

A Spectral Algorithm For Latent Junction Trees - Supplementary Material

Ankur P. Parikh, Le Song, Mariya Ishteva, Gabi Teodoru, Eric P. Xing

1 Discussion of Conditions for Observable Representation

The observable representation exists only if there exist transformations $\mathcal{F}_i = \mathbb{P}[\mathcal{O}_i|S_i]$ with rank $\tau_i := k_h \times |S_i|$ and $\mathbb{P}[\mathcal{O}_{i-}|S_i]$ also has rank τ_i (so that $\mathbb{P}(\mathcal{O}_i, \mathcal{O}_{i-})$ has rank τ_i). Thus, it is required that $\#\text{states}(\mathcal{O}_i) \geq \#\text{states}(S_i)$. This can either be achieved by either making \mathcal{O}_i consist of a few high dimensional observations, or many smaller dimensional ones. In the case when $\#\text{states}(\mathcal{O}_i) > \#\text{states}(S_i)$, we need to project \mathcal{F}_i to a lower dimensional space such that it can be inverted using a tensor \mathcal{U}_i . In this case, we define $\mathcal{F}_i := \mathbb{P}[\mathcal{O}_i|S_i] \times_{\mathcal{O}_i} \mathcal{U}_i$. Following this through the computation gives us that $\tilde{\mathcal{P}}(C_i) = \mathbb{P}(\mathcal{O}_i, \mathcal{O}_{i-}) \times_{\mathcal{O}_{i-}} (\mathbb{P}(\mathcal{O}_i, \mathcal{O}_{i-}) \times_{\mathcal{O}_i} \mathcal{U}_i)^{-1}$. A good choice of \mathcal{U}_i can be obtained by performing a singular value decomposition of the matricized version of $\mathbb{P}(\mathcal{O}_i, \mathcal{O}_{i-})$ (variables in \mathcal{O}_i are arranged to rows and those in \mathcal{O}_{i-} to columns).

For HMMs and latent trees, this rank condition can be expressed simply as requiring the conditional probability tables of the underlying model to not be rank-deficient. However, junction trees encode significantly more complex latent structures that introduce more subtle considerations. While we consider a general characterization of such models where the observable representation to be future work, here try to give some intuition on what types of latent structures the rank condition may fail.

First, the rank condition can fail is if there are not enough observed nodes/states, and thus $\#\text{states}(\mathcal{O}_i) < \tau_i$. Intuitively, this corresponds to a model where the latent space is too expressive and inherently represents an intractable model (e.g. a set of n binary variables connected by a hidden node with 2^n states is equivalent to a clique of size n).

However, there are more subtle considerations unique to non-tree models. In general, our method is not limited to non-triangulated graphs (see the factorial HMM in Figure 3 of the main paper), but the process of triangulation can introduce artificial dependencies that can lead to complications. Consider Figure 1(a) which shows a DAG and its corresponding junction tree. To construct $\tilde{\mathcal{P}}(C_{AD})$, we may set $\mathbb{P}[\mathcal{O}_i, \mathcal{O}_{-i}] = \mathbb{P}[D, F]$ based on the junction tree topology. However, in the original model before triangulation, $D \perp F$ because of the v-structure. As a result, $\mathbb{P}[D, F]$ does not have rank k_h and thus cannot be inverted. However, note that choosing $\mathbb{P}[\mathcal{O}_i, \mathcal{O}_{-i}] = \mathbb{P}[D, E]$ is valid.

Finally consider Figure 1 (b), ignoring the orange nodes for now and assuming the variables are binary. In this model, the hidden variables are largely redundant since integrating out A and B would simply give a chain. \mathcal{F}_{r_2} must be set to $\mathbb{P}[F|S_{r_2}] = \mathbb{P}[F|AB]$. If we think of A, B has just one large variable, then it is clear that $\mathbb{P}[F|AB] = \mathbb{P}[F|D]\mathbb{P}[D|AB]$. However, D only has two states while AB has 4, so $\mathbb{P}[F|AB]$ only has rank 2. Now, consider adding the orange node. In this case we could set \mathcal{F}_{r_2} to $\mathbb{P}[F, G|S_{r_2}] = \mathbb{P}[F, G|A, B]$ whose matricized version has rank 4. Note that once the orange node has been added, integrating out A and B no longer produces a simple chain, but a more complicated structure.

Thus, we believe that more rigorously characterizing the existence of the observable representation in more detail, may shed light on the ‘‘intrinsic’’ complexity/redundancy of latent variable models in the context of linear and tensor algebra.

2 Linear Systems to Improve Stability

As derived in Section 6 in the main paper, for non-root and non-leaf nodes:

$$\tilde{\mathcal{P}}(C_i) \times_{\mathcal{O}_i} \mathbb{P}(\mathcal{O}_i, \mathcal{O}_{i-}) = \mathbb{P}(\mathcal{O}_{i_1}, \mathcal{O}_{i_2}, \mathcal{O}_{i-}) \quad (1)$$

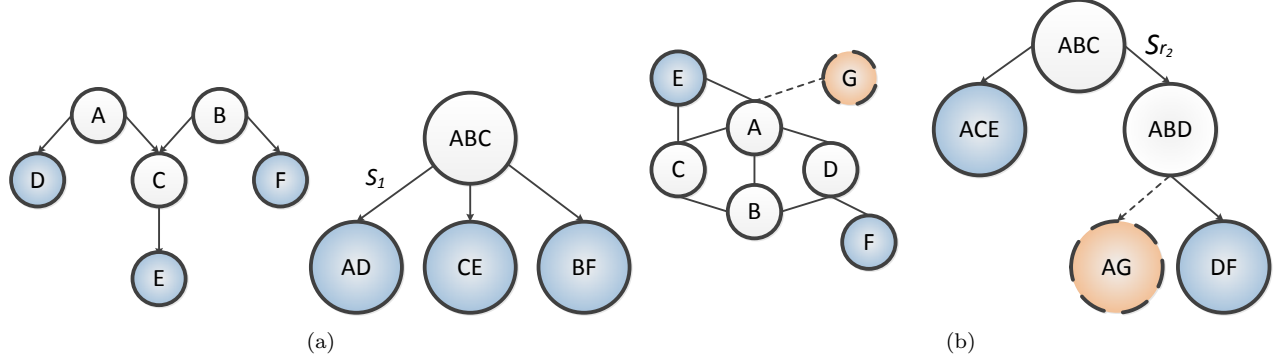


Figure 1: Models and their corresponding junction trees, where constructing the observable representation poses complications. See Section 1.

which then implies that

$$\tilde{\mathcal{P}}(C_i) = \mathbb{P}(\theta_{i_1}, \theta_{i_2}, \theta_{i_-}) \times_{\theta_{i_-}} \mathbb{P}(\theta_i, \theta_{i_-})^{-1} \quad (2)$$

However, there may be many choices of θ_{i_-} (which we denote with $\theta_{i_-}^{(1)}, \dots, \theta_{i_-}^{(N)}$) for which the above equality is true:

$$\begin{aligned} \tilde{\mathcal{P}}(C_i) \times_{\theta_i} \mathbb{P}(\theta_i, \theta_{i_-}^{(1)}) &= \mathbb{P}(\theta_{i_1}, \theta_{i_2}, \theta_{i_-}^{(1)}) \\ \tilde{\mathcal{P}}(C_i) \times_{\theta_i} \mathbb{P}(\theta_i, \theta_{i_-}^{(2)}) &= \mathbb{P}(\theta_{i_1}, \theta_{i_2}, \theta_{i_-}^{(2)}) \\ &\dots \\ \tilde{\mathcal{P}}(C_i) \times_{\theta_i} \mathbb{P}(\theta_i, \theta_{i_-}^{(N)}) &= \mathbb{P}(\theta_{i_1}, \theta_{i_2}, \theta_{i_-}^{(N)}) \end{aligned}$$

This defines an over-constrained linear system, and we can solve for $\tilde{\mathcal{P}}(C_i)$ using least squares. In the case where $\#\text{states}(\mathcal{O}_i) > \#\text{states}(S_i)$, and the projection tensor \mathbf{U}_i is needed, our system of equations becomes:

$$\begin{aligned} \tilde{\mathcal{P}}(C_i) \times_{\theta_i} (\mathbb{P}(\theta_i, \theta_{i_-}^{(1)}) \times_{\theta_i} \mathbf{U}_i) &= \mathbb{P}(\theta_{i_1}, \theta_{i_2}, \theta_{i_-}^{(1)}) \times_{\theta_{i_1}} \mathbf{U}_{i_1} \times_{\theta_{i_2}} \mathbf{U}_{i_2} \\ \tilde{\mathcal{P}}(C_i) \times_{\theta_i} (\mathbb{P}(\theta_i, \theta_{i_-}^{(2)}) \times_{\theta_i} \mathbf{U}_i) &= \mathbb{P}(\theta_{i_1}, \theta_{i_2}, \theta_{i_-}^{(2)}) \times_{\theta_{i_1}} \mathbf{U}_{i_1} \times_{\theta_{i_2}} \mathbf{U}_{i_2} \\ &\dots \\ \tilde{\mathcal{P}}(C_i) \times_{\theta_i} (\mathbb{P}(\theta_i, \theta_{i_-}^{(N)}) \times_{\theta_i} \mathbf{U}_i) &= \mathbb{P}(\theta_{i_1}, \theta_{i_2}, \theta_{i_-}^{(N)}) \times_{\theta_{i_1}} \mathbf{U}_{i_1} \times_{\theta_{i_2}} \mathbf{U}_{i_2} \end{aligned}$$

In this case, one good choice of \mathbf{U}_i is the top singular vectors of the matrix formed by the horizontal concatenation of $\mathbb{P}(\theta_i, \theta_{i_-}^{(1)}), \dots, \mathbb{P}(\theta_i, \theta_{i_-}^{(N)})$.

This linear system method allows for more robust estimation especially in smaller sample sizes (at the cost of more computation). It can be applied to the leaf case as well. One does not need to set up the linear system with all the valid choices; a subset is also acceptable.

3 Sample Complexity Theorem

We prove the sample complexity theorem.

3.1 Notation

For simplicity, the main text describes the algorithm in the context of a binary tree. However, in the sample complexity proof, we adopt a more general notation. Let C_i be a clique, and $C_{\alpha_i(1)}, \dots, C_{\alpha_i(\gamma_i)}$ denote its γ_i children. Let \mathbb{C} denote the set of all the cliques in the junction tree, and $|\mathbb{C}|$ denote its size. Define d_{\max} to be maximum degree of a tensor in the observable representation and e_{\max} to be maximum number of

observed variables in any tensor. Furthermore, let $\tau_i = |k_h| \times |S_i|$ (i.e. the number of states associated with the separator).

