

Spectral Learning of Latent Variable Models and its Interpretation as an Optimization Problem

– a CSML pizza talk –

Gabi Teodoru

joint work with Maneesh Sahani

Special thanks (in no particular order):

Jeff Beck, Lars Büsing, Bharath Sriperumbudur,
Byron Boots, YeeWhye Teh and Arthur Gretton

Outline

1. Spectral learning
 - Previous work
 - Optimization interpretation
 - Re-deriving some previous algorithms
2. Spectral learning for Latent Markov Lattices
 - The model
 - Learning: Three tricks
3. Future work

What is Spectral Learning?

- Spectral Learning is a learning algorithm for latent variable graphical models
 - More details coming up in a few slides
- Pros:
 - Consistent estimator of parameters; computable
 - Quick, one-shot learning
- Cons:
 - Statistically inefficient
 - Estimates a function of the parameters; additional steps are needed to recover the parameters
 - Numerical instability

Previous work

- Spectral algorithms:
 - Spectral learning for HMMs (Hsu et al, 2009)
 - Reduced-rank HMMs (Siddiqi et al, 2010)
 - Hilbert space embeddings of HMMs (Song et al, 2010)
 - Spectral learning of latent tree graphical models (Parikh et al, 2011)
- Related algorithms and ideas:
 - Subspace identification (SSID) for Linear Dynamical Systems (LDS) [aka Kalman Filter]
 - Observable Operator Models (OOM) / Predictive state representations (PSR)

Optimization interpretation: Motivation

- What previous work showed:
 - In the limit of infinite data, the algorithm identifies the true parameters
 - Finite sample bounds relate the finite sample size to prediction errors
- What previous work authors' claimed in talks:
 - Algorithm is convex, and therefore does not get stuck in local optima (unlike EM)
- What I wanted to see:
 - The cost function
 - (... they didn't have one)

Optimization interpretation

There is a common optimization interpretation for all previously mentioned spectral algorithms and SSID.

Pros of optimization interpretation

- Allows improving the cost functions
 - Adding regularizers
 - Using the Generalized Method of Moments
- Provides a way to use the data more efficiently
- Generalizes the method to other models

Constructing the optimization interpretation

- Step 1: Given a latent variable model parametrized by θ , pick a transformation of these parameters $\phi = \phi(\theta)$

- Step 2: Choose two sets of model statistics

$$\mathbb{E}_{\phi} [f_1(\mathbf{X})] \quad \mathbb{E}_{\phi} [f_2(\mathbf{X})]$$

- Step 3: Choose two functions

$$g_1(\cdot, \phi) \quad g_2(\cdot, \phi)$$

such that

$$g_1\left(\mathbb{E}_{\phi} [f_1(\mathbf{X})], \phi\right) = g_2\left(\mathbb{E}_{\phi} [f_1(\mathbf{X})], \phi\right) \forall \phi \in \Phi$$

Examples

$$g_1(\mathbb{E}_\phi[f_1(\mathbf{X})], \phi) = g_2(\mathbb{E}_\phi[f_1(\mathbf{X})], \phi) \quad \forall \phi \in \Phi$$

- LDS:** $f_1 = E[X_3 X_1^T]$ $f_2 = E[X_2 X_1^T]$
 $g_1(f_1, \phi) = f_1$ $g_2(f_2, \phi) = \Phi f_2$
- HMM:** $f_1 = P(X_{3..1})$ $f_2 = P(X_{2..1})$
 $g_1(f_1, \phi) = f_1$ $g_2(f_2, \phi) = \Phi f_2$
- SCFG:** $f_1 = P(X_{1..3})$ $f_2 = [P(X_{1..2}), P(X_{2..3})]$
 $g_1(f_1, \phi) = f_1$ $g_2(f_2, \phi) = \sum_k [f_2]_{i,k} \phi_k \delta_k$

