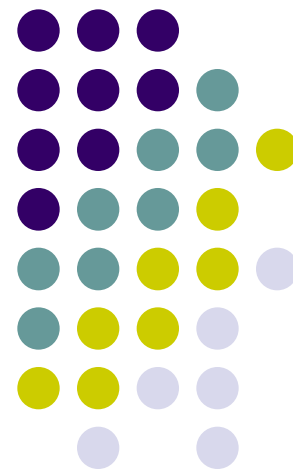
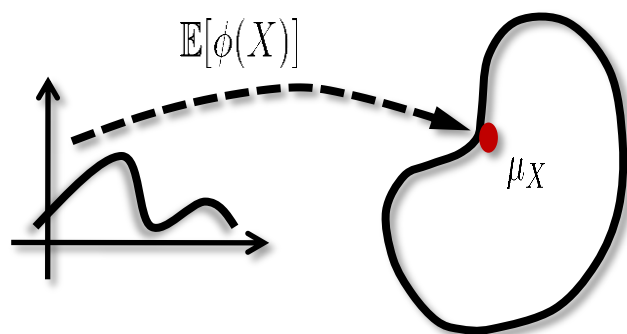


Probabilistic Graphical Models

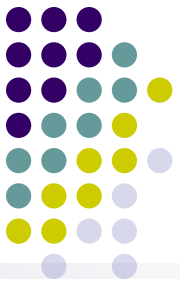
Spectral Algorithms for Graphical Models

Ankur Parikh

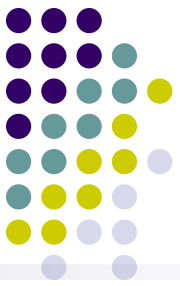
Lecture 21, April 3, 2013



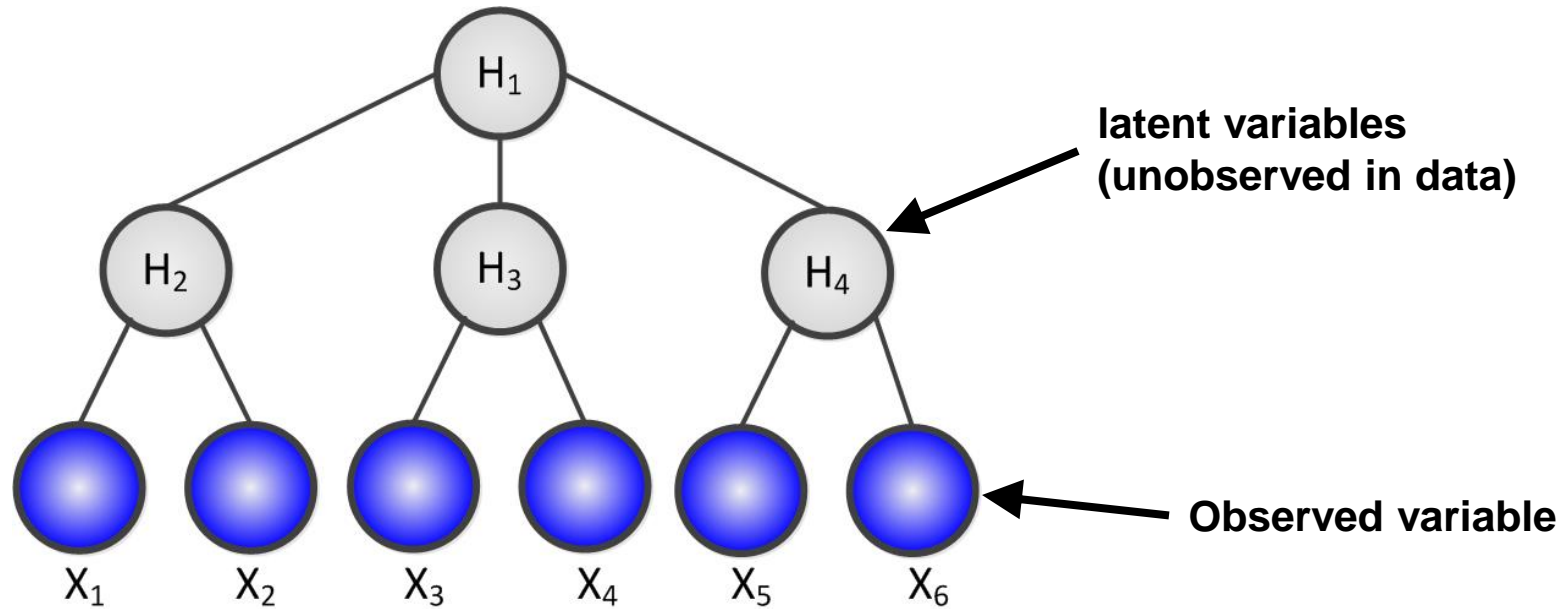
The Linear Algebra View of Latent Variable Models



- Linear algebra can provide a different perspective of graphical models
- Specifically, linear algebra for latent variable models (**Spectral Algorithms**)
- As we will see in this lecture, there is a deep connection between the notion of low rank models and latent variable models.
- This connection will allow us to derive local minima free learning algorithms for latent variable models (that do not use EM).



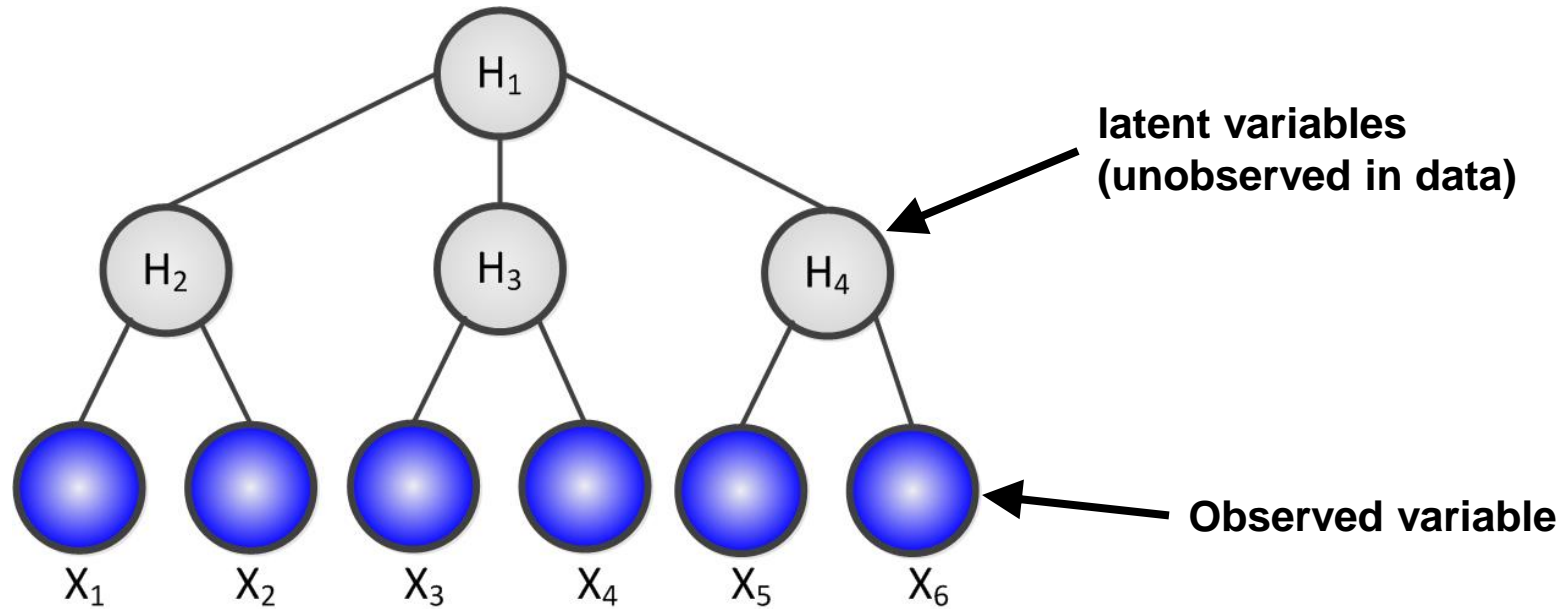
Consider the following Tree



$$\mathbb{P}[X_1, \dots, X_6, H_1, \dots, H_4] = \mathbb{P}[H_1] \prod_{i=2}^4 \mathbb{P}[H_i | H_{\pi(H_i)}] \prod_{i=1}^6 \mathbb{P}[X_i | H_{\pi(X_i)}]$$

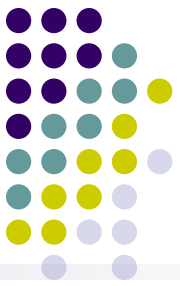


Learning Parameters (EM)



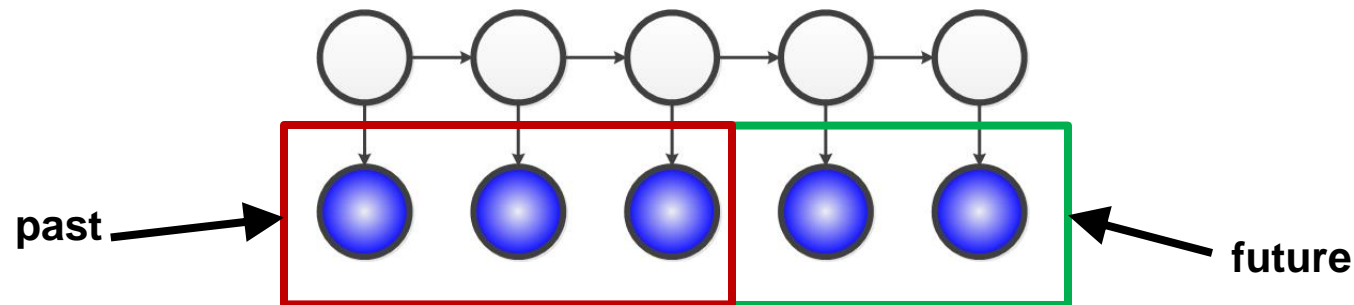
Since latent variables are not observed in the data, we have to use Expectation Maximization (EM) to learn parameters

- **Slow**
- **Local Minima**



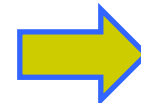
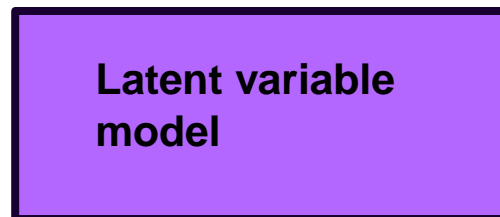
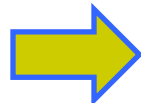
Focusing on Prediction

- In many applications that use latent variable models, the end task is not to recover the latent states, but rather to use the model for prediction among observed variables.
- Dynamical Systems – Predict future given past

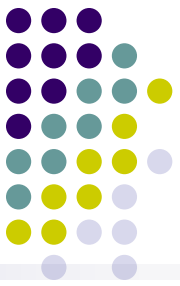


- Machine Translation

Apprentissage spectral est Awesome



Spectral Learning is Awesome

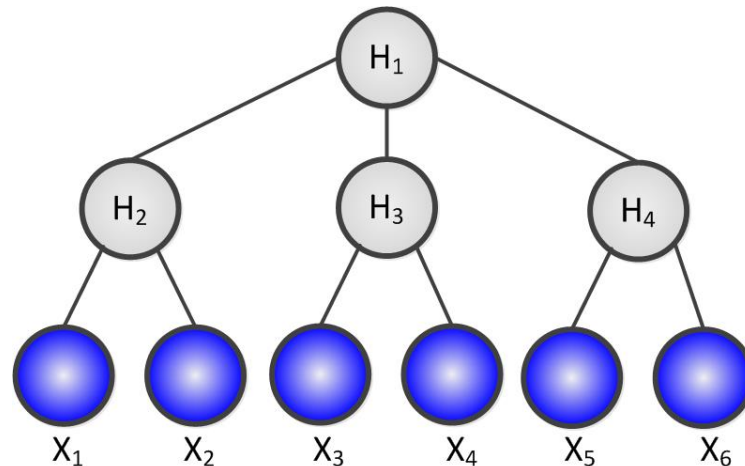


Focusing on Prediction

- As a result, for most of this lecture, we are concerned with quantities related to the observed variables:

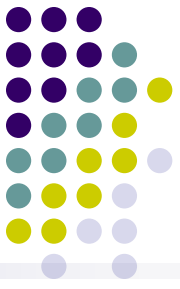
$$\mathcal{P}[X_1, X_2, X_3, X_4, X_5, X_6]$$

- We do not care about the latent variables explicitly.

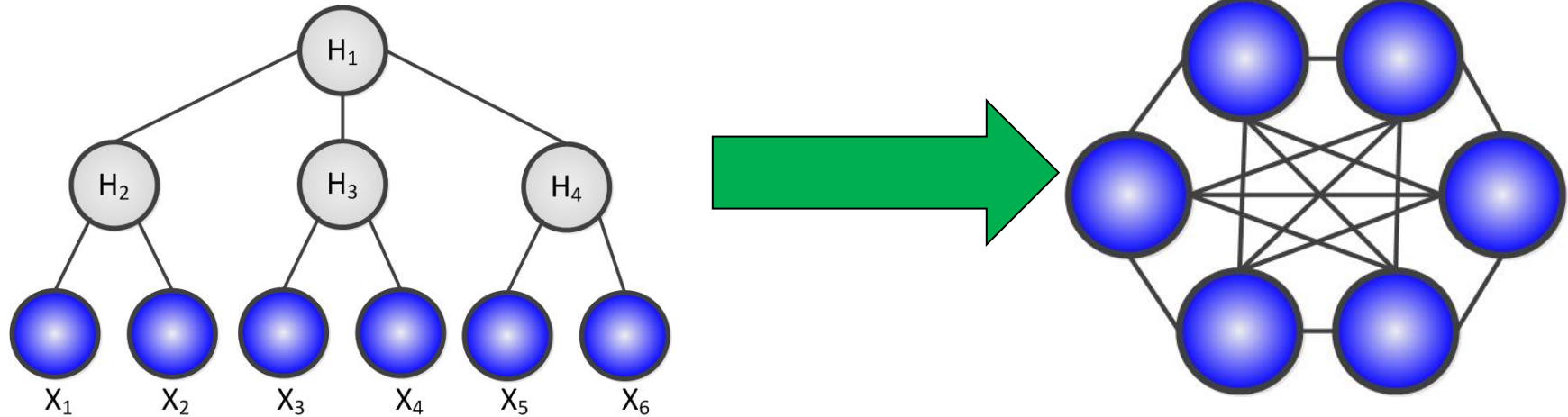


- **Do we still need EM to learn the parameters?**

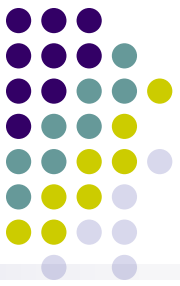
But if we don't care about the latent variables....