We can now write the transformed representation for junction trees with more than 3 neighbors as:

Root:

$$\tilde{\mathcal{P}}(C_i) = \mathcal{P}(C_i) \times_{S_{\alpha_i(1)}} \mathcal{F}_{\alpha_i(1)} \times \dots \times_{S_{\alpha_i(\gamma_i)}} \mathcal{F}_{\alpha_i(\gamma_i)}$$

Internal nodes:

$$\tilde{\mathcal{P}}(C_i) = \mathcal{P}(C_i) \times_{S_i} \mathcal{F}_i^\dagger \times_{S_{\alpha_i(1)}} \mathcal{F}_{\alpha_i(1)} \times \dots \times_{S_{\alpha_i(\gamma_i)}} \mathcal{F}_{\alpha_i(\gamma_i)}$$

Leaf:

$$\tilde{\mathcal{P}}(C_i) = \mathcal{P}(C_i) \times_{S_i} \mathcal{F}_i^\dagger$$

and the observable representation as:

$$\text{root: } \mathcal{P}(C_i) = \mathbb{P}(\theta_{\alpha_i(1)}, \dots, \theta_{\alpha_i(\gamma_i)}) \times_{\theta_{\alpha_i(1)}} \mathcal{U}_{\alpha_i(1)} \times \dots \times_{\theta_{\alpha_i(\gamma_i)}} \mathcal{U}_{\alpha_i(\gamma_i)}$$

$$\text{internal: } \mathcal{P}(C_i) = \mathbb{P}(\theta_{\alpha_i(1)}, \dots, \theta_{\alpha_i(\gamma_i)}, \theta_{-i}) \times_{\theta_{-i}} (\mathbb{P}(\theta_i, \theta_{-i}) \times_{\theta_i} \mathcal{U}_i)^\dagger \times_{\theta_{\alpha_i(1)}} \mathcal{U}_{\alpha_i(1)} \times \dots \times_{\theta_{\alpha_i(\gamma_i)}} \mathcal{U}_{\alpha_i(\gamma_i)}$$

$$\text{leaf: } \mathcal{P}(C_i) = \hat{\mathbb{P}}(R_i, \theta_{-i}) \times_{\theta_{-i}} (\hat{\mathbb{P}}(\theta_i, \theta_{-i}) \times_{\theta_i} \mathcal{U}_i)^\dagger$$

Sometimes when the multiplication indices are clear, we will omit it to make things simpler.

Rearranged version:

We will often find it convenient to rearrange the tensors into lower order tensors for the purpose of taking some norms (such as the spectral norm). We define $\mathbf{R}(\cdot)$ as the ‘‘rearranging’’ operation. For example, $\mathbf{R}(\hat{\mathbb{P}}(\theta_i, \theta_{-i}))$ is the matricized version of $\hat{\mathbb{P}}(\theta_i, \theta_{-i})$ with θ_i being mapped to the rows and θ_{-i} being mapped to the columns.

More generally, if $\mathcal{P}(C_i)$ (which has order d_i) then $\mathbf{R}(\mathcal{P}(C_i))$ is a rearrangement of $\mathcal{P}(C_i)$ such that it has order equal to the number of neighbors of C_i . The set of modes of $\mathcal{P}(C_i)$ corresponding to a single separator of C_i get mapped to a single mode in $\mathbf{R}(\mathcal{P}(C_i))$. We let $\mathbf{R}(S_{\alpha_i(j)})$ denote the mode that $S_{\alpha_i(j)}$ maps to in $\mathbf{R}(\mathcal{P}(C_i))$.

In the case of the root, $\mathbf{R}(\mathcal{P}(C_i))$ rearranges $\mathcal{P}(C_i)$ into a tensor of order γ_i . For other clique nodes, $\mathbf{R}(\mathcal{P}(C_i))$ rearranges $\mathcal{P}(C_i)$ into a tensor of order $\gamma_i + 1$.

Example: In Figure 2 of the main text $\mathcal{P}(C_{BCDE}) = \mathbb{P}(\oslash_2 B, D | \oslash_2 C, E)$. C_{BCDE} has 3 neighbors and so $\mathbf{R}(\mathcal{P}(C_{BCDE}))$ is of order 3 where each of the modes correspond to $\{B, C\}$, $\{B, D\}$, $\{C, E\}$ respectively.

This rearrangement can be applied to the other quantities. For example, $\mathbf{R}(\mathbb{P}(\theta_{\alpha_i(1)}, \dots, \theta_{\alpha_i(\gamma_i)}))$ is a rearranging of $\mathbb{P}(\theta_{\alpha_i(1)}, \dots, \theta_{\alpha_i(\gamma_i)})$ into a tensor of order γ_i where the $\theta_{\alpha_i(1)}, \dots, \theta_{\alpha_i(\gamma_i)}$ correspond to one mode each. $\mathbf{R}(\hat{\mathbb{P}}(\theta_i, \theta_{-i}))$ is the matricized version of $\hat{\mathbb{P}}(\theta_i, \theta_{-i})$ with θ_i being mapped to the rows and θ_{-i} being mapped to the columns. Similarly, $\mathbf{R}(\mathcal{F}_i)$ is the matricized version of \mathcal{F}_i and $\mathbf{R}(\mathcal{U}_i)$ is the matricized version of \mathcal{U}_i .

Thus, we can define the rearranged quantities:

Root:

$$\mathbf{R}(\tilde{\mathcal{P}}(C_i)) = \mathbf{R}(\mathcal{P}(C_i)) \times_{\mathbf{R}(S_{\alpha_i(1)})} \mathbf{R}(\mathcal{F}_{\alpha_i(1)}) \times \dots \times_{\mathbf{R}(S_{\alpha_i(\gamma_i)})} \mathbf{R}(\mathcal{F}_{\alpha_i(\gamma_i)})$$

$$\mathbf{R}(\tilde{\mathcal{P}}(C_i)) = \mathbf{R}(\mathcal{P}(C_i)) \times_{\mathbf{R}(S_i)} \mathbf{R}(\mathcal{F}_i^\dagger) \times_{\mathbf{R}(S_{\alpha_i(1)})} \mathcal{F}_{\alpha_i(1)} \times \dots \times_{\mathbf{R}(S_{\alpha_i(\gamma_i)})} \mathbf{R}(\mathcal{F}_{\alpha_i(\gamma_i)})$$

Leaf:

$$\mathbf{R}(\tilde{\mathcal{P}}(C_i)) = \mathbf{R}(\mathcal{P}(C_i)) \times_{\mathbf{R}(S_i)} \mathbf{R}(\mathcal{F}_i^\dagger)$$

and the observable representation as:

$$\begin{aligned}
\text{root: } \mathbf{R}(\mathcal{P}(C_i)) &= \mathbf{R}(\mathbb{P}(\mathcal{O}_{\alpha_i(1)}, \dots, \mathcal{O}_{\alpha_i(\gamma_i)})) \times_{\mathbf{R}(S_{\alpha_i(1)})} \mathbf{R}(\mathbf{U}_{\alpha_i(1)}) \times \dots \times_{\mathbf{R}(S_{\alpha_i(\gamma_i)})} \mathbf{R}(\mathbf{U}_{\alpha_i(\gamma_i)}) \\
\text{internal: } \mathbf{R}(\mathcal{P}(C_i)) &= \mathbf{R}(\mathbb{P}(\mathcal{O}_{S_{\alpha_i(1)}}, \dots, \mathcal{O}_{\alpha_i(\gamma_i)}, \mathcal{O}_{-i})) \times_{\mathbf{R}(S_{-i})} \mathbf{R}((\mathbb{P}(\mathcal{O}_i, \mathcal{O}_{-i}) \times_{S_i} \mathbf{U}_i)^\dagger) \\
&\quad \times_{\mathbf{R}(S_{\alpha_i(1)})} \mathbf{R}(\mathbf{U}_{\alpha_i(1)}) \times \dots \times_{\mathbf{R}(S_{\alpha_i(\gamma_i)})} \mathbf{R}(\mathbf{U}_{\alpha_i(\gamma_i)}) \\
\text{leaf: } \mathbf{R}(\mathcal{P}(C_i)) &= \mathbf{R}(\widehat{\mathbb{P}}(R_i, \mathcal{O}_{-i})) \times_{\mathbf{R}(\mathcal{O}_{-i})} \mathbf{R}((\widehat{\mathbb{P}}(\mathcal{O}_i, \mathcal{O}_{-i}) \times_{\mathcal{O}_i} \mathbf{U}_i)^\dagger)
\end{aligned} \tag{3}$$

Furthermore define the following:

$$\alpha = \min_{C_i \in \mathbb{C}} \sigma_{\tau_i}(\mathbb{P}[\mathcal{O}_i, \mathcal{O}_{-i}]) \tag{4}$$

$$\beta = \min_{C_i \in \mathbb{C}} \sigma_{\tau_i}(\mathcal{F}_i) \tag{5}$$

The proof is based on the technique of HKZ [2] but has differences due to the junction tree topology and higher order tensors.

3.2 Tensor Norms

We briefly define several tensor norms that are a generalization of matrix norms. For more information about matrix norms see [1].

Frobenius Norm: Just like for matrices, the frobenius norm of a tensor of order N is defined as:

$$\|\mathbf{T}\|_F = \sqrt{\sum_{i_1, \dots, i_N} \mathbf{T}(i_1, \dots, i_N)^2} \tag{6}$$

Elementwise One Norm: Similarly, the elementwise one norm of a tensor of order N is defined as:

$$\|\mathbf{T}\|_1 = \sum_{i_1, \dots, i_N} |\mathbf{T}(i_1, \dots, i_N)| \tag{7}$$

Spectral Norm: For tensors of order N the spectral norm [3] can be defined as

$$\|\|\mathbf{T}\|\|_2 = \sup_{\mathbf{v}_i \text{ s.t. } \|\mathbf{v}_i\|_2 \leq 1 \forall 1 \leq i \leq N} \mathbf{T} \times_N \mathbf{v}_N, \dots, \times_2 \mathbf{v}_2 \times_1 \mathbf{v}_1 \tag{8}$$

In our case, we will find it more convenient to use the rearranged spectral norm, which we define as:

$$\|\|\mathbf{T}\|\|_{2R} = \|\|\mathbf{R}(\mathbf{T})\|\|_2 \tag{9}$$

where $\mathbf{R}(\cdot)$ was defined in the previous section.

Induced One Norm: For matrices the induced one norm is defined as the max column sum of the matrix: $\|\|\mathbf{M}\|\|_1 = \sup_{\mathbf{v} \text{ s.t. } \|\mathbf{v}\|_1 \leq 1} \|\mathbf{M}\mathbf{v}\|_1$. We can generalize this to be the max slice sum of a tensor on a tensor where some modes are fixed and others are summed over. Let σ denote the modes that will be maxed over. Then:

$$\|\|\mathbf{T}\|\|_1^\sigma = \sup_{\mathbf{v}_i \text{ s.t. } \|\mathbf{v}_i\|_1 \leq 1, \forall 1 \leq i \leq |\sigma|} \|\mathbf{T} \times_{\sigma_1} \mathbf{v}_1, \dots, \times_{\sigma_{|\sigma|}} \mathbf{v}_{|\sigma|}\|_1 \tag{10}$$

In the Appendix, we prove several lemmas regarding these tensor norms.

