

# Variational inference (cont'd)

## Solving VB and Starting ExpFam

Helge Langseth and Thomas Dyhre Nielsen

Oct. 2018

# Introduction

## Day 1: Bayesian networks – Definition and inference

- Definition of Bayesian networks: Syntax and semantics
- Exact inference
- Approximate inference using MCMC

## Day 2: Variational inference – Introduction and basis

- Approximate inference through the *Kullback-Leibler divergence*
- *Variational Bayes*
- The *mean-field* approach to Variational Bayes

## Day 3: Variational Bayes – cont'd

- Solving the VB equations
- Introducing Exponential families

## Day 4: Scalable Variational Bayes

- Variational message passing
- Stochastic gradient ascent
- Stochastic variational inference

## Day 5: Current approaches and extensions

- Variational Auto Encoders
- Black Box variational inference
- Probabilistic Programming Languages

- Exact inference, in particular for Bayesian inference of latent-variable models, is computationally very expensive.
- Approximate inference using simulation, like MCMC, importance sampling, ... :
  - Have asymptotic guarantees
  - Have a non-deterministic behaviour
  - Can be memory-expensive in particular for “difficult models” where we need to store the full sample to approximate inference
  - Can be time-expensive (slow convergence properties for a “difficult” model and/or evidence).
- Approximate inference using projections,  $\hat{q} = \arg \min_{q \in \mathcal{Q}} \Delta(p; q)$ :
  - Behavior depends on choices of  $\mathcal{Q}$  and  $\Delta(\cdot; \cdot)$ .
  - We look mostly at the **variational objective**:  $\Delta(p; q) = \text{KL}(q||p)$ .
  - ... and make the **mean-field assumption**:  $q \in \mathcal{Q} \Leftrightarrow q(\mathbf{z}) = \prod_j q_j(z_j)$ .

- Notice how we can rearrange the terms to find:

$$\begin{aligned}\text{KL}(q(\mathbf{z} | \mathbf{x}) || p(\mathbf{z} | \mathbf{x})) &= \mathbb{E}_q \left[ \log \frac{q(\mathbf{z} | \mathbf{x})}{p(\mathbf{z}, \mathbf{x})} \right] + \log [p(\mathbf{x})] = -\mathcal{L}(q) + \log p(\mathbf{x}) \\ \log p(\mathbf{x}) &= \mathcal{L}(q) + \text{KL}(q(\mathbf{z} | \mathbf{x}) || p(\mathbf{z} | \mathbf{x}))\end{aligned}$$

... where the Evidence Lower Bound (ELBO) is  $\mathcal{L}(q) = -\mathbb{E}_q \left[ \log \frac{q(\mathbf{z} | \mathbf{x})}{p(\mathbf{z}, \mathbf{x})} \right]$ .

- Since  $\log p(\mathbf{x})$  is constant wrt.  $q$  and  $\text{KL}(q || p) \geq 0$  it follows:
  - We can minimize  $\text{KL}(q(\mathbf{z} | \mathbf{x}) || p(\mathbf{z} | \mathbf{x}))$  by maximizing  $\mathcal{L}(q)$ .
  - $\mathcal{L}(q)$  is a *lower bound* of the marginal data likelihood  $p(\mathbf{x})$ , hence the variational objective makes sense in a learning setting.

## Algorithm:

- We have observed  $\mathbf{X} = \mathbf{x}$ , and have access to the full joint  $p(\mathbf{z}, \mathbf{x})$ .
- We posit a *variational family* of distributions  $q_j(\cdot \mid \boldsymbol{\lambda}_j)$ , i.e., we choose the distributional form, while wanting to optimize the parameterization  $\boldsymbol{\lambda}_j$ .
- The posterior approximation is assumed to factorize according to the mean-field assumption, and we use the KL ( $q(\mathbf{z}) \parallel p(\mathbf{z} \mid \mathbf{x})$ ) as our objective.

## Algorithm:

Repeat until negligible improvement in terms of  $\mathcal{L}(q)$ :

- For each  $j$ :
  - Somehow choose  $\boldsymbol{\lambda}_j$  to maximize  $\mathcal{L}(q)$ , typically based on  $\{\boldsymbol{\lambda}_i\}_{i \neq j}$ .
- Calculate the new  $\mathcal{L}(q)$ .

## Solving the VB optimization

We will maximize  $\mathcal{L}(q) = \mathbb{E}_q \left[ \log \frac{p(\mathbf{z}, \mathbf{x})}{q(\mathbf{z})} \right]$  under the assumption that  $q(\cdot)$  factorizes.

Let us pick one  $j$ , utilize that  $q(\mathbf{z}) = q_j(z_j) \cdot q_{\neg j}(\mathbf{z}_{\neg j})$ , and assume  $q_{\neg j}(\cdot)$  is kept fixed.

$$\begin{aligned}\mathcal{L}(q) &= \mathbb{E}_q [\log p(\mathbf{z}, \mathbf{x})] - \mathbb{E}_q [\log q(\mathbf{z})] \\ &= \mathbb{E}_{q_j} \mathbb{E}_{q_{\neg j}} [\log p(\mathbf{z}, \mathbf{x})] - \mathbb{E}_{q_j} \mathbb{E}_{q_{\neg j}} [\log q(\mathbf{z})]\end{aligned}$$



We will maximize  $\mathcal{L}(q) = \mathbb{E}_q \left[ \log \frac{p(\mathbf{z}, \mathbf{x})}{q(\mathbf{z})} \right]$  under the assumption that  $q(\cdot)$  factorizes.

Let us pick one  $j$ , utilize that  $q(\mathbf{z}) = q_j(z_j) \cdot q_{\neg j}(\mathbf{z}_{\neg j})$ , and assume  $q_{\neg j}(\cdot)$  is kept fixed.

$$\begin{aligned}\mathcal{L}(q) &= \mathbb{E}_q [\log p(\mathbf{z}, \mathbf{x})] - \mathbb{E}_q [\log q(\mathbf{z})] \\ &= \mathbb{E}_{q_j} \mathbb{E}_{q_{\neg j}} [\log p(\mathbf{z}, \mathbf{x})] - \mathbb{E}_{q_j} \mathbb{E}_{q_{\neg j}} [\log q(\mathbf{z})]\end{aligned}$$

For the term  $\mathbb{E}_{q_{\neg j}} [\log p(\mathbf{z}, \mathbf{x})]$  we simply define  $f_j(z_j)$  so that

$$\log f_j(z_j) = \mathbb{E}_{q_{\neg j}} [\log p(\mathbf{z}, \mathbf{x})]$$

We will maximize  $\mathcal{L}(q) = \mathbb{E}_q \left[ \log \frac{p(\mathbf{z}, \mathbf{x})}{q(\mathbf{z})} \right]$  under the assumption that  $q(\cdot)$  factorizes.

Let us pick one  $j$ , utilize that  $q(\mathbf{z}) = q_j(z_j) \cdot q_{\neg j}(\mathbf{z}_{\neg j})$ , and assume  $q_{\neg j}(\cdot)$  is kept fixed.

$$\begin{aligned}\mathcal{L}(q) &= \mathbb{E}_q [\log p(\mathbf{z}, \mathbf{x})] - \mathbb{E}_q [\log q(\mathbf{z})] \\ &= \mathbb{E}_{q_j} \mathbb{E}_{q_{\neg j}} [\log p(\mathbf{z}, \mathbf{x})] - \mathbb{E}_{q_j} \mathbb{E}_{q_{\neg j}} [\log q(\mathbf{z})] \\ &= \mathbb{E}_{q_j} \log f_j(z_j) - \mathbb{E}_{q_j} \mathbb{E}_{q_{\neg j}} [\log q(\mathbf{z})]\end{aligned}$$

We will maximize  $\mathcal{L}(q) = \mathbb{E}_q \left[ \log \frac{p(\mathbf{z}, \mathbf{x})}{q(\mathbf{z})} \right]$  under the assumption that  $q(\cdot)$  factorizes.  
Let us pick one  $j$ , utilize that  $q(\mathbf{z}) = q_j(z_j) \cdot q_{\neg j}(\mathbf{z}_{\neg j})$ , and assume  $q_{\neg j}(\cdot)$  is kept fixed.

$$\begin{aligned}\mathcal{L}(q) &= \mathbb{E}_q [\log p(\mathbf{z}, \mathbf{x})] - \mathbb{E}_q [\log q(\mathbf{z})] \\ &= \mathbb{E}_{q_j} \mathbb{E}_{q_{\neg j}} [\log p(\mathbf{z}, \mathbf{x})] - \mathbb{E}_{q_j} \mathbb{E}_{q_{\neg j}} [\log q(\mathbf{z})] \\ &= \mathbb{E}_{q_j} \log f_j(z_j) - \mathbb{E}_{q_j} \mathbb{E}_{q_{\neg j}} [\log q(\mathbf{z})]\end{aligned}$$

For the other term, notice that  $\log q(\mathbf{z}) = \log q_j(z_j) + \log q_{\neg j}(\mathbf{z}_{\neg j})$  Therefore

$$\begin{aligned}\mathbb{E}_{q_j} \mathbb{E}_{q_{\neg j}} [\log q(\mathbf{z})] &= \mathbb{E}_{q_j} \mathbb{E}_{q_{\neg j}} [\log q_j(z_j) + \log q_{\neg j}(\mathbf{z}_{\neg j})] \\ &= \mathbb{E}_{q_j} [\log q_j(z_j)] + \mathbb{E}_{q_{\neg j}} [\log q_{\neg j}(\mathbf{z}_{\neg j})] \\ &= \mathbb{E}_{q_j} [\log q_j(z_j)] + c,\end{aligned}$$

because  $\mathbb{E}_{q_{\neg j}} [\log q_{\neg j}(\mathbf{z}_{\neg j})]$  is constant wrt.  $q_j(\cdot)$ .

We will maximize  $\mathcal{L}(q) = \mathbb{E}_q \left[ \log \frac{p(\mathbf{z}, \mathbf{x})}{q(\mathbf{z})} \right]$  under the assumption that  $q(\cdot)$  factorizes.

Let us pick one  $j$ , utilize that  $q(\mathbf{z}) = q_j(z_j) \cdot q_{\neg j}(\mathbf{z}_{\neg j})$ , and assume  $q_{\neg j}(\cdot)$  is kept fixed.

$$\begin{aligned}\mathcal{L}(q) &= \mathbb{E}_q [\log p(\mathbf{z}, \mathbf{x})] - \mathbb{E}_q [\log q(\mathbf{z})] \\ &= \mathbb{E}_{q_j} \mathbb{E}_{q_{\neg j}} [\log p(\mathbf{z}, \mathbf{x})] - \mathbb{E}_{q_j} \mathbb{E}_{q_{\neg j}} [\log q(\mathbf{z})] \\ &= \mathbb{E}_{q_j} \log f_j(z_j) - \mathbb{E}_{q_j} [\log q_j(z_j)] + c \\ &= -\text{KL}(q_j(z_j) || f_j(z_j)) + c\end{aligned}$$

We will maximize  $\mathcal{L}(q) = \mathbb{E}_q \left[ \log \frac{p(\mathbf{z}, \mathbf{x})}{q(\mathbf{z})} \right]$  under the assumption that  $q(\cdot)$  factorizes.  
Let us pick one  $j$ , utilize that  $q(\mathbf{z}) = q_j(z_j) \cdot q_{\neg j}(\mathbf{z}_{\neg j})$ , and assume  $q_{\neg j}(\cdot)$  is kept fixed.

$$\begin{aligned}\mathcal{L}(q) &= \mathbb{E}_q [\log p(\mathbf{z}, \mathbf{x})] - \mathbb{E}_q [\log q(\mathbf{z})] \\ &= \mathbb{E}_{q_j} \mathbb{E}_{q_{\neg j}} [\log p(\mathbf{z}, \mathbf{x})] - \mathbb{E}_{q_j} \mathbb{E}_{q_{\neg j}} [\log q(\mathbf{z})] \\ &= \mathbb{E}_{q_j} \log f_j(z_j) - \mathbb{E}_{q_j} [\log q_j(z_j)] + c \\ &= -\text{KL}(q_j(z_j) || f_j(z_j)) + c\end{aligned}$$

**We get the following result:**

The ELBO is maximized wrt.  $q_j$  by choosing

$$q_j(z_j) = \frac{1}{Z} \exp(\mathbb{E}_{q_{\neg j}} [\log p(\mathbf{z}, \mathbf{x})])$$

**... and made the following assumptions to get there:**

- Mean field:  $q(\mathbf{z}) = \prod_i q_i(z_i)$ , and specifically  $q(\mathbf{z}) = q_j(z_j) \cdot q_{\neg j}(\mathbf{z}_{\neg j})$ .
- We optimize wrt.  $q_j(\cdot)$ , while keeping  $q_{\neg j}(\cdot)$  fixed – i.e., we do coordinate ascent in probability distribution space.