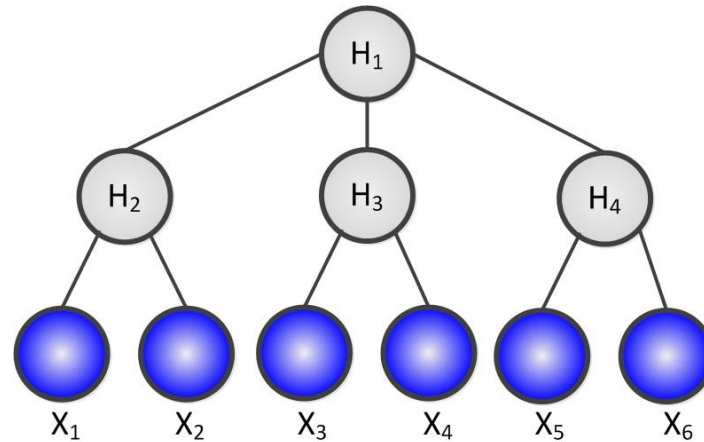
- Why don't we just integrate them out?
- Because integrating them out results in a clique ☹️



- The clique is the minimal I-map.



Marginal Does Not Factorize



$$\mathbb{P}[X_1, X_2, X_3, X_4, X_5, X_6] = \sum_{H_1, \dots, H_6} \prod_{i=2}^4 \mathbb{P}[H_i | H_{\pi(H_i)}] \prod_{i=1}^6 \mathbb{P}[X_i | H_{\pi(X_i)}]$$

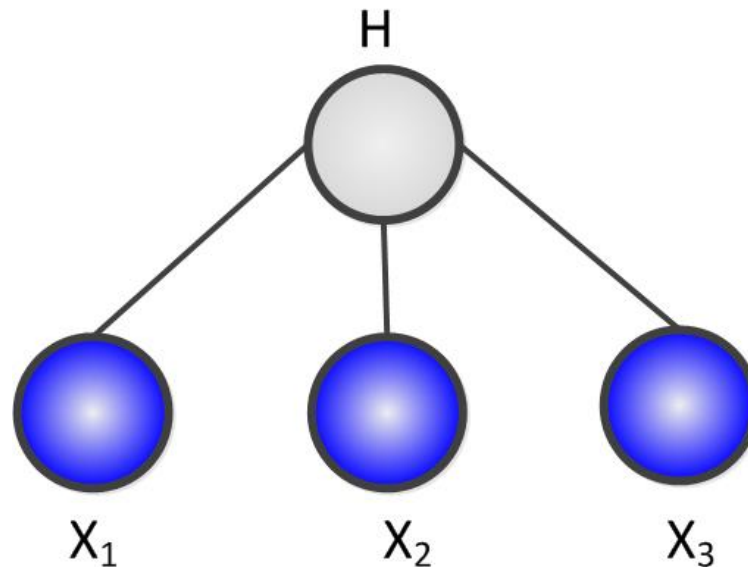
Does not factorize due to the outer sum (Can somewhat distribute the sum, but doesn't solve problem)

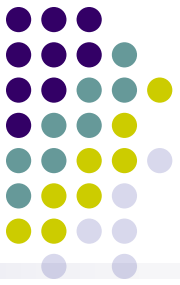
$$\mathbb{P}[X_1, X_2, X_3, X_4, X_5, X_6] = \sum_{H_1} \mathbb{P}[H_1] \prod_{i=2}^4 \sum_{H_i} \mathbb{P}[H_i | H_{\pi(H_i)}] \prod_{i=1}^6 \sum_{X_i} \mathbb{P}[X_i | H_{\pi(X_i)}]$$

But isn't a latent tree different from a clique?



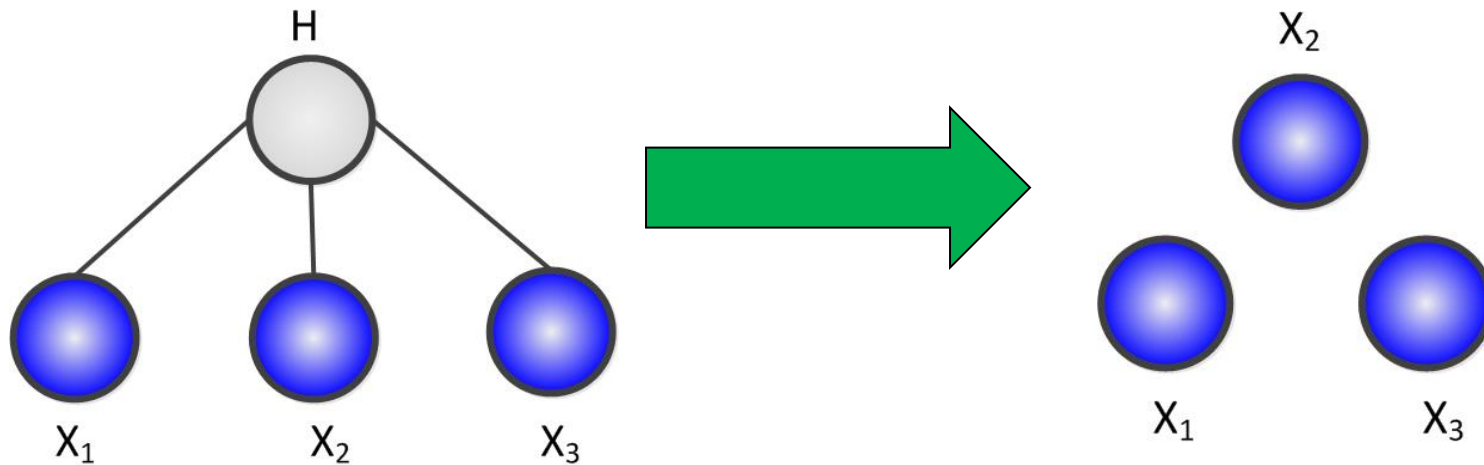
- It depends on the number of latent states.
- Consider the following model.

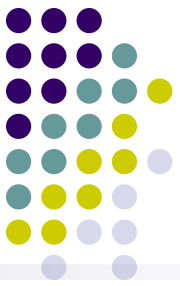




If H has only one state.....

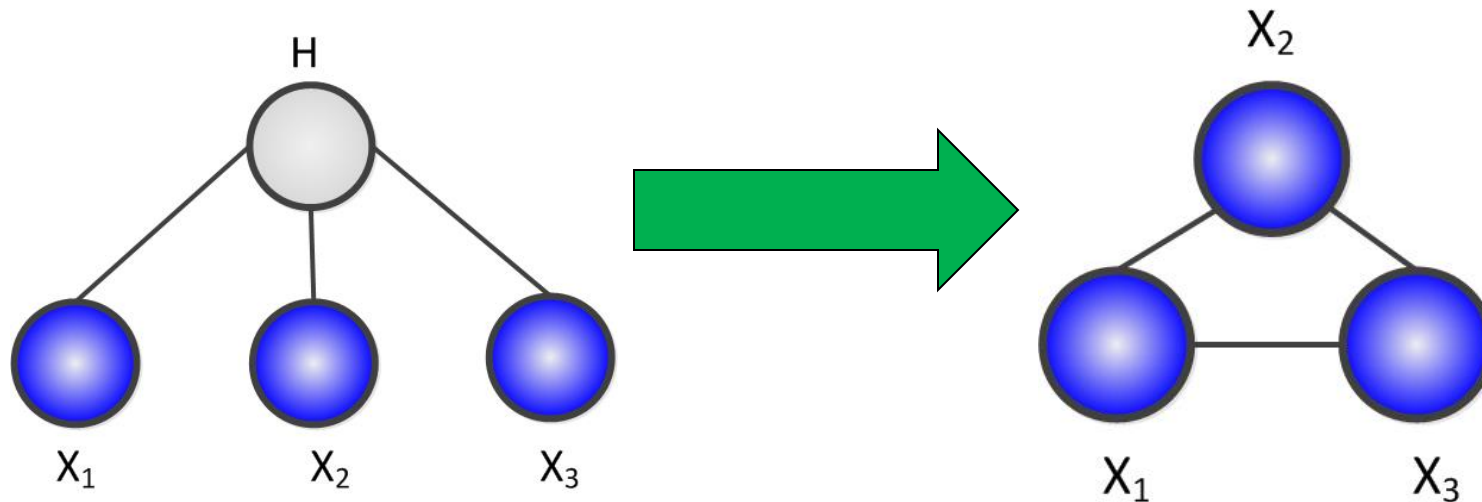
- Then the observed variables are independent!





What if H has many states?

- Let us say the observed variables each have m states.
- Then if H has m^3 states then the latent model can be exactly equivalent to a clique (depending on how parameters are set).



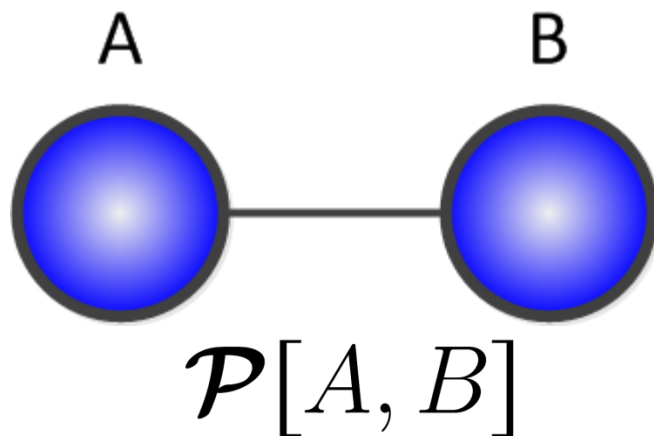
- But what about all the other cases?

The Question



- Under existing methods, latent models all require EM to learn regardless of the number of hidden states.
- However, is there a formulation of latent variable models where the difficulty of learning is a function of the number of latent states?
- This is the question that the *spectral view* will answer.

Graphical Models: The Linear Algebra View



A and B have m states each.

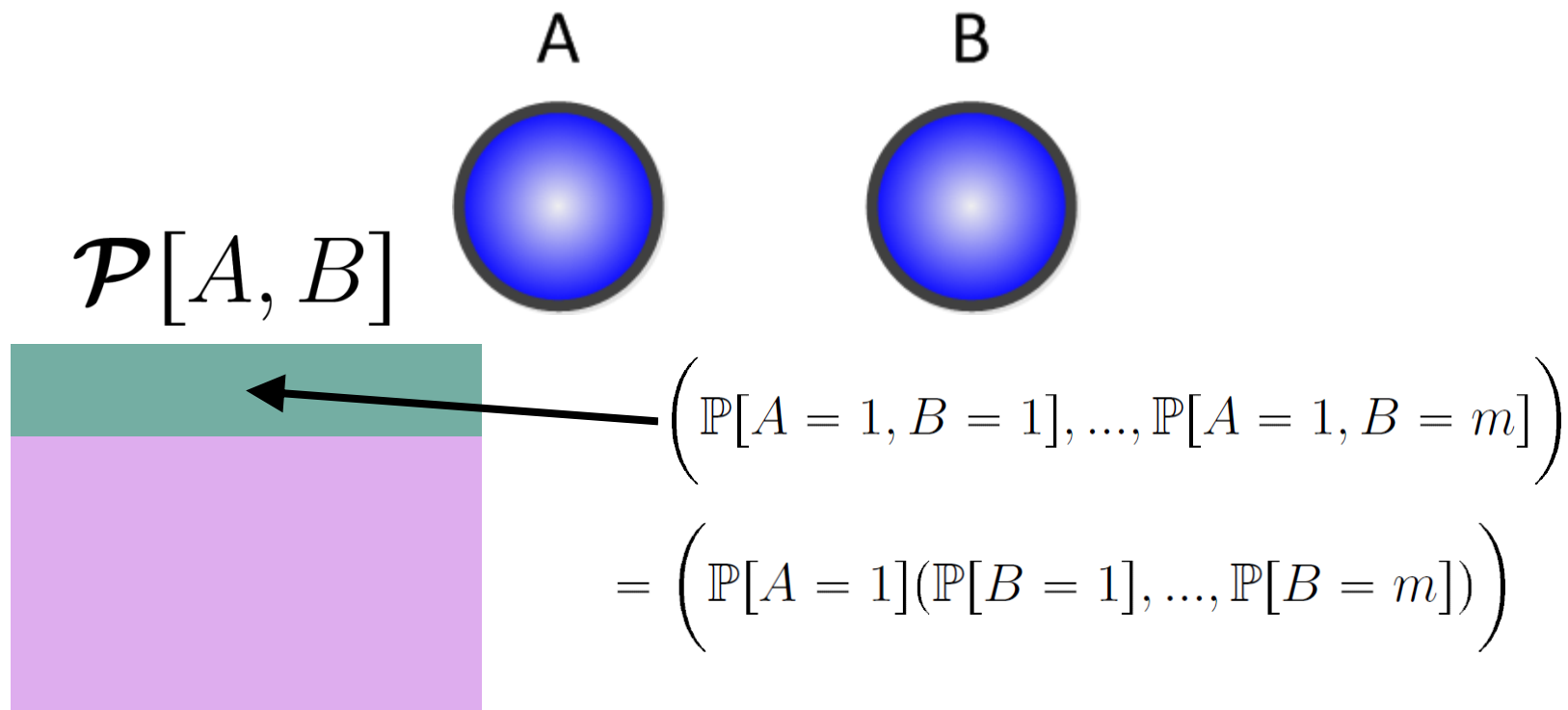


- In general, nothing we can say about the nature of this matrix.

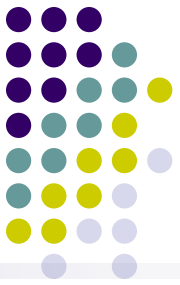
Independence: The Linear Algebra View



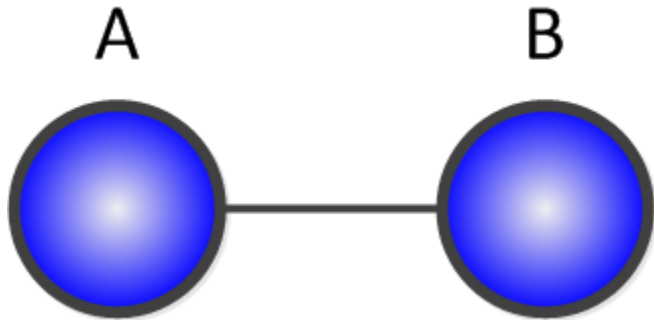
- What if we know A and B are independent?