3.3 Concentration Bounds

$$\epsilon(\mathcal{O}_1, \dots, \mathcal{O}_d) = \left\| \widehat{\mathbb{P}}(\mathcal{O}_1, \dots, \mathcal{O}_d) - \mathbb{P}(\mathcal{O}_1, \dots, \mathcal{O}_d) \right\|_F \tag{11}$$

$$\epsilon(\mathcal{O}_1, \dots, \mathcal{O}_{d-e}, \mathbf{o}_{d-e+1}, \dots, \mathbf{o}_d) = \left\| \widehat{\mathbb{P}}(\mathcal{O}_1, \dots, \mathcal{O}_{d-e}, \mathbf{o}_{d-e+1}, \dots, \mathbf{o}_d) - \mathbb{P}(\mathcal{O}_1, \dots, \mathcal{O}_{d-e}, \mathbf{o}_{d-e+1}, \dots, \mathbf{o}_d) \right\|_F \tag{12}$$

$1, \dots, d-e$ denote the $d-e$ non-evidence variables while $d-e+1, \dots, d$ denote the e evidence variables. d indicates the total number of modes of the tensor. As the number of samples N gets large, we expect these quantities to be small.

Lemma 1 (variant of HKZ [2]) *If the algorithm independently samples N of the observations, then with probability at least $1 - \delta$.*

$$\epsilon(\mathcal{O}_1, \dots, \mathcal{O}_d) \leq \sqrt{\frac{1}{N} \ln \frac{2|\mathbb{C}|}{\delta}} + \sqrt{\frac{1}{N}} \quad (13)$$

$$\sum_{\mathbf{o}_{d-e+1}, \dots, \mathbf{o}_d} \epsilon(\mathcal{O}_1, \dots, \mathcal{O}_{d-e}, \mathbf{o}_{d-e+1}, \dots, \mathbf{o}_d) \leq \sqrt{\frac{k_o^{\epsilon_{\max}}}{N} \ln \frac{2|\mathbb{C}|}{\delta}} + \sqrt{\frac{k_o^{\epsilon_{\max}}}{N}} \quad (14)$$

for **all** tuples $(\mathcal{O}_1, \dots, \mathcal{O}_d)$ that are used in the spectral algorithm.

The proof is the same as that of HKZ [2] except the union bound is taken over $2|\mathbb{C}|$ instead of 3 (since each transformed quantity in the spectral algorithm is composed of at most two such terms). The last bound can be made tighter, identical to HKZ, but for simplicity we do not pursue that approach here.

3.4 Singular Value Bounds

Basically this is the generalized version of Lemma 9 in HKZ [2], which is stated below for completeness:

Lemma 2 (variant of HKZ [2]) *Suppose $\epsilon(\mathcal{O}_i, \mathcal{O}_{-i}) \leq \varepsilon \times \sigma_{\tau_i}(\mathbb{P}(\mathcal{O}_i, \mathcal{O}_{-i}))$ for some $\varepsilon < 1/2$. Let $\varepsilon_0 = \varepsilon(\mathcal{O}_i, \mathcal{O}_{-i})^2 / ((1 - \varepsilon)\sigma_{\tau_i}(\mathbb{P}(\mathcal{O}_i, \mathcal{O}_{-i})))^2$. Then:*

1. $\varepsilon_0 < 1$
2. $\sigma_{\tau_i}(\widehat{\mathbb{P}}(\mathcal{O}_i, \mathcal{O}_{-i}) \times_{\mathcal{O}_i} \widehat{\mathbf{U}}_i) \geq (1 - \varepsilon_0)\sigma_{\tau_i}(\mathbb{P}(\mathcal{O}_i, \mathcal{O}_{-i}))$
3. $\sigma_{\tau_i}(\mathbb{P}(\mathcal{O}_i, \mathcal{O}_{-i}) \times_{\mathcal{O}_i} \widehat{\mathbf{U}}_i) \geq \sqrt{1 - \varepsilon_0}\sigma_{\tau_i}(\mathbb{P}(\mathcal{O}_i, \mathcal{O}_{-i}))$
4. $\sigma_{\tau_i}(\widehat{\mathbb{P}}(\mathcal{O}_i | S_i) \times_{\mathcal{O}_i} \widehat{\mathbf{U}}_i) \geq \sqrt{1 - \varepsilon_0}\sigma_{\tau_i}(\mathbb{P}(\mathcal{O}_i | S_i))$

It follows that if $\epsilon(\mathcal{O}_i, \mathcal{O}_{-i}) \leq \sigma_{\tau_i}(\mathbb{P}(\mathcal{O}_i, \mathcal{O}_{-i}))/3$ then this implies that $\varepsilon_0 \leq \frac{1/9}{4/9} = \frac{1}{4}$. It then follows that,

1. $\sigma_{\tau_i}(\widehat{\mathbb{P}}(\mathcal{O}_i, \mathcal{O}_{-i}) \times_{\mathcal{O}_i} \widehat{\mathbf{U}}_i) \geq \frac{3}{4}\sigma_{\tau_i}(\mathbb{P}(\mathcal{O}_i, \mathcal{O}_{-i}))$
2. $\sigma_{\tau_i}(\mathbb{P}(\mathcal{O}_i, \mathcal{O}_{-i}) \times_{\mathcal{O}_i} \widehat{\mathbf{U}}_i) \geq \frac{\sqrt{3}}{2}\sigma_{\tau_i}(\mathbb{P}(\mathcal{O}_i, \mathcal{O}_{-i}))$
3. $\sigma_{\tau_i}(\widehat{\mathbb{P}}(\mathcal{O}_i | S_i) \times_{\mathcal{O}_i} \widehat{\mathbf{U}}_i) \geq \frac{\sqrt{3}}{2}\sigma_{\tau_i}(\mathbb{P}(\mathcal{O}_i | S_i))$

3.5 Bounding the Transformed Quantities

Define,

1. **root:** $\check{\mathcal{P}}(C_i) := \mathbb{P}(\mathcal{O}_{\alpha_i(1)}, \dots, \mathcal{O}_{\alpha_i(\gamma_i)}) \times_{\mathcal{O}_{\alpha_i(1)}} \widehat{\mathbf{U}}_{\alpha_i(1)} \dots \times_{\mathcal{O}_{\alpha_i(\gamma_i)}} \widehat{\mathbf{U}}_{\alpha_i(\gamma_i)}$
2. **leaf:** $\check{\mathcal{P}}_{r_i}(C_i) = \mathbb{P}(R_i = r_i, \mathcal{O}_{-i}) \times_{\mathcal{O}_{-i}} (\mathbb{P}_{\widehat{\mathbf{U}}}(\mathcal{O}_i, \mathcal{O}_{-i}))^\dagger$
3. **internal:** $\check{\mathcal{P}}(C_i) = \mathbb{P}(\mathcal{O}_{S_{\alpha_i(1)}}, \dots, \mathcal{O}_{\alpha_i(\gamma_i)}, \mathcal{O}_{-i}) \times_{\mathcal{O}_{-i}} (\mathbb{P}_{\widehat{\mathbf{U}}}(\mathcal{O}_i, \mathcal{O}_{-i}))^\dagger \times_{\mathcal{O}_{\alpha_i(1)}} \widehat{\mathbf{U}}_{\alpha_i(1)} \times \dots \times_{\mathcal{O}_{\alpha_i(\gamma_i)}} \widehat{\mathbf{U}}_{\alpha_i(\gamma_i)}$

(which are the observable quantities with the true probabilities, but empirical $\widehat{\mathbf{U}}$'s). Similarly, we can define $\check{\mathcal{F}}_i = \mathbb{P}[\mathcal{O}_i | S_i] \times_{\mathcal{O}_i} \widehat{\mathbf{U}}_i$. We have also abbreviated $\mathbb{P}(\mathcal{O}_i, \mathcal{O}_{-i}) \times_{\mathcal{O}_i} \widehat{\mathbf{U}}_i$ with $\mathbb{P}_{\widehat{\mathbf{U}}}(\mathcal{O}_i, \mathcal{O}_{-i})$.

We seek to bound the following three quantities:

$$\delta_i^{\text{root}} := \left\| (\widehat{\mathcal{P}}(C_i) - \check{\mathcal{P}}(C_i)) \times_{\mathcal{O}_{\alpha_i(1)}} \check{\mathcal{F}}_{\alpha_i(1)}^{-1}, \dots, \times_{\mathcal{O}_{\alpha_i(\gamma_i)}} \check{\mathcal{F}}_{\alpha_i(\gamma_i)}^{-1} \right\|_1 \quad (15)$$

$$\delta_i^{\text{internal}} := \left\| \left\| (\widehat{\mathcal{P}}(C_i) - \check{\mathcal{P}}(C_i)) \times_{\mathcal{O}_{-i}} \check{\mathcal{F}}_i \times_{\mathcal{O}_{\alpha_i(1)}} \check{\mathcal{F}}_i^{-1}, \dots, \times_{\mathcal{O}_{\alpha_i(\gamma_i)}} \check{\mathcal{F}}_{\alpha_i(\gamma_i)}^{-1} \right\| \right\|_1^{S_i} \quad (16)$$

$$\xi_i^{\text{leaf}} := \sum_{r_i} \left\| (\widehat{\mathcal{P}}_{r_i}(C_i) - \check{\mathcal{P}}_{r_i}(C_i)) \times_{\mathcal{O}_{-i}} \check{\mathcal{F}}_i \right\|_1 \quad (17)$$

Lemma 3 If $\epsilon(\mathcal{O}_i, \mathcal{O}_{-i}) \leq \sigma_{\tau_i}(\mathbb{P}(\mathcal{O}_i, \mathcal{O}_{-i}))/3$ then

$$\delta_i^{root} \leq \frac{2^{d_{\max}} k_h^{d_{\max}} \epsilon(\mathcal{O}_{\alpha_i(1)}, \dots, \mathcal{O}_{\alpha_i(\gamma_i)})}{3^{d/2} \beta^d} \quad (18)$$

$$\delta_i^{internal} \leq \frac{2^{d_{\max}+3} k_h^{d_{\max}}}{3\sqrt{3}^d \beta^d} \left(\frac{\epsilon(\mathcal{O}_{\alpha_i(1)}, \dots, \mathcal{O}_{\alpha_i(\gamma_i)}, \mathcal{O}_{-i})}{\sigma_{\tau_i}(\mathbb{P}(\mathcal{O}_i, \mathcal{O}_{-i}))} + \frac{\epsilon(\mathcal{O}_i, \mathcal{O}_{-i})}{(\sigma_{\tau_i}(\mathbb{P}(\mathcal{O}_i, \mathcal{O}_{-i})))^2} \right) \quad (19)$$

$$\xi_i^{leaf} \leq \frac{8k_h^{d_{\max}}}{3} \left(\frac{\epsilon(\mathcal{O}_i, \mathcal{O}_{-i})}{\sigma_{\tau_i}(\mathbb{P}(\mathcal{O}_i, \mathcal{O}_{-i}))^2} + \frac{\sum_{r_i} \epsilon(R_i = r_i, \mathcal{O}_{-i})}{\sigma_{\tau_i}(\mathbb{P}_{\hat{\mathbf{u}}}(\mathcal{O}_i, \mathcal{O}_{-i}))} \right) \quad (20)$$

We define $\Delta = \max(\delta_i^{root}, \delta_i^{internal}, \xi_i^{leaf})$ (over all i).