The cost function

- Since

$$g_1\left(\mathbb{E}_\phi\left[f_1(\mathbf{X})\right],\phi\right) = g_2\left(\mathbb{E}_\phi\left[f_1(\mathbf{X})\right],\phi\right) \forall \phi \in \Phi$$

if we have finite samples from the model,

$$g_1\left(f_1(\mathbf{x}),\phi\right) \approx g_2\left(f_2(\mathbf{x}),\phi\right)$$

the dumbest thing we could possibly do is:

$$\min_{\phi} d\left(g_1\left(f_1(\mathbf{x}),\phi\right), g_2\left(f_2(\mathbf{x}),\phi\right)\right)$$

If the functions g_1, g_2 are linear in the 2nd parameter, and the distance is quadratic, then we have a convex optimization!

Adding bells and whistles

- Original cost function:

$$C(\mathbf{x}, \phi) = d\left(g_1\left(f_1(\mathbf{x}), \phi\right), g_2\left(f_2(\mathbf{x}), \phi\right)\right)$$

- With regularization:

$$C(\mathbf{x}, \phi) = d(g_1, g_2) + \Omega(\phi)$$

- Generalized Method of Moments:

$$C(\mathbf{x}, \phi) = (g_1 - g_2)^T \mathbf{W} (g_1 - g_2)$$

GMM iteratively re-estimates

$$\mathbf{W} = (g_1 - g_2)(g_1 - g_2)^T$$

Example 1: SSID for LDS

- Hidden RVs $H_t \in \mathbb{R}^m$, observables $X_t \in \mathbb{R}^n$, $m \leq n$
- Distribution: $H_t, X_t \sim \text{Normal}$
 - $H_t = \mathbf{T}H_{t-1} + \text{noise} \leftrightarrow E[H_t] = \mathbf{T}E[H_{t-1}]$
 - $X_t = \mathbf{O}H_t + \text{noise} \leftrightarrow E[X_t] = \mathbf{O}E[H_t]$
 - Assume $\mathbf{T}_{m \times m}$ and $\mathbf{O}_{n \times m}$ are full rank
- Special relationship and spectral coefficients:
 - $E[X_{t+1}] = \mathbf{O}\mathbf{T}E[H_t]$
 - $E[H_t] = \mathbf{O}^\dagger E[X_t]$
 - $E[X_{t+1}] = \mathbf{O}\mathbf{T}\mathbf{O}^\dagger E[X_t]$
 - $E[X_{t+1}] = \mathbf{\Phi}E[X_t] \leftrightarrow \mathbf{\Phi} = \mathbf{O}\mathbf{T}\mathbf{O}^\dagger$
- Cost function: $\|E[X_{t+1}] - \mathbf{\Phi}E[X_t]\|_2^2$

Example 1: SSID for LDS

- Minimization: $\Phi^* = \arg \min_{\Phi} \|E[X_{t+1}] - \Phi E[X_t]\|_F$
- Cost function is underdetermined!
 - Solution 1: $\Phi_t^* = \arg \min_{\Phi_t} \|E[X_{t+1}X_{t-1}^T] - \Phi_t E[X_tX_{t-1}^T]\|_F$
 - Solution 2: $\Phi^* = \arg \min_{\Phi} \sum_t \|E[X_{t+1}] - \Phi E[X_t]\|_F$
 - May be expressed also by horizontal concatenation:
$$\Phi^* = \arg \min_{\Phi} d\left(\begin{bmatrix} E[X_t] & E[X_{t+1}] & E[X_{t+2}] & \dots \end{bmatrix}, \Phi \begin{bmatrix} E[X_{t-1}] & E[X_t] & E[X_{t+1}] & \dots \end{bmatrix}\right)$$