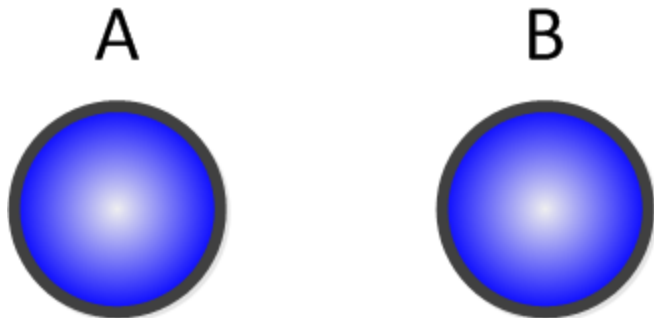
- Joint probability matrix is rank one, since all rows are multiples of one another!!



Independence and Rank

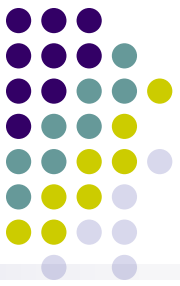


$\mathcal{P}[A, B]$ has rank m (at most)



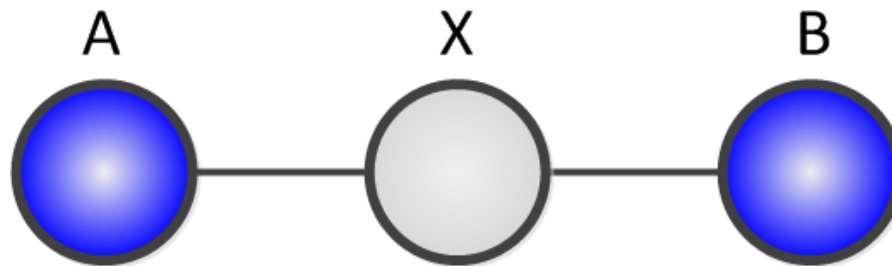
$\mathcal{P}[A, B]$ has rank 1

- What about rank in between 1 and m ?



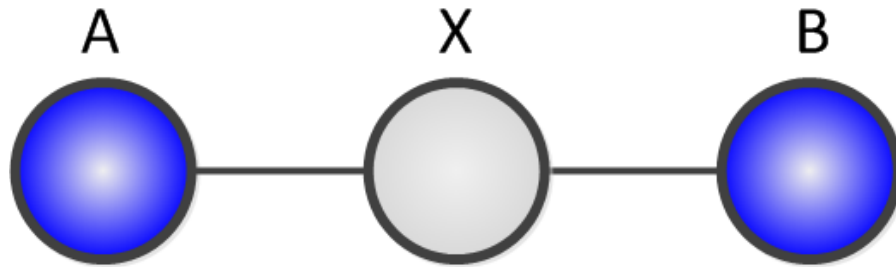
Low Rank Structure

- **A** and **B** are not marginally independent (They are only conditionally independent given **X**).




- Assume **X** has **k** states (while **A** and **B** have **m** states).
- Then, $\text{rank}(\mathcal{P}[A, B]) \leq k$
- Why?

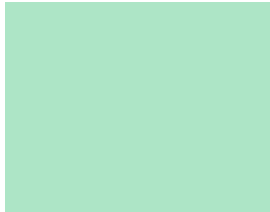
Low Rank Structure




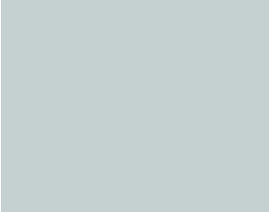
$$\mathcal{P}[A, B] = \mathcal{P}[A|X] \mathcal{P}(\emptyset|X) \mathcal{P}[B|X]^T$$

 $\text{rank} \leq k$

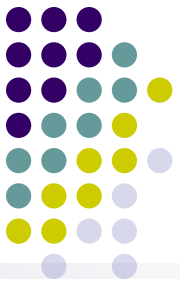
$=$

 $\text{rank} \leq k$

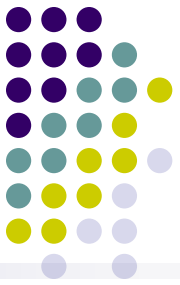
 $\text{rank} \leq k$

 $\text{rank} \leq k$

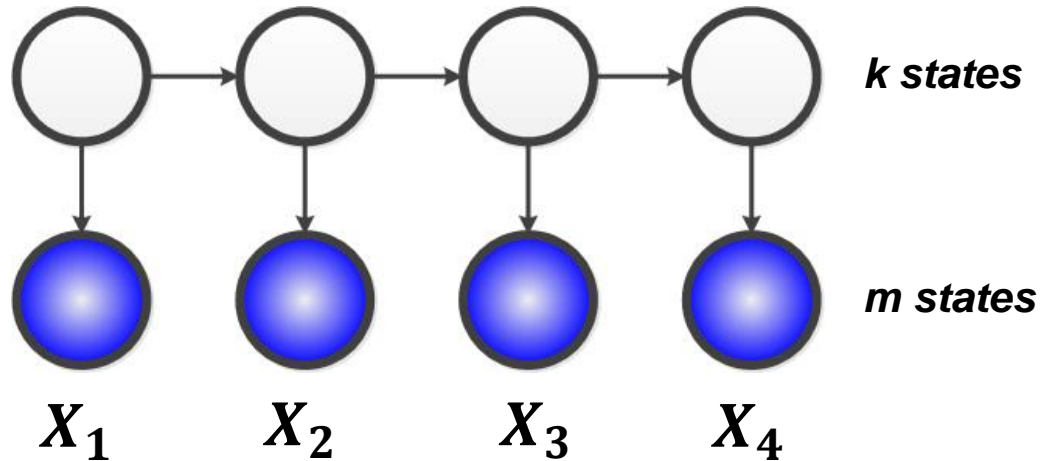
The Spectral View



- Latent variable models encode **low rank dependencies** among variables (*both marginal and conditional*)
- Use tools from linear algebra to exploit this structure.
 - Rank
 - Eigenvalues
 - SVD
 - Tensors
- In the rest of the lecture we will focus on **parameter learning** in latent variable models.



A More Interesting Example

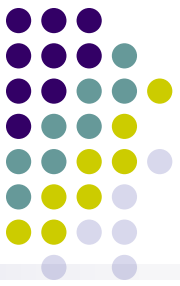


$$\mathcal{P}[X_{\{1,2\}}, X_{\{3,4\}}]$$

$\{X_1, X_2\}$

$\{X_3, X_4\}$

has rank k



Low Rank Matrices “Factorize”

$$M = RL$$

m by n

m by k k by n

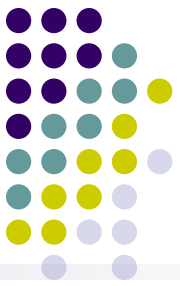
If M has rank k

We already know one factorization!!!

$$\mathcal{P}[X_{\{1,2\}}, X_{\{3,4\}}] = \mathcal{P}[X_{\{1,2\}} | H_2] \mathcal{P}[\bigoplus H_2] \mathcal{P}[X_{\{3,4\}} | H_2]^\top$$

Factor of 4 variables Factor of 3 variables \uparrow Factor of 3 variables

Factor of 1 variable



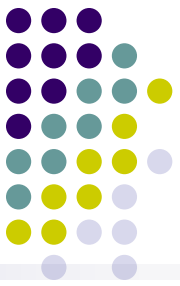
Alternate Factorizations

- The key insight is that this factorization is not unique.
- Consider Matrix Factorization. Can add any invertible transformation:

$$M = RL$$

$$M = RSS^{-1}L$$

- **The magic of spectral learning is that there exists an alternative factorization that only depends on observed variables!**



An Alternate Factorization

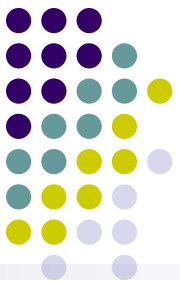
- Let us say we only want to factorize this matrix of 4 variables

$$\mathcal{P}[X_{\{1,2\}}, X_{\{3,4\}}]$$

such that it is product of matrices that contain at most three *observed* variables e.g.

$$\mathcal{P}[X_{\{1,2\}}, X_3]$$

$$\mathcal{P}[X_2, X_{\{3,4\}}]$$



An Alternate Factorization

- Note that

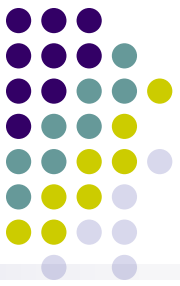
$$\mathcal{P}[X_{\{1,2\}}, X_3] = \underbrace{\mathcal{P}[X_{\{1,2\}}|H_2]}_{\text{green}} \underbrace{\mathcal{P}[\ominus H_2]}_{\text{green}} \underbrace{\mathcal{P}[X_3|H_2]}_{\text{red}}^\top$$

$$\mathcal{P}[X_2, X_{\{3,4\}}] = \underbrace{\mathcal{P}[X_2|H_2]}_{\text{red}} \underbrace{\mathcal{P}[\ominus H_2]}_{\text{red}} \underbrace{\mathcal{P}[X_{\{3,4\}}|H_2]}_{\text{green}}^\top$$

- Product of green terms (in some order) is

$$\mathcal{P}[X_{\{1,2\}}, X_{\{3,4\}}]$$

- Product of red terms (in some order) is $\mathcal{P}[X_2, X_3]$



An Alternate Factorization

$$\mathcal{P}[X_{\{1,2\}}, X_{\{3,4\}}] = \mathcal{P}[X_{\{1,2\}}, X_3] \mathcal{P}[X_2, X_3]^{-1} \mathcal{P}[X_2, X_{\{3,4\}}]$$

factor of 4 variables

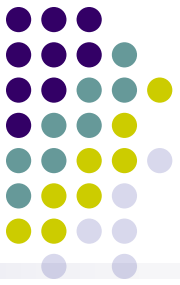
factor of 3 variables

factor of 3 variables

Advantage: Factors are only functions of observed variables! Can be directly computed from data without EM!!!!

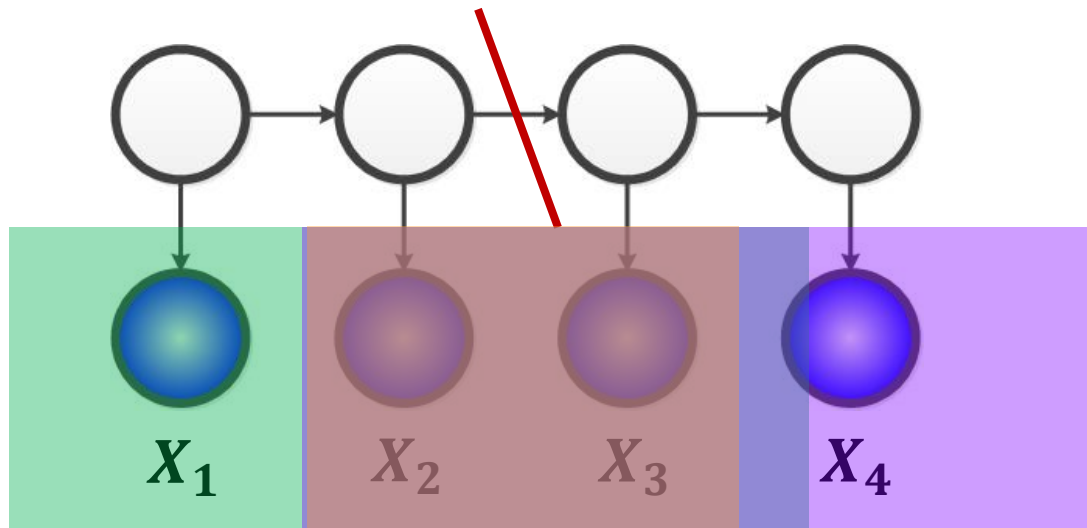
Caveat: Factors are no longer probability tables (do not have to be non-negative)

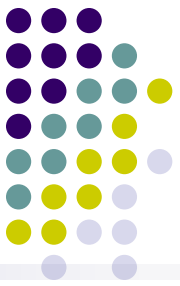
We will call this factorization the **observable factorization**.



Graphical Relationship

$$\mathcal{P}[X_{\{1,2\}}, X_{\{3,4\}}] = \mathcal{P}[X_{\{1,2\}}, X_3] \mathcal{P}[X_2, X_3]^{-1} \mathcal{P}[X_2, X_{\{3,4\}}]$$





Our Alternate Factorization

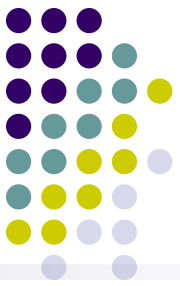
$$\mathcal{P}[X_{\{1,2\}}, X_{\{3,4\}}] = \mathcal{P}[X_{\{1,2\}}, X_3] \mathcal{P}[X_2, X_3]^{-1} \mathcal{P}[X_2, X_{\{3,4\}}]$$

factor of 4 variables

factor of 3 variables

factor of 4 variables

- It may not seem very amazing at the moment (we have only reduced the size of the factor by 1)
- What is cool is that every latent tree of N variables has such a factorization where:
 - All factors are of size 3 (for simplicity we only consider binary trees)
 - All factors are only functions of observed variables
- We will prove why in a few slides.

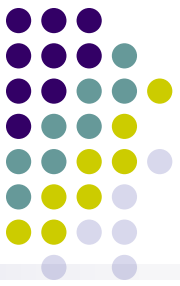


Where's the Catch?

- Before we said that if the number of latent states was very large then the model was equivalent to a clique.
- Where does that scenario enter in our factorization?

$$\mathcal{P}[X_{\{1,2\}}, X_{\{3,4\}}] = \mathcal{P}[X_{\{1,2\}}, X_3] \mathcal{P}[X_2, X_3]^{-1} \mathcal{P}[X_2, X_{\{3,4\}}]$$

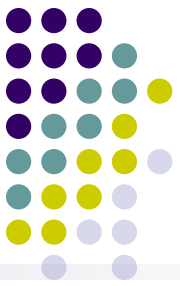
When does this inverse exist?



When Does the Inverse Exist

$$\mathcal{P}[X_2, X_3] = \mathcal{P}[X_2|H_2]\mathcal{P}[\ominus H_2]\mathcal{P}[X_3|H_2]^\top$$

- All the matrices on the right hand side must have full rank. (This is in general a requirement of spectral learning, although it can be somewhat relaxed, see Siddiqi et al. 2009)
- If m (number of observed states) $>$ k (number of hidden states), then the inverse cannot exist, but this situation is easily fixable (project onto lower dimensional space)
- But what happens when $k > m$?



