

---

# Supplementary Material

## Kernel Embeddings of Latent Tree Graphical Models

Le Song, Ankur P. Parikh, Eric P. Xing

---

Section A of the supplemental contains a brief introduction to Tensor Algebra while Section B contains details on the Spectral Algorithm.

### A Tensor Algebra

Here we give a brief introduction to tensor algebra (for more details, see [2]). 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  $\mathcal{T}(i_1, i_2, \dots, i_{n-1}, :, i_{n+1}, \dots, i_N)$ .

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  $\sum_{i_n=1}^{I_n} \mathcal{T}(i_1, \dots, i_n, \dots, i_N) \mathbf{A}(j, i_n)$ . For example, if  $\mathbf{A}$  and  $\mathbf{B}$  are matrices, then  $\mathbf{A} \times_1 \mathbf{B} = \mathbf{B}\mathbf{A}$  and  $\mathbf{A} \times_2 \mathbf{B}^\top = \mathbf{A}\mathbf{B}$ . 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  $N$ . Note that in general  $\mathcal{T} \bar{\times}_n \mathbf{a} = \text{squeeze}(\mathcal{T} \times_n \mathbf{a}^\top)$ .

### B Derivation of Spectral Algorithm

In this section, we provide a more detailed derivation of the spectral algorithm for (transformed) parameter learning. For simplicity of explanation, we will focus on latent tree structure where each internal node has exactly 3

neighbors. We can reroot the tree and redirect all the edges away from the root. For a variable  $X_s$ , we use  $\alpha_s$  to denote its sibling,  $\pi_s$  to denote its parent,  $\iota_s$  to denote its left child and  $\rho_s$  to denote its right child; the root node will have 3 children, we use  $\omega_s$  to denote the extra child. All the observed variables are leaves in the tree, and we will use  $\iota_s^*$ ,  $\rho_s^*$ ,  $\pi_s^*$  to denote an observed variable which is found by tracing in the direction from node  $s$  to its left child  $\iota_s$ , right child  $\rho_s$ , and its parent  $\pi_s$  respectively.  $s^*$  denotes any observed variable in the subtree rooted at node  $s$ .

Recall the transformed messages:

\* At leaf nodes,  $\tilde{m}_s = T_s^\top C_{s|\pi_s}^\top \phi(x_s) = (C_{s|\pi_s} \times_2 T_s^\top) \bar{\times}_1 \phi(x_s)$

\*\* At internal nodes,  $\tilde{m}_s = (C_{s^2|\pi_s} \times_1 T_{\iota_s}^{-1} \times_2 T_{\rho_s}^{-1} \times_3 T_s^\top) \bar{\times}_2 \tilde{m}_{\rho_s} \bar{\times}_1 \tilde{m}_{\iota_s}$

\*\*\* At the root,  $b_r = (C_{r^3} \times_1 T_{\iota_r}^{-1} \times_2 T_{\rho_r}^{-1} \times_3 T_{\omega_r}^{-1}) \bar{\times}_3 \tilde{m}_{\omega_s} \bar{\times}_2 \tilde{m}_{\rho_s} \bar{\times}_1 \tilde{m}_{\iota_s}$

Let  $\tilde{C}_{s^2|\pi_s} := C_{s^2|\pi_s} \times_1 T_{\iota_s}^{-1} \times_2 T_{\rho_s}^{-1} \times_3 T_s^\top$  and  $\tilde{C}_{r^3} := C_{r^3} \times_1 T_{\iota_r}^{-1} \times_2 T_{\rho_r}^{-1} \times_3 T_{\omega_r}^{-1}$ . We set  $T_s = (U_s^\top C_{s^*|\pi_s})^{-1}$ .  $U_s$  is chosen <sup>1</sup> to be the top  $d$  right singular vectors of  $C_{\pi_s^* s^*}$ , and therefore one can take the one-sided inverse of  $(C_{\pi_s^* s^*} U_s)$  assuming all latent variables have dimension  $d$ . For internal nodes we set  $s^*$  in  $(C_{\pi_s^* s^*} U_s)$  to  $\iota_s^*$  while for leaves we set  $s^*$  to  $s$ . We have the following observable representation that we derive in the following subsections:

\* At leaf nodes,  $\tilde{m}_s = (C_{\pi_s^* s} U_s)^\dagger C_{\pi_s^* s} \phi(x_s)$ .

\*\* At internal nodes,  $\tilde{C}_{s^2|\pi_s} = C_{\iota_s^* \rho_s^* \pi_s^*} \times_1 U_{\iota_s}^\top \times_2 U_{\rho_s}^\top \times_3 (C_{\pi_s^* \iota_s} U_s)^\dagger$ .

