# Coordinate Ascent Variational Inference
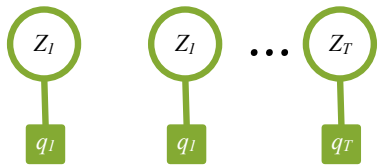
Matt Gormley
Lecture 17
Nov. 2, 2022
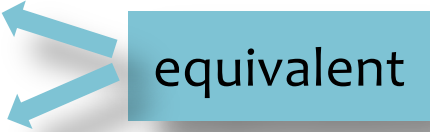
# Reminders

- **Lecture on Friday, Recitation on Monday**
- **Exam Rubrics and Exam Viewings**
- **Homework 4: MCMC**
  - **Out: Mon, Oct 24**
  - **Due: Fri, Nov 4 at 11:59pm**
- **Homework 5: Variational Inference**
  - **Out: Fri, Nov 5**
  - **Due: Wed, Nov 16 at 11:59pm**

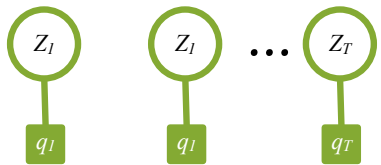# MEAN FIELD WITH GRADIENT ASCENT

# 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 ‖ p)

$$\hat{q}(\mathbf{z}) = \underset{q(\mathbf{z}) \in \mathcal{Q}}{\operatorname{argmin}} \, \mathsf{KL}(q(\mathbf{z}) \| p(\mathbf{z} \mid \mathbf{x}))$$

$$\hat{\theta} = \underset{\theta \in \Theta}{\operatorname{argmin}} \, \mathsf{KL}(q_\theta(\mathbf{z}) \| p_\alpha(\mathbf{z} \mid \mathbf{x}))$$

equivalent

5. *Optimization Algorithm*: pick your favorite {coordinate
   descent, gradient descent, etc.}

# 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 ‖ p)
$$\hat{\theta} = \underset{\theta}{\operatorname{argmin}} \, \mathsf{KL}(q_\theta(\mathbf{z}) \parallel p_\alpha(\mathbf{z} \mid \mathbf{x})) = \underset{\theta}{\operatorname{argmax}} \, \mathsf{ELBO}(q_\theta)$$
$$\mathsf{ELBO}(q_\theta) = E_{q_\theta(\mathbf{z})} \left[ \log p_\alpha(\mathbf{x}, \mathbf{z}) \right] - E_{q_\theta(\mathbf{z})} \left[ \log q_\theta(\mathbf{z}) \right]$$
$$\mathsf{ELBO}(q_\theta) = E_{q_\theta(\mathbf{z})} \left[ \log \tilde{p}_\alpha(\mathbf{z} \mid \mathbf{x}) \right] - E_{q_\theta(\mathbf{z})} \left[ \log q_\theta(\mathbf{z}) \right]$$

5. _Optimization Algorithm_: pick your favorite {coordinate
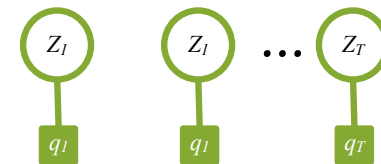   ascent, gradient ascent, etc.}

# 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 ∥ p)
$$\hat{\theta} = \operatorname*{argmin}_\theta \mathsf{KL}(q_\theta(\mathbf{z}) \parallel p_\alpha(\mathbf{z} \mid \mathbf{x})) = \operatorname*{argmax}_\theta \mathsf{ELBO}(q_\theta)$$
$$\mathsf{ELBO}(q_\theta) = E_{q_\theta(\mathbf{z})}\left[\log p_\alpha(\mathbf{x}, \mathbf{z})\right] - E_{q_\theta(\mathbf{z})}\left[\log q_\theta(\mathbf{z})\right]$$
$$\mathsf{ELBO}(q_\theta) = E_{q_\theta(\mathbf{z})}\left[\log \tilde{p}_\alpha(\mathbf{z} \mid \mathbf{x})\right] - E_{q_\theta(\mathbf{z})}\left[\log q_\theta(\mathbf{z})\right]$$