When $k > m$

- The inverse does exist. But it no longer satisfies the following property, which we used to derive the factorization

$$\mathcal{P}[X_2, X_3]^{-1} = (\mathcal{P}[X_3|H_2]^\top)^{-1} \mathcal{P}[\ominus H_2]^{-1} \mathcal{P}[X_2|H_2]^{-1}$$

- This is much more difficult to fix, and intuitively corresponds to how the problem becomes intractable if $k \gg m$.

Relationship to Original Factorization



- What is the relationship between the original factorization and the new factorization?

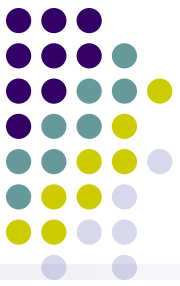
$$\underbrace{\mathcal{P}[X_{\{1,2\}}, X_{\{3,4\}}]}_M = \underbrace{\mathcal{P}[X_{\{1,2\}}|H_2]\mathcal{P}[\ominus H_2]}_R \underbrace{\mathcal{P}[X_{\{3,4\}}|H_2]}_L^\top$$

$$M = RL$$

$$M = RSS^{-1}L$$

Can I choose S to get the observable factorization?

Relationship to Original Factorization



- Let

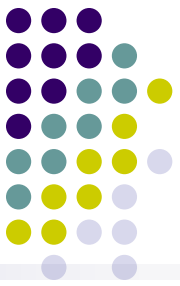
$$S := \mathcal{P}[X_3 | H_2]$$

$$\begin{aligned} \mathcal{P}[X_{\{1,2\}}, X_{\{3,4\}}] &= \mathcal{P}[X_{\{1,2\}}, X_3] \mathcal{P}[X_2, X_3]^{-1} \mathcal{P}[X_2, X_{\{3,4\}}] \\ &= \underline{LS} \qquad \qquad \qquad = \underline{S^{-1}R} \end{aligned}$$

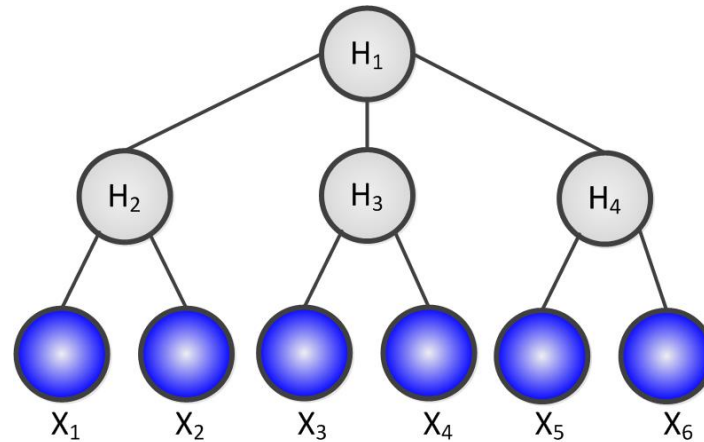
Overview of Spectral Learning Magic



- Represent original latent model as a decomposition of matrices (tensors in our case)
- Derive an alternate factorization that is only a function of observed variables (which we term the observable factorization).
- Learn parameters for alternate factorization directly from the data without using EM!
 - Local minima free and provably consistent.
- **For simplicity, assume number of hidden states is equal to the number of observed states and all CPTs have full rank.**



The Marginal Does Not Factorize

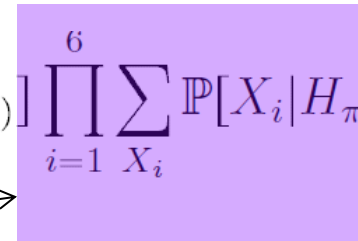


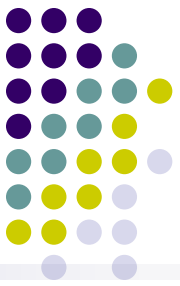
$$\mathbb{P}[X_1, X_2, X_3, X_4, X_5, X_6] = \sum_{H_1, \dots, H_6} \prod_{i=2}^4 \mathbb{P}[H_i | H_{\pi(H_i)}] \prod_{i=1}^6 \mathbb{P}[X_i | H_{\pi(X_i)}]$$

Does not factorize due to the outer sum. Can somewhat distribute the sum.

$$\mathbb{P}[X_1, X_2, X_3, X_4, X_5, X_6] = \sum_{H_1} \mathbb{P}[H_1] \prod_{i=2}^4 \sum_{H_i} \mathbb{P}[H_i | H_{\pi(H_i)}] \prod_{i=1}^6 \sum_{X_i} \mathbb{P}[X_i | H_{\pi(X_i)}]$$

Sum product reminds us
of matrix multiplication





Matrices are Insufficient

- Matrix factorization

$$M = RL$$

- For example, we can do the following:

$$\mathcal{P}[X_{\{1,2,3,4\}}, X_{\{5,6\}}] = \mathcal{P}[X_{\{1,2,3,4\}} | H_1] \mathcal{P}[\neg H_1] \mathcal{P}[X_{\{5,6\}} | H_1]^\top$$

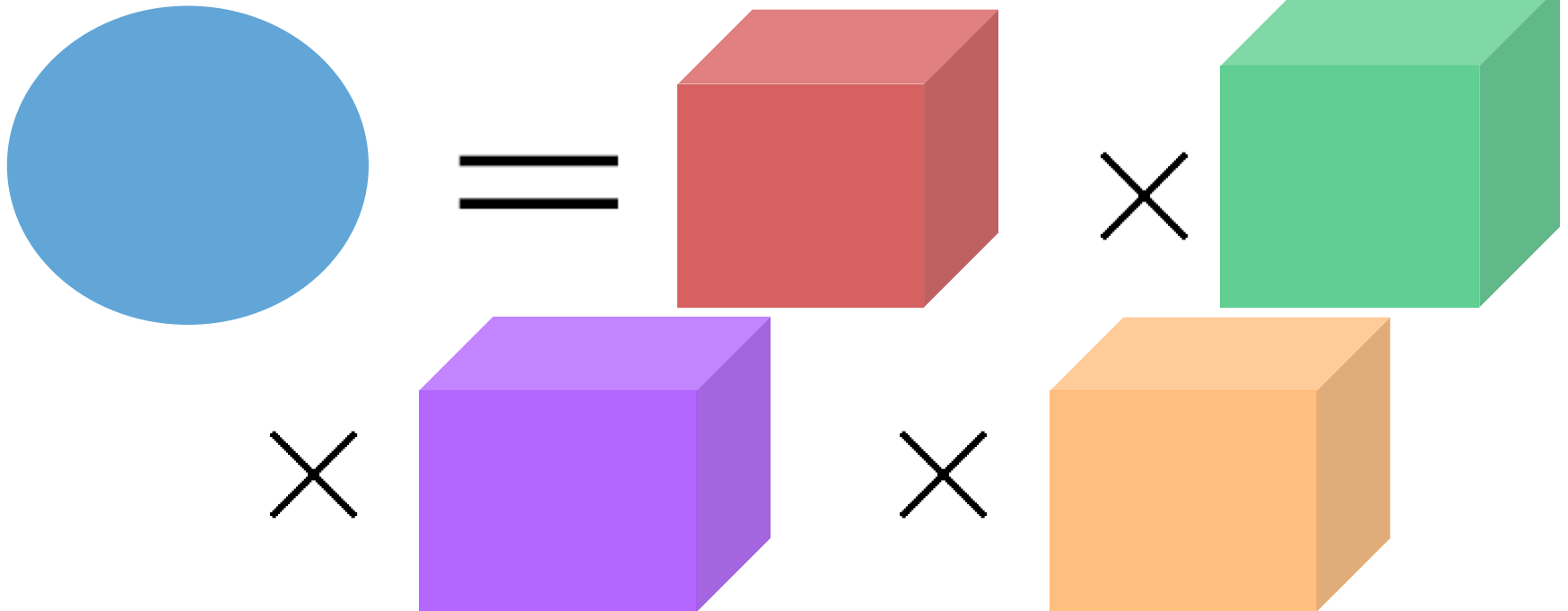
- But then there is nothing more we can do.

Tensor Factorization



- Express higher order tensors as products of lower order tensors.

$\mathcal{P}[X_1, X_2, X_3, X_4, X_5, X_6]$



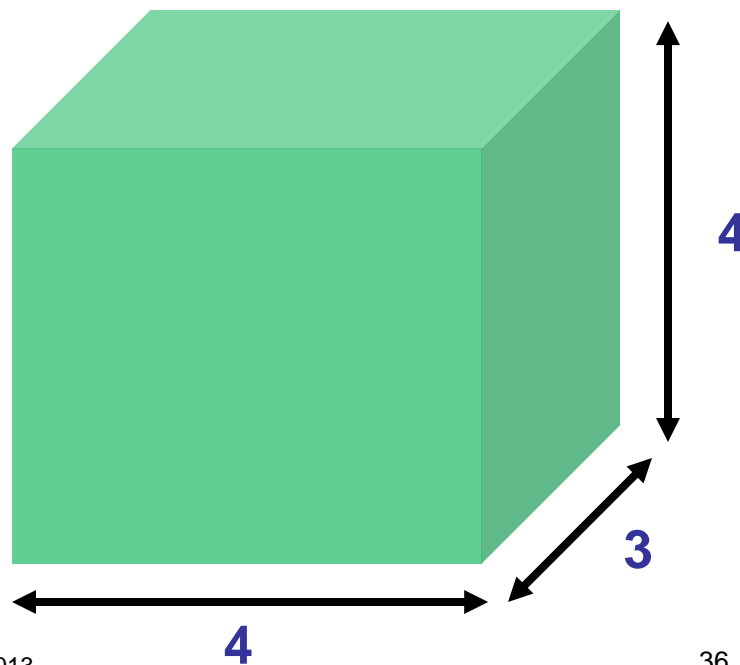
Tensors



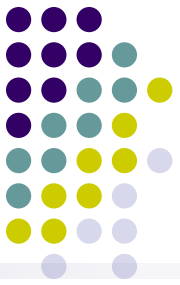
- Multidimensional arrays
- A Tensor of order \mathbf{N} has \mathbf{N} modes (\mathbf{N} indices):

$$\mathcal{T}(i_1, \dots, i_N)$$

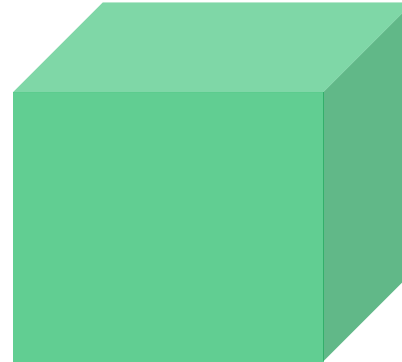
- Each mode is associated with a dimension. In the example,
 - Dimension of mode 1 is 4
 - Dimension of mode 2 is 3
 - Dimension of mode 3 is 4



Tensors



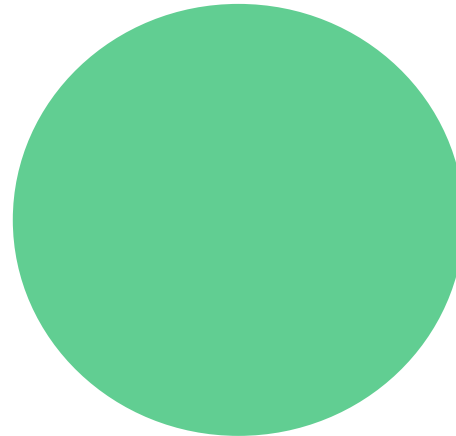
**1st order tensor
(vector)**



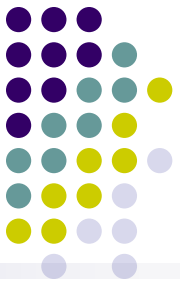
**3rd order tensor
(cube)**



**2nd order tensor
(matrix)**



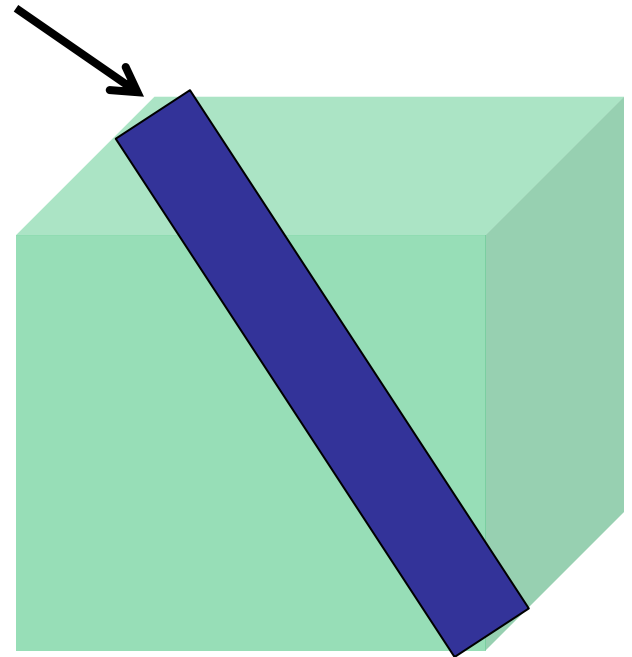
**Higher order
tensors**



Diagonal Tensors

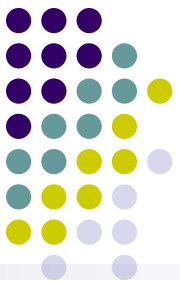
$\mathcal{P}[X]$

$$\mathcal{T}(i, j, k) = \begin{cases} \mathbb{P}[X = i] & \text{if } i = j = k \\ 0 & \text{otherwise} \end{cases}$$



