcal{E}|} \in \mathcal{E}$. This can be achieved by estimating $\mathbb{P}[x_i, x_{j_1}, \dots, x_{j_{|\mathcal{E}|}}]$ for each instantiation x_i . Then we make prediction by $\hat{x}_i = \arg\max_{x_i} \mathbb{P}[x_i, x_{j_1}, \dots, x_{j_{|\mathcal{E}|}}]$. We measure the prediction error using $\epsilon = |\hat{x}_i - x_i^*|$ where x_i^* is the true label.

We experiment with a varying number of query sizes. For each query size Q , we randomly pick Q stocks and predict the value of one stock conditioned on the other $Q - 1$, (and repeat for 50 trials). Over the entire range of query sizes, the advantage of latent tree approaches (Spectral, EM+/-) is clear over the Chow-Liu tree. Thus, the latent factors help better model the stock data in this case. Due to the small training sample size, the distinction between our method and EM is less clear.

²www.finance.yahoo.com

8 Conclusion

We have proposed a local-minimum-free spectral algorithm for latent tree models which only uses information from marginals of triples of observed variables irrespective of the maximum degree of the graph. Our algorithm is computationally efficient even for large sample sizes, and shows good performance in both synthetic and real datasets. There are many future directions: one can kernelize the method like Song et al. (2010) did for HMMs, or design spectral algorithms for loopy latent variable models.

Acknowledgements

APP is thankful for a NSF Graduate Fellowship and LS is supported by a Ray and Stephanie Lane Fellowship. This paper is also supported by NIH 1R01GM093156, NIH 1RC2HL101487, NSF IIS-0713379, NSF DBI-0546594, and an Alfred P. Sloan Fellowship to EPX.

References

- Blei, D., Ng, A., and Jordan, M. Latent dirichlet allocation. In *NIPS*, 2002.
- Choi, Myung J., Tan, Vicent Y., Anandkumar, Animashree, and Willsky, Alan S. Learning latent tree graphical models. In *arXiv:1009.2722v1*, 2010.
- Chow, C. and Liu, C. Approximating discrete probability distributions with dependence trees. *IEEE Transactions on Information Theory*, 14(3):462–467, 1968.
- Clark, A. Inference of haplotypes from pcr-amplified samples of diploid populations. *Molecular Biology and Evolution*, 7(2): 111–122, 1990.
- Dempster, A., Laird, N., and Rubin, D. Maximum likelihood from incomplete data via the em algorithm. *Journal of the Royal Statistical Society B*, 39(1):1–22, 1977.
- Hoff, Peter D., Raftery, Adrian E., and Handcock, Mark S. Latent space approaches to social network analysis. *JASA*, 97(460): 1090–1098, 2002.
- Hsu, D., Kakade, S., and Zhang, T. A spectral algorithm for learning hidden markov models. In *COLT*, 2009.
- Kolda, Tamara. and Bader, Brett. Tensor decompositions and applications. *SIAM Review*, 51(3):455–500, 2009.
- Mossel, E. and Roch, S. Learning nonsingular phylogenies and hidden markov models. *AOAP*, 16(2):583–614, 2006.
- Pearl, J. *Probabilistic reasoning in intelligent systems: networks of plausible inference*. Morgan Kaufmann, 1988.
- Rabiner, L. R. and Juang, B. H. An introduction to hidden Markov models. *IEEE ASSP Magazine*, 3(1):4–16, January 1986.
- Song, L., Boots, B., Siddiqi, S., Gordon, G., and Smola, A. Hilbert space embeddings of hidden markov models. In *ICML*, 2010.

Supplemental for Spectral Algorithm For Latent Tree Graphical Models

The supplemental contains 3 main things.

1. The first is network plots of the latent variable tree learned by [1] for the stock market data, and the Chow Liu tree to give a more intuitive explanation why latent variable trees can lead to better performance.
2. The second is a more detailed representation of the tensor representation where internal nodes are allowed to evidence variables.
3. The third is the proof of Theorem 1.

1 Latent Tree Structure for Stock Data

The latent tree structure learned by the algorithm by [1] is shown in Figure 1. The blue nodes are hidden nodes and the red nodes are observed. Note how integrating out some of these hidden nodes could lead to very large cliques. Thus it is not surprising why both our spectral method and EM perform better than Chow Liu. The Chow Liu Tree is shown in Figure 1. Note how it is forced to pick some of the observed variables as hubs even if latent variables may be more natural.

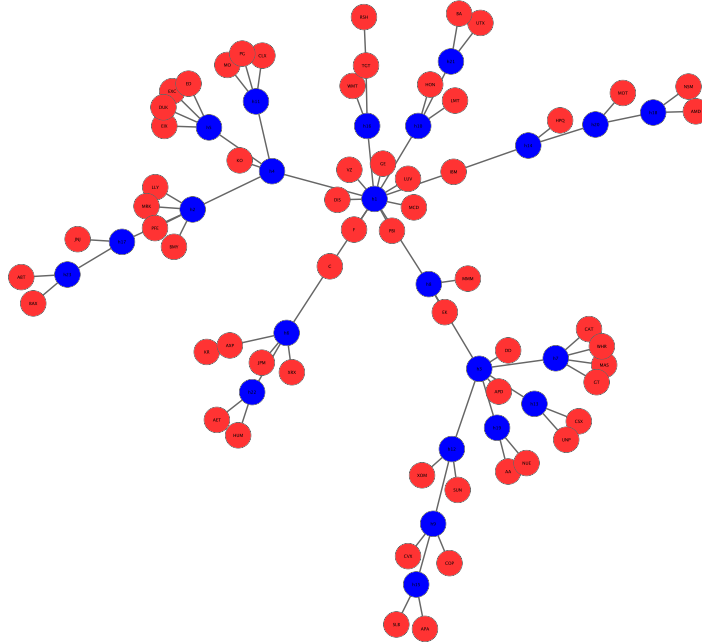


Figure 1: Latent variable tree learned by [1]. The hidden variables are in blue while the observed variables are in red. As one can see the hidden variables can model significantly more complex relationships among the observed variables.

2 More Detailed Information about Tensor Representation for LTMs

The computation of the marginal distribution of the observed variables can be expressed in terms of tensor multiplications. Basically, the information contained in each tensor will correspond to the information in a conditional probability table (CPT) of the model and the tensor multiplications implement the summations. However, there are multiple ways of rewriting the marginal distribution of the observed variables using tensor notation, and not all of them provide intuition or easy derivation to a spectral algorithm. In this section, we will derive a specific representation of the latent tree models which requires only tensors up to 3rd order and provides us a basis for deriving a spectral algorithm.

More specifically, we first select a latent or observed variable as the root node and sort the nodes in the tree in topological order. Then we associate the root node X_r with a vector \mathbf{r} related to the marginal probability of X_r . Depending on whether the root node is latent or observed, its entries are defined as

	X_r latent	X_r observed
$\mathbf{r}(k)$	$\mathbb{P}[X_r = k]$	$\delta_{kx_r} \mathbb{P}[x_r]$

where δ_{kx_r} is an indicator variable defined as $\delta_{kx_r} = 1$ if $k = x_r$ and 0 otherwise. Effectively, δ_{kx_r} sets all entries of \mathbf{r} to zero except the one corresponding to $\mathbb{P}[x_r]$.

