



10-418/10-618 Machine Learning for Structured Data

Machine Learning Department
School of Computer Science
Carnegie Mellon University



Variational Inference

Matt Gormley
Lecture 16
Oct. 31, 2022

Q&A

Q: The parameters of a K-dimensional Dirichlet(α) are a vector α of length K, so why are Dirichlet parameters sometimes given as a scalar? For example...

“We use a Dirichlet prior with parameter $\alpha = 0.1$.”

A: Great question!

A K-dimensional Dirichlet prior is said to be *symmetric* if all the values in the vector α are the same, i.e. for all k , $\alpha_k = c$ where c is a scalar constant.

We sometimes call this restricted version the *symmetric Dirichlet distribution*.

Reminders

- Exam Rubrics and Exam Viewings
- Homework 4: MCMC
 - Out: Mon, Oct 24
 - Due: Fri, Nov 3 at 11:59pm
- Homework 5: ~~MCMC~~ ^{VI}
 - ~~– Out: Mon, Oct 24~~
 - ~~– Due: Fri, Nov 3 at 11:59pm~~

Reminders

Happy Halloween!



SEMANTIC SEGMENTATION

Case Study: Image Segmentation

- Image segmentation (FG/BG) by modeling of interactions btw RVs
 - Images are noisy.
 - Objects occupy continuous regions in an image.

[Nowozin, Lampert 2012]



Input image



Pixel-wise separate
optimal labeling



Locally-consistent
joint optimal labeling

$$Y^* = \arg \max_{y \in \{0,1\}^n} \left[\overbrace{\sum_{i \in S} V_i(y_i, X)}^{\text{Unary Term}} + \overbrace{\sum_{i \in S} \sum_{j \in N_i} V_{i,j}(y_i, y_j)}^{\text{Pairwise Term}} \right].$$

© Eric Xing @ CMU, 2005-2015

Y : labels

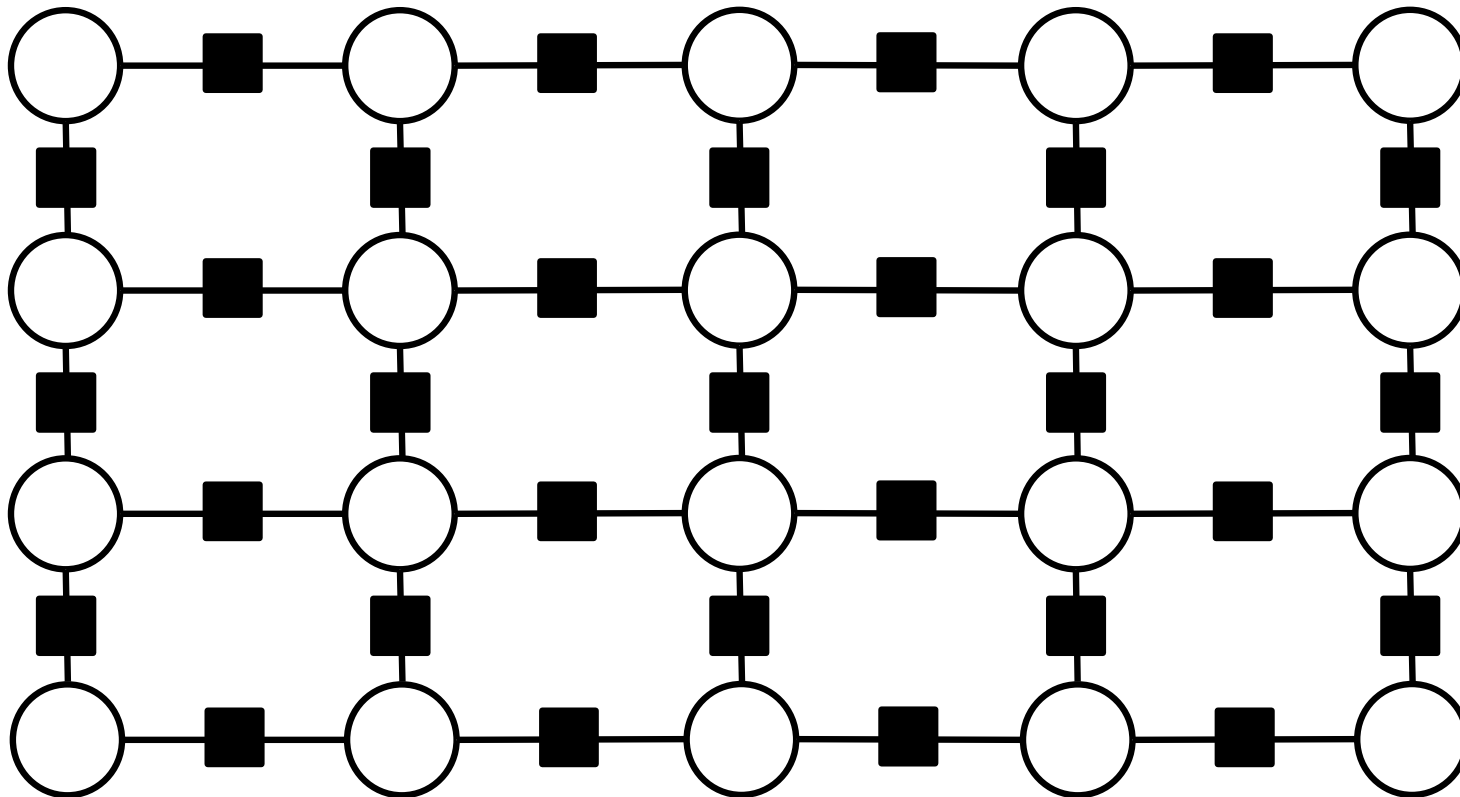
X : data (features)