Proof

Case δ_i^{root} :

For the root, $\hat{\mathcal{P}}(C_i) = \hat{\mathbb{P}}(\mathcal{O}_{\alpha_i(1)}, \dots, \mathcal{O}_{\alpha_i(\gamma_i)}) \times_{\mathcal{O}_{\alpha_i(1)}} \hat{\mathbf{u}}_{\alpha_i(1)} \times \dots \times_{\mathcal{O}_{\alpha_i(\gamma_i)}} \hat{\mathbf{u}}_{\alpha_i(\gamma_i)}$ and similarly $\check{\mathcal{P}}(C_i) = \mathbb{P}(\mathcal{O}_{\alpha_i(1)}, \dots, \mathcal{O}_{\alpha_i(\gamma_i)}) \times_{\mathcal{O}_{\alpha_i(1)}} \check{\mathbf{u}}_{\alpha_i(1)} \times \dots \times_{\mathcal{O}_{\alpha_i(\gamma_i)}} \check{\mathbf{u}}_{\alpha_i(\gamma_i)}$.

$$\begin{aligned} \delta_i^{root} &= \left\| \left(\hat{\mathbb{P}}(\mathcal{O}_{\alpha_i(1)}, \dots, \mathcal{O}_{\alpha_i(\gamma_i)}) \times_{\mathcal{O}_{\alpha_i(1)}} \hat{\mathbf{u}}_{\alpha_i(1)} \times \dots \times_{\mathcal{O}_{\alpha_i(\gamma_i)}} \hat{\mathbf{u}}_{\alpha_i(\gamma_i)} \right. \right. \\ &\quad \left. \left. - \mathbb{P}(\mathcal{O}_{\alpha_i(1)}, \dots, \mathcal{O}_{\alpha_i(\gamma_i)}) \times_{\mathcal{O}_{\alpha_i(1)}} \check{\mathbf{u}}_{\alpha_i(1)} \times \dots \times_{\mathcal{O}_{\alpha_i(\gamma_i)}} \check{\mathbf{u}}_{\alpha_i(\gamma_i)} \right) \times_{\mathcal{O}_{\alpha_i(1)}} \check{\mathcal{F}}_{\alpha_i(1)}^{-1}, \dots, \times_{\mathcal{O}_{\alpha_i(\gamma_i)}} \check{\mathcal{F}}_{\alpha_i(\gamma_i)}^{-1} \right\|_{1R} \\ &\leq k_h^{d_{\max}} \left\| \left(\hat{\mathbb{P}}(\mathcal{O}_{\alpha_i(1)}, \dots, \mathcal{O}_{\alpha_i(\gamma_i)}) - \mathbb{P}(\mathcal{O}_{\alpha_i(1)}, \dots, \mathcal{O}_{\alpha_i(\gamma_i)}) \right) \times_{\mathcal{O}_{\alpha_i(1)}} \hat{\mathbf{u}}_{\alpha_i(1)} \times \dots \times_{\mathcal{O}_{\alpha_i(\gamma_i)}} \hat{\mathbf{u}}_{\alpha_i(\gamma_i)} \right\|_{2R} \\ &\quad \times \left\| \check{\mathcal{F}}_{\alpha_i(1)}^{-1} \right\|_{2R} \dots \left\| \check{\mathcal{F}}_{\alpha_i(\gamma_i)}^{-1} \right\|_{2R} \\ &\leq \frac{k_h^{d_{\max}}}{\prod_{j=1}^{\gamma_i} \sigma_{\tau_{\alpha_i(j)}}(\check{\mathcal{F}}_{\alpha_i(j)})} \left\| \left(\hat{\mathbb{P}}(\mathcal{O}_{\alpha_i(1)}, \dots, \mathcal{O}_{\alpha_i(\gamma_i)}) - \mathbb{P}(\mathcal{O}_{\alpha_i(1)}, \dots, \mathcal{O}_{\alpha_i(\gamma_i)}) \right) \times_{\mathcal{O}_{\alpha_i(1)}} \hat{\mathbf{u}}_{\alpha_i(1)} \times \dots \times_{\mathcal{O}_{\alpha_i(\gamma_i)}} \hat{\mathbf{u}}_{\alpha_i(\gamma_i)} \right\|_{2R} \\ &\leq \frac{k_h^{d_{\max}}}{\prod_{j=1}^{\gamma_i} \sigma_{\tau_{\alpha_i(j)}}(\check{\mathcal{F}}_{\alpha_i(j)})} \left\| \left(\hat{\mathbb{P}}(\mathcal{O}_{\alpha_i(1)}, \dots, \mathcal{O}_{\alpha_i(\gamma_i)}) - \mathbb{P}(\mathcal{O}_{\alpha_i(1)}, \dots, \mathcal{O}_{\alpha_i(\gamma_i)}) \right) \right\|_{2R} \times \\ &\quad \left\| \hat{\mathbf{u}}_{\alpha_i(1)} \right\|_{2R} \left\| \hat{\mathbf{u}}_{\alpha_i(2)} \right\|_{2R} \dots \left\| \hat{\mathbf{u}}_{\alpha_i(\gamma_i)} \right\|_{2R} \\ &= \frac{k_h^{d_{\max}}}{\prod_{j=1}^{\gamma_i} \sigma_{\tau_{\alpha_i(j)}}(\check{\mathcal{F}}_{\alpha_i(j)})} \epsilon(\mathcal{O}_{\alpha_i(1)}, \dots, \mathcal{O}_{\alpha_i(\gamma_i)}) \end{aligned}$$

Between the first and second line we use Lemma 8 to convert from elementwise one norm to spectral norm, and Lemma 6 (submultiplicativity). Lemma 6 (submultiplicativity) is applied again to get to the second-to-last line. We also use the fact that $\left\| \hat{\mathbf{u}}_{\alpha_i(1)} \right\|_{2R} = 1$.

Thus, by Lemma 2

$$\delta_i^{root} \leq \frac{2^{d_{\max}} k_h^{d_{\max}} \epsilon(\mathcal{O}_{\alpha_i(1)}, \dots, \mathcal{O}_{\alpha_i(\gamma_i)})}{3^{d_{\max}/2} \beta^{d_{\max}}} \quad (21)$$

Case ξ_i^{leaf} :

We note that $\hat{\mathcal{P}}_{r_i}(C_i) = \hat{\mathbb{P}}(R_i = r_i, \mathcal{O}_{-i}) \times_{\mathcal{O}_{-i}} (\hat{\mathbb{P}}_{\hat{\mathbf{u}}}(\mathcal{O}_i, \mathcal{O}_{-i}))^\dagger$ and similarly $\check{\mathcal{P}}_{r_i}(C_i) = \mathbb{P}(R_i = r_i, \mathcal{O}_{-i}) \times_{\mathcal{O}_{-i}} (\mathbb{P}_{\check{\mathbf{u}}}(\mathcal{O}_i, \mathcal{O}_{-i}))^\dagger$.

Again, we use Lemma 8 to convert from the one norm to the spectral norm, and Lemma 6 for submulti-

plicativity.

$$\begin{aligned}
\xi_i^{leaf} &= \sum_{r_i} \left\| \left\| (\widehat{\mathcal{P}}_{r_i}(C_i) - \check{\mathcal{P}}_{r_i}(C_i)) \times_{\mathcal{O}_{-i}} \check{\mathcal{F}}_i \right\|_1^{S_i} \right\| \\
&= \sum_{r_i} \left\| \left\| (\widehat{\mathcal{P}}_{r_i}(C_i) - \check{\mathcal{P}}_{r_i}(C_i)) \times_{\mathcal{O}_{-i}} \widehat{\mathbf{u}}_i \right\|_1^{S_i} \right\| \|\mathbb{P}[\mathcal{O}_i | S_i]\|_1^{S_i} \\
&\leq \sum_{r_i} k_h^{d_{\max}} \left\| \left\| (\widehat{\mathcal{P}}_{r_i}(C_i) - \check{\mathcal{P}}_{r_i}(C_i)) \right\|_{2R} \right\| \left\| \widehat{\mathbf{u}}_i \right\|_{2R} \\
&\leq \sum_{r_i} k_h^{d_{\max}} \left\| \left\| (\widehat{\mathcal{P}}_{r_i}(C_i) - \check{\mathcal{P}}_{r_i}(C_i)) \right\|_{2R} \right\| \tag{22}
\end{aligned}$$

Note that $\|\mathbb{P}[\mathcal{O}_i | S_i]\|_1^{S_i} = 1$, and $\|\widehat{\mathbf{u}}_i\|_{2R} = 1$.

$$\begin{aligned}
\left\| \left\| \widehat{\mathcal{P}}_{r_i}(C_i) - \check{\mathcal{P}}_{r_i}(C_i) \right\|_{2R} \right\| &= \left\| \left\| \widehat{\mathbb{P}}(R_i = r_i, \mathcal{O}_{-i}) \times_{\mathcal{O}_{-i}} (\widehat{\mathbb{P}}\widehat{\mathbf{u}}(\mathcal{O}_i, \mathcal{O}_{-i}))^\dagger - \mathbb{P}(R_i = r_i, \mathcal{O}_{-i}) \times_{\mathcal{O}_{-i}} (\mathbb{P}\widehat{\mathbf{u}}(\mathcal{O}_i, \mathcal{O}_{-i}))^\dagger \right\|_{2R} \right\| \\
&\leq \left\| \left\| \mathbb{P}(R_i = r_i, \mathcal{O}_{-i}) \times_{\mathcal{O}_{-i}} (\widehat{\mathbb{P}}\widehat{\mathbf{u}}(\mathcal{O}_i, \mathcal{O}_{-i}))^\dagger - \mathbb{P}\widehat{\mathbf{u}}(\mathcal{O}_i, \mathcal{O}_{-i})^\dagger \right\|_{2R} \right\| + \\
&\quad \left\| \left\| (\widehat{\mathbb{P}}(R_i = r_i, \mathcal{O}_{-i}) - \mathbb{P}(R_i = r_i, \mathcal{O}_{-i})) \times_{\mathcal{O}_{-i}} \mathbb{P}\widehat{\mathbf{u}}(\mathcal{O}_i, \mathcal{O}_{-i})^\dagger \right\|_{2R} \right\| \\
&\leq \|\mathbb{P}(R_i = r_i, \mathcal{O}_{-i})\|_{2R} \left\| \left\| (\widehat{\mathbb{P}}\widehat{\mathbf{u}}(\mathcal{O}_i, \mathcal{O}_{-i}))^\dagger - \mathbb{P}\widehat{\mathbf{u}}(\mathcal{O}_i, \mathcal{O}_{-i})^\dagger \right\|_{2R} \right\| + \\
&\quad \|\mathbb{P}\widehat{\mathbf{u}}(\mathcal{O}_i, \mathcal{O}_{-i})^\dagger\|_{2R} \left\| \left\| (\widehat{\mathbb{P}}(R_i = r_i, \mathcal{O}_{-i}) - \mathbb{P}(R_i = r_i, \mathcal{O}_{-i})) \right\|_{2R} \right\| \\
&\leq \mathbb{P}(R_i = r_i) \frac{1 + \sqrt{5}}{2} \frac{\epsilon(\mathcal{O}_i, \mathcal{O}_{-i})}{\min(\sigma_{\tau_i}(\widehat{\mathbb{P}}\widehat{\mathbf{u}}(\mathcal{O}_i, \mathcal{O}_{-i})), \sigma_{\tau_{ij}}(\mathbb{P}\widehat{\mathbf{u}}(\mathcal{O}_i, \mathcal{O}_{-i})))^2} \\
&\quad + \frac{\epsilon(R_i = r_i, \mathcal{O}_{-i})}{\sigma_{\tau_i}(\mathbb{P}\widehat{\mathbf{u}}(\mathcal{O}_i, \mathcal{O}_{-i}))}
\end{aligned}$$