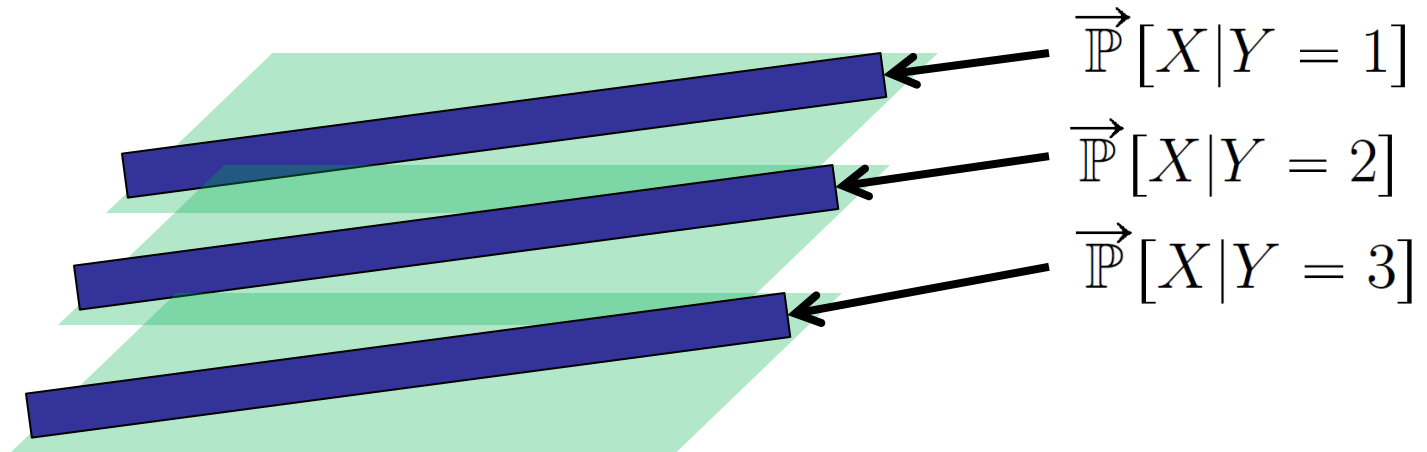
Notation for Diagonal Third Order Tensors:

$$\mathcal{P}[\textcircled{\diagup}_3 X]$$

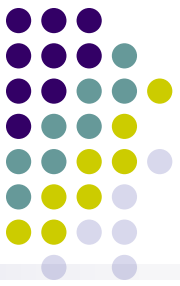


Partially Diagonal Tensors

$$\mathcal{T}(i, j, k) = \begin{cases} \mathbb{P}[X = i | Y = k] & \text{if } i = j \\ 0 & \text{otherwise} \end{cases}$$



Notation: $\mathcal{P}[\otimes X | Y]$



Tensor Vector Multiplication

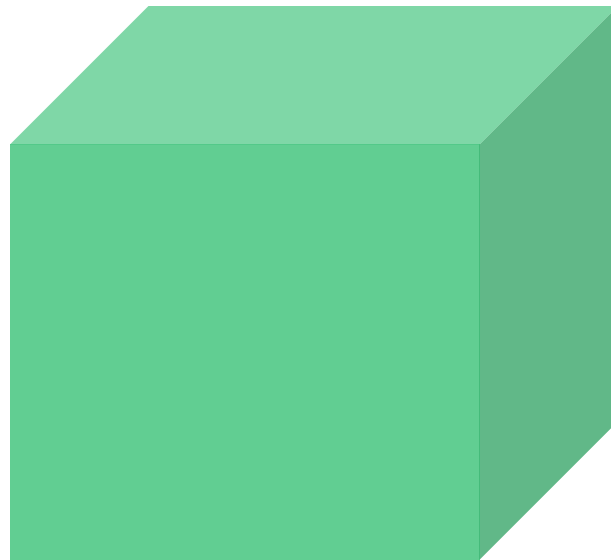
- Multiplying a 3rd order tensor by a vector produces a matrix

M



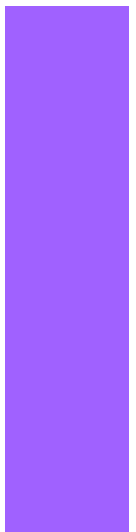
=

\mathcal{T}



\times_1

v



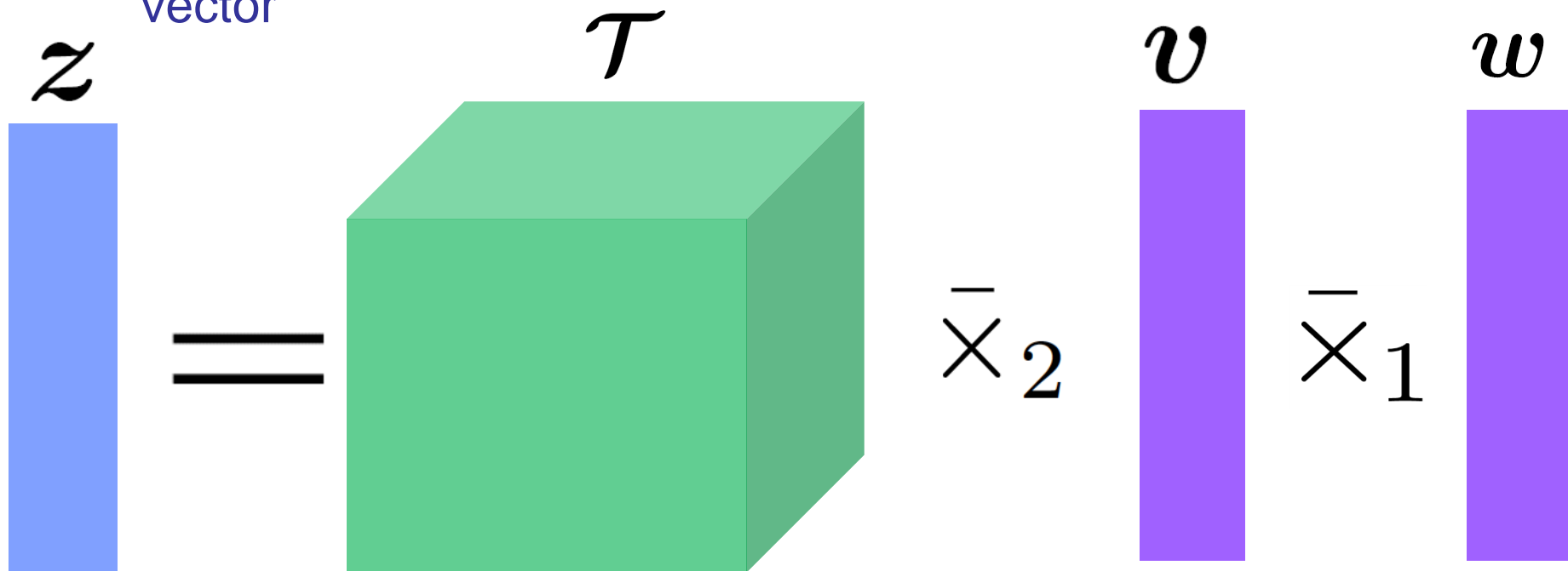
$$M(j, k) = \sum_i \mathcal{T}(i, j, k) v(i)$$

Tensor Vector Multiplication

Cont.

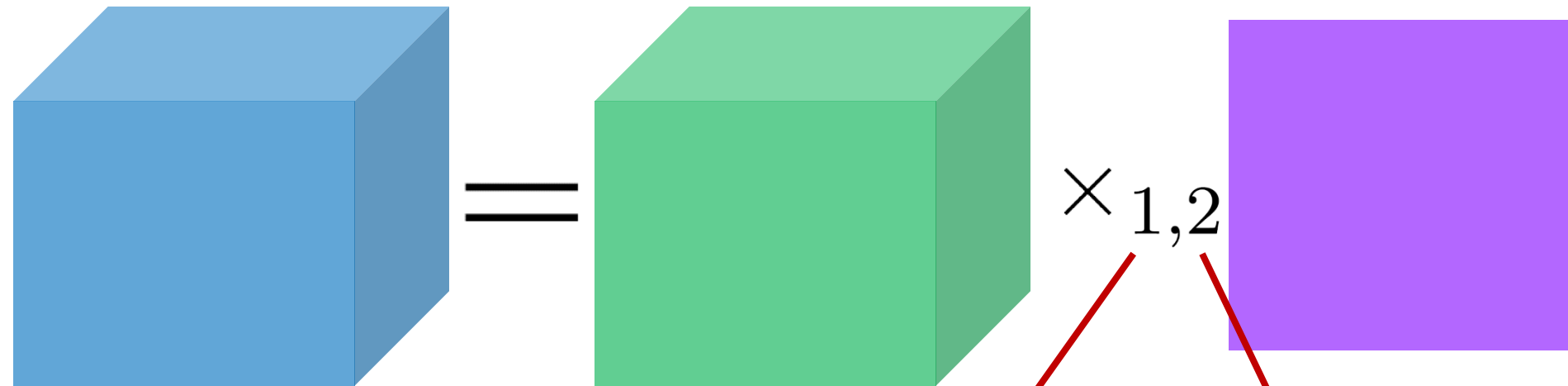
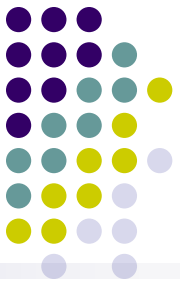


- Multiplying a 3rd order tensor by two vectors produces a vector



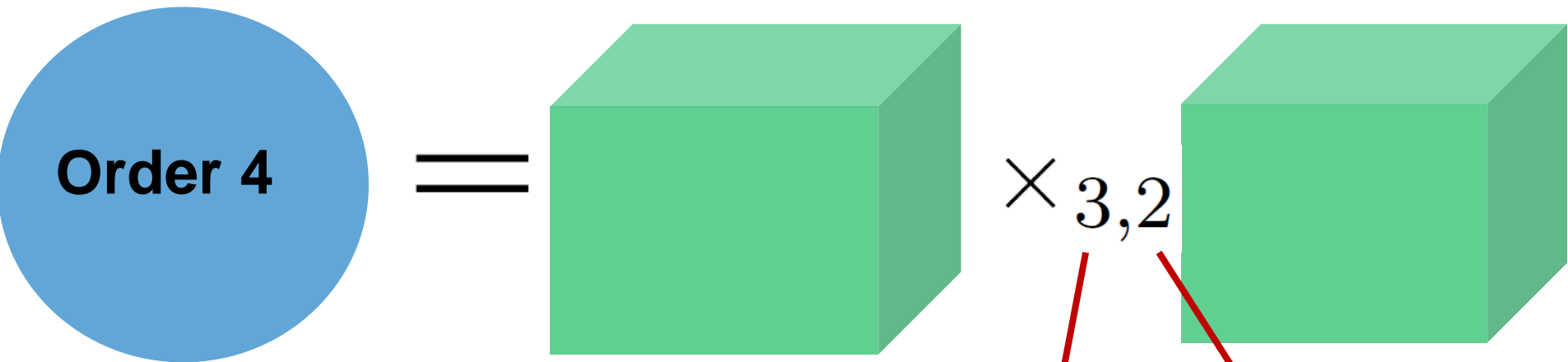
$$z(k) = \sum_i \left(\sum_j \mathcal{T}(i, j, k) v(i) \right) w(j) = \sum_{i,j} \mathcal{T}(i, j, k) v(i) w(j)$$

Tensor Matrix Multiplication



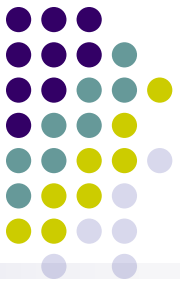
$$\mathbf{R}(i, j, k) = \sum_m \mathcal{T}(m, j, k) \mathbf{M}(i, m)$$

Tensor Tensor Multiplication

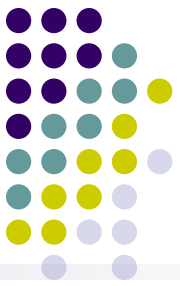


$$\mathbf{R}(i, j, k, l) = \sum_m \mathcal{T}_1(i, j, m) \mathcal{T}_2(k, m, l)$$

Tensor Tensor Multiplication



- Matrix multiplication is closed under multiplication (i.e. product of matrices is a matrix).
- Tensor multiplication is not closed under multiplication (e.g. product of two tensors of order n is a tensor of order $2n-2$)
- **This is key to developing the factorization - multiplying lower order tensors will result in a higher order tensor**



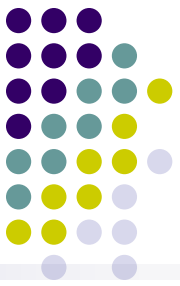
Cleaner Notation

- This is clunky

$$\mathcal{P}[X_1, X_2, X_3] = \mathcal{P}[X_1, X_2, H_1] \times_{3,2} \mathcal{P}[X_3|H_1]$$

- Instead use variable names instead of modes

$$\mathcal{P}[X_1, X_2, X_3] = \mathcal{P}[X_1, X_2, H_1] \times_{H_1} \mathcal{P}[X_3|H_1]$$



The Matricization View

- You can also interpret tensor-tensor (or tensor-matrix) multiplication as first matricizing the tensors, and performing normal matrix multiplication and then reshaping back into a tensor.
- Example:

$$\mathcal{P}[X_1, X_2, X_3] = \mathcal{P}[X_1, X_2, H_1] \times_{H_1} \mathcal{P}[X_3|H_1]$$

$$\mathcal{P}[X_{\{1,2\}}, X_3] = \mathcal{P}[X_{\{1,2\}}, H_1] \mathcal{P}[X_3|H_1]^\top$$

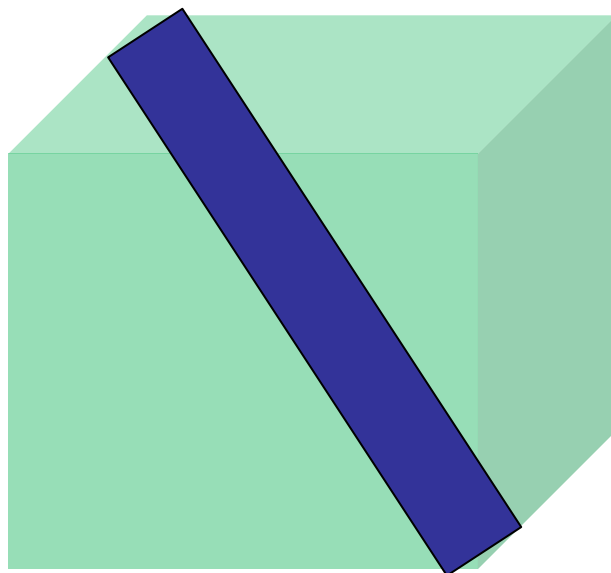
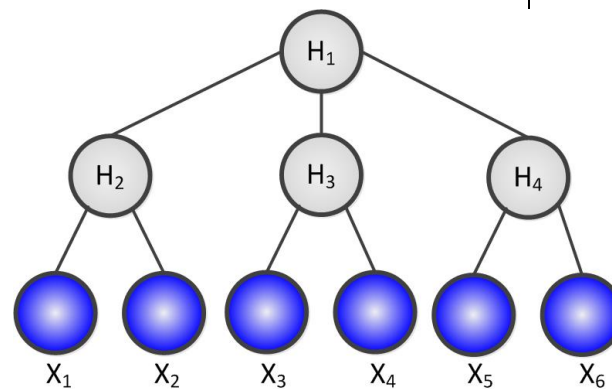