## Setup

- We have observed  $\mathbf{X} = \mathbf{x}$ , and can calculate the full joint  $p(\mathbf{z}, \mathbf{x})$ .
- The posterior approximation is assumed to factorize according to the mean-field assumption, and we use the KL ( $q(\mathbf{z}) || p(\mathbf{z} | \mathbf{x})$ ) as our objective.
- We posit a *variational family* of distributions  $q_j(z_j | \lambda_j)$ , i.e., we choose the distributional form, while wanting to **optimize the parameterization  $\lambda_j$** .
- The optimal  $\lambda_j$  **will** depend on  $\mathbf{x}$  – in fact  $\lambda_j$  encodes all the information about the other variables in the domain that  $Z_j$  is “aware of”.

## Setup

- We have observed  $\mathbf{X} = \mathbf{x}$ , and can calculate the full joint  $p(\mathbf{z}, \mathbf{x})$ .
- The posterior approximation is assumed to factorize according to the mean-field assumption, and we use the KL ( $q(\mathbf{z}) || p(\mathbf{z} | \mathbf{x})$ ) as our objective.
- We posit a *variational family* of distributions  $q_j(z_j | \lambda_j)$ , i.e., we choose the distributional form, while wanting to optimize the parameterization  $\lambda_j$ .
- The optimal  $\lambda_j$  **will** depend on  $\mathbf{x}$  – in fact  $\lambda_j$  encodes all the information about the other variables in the domain that  $Z_j$  is “aware of”.

## Algorithm

Repeat until negligible improvement in terms of  $\mathcal{L}(q)$ :

- For each  $j$ :
  - Calculate  $\mathbb{E}_{q_{-j}} [\log p(\mathbf{z}, \mathbf{x})]$  using current estimates for  $q_i(\cdot | \lambda_i)$ ,  $i \neq j$ .
  - Choose  $\lambda_j$  so that  $q_j(z_j | \lambda_j) \propto \exp(\mathbb{E}_{q_{-j}} [\log p(\mathbf{z}, \mathbf{x})])$ .
- Calculate the new  $\mathcal{L}(q)$ .

## Algorithm

Repeat until negligible improvement in terms of  $\mathcal{L}(q)$ :

- For each  $j$ :
  - Calculate  $\mathbb{E}_{q_{-j}} [\log p(\mathbf{z}, \mathbf{x})]$  using current estimates for  $q_i(\cdot | \boldsymbol{\lambda}_i)$ ,  $i \neq j$ .
  - Choose  $\boldsymbol{\lambda}_j$  so that  $q_j(z_j | \boldsymbol{\lambda}_j) \propto \exp(\mathbb{E}_{q_{-j}} [\log p(\mathbf{z}, \mathbf{x})])$ .
- Calculate the new  $\mathcal{L}(q)$ .

Calculating  $q_j(z_j | \boldsymbol{\lambda}_j)$ 

The update-rule can equivalently be expressed as

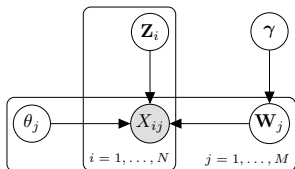
$$\log q_j(z_j | \boldsymbol{\lambda}_j) = \mathbb{E}_{q_{-j}} [\ln p(\mathbf{z}, \mathbf{x})] + c.$$

## Note!

- We only need to consider terms that share a factor with  $z_j$  – all other terms get absorbed into the constant  $c$ .
- $\rightsquigarrow$  need only reason about variables in the Markov blanket of  $Z_j$  – just as for Gibbs sampling!



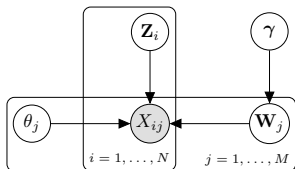
## The factor analysis model revisited and refined



- $X_{ij} \mid \{\mathbf{w}_j, \mathbf{z}_i, \theta_j\} \sim \mathcal{N}(\mathbf{w}_j^\top \mathbf{z}_i, 1/\theta_j)$
- $\mathbf{Z}_i \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_{D \times D})$
- $\mathbf{W}_j \sim \mathcal{N}(\mathbf{0}, \gamma^{-1} \mathbf{I}_{D \times D})$
- $\theta_j \sim \text{Gamma}(\theta_\theta, \theta_\theta)$
- $\gamma \sim \text{Gamma}(\alpha_\gamma, \beta_\gamma)$

- Num. of latent dim:  $D$
- Num. of data dim:  $M$
- Num. of data inst:  $N$

## The factor analysis model revisited and refined



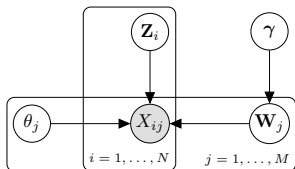
- $X_{ij} \mid \{\mathbf{w}_j, \mathbf{z}_i, \theta_j\} \sim \mathcal{N}(\mathbf{w}_j^T \mathbf{z}_i, 1/\theta_j)$
- $\mathbf{Z}_i \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_{D \times D})$
- $\mathbf{W}_j \sim \mathcal{N}(\mathbf{0}, \gamma^{-1} \mathbf{I}_{D \times D})$
- $\theta_j \sim \text{Gamma}(\theta_\theta, \theta_\theta)$
- $\gamma \sim \text{Gamma}(\alpha_\gamma, \beta_\gamma)$

- Num. of latent dim:  $D$
- Num. of data dim:  $M$
- Num. of data inst:  $N$

## The probability model

$$p(\cdot) = p(\gamma) \left[ \prod_{i=1}^N p(\mathbf{z}_i) \right] \left[ \prod_{j=1}^M p(\mathbf{w}_j \mid \gamma) p(\theta_j) \right] \left[ \prod_{i=1}^N \prod_{j=1}^M p(x_{ij} \mid \mathbf{w}_j, \mathbf{z}_i, \theta_j) \right]$$

## The factor analysis model revisited and refined



- $X_{ij} \mid \{\mathbf{w}_j, \mathbf{z}_i, \theta_j\} \sim \mathcal{N}(\mathbf{w}_j^T \mathbf{z}_i, 1/\theta_j)$
- $\mathbf{Z}_i \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_{D \times D})$
- $\mathbf{W}_j \sim \mathcal{N}(\mathbf{0}, \gamma^{-1} \mathbf{I}_{D \times D})$
- $\theta_j \sim \text{Gamma}(\theta_\theta, \theta_\theta)$
- $\gamma \sim \text{Gamma}(\alpha_\gamma, \beta_\gamma)$
- Num. of latent dim:  $D$
- Num. of data dim:  $M$
- Num. of data inst:  $N$

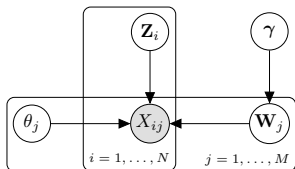
## The probability model

$$p(\cdot) = p(\gamma) \left[ \prod_{i=1}^N p(\mathbf{z}_i) \right] \left[ \prod_{j=1}^M p(\mathbf{w}_j \mid \gamma) p(\theta_j) \right] \left[ \prod_{i=1}^N \prod_{j=1}^M p(x_{ij} \mid \mathbf{w}_j, \mathbf{z}_i, \theta_j) \right]$$

## ...after taking the log

$$\log p(\cdot) = \log p(\gamma) + \sum_{i=1}^N \log p(\mathbf{z}_i) + \sum_{j=1}^M [\log p(\mathbf{w}_j \mid \gamma) + \log p(\theta_j)] + \sum_{i=1}^N \sum_{j=1}^M \log p(x_{ij} \mid \mathbf{w}_j, \mathbf{z}_i, \theta_j)$$

## The factor analysis model revisited and refined



- $X_{ij} \mid \{\mathbf{w}_j, \mathbf{z}_i, \theta_j\} \sim \mathcal{N}(\mathbf{w}_j^T \mathbf{z}_i, 1/\theta_j)$
- $\mathbf{Z}_i \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_{D \times D})$
- $\mathbf{W}_j \sim \mathcal{N}(\mathbf{0}, \gamma^{-1} \mathbf{I}_{D \times D})$
- $\theta_j \sim \text{Gamma}(\theta_\theta, \theta_\theta)$
- $\gamma \sim \text{Gamma}(\alpha_\gamma, \beta_\gamma)$
- Num. of latent dim:  $D$
- Num. of data dim:  $M$
- Num. of data inst:  $N$

## The probability model

$$p(\cdot) = p(\gamma) \left[ \prod_{i=1}^N p(\mathbf{z}_i) \right] \left[ \prod_{j=1}^M p(\mathbf{w}_j \mid \gamma) p(\theta_j) \right] \left[ \prod_{i=1}^N \prod_{j=1}^M p(x_{ij} \mid \mathbf{w}_j, \mathbf{z}_i, \theta_j) \right]$$

## The variational model

$$q(\cdot) = q(\gamma) \prod_{i=1}^N q(\mathbf{z}_i \mid \cdot) \prod_{j=1}^M q(\mathbf{w}_j \mid \cdot) q(\theta_j \mid \cdot)$$

We choose the variational distribution so that

$$\log q(\gamma \mid \cdot) = \mathbb{E}_{q_{-\gamma}} [\log p(\cdot)] + c$$

We choose the variational distribution so that

$$\log q(\gamma \mid \cdot) = \mathbb{E}_{q_{-\gamma}} [\log p(\cdot)] + c$$

## Recall

$$\begin{aligned} \log p(\cdot) \\ = \log p(\gamma) + \sum_{i=1}^N \log p(\mathbf{z}_i) + \sum_{j=1}^M [\log p(\mathbf{w}_j \mid \gamma) + \log p(\theta_j)] + \sum_{i=1}^N \sum_{j=1}^M \log p(x_{ij} \mid \mathbf{w}_j, \mathbf{z}_i, \theta_j) \end{aligned}$$

We choose the variational distribution so that

$$\log q(\gamma \mid \cdot) = \mathbb{E}_{q_{-\gamma}} [\log p(\cdot)] + c$$

## Recall

$$\begin{aligned} \log p(\cdot) \\ = \log p(\gamma) + \sum_{i=1}^N \log p(\mathbf{z}_i) + \sum_{j=1}^M [\log p(\mathbf{w}_j \mid \gamma) + \log p(\theta_j)] + \sum_{i=1}^N \sum_{j=1}^M \log p(x_{ij} \mid \mathbf{w}_j, \mathbf{z}_i, \theta_j) \end{aligned}$$

We choose the variational distribution so that

$$\log q(\gamma \mid \cdot) = \mathbb{E}_{q_{-\gamma}} [\log p(\cdot)] + c = \log p(\gamma) + \sum_{j=1}^M \mathbb{E}_{q_{-\gamma}} \log p(\mathbf{w}_j \mid \gamma) + c$$

## Recall

$$\begin{aligned} \log p(\cdot) \\ = \log p(\gamma) + \sum_{i=1}^N \log p(\mathbf{z}_i) + \sum_{j=1}^M [\log p(\mathbf{w}_j \mid \gamma) + \log p(\theta_j)] + \sum_{i=1}^N \sum_{j=1}^M \log p(x_{ij} \mid \mathbf{w}_j, \mathbf{z}_i, \theta_j) \end{aligned}$$



We choose the variational distribution so that

$$\log q(\gamma \mid \cdot) = \mathbb{E}_{q_{-\gamma}} [\log p(\cdot)] + c = \log p(\gamma) + \sum_{j=1}^M \mathbb{E}_{q_{-\gamma}} \log p(\mathbf{w}_j \mid \gamma) + c$$

## The gamma and multivariate normal

$$\log p(\gamma \mid \alpha_\gamma, \beta_\gamma) = \alpha_\gamma \log(\beta_\gamma) + (\alpha_\gamma - 1) \log(\gamma) - \beta_\gamma \cdot \gamma - \log(\Gamma(\alpha_\gamma))$$

$$\log p(\mathbf{w}_j \mid \gamma) = \log \mathcal{N}(\mathbf{w}_j \mid \mathbf{0}, \gamma^{-1} \mathbf{I}) = -\frac{D}{2} \log(2\pi) + \frac{D}{2} \log(\gamma) - \frac{\gamma}{2} \mathbf{w}_j^T \mathbf{w}_j$$

We choose the variational distribution so that

$$\log q(\gamma \mid \cdot) = \mathbb{E}_{q_{-\gamma}} [\log p(\cdot)] + c = \log p(\gamma) + \sum_{j=1}^M \mathbb{E}_{q_{-\gamma}} \log p(\mathbf{w}_j \mid \gamma) + c$$

## The gamma and multivariate normal

$$\log p(\gamma \mid \alpha_\gamma, \beta_\gamma) = \alpha_\gamma \log(\beta_\gamma) + (\alpha_\gamma - 1) \log(\gamma) - \beta_\gamma \cdot \gamma - \log(\Gamma(\alpha_\gamma))$$

$$\log p(\mathbf{w}_j \mid \gamma) = \log \mathcal{N}(\mathbf{w}_j \mid \mathbf{0}, \gamma^{-1} \mathbf{I}) = -\frac{D}{2} \log(2\pi) + \frac{D}{2} \log(\gamma) - \frac{\gamma}{2} \mathbf{w}_j^T \mathbf{w}_j$$