The last line follows from Eq. 35. We have also used the fact that $\|\mathbb{P}(R_i = r_i, \mathcal{O}_{-i})\|_{2R} \leq \|\mathbb{P}(R_i = r_i, \mathcal{O}_{-i})\|_F \leq \mathbb{P}[R_i = r]$ by Lemma 7. Thus, using Lemma 2,

$$\xi_i^{leaf} \leq \sum_{r_i} k_h^{d_{\max}} \left(\mathbb{P}(R_i = r_i) \frac{1 + \sqrt{5}}{2} \frac{\epsilon(\mathcal{O}_i, \mathcal{O}_{-i})}{\min(\sigma_{\tau_i}(\widehat{\mathbb{P}}\widehat{\mathbf{u}}(\mathcal{O}_i, \mathcal{O}_{-i})), \sigma_{\tau_i}(\mathbb{P}\widehat{\mathbf{u}}(\mathcal{O}_i, \mathcal{O}_{-i})))^2} + \frac{\epsilon(R_i = r_i, \mathcal{O}_{-i})}{\sigma_{\tau_i}(\mathbb{P}\widehat{\mathbf{u}}(\mathcal{O}_i, \mathcal{O}_{-i}))} \right) \tag{23}$$

$$\begin{aligned}
\xi_i^{leaf} &\leq \sum_{r_i} k_h^{d_{\max}} \left(\frac{1 + \sqrt{5}}{2} \frac{16\mathbb{P}(R_i = r_i)\epsilon(\mathcal{O}_i, \mathcal{O}_{-i})}{9\sigma_{\tau_i}(\mathbb{P}\widehat{\mathbf{u}}(\mathcal{O}_i, \mathcal{O}_{-i}))^2} + \frac{2\epsilon(R_i = r_i, \mathcal{O}_{-i})}{\sqrt{3}\sigma_{\tau_i}(\mathbb{P}\widehat{\mathbf{u}}(\mathcal{O}_i, \mathcal{O}_{-i}))} \right) \\
&\leq \frac{8k_h^{d_{\max}}}{3} \left(\frac{\epsilon(\mathcal{O}_i, \mathcal{O}_{-i})}{\sigma_{\tau_i}(\mathbb{P}(\mathcal{O}_i, \mathcal{O}_{-i}))^2} + \frac{\sum_{r_i} \epsilon(R_i = r_i, \mathcal{O}_{-i})}{\sigma_{\tau_i}(\mathbb{P}\widehat{\mathbf{u}}(\mathcal{O}_i, \mathcal{O}_{-i}))} \right) \tag{24}
\end{aligned}$$

3.5.1 $\delta_i^{internal}$

$$\check{\mathcal{P}}(C_i) = \mathbb{P}(\mathcal{O}_{\alpha_i(1)}, \dots, \mathcal{O}_{\alpha_i(\gamma_i)}, \mathcal{O}_{-i}) \times_{\mathcal{O}_{-i}} (\mathbb{P}\widehat{\mathbf{u}}(\mathcal{O}_i, \mathcal{O}_{-i}))^\dagger \times_{\mathcal{O}_{\alpha_i(1)}} \widehat{\mathbf{u}}_{\alpha_i(1)} \times \dots \times_{\mathcal{O}_{\alpha_i(\gamma_i)}} \widehat{\mathbf{u}}_{\alpha_i(\gamma_i)}.$$

$$\widehat{\mathcal{P}}(C_i) = \widehat{\mathbb{P}}(\mathcal{O}_{\alpha_i(1)}, \dots, \mathcal{O}_{\alpha_i(\gamma_i)}, \mathcal{O}_{-i}) \times_{\mathcal{O}_{-i}} (\widehat{\mathbb{P}}\widehat{\mathbf{u}}(\mathcal{O}_i, \mathcal{O}_{-i}))^\dagger \times_{\mathcal{O}_{\alpha_i(1)}} \widehat{\mathbf{u}}_{\alpha_i(1)} \times \dots \times_{\mathcal{O}_{\alpha_i(\gamma_i)}} \widehat{\mathbf{u}}_{\alpha_i(\gamma_i)}.$$

Again, we use Lemma 8 to convert from one norm to spectral norm and Lemma 6 for submultiplicativity.

$$\begin{aligned}
\delta_i^{internal} &= \left\| \left(\widehat{\mathcal{P}}(C_i) - \check{\mathcal{P}}(C_i) \right) \times_{\mathcal{O}_{-i}} \check{\mathcal{F}}_i \times_{\mathcal{O}_{\alpha_i(1)}} \check{\mathcal{F}}_{\alpha_i(1)}^{-1}, \dots, \times_{\mathcal{O}_{\alpha_i(\gamma_i)}} \check{\mathcal{F}}_{\alpha_i(\gamma_i)}^{-1} \right\|_1^{S_i} \\
&\leq \left\| \left(\widehat{\mathcal{P}}(C_i) - \check{\mathcal{P}}(C_i) \right) \times_{\mathcal{O}_{-i}} \widehat{\mathbf{u}}_i \times_{\mathcal{O}_{\alpha_i(1)}} \check{\mathcal{F}}_{\alpha_i(1)}^{-1}, \dots, \times_{\mathcal{O}_{\alpha_i(\gamma_i)}} \check{\mathcal{F}}_{\alpha_i(\gamma_i)}^{-1} \right\|_1^{S_i} \|\mathbb{P}[\mathcal{O}_i | S_i]\|_1^{S_i} \\
&\leq k_h^{d_{max}} \left\| \left(\widehat{\mathcal{P}}(C_i) - \check{\mathcal{P}}(C_i) \right) \right\|_{2R} \left\| \widehat{\mathbf{u}}_i \right\|_{2R} \left\| \check{\mathcal{F}}_{\alpha_i(1)}^{-1} \right\|_{2R} \dots \left\| \check{\mathcal{F}}_{\alpha_i(\gamma_i)}^{-1} \right\|_{2R} \\
&\leq \frac{k_h^{d_{max}}}{\prod_{j=1}^{\gamma_i} \sigma_{\tau_{\alpha_i(j)}}(\check{\mathcal{F}}_{\alpha_i(j)})} \left\| \left(\widehat{\mathcal{P}}(C_i) - \check{\mathcal{P}}(C_i) \right) \right\|_{2R}
\end{aligned} \tag{25}$$

Note that $\|\mathbb{P}[\mathcal{O}_i | S_i]\|_1^{S_i} = 1$, and $\left\| \widehat{\mathbf{u}}_i \right\|_{2R} = 1$.

$$\begin{aligned}
&\left\| \widehat{\mathcal{P}}(C_i) - \check{\mathcal{P}}(C_i) \right\|_{2R} \\
&= \left\| \widehat{\mathbb{P}}(\mathcal{O}_{\alpha_i(1)}, \dots, \mathcal{O}_{\alpha_i(\gamma_i)}, \mathcal{O}_{-i}) \times_{\mathcal{O}_{-i}} (\widehat{\mathbb{P}}_{\widehat{\mathbf{u}}}(\mathcal{O}_i, \mathcal{O}_{-i}))^\dagger \times_{\mathcal{O}_{\alpha_i(1)}} \widehat{\mathbf{u}}_{\alpha_i(1)} \times \dots \times_{\mathcal{O}_{\alpha_i(\gamma_i)}} \widehat{\mathbf{u}}_{\alpha_i(\gamma_i)} - \right. \\
&\quad \left. \mathbb{P}(\mathcal{O}_{\alpha_i(1)}, \dots, \mathcal{O}_{\alpha_i(\gamma_i)}, \mathcal{O}_{-i}) \times_{\mathcal{O}_{-i}} (\mathbb{P}_{\widehat{\mathbf{u}}}(\mathcal{O}_i, \mathcal{O}_{-i}))^\dagger \times_{\mathcal{O}_{\alpha_i(1)}} \widehat{\mathbf{u}}_{\alpha_i(1)} \times \dots \times_{\mathcal{O}_{\alpha_i(\gamma_i)}} \widehat{\mathbf{u}}_{\alpha_i(\gamma_i)} \right\|_{2R} \\
&\leq \left\| (\widehat{\mathbb{P}}(\mathcal{O}_{\alpha_i(1)}, \dots, \mathcal{O}_{\alpha_i(\gamma_i)}, \mathcal{O}_{-i}) - \mathbb{P}(\mathcal{O}_{\alpha_i(1)}, \dots, \mathcal{O}_{\alpha_i(\gamma_i)}, \mathcal{O}_{-i})) \times_{\mathcal{O}_{-i}} \mathbb{P}_{\widehat{\mathbf{u}}}(\mathcal{O}_i, \mathcal{O}_{-i})^\dagger \right. \\
&\quad \left. \times_{\mathcal{O}_{\alpha_i(1)}} \widehat{\mathbf{u}}_{\alpha_i(1)} \times \dots \times_{\mathcal{O}_{\alpha_i(\gamma_i)}} \widehat{\mathbf{u}}_{\alpha_i(\gamma_i)} \right\|_{2R} \\
&\quad + \left\| \mathbb{P}(\mathcal{O}_{\alpha_i(1)}, \dots, \mathcal{O}_{\alpha_i(\gamma_i)}, \mathcal{O}_{-i}) \times_{\mathcal{O}_{-i}} ((\widehat{\mathbb{P}}_{\widehat{\mathbf{u}}}(\mathcal{O}_i, \mathcal{O}_{-i}))^\dagger - (\mathbb{P}_{\widehat{\mathbf{u}}}(\mathcal{O}_i, \mathcal{O}_{-i}))^\dagger) \right. \\
&\quad \left. \times_{\mathcal{O}_{\alpha_i(1)}} \widehat{\mathbf{u}}_{\alpha_i(1)} \times \dots \times_{\mathcal{O}_{\alpha_i(\gamma_i)}} \widehat{\mathbf{u}}_{\alpha_i(\gamma_i)} \right\|_{2R} \\
&\leq \left\| (\widehat{\mathbb{P}}(\mathcal{O}_{\alpha_i(1)}, \dots, \mathcal{O}_{\alpha_i(\gamma_i)}, \mathcal{O}_{-i}) - \mathbb{P}(\mathcal{O}_{\alpha_i(1)}, \dots, \mathcal{O}_{\alpha_i(\gamma_i)}, \mathcal{O}_{-i})) \right\|_{2R} \left\| \mathbb{P}_{\widehat{\mathbf{u}}}(\mathcal{O}_i, \mathcal{O}_{-i})^\dagger \right\|_{2R} \left\| \widehat{\mathbf{u}}_{\alpha_i(1)} \right\|_{2R} \dots \left\| \widehat{\mathbf{u}}_{\alpha_i(\gamma_i)} \right\|_{2R} \\
&\quad + \left\| \mathbb{P}(\mathcal{O}_{\alpha_i(1)}, \dots, \mathcal{O}_{\alpha_i(\gamma_i)}, \mathcal{O}_{-i}) \right\|_{2R} \left\| (\widehat{\mathbb{P}}_{\widehat{\mathbf{u}}}(\mathcal{O}_i, \mathcal{O}_{-i}))^\dagger - \mathbb{P}_{\widehat{\mathbf{u}}}(\mathcal{O}_i, \mathcal{O}_{-i})^\dagger \right\|_{2R} \left\| \widehat{\mathbf{u}}_{\alpha_i(1)} \right\|_{2R} \dots \left\| \widehat{\mathbf{u}}_{\alpha_i(\gamma_i)} \right\|_{2R} \\
&\leq \frac{\epsilon(\mathcal{O}_{\alpha_i(1)}, \dots, \mathcal{O}_{\alpha_i(\gamma_i)}, \mathcal{O}_{-i})}{\sigma_{\tau_i}(\mathbb{P}_{\widehat{\mathbf{u}}}(\mathcal{O}_i, \mathcal{O}_{-i}))} + \frac{1 + \sqrt{5}}{2} \frac{\epsilon(\mathcal{O}_i, \mathcal{O}_{-i})}{\min(\sigma_{\tau_i}(\mathbb{P}_{\widehat{\mathbf{u}}}(\mathcal{O}_i, \mathcal{O}_{-i})), \sigma_{\tau_i}(\widehat{\mathbb{P}}_{\widehat{\mathbf{u}}}(\mathcal{O}_i, \mathcal{O}_{-i})))^2}
\end{aligned}$$