Example 1: SSID for LDS

- How are the LDS spectral parameters useful?
 - Parameters are identifiable up to a transformation:

$$E[H_t] = \mathbf{T}E[H_{t-1}]$$

$$E[X_t] = \mathbf{O}E[H_t]$$

Example 1: SSID for LDS

- How are the LDS spectral parameters useful?
 - Parameters are identifiable up to a transformation:

$$E[\mathbf{S}H_t] = \mathbf{S}\mathbf{T}\mathbf{S}^{-1}E[\mathbf{S}H_{t-1}]$$

$$E[X_t] = \mathbf{O}\mathbf{S}^{-1}E[\mathbf{S}H_t]$$

Example 1: SSID for LDS

- How are the LDS spectral parameters useful?
 - Parameters are identifiable up to a transformation:

$$E[\mathbf{S}H_t] = \mathbf{S}\mathbf{T}\mathbf{S}^\dagger E[\mathbf{S}H_{t-1}]$$

$$E[X_t] = \mathbf{O}\mathbf{S}^\dagger E[\mathbf{S}H_t]$$

Example 1: SSID for LDS

- How are the LDS spectral parameters useful?
 - Parameters are identifiable up to a transformation:

$$\left. \begin{aligned} E[\mathbf{S}H_t] &= \mathbf{S}\mathbf{T}\mathbf{S}^\dagger E[\mathbf{S}H_{t-1}] \\ E[X_t] &= \mathbf{O}\mathbf{S}^\dagger E[\mathbf{S}H_t] \end{aligned} \right\} \rightarrow \begin{cases} E[H'_t] = \mathbf{S}\mathbf{T}\mathbf{S}^\dagger E[H'_{t-1}] \\ E[X_t] = \mathbf{O}\mathbf{S}^\dagger E[H'_t] \end{cases}$$

- So what were the spectral coefficients we learned?

$$\begin{aligned} \Phi &= \mathbf{O}\mathbf{T}\mathbf{O}^\dagger & \Phi^\dagger &= \mathbf{O}\mathbf{T}^{-1}\mathbf{O}^\dagger \\ \Phi &= \mathbf{O}\mathbf{T}(\mathbf{T}\mathbf{T}^{-1})\mathbf{O}^\dagger & \Phi^\dagger &= \mathbf{O}(\mathbf{O}\mathbf{T})^\dagger \\ \Phi &= (\mathbf{O}\mathbf{T})\mathbf{T}(\mathbf{O}\mathbf{T})^\dagger \end{aligned}$$

- The spectral coefficients matrix and its pseudoinverse provide the corresponding transformations of the model parameters

Example 1: SSID for LDS

- Hidden RVs $H_t \in \mathbb{R}^m$, observables $X_t \in \mathbb{R}^n$, $m \leq n$
- Distribution: $H_t, X_t \sim \text{Normal}$
 - $H_t = \mathbf{T}H_{t-1} + \text{noise} \leftrightarrow E[H_t] = \mathbf{T}E[H_{t-1}]$
 - $X_t = \mathbf{O}H_t + \text{noise} \leftrightarrow E[X_t] = \mathbf{O}E[H_t]$
 - Assume $\mathbf{T}_{m \times m}$ and $\mathbf{O}_{n \times m}$ are full rank
- Special relationship and spectral coefficients:
 - $E[X_{t+1}] = \mathbf{O}\mathbf{T}E[H_t]$
 - $E[H_t] = \mathbf{O}^\dagger E[X_t]$
 - $E[X_{t+1}] = \mathbf{O}\mathbf{T}\mathbf{O}^\dagger E[X_t]$
 - $E[X_{t+1}] = \mathbf{\Phi}E[X_t] \leftrightarrow \mathbf{\Phi} = \mathbf{O}\mathbf{T}\mathbf{O}^\dagger$
- Cost function: $\|E[X_{t+1}] - \mathbf{\Phi}E[X_t]\|_F$