We choose the variational distribution so that

$$\begin{aligned}\log q(\gamma \mid \cdot) &= \mathbb{E}_{q_{\neg\gamma}} [\log p(\cdot)] + c = \log p(\gamma) + \sum_{j=1}^M \mathbb{E}_{q_{\neg j}} \log p(\mathbf{w}_j \mid \gamma) + c \\ &= (\alpha_\gamma - 1) \log(\gamma) - \beta_\gamma \cdot \gamma + \frac{D \cdot M}{2} \log(\gamma) - \frac{\gamma}{2} \sum_{j=1}^M \mathbb{E}(\mathbf{w}_j^T \mathbf{w}_j) + c\end{aligned}$$

## The gamma and multivariate normal

$$\begin{aligned}\log p(\gamma \mid \alpha_\gamma, \beta_\gamma) &= \alpha_\gamma \log(\beta_\gamma) + (\alpha_\gamma - 1) \log(\gamma) - \beta_\gamma \cdot \gamma - \log(\Gamma(\alpha_\gamma)) \\ \log p(\mathbf{w}_j \mid \gamma) &= \log \mathcal{N}(\mathbf{w}_j \mid \mathbf{0}, \gamma^{-1} \mathbf{I}) = -\frac{D}{2} \log(2\pi) + \frac{D}{2} \log(\gamma) - \frac{\gamma}{2} \mathbf{w}_j^T \mathbf{w}_j\end{aligned}$$

We choose the variational distribution so that

$$\begin{aligned}\log q(\gamma \mid \cdot) &= \mathbb{E}_{q_{-\gamma}} [\log p(\cdot)] + c = \log p(\gamma) + \sum_{j=1}^M \mathbb{E}_{q_{-\gamma}} \log p(\mathbf{w}_j \mid \gamma) + c \\&= (\alpha_\gamma - 1) \log(\gamma) - \beta_\gamma \cdot \gamma + \frac{D \cdot M}{2} \log(\gamma) - \frac{\gamma}{2} \sum_{j=1}^M \mathbb{E}(\mathbf{w}_j^T \mathbf{w}_j) + c \\&= \left[ \alpha_\gamma - 1 + \frac{D \cdot M}{2} \right] \log(\gamma) - \left[ \beta_\gamma + \frac{1}{2} \sum_{j=1}^M \mathbb{E}(\mathbf{w}_j^T \mathbf{w}_j) \right] \cdot \gamma + c\end{aligned}$$

We choose the variational distribution so that

$$\begin{aligned}\log q(\gamma \mid \cdot) &= \mathbb{E}_{q_{-\gamma}} [\log p(\cdot)] + c = \log p(\gamma) + \sum_{j=1}^M \mathbb{E}_{q_{-\gamma}} \log p(\mathbf{w}_j \mid \gamma) + c \\&= (\alpha_\gamma - 1) \log(\gamma) - \beta_\gamma \cdot \gamma + \frac{D \cdot M}{2} \log(\gamma) - \frac{\gamma}{2} \sum_{j=1}^M \mathbb{E}(\mathbf{w}_j^T \mathbf{w}_j) + c \\&= \left[ \alpha_\gamma - 1 + \frac{D \cdot M}{2} \right] \log(\gamma) - \left[ \beta_\gamma + \frac{1}{2} \sum_{j=1}^M \mathbb{E}(\mathbf{w}_j^T \mathbf{w}_j) \right] \cdot \gamma + c\end{aligned}$$

Recall the (generic) gamma distribution

$$\log p(x \mid \alpha, \beta) = \alpha \log(\beta) + (\alpha - 1) \log(x) - \beta \cdot x - \log(\Gamma(\alpha))$$

We choose the variational distribution so that

$$\begin{aligned}\log q(\gamma \mid \cdot) &= \mathbb{E}_{q_{-\gamma}} [\log p(\cdot)] + c = \log p(\gamma) + \sum_{j=1}^M \mathbb{E}_{q_{-\gamma}} \log p(\mathbf{w}_j \mid \gamma) + c \\&= (\alpha_{\gamma} - 1) \log(\gamma) - \beta_{\gamma} \cdot \gamma + \frac{D \cdot M}{2} \log(\gamma) - \frac{\gamma}{2} \sum_{j=1}^M \mathbb{E}(\mathbf{w}_j^T \mathbf{w}_j) + c \\&= \left[ \alpha_{\gamma} - 1 + \frac{D \cdot M}{2} \right] \log(\gamma) - \left[ \beta_{\gamma} + \frac{1}{2} \sum_{j=1}^M \mathbb{E}(\mathbf{w}_j^T \mathbf{w}_j) \right] \cdot \gamma + c\end{aligned}$$

Recall the (generic) gamma distribution

$$\log p(x \mid \alpha, \beta) = \alpha \log(\beta) + (\alpha - 1) \log(x) - \beta \cdot x - \log(\Gamma(\alpha))$$

We choose the variational distribution so that

$$\begin{aligned}
 \log q(\gamma | \cdot) &= \mathbb{E}_{q_{-\gamma}} [\log p(\cdot)] + c = \log p(\gamma) + \sum_{j=1}^M \mathbb{E}_{q_{-\gamma}} \log p(\mathbf{w}_j | \gamma) + c \\
 &= (\alpha_\gamma - 1) \log(\gamma) - \beta_\gamma \cdot \gamma + \frac{D \cdot M}{2} \log(\gamma) - \frac{\gamma}{2} \sum_{j=1}^M \mathbb{E}(\mathbf{w}_j^T \mathbf{w}_j) + c \\
 &= \left[ \alpha_\gamma - 1 + \frac{D \cdot M}{2} \right] \log(\gamma) - \left[ \beta_\gamma + \frac{1}{2} \sum_{j=1}^M \mathbb{E}(\mathbf{w}_j^T \mathbf{w}_j) \right] \cdot \gamma + c
 \end{aligned}$$

**Thus** we see that  $q(\gamma | \cdot)$  is gamma distributed with

- shape parameter:  $\alpha \leftarrow \alpha_\gamma + \frac{DM}{2}$
- rate parameter:  $\beta \leftarrow \beta_\gamma + \frac{1}{2} \sum_{j=1}^M \mathbb{E}_{q(\mathbf{w}_j)} [\mathbf{w}_j^T \mathbf{w}_j]$

Recall the (generic) gamma distribution

$$\log p(x | \alpha, \beta) = \alpha \log(\beta) + (\alpha - 1) \log(x) - \beta \cdot x - \log(\Gamma(\alpha))$$

We choose the variational distribution so that

$$\begin{aligned}
 \log q(\gamma | \cdot) &= \mathbb{E}_{q_{\neg\gamma}} [\log p(\cdot)] + c = \log p(\gamma) + \sum_{j=1}^M \mathbb{E}_{q_{\neg j}} \log p(\mathbf{w}_j | \gamma) + c \\
 &= (\alpha_\gamma - 1) \log(\gamma) - \beta_\gamma \cdot \gamma + \frac{D \cdot M}{2} \log(\gamma) - \frac{\gamma}{2} \sum_{j=1}^M \mathbb{E}(\mathbf{w}_j^T \mathbf{w}_j) + c \\
 &= \left[ \alpha_\gamma - 1 + \frac{D \cdot M}{2} \right] \log(\gamma) - \left[ \beta_\gamma + \frac{1}{2} \sum_{j=1}^M \mathbb{E}(\mathbf{w}_j^T \mathbf{w}_j) \right] \cdot \gamma + c
 \end{aligned}$$

**Thus** we see that  $q(\gamma | \cdot)$  is gamma distributed with

- shape parameter:  $\alpha \leftarrow \alpha_\gamma + \frac{DM}{2}$
- rate parameter:  $\beta \leftarrow \beta_\gamma + \frac{1}{2} \sum_{j=1}^M \mathbb{E}_{q(\mathbf{w}_j)} [\mathbf{w}_j^T \mathbf{w}_j]$

Calculation of  $\mathbb{E}_{q(\mathbf{w}_j)} [\mathbf{w}_j^T \mathbf{w}_j]$

$$\mathbb{E}_{q(\mathbf{w}_j)} [\mathbf{w}_j^T \mathbf{w}_j] = \sum_{d=1}^D \text{Var}_{q(\mathbf{w}_j)} [\mathbf{w}_{jd}] + \sum_{d=1}^D (\mathbb{E}_{q(\mathbf{w}_j)} [\mathbf{w}_{j,d}])^2$$



We choose the variational distribution so that

$$\begin{aligned}
 \log q(\gamma | \cdot) &= \mathbb{E}_{q_{-\gamma}} [\log p(\cdot)] + c = \log p(\gamma) + \sum_{j=1}^M \mathbb{E}_{q_{-\gamma}} \log p(\mathbf{w}_j | \gamma) + c \\
 &= (\alpha_\gamma - 1) \log(\gamma) - \beta_\gamma \cdot \gamma + \frac{D \cdot M}{2} \log(\gamma) - \frac{\gamma}{2} \sum_{j=1}^M \mathbb{E}(\mathbf{w}_j^T \mathbf{w}_j) + c \\
 &= \left[ \alpha_\gamma - 1 + \frac{D \cdot M}{2} \right] \log(\gamma) - \left[ \beta_\gamma + \frac{1}{2} \sum_{j=1}^M \mathbb{E}(\mathbf{w}_j^T \mathbf{w}_j) \right] \cdot \gamma + c
 \end{aligned}$$

**Thus** we see that  $q(\gamma | \cdot)$  is gamma distributed with

- shape parameter:  $\alpha \leftarrow \alpha_\gamma + \frac{DM}{2}$
- rate parameter:  $\beta \leftarrow \beta_\gamma + \frac{1}{2} \sum_{j=1}^M \mathbb{E}_{q(\mathbf{w}_j)} [\mathbf{w}_j^T \mathbf{w}_j]$

Compare this to the Gibbs sampler:

- shape parameter:  $\alpha \leftarrow \alpha_\gamma + \frac{DM}{2}$
- rate parameter:  $\beta \leftarrow \beta_\gamma + \frac{1}{2} \sum_{j=1}^M \mathbf{w}_j^T \mathbf{w}_j$

**VB uses posterior expectations where Gibbs uses samples!**

We choose the variational distribution so that

$$\log q(\mathbf{w}_j) = \mathbb{E}_{q-\mathbf{w}_j}[\log p(\cdot)] + c$$

We choose the variational distribution so that

$$\log q(\mathbf{w}_j) = \mathbb{E}_{q-\mathbf{w}_j}[\log p(\cdot)] + c$$

## Recall

$$\begin{aligned} \log p(\cdot) \\ = \log p(\gamma) + \sum_{i=1}^N \log p(\mathbf{z}_i) + \sum_{j=1}^M [\log p(\mathbf{w}_j | \gamma) + \log p(\theta_j)] + \sum_{i=1}^N \sum_{j=1}^M \log p(x_{ij} | \mathbf{w}_j, \mathbf{z}_i, \theta_j) \end{aligned}$$

We choose the variational distribution so that

$$\log q(\mathbf{w}_j) = \mathbb{E}_{q-\mathbf{w}_j}[\log p(\cdot)] + c$$

## Recall

$$\begin{aligned} \log p(\cdot) \\ = \log p(\gamma) + \sum_{i=1}^N \log p(\mathbf{z}_i) + \sum_{j=1}^M [\log p(\mathbf{w}_j | \gamma) + \log p(\theta_j)] + \sum_{i=1}^N \sum_{j=1}^M \log p(x_{ij} | \mathbf{w}_j, \mathbf{z}_i, \theta_j) \end{aligned}$$

We choose the variational distribution so that

$$\log q(\mathbf{w}_j) = \mathbb{E}_{q-\mathbf{w}_j}[\log p(\cdot)] + c = \mathbb{E} \log p(\mathbf{w}_j | \gamma) + \sum_{i=1}^N \mathbb{E} \log p(x_{ij} | \mathbf{w}_j, \mathbf{z}_i, \theta_j) + c$$

## Recall

$$\begin{aligned} \log p(\cdot) \\ = \log p(\gamma) + \sum_{i=1}^N \log p(\mathbf{z}_i) + \sum_{j=1}^M [\log p(\mathbf{w}_j | \gamma) + \log p(\theta_j)] + \sum_{i=1}^N \sum_{j=1}^M \log p(x_{ij} | \mathbf{w}_j, \mathbf{z}_i, \theta_j) \end{aligned}$$

We choose the variational distribution so that

$$\log q(\mathbf{w}_j) = \mathbb{E}_{q-\mathbf{w}_j}[\log p(\cdot)] + c = \mathbb{E} \log p(\mathbf{w}_j \mid \gamma) + \sum_{i=1}^N \mathbb{E} \log p(x_{ij} \mid \mathbf{w}_j, \mathbf{z}_i, \theta_j) + c$$

## The (multivariate) normal distribution

$$\log p(\mathbf{w}_j \mid \gamma) = \log \mathcal{N}(\mathbf{0}, \gamma^{-1} \mathbf{I}) = -\frac{D}{2} \log(2\pi) + \frac{D}{2} \log(\gamma) - \frac{\gamma}{2} \mathbf{w}_j^T \mathbf{w}_j$$

$$\log p(x_{ij} \mid \mathbf{w}_j, \mathbf{z}_i, \theta_j) = \log \mathcal{N}(\mathbf{w}_j^T \mathbf{z}_i, \theta_j) = -\frac{1}{2} \log(2\pi) + \frac{1}{2} \log(\theta_j) - \frac{\theta_j}{2} (x_{ij} - \mathbf{w}_j^T \mathbf{z}_i)^2$$