The last line follows from Eq. 35. We have also used the fact that

$\left\| \mathbb{P}(\mathcal{O}_{\alpha_i(1)}, \dots, \mathcal{O}_{\alpha_i(\gamma_i)}, \mathcal{O}_{-i}) \right\|_{2R} \leq \|\mathbb{P}(\mathcal{O}_{\alpha_i(1)}, \dots, \mathcal{O}_{\alpha_i(\gamma_i)}, \mathcal{O}_{-i})\|_F \leq 1$ via Lemma 7. Using Lemma 2,

$$\delta_i^{internal} \leq \frac{(2k_h)^{d_{max}}}{(\beta\sqrt{3})^{d_{max}}} \left(\frac{\epsilon(\mathcal{O}_{\alpha_i(1)}, \dots, \mathcal{O}_{\alpha_i(\gamma_i)}, \mathcal{O}_{-i})}{\sigma_{\tau_i}(\mathbb{P}_{\widehat{\mathbf{u}}}(\mathcal{O}_i, \mathcal{O}_{-i}))} + \frac{1 + \sqrt{5}}{2} \frac{\epsilon(\mathcal{O}_i, \mathcal{O}_{-i})}{\min(\sigma_{\tau_i}(\mathbb{P}_{\widehat{\mathbf{u}}}(\mathcal{O}_i, \mathcal{O}_{-i})), \sigma_{\tau_i}(\widehat{\mathbb{P}}_{\widehat{\mathbf{u}}}(\mathcal{O}_i, \mathcal{O}_{-i})))^2} \right) \tag{26}$$

$$\begin{aligned}
\delta_i^{internal} &\leq \frac{(2k_h)^{d_{max}}}{(\beta\sqrt{3})^{d_{max}}} \left(\frac{2\epsilon(\mathcal{O}_{\alpha_i(1)}, \dots, \mathcal{O}_{\alpha_i(\gamma_i)}, \mathcal{O}_{-i})}{\sqrt{3}\sigma_{\tau_i}(\mathbb{P}(\mathcal{O}_i, \mathcal{O}_{-i}))} + \frac{8\epsilon(\mathcal{O}_i, \mathcal{O}_{-i})}{3(\sigma_{\tau_i}(\mathbb{P}(\mathcal{O}_i, \mathcal{O}_{-i})))^2} \right) \\
&\leq \frac{8(2k_h)^{d_{max}}}{3(\beta\sqrt{3})^{d_{max}}} \left(\frac{\epsilon(\mathcal{O}_{\alpha_i(1)}, \dots, \mathcal{O}_{\alpha_i(\gamma_i)}, \mathcal{O}_{-i})}{\sigma_{\tau_i}(\mathbb{P}(\mathcal{O}_i, \mathcal{O}_{-i}))} + \frac{\epsilon(\mathcal{O}_i, \mathcal{O}_{-i})}{(\sigma_{\tau_i}(\mathbb{P}(\mathcal{O}_i, \mathcal{O}_{-i})))^2} \right)
\end{aligned} \tag{27}$$

■

3.6 Bounding the Propagation of Error

We now show all these errors propagate on the junction tree. For this section, assume the clique nodes are numbered $1, 2, \dots, |\mathbb{C}|$ in breadth first order (such that 1 is the root). Moreover let $\Phi_{1:c}(\mathcal{C})$ be the transformed

factors accumulated so far if we computed the joint probability from the root down (instead of the bottom up). For example,

$$\Phi_{1:1}(\mathcal{C}) = \mathcal{P}(\mathcal{C}_1) \quad (28)$$

$$\Phi_{1:2}(\mathcal{C}) = \mathcal{P}(\mathcal{C}_1) \times_{S_2} \mathcal{P}(\mathcal{C}_2) \quad (29)$$

$$\Phi_{1:2}(\mathcal{C}) = \mathcal{P}(\mathcal{C}_1) \times_{S_2} \mathcal{P}(\mathcal{C}_2) \times_{S_3} \mathcal{P}(\mathcal{C}_3) \quad (30)$$

(Note how this is very computationally inefficient: the tensors get very large. However, it is useful for proving statistical properties). Then the modes of $\Phi_{1:c}$ can be partitioned into mode groups, $\mathcal{M}_1, \dots, \mathcal{M}_{d_c}$ (where each mode group consists of the variables on the corresponding separator edge). We now prove the following lemma,

Lemma 4 Define $\Delta = \max(\delta_i^{root}, \delta_i^{internal}, \xi_i^{leaf})$. Then,

$$\sum_{\mathbf{x}_{1:c}} \|(\check{\Phi}_{1:c}(\mathcal{C}) - \widehat{\Phi}_{1:c}(\mathcal{C})) \times_{\mathcal{M}_1} \check{\mathcal{F}}_1^{-1}, \dots, \times_{\mathcal{M}_{d_c}} \check{\mathcal{F}}_{d_c}^{-1}\|_1 \leq (1 + \Delta)^{c-1} \delta_i^{root} + (1 + \Delta)^{c-1} - 1 \quad (31)$$

$\mathbf{x}_{1:c}$ is all the observed variables in cliques 1 through c . Note that when $c = |\mathcal{C}|$ then this implies that $\widehat{\Phi}_{1:c}(\mathcal{C}) = \widehat{\mathbb{P}}[x_1, \dots, x_O]$ and thus,

$$\sum_{\mathbf{x}} |\widehat{\mathbb{P}}[x_1, \dots, x_O] - \mathbb{P}[x_1, \dots, x_O]| \leq (1 + \Delta)^{|\mathcal{C}|-1} \delta_i^{root} + (1 + \Delta)^{|\mathcal{C}|-1} - 1 \quad (32)$$

Proof The proof is by induction on c . The base case follows trivially from the definition of δ_i^{root} . For the induction step, assume the claim holds for $c \geq 1$. Then we prove it holds for $c + 1$.

$$\begin{aligned} & \sum_{\mathbf{x}_{1:c+1}} \|(\check{\Phi}_{1:c+1}(\mathcal{C}) - \widehat{\Phi}_{1:c+1}(\mathcal{C})) \times_{\mathcal{M}_1} \check{\mathcal{F}}_1^{-1}, \dots, \times_{\mathcal{M}_{d_{c+1}}} \check{\mathcal{F}}_{d_{c+1}}^{-1}\|_1 \\ = & \sum_{\mathbf{x}_{1:c+1}} \left\| \left(\check{\Phi}_{1:c}(\mathcal{C}) \times (\widehat{\mathcal{P}}(\mathcal{C}_{c+1}) - \check{\mathcal{P}}(\mathcal{C}_{c+1})) + (\check{\Phi}_{1:c}(\mathcal{C}) - \widehat{\Phi}_{1:c}(\mathcal{C})) \times (\widehat{\mathcal{P}}(\mathcal{C}_{c+1}) - \check{\mathcal{P}}(\mathcal{C}_{c+1})) + (\check{\Phi}_{1:c}(\mathcal{C}) - \widehat{\Phi}_{1:c}(\mathcal{C})) \times \check{\mathcal{P}}(\mathcal{C}_{c+1}) \right) \right. \\ & \left. \times_{\mathcal{M}_1} \check{\mathcal{F}}_1^{-1}, \dots, \times_{\mathcal{M}_{d_{c+1}}} \check{\mathcal{F}}_{d_{c+1}}^{-1} \right\|_1 \\ \leq & \sum_{\mathbf{x}_{1:c+1}} \|(\check{\Phi}_{1:c}(\mathcal{C}) \times (\widehat{\mathcal{P}}(\mathcal{C}_{c+1}) - \check{\mathcal{P}}(\mathcal{C}_{c+1}))) \times_{\mathcal{M}_1} \check{\mathcal{F}}_1^{-1}, \dots, \times_{\mathcal{M}_{d_{c+1}}} \check{\mathcal{F}}_{d_{c+1}}^{-1}\|_1 + \\ & \sum_{\mathbf{x}_{1:c+1}} \|((\check{\Phi}_{1:c}(\mathcal{C}) - \widehat{\Phi}_{1:c}(\mathcal{C})) \times (\widehat{\mathcal{P}}(\mathcal{C}_{c+1}) - \check{\mathcal{P}}(\mathcal{C}_{c+1}))) \times_{\mathcal{M}_1} \check{\mathcal{F}}_1^{-1}, \dots, \times_{\mathcal{M}_{d_{c+1}}} \check{\mathcal{F}}_{d_{c+1}}^{-1}\|_1 + \\ & \sum_{\mathbf{x}_{1:c+1}} \|((\check{\Phi}_{1:c}(\mathcal{C}) - \widehat{\Phi}_{1:c}(\mathcal{C})) \times \check{\mathcal{P}}(\mathcal{C}_{c+1})) \times_{\mathcal{M}_1} \check{\mathcal{F}}_1^{-1}, \dots, \times_{\mathcal{M}_{d_{c+1}}} \check{\mathcal{F}}_{d_{c+1}}^{-1}\|_1 \end{aligned}$$

Now we consider two cases, when $c + 1$ is a leaf clique and when it is an internal clique.