S : pixels

N_i : neighbors of pixel i

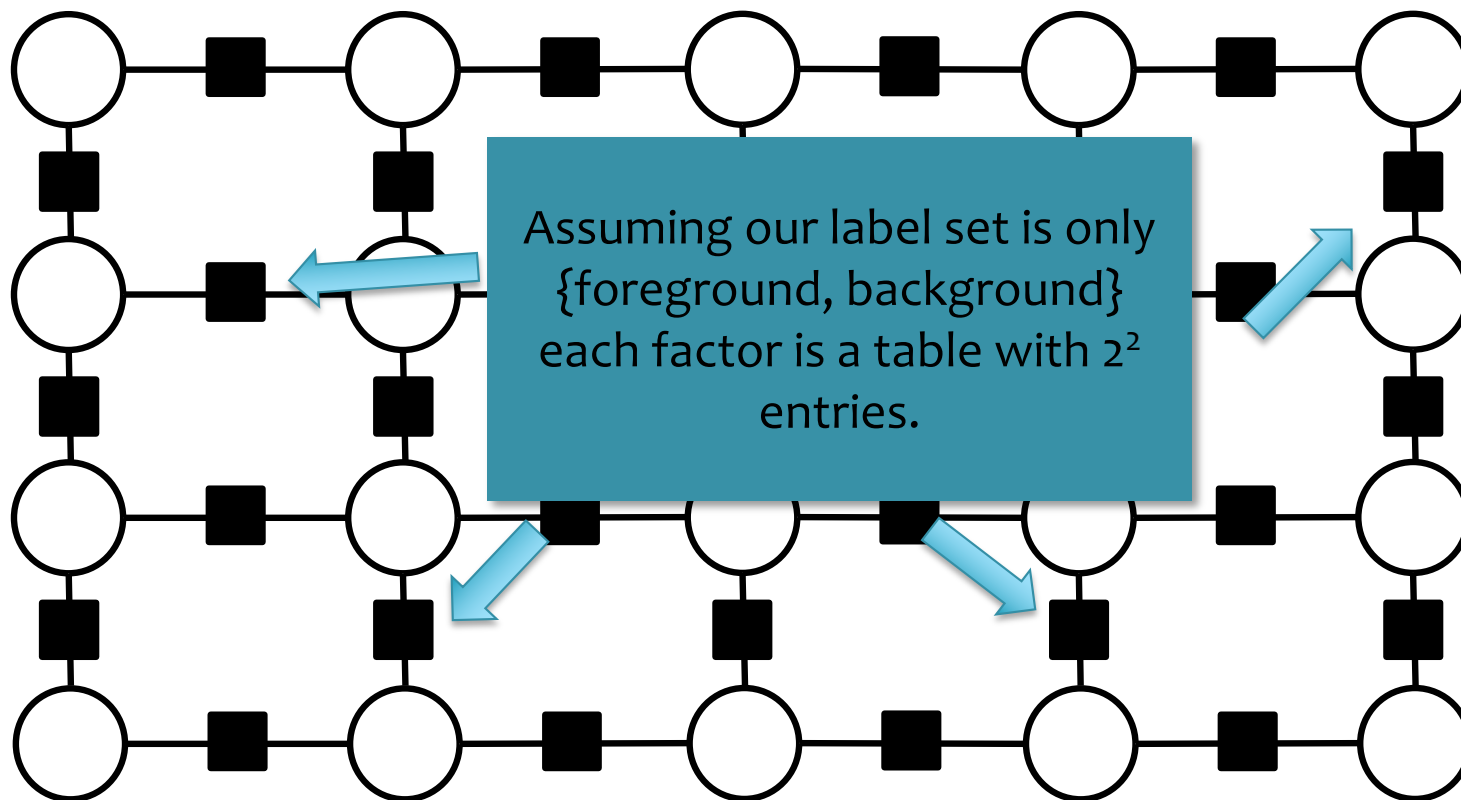
Grid CRF

- Suppose we want to image segmentation using a grid model



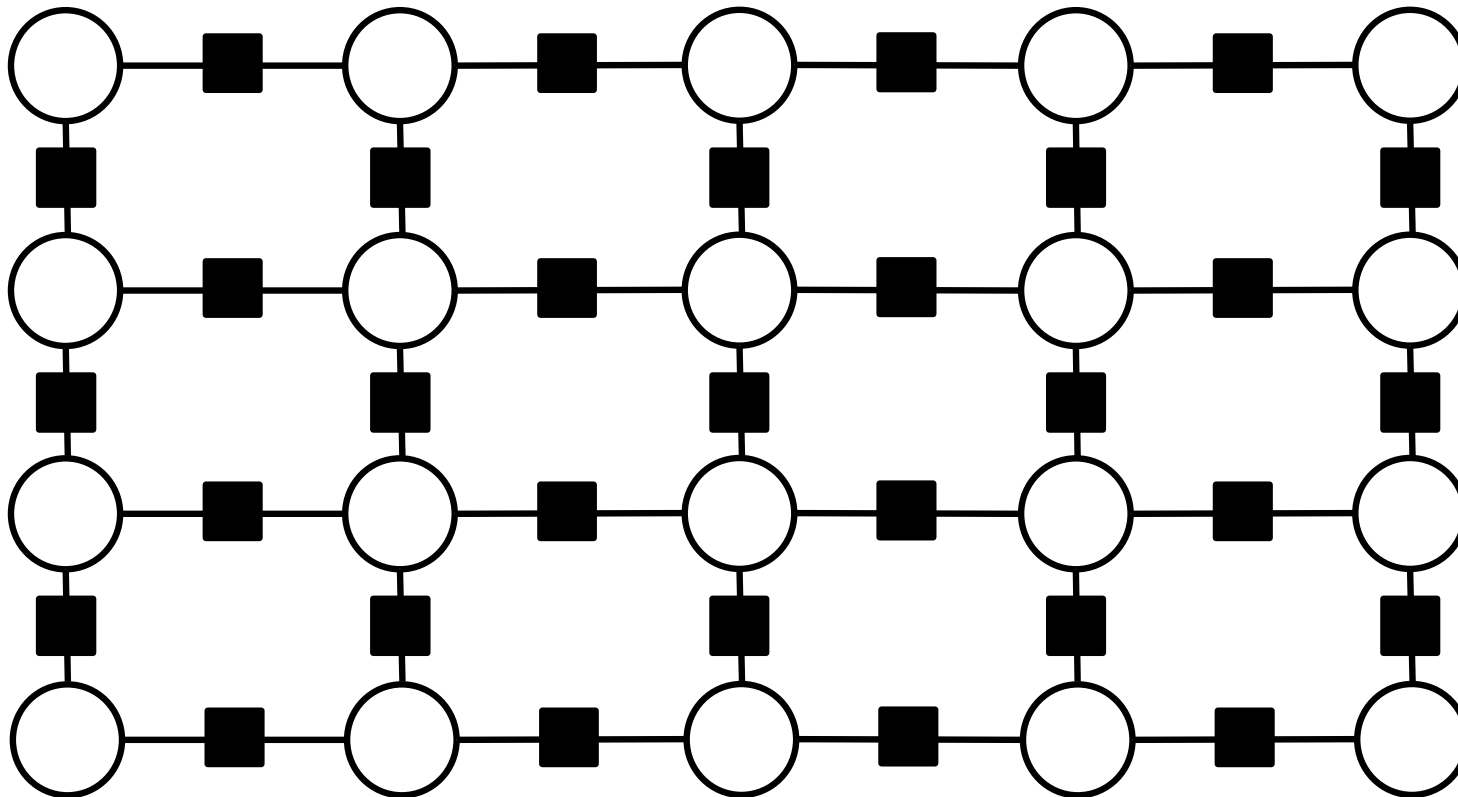
Grid CRF

- Suppose we want to image segmentation using a grid model



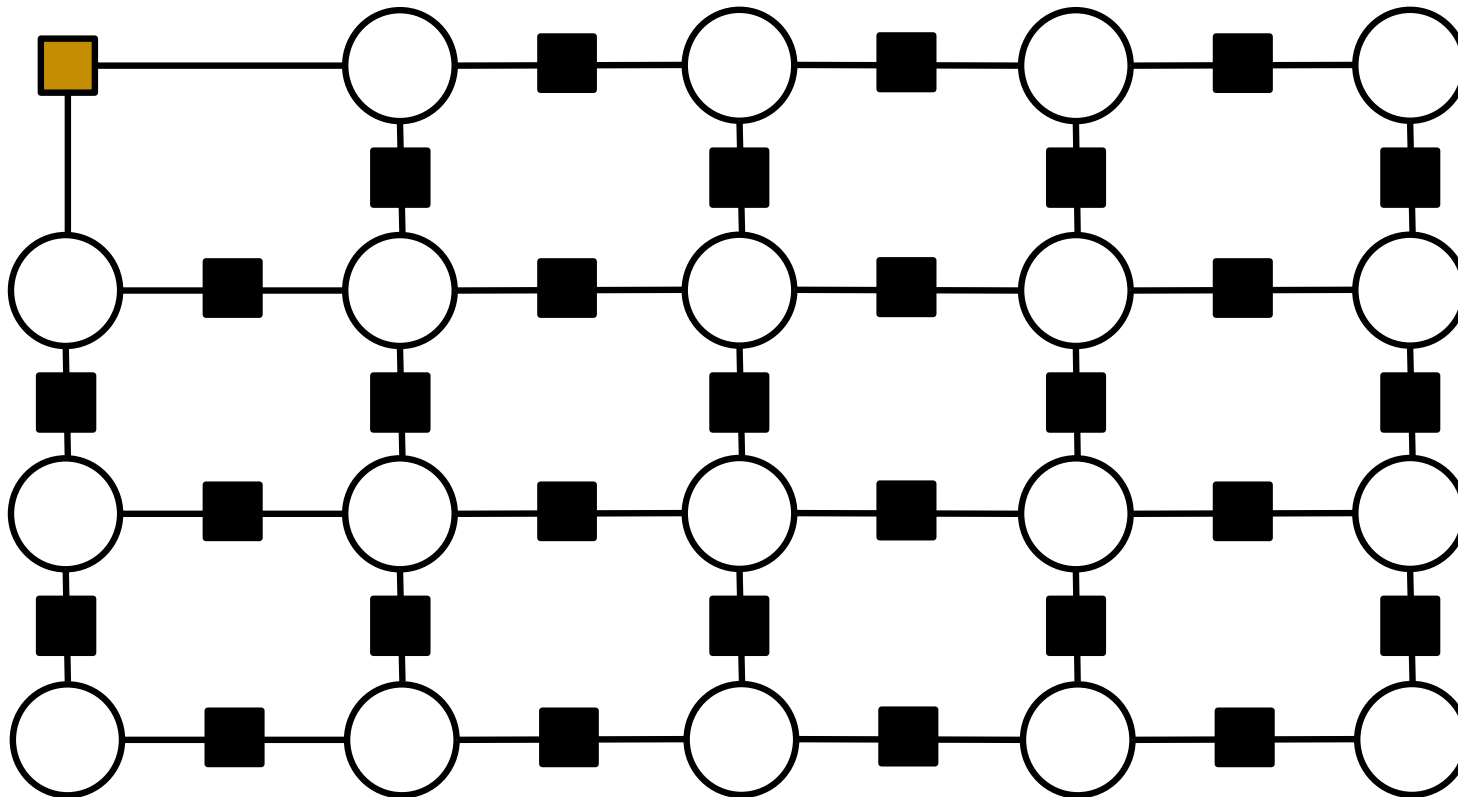
Grid CRF

- Suppose we want to image segmentation using a grid model
- What happens when we run variable elimination?



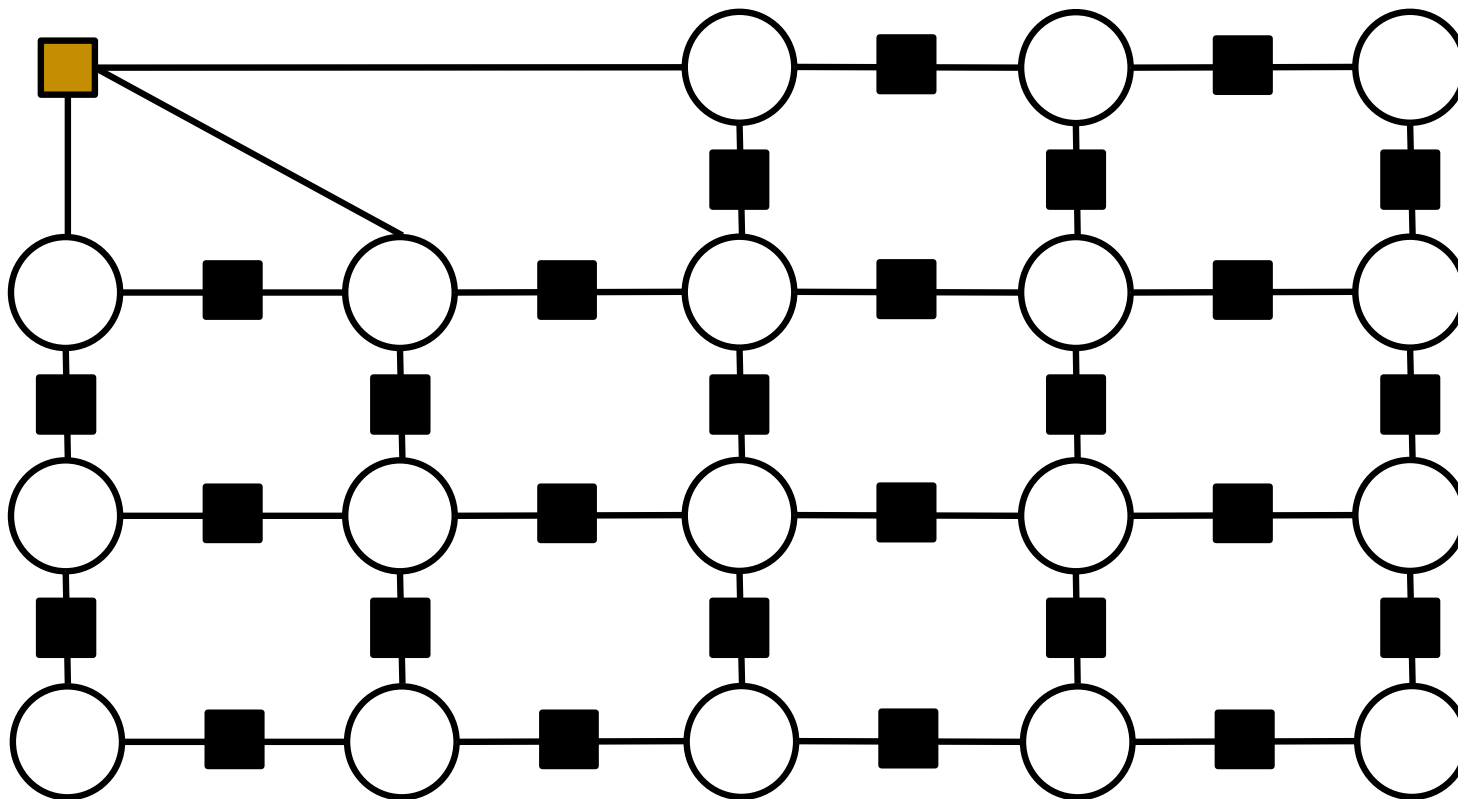
Grid CRF

- Suppose we want to image segmentation using a grid model
- What happens when we run variable elimination?



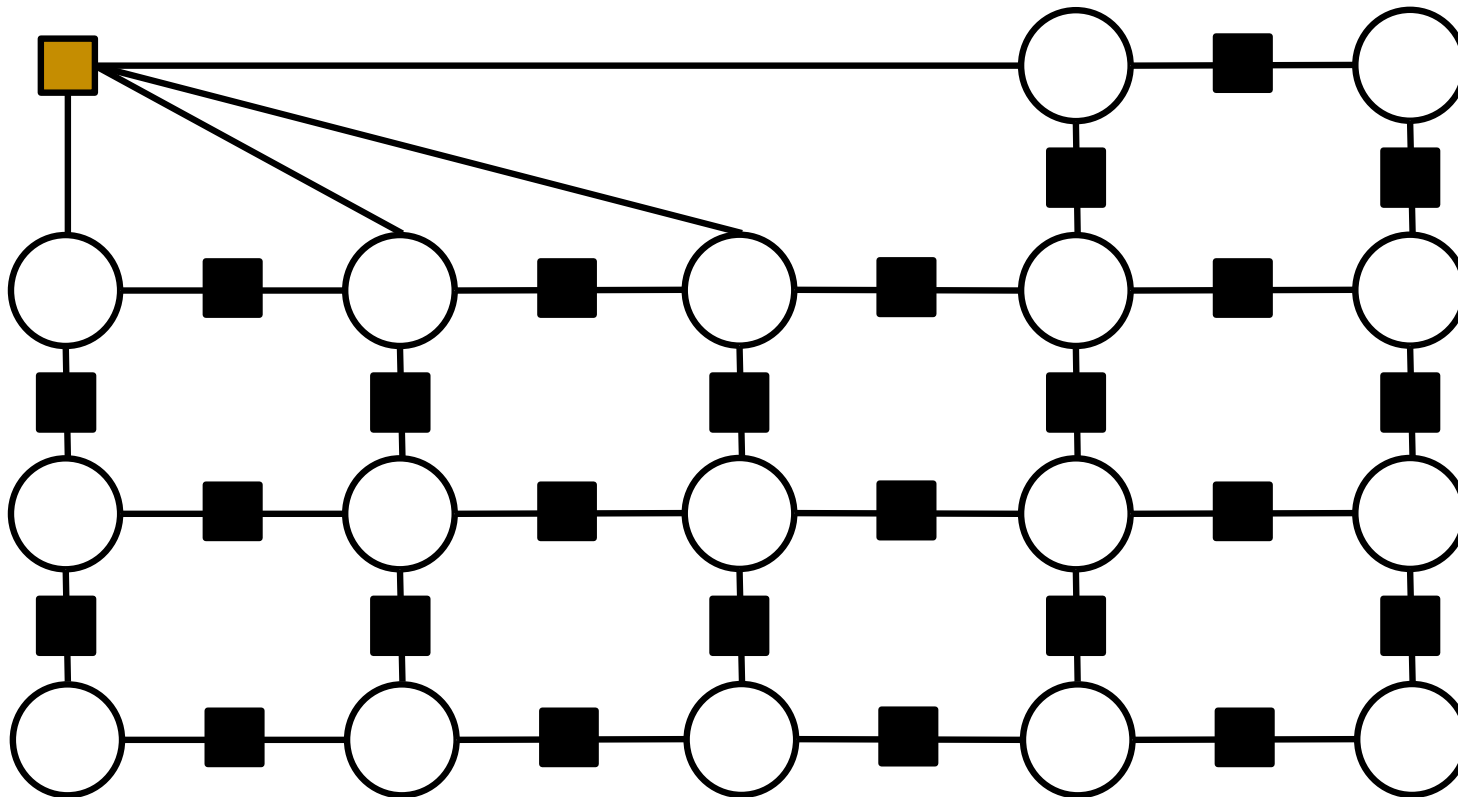
Grid CRF

- Suppose we want to image segmentation using a grid model
- What happens when we run variable elimination?