Adding the SVD

- Hidden RVs $H_t \in \mathbb{R}^m$, observables $X_t \in \mathbb{R}^n$, $m \leq n$
- Distribution: $H_t, X_t \sim \text{Normal}$
 - $H_t = \mathbf{T}H_{t-1} + \text{noise} \Leftrightarrow E[H_t] = \mathbf{T}E[H_{t-1}]$
 - $X_t = \mathbf{O}H_t + \text{noise} \Leftrightarrow E[X_t] = \mathbf{O}E[H_t]$
 - Assume $\mathbf{T}_{m \times m}$ and $\mathbf{O}_{n \times m}$ are full rank
- $\mathbf{\Phi} = \mathbf{O}\mathbf{T}\mathbf{O}^\dagger$ implies that $\mathbf{\Phi}$ has rank m
- Make this a constraint of the optimization:

$$\begin{cases} \min_{\mathbf{\Phi}} d(E[X_{t+1}], \mathbf{\Phi}E[X_t]) \\ \text{s.t. rank}(\mathbf{\Phi}) \leq m \end{cases}$$

Adding the SVD

$$\begin{cases} \min_{\Phi} \|E[X_{t+1}] - \Phi E[X_t]\|_F \\ \text{s.t. rank}(\Phi) \leq m \end{cases}$$

is equivalent to:

$$\min_{\Phi} \|\mathbf{U}^T E[X_{t+1}] - \Phi \mathbf{U}^T E[X_t]\|_F$$

where \mathbf{U} is the matrix of the first m -left singular vectors of

$$E[X_t X_{t-1}^T] \quad \text{or} \quad [E[X_t] \ E[X_{t+1}] \ E[X_{t+2}] \ \dots]$$

Example 2: Spectral Learning of HMMs

- Hidden $H_t \in \{1..m\}$, observables $X_t \in \{1..n\}$, $m \leq n$
- Parameters: \mathbf{T} , \mathbf{O}
 - $\mathbf{p}_{H_{t+1}} = \mathbf{T}\mathbf{p}_{H_t}$
 - $\mathbf{p}_{X_t} = \mathbf{O}\mathbf{p}_{H_t}$
 - Assume $\mathbf{T}_{m \times m}$ and $\mathbf{O}_{n \times m}$ are full rank
- Special relationship and spectral coefficients:
 - $\mathbf{p}_{X_{t+1}} = \mathbf{O}\mathbf{T}\mathbf{p}_{H_t}$
 - $\mathbf{p}_{H_t} = \mathbf{O}^\dagger \mathbf{p}_{X_t}$
 - $\mathbf{p}_{X_{t+1}} = \mathbf{O}\mathbf{T}\mathbf{p}_{H_t} \rightarrow \mathbf{p}_{X_{t+1}} = \mathbf{O}\mathbf{T}\mathbf{O}^\dagger \mathbf{p}_{X_t}$
 - $\mathbf{p}_{X_{t+1}} = \mathbf{\Phi}\mathbf{p}_{X_t} \leftrightarrow \mathbf{\Phi} = \mathbf{O}\mathbf{T}\mathbf{O}^\dagger$
- Spectral coefficients are a transform of a conditional probability:
 - Useless for computing the joint distribution $p(x_{1..t})$

Example 2: Spectral Learning of HMMs

- Use probability matrices: $\left[\mathbf{P}_{X,Y} \right]_{x,y} = P(X = x, Y = y)$
- We had from the previous slide:
 - $\mathbf{p}_{H_{t+1}} = \mathbf{T}\mathbf{p}_{H_t}$
 - $\mathbf{p}_{X_t} = \mathbf{O}\mathbf{p}_{H_t}$
 - Assume $\mathbf{T}_{m \times m}$ and $\mathbf{O}_{n \times m}$ are full rank, with $m \leq n$
- Special relationship and spectral coefficients:
 - $\mathbf{p}_{X_{t+1}} = \mathbf{O}\mathbf{T}\mathbf{p}_{H_t}$
 - $\mathbf{p}_{H_t} = \mathbf{O}^\dagger \mathbf{p}_{X_t}$
 - $\mathbf{p}_{X_{t+1}} = \mathbf{\Phi}\mathbf{p}_{X_t} \iff \mathbf{\Phi} = \mathbf{O}\mathbf{T}\mathbf{O}^\dagger$