We choose the variational distribution so that

$$\log q(\mathbf{w}_j) = \mathbb{E}_{q-\mathbf{w}_j}[\log p(\cdot)] + c = \mathbb{E} \log p(\mathbf{w}_j \mid \gamma) + \sum_{i=1}^N \mathbb{E} \log p(x_{ij} \mid \mathbf{w}_j, \mathbf{z}_i, \theta_j) + c$$

## The (multivariate) normal distribution

$$\log p(\mathbf{w}_j \mid \gamma) = \log \mathcal{N}(\mathbf{0}, \gamma^{-1} \mathbf{I}) = -\frac{D}{2} \log(2\pi) + \frac{D}{2} \log(\gamma) - \frac{\gamma}{2} \mathbf{w}_j^T \mathbf{w}_j$$

$$\log p(x_{ij} \mid \mathbf{w}_j, \mathbf{z}_i, \theta_j) = \log \mathcal{N}(\mathbf{w}_j^T \mathbf{z}_i, \theta_j) = -\frac{1}{2} \log(2\pi) + \frac{1}{2} \log(\theta_j) - \frac{\theta_j}{2} (x_{ij} - \mathbf{w}_j^T \mathbf{z}_i)^2$$

We choose the variational distribution so that

$$\begin{aligned}\log q(\mathbf{w}_j) &= \mathbb{E}_{q-\mathbf{w}_j}[\log p(\cdot)] + c = \mathbb{E} \log p(\mathbf{w}_j \mid \gamma) + \sum_{i=1}^N \mathbb{E} \log p(x_{ij} \mid \mathbf{w}_j, \mathbf{z}_i, \theta_j) + c \\ &= -\frac{\mathbb{E}(\gamma)}{2} \mathbf{w}_j^T \mathbf{w}_j - \frac{\mathbb{E}(\theta_j)}{2} \sum_{i=1}^N \mathbb{E}(x_{ij} - \mathbf{w}_j^T \mathbf{z}_i)^2 + c\end{aligned}$$

## The (multivariate) normal distribution

$$\begin{aligned}\log p(\mathbf{w}_j \mid \gamma) &= \log \mathcal{N}(\mathbf{0}, \gamma^{-1} \mathbf{I}) = -\frac{D}{2} \log(2\pi) + \frac{D}{2} \log(\gamma) - \frac{\gamma}{2} \mathbf{w}_j^T \mathbf{w}_j \\ \log p(x_{ij} \mid \mathbf{w}_j, \mathbf{z}_i, \theta_j) &= \log \mathcal{N}(\mathbf{w}_j^T \mathbf{z}_i, \theta_j) = -\frac{1}{2} \log(2\pi) + \frac{1}{2} \log(\theta_j) - \frac{\theta_j}{2} (x_{ij} - \mathbf{w}_j^T \mathbf{z}_i)^2\end{aligned}$$



We choose the variational distribution so that

$$\begin{aligned}\log q(\mathbf{w}_j) &= \mathbb{E}_{q-\mathbf{w}_j}[\log p(\cdot)] + c = \mathbb{E} \log p(\mathbf{w}_j \mid \gamma) + \sum_{i=1}^N \mathbb{E} \log p(x_{ij} \mid \mathbf{w}_j, \mathbf{z}_i, \theta_j) + c \\ &= -\frac{\mathbb{E}(\gamma)}{2} \mathbf{w}_j^T \mathbf{w}_j - \frac{\mathbb{E}(\theta_j)}{2} \sum_{i=1}^N \mathbb{E}(x_{ij} - \mathbf{w}_j^T \mathbf{z}_i)^2 + c\end{aligned}$$

## Expanding the square

$$\begin{aligned}(x_{ij} - \mathbf{w}_j^T \mathbf{z}_i)^2 &= x_{ij}^2 - 2x_{ij} \mathbf{w}_j^T \mathbf{z}_i + \mathbf{w}_j^T \mathbf{z}_i \mathbf{z}_i^T \mathbf{w}_j \\ &= x_{ij}^2 - 2x_{ij} \mathbf{w}_j^T \mathbf{z}_i + \mathbf{w}_j^T \mathbf{z}_i \mathbf{z}_i^T \mathbf{w}_j\end{aligned}$$

We choose the variational distribution so that

$$\begin{aligned}\log q(\mathbf{w}_j) &= \mathbb{E}_{q-\mathbf{w}_j}[\log p(\cdot)] + c = \mathbb{E} \log p(\mathbf{w}_j \mid \gamma) + \sum_{i=1}^N \mathbb{E} \log p(x_{ij} \mid \mathbf{w}_j, \mathbf{z}_i, \theta_j) + c \\&= -\frac{\mathbb{E}(\gamma)}{2} \mathbf{w}_j^T \mathbf{w}_j - \frac{\mathbb{E}(\theta_j)}{2} \sum_{i=1}^N \mathbb{E}(x_{ij} - \mathbf{w}_j^T \mathbf{z}_i)^2 + c \\&= -\frac{\mathbb{E}(\gamma)}{2} \mathbf{w}_j^T \mathbf{w}_j - \mathbb{E}(\theta_j) \frac{1}{2} \sum_{i=1}^N [-2x_{ij} \mathbf{w}_j^T \mathbb{E}(\mathbf{z}_i) + \mathbf{w}_j^T \mathbb{E}(\mathbf{z}_i \mathbf{z}_i^T) \mathbf{w}_j] + c\end{aligned}$$

## Expanding the square

$$\begin{aligned}(x_{ij} - \mathbf{w}_j^T \mathbf{z}_i)^2 &= x_{ij}^2 - 2x_{ij} \mathbf{w}_j^T \mathbf{z}_i + \mathbf{w}_j^T \mathbf{z}_i \mathbf{z}_i^T \mathbf{w}_j \\&= x_{ij}^2 - 2x_{ij} \mathbf{w}_j^T \mathbf{z}_i + \mathbf{w}_j^T \mathbf{z}_i \mathbf{z}_i^T \mathbf{w}_j\end{aligned}$$

We choose the variational distribution so that

$$\begin{aligned}
 \log q(\mathbf{w}_j) &= \mathbb{E}_{q-\mathbf{w}_j}[\log p(\cdot)] + c = \mathbb{E} \log p(\mathbf{w}_j \mid \gamma) + \sum_{i=1}^N \mathbb{E} \log p(x_{ij} \mid \mathbf{w}_j, \mathbf{z}_i, \theta_j) + c \\
 &= -\frac{\mathbb{E}(\gamma)}{2} \mathbf{w}_j^T \mathbf{w}_j - \frac{\mathbb{E}(\theta_j)}{2} \sum_{i=1}^N \mathbb{E}(x_{ij} - \mathbf{w}_j^T \mathbf{z}_i)^2 + c \\
 &= -\frac{\mathbb{E}(\gamma)}{2} \mathbf{w}_j^T \mathbf{w}_j - \mathbb{E}(\theta_j) \frac{1}{2} \sum_{i=1}^N [-2x_{ij} \mathbf{w}_j^T \mathbb{E}(\mathbf{z}_i) + \mathbf{w}_j^T \mathbb{E}(\mathbf{z}_i \mathbf{z}_i^T) \mathbf{w}_j] + c \\
 &= -\frac{1}{2} \mathbf{w}_j^T \left[ \mathbb{E}(\gamma) \mathbf{I} + \mathbb{E}(\theta_j) \sum_{i=1}^N \mathbb{E}(\mathbf{z}_i \mathbf{z}_i^T) \right] \mathbf{w}_j + \mathbf{w}_j^T \mathbb{E}(\theta_j) \sum_{i=1}^N x_{ij} \mathbb{E}(\mathbf{z}_i) + c
 \end{aligned}$$

We choose the variational distribution so that

$$\begin{aligned}\log q(\mathbf{w}_j) &= \mathbb{E}_{q-\mathbf{w}_j}[\log p(\cdot)] + c = \mathbb{E} \log p(\mathbf{w}_j \mid \gamma) + \sum_{i=1}^N \mathbb{E} \log p(x_{ij} \mid \mathbf{w}_j, \mathbf{z}_i, \theta_j) + c \\ &= -\frac{1}{2} \mathbf{w}_j^T \left[ \mathbb{E}(\gamma) \mathbf{I} + \mathbb{E}(\theta_j) \sum_{i=1}^N \mathbb{E}(\mathbf{Z}_i \mathbf{Z}_i^T) \right] \mathbf{w}_j + \mathbf{w}_j^T \mathbb{E}(\theta_j) \sum_{i=1}^N x_{ij} \mathbb{E}(\mathbf{Z}_i) + c\end{aligned}$$

We choose the variational distribution so that

$$\begin{aligned}\log q(\mathbf{w}_j) &= \mathbb{E}_{q-\mathbf{w}_j}[\log p(\cdot)] + c = \mathbb{E} \log p(\mathbf{w}_j \mid \gamma) + \sum_{i=1}^N \mathbb{E} \log p(x_{ij} \mid \mathbf{w}_j, \mathbf{z}_i, \theta_j) + c \\ &= -\frac{1}{2} \mathbf{w}_j^T \left[ \mathbb{E}(\gamma) \mathbf{I} + \mathbb{E}(\theta_j) \sum_{i=1}^N \mathbb{E}(\mathbf{Z}_i \mathbf{Z}_i^T) \right] \mathbf{w}_j + \mathbf{w}_j^T \mathbb{E}(\theta_j) \sum_{i=1}^N x_{ij} \mathbb{E}(\mathbf{Z}_i) + c\end{aligned}$$

Recall the (generic) multivariate normal distribution

$$\begin{aligned}\log p(\mathbf{y} \mid \boldsymbol{\mu}, \mathbf{Q}) &= \log \mathcal{N}(\mathbf{y} \mid \boldsymbol{\mu}, \mathbf{Q}^{-1}) = -\frac{D}{2} \log(2\pi) + \frac{1}{2} \log |\mathbf{Q}| - \frac{1}{2} (\mathbf{y} - \boldsymbol{\mu})^T \mathbf{Q} (\mathbf{y} - \boldsymbol{\mu}) \\ &= -\frac{D}{2} \log(2\pi) + \frac{1}{2} \log |\mathbf{Q}| - \frac{1}{2} \mathbf{y}^T \mathbf{Q} \mathbf{y} + \mathbf{y}^T \mathbf{Q} \boldsymbol{\mu} - \frac{1}{2} \boldsymbol{\mu}^T \mathbf{Q} \boldsymbol{\mu}\end{aligned}$$

We choose the variational distribution so that

$$\begin{aligned}\log q(\mathbf{w}_j) &= \mathbb{E}_{q-\mathbf{w}_j}[\log p(\cdot)] + c = \mathbb{E} \log p(\mathbf{w}_j \mid \gamma) + \sum_{i=1}^N \mathbb{E} \log p(x_{ij} \mid \mathbf{w}_j, \mathbf{z}_i, \theta_j) + c \\ &= -\frac{1}{2} \mathbf{w}_j^\top \left[ \mathbb{E}(\gamma) \mathbf{I} + \mathbb{E}(\theta_j) \sum_{i=1}^N \mathbb{E}(\mathbf{Z}_i \mathbf{Z}_i^\top) \right] \mathbf{w}_j + \mathbf{w}_j^\top \mathbb{E}(\theta_j) \sum_{i=1}^N x_{ij} \mathbb{E}(\mathbf{Z}_i) + c\end{aligned}$$

**Thus** we see that  $q(\mathbf{w}_j \mid \cdot)$  is normally distributed with

- precision  $\mathbf{Q} \leftarrow \mathbb{E}(\gamma) \mathbf{I} + \mathbb{E}(\theta_j) \sum_{i=1}^N \mathbb{E}(\mathbf{Z}_i \mathbf{Z}_i^\top)$
- mean  $\boldsymbol{\mu} \leftarrow \mathbf{Q}^{-1} \left[ \mathbb{E}(\theta_j) \sum_{i=1}^N x_{ij} \mathbb{E}(\mathbf{Z}_i) \right]$

Recall the (generic) multivariate normal distribution

$$\begin{aligned}\log p(\mathbf{y} \mid \boldsymbol{\mu}, \mathbf{Q}) &= \log \mathcal{N}(\mathbf{y} \mid \boldsymbol{\mu}, \mathbf{Q}^{-1}) = -\frac{D}{2} \log(2\pi) + \frac{1}{2} \log |\mathbf{Q}| - \frac{1}{2} (\mathbf{y} - \boldsymbol{\mu})^\top \mathbf{Q} (\mathbf{y} - \boldsymbol{\mu}) \\ &= -\frac{D}{2} \log(2\pi) + \frac{1}{2} \log |\mathbf{Q}| - \frac{1}{2} \mathbf{y}^\top \mathbf{Q} \mathbf{y} + \mathbf{y}^\top \mathbf{Q} \boldsymbol{\mu} - \frac{1}{2} \boldsymbol{\mu}^\top \mathbf{Q} \boldsymbol{\mu}\end{aligned}$$