Next, we associate each internal node X_i with a 3rd order tensor \mathcal{T}_i related the conditional probability table between X_i and its parent X_{π_i} . This tensor is diagonal in its 2nd and 3rd mode, and hence its nonzero entries can be accessed by

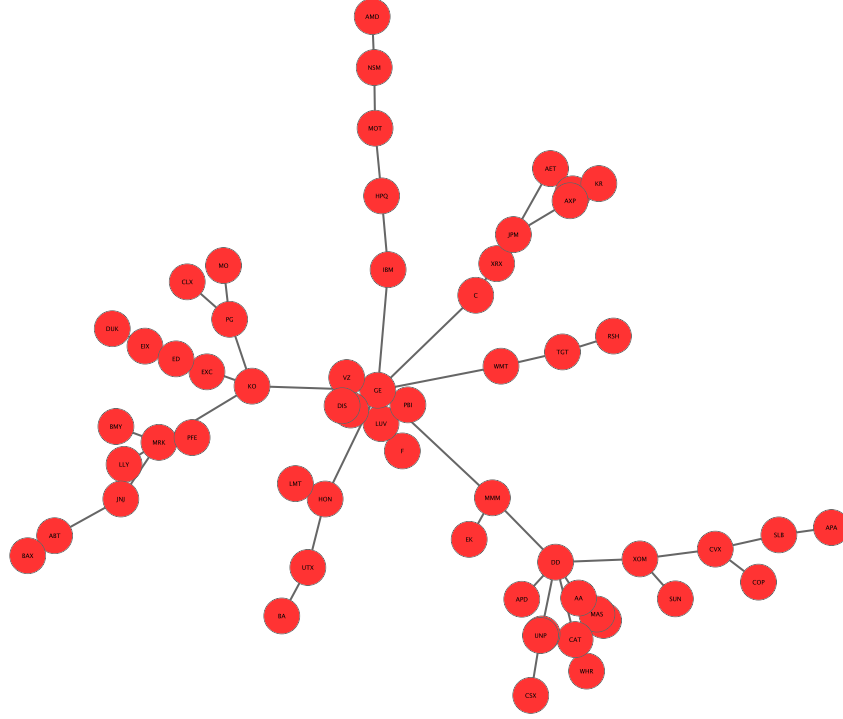


Figure 2: Tree learned by chow liu algorithm over only observed variables. Note how it is forced to pick some of the observed variables as hubs even if latent variables may be more natural.

two indices k and l . Depending on whether the internal node and its parent are latent or observed variables, the nonzero entries of \mathcal{T}_i are defined as

$\mathcal{T}_i(k, l)$	X_{π_i} latent	X_{π_i} observed
X_i latent	$\mathbb{P}[X_i = k X_{\pi_i} = l]$	$\delta_{l x_{\pi_i}} \mathbb{P}[X_i = k x_{\pi_i}]$
X_i observed	$\delta_{k x_i} \mathbb{P}[x_i X_{\pi_i} = l]$	$\delta_{k x_i} \delta_{l x_{\pi_i}} \mathbb{P}[x_i x_{\pi_i}]$

where $\delta_{k x_i}$ and $\delta_{l x_{\pi_i}}$ are also indicator variables. Effectively, the indicator variables zero out further entries in \mathcal{T}_i for those values that are not equal to the actual observation.

Last, we associate each leaf node x_i , which is always observed, with a diagonal matrix \mathbf{M}_i related to the likelihood of x_i . Depending on whether the parent of x_i is latent or observed, the diagonal entries of \mathbf{M}_i are defined as

	X_{π_i} latent	X_{π_i} observed
$\mathbf{M}_i(k, k)$	$\mathbb{P}[x_i X_{\pi_i} = k]$	$\delta_{k x_{\pi_i}} \mathbb{P}[x_i x_{\pi_i}]$

Let \mathbf{M}_i defined above be the messages passed from the leaf nodes to their parents. We can show that the marginal probability of the observed variables can be computed recursively using a message passing algorithm: each node in the tree sends message to its parent according to the reverse topological order of the nodes, and the final messages are aggregated in the root to give the desired quantity.

More formally, the outgoing message from an internal node X_i to its parent can be computed as

$$\mathbf{M}_i = \mathcal{T}_i \bar{\times}_1 (\mathbf{M}_{j_1} \mathbf{M}_{j_2} \dots \mathbf{M}_{j_J} \mathbf{1}_i) \quad (1)$$

where $j_1, j_2, \dots, j_J \in \chi_i$ range over all children of X_i ($J = |\chi_i|$). The $\mathbf{1}_i$ is a vector of all ones with suitable size, and it is used to reduce the incoming messages (all are diagonal matrices) to a single vector. The computation in (1) essentially implements the message update we often see in ordinary message passing algorithm, *i.e.*,

$$m_i[x_{\pi_i}] = \sum_{x_i} \mathbb{P}[x_i | x_{\pi_i}] m_{j_1}[x_i] \dots m_{j_J}[x_i], \quad (2)$$

where $m_j[x_i]$ represents incoming messages to X_i (or intermediate results of the marginalization operation by summing out all latent variables in subtree \mathcal{T}_j). The $\mathbf{M}_{j_1} \mathbf{M}_{j_2} \dots \mathbf{M}_{j_J} \mathbf{1}_i$ corresponds to aggregating all incoming messages $m_{j_1}[x_i] \dots m_{j_J}[x_i]$, and the $\mathcal{T}_i \bar{\times}_1 *$ corresponds to $\sum_{x_i} \mathbb{P}[x_i | x_{\pi_i}] *$.

At the root node, all incoming messages are combined to produce the final joint probability, *i.e.*,

$$\mathbb{P}[x_1, \dots, x_O] = \mathbf{r}^\top (\mathbf{M}_{j_1} \mathbf{M}_{j_2} \dots \mathbf{M}_{j_J} \mathbf{1}_r). \quad (3)$$

Here $\mathbf{r}^\top *$ basically implements the operation $\sum_{x_r} \mathbb{P}[x_r] *$, which sums out the root variable.

3 Notation for Proof of Theorem 1

We now proceed to prove Theorem 1. $\|\cdot\|_2$ refers to spectral norm for matrices and tensors (but normal euclidean norm for vectors). $\|\cdot\|_1$ refers to induced 1 norm for matrices and tensors (max column sum), (but normal l1 norm for vectors). $\|\cdot\|_F$ refers to Frobenius norm.

The tensor spectral norm (for 3 dimensions) is defined in [4]:

$$\|\mathcal{T}\|_2 = \sup_{\|v_i\|_2 \leq 1} \mathcal{T} \bar{\times}_3 v_3 \bar{\times}_2 v_2 \bar{\times}_1 v_1 \quad (4)$$

We will define the induced 1-norm of a tensor as $\|\mathcal{T}\|_{1,1} = \sup_{\|v\|_1 \leq 1} \|\mathcal{T} \bar{\times}_1 v\|_1$ using the ℓ_1 norm of a matrix (*i.e.*, $\|A\|_1 = \sup_{\|v\|_1 \leq 1} \|Av\|_1$).

For more information about matrix norms see [2].

In general, we suppress the actual ssubscripts/superscripts on \mathbf{U} and \mathbf{O} . It is implied that \mathbf{U} and \mathbf{O} can often be different depending on the transform being considered. However, this makes the notation very messy. It will generally be clear from context which \mathbf{U} and \mathbf{O} are being referred to. When it is not we will arbitrarily index them 1, 2, ..., so that it is clear which corresponds to which.

In general, for simplicity of exposition, we assume that all internal nodes in the tree are unobserved, and all leaves are observed (since this is the hardest case).

The proof generally follows the technique of HKZ [3], but has key differences due to the tree topology instead of the HMM.

We define $\tilde{\mathbf{M}}_i = (\hat{\mathbf{U}}^T \mathbf{O})^{-1} \mathbf{M}_i (\hat{\mathbf{U}}^T \mathbf{O})$. Then as long as $(\hat{\mathbf{U}}^T \mathbf{O})$ is invertible, $(\hat{\mathbf{U}}^T \mathbf{O})^{-1} \tilde{\mathbf{M}}_i (\hat{\mathbf{U}}^T \mathbf{O}) = \mathbf{M}_i$. (We admit this is a slight abuse of notation, since $\tilde{\mathbf{M}}_i$ was previously defined to be $(\mathbf{U}^T \mathbf{O})^{-1} \mathbf{M}_i (\mathbf{U}^T \mathbf{O})$, but as long as $(\hat{\mathbf{U}}^T \mathbf{O})$ is invertible it doesn't really matter whether it equals $(\mathbf{U}^T \mathbf{O})$ or not for the purposes of this proof). The other quantities are defined similarly.

We seek to prove the following theorem:

Theorem 1 *Pick any $\epsilon > 0, \delta < 1$. Let*

$$N \geq O \left(\frac{1}{\epsilon^2} \left(\frac{(d_{max} m)^{\ell+1} S_O}{\min_i \sigma_{SH}(\mathbf{O}_i)^2 \min_{i \neq j} \sigma_{SH}(\mathbf{P}_{i,j})^4} \right) \right) \log \frac{|\mathcal{O}|}{\delta} \quad (5)$$

Then with probability $1 - \delta$

$$\sum_{x_1, \dots, x_O} \left| \hat{\mathbb{P}}[x_1, \dots, x_O] - \mathbb{P}[x_1, \dots, x_O] \right| \leq \epsilon \quad (6)$$

In many cases, if the frequency of the observation symbols follow certain distributions, than the dependence on S_O can be removed as showed in HKZ [3]. That observation can easily be incorporated into our theorem if desired.