Case 1: Internal Clique

Note here the summation over $\mathbf{x}_{1:c+1}$ is irrelevant since there is no evidence. We use Lemma 5 to break up the three terms:

The first term,

$$\begin{aligned} & \|(\check{\Phi}_{1:c}(\mathcal{C}) \times (\widehat{\mathcal{P}}(\mathcal{C}_{c+1}) - \check{\mathcal{P}}(\mathcal{C}_{c+1}))) \times_{\mathcal{M}_1} \check{\mathcal{F}}_1^{-1}, \dots, \times_{\mathcal{M}_{d_{c+1}}} \check{\mathcal{F}}_{d_{c+1}}^{-1}\|_1 \\ \leq & \left\| \left((\widehat{\mathcal{P}}(\mathcal{C}_{c+1}) - \check{\mathcal{P}}(\mathcal{C}_{c+1})) \times_{\mathcal{O}_{-(c+1)}} \check{\mathcal{F}}_{c+1} \times_{\mathcal{O}_{\alpha_{c+1}(1)}} \check{\mathcal{F}}_{\alpha_{c+1}(1)}^{-1}, \dots, \times_{\mathcal{O}_{\alpha_{c+1}(\gamma_{c+1})}} \check{\mathcal{F}}_{\alpha_{c+1}(\gamma_{c+1})}^{-1} \right) \right\|_1^{S_i} \\ \times & \| \check{\Phi}_{1:c}(\mathcal{C}) \times_{\mathcal{M}_1} \check{\mathcal{F}}_1^{-1}, \dots, \times_{\mathcal{M}_{d_c}} \check{\mathcal{F}}_{d_c}^{-1} \|_1 \\ \leq & \Delta \times 1 \end{aligned}$$

The first term above is simply $\delta_i^{internal} \leq \Delta$ while the second equals one since it is a joint distribution.

Now for the second term,

$$\begin{aligned}
& \|(\check{\Phi}_{1:c}(\mathcal{C}) - \widehat{\Phi}_{1:c}(\mathcal{C})) \times (\widehat{\mathcal{P}}(\mathcal{C}_{c+1}) - \check{\mathcal{P}}(\mathcal{C}_{c+1})) \times_{\mathcal{M}_1} \check{\mathcal{F}}_1^{-1}, \dots, \times_{\mathcal{M}_{d_{c+1}}} \check{\mathcal{F}}_{d_{c+1}}^{-1}\|_1 \\
& \leq \|(\check{\Phi}_{1:c}(\mathcal{C}) - \widehat{\Phi}_{1:c}(\mathcal{C})) \times_{\mathcal{M}_1} \check{\mathcal{F}}_1^{-1}, \dots, \times_{\mathcal{M}_{d_c}} \check{\mathcal{F}}_{d_c}^{-1}\|_1 \times \\
& \quad \left\| \left(\widehat{\mathcal{P}}(\mathcal{C}_{c+1}) - \check{\mathcal{P}}(\mathcal{C}_{c+1}) \right) \times_{\mathcal{O}_{-(c+1)}} \check{\mathcal{F}}_{c+1} \times_{\mathcal{O}_{\alpha_{c+1}(1)}} \check{\mathcal{F}}_{\alpha_{c+1}(1)}^{-1}, \dots, \times_{\mathcal{O}_{\alpha_{c+1}(\gamma_{c+1})}} \check{\mathcal{F}}_{\alpha_{c+1}(\gamma_{c+1})}^{-1} \right\|_1^{S_{c+1}} \\
& \leq ((1 + \Delta)^{c-1} \delta^{root} + (1 + \Delta)^{c-1} - 1) \times \delta_{c+1}^{internal} \quad (\text{via induction hypothesis}) \\
& \leq ((1 + \Delta)^{c-1} \delta^{root} + (1 + \Delta)^{c-1} - 1) \times \Delta
\end{aligned}$$

The third term,

$$\begin{aligned}
& \|((\check{\Phi}_{1:c}(\mathcal{C}) - \widehat{\Phi}_{1:c}(\mathcal{C})) \times \check{\mathcal{P}}(\mathcal{C}_{c+1})) \times_{\mathcal{M}_1} \check{\mathcal{F}}_1^{-1}, \dots, \times_{\mathcal{M}_{d_{c+1}}} \check{\mathcal{F}}_{d_{c+1}}^{-1}\|_1 \\
& \leq \|(\check{\Phi}_{1:c}(\mathcal{C}) - \widehat{\Phi}_{1:c}(\mathcal{C})) \times_{\mathcal{M}_1} \check{\mathcal{F}}_1^{-1}, \dots, \times_{\mathcal{M}_{d_c}} \check{\mathcal{F}}_{d_c}^{-1}\|_1 \times \\
& \quad \left\| \check{\mathcal{P}}(\mathcal{C}_{c+1}) \times_{\mathcal{O}_{-(c+1)}} \check{\mathcal{F}}_{c+1} \times_{\mathcal{O}_{\alpha_{c+1}(1)}} \check{\mathcal{F}}_{\alpha_{c+1}(1)}^{-1}, \dots, \times_{\mathcal{O}_{\alpha_{c+1}(\gamma_{c+1})}} \check{\mathcal{F}}_{\alpha_{c+1}(\gamma_{c+1})}^{-1} \right\|_1^{S_{c+1}} \\
& \leq ((1 + \Delta)^{c-1} \delta^{root} + (1 + \Delta)^{c-1} - 1) \times 1
\end{aligned}$$

Case 2: Leaf Clique

Again we use Lemma 5 to break up the three terms:

The first term,

$$\begin{aligned}
& \sum_{\mathbf{x}_{1:c+1}} \|(\check{\Phi}_{1:c}(\mathcal{C}) \times (\widehat{\mathcal{P}}_{r_i}(\mathcal{C}_{c+1}) - \check{\mathcal{P}}_{r_i}(\mathcal{C}_{c+1}))) \times_{\mathcal{M}_1} \check{\mathcal{F}}_1^{-1}, \dots, \times_{\mathcal{M}_{d_{c+1}}} \check{\mathcal{F}}_{d_{c+1}}^{-1}\|_1 \\
& \leq \sum_{\mathbf{x}_{1:c+1}} \|(\check{\Phi}_{1:c}(\mathcal{C}) \times_{\mathcal{M}_1} \check{\mathcal{F}}_1^{-1}, \dots, \times_{\mathcal{M}_{d_c}} \check{\mathcal{F}}_{d_c}^{-1})\|_1 \left\| (\widehat{\mathcal{P}}_{r_i}(\mathcal{C}_{c+1}) - \check{\mathcal{P}}_{r_i}(\mathcal{C}_{c+1})) \times_{\mathcal{O}_{-(c+1)}} \check{\mathcal{F}}_{c+1} \right\|_1^{S_{c+1}} \\
& \leq \sum_{\mathbf{x}_{1:c}} \|(\check{\Phi}_{1:c}(\mathcal{C}) \times_{\mathcal{M}_1} \check{\mathcal{F}}_1^{-1}, \dots, \times_{\mathcal{M}_{d_c}} \check{\mathcal{F}}_{d_c}^{-1})\|_1 \sum_{\mathbf{x}_{c+1}} \left\| (\widehat{\mathcal{P}}_{r_i}(\mathcal{C}_{c+1}) - \check{\mathcal{P}}_{r_i}(\mathcal{C}_{c+1})) \times_{\mathcal{O}_{-(c+1)}} \check{\mathcal{F}}_{c+1} \right\|_1^{S_{c+1}} \\
& \leq 1 \times \Delta
\end{aligned}$$

The first term above equals 1 because it is a joint distribution and the second is the bound on the transformed quantity we had proved earlier (since $r_i = \mathbf{x}_{c+1}$).

The second term,

$$\begin{aligned}
& \sum_{\mathbf{x}_{1:c+1}} \|(\check{\Phi}_{1:c}(\mathcal{C}) - \widehat{\Phi}_{1:c}(\mathcal{C})) \times (\widehat{\mathcal{P}}_{r_i}(\mathcal{C}_{c+1}) - \check{\mathcal{P}}_{r_i}(\mathcal{C}_{c+1})) \times_{\mathcal{M}_1} \check{\mathcal{F}}_1^{-1}, \dots, \times_{\mathcal{M}_{d_{c+1}}} \check{\mathcal{F}}_{d_{c+1}}^{-1}\|_1 \\
& \leq \sum_{\mathbf{x}_{1:c}} \|(\check{\Phi}_{1:c}(\mathcal{C}) - \widehat{\Phi}_{1:c}(\mathcal{C})) \times_{\mathcal{M}_1} \check{\mathcal{F}}_1^{-1}, \dots, \times_{\mathcal{M}_{d_c}} \check{\mathcal{F}}_{d_c}^{-1}\|_1 \sum_{\mathbf{x}_{c+1}} \left\| (\widehat{\mathcal{P}}_{r_i}(\mathcal{C}_{c+1}) - \check{\mathcal{P}}_{r_i}(\mathcal{C}_{c+1})) \times_{\mathcal{O}_{-(c+1)}} \check{\mathcal{F}}_{c+1} \right\|_1^{S_{c+1}} \\
& \leq ((1 + \Delta)^{c-1} \delta^{root} + (1 + \Delta)^{c-1} - 1) \times \xi_{c+1}^{leaf} \\
& \leq ((1 + \Delta)^{c-1} \delta^{root} + (1 + \Delta)^{c-1} - 1) \times \Delta
\end{aligned}$$

The third term,

$$\begin{aligned}
& \sum_{\mathbf{x}_{1:c+1}} \|((\check{\Phi}_{1:c}(\mathcal{C}) - \widehat{\Phi}_{1:c}(\mathcal{C})) \times \check{\mathcal{P}}_{r_i}(\mathcal{C}_{c+1})) \times_{\mathcal{M}_1} \check{\mathcal{F}}_1^{-1}, \dots, \times_{\mathcal{M}_{d_{c+1}}} \check{\mathcal{F}}_{d_{c+1}}^{-1}\|_1 \\
& \leq \sum_{\mathbf{x}_{1:c}} \|(\check{\Phi}_{1:c}(\mathcal{C}) - \widehat{\Phi}_{1:c}(\mathcal{C})) \times_{\mathcal{M}_1} \check{\mathcal{F}}_1^{-1}, \dots, \times_{\mathcal{M}_{d_c}} \check{\mathcal{F}}_{d_c}^{-1}\|_1 \sum_{\mathbf{x}_{c+1}} \left\| \check{\mathcal{P}}_{r_i}(\mathcal{C}_{c+1}) \times_{\mathcal{O}_{-(c+1)}} \check{\mathcal{F}}_{c+1} \right\|_1^{S_{c+1}} \\
& \leq ((1 + \Delta)^{c-1} \delta^{root} + (1 + \Delta)^{c-1} - 1) \times 1
\end{aligned}$$

Combining these terms proves the induction step. ■

3.7 Putting it all together

We use the fact from HKZ [2] that $(1 + a/t)^t \leq 1 + 2a$ for $a \leq 1/2$. Now Δ is the main source of error. We set $\Delta \leq O(\epsilon_{total}/|\mathbb{C}|)$.