We choose the variational distribution so that

$$\begin{aligned}\log q(\mathbf{w}_j) &= \mathbb{E}_{q-\mathbf{w}_j}[\log p(\cdot)] + c = \mathbb{E} \log p(\mathbf{w}_j \mid \gamma) + \sum_{i=1}^N \mathbb{E} \log p(x_{ij} \mid \mathbf{w}_j, \mathbf{z}_i, \theta_j) + c \\ &= -\frac{1}{2} \mathbf{w}_j^T \left[ \mathbb{E}(\gamma) \mathbf{I} + \mathbb{E}(\theta_j) \sum_{i=1}^N \mathbb{E}(\mathbf{Z}_i \mathbf{Z}_i^T) \right] \mathbf{w}_j + \mathbf{w}_j^T \mathbb{E}(\theta_j) \sum_{i=1}^N x_{ij} \mathbb{E}(\mathbf{Z}_i) + c\end{aligned}$$

**Thus** we see that  $q(\mathbf{w}_j \mid \cdot)$  is normally distributed with

- precision  $\mathbf{Q} \leftarrow \mathbb{E}(\gamma) \mathbf{I} + \mathbb{E}(\theta_j) \sum_{i=1}^N \mathbb{E}(\mathbf{Z}_i \mathbf{Z}_i^T)$
- mean  $\boldsymbol{\mu} \leftarrow \mathbf{Q}^{-1} \left[ \mathbb{E}(\theta_j) \sum_{i=1}^N x_{ij} \mathbb{E}(\mathbf{Z}_i) \right]$

Calculation of  $\mathbb{E}(\mathbf{Z}_i \mathbf{Z}_i^T)$

$$\mathbb{E}(\mathbf{Z}_i \mathbf{Z}_i^T) = \text{Cov}(\mathbf{Z}_i) + \mathbb{E}(\mathbf{Z}_i) \mathbb{E}(\mathbf{Z}_i)^T$$

We choose the variational distribution so that

$$\begin{aligned}\log q(\mathbf{w}_j) &= \mathbb{E}_{q-\mathbf{w}_j}[\log p(\cdot)] + c = \mathbb{E} \log p(\mathbf{w}_j \mid \gamma) + \sum_{i=1}^N \mathbb{E} \log p(x_{ij} \mid \mathbf{w}_j, \mathbf{z}_i, \theta_j) + c \\ &= -\frac{1}{2} \mathbf{w}_j^T \left[ \mathbb{E}(\gamma) \mathbf{I} + \mathbb{E}(\theta_j) \sum_{i=1}^N \mathbb{E}(\mathbf{Z}_i \mathbf{Z}_i^T) \right] \mathbf{w}_j + \mathbf{w}_j^T \mathbb{E}(\theta_j) \sum_{i=1}^N x_{ij} \mathbb{E}(\mathbf{Z}_i) + c\end{aligned}$$

**Thus** we see that  $q(\mathbf{w}_j \mid \cdot)$  is normally distributed with

- precision  $\mathbf{Q} \leftarrow \mathbb{E}(\gamma) \mathbf{I} + \mathbb{E}(\theta_j) \sum_{i=1}^N \mathbb{E}(\mathbf{Z}_i \mathbf{Z}_i^T)$
- mean  $\boldsymbol{\mu} \leftarrow \mathbf{Q}^{-1} \left[ \mathbb{E}(\theta_j) \sum_{i=1}^N x_{ij} \mathbb{E}(\mathbf{Z}_i) \right]$

Compare this to the Gibbs sampler:

- $\mathbf{Q} \leftarrow \gamma \mathbf{I} + \theta_j \sum_{i=1}^N \mathbf{z}_i \mathbf{z}_i^T$
- $\boldsymbol{\mu} \leftarrow \mathbf{Q}^{-1} \left[ \theta_j \sum_{i=1}^N x_{ij} \mathbf{z}_i \right]$

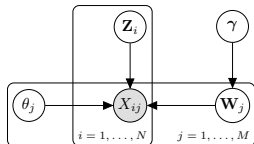
Once again, the only difference between VB and Gibbs is that where VB uses posterior expectations, Gibbs uses samples.



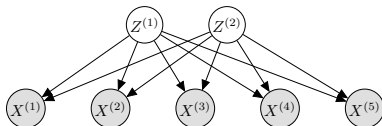
## Data

100 data points was randomly sampled from a 5-dim multivariate Gaussian distribution.

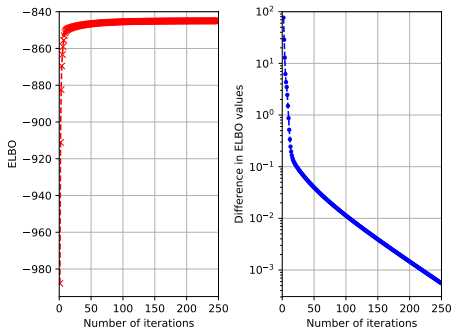
## Global model



## Local model



## Monitoring convergence



## Data

100 data points was randomly sampled from a 5-dim multivariate Gaussian distribution.

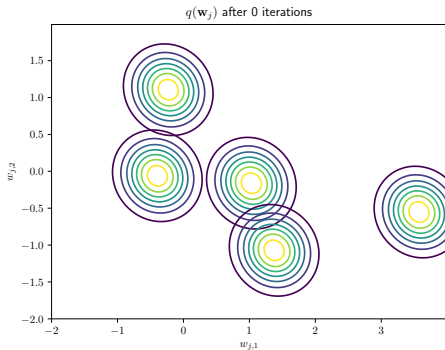
## Global model



## Local model



## Variational posteriors



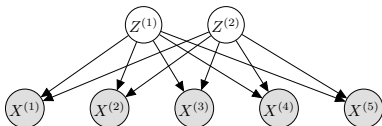
## Data

100 data points was randomly sampled from a 5-dim multivariate Gaussian distribution.

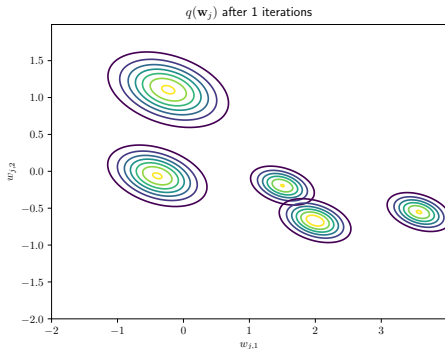
## Global model



## Local model



## Variational posteriors



## Data

100 data points were randomly sampled from a 5-dim multivariate Gaussian distribution.

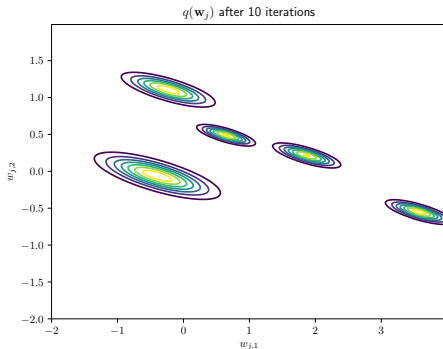
## Global model



## Local model



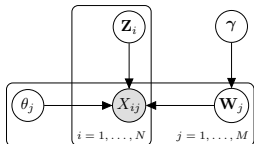
## Variational posteriors



## Data

100 data points was randomly sampled from a 5-dim multivariate Gaussian distribution.

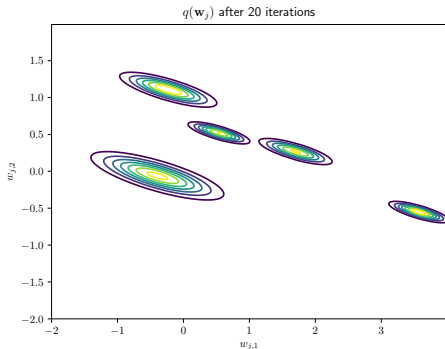
## Global model



## Local model



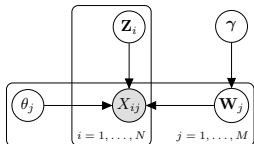
## Variational posteriors



## Data

100 data points was randomly sampled from a 5-dim multivariate Gaussian distribution.

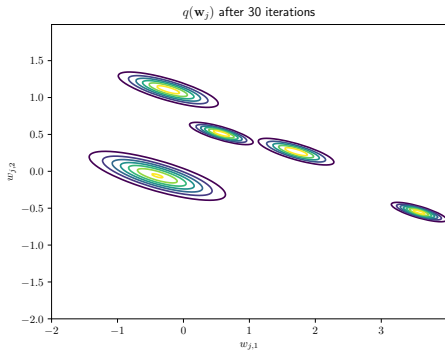
## Global model



## Local model



## Variational posteriors



## Data

100 data points was randomly sampled from a 5-dim multivariate Gaussian distribution.

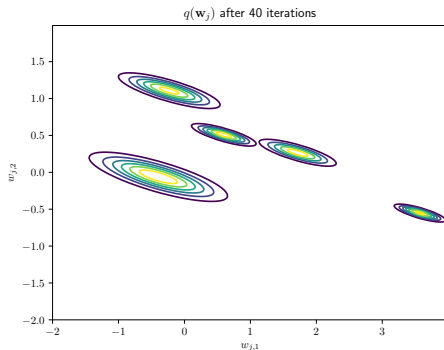
## Global model



## Local model



## Variational posteriors



## Data

100 data points was randomly sampled from a 5-dim multivariate Gaussian distribution.

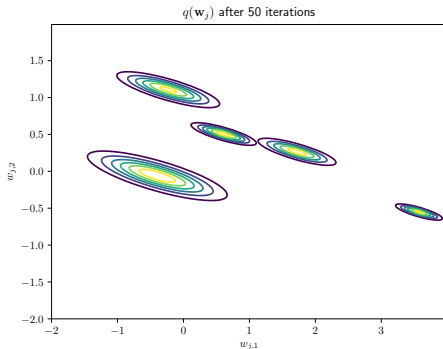
## Global model



## Local model



## Variational posteriors





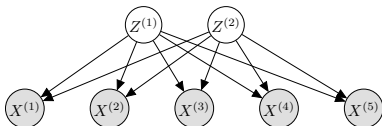
## Data

100 data points was randomly sampled from a 5-dim multivariate Gaussian distribution.

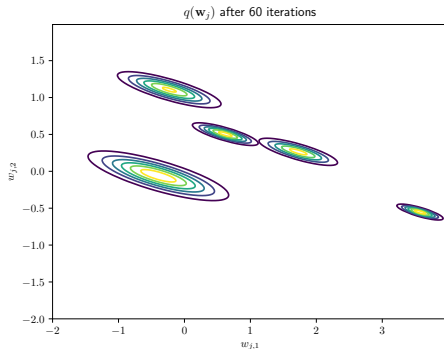
## Global model



## Local model



## Variational posteriors



## Data

100 data points was randomly sampled from a 5-dim multivariate Gaussian distribution.

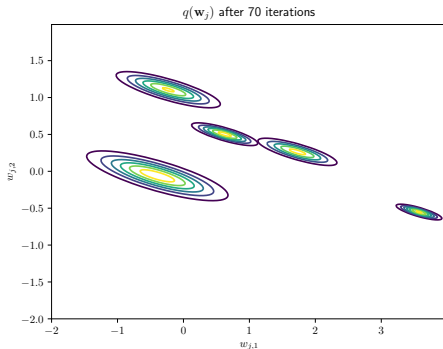
## Global model



## Local model



## Variational posteriors



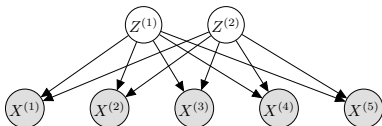
## Data

100 data points was randomly sampled from a 5-dim multivariate Gaussian distribution.

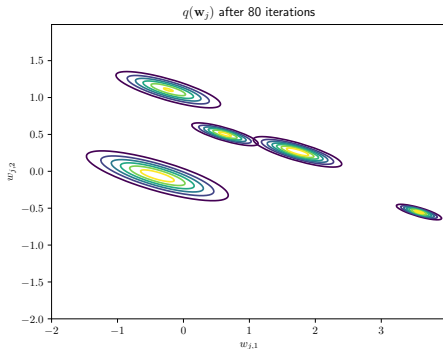
## Global model



## Local model



## Variational posteriors



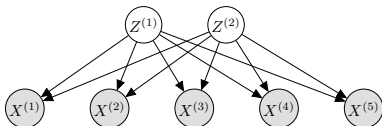
## Data

100 data points was randomly sampled from a 5-dim multivariate Gaussian distribution.

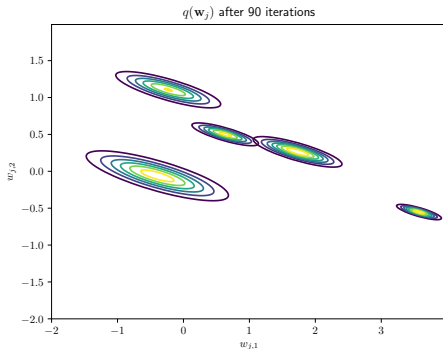
## Global model



## Local model



## Variational posteriors



## Data

100 data points was randomly sampled from a 5-dim multivariate Gaussian distribution.