\*\*\* At the root,  $\tilde{C}_{r^3} = C_{\iota_r^* \rho_r^* \omega_r^*} \times_1 U_{\iota_r}^\top \times_2 U_{\rho_r}^\top \times_3 U_{\omega_r}^\top$

## B.1 Root

Recall that

$$\tilde{C}_{r^3} = C_{r^3} \times_1 T_{\iota_r}^{-1} \times_2 T_{\rho_r}^{-1} \times_3 T_{\omega_r}^{-1} \quad (5)$$

$$= C_{r^3} \times_1 U_{\iota_r}^\top C_{\iota_r^*|r} \times_2 U_{\rho_r}^\top C_{\rho_r^*|r} \times_3 U_{\omega_r}^\top C_{\omega_r^*|r} \quad (6)$$

$$= C_{r^3} \times_1 C_{\iota_r^*|r} \times_2 C_{\rho_r^*|r} \times_3 C_{\omega_r^*|r} \times_1 U_{\iota_r}^\top \times_2 U_{\rho_r}^\top \times_3 U_{\omega_r}^\top \quad (7)$$

where  $T_s^{-1} = U_s^\top C_{s^*|\pi_s}$ .

We first prove that  $C_{r^3} \times_1 C_{\iota_r^*|r} \times_2 C_{\rho_r^*|r} \times_3 C_{\omega_r^*|r} = C_{\iota_r^* \rho_r^* \omega_r^*}$ : Consider any  $f, g, h \in \mathcal{F}$ . Then,

$$C_{r^3} \times_1 C_{\iota_r^*|r} \times_2 C_{\rho_r^*|r} \times_3 C_{\omega_r^*|r} \bar{\times}_3 h \bar{\times}_2 g \bar{\times}_1 f \quad (8)$$

$$= \langle f \otimes g \otimes h, C_{r^3} \times_1 C_{\iota_r^*|r} \times_2 C_{\rho_r^*|r} \times_3 C_{\omega_r^*|r} \rangle \quad (9)$$

$$= \mathbb{E}_{X_r} \left[ \langle C_{\iota_r^*|r}^\top f, \phi(X_r) \rangle \langle C_{\rho_r^*|r}^\top g, \phi(X_r) \rangle \langle C_{\omega_r^*|r}^\top h, \phi(X_r) \rangle \right] \quad (10)$$

$$= \mathbb{E}_{X_r} \left[ \langle f, C_{\iota_r^*|r} \phi(X_r) \rangle \langle g, C_{\rho_r^*|r} \phi(X_r) \rangle \langle h, C_{\omega_r^*|r} \phi(X_r) \rangle \right] \quad (11)$$

$$= \mathbb{E}_{X_r} \left[ \mathbb{E}_{X_{\rho_r^*} | X_r} [f(X_{\iota_r^*})] \mathbb{E}_{X_{\rho_r^*} | X_r} [g(X_{\rho_r^*})] \mathbb{E}_{X_{\omega_r^*} | X_r} [h(X_{\omega_r^*})] \right] \quad (12)$$

$$= \mathbb{E}_{X_{\iota_r^*}, X_{\rho_r^*}, X_{\omega_r^*}} [f(X_{\iota_r^*}) g(X_{\rho_r^*}) h(X_{\omega_r^*})] \quad (13)$$

$$= \langle f \otimes g \otimes h, \mathbb{E}_{X_{\iota_r^*}, X_{\rho_r^*}, X_{\omega_r^*}} [\phi(X_{\iota_r^*}) \otimes \phi(X_{\rho_r^*}) \otimes \phi(X_{\omega_r^*})] \rangle \quad (14)$$

$$= C_{\iota_r^* \rho_r^* \omega_r^*} \bar{\times}_3 h \bar{\times}_2 g \bar{\times}_1 f \quad (15)$$

Combining this result with Eq. 7 gives,

$$\tilde{C}_{r^3} = C_{r^3} \times_1 T_{\iota_r}^{-1} \times_2 T_{\rho_r}^{-1} \times_3 T_{\omega_r}^{-1} = C_{\iota_r^* \rho_r^* \omega_r^*} \times_1 U_{\iota_r}^\top \times_2 U_{\rho_r}^\top \times_3 U_{\omega_r}^\top \quad (16)$$

<sup>1</sup>This is not the only valid choice of  $U_s$  but will generally result in better performance. See [1, 3] for more details.