Example 2: Spectral Learning of HMMs

- Use probability matrices: $\left[\mathbf{P}_{X,Y} \right]_{x,y} = P(X = x, Y = y)$
- We had from the previous slide:
 - $\mathbf{p}_{H_{t+1}} = \mathbf{T} \mathbf{p}_{H_t}$
 - $\mathbf{p}_{X_t} = \mathbf{O} \mathbf{p}_{H_t}$
 - Assume $\mathbf{T}_{m \times m}$ and $\mathbf{O}_{n \times m}$ are full rank, with $m \leq n$
- Special relationship and spectral coefficients:
 - $\left[\mathbf{P}_{X_{t+1}, X_t} \right]_{\cdot, x_t} = \mathbf{O} \mathbf{T} \text{diag}([\mathbf{O}]_{x_t, \cdot}) \mathbf{p}_{H_t}$
 - $\mathbf{p}_{H_t} = \mathbf{O}^\dagger \mathbf{p}_{X_t}$
 - $\left[\mathbf{P}_{X_{t+1}, X_t} \right]_{\cdot, x_t} = \mathbf{O} \mathbf{T} \text{diag}([\mathbf{O}]_{x_t, \cdot}) \mathbf{O}^\dagger \mathbf{p}_{X_t}$
 - $\left[\mathbf{P}_{X_{t+1}, X_t} \right]_{\cdot, x_t} = \Phi_{x_t} \mathbf{p}_{X_t} \iff \Phi_{x_t} = \mathbf{O} \mathbf{T} \text{diag}[\mathbf{O}_{x_t, \cdot}] \mathbf{O}^\dagger$
- This can be used to compute the joint, since:

$$\left[\mathbf{P}_{X_{t+1}, X_t, X_{\text{past}}} \right]_{\cdot, x_t, \mathbf{x}_{\text{past}}} = \Phi_{x_t} \left[\mathbf{P}_{X_t, X_{\text{past}}} \right]_{\cdot, \mathbf{x}_{\text{past}}}$$

Example 2: Spectral Learning of HMMs

$$\left\{ \begin{array}{l} \min_{\Phi_{x_t}} \left\| \left[\mathbf{P}_{X_{t+1}, X_t, X_{\text{past}}} \right]_{\bullet, x_t, \mathbf{x}_{\text{past}}} - \Phi_{x_t} \left[\mathbf{P}_{X_t, X_{\text{past}}} \right]_{\bullet, \mathbf{x}_{\text{past}}} \right\|_F \\ \text{s.t. } \text{rank}(\Phi_{x_t}) \leq m \end{array} \right.$$

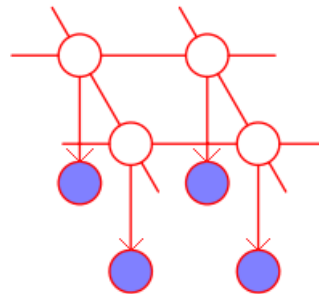
Spectral learning for other models

$$C(\mathbf{x}, \phi) = d\left(g_1\left(f_1(\mathbf{x}), \phi\right), g_2\left(f_2(\mathbf{x}), \phi\right)\right)$$

- $f_{1,2}(\mathbf{x})$ are Hilbert space embeddings
 - Hilbert space embeddings of HMMs (Song et al, 2010)
 - $f_{1,2}(\mathbf{x})$ are certain sets of joint probabilities in a tree chosen such that linear relationships hold
 - Spectral learning of latent tree graphical models (Parikh et al, '11)
- Other models on which I've tried spectral learning:
 - PCFGs (can only learn from strings of length 1..3)
 - Finite State Transducers (adds an extra layer of hidden variables; optimization becomes bilinear least squares)
 - Markov Random Fields (that's the segue...)