4 Concentration Bounds

$$\epsilon_i = \|\hat{\mathbf{P}}_i - \mathbf{P}_i\|_F \quad (7)$$

$$\epsilon_{i,j} = \|\hat{\mathbf{P}}_{i,j} - \mathbf{P}_{i,j}\|_F \quad (8)$$

$$\epsilon_{x,i,j} = \|\hat{\mathbf{P}}_{x,i,j} - \mathbf{P}_{x,i,j}\|_F \quad (9)$$

$$\epsilon_{i,j,k} = \|\hat{\mathbf{P}}_{i,j,k} - \mathbf{P}_{i,j,k}\|_F \quad (10)$$

(x denotes a fixed element while i, j, k are over indices).

As the number of samples N gets large, we expect these quantities to be small.

Lemma 1 (variant of HKZ [3]) *If the algorithm independently samples N observation triples from the tree, then with*

probability at least $1 - \eta$.

$$\epsilon_i \leq \sqrt{\frac{C}{N} \ln \frac{|\mathcal{O}|}{\eta}} + \sqrt{\frac{1}{N}} \quad (11)$$

$$\epsilon_{i,j} \leq \sqrt{\frac{C}{N} \ln \frac{|\mathcal{O}|}{\eta}} + \sqrt{\frac{1}{N}} \quad (12)$$

$$\epsilon_{i,j,k} \leq \sqrt{\frac{C}{N} \ln \frac{|\mathcal{O}|}{\eta}} + \sqrt{\frac{1}{N}} \quad (13)$$

$$\max_x \epsilon_{i,x,j} \leq \sqrt{\frac{C}{N} \ln \frac{|\mathcal{O}|}{\eta}} + \sqrt{\frac{1}{N}} \quad (14)$$

where C is some constant (from the union bound over $O(V^3)$). (V is the total number of observed variables in the tree). The proof is the same as that of HKZ [3] except the union bound is larger. The last bound can be made tighter, identical to HKZ, but for simplicity we do not pursue that approach here.

5 Eigenvalue Bounds

Basically this is Lemma 9 in HKZ [3], which is stated below for completeness:

Lemma 2 Suppose $\epsilon_{i,j} \leq \varepsilon \times \sigma_{S_H}(\mathbf{P}_{i,j})$ for some $\varepsilon < 1/2$. Let $\varepsilon_0 = \epsilon_{i,j}^2 / ((1 - \varepsilon)\sigma_{S_H}(\mathbf{P}_{i,j}))^2$. Then:

1. $\varepsilon_0 < 1$
2. $\sigma_{S_H}(\hat{\mathbf{U}}^T \hat{\mathbf{P}}_{i,j}) \geq (1 - \varepsilon)\sigma_{S_H}(\mathbf{P}_{i,j})$
3. $\sigma_{S_H}(\hat{\mathbf{U}}^T \mathbf{P}_{i,j}) \geq \sqrt{1 - \varepsilon_0} \sigma_{S_H}(\mathbf{P}_{i,j})$
4. $\sigma_{S_H}(\hat{\mathbf{U}}^T \mathbf{O}) \geq \sqrt{1 - \varepsilon} \sigma_{S_H}(\mathbf{O})$

The proof is in HKZ [3].

6 Bounding the Transformed Quantities

If Lemma 2 holds then $(\hat{\mathbf{U}}^T \mathbf{O})$ is invertible. Thus, if we define $\tilde{\mathbf{M}}_i = (\hat{\mathbf{U}}^T \mathbf{O})^{-1} \mathbf{M}_i (\hat{\mathbf{U}}^T \mathbf{O})$. Then clearly, $(\hat{\mathbf{U}}^T \mathbf{O})^{-1} \tilde{\mathbf{M}}_i (\hat{\mathbf{U}}^T \mathbf{O}) = \mathbf{M}_i$. (We admit this is a slight abuse of notation, since $\tilde{\mathbf{M}}_i$ is previously defined to be $(\mathbf{U}^T \mathbf{O})^{-1} \mathbf{M}_i (\mathbf{U}^T \mathbf{O})$, but as long as $(\hat{\mathbf{U}}^T \mathbf{O})$ is invertible it doesn't really matter whether it equals $(\mathbf{U}^T \mathbf{O})$ or not for the purposes of this proof). The other quantities are defined similarly.

We seek to bound the following four quantities:

$$\delta_{one}^i = \|(\hat{\mathbf{U}}^T \mathbf{O})(\hat{\mathbf{1}}_i - \tilde{\mathbf{1}}_i)\|_1 \quad (15)$$

$$\gamma_i = \|(\hat{\mathcal{T}}_i - \tilde{\mathcal{T}}_i) \times_1 (\hat{\mathbf{U}}_1^T \mathbf{O}_1)^{-1} \times_2 (\hat{\mathbf{U}}_2^T \mathbf{O}_2) \times_3 (\hat{\mathbf{U}}_3^T \mathbf{O}_3)^{-1}\|_2 \quad (16)$$

$$\delta_{root} = \|(\hat{\mathbf{r}} - \tilde{\mathbf{r}})^T (\hat{\mathbf{U}}^T \mathbf{O})^{-1}\|_\infty \quad (17)$$

$$\Delta_i = \sum_{\mathbf{x}_i} \|(\hat{\mathbf{U}}_1^T \mathbf{O}_1)(\hat{\mathbf{M}}_i - \tilde{\mathbf{M}}_i)(\hat{\mathbf{U}}_2^T \mathbf{O}_2)^{-1}\|_1 \quad (18)$$

Here \mathbf{x}_i denotes all observations that are in the subtree of node i (since i may be hidden or observed). Sometimes we like to distinguish between when i is observed and when i is hidden. Thus, we sometimes refer to the quantity Δ_i^{obs} and Δ_i^{hidden} for when i is observed or hidden respectively.

Again note that the numbering in $(\hat{\mathbf{U}}_1^T \mathbf{O}_1)$ and $(\hat{\mathbf{U}}_2^T \mathbf{O}_2)$ is just there to avoid confusion in the same equation (In reality there are many \mathbf{U} 's and \mathbf{O} 's).

Lemma 3 Assume $\epsilon_{i,j} \leq \sigma_{S_H}(\mathbf{P}_{i,j})/3$ for all $i \neq j$. Then

$$\delta_{root} \leq \frac{2\epsilon_r}{\sqrt{3}\sigma_{S_H}(\mathbf{O})} \quad (19)$$

$$\delta_{one}^i \leq 4\sqrt{S_H} \left(\frac{\epsilon_{i,j}}{\sigma_{S_H}(\mathbf{P}_{i,j})^2} + \frac{\epsilon_i}{\sqrt{3}\sigma_{S_H}(\mathbf{P}_{i,j})} \right) \quad (20)$$

$$\gamma_i \leq \frac{4\sqrt{S_H}}{\sigma_{S_H}(\mathbf{O})} \left(\frac{\epsilon_{m,j}}{\sigma_{S_H}(\mathbf{P}_{i,j})} + \frac{\epsilon_{m,j,k}}{\sqrt{3}\sigma_{S_H}(\mathbf{P}_{i,j})} \right) \quad (21)$$

$$\Delta_i^{hidden} \leq \left((1 + \gamma_i) \prod_{k=1}^J (1 + \Delta_{j_k}) \delta_{one}^i + (1 + \gamma_i) m \prod_{k=1}^J (1 + \Delta_{j_k}) - m \right) \quad (22)$$

$$\Delta_i^{obs} \leq 4 \frac{\sqrt{S_H}}{\sigma_{S_H}(\mathbf{O})} \left(\frac{\epsilon_{j,i}}{(\sigma_{S_H}(\mathbf{P}_{j,i}))^2} + \frac{\sum_{x_i} \epsilon_{m,x_i,j}}{\sqrt{3}\sigma_{S_H}(\mathbf{P}_{j,i})} \right) \quad (23)$$

The main challenge in this part is Δ_v and γ_v^{hidden} . The rest are similar to HKZ. However, we go through the other bounds to be more explicit about some of the properties used, since sometimes we have used different norms etc.