Can reshape this back into a tensor

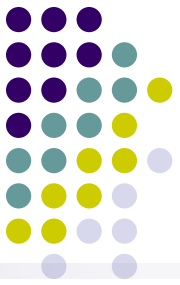
The Tensor Form of The Root



$$\mathcal{P}(\otimes_3 H_1)$$



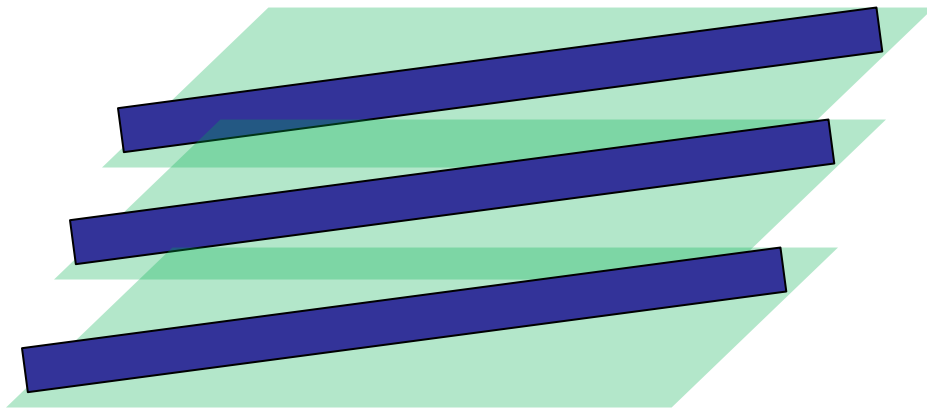
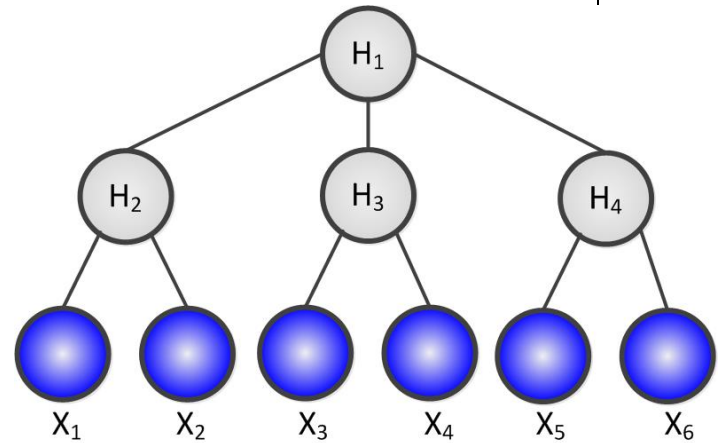
The Tensor Form of the Internal Nodes



$$\mathcal{P}(\otimes H_2 | H_1)$$

$$\mathcal{P}(\otimes H_3 | H_1)$$

$$\mathcal{P}(\otimes H_4 | H_1)$$



The Tensor Form of the Factorization-Leaves



$$\mathcal{P}(X_1|H_2)$$

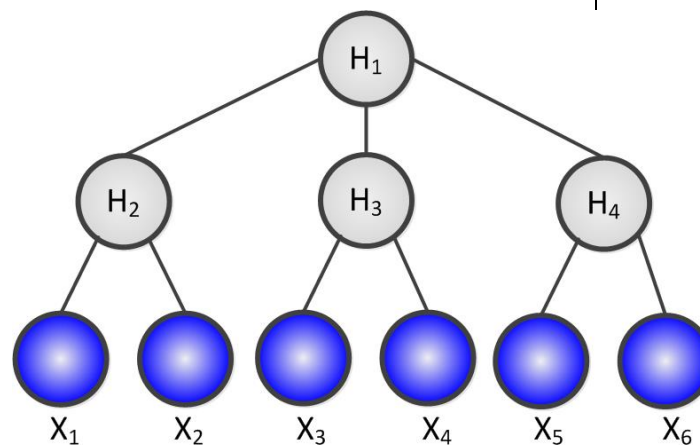
$$\mathcal{P}(X_2|H_2)$$

$$\mathcal{P}(X_3|H_3)$$

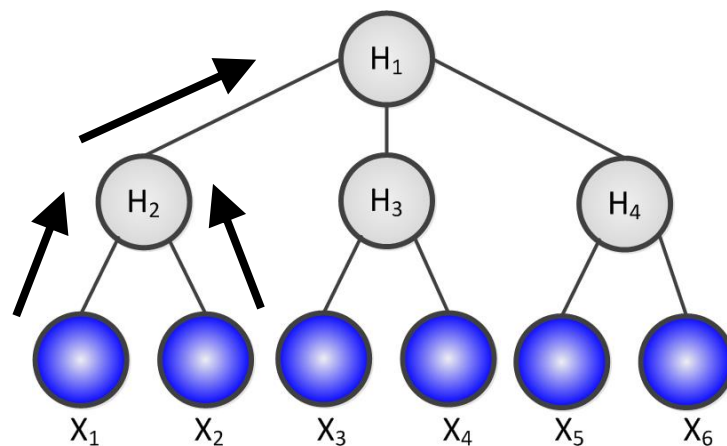
$$\mathcal{P}(X_4|H_3)$$

$$\mathcal{P}(X_5|H_4)$$

$$\mathcal{P}(X_6|H_4)$$



Tensor Message Passing

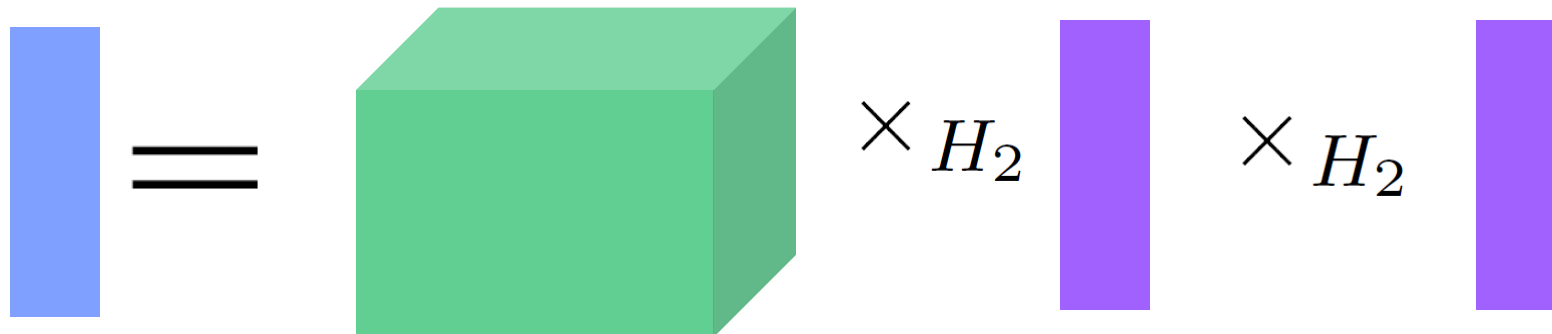


$$\mathcal{P}[x_1, x_2 | H_1]$$

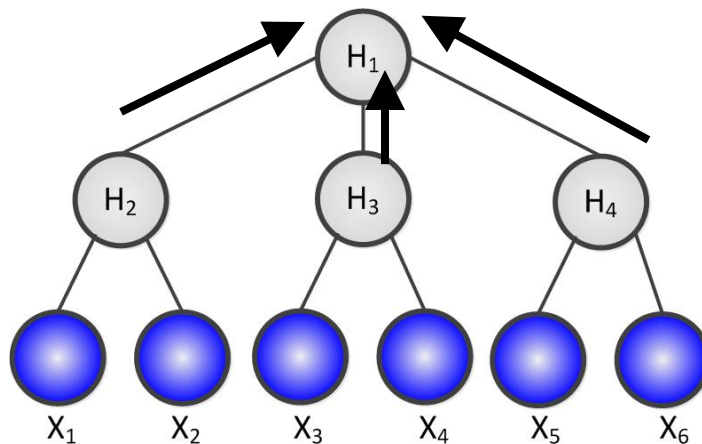
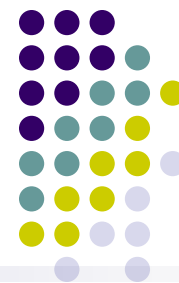
$$\mathcal{P}(\ominus H_2 | H_1)$$

$$\mathcal{P}[x_1 | H_2]$$

$$\mathcal{P}[x_2 | H_2]$$

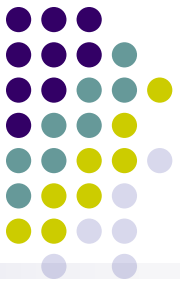


At the Root



$$\mathbb{P}[\text{evidence}] = \mathcal{P}(\bigoplus_3 H_1) \times_{H_1} \mathcal{P}[x_1, x_2 | H_1] \times_{H_1} \mathcal{P}[x_3, x_4 | H_1] \times_{H_1} \mathcal{P}[x_5, x_6 | H_1]$$

To Compute the Marginal Tensor

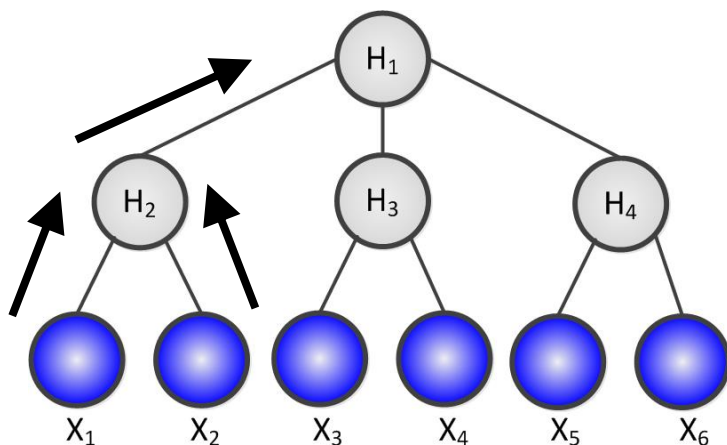


- What if we want to compute the whole marginal tensor?

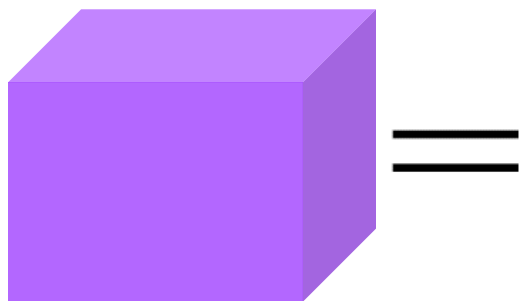
$$\mathcal{P}[X_1, X_2, X_3, X_4, X_5, X_6]$$

- Run message passing, except don't incorporate evidence at the leaves

To Compute the Marginal Tensor



$$\mathcal{P}[X_1, X_2|H_1]$$



$=$

$$\mathcal{P}(\oslash H_2|H_1)$$



\times_{H_2}

$$\mathcal{P}(X_1|H_2)$$



\times_{H_2}

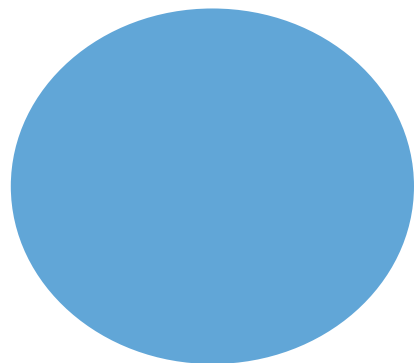
$$\mathcal{P}(X_2|H_2)$$



At the Root

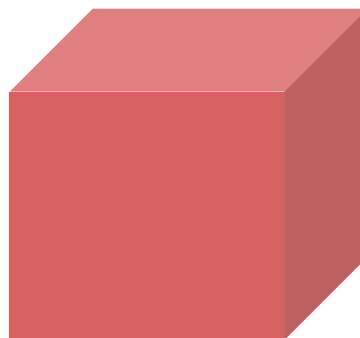


$\mathcal{P}[X_1, X_2, X_3, X_4, X_5, X_6]$

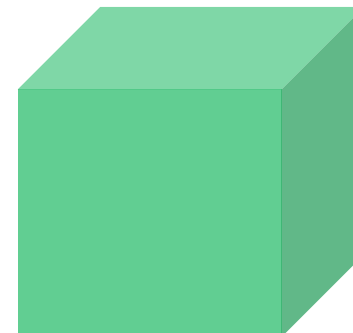


$=$

$\mathcal{P}(\bigoplus_3 H_1)$

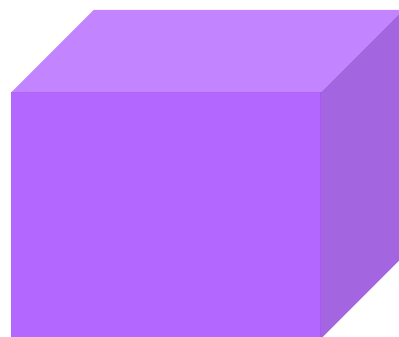


$\mathcal{P}[X_1, X_2|H_1]$



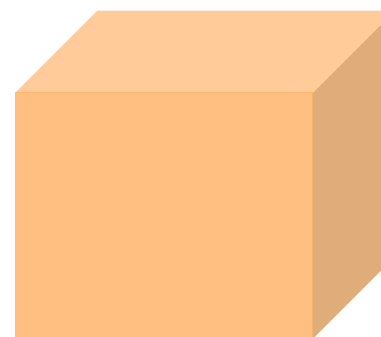
$\times H_1$

$\times H_1$



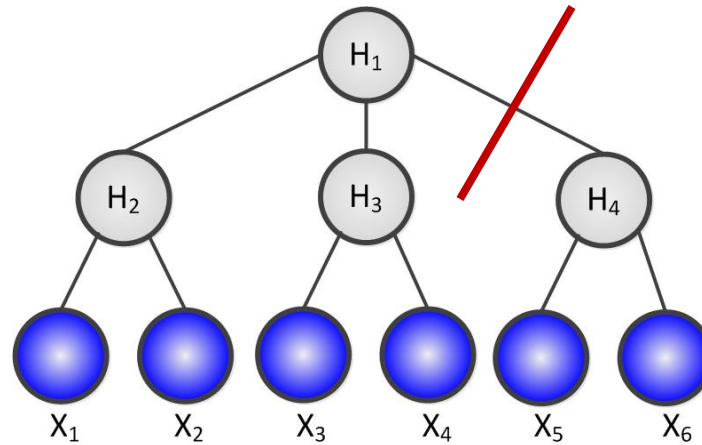
$\mathcal{P}[X_3, X_4|H_1]$

$\times H_1$



$\mathcal{P}[X_5, X_6|H_1]$

Constructing a Different Factorization



Like we did before

$$\mathcal{P}[X_{\{1,2,3,4\}}, X_{\{5,6\}}] = \mathcal{P}[X_{\{1,2,3,4\}}, X_5] \mathcal{P}[X_4, X_5]^{-1} \mathcal{P}[X_4, X_{\{5,6\}}]$$