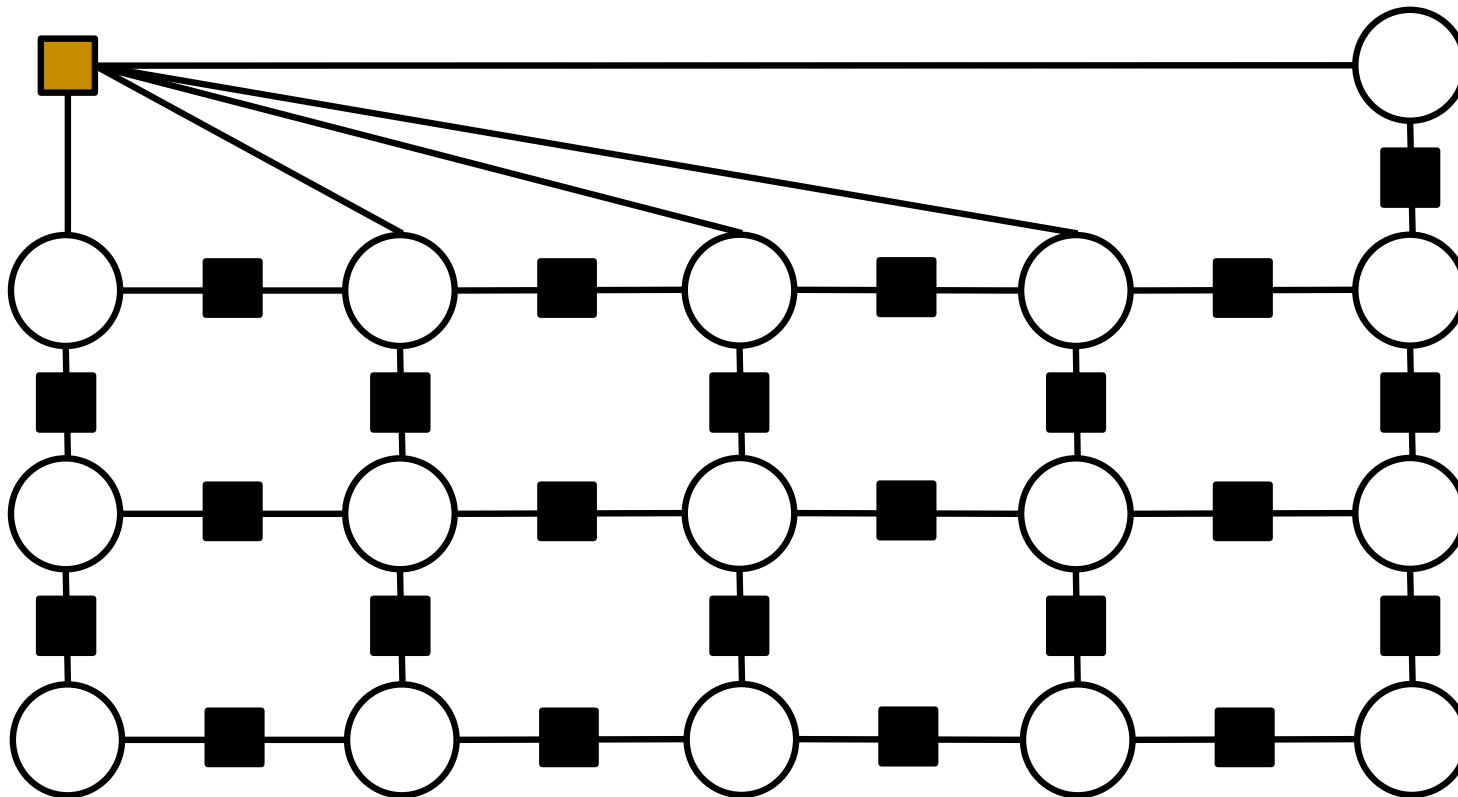
Grid CRF

- Suppose we want to image segmentation using a grid model
- What happens when we run variable elimination?



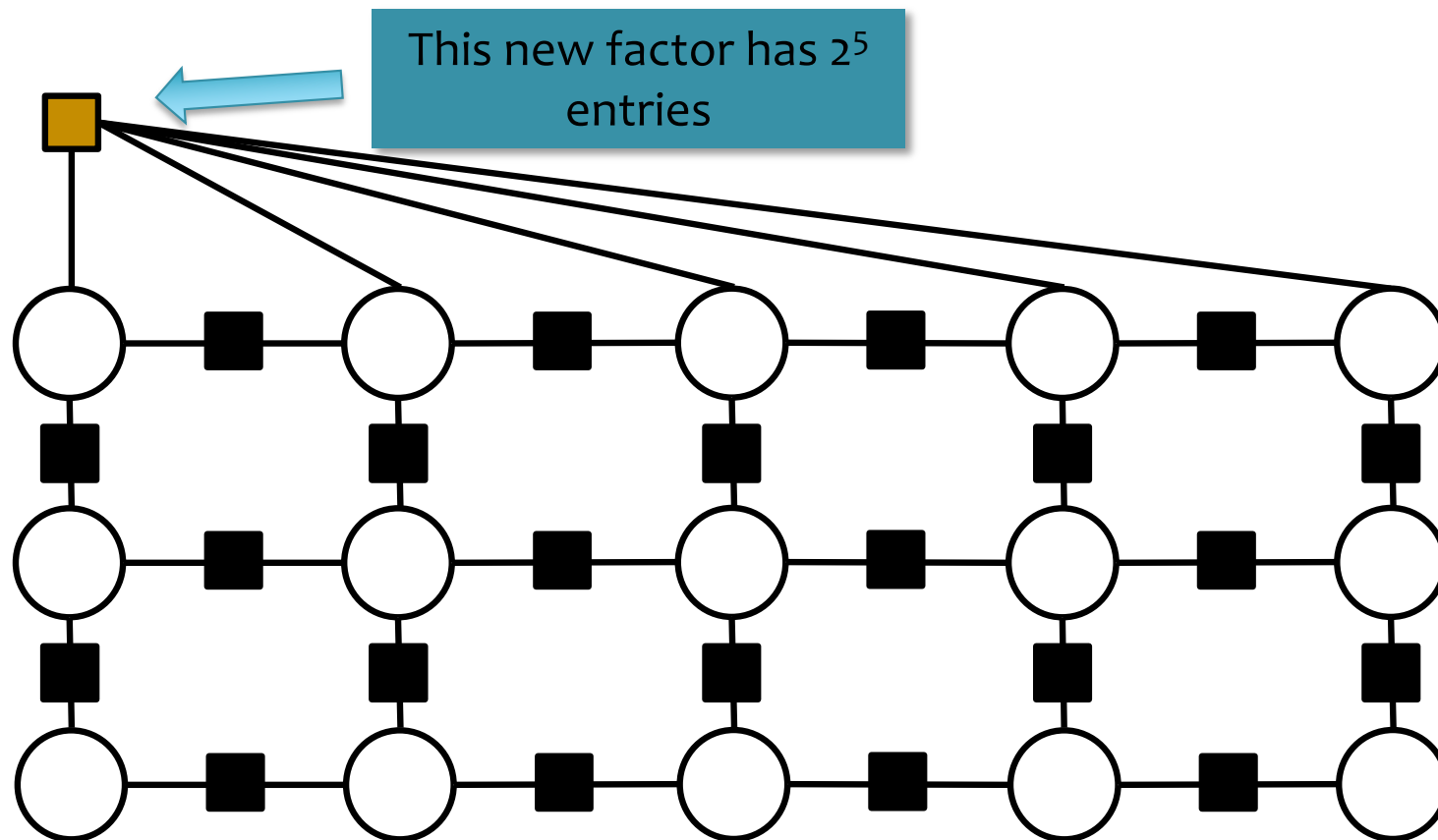
Grid CRF

- Suppose we want to image segmentation using a grid model
- What happens when we run variable elimination?



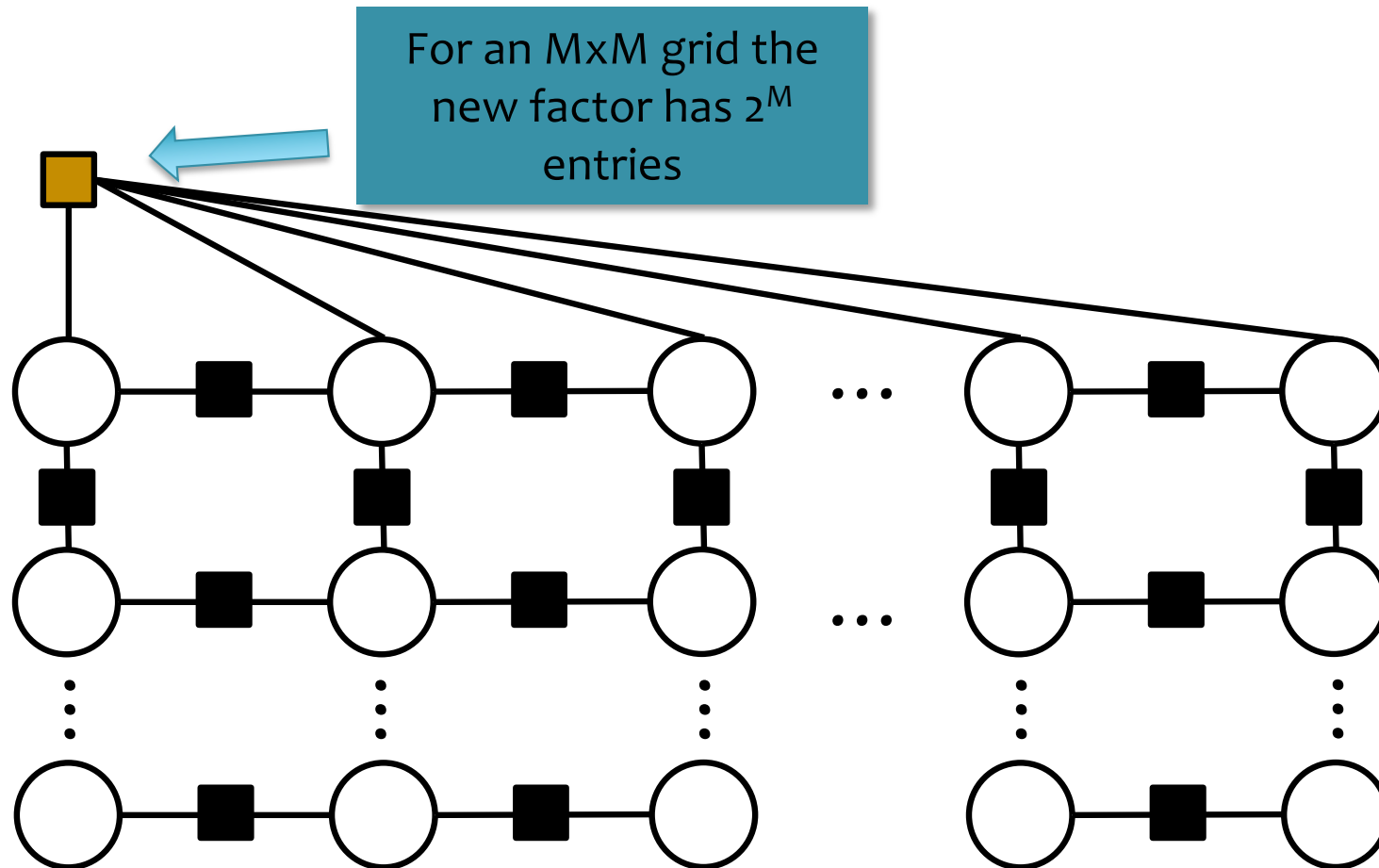
Grid CRF

- Suppose we want to image segmentation using a grid model
- What happens when we run variable elimination?