Note that

$$\Delta \leq \frac{2^{d_{\max}+3} k_h^{d_{\max}}}{3\sqrt{3}^{d_{\max}} \beta^{d_{\max}}} \left(\frac{\sum_{\mathbf{o}_{d-e+1}, \dots, \mathbf{o}_d} \epsilon(\theta_1, \dots, \theta_{d-e}, \mathbf{o}_{d-e+1}, \dots, \mathbf{o}_d)}{\alpha} + \frac{\sum_{\mathbf{o}_{d-e+1}, \dots, \mathbf{o}_d} \epsilon(\theta_1, \dots, \theta_{d-e}, \mathbf{o}_{d-e+1}, \dots, \mathbf{o}_d)}{\alpha^2} \right)$$

This gives,

$$\frac{2^{d_{\max}+3} k_h^{d_{\max}}}{3\sqrt{3}^{d_{\max}} \beta^{d_{\max}}} \left(\frac{\sum_{\mathbf{o}_{d-e+1}, \dots, \mathbf{o}_d} \epsilon(\theta_1, \dots, \theta_{d-e}, \mathbf{o}_{d-e+1}, \dots, \mathbf{o}_d)}{\alpha} + \frac{\sum_{\mathbf{o}_{d-e+1}, \dots, \mathbf{o}_d} \epsilon(\theta_1, \dots, \theta_{d-e}, \mathbf{o}_{d-e+1}, \dots, \mathbf{o}_d)}{\alpha^2} \right) \leq K \epsilon_{total}/|\mathbb{C}|$$

where K is some constant.

This implies,

$$\sum_{\mathbf{o}_{d-e+1}, \dots, \mathbf{o}_d} \epsilon(\theta_1, \dots, \theta_{d-e}, \mathbf{o}_{d-e+1}, \dots, \mathbf{o}_d) \leq K \frac{3^{d_{\max}/2+1} \epsilon_{total} \alpha^2 \beta^{d_{\max}}}{2^{d_{\max}+3} k_h^{d_{\max}} |\mathbb{C}|} \quad (33)$$

Now using the concentration bound (Lemma 1) will give,

$$K \frac{3^{d_{\max}/2+1} \epsilon_{total} \alpha^2 \beta^{d_{\max}}}{2^{d_{\max}+3} k_h^{d_{\max}} |\mathbb{C}|} \leq \sqrt{\frac{k_o^{e_{\max}}}{N} \ln \frac{2|\mathbb{C}|}{\delta}} + \sqrt{\frac{k_o^{e_{\max}}}{N}}$$

Solving for N :

$$N \geq O \left(\left(\frac{4k_h^2}{3\beta^2} \right)^{d_{\max}} \frac{k_o^{e_{\max}} \ln \frac{|\mathbb{C}|}{\delta} |\mathbb{C}|^2}{\epsilon_{total}^2 \alpha^4} \right) \quad (34)$$

and this completes the proof.

4 Appendix

4.1 Matrix Perturbation Bounds

This is Theorem 3.8 from pg. 143 in Stewart and Sun, 1990 [4]. Let $A \in \mathbf{R}^{m \times n}$, with $m \geq n$ and let $\tilde{A} = A + E$. Then

$$\left\| \tilde{A}^+ - A^+ \right\|_2 \leq \frac{1 + \sqrt{5}}{2} \max(\|A^+\|_2^2, \|\tilde{A}\|_2^2) \|E\|_2 \quad (35)$$

4.2 Tensor Norm Bounds

For matrices it is true that $\|\mathbf{M}\mathbf{v}\|_1 \leq \|\mathbf{M}\|_1 \|\mathbf{v}\|_1$. We prove the generalization to tensors.

Lemma 5 *Let \mathbf{T}_1 and \mathbf{T}_2 be tensors σ a set of (labeled) modes.*

$$\|\mathbf{T}_1 \times_{\sigma} \mathbf{T}_2\|_1 \leq \|\mathbf{T}_2\|_1^{\sigma} \|\mathbf{T}_1\|_1 \quad (36)$$

Proof

$$\begin{aligned}
\|\mathbf{T}_1 \times_{\sigma} \mathbf{T}_2\|_1 &= \sum_{\mathbf{i}_{1:N} \setminus \sigma} \sum_{\mathbf{j}_{1:M} \setminus \sigma} \sum_{\mathbf{x}} \mathbf{T}_1(\mathbf{i}_{1:N} \setminus \sigma, \sigma = \mathbf{x}) \mathbf{T}_2(\mathbf{j}_{1:M} \setminus \sigma, \sigma = \mathbf{x}) \\
&= \sum_{\mathbf{i}_{1:N} \setminus \sigma} \sum_{\mathbf{x}} \sum_{\mathbf{j}_{1:M} \setminus \sigma} \mathbf{T}_1(\mathbf{i}_{1:N} \setminus \sigma, \sigma = \mathbf{x}) \mathbf{T}_2(\mathbf{j}_{1:M} \setminus \sigma, \sigma = \mathbf{x}) \\
&= \sum_{\mathbf{i}_{1:N} \setminus \sigma} \sum_{\sigma} \mathbf{T}_1(\mathbf{i}_{1:N} \setminus \sigma, \sigma = \mathbf{x}) \sum_{\mathbf{j}_{1:M} \setminus \sigma} \mathbf{T}_2(\mathbf{j}_{1:M} \setminus \sigma, \sigma = \mathbf{x}) \\
&\leq \max_{\mathbf{x}} \left(\sum_{\mathbf{j}_{1:M} \setminus \sigma} \mathbf{T}_2(\mathbf{j}_{1:M} \setminus \sigma, \sigma = \mathbf{x}) \right) \|\mathbf{T}_1\|_1 \\
&= \|\mathbf{T}_2\|_1^{\sigma} \|\mathbf{T}_1\|_1
\end{aligned}$$

■

We prove a restricted analog of the fact that spectral norm is submultiplicative for matrices i.e. $\|AB\|_2 \leq \|A\|_2 \|B\|_2$.

Lemma 6 Let \mathbf{T} be a tensor of order N and let \mathbf{M} be a matrix. Then,

$$\|\mathbf{T} \times_1 \mathbf{M}\|_2 \leq \|\mathbf{T}\|_2 \|\mathbf{M}\|_2 \quad (37)$$

Proof

$$\begin{aligned}
\|\mathbf{T} \times_1 \mathbf{M}\|_2 &= \sup_{\mathbf{v}_m, \mathbf{v}_2, \dots, \mathbf{v}_N} \sum_{i_1, \dots, i_N, m} \mathbf{T}(i_1, i_2, \dots, i_N) \mathbf{M}(i_1, m) \mathbf{v}_m(m) \mathbf{v}_{i_2}(i_2) \mathbf{v}_{i_3}(i_3) \dots \mathbf{v}_{i_n}(i_n) \\
&= \sup_{\mathbf{v}_m, \mathbf{v}_2, \dots, \mathbf{v}_N} \sum_{i_2, \dots, i_N} \sum_{i_1} \mathbf{T}(i_1, i_2, \dots, i_N) \sum_m \mathbf{M}(i_1, m) \mathbf{v}_m(m) \\
&\leq \sup_{\mathbf{v}_m, \mathbf{v}_2, \dots, \mathbf{v}_N} \left(\sup_{i_1} \left\| \sum_m \mathbf{M}(i_1, m) \mathbf{v}_m(m) \right\|_2 \right) \times \\
&\quad \sum_{i_1, \dots, i_N} \mathbf{T}(i_1, i_2, \dots, i_N) \frac{1}{\|\sum_m \mathbf{M}(i_1, m) \mathbf{v}_m(m)\|_2} \left(\sum_m \mathbf{M}(i_1, m) \mathbf{v}_m(m) \right) \mathbf{v}_{i_2}(i_2) \mathbf{v}_{i_3}(i_3) \dots \mathbf{v}_{i_n}(i_n) \\
&\leq \|\mathbf{M}\|_2 \|\mathbf{T}\|_2
\end{aligned}$$

■

Lemma 7 Let \mathbf{T} be a tensor of order N . Then,

$$\|\mathbf{T}\|_2 \leq \|\mathbf{T} \times_1 \mathbf{v}_1, \dots, \times_{n-1} \mathbf{v}_{n-1}\|_F \leq \|\mathbf{T} \times_1 \mathbf{v}_1, \dots, \times_{n-2} \mathbf{v}_{n-2}\|_F \leq \dots \leq \|\mathbf{T}\|_F \quad (38)$$

Proof It suffices to show that $\sup_{\mathbf{v}, s.t. \|\mathbf{v}\| \leq 1} \|\mathbf{T} \times_1 \mathbf{v}\|_F \leq \|\mathbf{T}\|_F$. By submultiplicativity of the Frobenius norm: $\|\mathbf{T} \times_1 \mathbf{v}\|_F \leq \|\mathbf{T}\|_F \|\mathbf{v}\|_F \leq \|\mathbf{T}\|_F$, since $\|\mathbf{v}\|_F = \|\mathbf{v}\|_2 \leq 1$. ■

Lemma 8 Let \mathbf{T} be a tensor of order N , where each mode is of dimension k . Then,

$$\|\mathbf{T}\|_1^{\sigma} \leq k^N \|\mathbf{T}\|_2 \quad (39)$$

For any σ .

Proof We simply prove this for $\sigma = \emptyset$ (which corresponds to elementwise one norm) since $\|\mathbf{T}\|_1^{\sigma_1} \leq \|\mathbf{T}\|_1^{\sigma_2}$ if $\sigma_2 \subseteq \sigma_1$. Note that $\|\mathbf{T}\|_1 \leq k^N \max(|\mathbf{T}|)$ (where $\max(\mathbf{T})$ is the maximum element of $|\mathbf{T}|$). Similarly, $\max(|\mathbf{T}|) \leq \|\mathbf{T}\|_2$ which implies that $\|\mathbf{T}\|_1 \leq k^N \|\mathbf{T}\|_2$. ■

References

- [1] R.A. Horn and C.R. Johnson. *Matrix analysis*. Cambridge Univ Pr, 1990.
- [2] D. Hsu, S. Kakade, and T. Zhang. A spectral algorithm for learning hidden Markov models. In *Proc. Annual Conf. Computational Learning Theory*, 2009.
- [3] 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.
- [4] GW Stewart and J. Sun. *Matrix Perturbation Theory*. Academic Press, 1990.