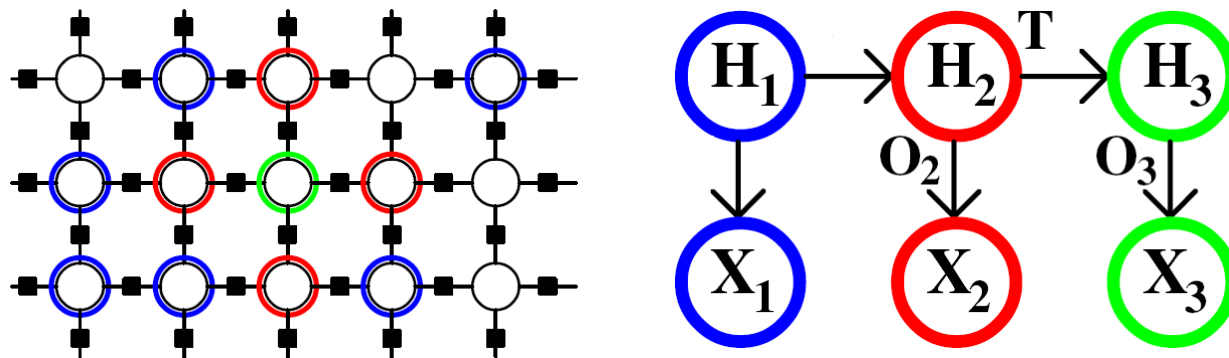
MRFs and Latent Markov Lattices

- Hidden $H_i \in \{1..m\}$, observables $X_i \in \{1..n\}$, $m \leq n$
- Distribution:
$$P(\mathbf{x}) \propto \sum_{\mathbf{h}} \left(\prod_i P(x_i | h_i) \right) \left(\prod_{(i,j) \in \text{edges}} \theta_{h_i, h_j} \right)$$
- Parameters:
 - $\theta_{ab} = \theta_{ba} \geq 0$
 - $P(X_i = a | H_i = b) = [\mathbf{O}]_{ab}$. \mathbf{O} is n -by- m and full rank
- For concreteness, we'll look at 2D square lattices:



Learning LMLs: Trick #1: HMM equivalence

- Construct a new set of variables:
 - $\mathbb{H}_3 = H_i, \mathbb{X}_3 = X_i$
 - $\mathbb{H}_2 = \text{Markov Boundary}(H_i), \mathbb{X}_2 = \text{child}(\mathbb{H}_2)$
 - $\mathbb{H}_1 = H_\sigma, \mathbb{H}_1 \cap (\mathbb{H}_2 \cup \mathbb{H}_3) = \emptyset, \mathbb{X}_1 = X_\sigma$
- This is a Hidden Markov Model



- Parameters:
 - $\mathbf{O}_3 = \mathbf{O}$
 - $[\mathbf{O}_2]_{\sigma, \tau} = P(\mathbb{X}_2 = \sigma \mid \mathbb{H}_2 = \tau), \mathbf{O}_{2^{n^4} \times m^4}$
 $= P(X_{2,1} = \sigma_1, X_{2,2} = \sigma_2, \dots \mid H_{2,1} = \tau_1, H_{2,2} = \tau_2, \dots)$
 - $[\mathbf{T}]_{\sigma, i} = P(\mathbb{H}_3 = i \mid \mathbb{H}_2 = \tau)$. Assume $\mathbf{T}_{m \times m^4}$ is full rank

Learning LMLs: Trick #2: Tensor re-slicing

- Solve:

- $\min_{\Phi_{\mathbf{x}_2}} d \left(\left[\mathbf{P}_{\mathbb{X}_3, \mathbb{X}_2, \mathbb{X}_1} \right]_{\cdot, \mathbf{x}_2, \cdot}, \Phi_{\mathbf{x}_2} \left[\mathbf{P}_{\mathbb{X}_2, \mathbb{X}_1} \right]_{\cdot, \cdot} \right) \text{ s.t. } \text{rank}(\Phi_{\mathbf{x}_2}) \leq m$

- $\Phi_{\mathbf{x}_2} = \mathbf{O}_3 \mathbf{T} \text{diag} \left(\left[\mathbf{O}_2 \right]_{\mathbf{x}_2, \cdot} \right) \mathbf{O}_2^\dagger$

- This spectral parameterization is useless!