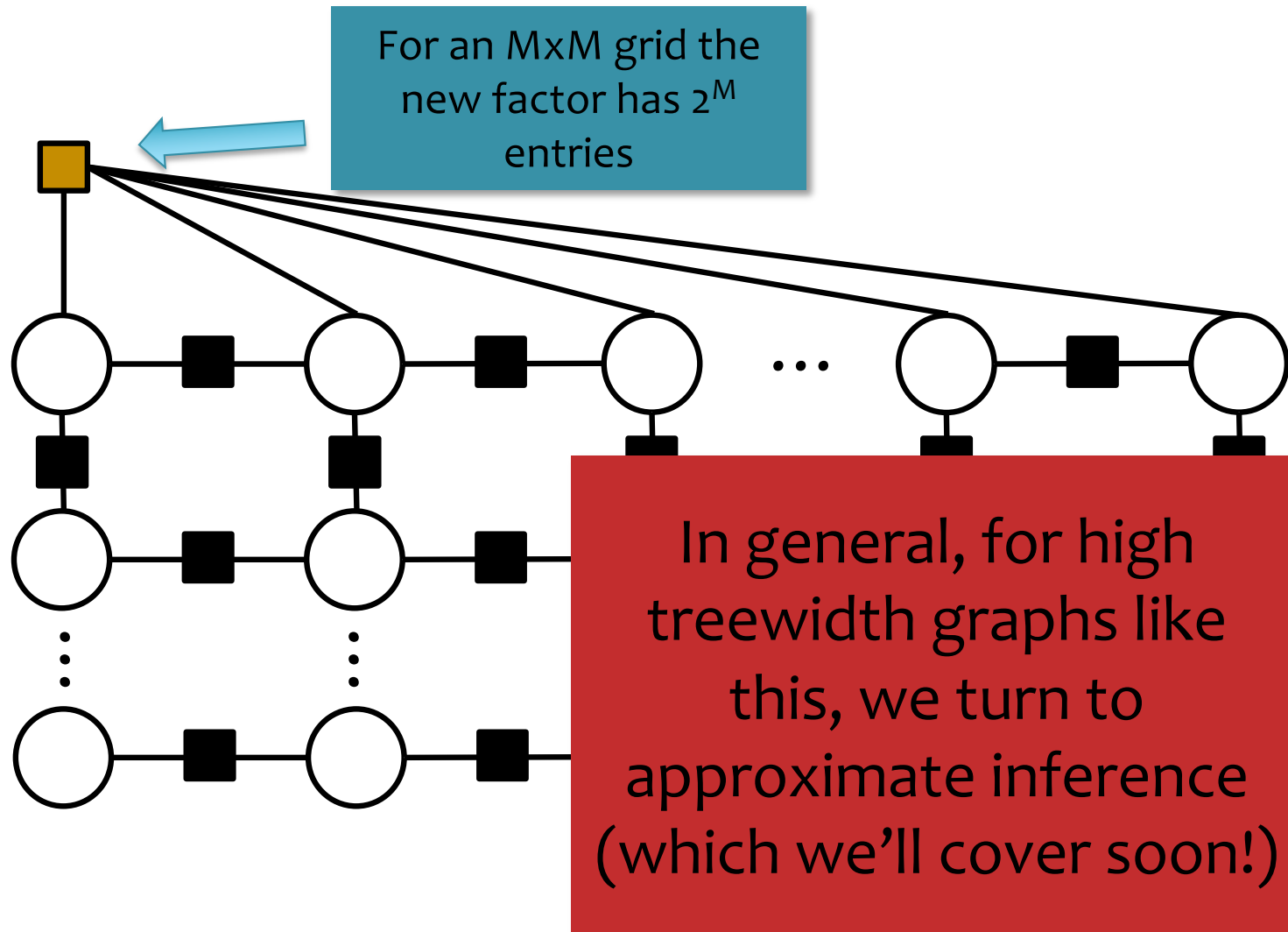
Grid CRF

- Suppose we want to image segmentation using a grid model
- What happens when we run variable elimination?



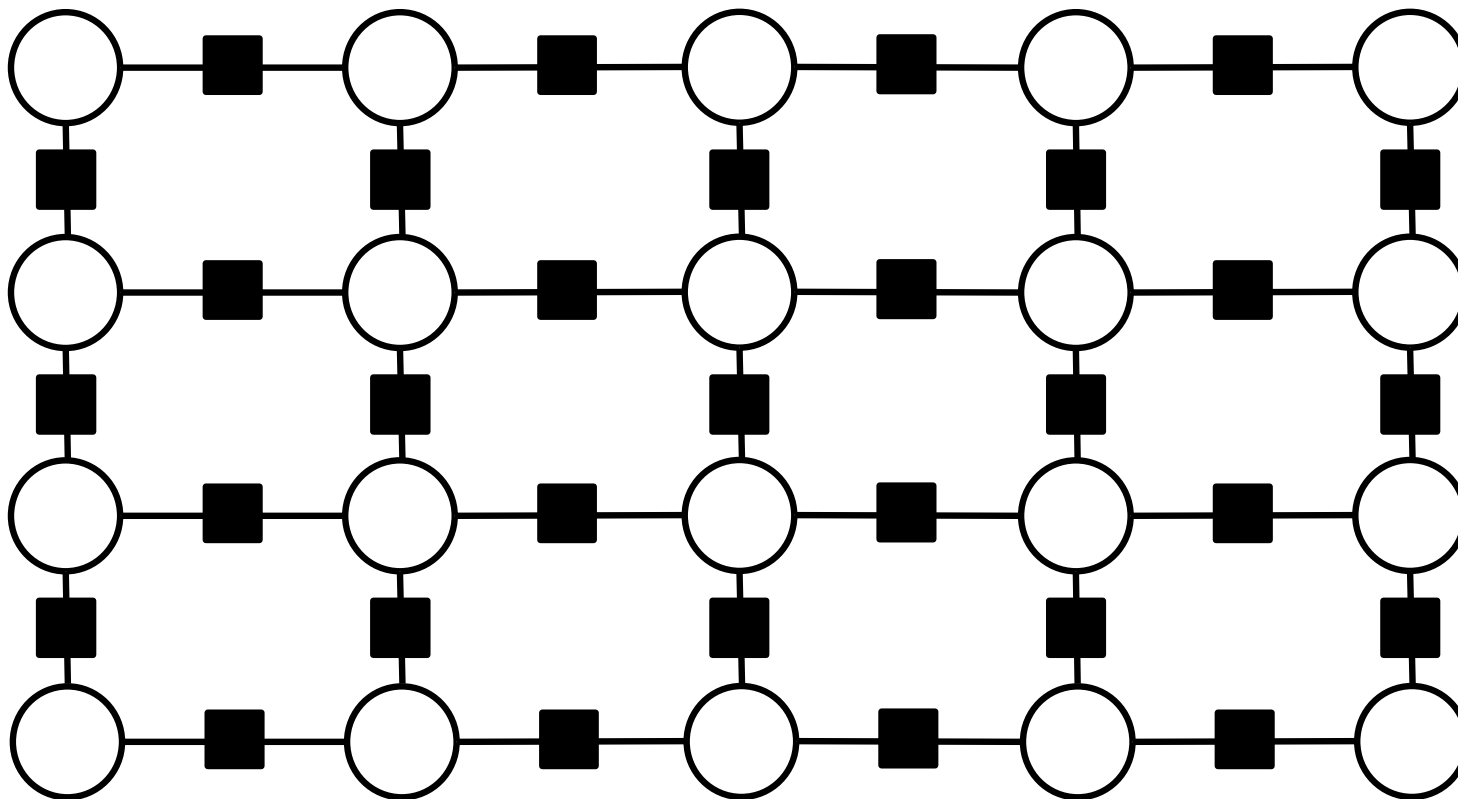
Grid CRF

- Suppose we want to image segmentation using a grid model
- What happens when we run variable elimination?



Grid CRF

- Suppose we want to image segmentation using a grid model
- What happens when we run variable elimination?
- Can we instead run belief propagation to do exact inference?

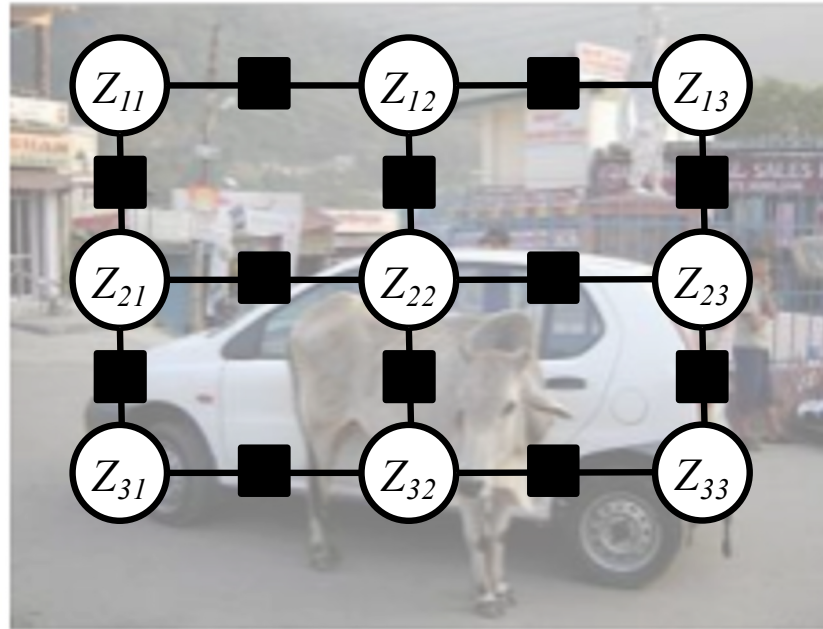


HIGH-LEVEL INTRO TO VARIATIONAL INFERENCE

Variational Inference

Problem:

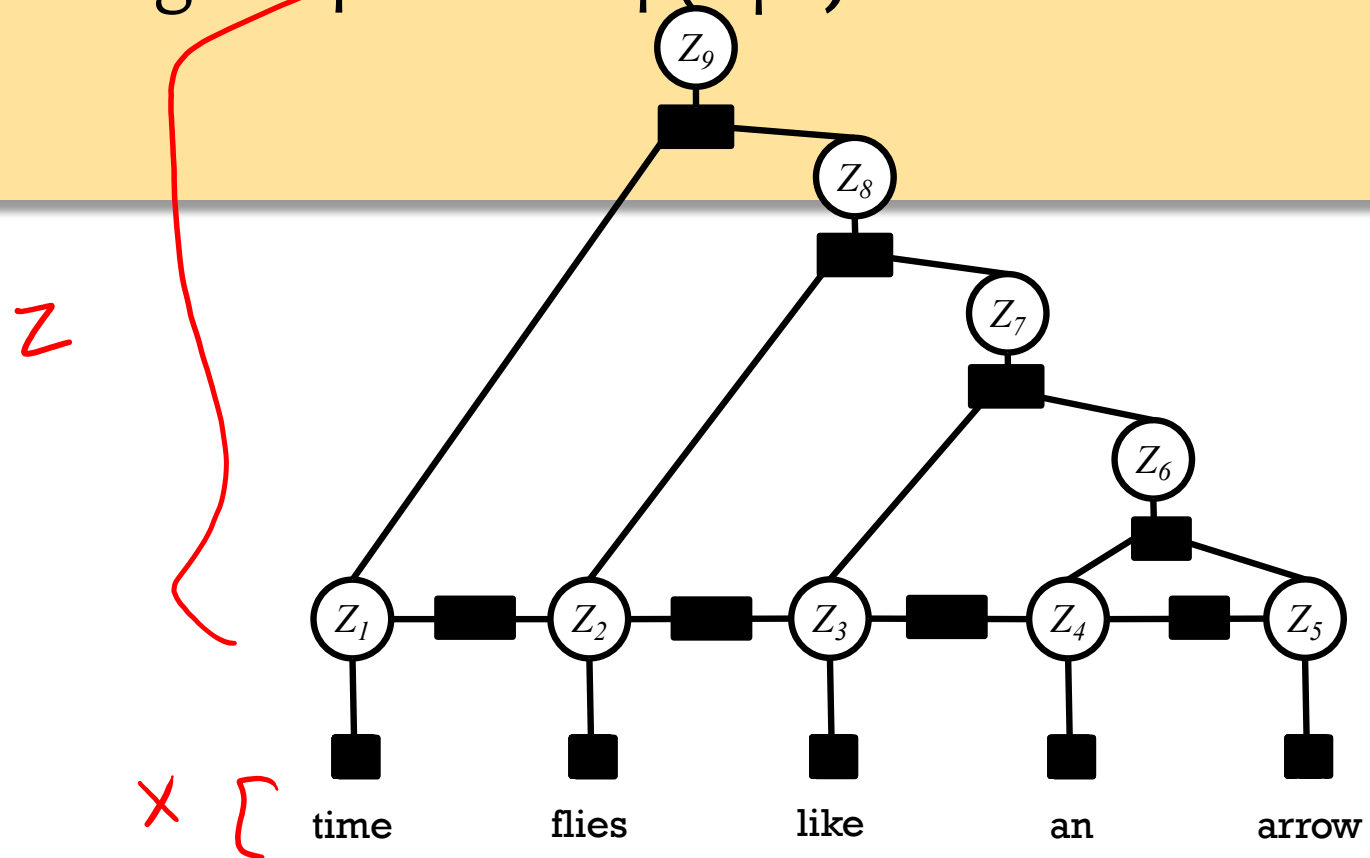
- For observed variables \mathbf{x} and latent variables \mathbf{z} , estimating the posterior $p(\mathbf{z} \mid \mathbf{x})$ is intractable



Variational Inference

Problem:

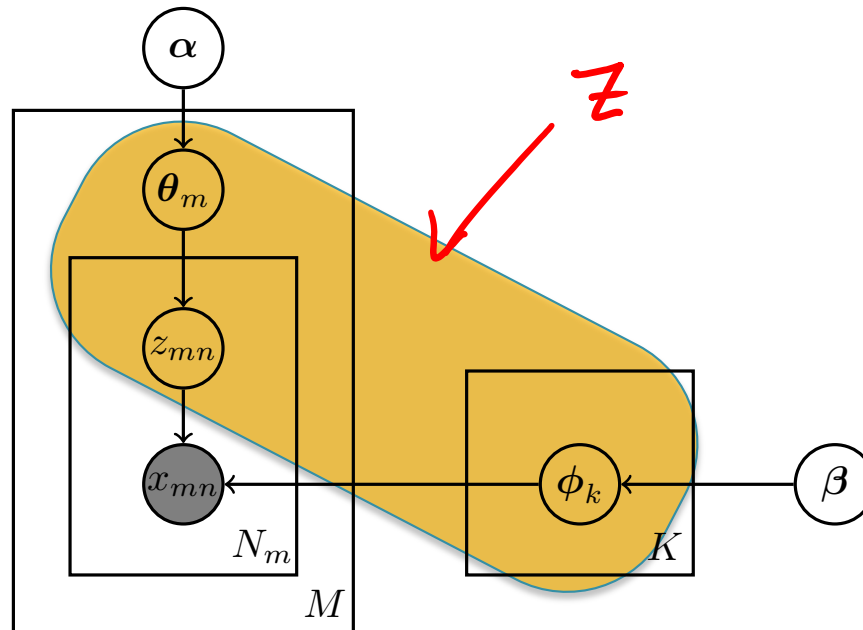
- For observed variables \mathbf{x} and latent variables \mathbf{z} , estimating the posterior $p(\mathbf{z} \mid \mathbf{x})$ is intractable