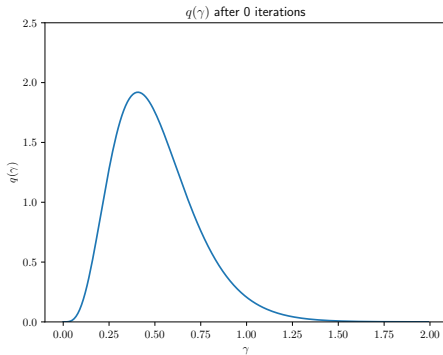
## Global model



## Local model



## Variational posteriors



## Data

100 data points was randomly sampled from a 5-dim multivariate Gaussian distribution.

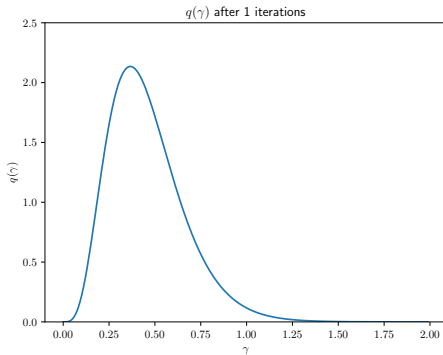
## Global model



## Local model



## Variational posteriors



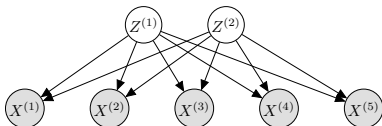
## Data

100 data points was randomly sampled from a 5-dim multivariate Gaussian distribution.

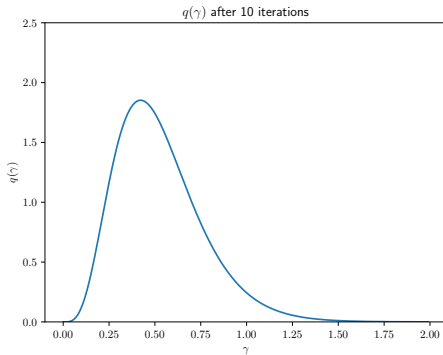
## Global model



## Local model



## Variational posteriors



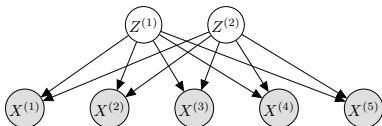
## Data

100 data points was randomly sampled from a 5-dim multivariate Gaussian distribution.

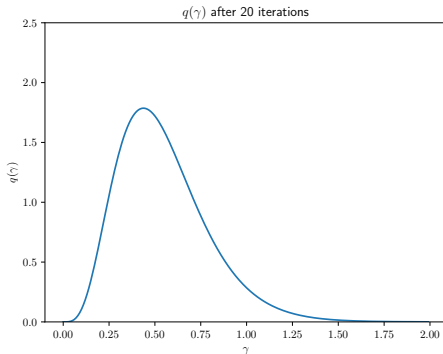
## Global model



## Local model



## Variational posteriors

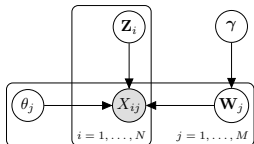




## Data

100 data points was randomly sampled from a 5-dim multivariate Gaussian distribution.

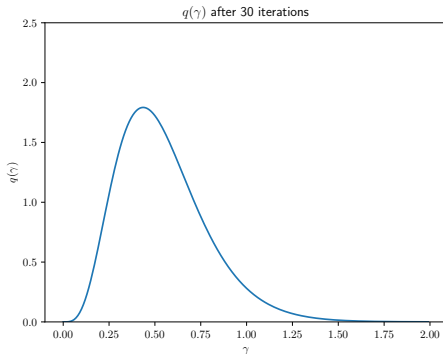
## Global model



## Local model



## Variational posteriors



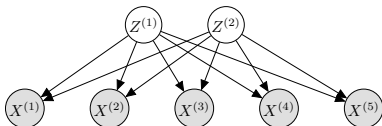
## Data

100 data points was randomly sampled from a 5-dim multivariate Gaussian distribution.

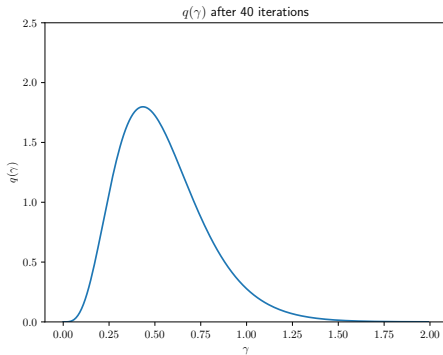
## Global model



## Local model



## Variational posteriors



## Data

100 data points was randomly sampled from a 5-dim multivariate Gaussian distribution.

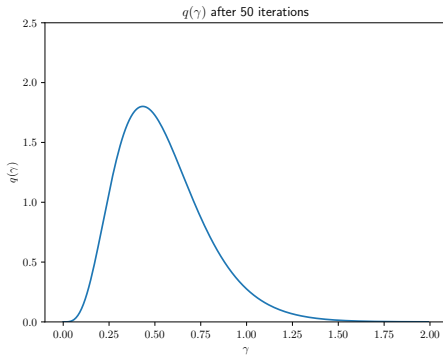
## Global model



## Local model



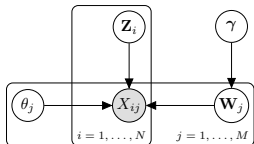
## Variational posteriors



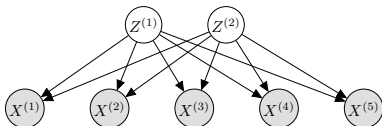
## Data

100 data points was randomly sampled from a 5-dim multivariate Gaussian distribution.

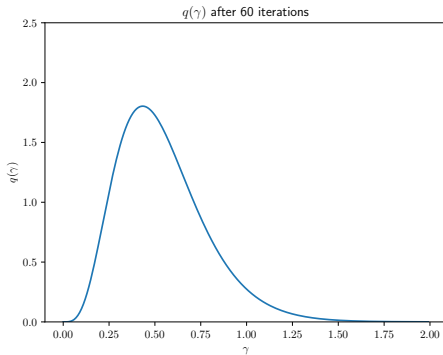
## Global model



## Local model



## Variational posteriors



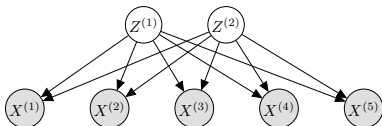
## Data

100 data points was randomly sampled from a 5-dim multivariate Gaussian distribution.

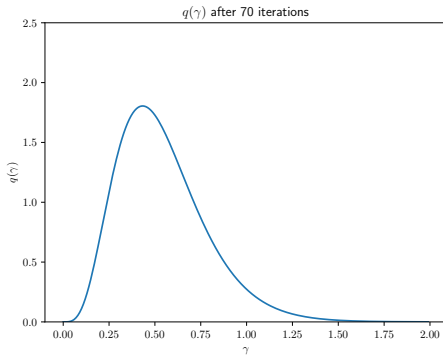
## Global model



## Local model



## Variational posteriors



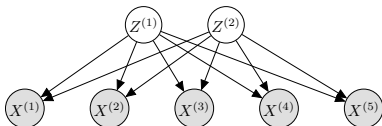
## Data

100 data points was randomly sampled from a 5-dim multivariate Gaussian distribution.

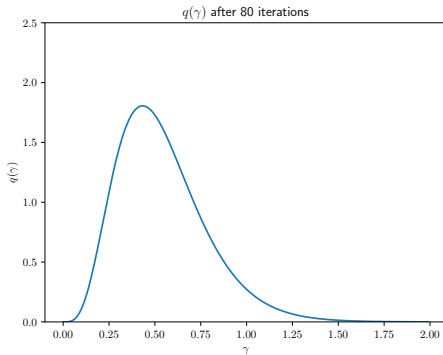
## Global model



## Local model



## Variational posteriors



## Data

100 data points was randomly sampled from a 5-dim multivariate Gaussian distribution.

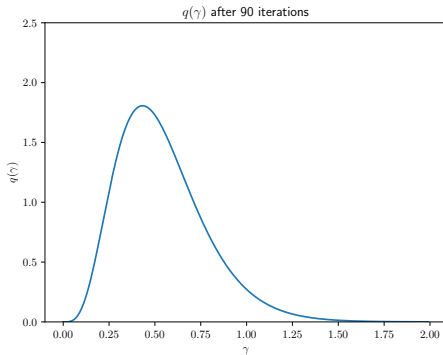
## Global model



## Local model



## Variational posteriors



## Data

100 data points was randomly sampled from a 5-dim multivariate Gaussian distribution.

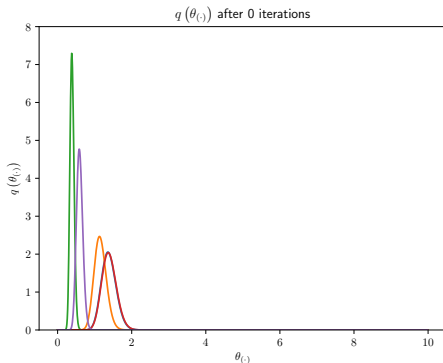
### Global model



### Local model



## Variational posteriors





## Data

100 data points was randomly sampled from a 5-dim multivariate Gaussian distribution.

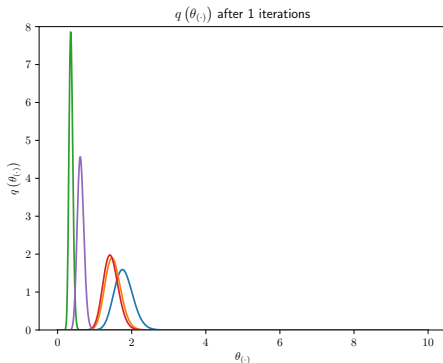
### Global model



### Local model



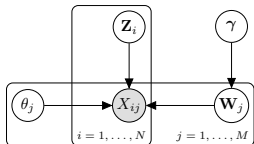
## Variational posteriors



## Data

100 data points was randomly sampled from a 5-dim multivariate Gaussian distribution.

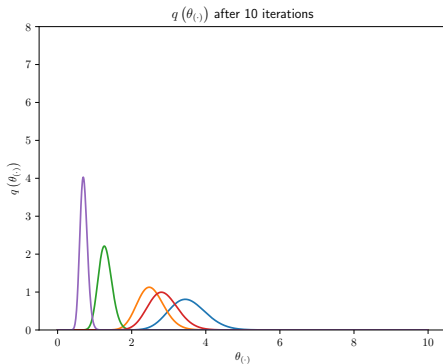
### Global model



### Local model



## Variational posteriors



## Data

100 data points was randomly sampled from a 5-dim multivariate Gaussian distribution.

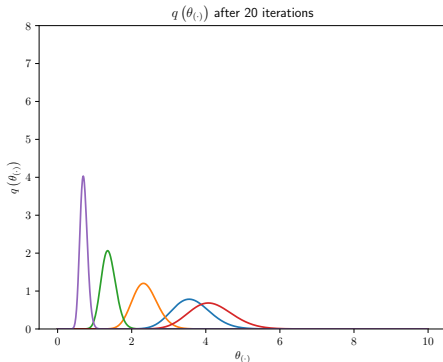
## Global model



## Local model



## Variational posteriors



## Data

100 data points was randomly sampled from a 5-dim multivariate Gaussian distribution.

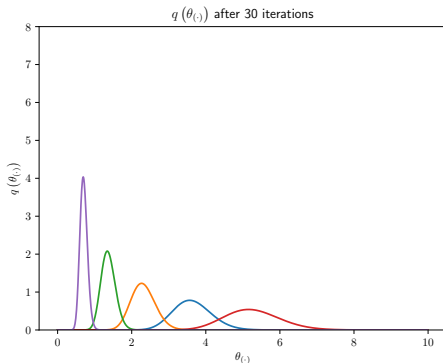
### Global model



### Local model



## Variational posteriors



## Data

100 data points was randomly sampled from a 5-dim multivariate Gaussian distribution.

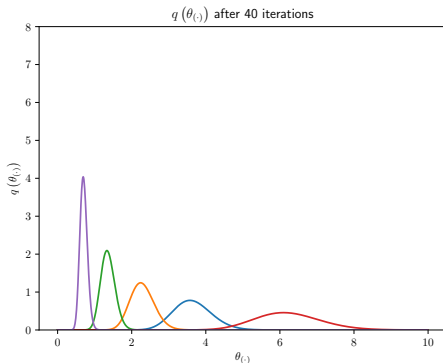
## Global model



## Local model



## Variational posteriors



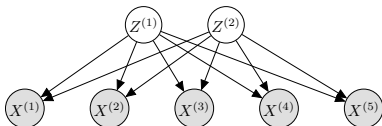
## Data

100 data points was randomly sampled from a 5-dim multivariate Gaussian distribution.