## B.2 Leaf

Recall that  $\tilde{m}_s = T_s^\top C_{s|\pi_s}^\top \phi(x_s)$  and  $T_s = (U_s^\top C_{s^*|\pi_s})^{-1}$ . However since  $s$  is a leaf we can set  $s^* = s$ . Consider expanding the related quantity  $\tilde{m}_s^\top (U_s^\top C_{s\pi_s^*})$ :

$$\tilde{m}_s^\top (U_s^\top C_{s\pi_s^*}) = \phi^\top(x_s) C_{s|\pi_s} (U_s^\top C_{s|\pi_s})^{-1} (U_s^\top C_{s\pi_s^*}) \quad (17)$$

$$= \phi^\top(x_s) C_{s|\pi_s} (U_s^\top C_{s|\pi_s})^{-1} (U_s^\top C_{s|\pi_s} C_{\pi_s^2} C_{\pi_s^*|\pi_s}^\top) \quad (18)$$

$$= \phi^\top(x_s) C_{s|\pi_s} (U_s^\top C_{s|\pi_s})^{-1} (U_s^\top C_{s|\pi_s}) (C_{\pi_s^2} C_{\pi_s^*|\pi_s}^\top) \quad (19)$$

$$= \phi^\top(x_s) C_{s|\pi_s} C_{\pi_s^2} C_{\pi_s^*|\pi_s}^\top \quad (20)$$

$$= \phi(x_s)^\top C_{s\pi_s^*} \quad (21)$$

where we have used the fact that  $C_{s|\pi_s} C_{\pi_s^2} C_{\pi_s^*|\pi_s}^\top = C_{s\pi_s^*}$  (which is proved using the same technique as used in Section B.1).

This implies that  $\tilde{m}_s = (C_{\pi_s^*} U_s)^\dagger C_{\pi_s^*} \phi(x_s) = C_{s\pi_s^*} (U_s^\top C_{s\pi_s^*})^\dagger \bar{\times}_1 \phi(x_s)$ . We choose  $U_s$  to be the top  $d$  right singular vectors of  $C_{\pi_s^*}$ , and therefore the one-sided inverse exists (since all latent variables are assumed to have dimension  $d$ ).

## B.3 Intermediate Node

Recall that  $T_s = (U_s^\top C_{s^*|\pi_s})^{-1}$  and  $\tilde{C}_{s^2|\pi_s} = C_{s^2|\pi_s} \times_1 T_{\iota_s}^{-1} \times_2 T_{\rho_s}^{-1} \times_3 T_s^\top$ . Thus,

$$\tilde{C}_{s^2|\pi_s} = C_{s^2|\pi_s} \times_1 U_{\iota_s}^\top C_{\iota_s^*|s} \times_2 U_{\rho_s}^\top C_{\rho_s^*|s} \times_3 (C_{s|\pi_s}^\top U_s)^{-1} \quad (22)$$

Consider expanding the quantity  $\tilde{C}_{s^2|\pi_s} \times_3 (C_{\pi_s^*} U_s)$ :

$$\tilde{C}_{s^2|\pi_s} \times_3 (C_{\pi_s^*} U_s) = C_{s^2|\pi_s} \times_1 U_{\iota_s}^\top C_{\iota_s^*|s} \times_2 U_{\rho_s}^\top C_{\rho_s^*|s} \times_3 (C_{s|\pi_s}^\top U_s)^{-1} \times_3 (C_{\pi_s^*} U_s) \quad (23)$$

$$= C_{s^2|\pi_s} \times_1 C_{\iota_s^*|s} \times_2 C_{\rho_s^*|s} \times_3 (C_{\pi_s^*} U_s) (C_{s|\pi_s}^\top U_s)^{-1} \times_1 U_{\iota_s}^\top \times_2 U_{\rho_s}^\top \quad (24)$$

$$= C_{s^2|\pi_s} \times_1 C_{\iota_s^*|s} \times_2 C_{\rho_s^*|s} \times_3 (C_{\pi_s^*|\pi_s} C_{\pi_s^2} C_{s|\pi_s}^\top U_s) (C_{s|\pi_s}^\top U_s)^{-1} \times_1 U_{\iota_s}^\top \times_2 U_{\rho_s}^\top \quad (25)$$

$$= C_{s^2|\pi_s} \times_1 C_{\iota_s^*|s} \times_2 C_{\rho_s^*|s} \times_3 (C_{\pi_s^*|\pi_s} C_{\pi_s^2}) (C_{s|\pi_s}^\top U_s) (C_{s|\pi_s}^\top U_s)^{-1} \times_1 U_{\iota_s}^\top \times_2 U_{\rho_s}^\top \quad (26)$$

$$= C_{s^2|\pi_s} \times_1 C_{\iota_s^*|s} \times_2 C_{\rho_s^*|s} \times_3 (C_{\pi_s^*|\pi_s} C_{\pi_s^2}) \times_1 U_{\iota_s}^\top \times_2 U_{\rho_s}^\top \quad (27)$$

$$= C_{\iota_s^*, \rho_s^*, \pi_s^*} \times_1 U_{\iota_s}^\top \times_2 U_{\rho_s}^\top \quad (28)$$