Variational Inference

Problem:

- For observed variables \mathbf{x} and latent variables \mathbf{z} , estimating the posterior $p(\mathbf{z} \mid \mathbf{x})$ is intractable
- For training data \mathbf{x} and parameters \mathbf{z} , estimating the posterior $p(\mathbf{z} \mid \mathbf{x})$ is intractable



Variational Inference

Problem:

- For observed variables \mathbf{x} and latent variables \mathbf{z} , estimating the posterior $p(\mathbf{z} \mid \mathbf{x})$ is intractable
- For training data \mathbf{x} and parameters \mathbf{z} , estimating the posterior $p(\mathbf{z} \mid \mathbf{x})$ is intractable

Solution:

- Approximate $p(\mathbf{z} \mid \mathbf{x})$ with a simpler $q(\mathbf{z})$
- Typically $q(\mathbf{z})$ has more independence assumptions than $p(\mathbf{z} \mid \mathbf{x})$ – fine b/c $q(\mathbf{z})$ is tuned for a specific \mathbf{x}
- **Key idea:** pick a single $q(\mathbf{z})$ from some family Q that best approximates $p(\mathbf{z} \mid \mathbf{x})$

Variational Inference

Terminology:

- $q(\mathbf{z})$: the **variational approximation**
- Q : the **variational family**
- Usually $q_{\theta}(\mathbf{z})$ is parameterized by some θ called **variational parameters**
- Usually $p_{\alpha}(\mathbf{z} \mid \mathbf{x})$ is parameterized by some fixed α – we'll call them the parameters

Example Algorithms:

- mean-field variational inference
- loopy belief propagation
- tree-reweighted belief propagation
- expectation propagation

Variational Inference

Is this trivial?

- Note: We are not defining a new distribution simple $q_{\theta}(\mathbf{z} \mid \mathbf{x})$, there is one simple $q_{\theta}(\mathbf{z})$ for each $p_{\alpha}(\mathbf{z} \mid \mathbf{x})$
- Consider the MCMC equivalent of this:
 - you could draw samples $\mathbf{z}^{(i)} \sim p(\mathbf{z} \mid \mathbf{x})$
 - then train some simple $q_{\theta}(\mathbf{z})$ on $\mathbf{z}^{(1)}, \mathbf{z}^{(2)}, \dots, \mathbf{z}^{(N)}$
 - hope that the sample adequately represents the posterior for the given \mathbf{x}
- How is VI different from this?
 - VI doesn't require sampling
 - VI is fast and deterministic
 - Why? b/c we choose an objective function (KL divergence) that defines which q_{θ} best approximates p_{α} , and exploit the special structure of q_{θ} to optimize it

Variational Inference

V.I. offers a new design decision

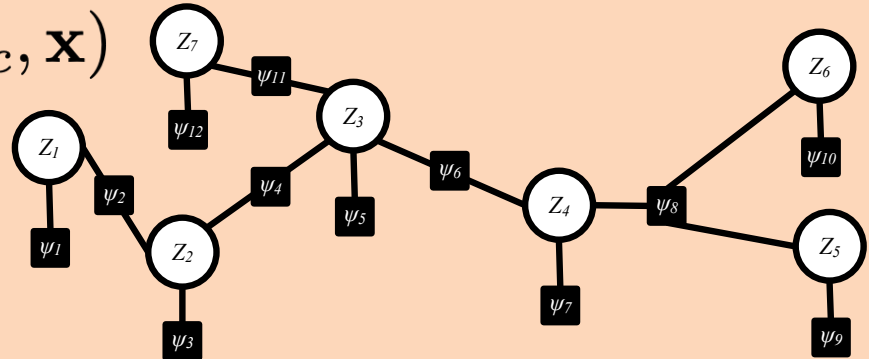
- Choose the distribution $p_{\alpha}(\mathbf{z} \mid \mathbf{x})$ that you really want, i.e. don't just simplify it to make it computationally convenient
- Then design the structure of another distribution $q_{\theta}(\mathbf{z})$ such that V.I. is efficient

TYPES OF VARIATIONAL APPROXIMATIONS

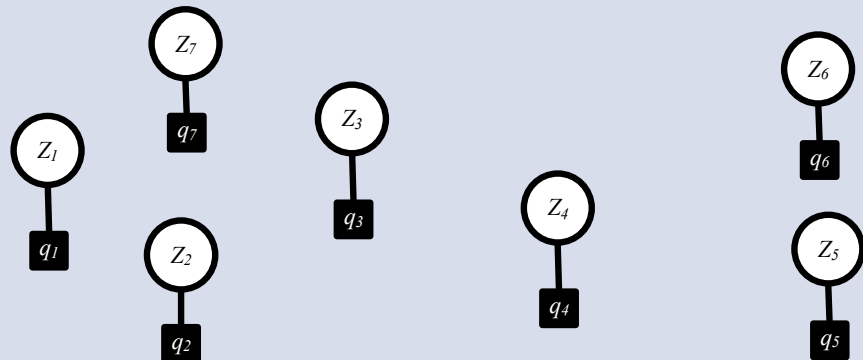
Mean Field Approximation

The **mean field approximation** assumes our variational approximation $q_{\theta}(\mathbf{z})$ treats each variable as independent

$$p_{\alpha}(\mathbf{z} \mid \mathbf{x}) = \frac{1}{Z(\mathbf{x})} \prod_{c \in \mathcal{C}} \psi_c(\mathbf{z}_c, \mathbf{x})$$



$$q_{\theta}(\mathbf{z}) = \prod_{t=1}^T q_t(z_t)$$

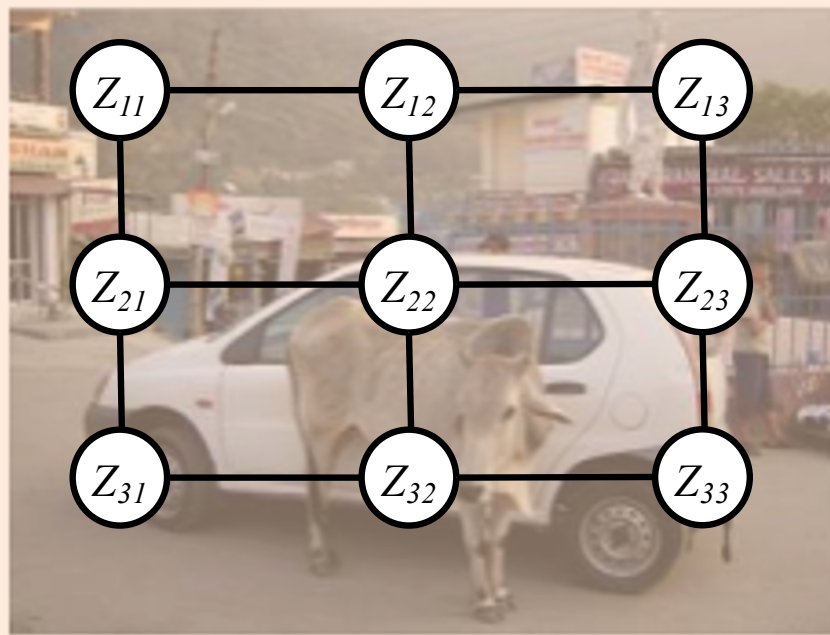


Mean Field Approximation

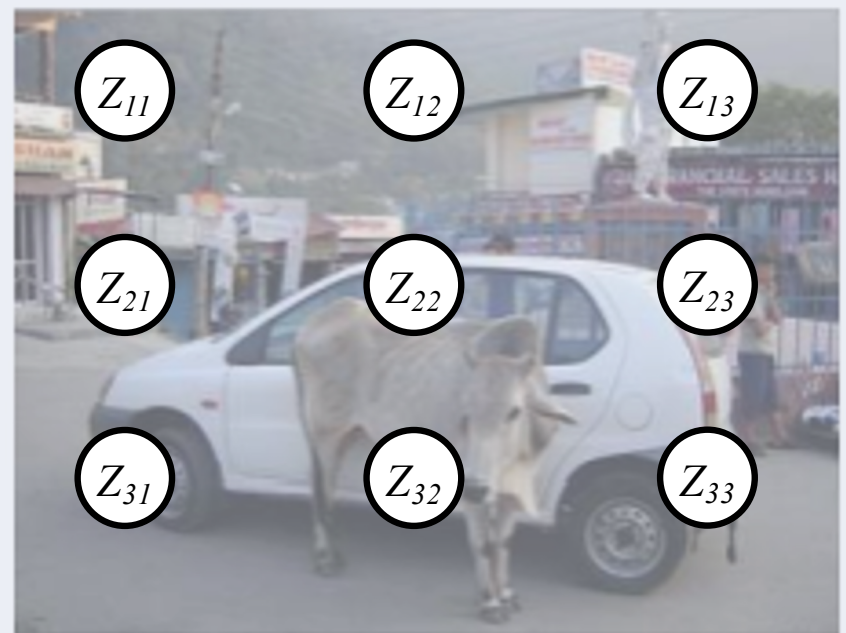
The **mean field approximation** assumes our variational approximation $q_{\theta}(\mathbf{z})$ treats each variable as independent

Ising Model


$$p_{\alpha}(\mathbf{z} \mid \mathbf{x}) = \frac{1}{Z(\mathbf{x})} \prod_{c \in \mathcal{C}} \psi_c(\mathbf{z}_c, \mathbf{x})$$



$$q_{\theta}(\mathbf{z}) = \prod_{t=1}^T q_t(z_t)$$

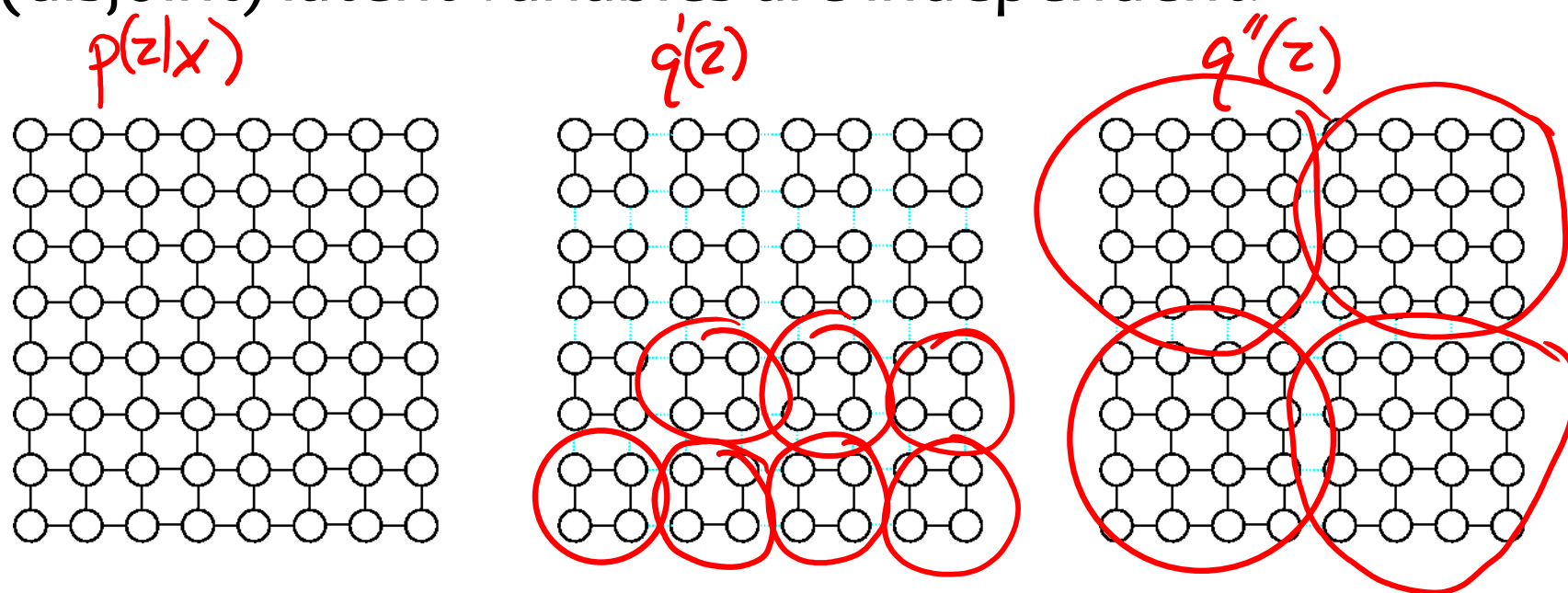


Structured Mean Field

- If q is not a mean-field approximation, but decomposes over “blocks” of variables, then we have the **Structured Mean Field algorithm**