5. *Optimization Algorithm*: gradient ascent

# Mean Field w/Gradient Ascent

- **Note**: GA does local maximization, but ELBO is generally non-convex
- **Algorithm**:
  - Initialize θ
  - while not converged:
  $$\theta \leftarrow \theta + \gamma \nabla_\theta \text{ELBO}(q_\theta)$$
- **Gradient of ELBO**:

$$\nabla_\theta \text{ELBO}(q_\theta) = \nabla_\theta \mathbb{E}_{q_\theta}[\log p_\alpha(x, z)] - \nabla_\theta \mathbb{E}_{q_\theta}[\log q_\theta(z)]$$

$$= \cdots$$

$$= \cdots$$
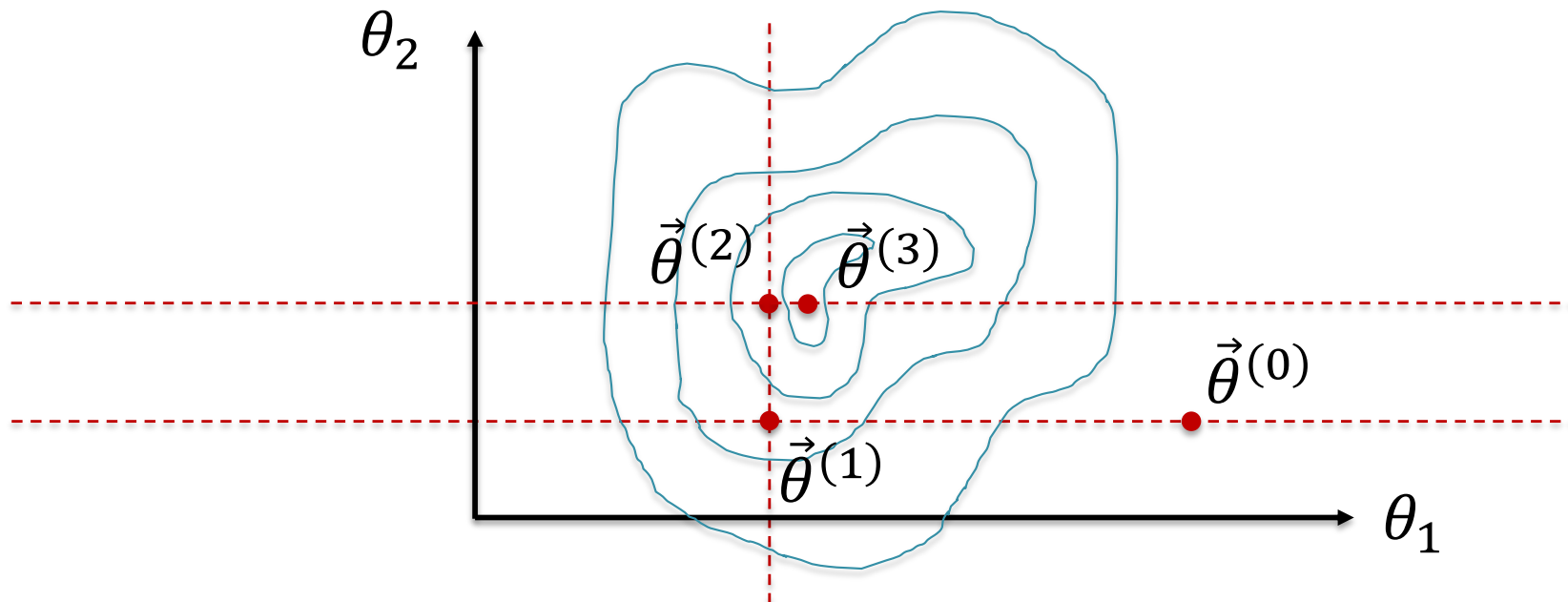
$$= \text{easy b/c of a simple } q_\theta$$

HW5?

# BACKGROUND: BLOCK COORDINATE DESCENT

# Coordinate Descent

- Goal: minimize some objective

$$\vec{\theta}^* = \underset{\vec{\theta}}{\mathrm{argmin}}\, J(\vec{\theta})$$

- Idea: iteratively pick one variable and minimize the objective w.r.t. just that one variable, *keeping all the others fixed.*

# Block Coordinate Descent

- Goal: minimize some objective (with 2 blocks)

$$\vec{\alpha}^*, \vec{\beta}^* = \underset{\vec{\alpha}, \vec{\beta}}{\arg\min} J(\vec{\alpha}, \vec{\beta})$$

- Idea: iteratively pick one *block* of variables ($\vec{\alpha}$ or $\vec{\beta}$) and minimize the objective w.r.t. that block, keeping the other(s) fixed.