where in the last line we have claimed that  $C_{\iota_s^*, \rho_s^*, \pi_s^*} = C_{s^2|\pi_s} \times_1 C_{\iota_s^*|s} \times_2 C_{\rho_s^*|s} \times_3 C_{\pi_s^*|\pi_s} C_{\pi_s^2}$ . To prove this assertion, first consider the  $C_{s^2|\pi_s} \times_1 C_{\iota_s^*|s} \times_2 C_{\rho_s^*|s}$  part. For any  $f, g \in \mathcal{F}$ :

$$\langle f \otimes g, C_{s^2|\pi_s} \times_1 C_{\iota_s^*|s} \times_2 C_{\rho_s^*|s} \bar{\times}_3 \phi(x_{\pi_s}) \rangle = \langle (C_{\iota_s^*|s}^\top f) \otimes (C_{\rho_s^*|s}^\top g), C_{s^2|\pi_s} \bar{\times}_3 \phi(x_{\pi_s}) \rangle \quad (29)$$

$$= \langle (C_{\iota_s^*|s}^\top f) \otimes (C_{\rho_s^*|s}^\top g), \mathbb{E}_{X_s|x_{\pi_s}} [\phi(X_s) \otimes \phi(X_s)] \rangle \quad (30)$$

$$= \mathbb{E}_{X_s|x_{\pi_s}} \left[ \langle (C_{\iota_s^*|s}^\top f) \otimes (C_{\rho_s^*|s}^\top g), \phi(X_s) \otimes \phi(X_s) \rangle \right] \quad (31)$$

$$= \mathbb{E}_{X_s|x_{\pi_s}} \left[ \langle f, C_{\iota_s^*|s} \phi(X_s) \rangle \langle g, C_{\rho_s^*|s} \phi(X_s) \rangle \right] \quad (32)$$

$$= \mathbb{E}_{X_s|x_{\pi_s}} \left[ \mathbb{E}_{X_{\iota_s^*}|X_s} [f(X_{\iota_s^*})] \mathbb{E}_{X_{\rho_s^*}|X_s} [g(X_{\rho_s^*})] \right] \quad (33)$$

$$= \mathbb{E}_{\iota_s^*, \rho_s^*|x_{\pi_s}} [f(X_{\iota_s^*}) g(X_{\rho_s^*})] \quad (34)$$

$$= \langle f \otimes g, C_{\iota_s^*, \rho_s^*|\pi_s} \bar{\times}_3 \phi(x_{\pi_s}) \rangle \quad (35)$$

Thus,  $\mathcal{C}_{\iota_s^* \rho_s^* | \pi_s} = \mathcal{C}_{s^2 | \pi_s} \times_1 \mathcal{C}_{\iota_s^* | s} \times_2 \mathcal{C}_{\rho_s^* | s}$ . We can then conclude (using a similar derivation to that in Section B.1) that  $\mathcal{C}_{\iota_s^*, \rho_s^*, \pi_s^*} = \mathcal{C}_{\iota_s^* \rho_s^* | \pi_s} \times_3 \mathcal{C}_{\pi_s^* | \pi_s} \mathcal{C}_{\pi_s^2}$ . Thus,

$$\mathcal{C}_{\iota_s^*, \rho_s^*, \pi_s^*} = \mathcal{C}_{s^2 | \pi_s} \times_1 \mathcal{C}_{\iota_s^* | s} \times_2 \mathcal{C}_{\rho_s^* | s} \times_3 \mathcal{C}_{\pi_s^* | \pi_s} \mathcal{C}_{\pi_s^2} \quad (36)$$

Now, returning to Eq. 28 we get that

$$\tilde{\mathcal{C}}_{s^2 | \pi_s} = \mathcal{C}_{\iota_s^*, \rho_s^*, \pi_s^*} \times_1 U_{\iota_s^*}^\top \times_2 U_{\rho_s^*}^\top \times_3 (\mathcal{C}_{\pi_s^* s^*} U_s)^\dagger \quad (37)$$

where one valid choice for  $s^*$  is  $\iota_s^*$ .  $U_s$  is chosen to be the top  $d$  right singular vectors of  $\mathcal{C}_{\pi_s^* \iota_s^*}$ , and therefore one can take a one-sided inverse of  $(\mathcal{C}_{\pi_s^* s^*} U_s)$  (assuming all latent variables have dimension  $d$ ).

## References

- [1] D. Hsu, S. Kakade, and T. Zhang. A spectral algorithm for learning hidden markov models. In *COLT*, 2009.
- [2] Tamara Kolda and Brett Bader. Tensor decompositions and applications. *SIAM Review*, 51(3):455–500, 2009.
- [3] A. Parikh, L. Song, and E. Xing. A spectral algorithm for latent tree graphical models. In *ICML*, 2011.