The Tensor View



“Matricized” Way

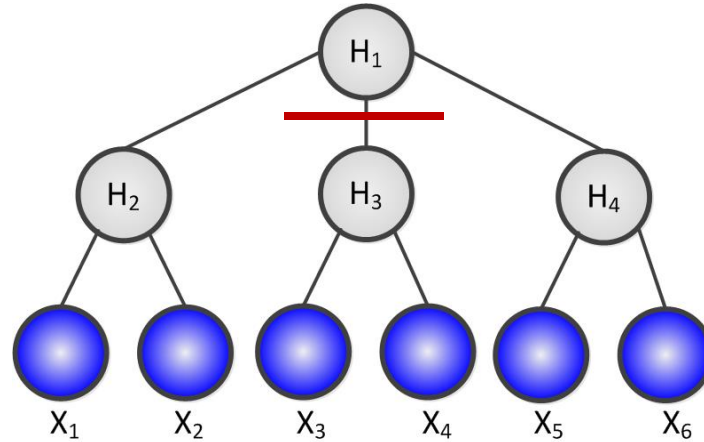
$$\mathcal{P}[X_{\{1,2,3,4\}}, X_{\{5,6\}}] = \mathcal{P}[X_{\{1,2,3,4\}}, X_5] \mathcal{P}[X_4, X_5]^{-1} \mathcal{P}[X_4, X_{\{5,6\}}]$$

“Tensor” Way

$$\begin{aligned} & \mathcal{P}[X_1, X_2, X_3, X_4, X_5, X_6] \\ = & \mathcal{P}[X_1, X_2, X_3, X_4, X_5] \times_{X_5} \mathcal{P}[X_4, X_5]^{-1} \times_{X_4} \mathcal{P}[X_4, X_5, X_6] \end{aligned}$$

Decompose this recursively

Constructing a Different Factorization

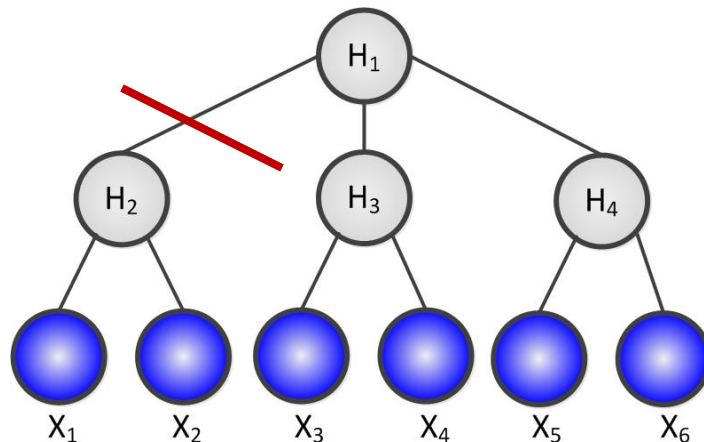


$$\mathcal{P}[X_{\{1,2,5\}}, X_{\{3,4\}}] = \mathcal{P}[X_{\{1,2,5\}}, X_3] \mathcal{P}[X_5, X_3]^{-1} \mathcal{P}[X_5, X_{\{3,4\}}]$$

$$\begin{aligned} & \mathcal{P}[X_1, X_2, X_3, X_4, X_5] \\ = & \underline{\mathcal{P}[X_1, X_2, X_3, X_5]} \times_{X_3} \mathcal{P}[X_5, X_3]^{-1} \times_{X_5} \mathcal{P}[X_3, X_4, X_5] \end{aligned}$$

Decompose this recursively

Constructing a Different Factorization



$$\mathcal{P}[X_{\{1,2\}}, X_{\{3,5\}}] = \mathcal{P}[X_{\{1,2\}}, X_3] \mathcal{P}[X_1, X_3]^{-1} \mathcal{P}[X_1, X_{\{3,5\}}]$$

$$\begin{aligned} & \mathcal{P}[X_1, X_2, X_3, X_5] \\ = & \mathcal{P}[X_1, X_2, X_3] \times_{X_3} \mathcal{P}[X_1, X_3]^{-1} \times_{X_1} \mathcal{P}[X_1, X_3, X_5] \end{aligned}$$



Our Observable Factorization

$$\begin{aligned} & \mathcal{P}[X_1, X_2, X_3, X_4, X_5, X_6] \\ = & \mathcal{P}[X_1, X_2, X_3, X_4, X_5] \times_{X_5} \mathcal{P}[X_4, X_5]^{-1} \times_{X_4} \mathcal{P}[X_4, X_5, X_6] \end{aligned}$$

$$\begin{aligned} & \mathcal{P}[X_1, X_2, X_3, X_4, X_5] \\ = & \mathcal{P}[X_1, X_2, X_3, X_5] \times_{X_3} \mathcal{P}[X_5, X_3]^{-1} \times_{X_5} \mathcal{P}[X_3, X_4, X_5] \end{aligned}$$

$$\begin{aligned} & \mathcal{P}[X_1, X_2, X_3, X_5] \\ = & \mathcal{P}[X_1, X_2, X_3] \times_{X_3} \mathcal{P}[X_1, X_3]^{-1} \times_{X_1} \mathcal{P}[X_1, X_3, X_5] \end{aligned}$$

All Third Order Tensors

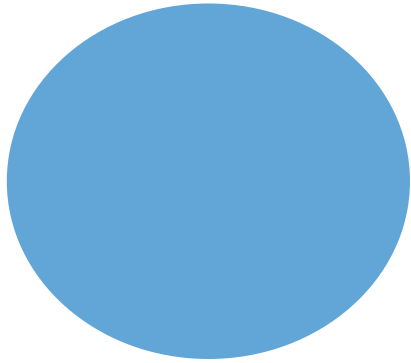
Observable Factorization



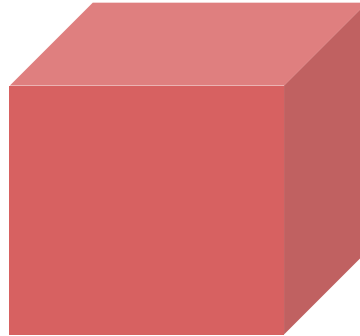
$\mathcal{P}[X_1, X_2, X_3, X_4, X_5, X_6]$

$\mathcal{P}[X_1, X_3, X_5]$

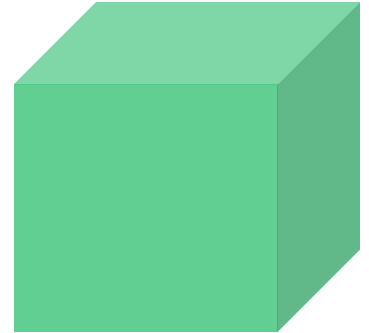
$\mathcal{P}[X_1, X_2, X_3] \times_{X_3} \mathcal{P}[X_1, X_3]^{-1}$



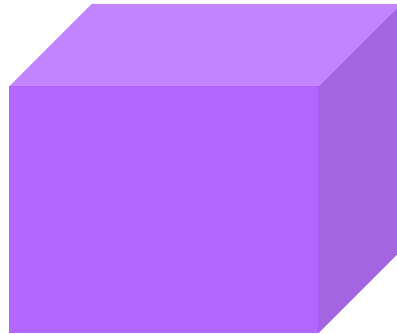
$=$



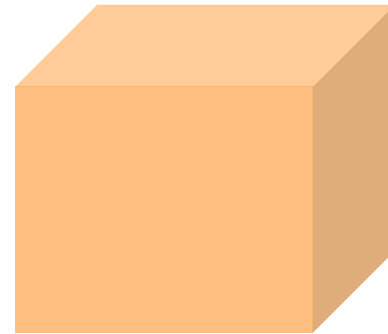
\times



\times



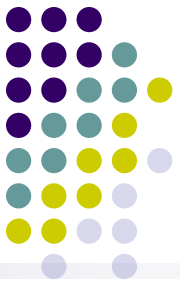
\times



$\mathcal{P}[X_5, X_3]^{-1} \times_{X_5} \mathcal{P}[X_3, X_4, X_5]$

$\mathcal{P}[X_4, X_5]^{-1} \times_{X_4} \mathcal{P}[X_4, X_5, X_6]$

Can We Decompose the Third Order Tensors Further?



$$\mathcal{P}[X_1, X_3, X_5] \quad ?$$

Not to lower order tensors. Multiplying second order tensors doesn't increase the tensor order.



The Observable Factorization

Root

$$\mathcal{P}(\emptyset_3 H_1)$$



$$\mathcal{P}[X_1, X_3, X_5]$$

Internal Nodes
(Example)

$$\mathcal{P}(\emptyset H_2 | H_1)$$



$$\mathcal{P}[X_1, X_2, X_3] \times_{X_3} \mathcal{P}[X_1, X_3]^{-1}$$

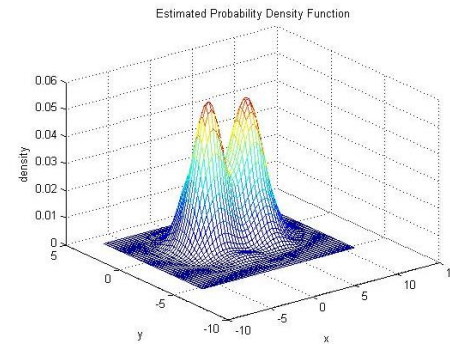
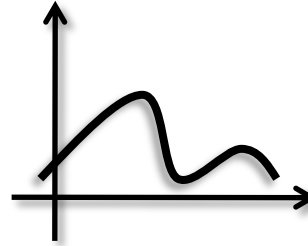
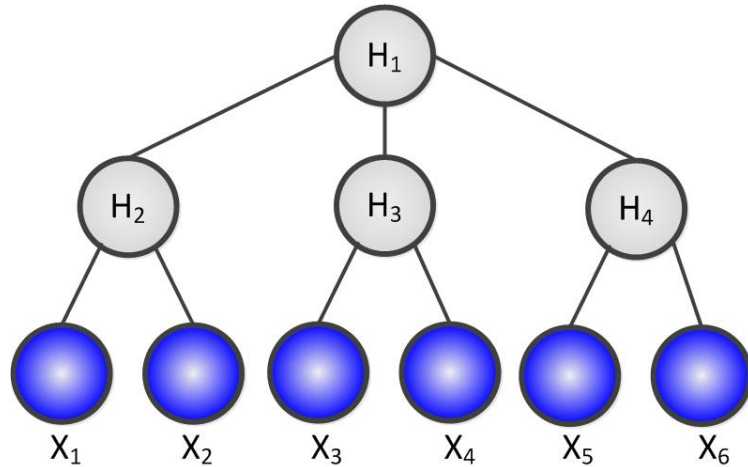
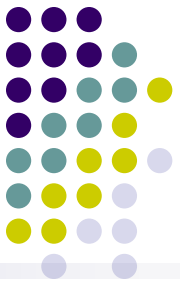
Leaves (Example)

$$\mathcal{P}(X_1 | H_2)$$



Identity Matrix

How To Deal with (Nonparametric) Continuous Variables?



Hilbert Space Embeddings!!!!

How To Deal with (Nonparametric) Continuous Variables?



Root

$$\mathcal{P}(\emptyset_3 H_1)$$



$$\mathcal{C}_{1,3,5}$$

Internal Nodes
(Example)

$$\mathcal{P}(\emptyset H_2 | H_1)$$



$$\mathcal{C}_{1,2,3} \times_3 \mathcal{C}_{1,3}^{-1}$$

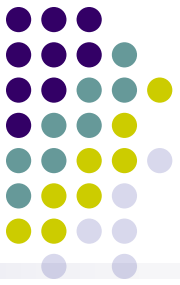
Leaves (Example)

$$\mathcal{P}(X_1 | H_2)$$



Identity Operator

(something very similar to this is true)



Comparing the Factorizations

Traditional (CPT)

**Requires EM to compute
(local minima, slow)**

**No theoretical
guarantees**

Aims to Find MLE

**Allows for inference
among latent variables**

Observable

Very fast, local minima free

Provably consistent

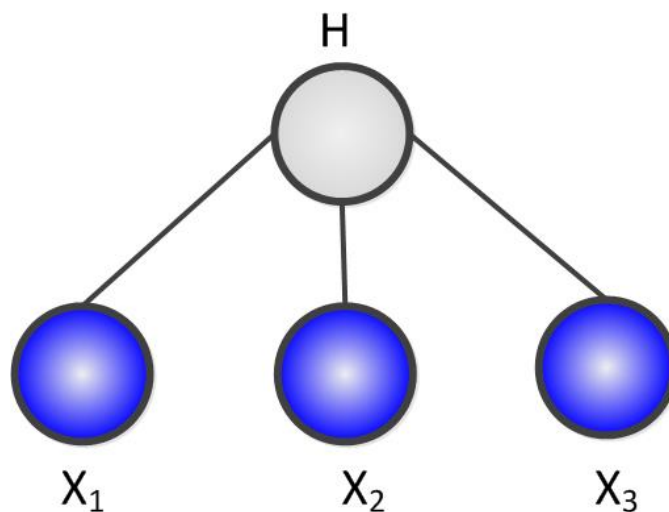
**Less statistically efficient
(does not aim to find MLE)**

**Does not allow for inference
among latent variables**

What If I Want to Extract the Actual Latent Parameters?

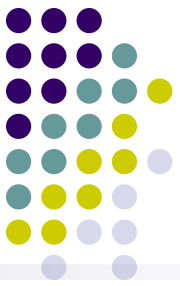


- Let's consider an even simpler example.



- Assume all the conditional probability tables are the same i.e.