6.1 δ_{root}

We note that $\hat{\mathbf{r}} = \hat{\mathbf{U}}^T \hat{\mathbf{P}}_i$ and similarly $\tilde{\mathbf{r}} = \hat{\mathbf{U}}^T \mathbf{P}_i$.

$$\delta_{root} = \|(\hat{\mathbf{r}} - \tilde{\mathbf{r}})^T (\hat{\mathbf{U}}^T \mathbf{O})^{-1}\|_{\infty} \leq \|\hat{\mathbf{U}}\|_2 \|\hat{\mathbf{P}}_i - \mathbf{P}_i\|_2 \|(\hat{\mathbf{U}}^T \mathbf{O})^{-1}\|_2 \quad (24)$$

$$\leq \|\hat{\mathbf{P}}_i - \mathbf{P}_i\|_2 \|(\hat{\mathbf{U}}^T \mathbf{O})^{-1}\|_2 \leq \frac{\epsilon_i}{\sigma_{S_H}(\hat{\mathbf{U}}^T \mathbf{O})} \quad (25)$$

The first inequality follows from the relationship between ℓ_{∞} and ℓ_2 norm and submultiplicativity. The second follows from a matrix perturbation bound given in Lemma 89. We also use the fact that since $\hat{\mathbf{U}}$ is orthonormal it has spectral norm 1.

Assuming that $\epsilon_{i,j} \leq \sigma_{S_H}(\mathbf{P}_{i,j})/3$ gives $\delta_{root} \leq \frac{2\epsilon_r}{\sqrt{3}\sigma_{S_H}(\mathbf{O})}$ by Lemma 2.

6.2 δ_{one}^i

$$\delta_{one}^i = \|(\hat{\mathbf{U}}^T \mathbf{O})(\hat{\mathbf{1}}_i - \tilde{\mathbf{1}}_i)\|_1 \leq \sqrt{S_H} \|\hat{\mathbf{U}}\|_2 \|\mathbf{O}\|_2 \|\hat{\mathbf{1}}_i - \tilde{\mathbf{1}}_i\|_2 \quad (26)$$

$$= \sqrt{S_H} \|\hat{\mathbf{1}}_i - \tilde{\mathbf{1}}_i\|_2 = \sqrt{S_H} \|\hat{\mathbf{1}}_i^T - \tilde{\mathbf{1}}_i^T\|_2 \quad (27)$$

Here we have converted ℓ_1 norm to ℓ_2 norm, used submultiplicativity, the fact that $\hat{\mathbf{U}}$ is orthonormal so has spectral norm 1, and that \mathbf{O} is a conditional probability matrix and therefore also has spectral norm 1.

We note that $\hat{\mathbf{1}}_i^T = (\hat{\mathbf{U}}^T \hat{\mathbf{P}}_{m,j})^+ \hat{\mathbf{P}}_j$ and similarly $\tilde{q}_{one} = (\hat{\mathbf{U}}^T \mathbf{P}_{j,i})^+ \mathbf{P}_i$, where m and j are a particular pair of observations described in the main paper.

$$\|\hat{\mathbf{1}}_i^T - \tilde{\mathbf{1}}_i^T\|_2 = \|(\hat{\mathbf{U}}^T \hat{\mathbf{P}}_{m,j})^+ \hat{\mathbf{P}}_j - (\hat{\mathbf{U}}^T \mathbf{P}_{m,j})^+ \mathbf{P}_j\|_2 \quad (28)$$

$$= \|(\hat{\mathbf{U}}^T \hat{\mathbf{P}}_{m,j})^+ \hat{\mathbf{P}}_j - (\hat{\mathbf{U}}^T \mathbf{P}_{m,j})^+ \hat{\mathbf{P}}_j + (\hat{\mathbf{U}}^T \mathbf{P}_{m,j})^+ \hat{\mathbf{P}}_j - (\hat{\mathbf{U}}^T \mathbf{P}_{m,j})^+ \mathbf{P}_j\|_2 \quad (29)$$

$$\leq \|(\hat{\mathbf{U}}^T \hat{\mathbf{P}}_{m,j})^+ \hat{\mathbf{P}}_j - (\hat{\mathbf{U}}^T \mathbf{P}_{m,j})^+ \hat{\mathbf{P}}_j\|_2 + \|(\hat{\mathbf{U}}^T \hat{\mathbf{P}}_{m,j})^+ \hat{\mathbf{P}}_j - (\hat{\mathbf{U}}^T \mathbf{P}_{m,j})^+ \mathbf{P}_j\|_2 \quad (30)$$

$$\leq \|(\hat{\mathbf{U}}^T \hat{\mathbf{P}}_{m,j})^+ - (\hat{\mathbf{U}}^T \mathbf{P}_{m,j})^+\|_2 \|\hat{\mathbf{P}}_j\|_1 + \|(\hat{\mathbf{U}}^T \hat{\mathbf{P}}_{m,j})^+ - (\hat{\mathbf{U}}^T \mathbf{P}_{m,j})^+\|_2 \|\hat{\mathbf{P}}_j - \mathbf{P}_j\|_2 \quad (31)$$

$$\leq \frac{1 + \sqrt{5}}{2} \times \frac{\epsilon_{m,j}}{\min(\sigma_{S_H}(\hat{\mathbf{P}}_{m,j}), \sigma_{S_H}(\mathbf{P}_{m,j}^T \hat{\mathbf{U}}))^2} + \frac{\epsilon_j}{\sigma_{S_H}(\mathbf{P}_{m,j}^T \hat{\mathbf{U}})} \quad (32)$$

where we have used the triangle inequality in the first inequality and the submultiplicative property of matrix norms in the second. The last inequality follows by matrix perturbation bounds. Thus using the assumption that $\epsilon_{i,j} \leq \sigma_{S_H}(\mathbf{P}_{i,j})/3$, we get that

$$\delta_{one} \leq 4\sqrt{S_H} \left(\frac{\epsilon_{i,j}}{\sigma_{S_H}(\mathbf{P}_{i,j})^2} + \frac{\epsilon_i}{\sqrt{3}\sigma_{S_H}(\mathbf{P}_{i,j})} \right) \quad (33)$$

6.3 Tensor

We will define the induced norm of a tensor as $\|\mathcal{T}\|_{1,1} = \sup_{\|v\|_1 \leq 1} \|\mathcal{T} \bar{\times}_1 v\|_1$ using the ℓ_1 norm of a matrix (*i.e.*, $\|A\|_1 = \sup_{\|v\|_1 \leq 1} \|Av\|_1$).

Recall that $\tilde{\mathcal{T}}_i = \mathcal{T}_i \times_1 (\hat{U}_1^T \mathbf{O}_1) \times_2 (\hat{U}_2^T \mathbf{O}_2)^{-1} \times_3 (\hat{U}_3^T \mathbf{O}_3) = \mathbf{P}_{m,j,k} \times_1 \hat{U}_1^T \times_2 (\mathbf{P}_{l,j} \hat{U}_2)^+ \times_3 \hat{U}_3^T$. Similarly, $\hat{\mathcal{T}}_i = \mathbf{P}_{m,j,k} \times_1 \hat{U}_1^T \times_2 (\mathbf{P}_{l,j} \hat{U}_2)^+ \times_3 \hat{U}_3^T$.

$$\|(\hat{\mathcal{T}}_i - \tilde{\mathcal{T}}_i) \times_1 (\hat{U}_1^T \mathbf{O}_1)^{-1} \times_2 (\hat{U}_2^T \mathbf{O}_2) \times_3 (\hat{U}_3^T \mathbf{O}_3)^{-1}\|_{1,1} \leq \frac{\sqrt{S_H}}{\sigma_{S_H}(\mathbf{O})} \|\hat{\mathcal{T}}_i - \tilde{\mathcal{T}}_i\|_2 \quad (34)$$

This is because both \hat{U} and \mathbf{O} have spectral norm one and the $\sqrt{S_H}$ factor is the cost of converting from 1 norm to spectral norm.