- Solve:

- $\min_{\Phi_{\mathbf{x}_3}} d \left(\left[\mathbf{P}_{\mathbb{X}_3, \mathbb{X}_2, \mathbb{X}_1} \right]_{\mathbf{x}_3, \cdot, \cdot}, \Phi_{\mathbf{x}_3} \left[\mathbf{P}_{\mathbb{X}_2, \mathbb{X}_1} \right]_{\cdot, \cdot} \right) \text{ s.t. } \text{rank}(\Phi_{\mathbf{x}_3}) \leq m$

- $\Phi_{\mathbf{x}_3} = \mathbf{O}_2 \text{diag} \left(\left[\mathbf{O}_3 \right]_{\mathbf{x}_3, \cdot} \mathbf{T} \right) \mathbf{O}_2^\dagger$

- This is somewhat more useful...

- Simply a re-slicing of the tensor above

Learning LMLs: Trick #3: Eigendecomposition

- Warning: Hand-wavy math
- The eigenvectors of the spectral parameters:

$$\Phi_{\mathbf{x}_3} = \mathbf{O}_2 \text{diag}([\mathbf{O}_3]_{\mathbf{x}_3}, \mathbf{T}) \mathbf{O}_2^\dagger$$

are the columns of \mathbf{O}_2 , up to scaling

- due to model symmetries, the matrix $\Phi_{\mathbf{x}_3}$ is guaranteed to have repeated eigenvalues
- also due to model symmetries: the eigenvectors we can recuperate are sufficient to identify the model

Computing the MRF model parameters

- Compute $\mathbf{O} = \mathbf{O}_3$ using its relationships to \mathbf{O}_2 :
 - $[\mathbf{O}_2]_{\sigma,\tau} = \prod_i [\mathbf{O}]_{\sigma_i,\tau_i}$ or $[\mathbf{O}_2]_{[i,i,i,i],[j,j,j,j]} = ([\mathbf{O}]_{i,j})^4$
 - $\sum_{\sigma: \sigma_i=a} [\mathbf{O}_2]_{\sigma,\tau} = [\mathbf{O}]_{a,\tau_i}$
- May be computed by multiple methods:
 - minimize least squares in log space (linear)
 - exponentiation
 - marginalization
 - minimize least squares
- Compute then the remaining unknown eigenvectors of \mathbf{O}_2 from the now known \mathbf{O}

Computing the MRF model parameters

- Using the original slicing of the spectral parameters, compute \mathbf{T} :

$$\Phi_{\mathbf{x}_2} = \mathbf{O}_3 \mathbf{T} \text{diag} \left([\mathbf{O}_2]_{\mathbf{x}_2, \cdot} \right) \mathbf{O}_2^\dagger$$

- Compute θ_{ab} (up to proportionality):

$$[\mathbf{T}]_{i,\sigma} = Z_\sigma^{-1} \prod_j \theta_{i,\sigma_j} \quad \text{or} \quad [\mathbf{T}]_{i,[j,j,j,j]} = Z_\sigma^{-1} \theta_{i,\sigma_j}^4$$

- May be computed by multiple methods:
 - minimize least squares in log space (linear)
 - exponentiation (and some divisions)
 - minimize least squares

Future work

- Spectral learning
 - Other models
 - More unknowns (bilinear least squares, etc.)
 - Non-quadratic cost functions
- Spectral Learning of MRFs:
 - Make it work in practice
 - Also make it work for non-symmetric MRFs
 - Run experiments
 - Cry