  $Init \quad \alpha, \beta$

  **while** not converged:

  $$\vec{\alpha} = \underset{\vec{\alpha}}{\arg\min} J(\vec{\alpha}, \vec{\beta})$$

  $$\vec{\beta} = \underset{\vec{\beta}}{\arg\min} J(\vec{\alpha}, \vec{\beta})$$

# Block Coordinate Descent

- Goal: minimize some objective (with T blocks)

$$\boldsymbol{\alpha}_1, \ldots, \boldsymbol{\alpha}_T = \operatorname*{argmin}_{\boldsymbol{\alpha}_1} \cdots \operatorname*{argmin}_{\boldsymbol{\alpha}_T} J(\boldsymbol{\alpha}_1, \ldots, \boldsymbol{\alpha}_T)$$

- Idea: iteratively pick one *block* of variables (e.g. the vector $\boldsymbol{\alpha}_t$) and minimize the objective w.r.t. that block, keeping the other(s) fixed.

Init, $\alpha_1, \ldots, \alpha_T$

    **while** not converged:

      **for** $t = 1, \ldots, T$ :

$$\boldsymbol{\alpha}_t = \operatorname*{argmin}_{\boldsymbol{\alpha}_t} J(\boldsymbol{\alpha}_1, \ldots, \boldsymbol{\alpha}_T)$$

# COORDINATE ASCENT VARIATIONAL INFERENCE (CAVI)

# 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 ∥ p)
$$\hat{\theta} = \underset{\theta}{\mathrm{argmin}}\, \mathsf{KL}(q_\theta(\mathbf{z}) \,\|\, p_\alpha(\mathbf{z} \mid \mathbf{x})) = \underset{\theta}{\mathrm{argmax}}\, \mathsf{ELBO}(q_\theta)$$
$$\mathsf{ELBO}(q_\theta) = E_{q_\theta(\mathbf{z})}\left[\log p_\alpha(\mathbf{x}, \mathbf{z})\right] - E_{q_\theta(\mathbf{z})}\left[\log q_\theta(\mathbf{z})\right]$$
$$\mathsf{ELBO}(q_\theta) = E_{q_\theta(\mathbf{z})}\left[\log \tilde{p}_\alpha(\mathbf{z} \mid \mathbf{x})\right] - E_{q_\theta(\mathbf{z})}\left[\log q_\theta(\mathbf{z})\right]$$

5. *Optimization Algorithm*: coordinate ascent
   i.e. pick the best $q_t(z_t)$ based on the other $\{ q_s(z_s) \}_{s \neq t}$ being fixed

Choosing coordinate descent here yields the Coordinate Ascent Variational Inference (CAVI) algorithm

# CAVI Algorithm

Coordinate Ascent Variational Inference (CAVI)

- here we assume a **mean field** approximation
- application of **coordinate ascent** to maximization of ELBO
- converges to a **local optimum** of the **nonconvex** ELBO objective

$$q_{\neg t}(z_{\neg t}) = \prod_{s:\, s \neq t} q_s(z_s)$$

1: **procedure** CAVI($p_\alpha$)

2:      Let $q_\theta(\mathbf{z}) = \prod_{t=1}^{T} q_t(z_t)$          ▷ Mean field approx.

3:      **while** ELBO($q_\theta$) has not converged **do**

4:          **for** $t \in \{1, \dots, T\}$ **do**          ▷ For each variable

5:              Set $q_t(z_t) \propto \exp(E_{q_{\neg t}}[\log p_\alpha(z_t \mid z_{\neg t}, x)])$

6:              while keeping all $\{q_s(\cdot)\}_{s \neq t}$ fixed

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

8:      **return** $q_\theta$

# CAVI Algorithm

Coordinate ~~here~~ ... (CAVI)
- here ... ation
- appl ... imization
- conv ... **nconvex**

1: **procedure** CAVI($p_\alpha$)
2:     Let $q_\theta(\mathbf{z}) = \prod_{t=1}^{T} q_t(z_t)$                    ▷ Mean field approx.
3:     **while** ELBO($q_\theta$) has not converged **do**
4:         **for** $t \in \{1, \ldots, T\}$ **do**                    ▷ For each variable
5:             Set $q_t(z_t) \propto \exp(E_{q_{\neg t}}[\log p_\alpha(z_t \mid z_{\neg t}, x)])$
6:             while keeping all $\{q_s(\cdot)\}_{s \neq t}$ fixed
7:         Compute ELBO($q_\theta$) $= E_{q_\theta(\mathbf{z})}[\log p_\alpha(\mathbf{x}, \mathbf{z})] - E_{q_\theta(\mathbf{z})}[\log q_\theta(\mathbf{z})]$
8:     **return** $q_\theta$