$$\|\hat{\mathcal{T}}_i - \tilde{\mathcal{T}}_i\|_2 = \|\hat{\mathbf{P}}_{m,j,k} \times_1 \hat{U}_1^T \times_2 (\hat{\mathbf{P}}_{l,j} \hat{U}_2)^+ \times_3 \hat{U}_3^T - \mathbf{P}_{m,j,k} \times_1 \hat{U}_1^T \times_2 (\mathbf{P}_{l,j} \hat{U}_2)^+ \times_3 \hat{U}_3^T\|_2 \quad (35)$$

$$= \|\hat{\mathbf{P}}_{m,j,k} \times_1 \hat{U}_1^T \times_2 \hat{U}_2^T \times_3 (\hat{\mathbf{P}}_{l,j} \hat{U}_3)^+ - \hat{\mathbf{P}}_{m,j,k} \times_1 \hat{U}_1^T \times_2 \hat{U}_2^T \times_3 (\mathbf{P}_{l,j} \hat{U}_3)^+\|_2 \quad (36)$$

$$+ \|\hat{\mathbf{P}}_{m,j,k} \times_1 \hat{U}_1^T \times_2 (\mathbf{P}_{l,j} \hat{U}_2)^+ \times_3 \hat{U}_3^T - \mathbf{P}_{m,j,k} \times_1 \hat{U}_1^T \times_2 (\mathbf{P}_{l,j} \hat{U}_3)^+ \times_3 \hat{U}_3^T\|_2 \quad (37)$$

$$= \|\hat{\mathbf{P}}_{m,j,k} \times_1 \hat{U}_1^T \times_2 ((\hat{\mathbf{P}}_{l,j} \hat{U}_2)^+ - (\mathbf{P}_{l,j} \hat{U}_3)^+) \times_3 \hat{U}_3^T\|_2 \quad (38)$$

$$+ \|(\hat{\mathbf{P}}_{l,j,k} \times_1 \hat{U}_1^T \times_3 \hat{U}_3^T - \mathbf{P}_{l,j,k} \times_1 \hat{U}_1^T \times_3 \hat{U}_3^T) \times_2 (\mathbf{P}_{l,j} \hat{U}_2)^+\|_2 \quad (39)$$

$$= \|\hat{\mathbf{P}}_{m,j,k}\|_2 \frac{1+\sqrt{5}}{2} \frac{\epsilon_{l,j}}{\min(\sigma_{S_H}(\hat{\mathbf{P}}_{l,j}), \sigma_{S_H}(\hat{U}^T \mathbf{P}_{l,j}))} + \frac{\epsilon_{m,j,k}}{\sigma_{S_H}(\hat{U}^T \mathbf{P}_{l,j})} \quad (40)$$

It is clear that $\|\hat{\mathbf{P}}_{m,j,k}\|_2 \leq \|\hat{\mathbf{P}}_{m,j,k}\|_F \leq 1$.

Using the fact that $\epsilon_{i,j} \leq \sigma_{S_H}(\mathbf{P}_{i,j})/3$ gives us the following bound:

$$\gamma_v \leq \frac{4\sqrt{S_H}}{\sigma_{S_H}(\mathbf{O})} \left(\frac{\epsilon_{i,j}}{\sigma_{S_H}(\mathbf{P}_{i,j})} + \frac{\epsilon_{i,j,k}}{\sqrt{3}\sigma_{S_H}(\mathbf{P}_{i,j})} \right) \quad (41)$$

6.4 Bounding Δ_i

We now seek to bound $\Delta_i = \sum_{\mathbf{x}_i} \|(\hat{U}_1^T \mathbf{O}_1)(\hat{\mathbf{M}}_i - \tilde{\mathbf{M}}_i)(\hat{U}_2^T \mathbf{O}_2)^{-1}\|_1$. There are two cases: either i is a leaf or it is not.

6.4.1 i is leaf node

In this case our proof simply follows from HKZ [3] and is repeated here for convenience.

$$\|(\hat{U}_1^T \mathbf{O}_1)(\hat{\mathbf{M}} - \tilde{\mathbf{M}})(\hat{U}_2^T \mathbf{O}_2)^{-1}\|_1 \leq \sqrt{S_H} \|\mathbf{O}_1\|_1 \|(\hat{\mathbf{M}}_i - \tilde{\mathbf{M}}_i)(\hat{U}_2^T \mathbf{O}_2)^{-1}\|_2 \quad (42)$$

$$\leq \sqrt{S_H} \frac{\|\hat{\mathbf{M}}_i - \tilde{\mathbf{M}}_i\|_2}{\sigma_{S_H}(\hat{U}^T \mathbf{O})} \quad (43)$$

Note that $\hat{\mathbf{M}}_i = (\hat{\mathbf{P}}_{j,i} \hat{U}_1)^{-1} \hat{\mathbf{P}}_{m,x_i,j} \hat{U}_2$ and $\tilde{\mathbf{M}}_i = (\mathbf{P}_{j,i} \hat{U}_1)^{-1} \mathbf{P}_{m,x_i,j} \hat{U}_2$.

$$\|\hat{\mathbf{M}}_i - \tilde{\mathbf{M}}_i\|_2 = \|(\hat{\mathbf{P}}_{j,i} \hat{U}_1)^{-1} \hat{\mathbf{P}}_{m,x_i,j} \hat{U}_2 - (\mathbf{P}_{j,i} \hat{U}_1)^{-1} \mathbf{P}_{m,x_i,j} \hat{U}_2\|_2 \quad (44)$$

$$= \|(\hat{\mathbf{P}}_{j,i} \hat{U}_1)^{-1} \hat{\mathbf{P}}_{m,x_i,j} \hat{U}_2 - (\mathbf{P}_{j,i} \hat{U}_1)^{-1} \hat{\mathbf{P}}_{m,x_i,j} \hat{U}_2 + (\mathbf{P}_{j,i} \hat{U}_1)^{-1} \hat{\mathbf{P}}_{m,x_i,j} \hat{U}_2 - \hat{U}_1^T \mathbf{P}_{x_i,j} (\hat{U}_2^T \mathbf{P}_{i,j})^{-1}\|_2 \quad (45)$$

$$\leq \|((\hat{\mathbf{P}}_{j,i} \hat{U}_1)^{-1} - (\mathbf{P}_{j,i} \hat{U}_1)^{-1}) \hat{\mathbf{P}}_{m,x_i,j} \hat{U}_2\|_2 + \|(\mathbf{P}_{j,i} \hat{U}_1)^{-1} (\hat{\mathbf{P}}_{m,x_i,j} \hat{U}_2 - \mathbf{P}_{m,x_i,j} \hat{U}_2)\|_2 \quad (46)$$

$$\leq \|\hat{\mathbf{P}}_{m,x_i,j}\|_2 \frac{1+\sqrt{5}}{2} \frac{\epsilon_{j,i}}{\min(\sigma_{S_H}(\hat{\mathbf{P}}_{j,i}), \sigma_{S_H}(\hat{U}^T \mathbf{P}_{j,i}))} + \frac{\epsilon_{m,x_i,j}}{\sigma_{S_H}(\hat{U}^T \mathbf{P}_{j,i})} \quad (47)$$

$$\leq \mathbb{P}[x_i = x] \frac{1+\sqrt{5}}{2} \frac{\epsilon_{j,i}}{\min(\sigma_{S_H}(\hat{\mathbf{P}}_{j,i}), \sigma_{S_H}(\mathbf{P}_{j,i} \hat{U}))^2} + \frac{\epsilon_{m,x_i,j}}{\sigma_{S_H}(\mathbf{P}_{j,i} \hat{U})} \quad (48)$$

where the first inequality follows from the triangle inequality, and the second uses matrix perturbation bounds (and the fact that spectral norm of \hat{U} is 1).

The final inequality follows from the fact that spectral norm is less than frobenius norm which is less than l1 norm:

$$\|\hat{\mathbf{P}}_{m,x_i,j}\| \leq \sqrt{\sum_{m,j} [\hat{\mathbf{P}}_{m,x_i,j}]_{m,j}^2} \leq \sum_{m,j} [\mathbf{P}_{m,x_i,j}]_{m,j} \leq \mathbb{P}[x_i = x] \quad (49)$$

The first inequality follows from relation between 1 operator norm and 2 operator norm. Because \mathbf{O} is a conditional probability matrix $\|\mathbf{O}\|_1 = 1$ (i.e. the max column sum is 1).