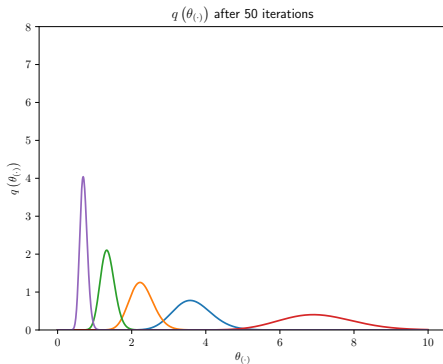
## Global model



## Local model



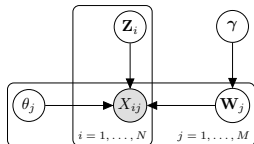
## Variational posteriors



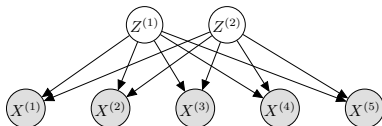
## Data

100 data points was randomly sampled from a 5-dim multivariate Gaussian distribution.

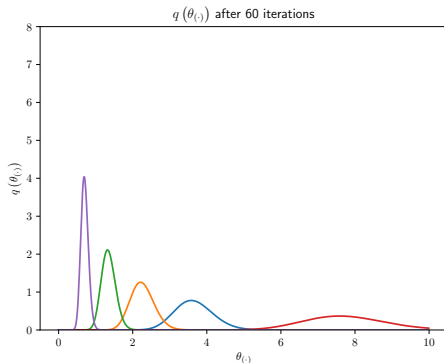
## Global model



## Local model



## Variational posteriors



## Data

100 data points was randomly sampled from a 5-dim multivariate Gaussian distribution.

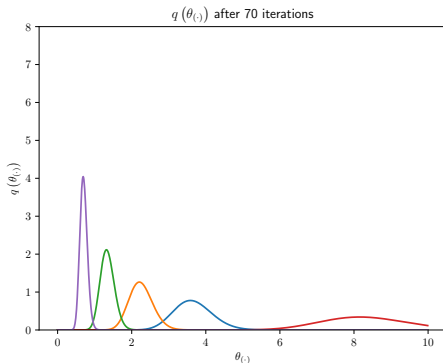
### Global model



### Local model



## Variational posteriors

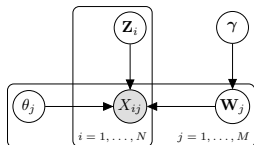




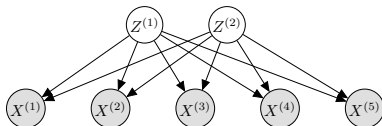
## Data

100 data points was randomly sampled from a 5-dim multivariate Gaussian distribution.

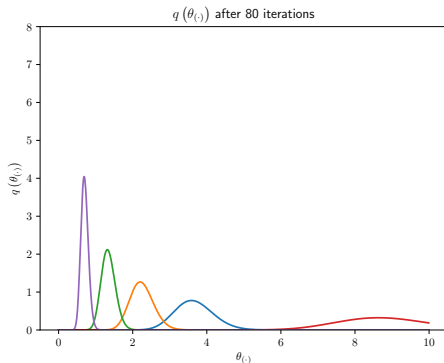
## Global model



## Local model



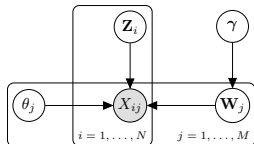
## Variational posteriors



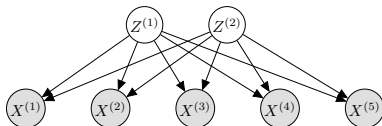
## Data

100 data points was randomly sampled from a 5-dim multivariate Gaussian distribution.

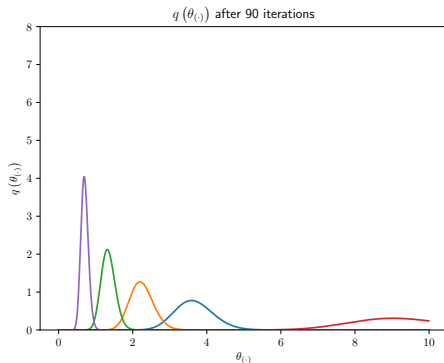
## Global model



## Local model



## Variational posteriors



## Data

100 data points was randomly sampled from a 5-dim multivariate Gaussian distribution.

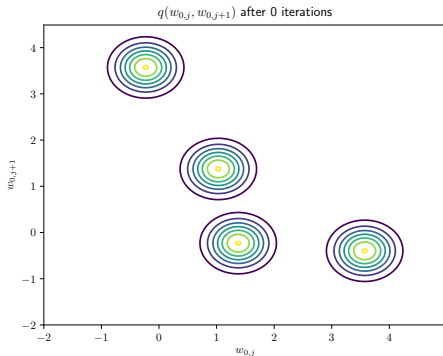
## Global model



## Local model



## Variational posteriors



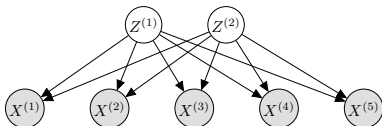
## Data

100 data points was randomly sampled from a 5-dim multivariate Gaussian distribution.

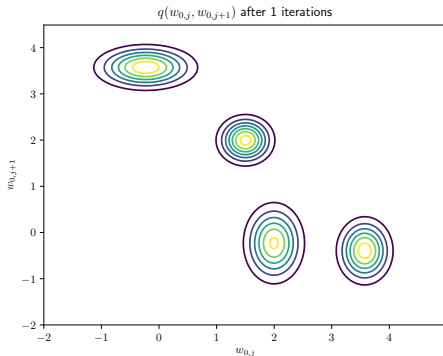
## Global model



## Local model



## Variational posteriors



## Data

100 data points was randomly sampled from a 5-dim multivariate Gaussian distribution.

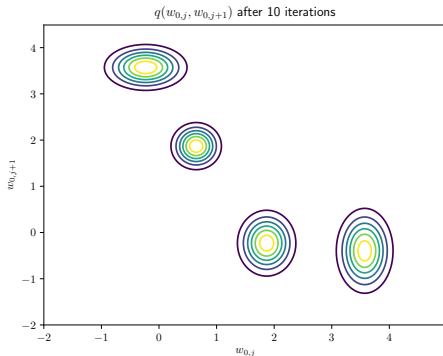
## Global model



## Local model



## Variational posteriors



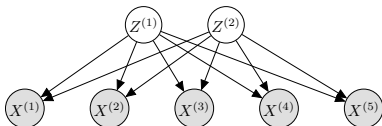
## Data

100 data points was randomly sampled from a 5-dim multivariate Gaussian distribution.

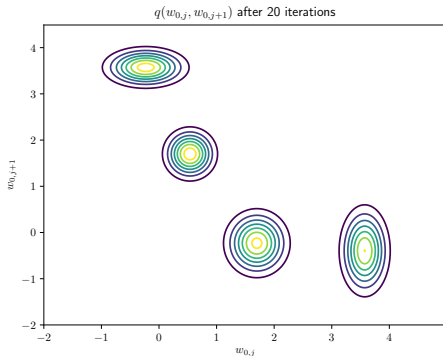
## Global model



## Local model



## Variational posteriors



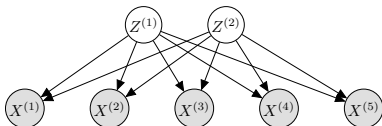
## Data

100 data points was randomly sampled from a 5-dim multivariate Gaussian distribution.

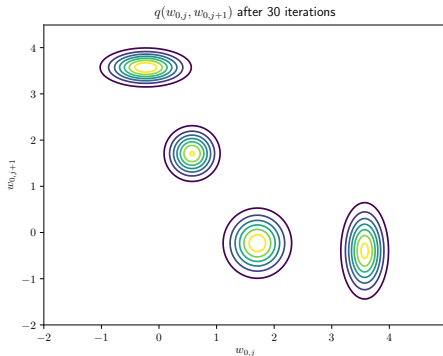
## Global model



## Local model



## Variational posteriors



## Data

100 data points was randomly sampled from a 5-dim multivariate Gaussian distribution.

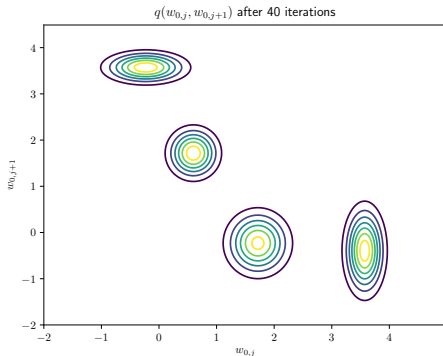
## Global model



## Local model



## Variational posteriors





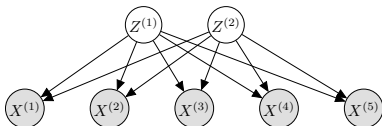
## Data

100 data points was randomly sampled from a 5-dim multivariate Gaussian distribution.

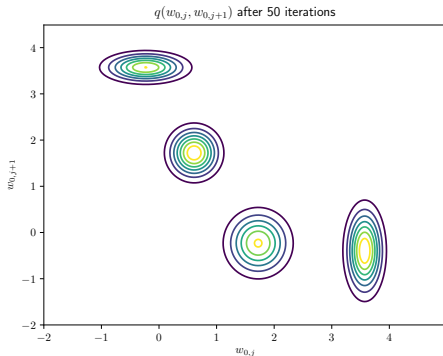
## Global model



## Local model



## Variational posteriors



## Data

100 data points was randomly sampled from a 5-dim multivariate Gaussian distribution.

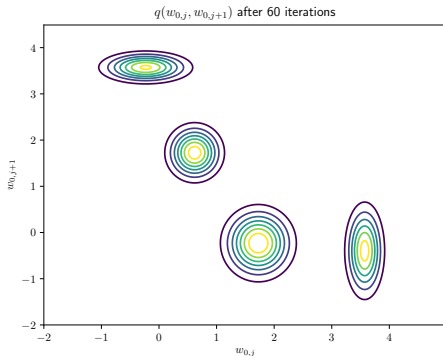
## Global model



## Local model



## Variational posteriors



## Data

100 data points was randomly sampled from a 5-dim multivariate Gaussian distribution.

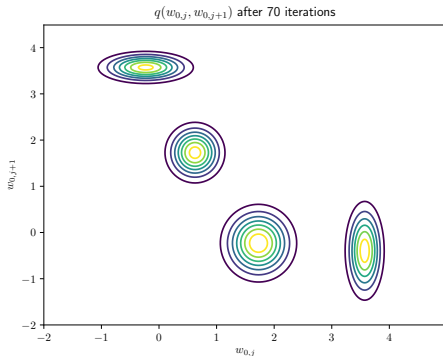
## Global model



## Local model



## Variational posteriors



## Data

100 data points was randomly sampled from a 5-dim multivariate Gaussian distribution.

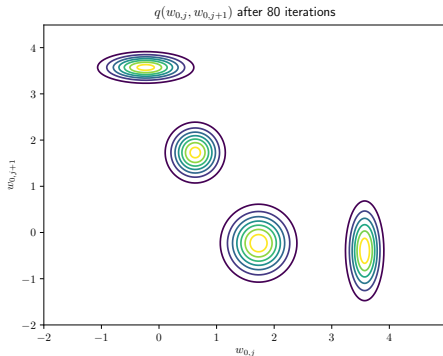
## Global model



## Local model



## Variational posteriors



## Data

100 data points was randomly sampled from a 5-dim multivariate Gaussian distribution.

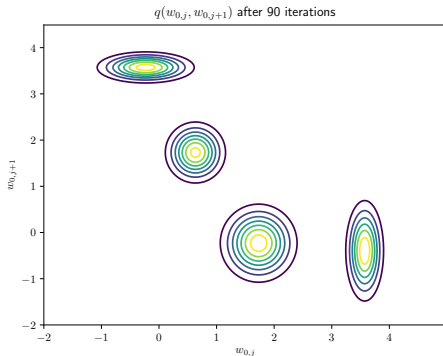
## Global model



## Local model



## Variational posteriors



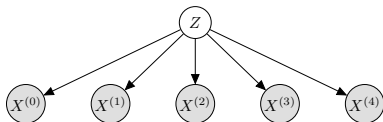
## Data

100 data points was randomly sampled from a 5-dim multivariate Gaussian distribution.

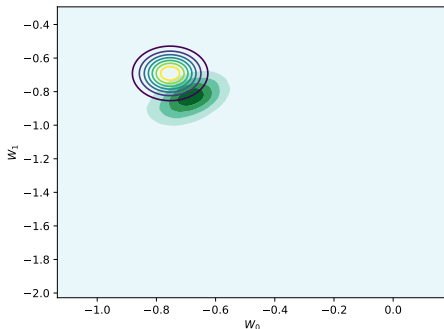
## Global model



## Local model



## Comparison with Gibbs sampling



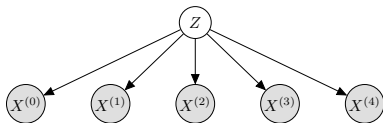
## Data

100 data points was randomly sampled from a 5-dim multivariate Gaussian distribution.