V.I.
- Connection to related algorithms:
 - This is analogous to **Blocked** Gibbs Sampling
 - This is analogous to **Generalized** Belief Propagation
 - The names here (**Structured**, **Blocked**, **Generalized**) are different b/c they were invented by different people and no-one thought to rename them all “Blocked”

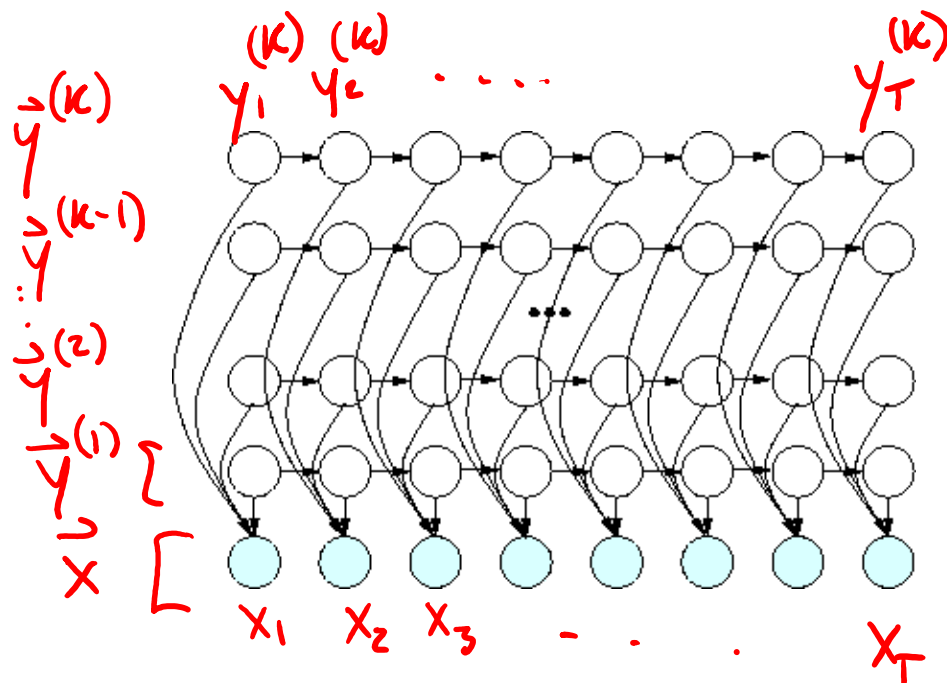
Structured Mean Field

- We can also apply more general forms of mean field approximations (involving clusters) to the Ising model:
- Instead of making all latent variables independent (i.e. naïve mean field, previous example), clusters of (disjoint) latent variables are independent.

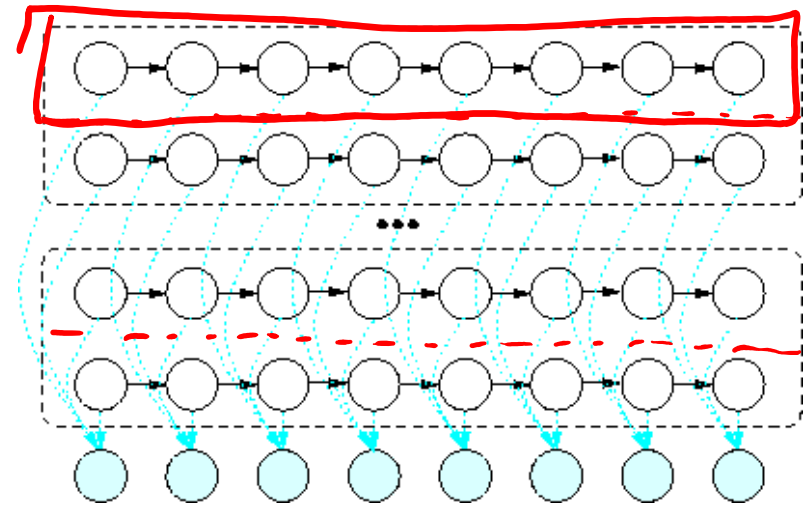


Structured Mean Field

- For a factorial HMM, we could decompose into chains



$$\frac{p(y|x)}{Z}$$



$$q(y)$$

Collapsed vs. Uncollapsed V.I.

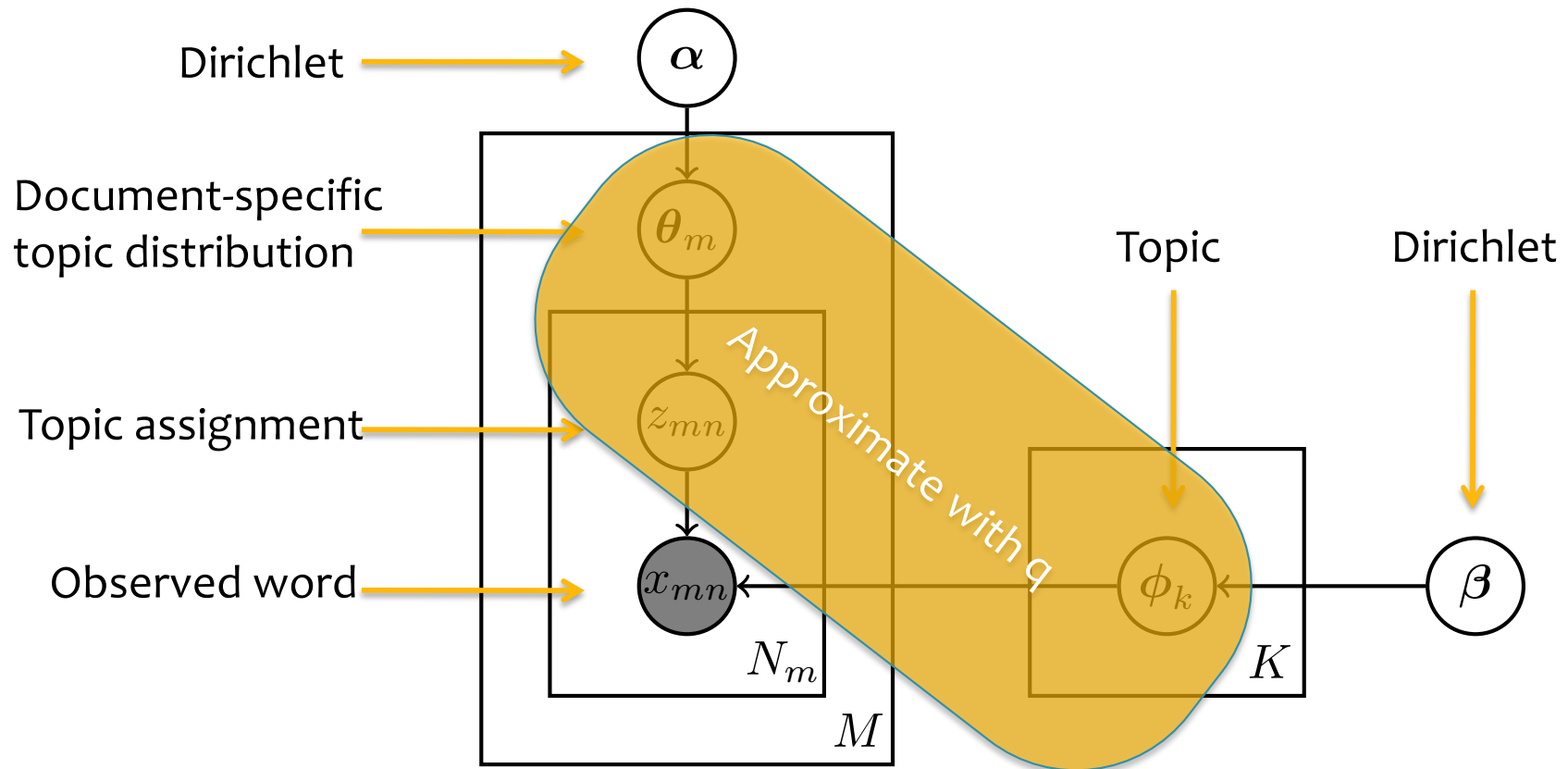
Just as we had collapsed and uncollapsed
Gibbs samplers for LDA...

... we can have collapsed and uncollapsed
variational inference for LDA

Collapsed vs. Uncollapsed V.I.

Latent Dirichlet Allocation (LDA)

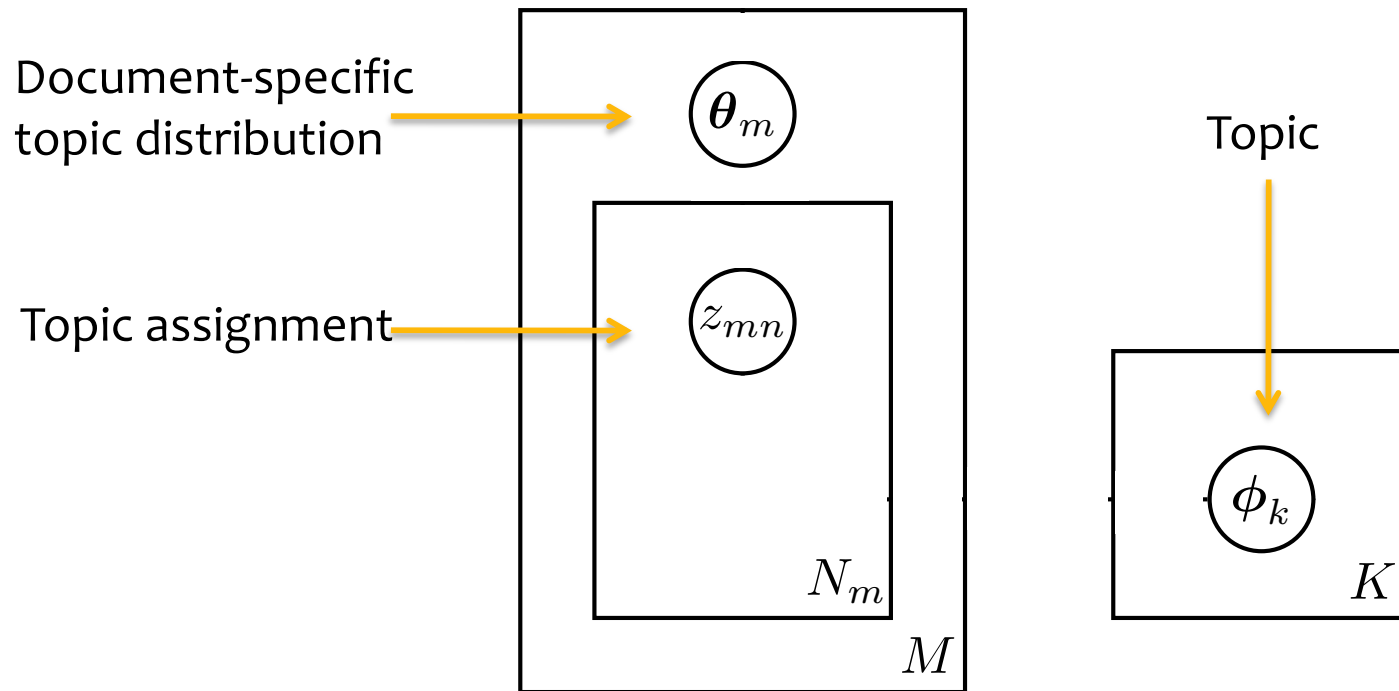
- Uncollapsed Variational Inference, aka. Explicit V.I. (original distribution)



Collapsed vs. Uncollapsed V.I.