Unlike Gibbs Sampling:
- we compute an entire distribution (instead of sampling a value)
- we condition on variable marginals (instead of on variable assignment)
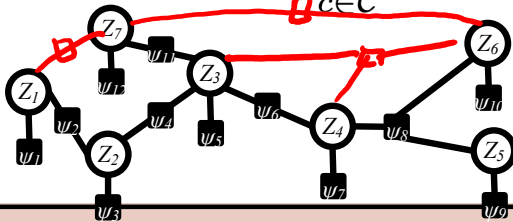
# Variational Inference

**Whiteboard**

– Computing marginals from a trained mean field approximation
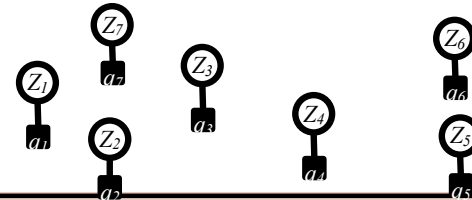
# EXAMPLE: CAVI FOR DISCRETE FACTOR GRAPH

# CAVI for a Discrete Factor Graph

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

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



1: **procedure** CAVI($p_\alpha$)

2:     Let $q_\theta(\mathbf{z}) = \prod_{t=1}^{T} q_t(z_t)$      ▷ Mean field approx.

3:     **while** ELBO($q_\theta$) has not converged **do**

4:         **for** $t \in \{1, \dots, T\}$ **do**      ▷ For each variable

5:             Set $q_t(z_t) \propto \exp(E_{q_{\neg t}}[\log p_\alpha(z_t \mid z_{\neg t}, x)])$

6:             while keeping all $\{q_s(\cdot)\}_{s \neq t}$ fixed

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

8:     **return** $q_\theta$

$$\Rightarrow q_t(z_t) \propto \exp\left( \sum_{\mathbf{z}_{\mathrm{MB}(z_t)}} \prod_{s \in \mathrm{MB}(z_t)} q_s(z_s) \log \prod_{c \in N(z_t)} \psi_c(\mathbf{z}_c) \right)$$

efficiently computed assuming number of neighbors N($z_t$) is not too large

20

# CAVI as Message Passing



Case 1: One Neighbor

$q_1(z_1)$

$\log(\psi_c(z_c))$

$q_2(z_2)$

| p | exp(0.08 + 0.16)/Z |
|---|---|
| d | exp(2.4 + 0)/Z |
| n | exp(0.8 + 0.2)/Z |

| | v | n |
|---|---|---|
| p | 0.1 | 8 |
| d | 3 | 0 |
| n | 1 | 1 |

| v | 0.8 |
|---|---|
| n | 0.2 |

$Z_1$ — $\psi_{12}$ — $Z_2$

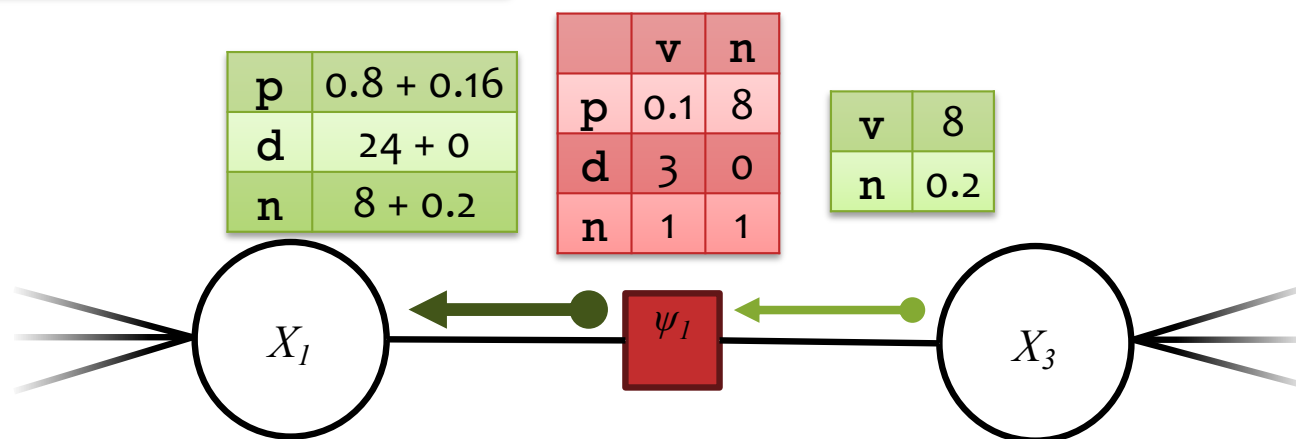CAVI message passing differs from BP in several ways:
- the beliefs are normalized (i.e. beliefs = marginals)
- no messages **to** factors (i.e. all messages are directly to a variable)
- matrix-vector product is exponentiated and normalized