Using the fact that $\epsilon_{i,j} \leq \sigma_{S_H}(\mathbf{P}_{i,j})/3$ gives us the following bound:

$$\Delta_{i,x} \leq 4 \frac{\sqrt{S_H}}{\sigma_{S_H}(\mathbf{O})} \left(\mathbb{P}[x_i = x] \frac{\epsilon_{m,x_i,j}}{(\sigma_{S_H}(\mathbf{P}_{j,i}))^2} + \frac{\epsilon_{m,x_i,j}}{\sqrt{3}\sigma_{S_H}(\mathbf{P}_{j,i})} \right) \quad (50)$$

Summing over v would give

$$\Delta_i \leq 4 \frac{\sqrt{S_H}}{\sigma_{S_H}(\mathbf{O})} \left(\frac{\epsilon_{j,i}}{(\sigma_{S_H}(\mathbf{P}_{j,i}))^2} + \frac{\sum_{x_i} \epsilon_{m,x_i,j}}{\sqrt{3}\sigma_{S_H}(\mathbf{P}_{j,i})} \right) \quad (51)$$

6.4.2 i is not a leaf node

Let $\hat{\mathbf{m}}_{J:1} = \hat{\mathbf{M}}_J \dots \hat{\mathbf{M}}_1 \hat{\mathbf{1}}_i$ and $\tilde{\mathbf{m}}_{J:1} = \tilde{\mathbf{M}}_J \dots \tilde{\mathbf{M}}_1 \tilde{\mathbf{1}}_i$

$$\sum_{\mathbf{x}_i} \|(\hat{\mathbf{U}}_2^T \mathbf{O}_2)(\hat{\mathbf{M}}_i - \tilde{\mathbf{M}}_i)(\hat{\mathbf{U}}_3^T \mathbf{O}_3)^{-1}\|_1 \quad (52)$$

$$= \sum_{\mathbf{x}_i} \|(\hat{\mathbf{U}}_2^T \mathbf{O}_2)(\hat{\mathcal{T}}_i \times_1 \hat{\mathbf{M}}_u \dots \hat{\mathbf{M}}_1 \hat{\mathbf{1}}_i - \tilde{\mathcal{T}}_i \times_1 \tilde{\mathbf{M}}_u \dots \tilde{\mathbf{M}}_1 \tilde{\mathbf{1}}_i)(\hat{\mathbf{U}}_3^T \mathbf{O}_3)^{-1}\|_1 \quad (53)$$

$$= \sum_{\mathbf{x}_i} \|(\hat{\mathbf{U}}_2^T \mathbf{O}_2)(\hat{\mathcal{T}}_i \bar{\times}_1 \hat{\mathbf{m}}_{J:1} - \tilde{\mathcal{T}}_i \bar{\times}_1 \tilde{\mathbf{m}}_{J:1})(\hat{\mathbf{U}}_3^T \mathbf{O}_3)^{-1}\|_1 \quad (54)$$

$$= \sum_{\mathbf{x}_i} \|(\hat{\mathbf{U}}^T \mathbf{O}) \left((\hat{\mathcal{T}}_i - \tilde{\mathcal{T}}_i) \bar{\times}_1 \tilde{\mathbf{m}}_{J:1} + (\hat{\mathcal{T}}_i - \tilde{\mathcal{T}}_i) \bar{\times}_1 (\hat{\mathbf{m}}_{J:1} - \tilde{\mathbf{m}}_{J:1}) + \tilde{\mathcal{T}}_i \bar{\times}_1 (\hat{\mathbf{m}}_{J:1} - \tilde{\mathbf{m}}_{J:1}) \right) (\hat{\mathbf{U}}_3^T \mathbf{O}_3)^{-1}\|_1 \quad (55)$$

$$\leq \sum_{\mathbf{x}_i} \|(\hat{\mathcal{T}}_i - \tilde{\mathcal{T}}_i) \times_1 (\hat{\mathbf{U}}_1^T \mathbf{O}_1)^{-1} \times_2 (\hat{\mathbf{U}}_2^T \mathbf{O}_2) \times_3 (\hat{\mathbf{U}}_3^T \mathbf{O}_3)^{-1}\|_{1,1} \|(\hat{\mathbf{U}}_1^T \mathbf{O}_1) \tilde{\mathbf{m}}_{J:1}\|_1 \quad (56)$$

$$+ \sum_{\mathbf{x}_i} \|(\hat{\mathbf{U}}_1^T \mathbf{O}_1)(\hat{\mathbf{m}}_{J:1} - \tilde{\mathbf{m}}_{J:1})\|_1 \|(\hat{\mathcal{T}}_i - \tilde{\mathcal{T}}_i) \times_1 (\hat{\mathbf{U}}_1^T \mathbf{O}_1)^{-1} \times_2 (\hat{\mathbf{U}}_2^T \mathbf{O}_2) \times_3 (\hat{\mathbf{U}}_3^T \mathbf{O}_3)^{-1}\|_{1,1} \quad (57)$$

$$+ \sum_{\mathbf{x}_i} \|\tilde{\mathcal{T}}_i \times_1 (\hat{\mathbf{U}}_1^T \mathbf{O}_1)^{-1} \times_2 (\hat{\mathbf{U}}_2^T \mathbf{O}_2) \times_3 (\hat{\mathbf{U}}_3^T \mathbf{O}_3)^{-1}\|_{1,1} \|(\hat{\mathbf{U}}_1^T \mathbf{O}_1)(\hat{\mathbf{m}}_{J:1} - \tilde{\mathbf{m}}_{J:1})\|_1 \quad (58)$$

First term is bounded by:

$$\left\| (\hat{\mathcal{T}}_i - \tilde{\mathcal{T}}_i) \times_1 (\hat{\mathbf{U}}_1^T \mathbf{O}_1)^{-1} \times_2 (\hat{\mathbf{U}}_2^T \mathbf{O}_2) \times_3 (\hat{\mathbf{U}}_3^T \mathbf{O}_3)^{-1} \right\|_{1,1} S_H \leq S_H \gamma_i \quad (59)$$

Second term is bounded by (since $\|Q\|_{1,1} \leq \|Q\|_1$):

$$\sum_x \|(\hat{\mathbf{U}}_1^T \mathbf{O}_1)(\hat{\mathbf{m}}_{J:1} - \tilde{\mathbf{m}}_{J:1})\|_1 \|(\hat{\mathcal{T}}_i - \tilde{\mathcal{T}}_i) \times_1 (\hat{\mathbf{U}}_1^T \mathbf{O}_1)^{-1} \times_2 (\hat{\mathbf{U}}_2^T \mathbf{O}_2) \times_3 (\hat{\mathbf{U}}_3^T \mathbf{O}_3)^{-1}\|_{1,1} \quad (60)$$

$$\leq \gamma_i \sum_{\mathbf{x}_i} \|(\hat{\mathbf{U}}_1^T \mathbf{O}_1)(\hat{\mathbf{m}}_{J:1} - \tilde{\mathbf{m}}_{J:1})\|_1 \quad (61)$$

Third Term is bounded by:

$$\|\tilde{\mathcal{T}}_i \times_1 (\hat{\mathbf{U}}_1^T \mathbf{O}_1)^{-1} \times_2 (\hat{\mathbf{U}}_2^T \mathbf{O}_2) \times_3 (\hat{\mathbf{U}}_3^T \mathbf{O}_3)^{-1}\|_{1,1} \sum_{\mathbf{x}_i} \|(\hat{\mathbf{U}}_1^T \mathbf{O}_1)(\hat{\mathbf{m}}_{J:1} - \tilde{\mathbf{m}}_{J:1})\|_1 \leq \sum_{\mathbf{x}_i} \|(\hat{\mathbf{U}}_1^T \mathbf{O}_1)(\hat{\mathbf{m}}_{J:1} - \tilde{\mathbf{m}}_{J:1})\|_1 \quad (62)$$

In the next section, we will see that $\sum_{\mathbf{x}_i} \|(\hat{\mathbf{U}}^T \mathbf{O})(\hat{\mathbf{m}}_{J:1} - \tilde{\mathbf{m}}_{J:1})\|_1 \leq \left(\prod_{k=1}^J (1 + \Delta_{j_k}) \delta_{one}^i + S_H \prod_{k=1}^J (1 + \Delta_{j_k}) - S_H \right)$. So the overall bound is

$$\Delta_i \leq \left((1 + \gamma_i) \prod_{k=1}^J (1 + \Delta_{j_k}) \delta_{one}^i + (1 + \gamma_i) S_H \prod_{k=1}^J (1 + \Delta_{j_k}) - S_H \right). \quad (63)$$

(where j_1, \dots, j_J are children of node i).