Latent Dirichlet Allocation (LDA)

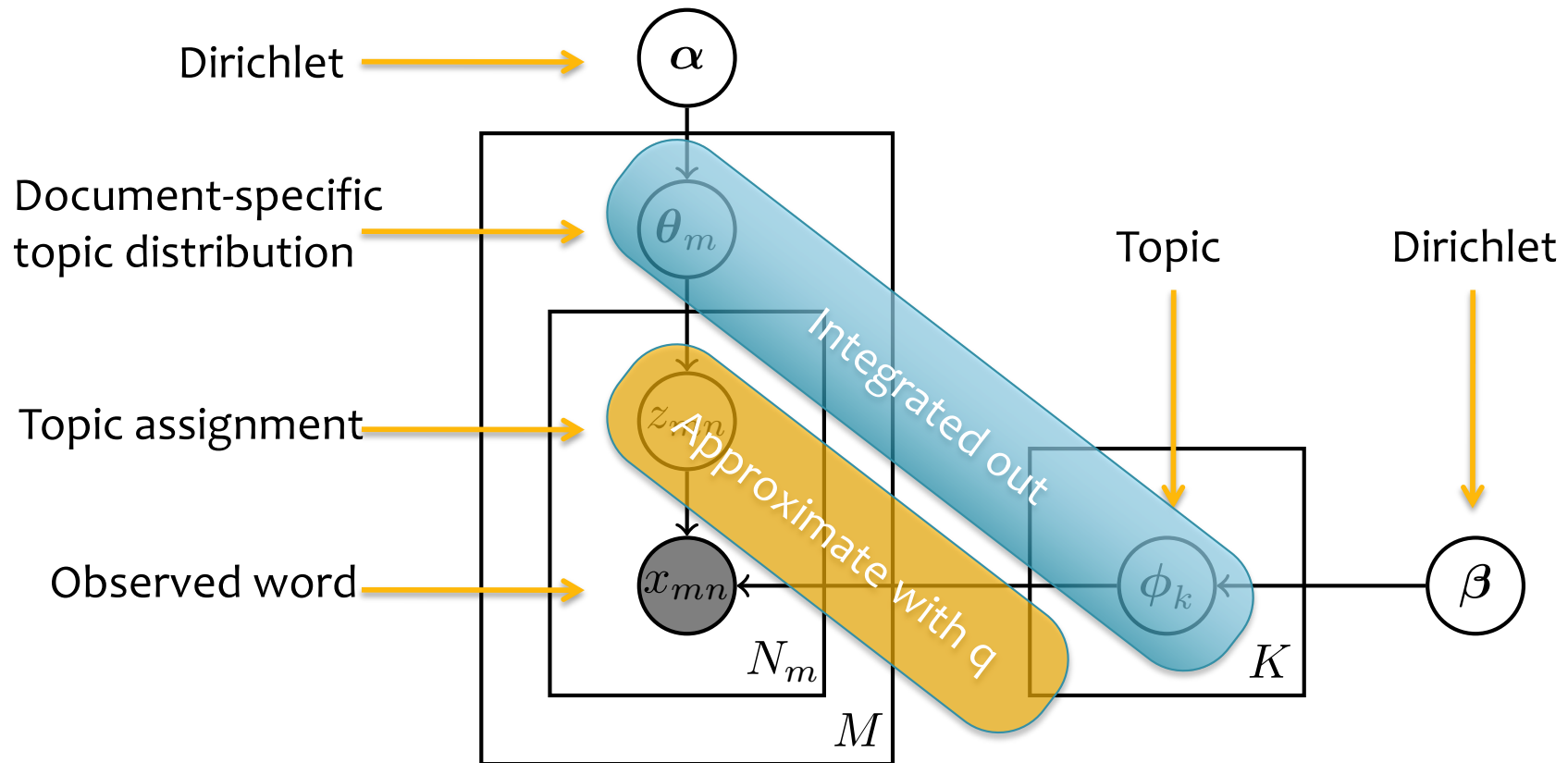
- Uncollapsed Variational Inference, aka. Explicit V.I.
(mean field variational approximation)



Collapsed vs. Uncollapsed V.I.

Latent Dirichlet Allocation (LDA)

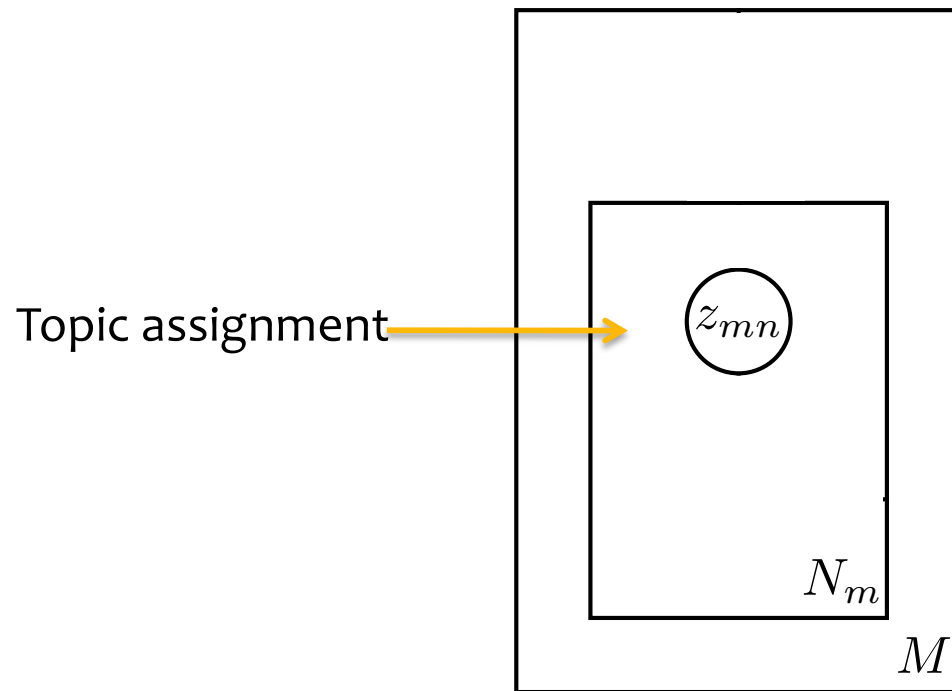
- Collapsed Variational Inference (original distribution)



Collapsed vs. Uncollapsed V.I.

Latent Dirichlet Allocation (LDA)

- Collapsed Variational Inference
(mean field variational approximation)



MEAN FIELD VARIATIONAL INFERENCE

Side Note

Contrast of three variational inference techniques:

1. Mean field variational inference minimizes $KL(q \parallel p)$
2. Expectation propagation minimizes $KL(p \parallel q)$
3. Loopy Belief Propagation minimizes the Bethe Free Energy

We are focused here on $KL(q \parallel p)$

KL Divergence

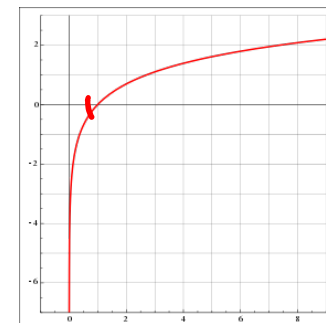
- Definition: for two distributions $q(x)$ and $p(x)$ over $x \in \mathcal{X}$, the **KL Divergence** is:

$$KL(q||p) = E_{q(x)} \left[\log \frac{q(x)}{p(x)} \right] = \begin{cases} \sum_x q(x) \log \frac{q(x)}{p(x)} & \text{discrete } x \\ \int_x q(x) \log \frac{q(x)}{p(x)} dx & \text{cont. } x \end{cases}$$

$$KL(p||q) = E_{p(x)} \left[\log \left(\frac{p(x)}{q(x)} \right) \right]$$

- Properties:
 - $KL(q || p)$ measures the **proximity** of two distributions q and p
 - KL is **not** symmetric: $KL(q || p) \neq KL(p || q)$
 - KL is minimized when $q(x) = p(x)$ for all $x \in \mathcal{X}$

KL Divergence

$$KL(q||p) = E_{q(x)} \left[\log \frac{q(x)}{p(x)} \right]$$


Understanding the Behavior of KL as an objective function

Example 1: Keeping all else constant, consider the effect of a particular x' on $KL(q || p)$

x'	$q(x')$	$p(x')$	$q(x') \log(q(x')/p(x'))$	effect on $KL(q p)$
1 →	0.9	0.9	0	no increase
2 →	0.9	0.1	1.97	big increase
3 →	0.1	0.9	-0.21	little decrease
4 →	0.1	0.1	0	little decrease

KL **does** insist on good approximations for values that have **high** probability in q

KL **does not** insist on good approximations for values that have **low** probability in q

Example 2: Which q distribution minimizes $KL(q || p)$?

$$p = \begin{bmatrix} 0.7 \\ 0.2 \\ 0.1 \end{bmatrix}$$

$$q^{(1)} = \begin{bmatrix} 1/3 \\ 1/3 \\ 1/3 \end{bmatrix}$$

$$q^{(2)} = \begin{bmatrix} 0.7 \\ 0.2 \\ 0.1 \end{bmatrix}$$

$$q^{(3)} = \begin{bmatrix} 0.1 \\ 0.1 \\ 0.1 \end{bmatrix}$$

Q: If we're minimizing KL, why not return $q^{(3)}$?
A: Because it's not a distribution!

$q^{(2)}$ minimizes KL

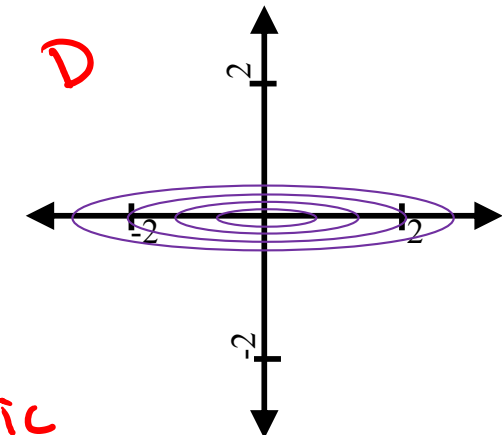
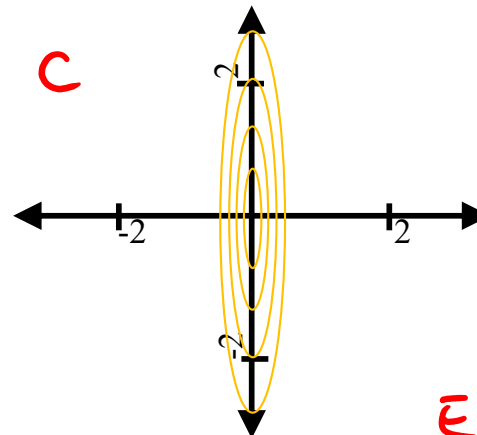
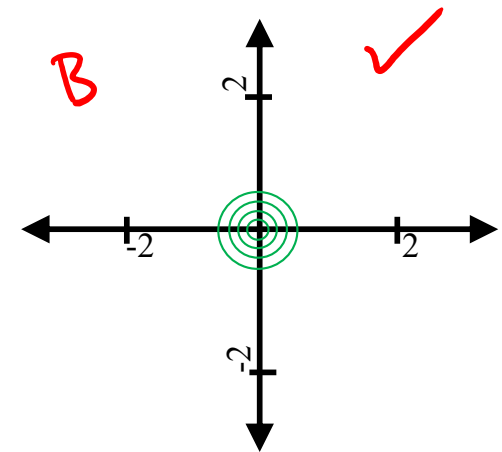
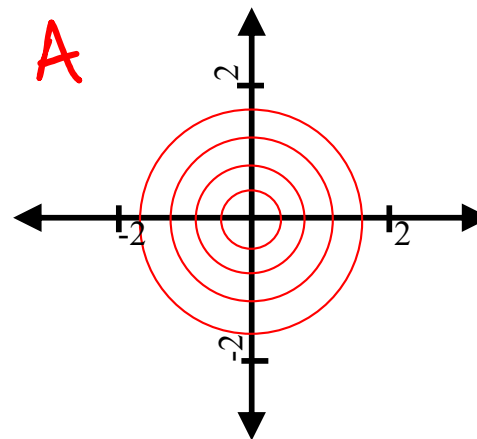
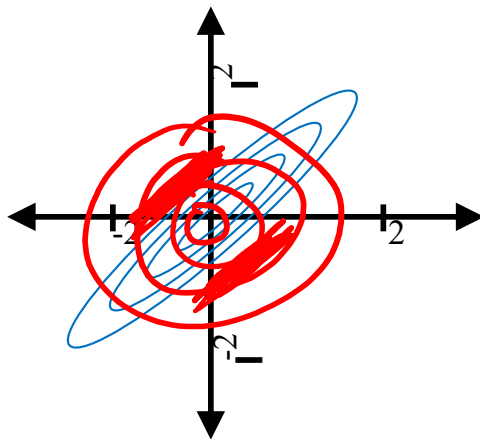
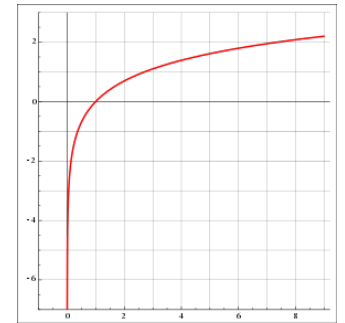
$KL(q||p) = E_{q(x)} \left[\log \frac{q(x)}{p(x)} \right]$ KL Divergence

Understanding the Behavior of KL as an objective function

Example 3: Which q distribution minimizes $KL(q || p)$?