$$\Rightarrow q_t(z_t) \propto \exp\left( \sum_{\mathbf{z}_{\text{MB}(z_t)}} \prod_{s \in \text{MB}(z_t)} q_s(z_s) \log \prod_{c \in N(z_t)} \psi_c(\mathbf{z}_c) \right)$$

# Sum-Product Belief Propagation

**Factor Message**

| | | |
|---|---|---|
| **p** | 0.8 + 0.16 | |
| **d** | 24 + 0 | |
| **n** | 8 + 0.2 | |

| | **v** | **n** |
|---|---|---|
| **p** | 0.1 | 8 |
| **d** | 3 | 0 |
| **n** | 1 | 1 |

| | |
|---|---|
| **v** | 8 |
| **n** | 0.2 |

$X_1$   $\psi_1$   $X_3$

$$\mu_{\alpha \to i}(x_i) = \sum_{\boldsymbol{x_\alpha} : \boldsymbol{x_\alpha}[i] = x_i} \psi_\alpha(\boldsymbol{x_\alpha}) \prod_{j \in \mathcal{N}(\alpha) \setminus i} \mu_{j \to \alpha}(\boldsymbol{x_\alpha}[i])$$

# CAVI as Message Passing

$\log \psi_c(z_c)$

|   |   | v | n |
|---|---|---|---|
| p | 0.08 + 0.2 | 0.1 | 8 |
| d | 2.4 + 0 | 3 | 0 |
| n | 0.8 + 0.2 | 1 | 1 |

| v | 0.8 |
|---|---|
| n | 0.2 |

| p | exp(0.28*3)/Z |
|---|---|
| d | exp(2.4*4)/Z |
| n | exp(1*1)/Z |

$Z_1$  $\psi_{12}$  $Z_2$

|   |   | a | d |
|---|---|---|---|
| p | 1 + 2 | 2 | 4 |
| d | 1.5 + 2.5 | 3 | 5 |
| n | 0.5 + 0.5 | 1 | 1 |

| a | 0.5 |
|---|---|
| d | 0.5 |

$\psi_{13}$  $Z_3$

$$\Rightarrow q_t(z_t) \propto \exp\left( \sum_{\mathbf{z}_{\mathrm{MB}(z_t)}} \prod_{s \in \mathrm{MB}(z_t)} q_s(z_s) \log \prod_{c \in N(z_t)} \psi_c(\mathbf{z}_c) \right)$$

23

# CAVI as Message Passing

**Case 2: Two Neighbors**

$\mu_{2\to1}(z_1)$

| | |
|---|---|
| p | 0.08 + 0.2 |
| d | 2.4 + 0 |
| n | 0.8 + 0.2 |

| | v | n |
|---|---|---|
| p | 0.1 | 8 |
| d | 3 | 0 |
| n | 1 | 1 |

| | |
|---|---|
| v | 0.8 |
| n | 0.2 |

| | |
|---|---|
| p | exp(0.28*3)/Z |
| d | exp(2.4*4)/Z |
| n | exp(1*1)/Z |

$Z_1$ — $\psi_{12}$ — $Z_2$

$\mu_{3\to1}(z_1)$

| | |
|---|---|
| p | 1 + 2 |
| d | 1.5 + 2.5 |
| n | 0.5 + 0.5 |

| | a | d |
|---|---|---|

| | |
|---|---|
| a | 0.5 |

For a **pairwise MRF,** we have the following simplified the update rules:

$$\mu_{s\to t}(z_t) = \sum_{z_s} q_s(z_s)\psi_{s,t}(z_s, z_t)$$

$$q_t(z_t) \propto \exp\left(\prod_{s\in\mathrm{MB}(z_t)} \mu_{s\to t}(z_t)\right)$$

# Variational Inference

## *Whiteboard*

- Computing the CAVI update ✓
  - Multinomial full conditionals
- Example: two variable factor graph
  - Joint distribution
  - Mean Field Variational Inference
  - Gibbs Sampling

*Q2: what questions do you have?*

# Q&A