6.5 Bounding $\sum_{\mathbf{x}_i} \|(\hat{\mathbf{U}}^T \mathbf{O})(\widehat{\mathbf{m}}_{J:1} - \tilde{\mathbf{m}}_{J:1})\|_1$

Lemma 4

$$\sum_{\mathbf{x}_i} \|(\hat{\mathbf{U}}^T \mathbf{O})(\widehat{\mathbf{m}}_{J:1} - \tilde{\mathbf{m}}_{J:1})\|_1 \leq \prod_{k=1}^J (1 + \Delta_{j_k}) \delta_{one}^i + S_H \prod_{k=1}^J (1 + \Delta_{j_k}) - S_H \quad (64)$$

(where j_1, \dots, j_J are children of node i).

The proof is by induction. Base case: $\|(\hat{\mathbf{U}}^T \mathbf{O})(\widehat{\mathbf{1}}_i - \tilde{\mathbf{1}}_i)\|_1 \leq \delta_{one}^i$, by definition of δ_{one}^i .

Inductive step: Let us say claim holds up until $u-1$. We show it holds for u . Thus

$$\sum_{\mathbf{x}_i} \|(\hat{\mathbf{U}}^T \mathbf{O})(\widehat{\mathbf{m}}_{(u-1):1} - \tilde{\mathbf{m}}_{(u-1):1})\|_1 \leq \prod_{k=1}^{u-1} (1 + \Delta_{j_k}) \delta_{one}^i + S_H \prod_{k=1}^{u-1} (1 + \Delta_{j_k}) - S_H \quad (65)$$

We now decompose the sum over x as

$$\begin{aligned} & \sum_{\mathbf{x}_{u:1}} \|(\hat{\mathbf{U}}^T \mathbf{O})(\widehat{\mathbf{m}}_{u:1} - \tilde{\mathbf{m}}_{u:1})\|_1 \\ &= \sum_{\mathbf{x}_{u:1}} \|(\hat{\mathbf{U}}^T \mathbf{O}) \left((\widehat{\mathbf{M}}_u - \tilde{\mathbf{M}}_u) \tilde{\mathbf{m}}_{(u-1):1} + (\widehat{\mathbf{M}}_u - \tilde{\mathbf{M}}_u)(\widehat{\mathbf{m}}_{(u-1):1} - \tilde{\mathbf{m}}_{(u-1):1}) + (\widehat{\mathbf{m}}_{(u-1):1} - \tilde{\mathbf{m}}_{(u-1):1}) \right)\|_1 \end{aligned} \quad (66)$$

Using the triangle inequality, we get

$$\sum_{\mathbf{x}_{u:1}} \|(\hat{\mathbf{U}}_2^T \mathbf{O}_2)(\widehat{\mathbf{M}}_u - \tilde{\mathbf{M}}_u)(\hat{\mathbf{U}}_1^T \mathbf{O}_1)^{-1}\|_1 \|(\hat{\mathbf{U}}_1^T \mathbf{O}_1) \tilde{\mathbf{m}}_{(u-1):1}\|_1 \quad (67)$$

$$+ \sum_{\mathbf{x}_{u:1}} \|(\hat{\mathbf{U}}_2^T \mathbf{O}_2)(\widehat{\mathbf{M}}_u - \tilde{\mathbf{M}}_u)(\hat{\mathbf{U}}_1^T \mathbf{O}_1)^{-1}\|_1 \|(\hat{\mathbf{U}}_1^T \mathbf{O}_1)(\widehat{\mathbf{m}}_{(u-1):1} - \tilde{\mathbf{m}}_{(u-1):1})\|_1 \quad (68)$$

$$+ \sum_{\mathbf{x}_{u:1}} \|(\hat{\mathbf{U}}^T \mathbf{O}) \tilde{\mathbf{M}}_u (\hat{\mathbf{U}}^T \mathbf{O})^{-1}\|_1 \|(\hat{\mathbf{U}}^T \mathbf{O})(\widehat{\mathbf{m}}_{(u-1):1} - \tilde{\mathbf{m}}_{(u-1):1})\|_1 \quad (69)$$

Again we are just numbering the \mathbf{U} 's and \mathbf{O} 's for clarity to see which corresponds with which. They are omitted in the actual theorem statements since we will take minimums etc. at the end.

We now must bound these terms. First term:

$$\sum_{\mathbf{x}_u} \|(\hat{\mathbf{U}}_2^T \mathbf{O}_2)(\widehat{\mathbf{M}}_u - \tilde{\mathbf{M}}_u)(\hat{\mathbf{U}}_2^T \mathbf{O}_2)^{-1}\|_1 \sum_{\mathbf{x}_{1:u-1}} \|(\hat{\mathbf{U}}_1^T \mathbf{O}_1) \tilde{\mathbf{m}}_{(u-1):1}\|_1 \leq \Delta_u \sum_{\mathbf{x}_{(u-1):1}} \|\tilde{\mathbf{m}}_{(u-1):1}(\hat{\mathbf{U}}^T \mathbf{O})\|_1 \leq S_H \Delta_u \quad (70)$$

since $\Delta_u = \|(\hat{\mathbf{U}}_2^T \mathbf{O}_2)(\widehat{\mathbf{M}}_u - \tilde{\mathbf{M}}_u)(\hat{\mathbf{U}}_1^T \mathbf{O}_1)^{-1}\|_1$. Second term can be bounded by inductive hypothesis:

$$\sum_{\mathbf{x}_{u:1}} \|(\hat{\mathbf{U}}_2^T \mathbf{O}_2)(\widehat{\mathbf{M}}_u - \tilde{\mathbf{M}}_u)(\hat{\mathbf{U}}_1^T \mathbf{O}_1)^{-1}\|_1 \|(\hat{\mathbf{U}}_1^T \mathbf{O}_1)(\widehat{\mathbf{m}}_{(u-1):1} - \tilde{\mathbf{m}}_{(u-1):1})\|_1 \leq \Delta_i \left(\prod_{k=1}^{u-1} (1 + \Delta_{j_k}) \delta_{one}^i + S_H \prod_{k=1}^{u-1} (1 + \Delta_{j_k}) - S_H \right) \quad (71)$$

The third term is bounded by observing that $(\hat{\mathbf{U}}^T \mathbf{O}) \tilde{\mathbf{M}}_u (\hat{\mathbf{U}}^T \mathbf{O})^{-1} = \text{diag}(Pr[\mathbf{x}_u | \text{Parent}])$. Thus it is diagonal, and $\sum_{\mathbf{x}_{u:1}} Pr[\mathbf{x} | \text{Parent}]$ has max row or column sum as 1. This means that the third term is bounded by the inductive hypothesis as well:

$$\sum_{\mathbf{x}_{u:1}} \|(\hat{\mathbf{U}}^T \mathbf{O}) \tilde{\mathbf{M}}_u (\hat{\mathbf{U}}^T \mathbf{O})^{-1}\|_1 \|(\hat{\mathbf{U}}^T \mathbf{O})(\widehat{\mathbf{m}}_{(u-1):1} - \tilde{\mathbf{m}}_{(u-1):1})\|_1 \leq \left(\prod_{k=1}^{u-1} (1 + \Delta_{j_k}) \delta_{one}^i + S_H \prod_{k=1}^{u-1} (1 + \Delta_{j_k}) - S_H \right) \quad (72)$$

7 Bounding the propagation of error in tree

We now wrap up the proof based on the approach of HKZ[3].