$$p(\mathbf{x}) = \mathcal{N}(\boldsymbol{\mu} = [0, 0]^T, \boldsymbol{\Sigma})$$

$$q(x_1, x_2) = \mathcal{N}_1(x_1 | \mu_1, \sigma_1^2) \mathcal{N}_2(x_2 | \mu_2, \sigma_2^2)$$




E = toxic

Two Cases for Intractability

- Case 1:

given a **joint distribution** $p(x, z)$

$$\Rightarrow p(z \mid x) = \frac{p(x, z)}{p(x)}$$




we assume
 $p(x)$ is
intractable

- Case 2:

give **factor graph** and potentials

$$\Rightarrow p(z \mid x) = \frac{\tilde{p}(x, z)}{Z(x)}$$

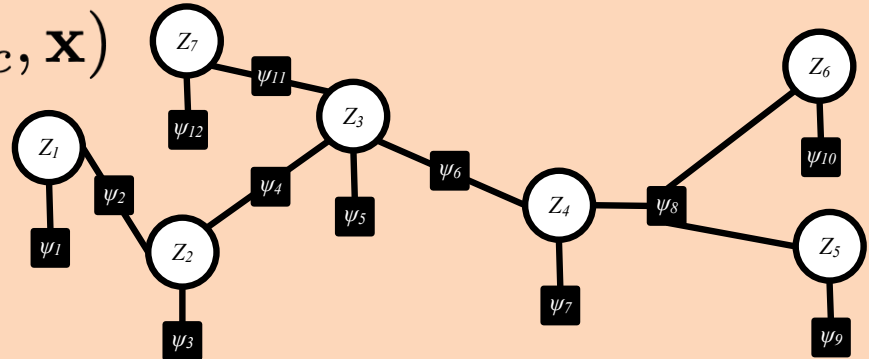


we assume
 $Z(x)$ is
intractable

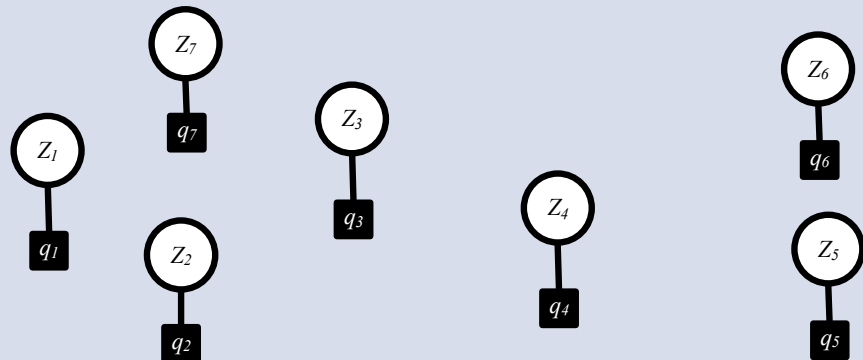
Mean Field Approximation

The **mean field approximation** assumes our variational approximation $q_{\theta}(\mathbf{z})$ treats each variable as independent

$$p_{\alpha}(\mathbf{z} \mid \mathbf{x}) = \frac{1}{Z(\mathbf{x})} \prod_{c \in \mathcal{C}} \psi_c(\mathbf{z}_c, \mathbf{x})$$



$$q_{\theta}(\mathbf{z}) = \prod_{t=1}^T q_t(z_t)$$

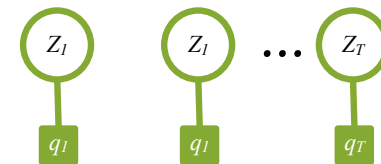


Mean Field V.I. Overview

1. Goal: estimate $p_\alpha(\mathbf{z} \mid \mathbf{x})$
we assume this is intractable to compute exactly
2. Idea: approximate with another distribution $q_\theta(\mathbf{z}) \approx p_\alpha(\mathbf{z} \mid \mathbf{x})$ for each \mathbf{x}

3. Mean Field: assume $q_\theta(\mathbf{z}) = \prod_t q_t(z_t; \theta)$

i.e., we decompose over variables



other choices for the decomposition of $q_\theta(\mathbf{z})$ give rise to “structured mean field”

4. Optimization Problem: pick the q that minimizes $KL(q \parallel p)$

$$\hat{q}(\mathbf{z}) = \operatorname{argmin}_{q(\mathbf{z}) \in \mathcal{Q}} KL(q(\mathbf{z}) \parallel p(\mathbf{z} \mid \mathbf{x}))$$

$$\hat{\theta} = \operatorname{argmin}_{\theta \in \Theta} KL(q_\theta(\mathbf{z}) \parallel p_\alpha(\mathbf{z} \mid \mathbf{x}))$$

equivalent

5. Optimization Algorithm: coordinate descent

i.e. pick the best $q_t(z_t)$ based on the other $\{q_s(z_s)\}_{s \neq t}$ being fixed

or gradient descent
SGD
Adam

Optimizing KL Divergence

- Question: How do we minimize KL?

$$\hat{\theta} = \operatorname{argmin}_{\theta \in \Theta} \underbrace{\text{KL}(q_{\theta}(\mathbf{z}) \parallel p_{\alpha}(\mathbf{z} \mid \mathbf{x}))}_{\text{KL Divergence}}$$

- Answer #1: Oh no! We can't even compute this KL.

Why we can't compute KL...

$$\begin{aligned} \text{KL}(q(\mathbf{z}) \parallel p(\mathbf{z} \mid \mathbf{x})) &= E_{q(\mathbf{z})} \left[\log \left(\frac{q(\mathbf{z})}{p(\mathbf{z} \mid \mathbf{x})} \right) \right] \\ &= E_{q(\mathbf{z})} [\log q(\mathbf{z})] - E_{q(\mathbf{z})} [\log p(\mathbf{z} \mid \mathbf{x})] \\ &= E_{q(\mathbf{z})} [\log q(\mathbf{z})] - E_{q(\mathbf{z})} [\log p(\mathbf{x}, \mathbf{z})] + E_{q(\mathbf{z})} [\log p(\mathbf{x})] \\ &= E_{q(\mathbf{z})} [\log q(\mathbf{z})] - E_{q(\mathbf{z})} [\log p(\mathbf{x}, \mathbf{z})] + \underbrace{\log p(\mathbf{x})}_{\text{we assumed this is intractable to compute!}} \end{aligned}$$

Handwritten red note: $\neq \frac{p(\mathbf{x}, \mathbf{z})}{p(\mathbf{x})}$

we have the same problem
with an intractable data
likelihood $p(\mathbf{x})$ or an intractable
partition function $Z(\mathbf{x})$

we assumed this
is intractable to
compute!

Optimizing KL Divergence

- Question: How do we minimize KL?

$$\hat{\theta} = \operatorname{argmin}_{\theta \in \Theta} \underbrace{\text{KL}(q_{\theta}(\mathbf{z}) \parallel p_{\alpha}(\mathbf{z} \mid \mathbf{x}))}_{\text{KL Divergence}}$$

- Answer #1: Oh no! We can't even compute this KL.

Why we can't compute KL...

$$\begin{aligned} \text{KL}(q(\mathbf{z}) \parallel p(\mathbf{z} \mid \mathbf{x})) &= E_{q(\mathbf{z})} \left[\log \left(\frac{q(\mathbf{z})}{p(\mathbf{z} \mid \mathbf{x})} \right) \right] \\ &= E_{q(\mathbf{z})} [\log q(\mathbf{z})] - E_{q(\mathbf{z})} [\log p(\mathbf{z} \mid \mathbf{x})] \\ &= E_{q(\mathbf{z})} [\log q(\mathbf{z})] - E_{q(\mathbf{z})} [\log \tilde{p}(\mathbf{z} \mid \mathbf{x})] + E_{q(\mathbf{z})} [\log Z(\mathbf{x})] \\ &= E_{q(\mathbf{z})} [\log q(\mathbf{z})] - E_{q(\mathbf{z})} [\log \tilde{p}(\mathbf{z} \mid \mathbf{x})] + \log Z(\mathbf{x}) \end{aligned}$$

we have the same problem
with an intractable data
likelihood $p(\mathbf{x})$ or an intractable
partition function $Z(\mathbf{x})$

we assumed this
is intractable to
compute!

Optimizing KL Divergence

- Question: How do we minimize KL?

$$\hat{\theta} = \operatorname{argmin}_{\theta \in \Theta} \underbrace{\text{KL}(q_{\theta}(\mathbf{z}) \parallel p_{\alpha}(\mathbf{z} \mid \mathbf{x}))}_{\text{KL Divergence}}$$

- Answer #2: We don't need to compute this KL
We can instead maximize the ELBO (i.e. **E**vidence **L**ower **B**ound)

$$\text{ELBO}(q_{\theta}) = E_{q_{\theta}(\mathbf{z})} [\log p_{\alpha}(\mathbf{x}, \mathbf{z})] - E_{q_{\theta}(\mathbf{z})} [\log q_{\theta}(\mathbf{z})]$$

The ELBO for a DGM

Here is why...

$$\begin{aligned} \theta &= \operatorname{argmin}_{\theta} \text{KL}(q_{\theta}(\mathbf{z}) \parallel p_{\alpha}(\mathbf{z} \mid \mathbf{x})) \\ &= \operatorname{argmin}_{\theta} E_{q_{\theta}(\mathbf{z})} [\log q_{\theta}(\mathbf{z})] - E_{q_{\theta}(\mathbf{z})} [\log p_{\alpha}(\mathbf{x}, \mathbf{z})] + \underbrace{\log p_{\alpha}(\mathbf{x})}_{\text{dropping the intractable term gives the ELBO}} \\ &= \operatorname{argmin}_{\theta} E_{q_{\theta}(\mathbf{z})} [\log q_{\theta}(\mathbf{z})] - E_{q_{\theta}(\mathbf{z})} [\log p_{\alpha}(\mathbf{x}, \mathbf{z})] \\ &= \operatorname{argmax}_{\theta} \text{ELBO}(q_{\theta}) \end{aligned}$$

Optimizing KL Divergence

- Question: How do we minimize KL?

$$\hat{\theta} = \operatorname{argmin}_{\theta \in \Theta} \underbrace{\text{KL}(q_{\theta}(\mathbf{z}) \parallel p_{\alpha}(\mathbf{z} \mid \mathbf{x}))}_{\text{KL Divergence}}$$

- Answer #2: We don't need to compute this KL
We can instead maximize the ELBO (i.e. **E**vidence **L**ower **B**ound)

$$\text{ELBO}(q_{\theta}) = E_{q_{\theta}(\mathbf{z})} [\log \tilde{p}_{\alpha}(\mathbf{z} \mid \mathbf{x})] - E_{q_{\theta}(\mathbf{z})} [\log q_{\theta}(\mathbf{z})]$$

Here is why... The ELBO for a UGM

$$\theta = \operatorname{argmin}_{\theta} \text{KL}(q_{\theta}(\mathbf{z}) \parallel p_{\alpha}(\mathbf{z} \mid \mathbf{x}))$$

$$= \frac{\tilde{p}(z|x)}{Z(x)}$$

$$= \operatorname{argmin}_{\theta} E_{q_{\theta}(\mathbf{z})} [\log q_{\theta}(\mathbf{z})] - E_{q_{\theta}(\mathbf{z})} [\log \tilde{p}_{\alpha}(\mathbf{z} \mid \mathbf{x})] + \underbrace{\log Z_{\alpha}(\mathbf{x})}_{\text{intractable term}}$$

$$= \operatorname{argmin}_{\theta} E_{q_{\theta}(\mathbf{z})} [\log q_{\theta}(\mathbf{z})] - E_{q_{\theta}(\mathbf{z})} [\log \tilde{p}_{\alpha}(\mathbf{z} \mid \mathbf{x})]$$

$$= \operatorname{argmax}_{\theta} \text{ELBO}(q_{\theta})$$

dropping the
intractable term
gives the ELBO

ELBO as Objective Function

What does maximizing $\text{ELBO}(q_\theta)$ accomplish?

$$\text{ELBO}(q_\theta) = E_{q_\theta(\mathbf{z})} [\log p_\alpha(\mathbf{x}, \mathbf{z})] - E_{q_\theta(\mathbf{z})} [\log q_\theta(\mathbf{z})]$$

1. The first expectation is high if q_θ puts probability mass on the same values of \mathbf{z} that p_α puts probability mass

2. The second term is the entropy of q_θ and the entropy will be high if q_θ spreads its probability mass evenly

ELBO as lower bound

- For a DGM:
 - $\text{ELBO}(q)$ is a lower bound for $\log p(x)$
- For a UGM:
 - $\text{ELBO}(q)$ is a lower bound for $\log Z(x)$

Takeaway: in variational inference, we find the q that gives the **tightest bound** on the normalization constant for $p(z \mid x)$

Variational Inference

Whiteboard

- Evidence Lower Bound (ELBO)
- ELBO's relation to $\log p(x)$