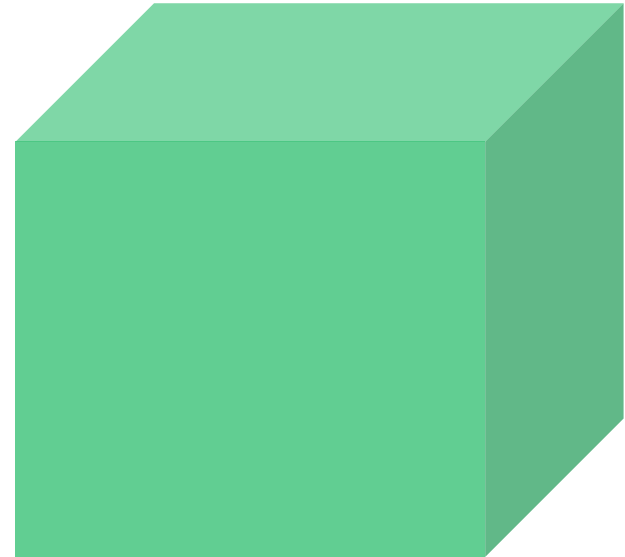
$$O := \mathcal{P}[X_1|H] = \mathcal{P}[X_2|H] = \mathcal{P}[X_3|H]$$



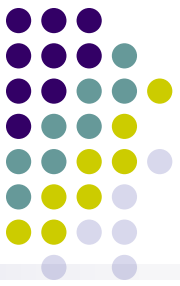
Extracting Latent Parameters

- We only have three observed variables, so one observable parameterization is just their joint:

$$\mathcal{P}[X_3, X_2, X_1]$$



$\mathcal{P}[X_3, x_2, X_1]$ is a slice of the tensor



Extracting Latent Parameters

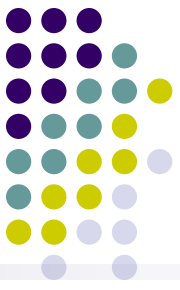
$$\mathcal{P}[X_3, x_2, X_1] = \mathbf{O} \text{diag}(O_{x_2}) P(\otimes H_1) \mathbf{O}^\top$$



I want to extract this....

Looks kind of like an eigenvalue decomposition, but not quite.....

$$\mathbf{M} = \mathbf{S} \mathbf{\Lambda} \mathbf{S}^{-1}$$



What about the Related Quantity?

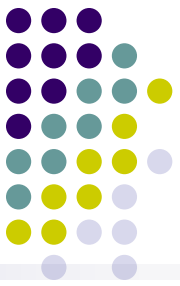
$$\mathcal{P}[X_3, x_2, X_1] = \mathbf{O} \text{diag}(O_{x_2}) P(\emptyset H_1) \mathbf{O}^\top$$

$$\mathcal{P}[X_3, x_2, X_1] = \mathbf{O} \text{diag}(O_{x_2}) \mathbf{O}^{-1} \mathbf{O} P(\emptyset H_1) \mathbf{O}^\top$$



This looks better!!!!

$$\mathbf{B}_{x_2} := \mathbf{O} \text{diag}(O_{x_2}) \mathbf{O}^{-1} = \mathcal{P}[X_3, x_2, X_1] \mathcal{P}[X_2, X_1]^{-1}$$

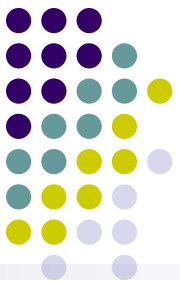


Eigenvalue Decomposition

Run Eigenvalue Decomposition on this

$$B_{x_2} := \mathbf{O} \operatorname{diag}(O_{x_2}) \mathbf{O}^{-1} = \mathcal{P}[X_3, x_2, X_1] \mathcal{P}[X_2, X_1]^{-1}$$

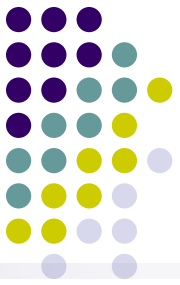
Eigenvectors correspond to columns of \mathbf{O}
(*need to be normalized*)



Comments

- Same technique is also the basis of recent advances in spectral LDA (Anandkumar et al. 2012)
- Not stable in practice in its naïve form.
- Anandkumar, Hsu, Kakade (2012) recently propose a tensor power iteration method that works better in practice (still current area of research)

Case Study: Supervised Parsing



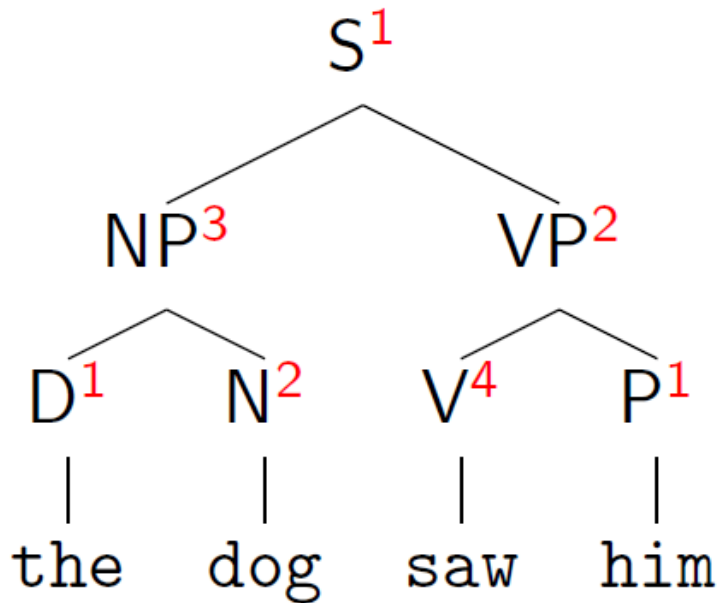
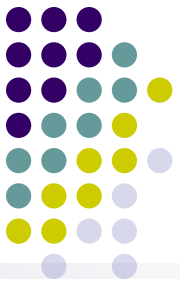
- Cohen et al. 2012 derived a spectral learning method for latent PCFGs.
- Latent PCFGs (Matsuzaki et al, Petrov et al.) are Probabilistic Context Free Grammars with additional latent nodes.
- The challenge in their setting is that the tree structure changes on every example.

Case Study: Supervised Parsing



- Cohen et al. 2012 derived a spectral learning method for latent PCFGs.
- Latent PCFGs (Matsuzaki et al, Petrov et al.) are Probabilistic Context Free Grammars with additional latent nodes.
- The challenge in their setting is that the tree structure changes on every example.

Latent PCFG Supervised Parsing

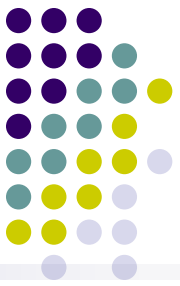


$$\begin{aligned} & p(\text{tree}, 1\ 3\ 1\ 2\ 2\ 4\ 1) \\ &= \pi(S^1) \times \\ & \quad t(S^1 \rightarrow NP^3\ VP^2 | S^1) \times \\ & \quad t(NP^3 \rightarrow D^1\ N^2 | NP^3) \times \\ & \quad t(VP^2 \rightarrow V^4\ P^1 | VP^2) \times \\ & \quad q(D^1 \rightarrow \text{the} | D^1) \times \\ & \quad q(N^2 \rightarrow \text{dog} | N^2) \times \\ & \quad q(V^4 \rightarrow \text{saw} | V^4) \times \\ & \quad q(P^1 \rightarrow \text{him} | P^1) \end{aligned}$$

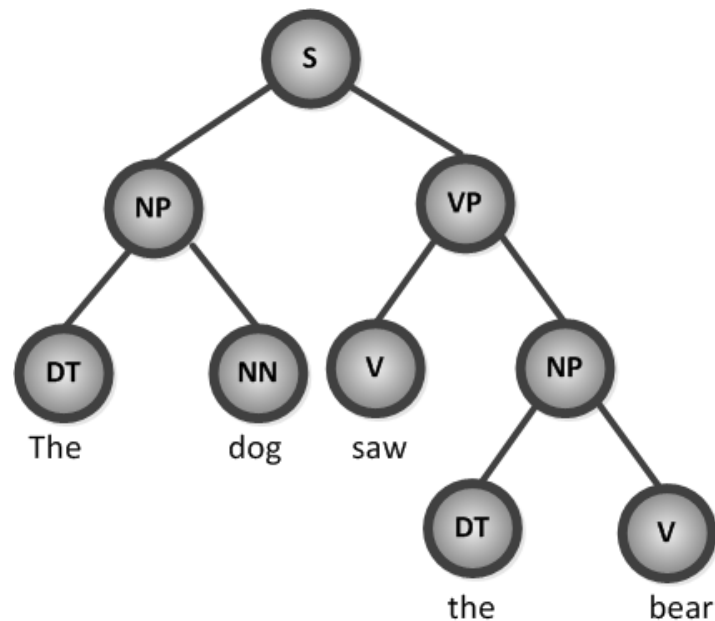
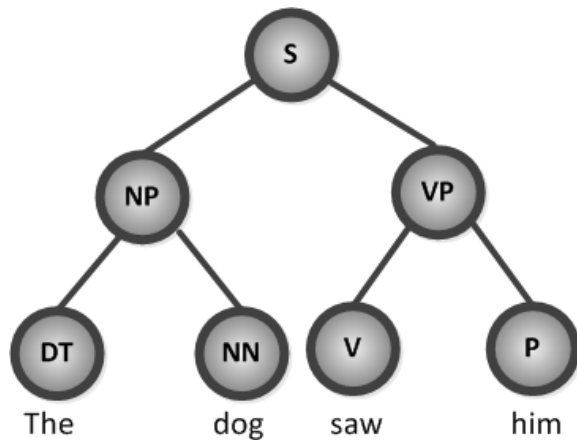
$$p(\text{tree}) = \sum_{h_1 \dots h_7} p(\text{tree}, h_1\ h_2\ h_3\ h_4\ h_5\ h_6\ h_7)$$

(from Karl Stratos ACL 2012 presentation)

Different Examples Have Different Trees



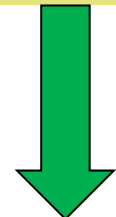
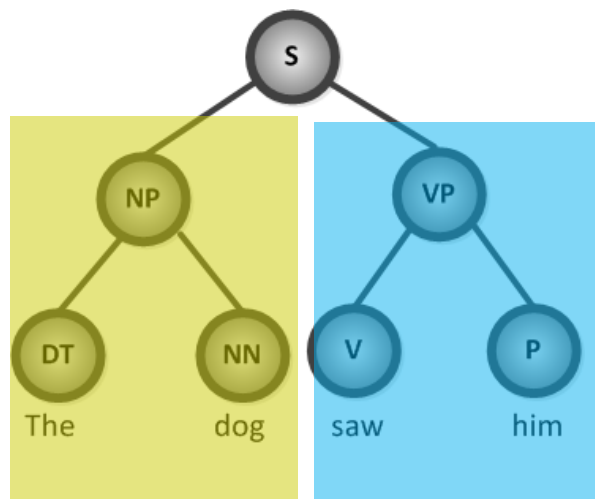
- If all examples had the same tree then the method would be very similar to spectral learning for trees.
- However, different sentences have different parse trees.



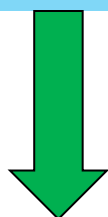


Construct Features of Subtrees

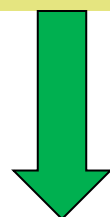
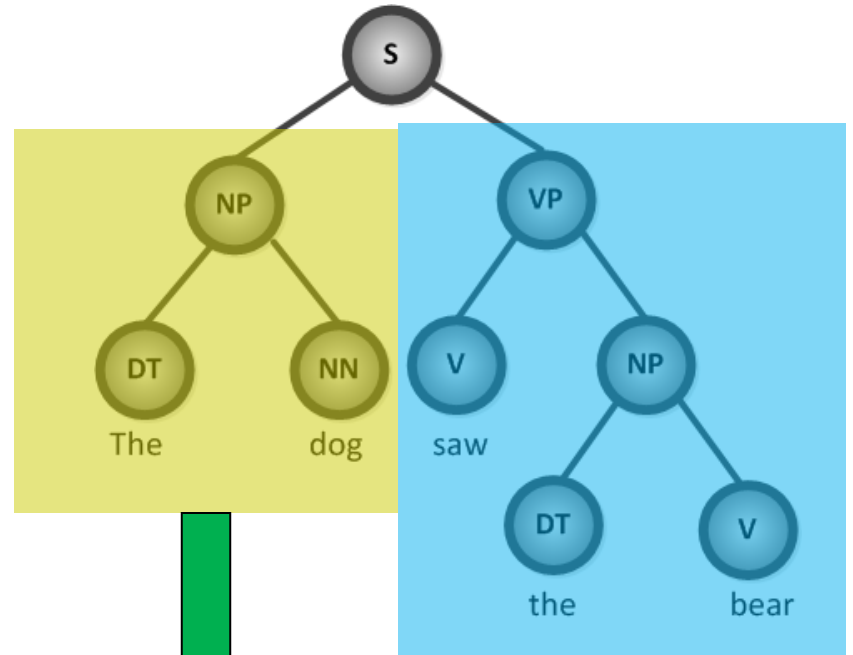
- Consider the root



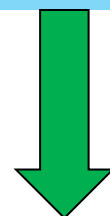
$$\phi_L^{(1)}$$



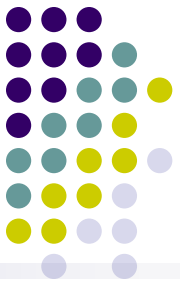
$$\phi_R^{(1)}$$



$$\phi_L^{(2)}$$



$$\phi_R^{(2)}$$



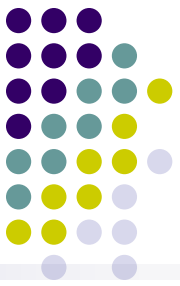
Use Features of Subtrees

- In our previous case the parameter for the root looked something like this (one observation from each subtree):

$$\mathcal{P}[X_1, X_3, X_5]$$

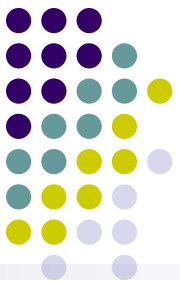
- For Cohen et al. 2012, it will look something like

$$\mathbb{E}[\phi_L \phi_R^\top]$$



Empirical Results

- (Results are unpublished – will be presented at NAACL 2013).
- In general, the algorithm performs comparably with EM in accuracy.
- However, it is around 20x faster.



References

- Song, L., Ishteva, M. Parikh, A.P., Park, H. Xing, E.P. **Under Review (2013)**
- Mossel, E. and Roch, S. **Learning nonsingular phylogenies and hidden markov models**. Annals of Applied Probability, 16(2):583–614, 2006.
- Hsu, D., Kakade, S., and Zhang, T. **A Spectral Algorithm for Learning Hidden Markov Models**. Conference on Learning Theory, 2009.
- Siddiqi, S., Boots, B. Gordon, G., **Reduced Rank Hidden Markov Models**, Artificial Intelligence and Statistics (AISTATS), 2009.
- Parikh, A.P. Song, L., and Xing, E.P. **A Spectral Algorithm for Latent Tree Graphical Models**, International Conference of Machine Learning (ICML), 2011.
- Song, L., Boots, B., Siddiqi, S., Gordon, G., and Smola, A. **Hilbert space embeddings of Hidden Markov Models**. International Conference of Machine Learning (ICML), 2010..
- S. Cohen, K. Stratos, M. Collins, D. Foster, L. Ungar. **Spectral Learning of Latent-Variable PCFGs**. Association of Computational Linguistics (ACL) 2012.