Lemma 5

$$\sum_{x_1, \dots, x_O} \left| \widehat{\mathbb{P}}[x_1, \dots, x_O] - \mathbb{P}[x_1, \dots, x_O] \right| \leq S_H \delta_{root} + (1 + \delta_{root}) \left(\prod_{k=1}^J (1 + \Delta_{j_k}) \delta_{one}^r + S_H \prod_{k=1}^J (1 + \Delta_{j_k}) - S_H \right) \quad (73)$$

$$\sum_{x_1, \dots, x_O} \left| \widehat{\mathbb{P}}[x_1, \dots, x_O] - \mathbb{P}[x_1, \dots, x_O] \right| = \sum_{x_1, \dots, x_O} \left| \widehat{\mathbf{r}}^T \widehat{\mathbf{M}}_{j_1} \dots \widehat{\mathbf{M}}_{j_J} \widehat{\mathbf{1}}_r - \tilde{\mathbf{r}}^T \tilde{\mathbf{M}}_{j_1} \dots \tilde{\mathbf{M}}_{j_J} \tilde{\mathbf{1}}_r \right| \quad (74)$$

$$\leq \sum_{x_1, \dots, x_O} \left| (\widehat{\mathbf{r}} - \tilde{\mathbf{r}})^T (\widehat{\mathbf{U}}^T \mathbf{O})^{-1} (\widehat{\mathbf{U}}^T \mathbf{O}) (\tilde{\mathbf{M}}_{J:1} \tilde{\mathbf{1}}) \right| \quad (75)$$

$$+ \sum_{x_1, \dots, x_O} \left| (\widehat{\mathbf{r}} - \tilde{\mathbf{r}})^T (\widehat{\mathbf{U}}^T \mathbf{O})^{-1} (\widehat{\mathbf{U}}^T \mathbf{O}) (\widehat{\mathbf{M}}_{J:1} \widehat{\mathbf{1}}_r - \tilde{\mathbf{M}}_{J:1} \tilde{\mathbf{1}}_r) \right| \quad (76)$$

$$+ \sum_{x_1, \dots, x_O} \left| \tilde{\mathbf{r}}^T (\widehat{\mathbf{U}}^T \mathbf{O})^{-1} (\widehat{\mathbf{U}}^T \mathbf{O}) (\widehat{\mathbf{M}}_{J:1} \widehat{\mathbf{1}}_i - \tilde{\mathbf{M}}_{J:1} \tilde{\mathbf{1}}) \right| \quad (77)$$

The first sum is bounded using Holder inequality and noting that the first term is a conditional probability (of all observed variables conditioned on the root)

$$\sum_{x_1, \dots, x_O} \left| (\widehat{\mathbf{r}} - \tilde{\mathbf{r}})^T (\widehat{\mathbf{U}}^T \mathbf{O})^{-1} (\widehat{\mathbf{U}}^T \mathbf{O}) (\tilde{\mathbf{M}}_{J:1} \tilde{\mathbf{1}}) \right| \quad (78)$$

$$\leq \sum_{x_1, \dots, x_O} \|(\widehat{\mathbf{r}} - \tilde{\mathbf{r}})^T (\widehat{\mathbf{U}}^T \mathbf{O})^{-1}\|_\infty \|(\widehat{\mathbf{U}}^T \mathbf{O}) (\tilde{\mathbf{M}}_{J:1} \tilde{\mathbf{1}})\|_1 \leq S_H \delta_{root} \quad (79)$$

The second sum is bounded by another application of Holder's inequality (and Lemma XX):

$$\sum_{x_1, \dots, x_O} \left| (\widehat{\mathbf{r}} - \tilde{\mathbf{r}})^T (\widehat{\mathbf{U}}^T \mathbf{O})^{-1} (\widehat{\mathbf{U}}^T \mathbf{O}) (\widehat{\mathbf{M}}_{J:1} \widehat{\mathbf{1}}_r - \tilde{\mathbf{M}}_{J:1} \tilde{\mathbf{1}}_r) \right| \quad (80)$$

$$\leq \sum_{x_1, \dots, x_O} \|(\widehat{\mathbf{r}} - \tilde{\mathbf{r}})^T (\widehat{\mathbf{U}}^T \mathbf{O})^{-1}\|_\infty \|(\widehat{\mathbf{U}}^T \mathbf{O}) (\widehat{\mathbf{M}}_{J:1} \widehat{\mathbf{1}}_r - \tilde{\mathbf{M}}_{J:1} \tilde{\mathbf{1}}_r)\|_1 \quad (81)$$

$$\leq \delta_{root} \left(\prod_{k=1}^J (1 + \Delta_{j_k}) \delta_{one}^r + S_H \prod_{k=1}^J (1 + \Delta_{j_k}) - S_H \right) \quad (82)$$

The third sum is also bounded by Holder's Inequality and previous lemmas and noting that $\tilde{\mathbf{r}}^T (\mathbf{U}^T \mathbf{O})^{-1} = \mathbb{P}[R = r]$:

$$\sum_{x_1, \dots, x_O} \left| \tilde{\mathbf{r}}^T (\widehat{\mathbf{U}}^T \mathbf{O})^{-1} (\widehat{\mathbf{U}}^T \mathbf{O}) (\widehat{\mathbf{M}}_{J:1} \widehat{\mathbf{1}}_i - \tilde{\mathbf{M}}_{J:1} \tilde{\mathbf{1}}) \right| \quad (83)$$

$$\leq \sum_{x_1, \dots, x_O} \|\tilde{\mathbf{r}}^T (\widehat{\mathbf{U}}^T \mathbf{O})^{-1}\|_\infty \|(\widehat{\mathbf{U}}^T \mathbf{O}) (\widehat{\mathbf{M}}_{J:1} \widehat{\mathbf{1}}_i - \tilde{\mathbf{M}}_{J:1} \tilde{\mathbf{1}})\|_1 \quad (84)$$

$$\leq \left(\prod_{k=1}^J (1 + \Delta_{j_k}) \delta_{one}^r + S_H \prod_{k=1}^J (1 + \Delta_{j_k}) - S_H \right) \quad (85)$$

Combining these bounds gives us the desired solution.

8 Putting it all together

We seek for

$$\sum_{x_1, \dots, x_O} \left| \widehat{\mathbb{P}}[x_1, \dots, x_O] - \mathbb{P}[x_1, \dots, x_O] \right| \leq \epsilon \quad (86)$$

Using the fact that for $a < .5$, $(1 + a/t)^t \leq 1 + 2a$, we get that $\Delta_{j_k} \leq O(\epsilon/mJ)$. However, Δ_j is defined recursively, and thus the error accumulates exponential in the longest path of hidden nodes. For example, $\Delta_i^{obs} \leq O(\frac{\epsilon}{d_{max} m^\ell})$ where ℓ is the longest path of hidden nodes. Tracing this back through will gives the result:

Pick any $\epsilon > 0, \delta < 1$. Let

$$N \geq O \left(\frac{1}{\epsilon^2} \left(\frac{(d_{max} m)^{\ell+1} S_O}{\min_i \sigma_{S_H}(\mathbf{O}_i)^2 \min_{i \neq j} \sigma_{S_H}(P_{i,j})^4} \right) \right) \log \frac{\mathcal{O}}{\delta} \quad (87)$$

Then with probability $1 - \delta$

$$\sum_{x_1, \dots, x_O} \left| \widehat{\mathbb{P}}[x_1, \dots, x_O] - \mathbb{P}[x_1, \dots, x_O] \right| \leq \epsilon \quad (88)$$

In many cases, if the frequency of the observation symbols follow certain distributions, than the dependence on S_O can be removed as showed in HKZ [3].

9 Appendix

9.1 Matrix Perturbation Bounds

This is Theorem 3.8 from pg. 143 in Stewart and Sun, 1990 [5]. Let $A \in \mathbf{R}^{m \times n}$, with $m \geq n$ and let $\tilde{A} = A + E$. Then

$$\|\tilde{A}^+ - A^+\|_2 \leq \frac{1 + \sqrt{5}}{2} \max(\|A^+\|_2^2, \|\tilde{A}\|_2^2) \|E\|_2 \quad (89)$$

References

- [1] Myung J. Choi, Vicent Y. Tan, Animashree Anandkumar, and Alan S. Willsky. Learning latent tree graphical models. In *arXiv:1009.2722v1*, 2010.
- [2] R.A. Horn and C.R. Johnson. *Matrix analysis*. Cambridge Univ Pr, 1990.
- [3] D. Hsu, S. Kakade, and T. Zhang. A spectral algorithm for learning hidden markov models. In *Proc. Annual Conf. Computational Learning Theory*, 2009.
- [4] N.H. Nguyen, P. Drineas, and T.D. Tran. Tensor sparsification via a bound on the spectral norm of random tensors. *Arxiv preprint arXiv:1005.4732*, 2010.
- [5] G.W. Stewart and J. Sun. *Matrix perturbation theory*, volume 175. Academic press New York, 1990.