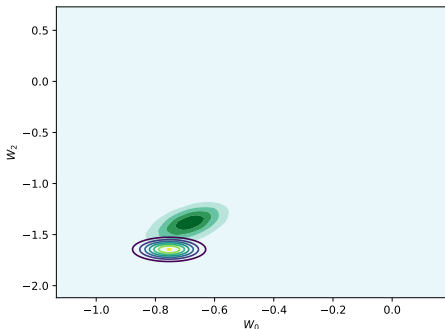
## Global model



## Local model



## Comparison with Gibbs sampling



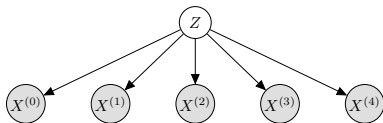
## Data

100 data points was randomly sampled from a 5-dim multivariate Gaussian distribution.

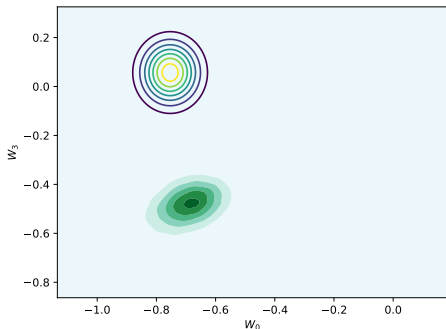
## Global model



## Local model



## Comparison with Gibbs sampling

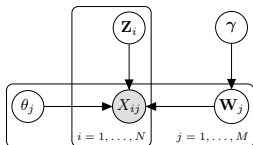




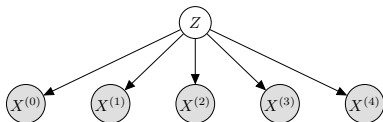
## Data

100 data points were randomly sampled from a 5-dim multivariate Gaussian distribution.

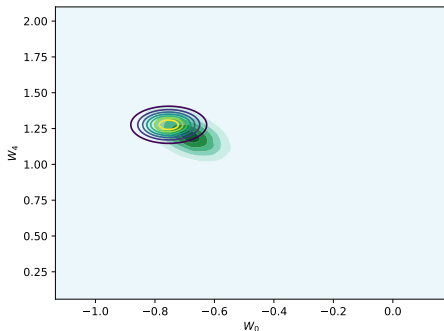
## Global model



## Local model



## Comparison with Gibbs sampling



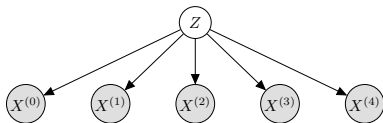
## Data

100 data points was randomly sampled from a 5-dim multivariate Gaussian distribution.

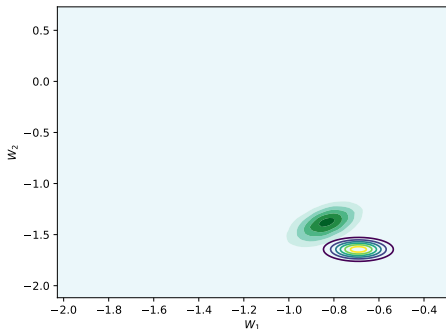
## Global model



## Local model



## Comparison with Gibbs sampling



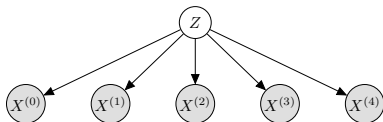
## Data

100 data points was randomly sampled from a 5-dim multivariate Gaussian distribution.

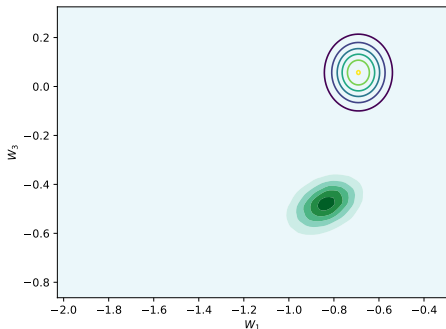
## Global model



## Local model



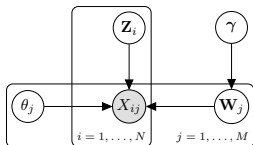
## Comparison with Gibbs sampling



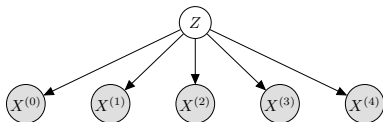
## Data

100 data points was randomly sampled from a 5-dim multivariate Gaussian distribution.

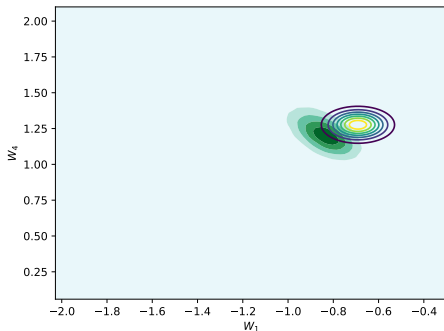
## Global model



## Local model



## Comparison with Gibbs sampling



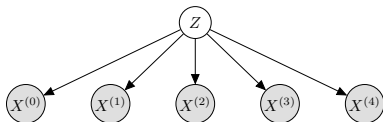
## Data

100 data points was randomly sampled from a 5-dim multivariate Gaussian distribution.

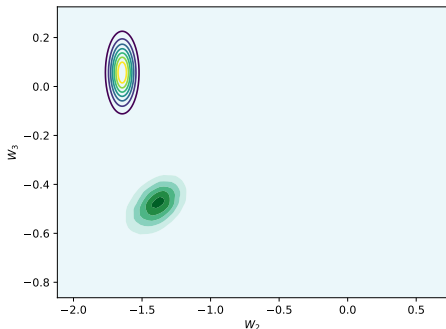
## Global model



## Local model



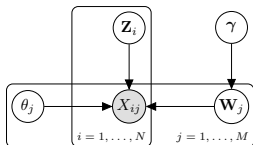
## Comparison with Gibbs sampling



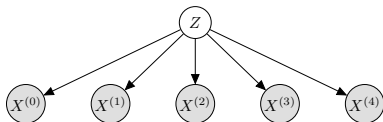
## Data

100 data points was randomly sampled from a 5-dim multivariate Gaussian distribution.

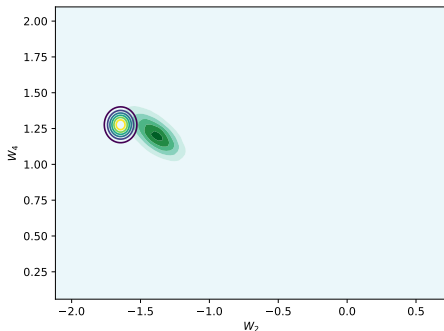
## Global model



## Local model



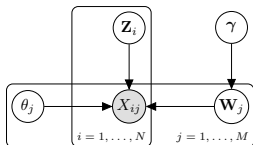
## Comparison with Gibbs sampling



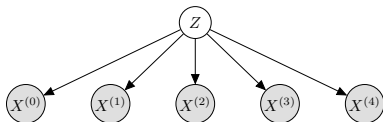
## Data

100 data points was randomly sampled from a 5-dim multivariate Gaussian distribution.

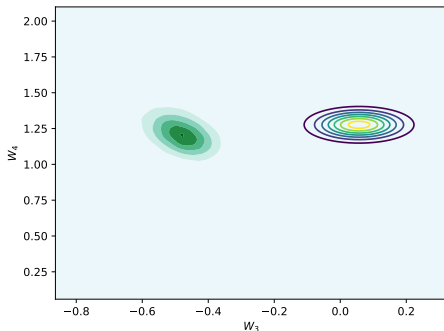
## Global model



## Local model



## Comparison with Gibbs sampling



Not seen from the plot, **but** the results strongly dependent on the VI initialization.

## Algorithm:

- We have observed  $\mathbf{X} = \mathbf{x}$ , and have access to the full joint  $p(\mathbf{z}, \mathbf{x})$ .
- We posit a *variational family* of distributions  $q_j(\cdot \mid \boldsymbol{\lambda}_j)$ , i.e., we choose the distributional form, while wanting to optimize the parameterization  $\boldsymbol{\lambda}_j$ .
- The posterior approximation is assumed to factorize according to the mean-field assumption, and we use the KL ( $q(\mathbf{z}) \parallel p(\mathbf{z} \mid \mathbf{x})$ ) as our objective.

## Algorithm:

Repeat until negligible improvement in terms of  $\mathcal{L}(q)$ :

- For each  $j$ :
  - Calculate  $\mathbb{E}_{q_{-j}} [\log p(\mathbf{z}, \mathbf{x})]$  using current estimates for  $q_i(\cdot \mid \boldsymbol{\lambda}_i)$ ,  $i \neq j$ .
  - Choose  $\boldsymbol{\lambda}_j$  so that  $q_j(z_j \mid \boldsymbol{\lambda}_j) \propto \exp (\mathbb{E}_{q_{-j}} [\log p(\mathbf{z}, \mathbf{x})])$ .
- Calculate the new  $\mathcal{L}(q)$ .



## Algorithm:

- We have observed  $\mathbf{X} = \mathbf{x}$ , and have access to the full joint  $p(\mathbf{z}, \mathbf{x})$ .
- We posit a *variational family* of distributions  $q_j(\cdot \mid \boldsymbol{\lambda}_j)$ , i.e., we choose the distributional form, while wanting to optimize the parameterization  $\boldsymbol{\lambda}_j$ .
- The posterior approximation is assumed to factorize according to the mean-field assumption, and we use the KL ( $q(\mathbf{z}) \parallel p(\mathbf{z} \mid \mathbf{x})$ ) as our objective.

## Algorithm:

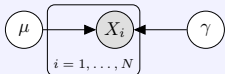
Repeat until negligible improvement in terms of  $\mathcal{L}(q)$ :

- For each  $j$ :
  - Calculate  $\mathbb{E}_{q_{-j}} [\log p(\mathbf{z}, \mathbf{x})]$  using current estimates for  $q_i(\cdot \mid \boldsymbol{\lambda}_i)$ ,  $i \neq j$ .
  - Choose  $\boldsymbol{\lambda}_j$  so that  $q_j(z_j \mid \boldsymbol{\lambda}_j) \propto \exp(\mathbb{E}_{q_{-j}} [\log p(\mathbf{z}, \mathbf{x})])$ .
- Calculate the new  $\mathcal{L}(q)$ .

As we just realized, calculations of  $\mathbb{E}_{q_{-j}} [\log p(\mathbf{z}, \mathbf{x})]$  and  $\mathcal{L}(q)$  are quite tedious – and apparently must be done separately for each model we make.

This **harms the applicability** of variational inference, even under the **quite restrictive** mean field assumption.

## Code Task: VB for the a simple Gaussian model



- $X_i \mid \{\mu, \gamma\} \sim \mathcal{N}(\mu, 1/\gamma)$
- $\mu \sim \mathcal{N}(0, \tau)$
- $\gamma \sim \text{Gamma}(\alpha, \beta)$

In this task you need to use mean-field, and look for  $q(\mu, \gamma) = q(\mu) \cdot q(\gamma)$  that best approximates  $p(\mu, \tau \mid x_1, \dots, x_N)$  wrt. the VB measure  $\text{KL}(q||p)$ .

- Calculate the update rules for  $q(\mu)$  and  $q(\gamma)$ .
  - Hint:  $q(\mu)$  is Gaussian with mean  $\nu^*$  and precision  $\tau^*$ ;  $q(\gamma)$  is Gamma-distributed with parameters  $\alpha^*$  and  $\beta^*$ .
- Implement the update rules you find in the notebook

`students_simple_model.ipynb`

It may be useful to recall the definition of pdfs:

- Gamma:  $\log p(x \mid \alpha, \beta) = \alpha \log(\beta) + (\alpha - 1) \log(x) - \beta \cdot x - \log(\Gamma(\alpha))$ .
- Gauss:  $\log p(x \mid \mu, 1/\gamma) = -\frac{1}{2} \log(2\pi) + \frac{1}{2} \log(\gamma) - \frac{\gamma}{2} (x - \mu)^2$ .

# The Exponential Family of distributions

## Positives with the Exponential Family:

- It is the only family of distributions with finite-sized sufficient statistics\*;
- It is the only family of distributions that has conjugate priors;
- It simplifies the operations of variational inference;
- It has simple mathematical procedures for calculating moments, MLEs, Bayesian posteriors, . . .

\* Under certain regularity conditions. . .

## Negatives:

- Standard distributions are defined using a new parameterization. Can be “unnatural”.
- While all the calculations we will do are in principle simple, they are sometimes a bit abstract – due to the massive generality of the construction.

Consider a univariate distribution  $f_X(x | \theta)$ , written as:

$$f_X(x | \theta) = \exp(h(x) + \boldsymbol{\eta}(\theta)^\top \mathbf{t}(x) - A(\theta))$$

Here we define:

- $h(x)$ : log base measure
- $\boldsymbol{\eta}(\theta)$ : the natural parameters
- $\mathbf{t}(x)$ : the sufficient statistics
- $A(\theta)$ : the log partition function

Consider a univariate distribution  $f_X(x | \theta)$ , written as:

$$f_X(x | \theta) = \exp(h(x) + \boldsymbol{\eta}(\theta)^\top \mathbf{t}(x) - A(\theta))$$

Here we define:

- $h(x)$ : log base measure
- $\boldsymbol{\eta}(\theta)$ : the natural parameters
- $\mathbf{t}(x)$ : the sufficient statistics
- $A(\theta)$ : the log partition function

**Members:**

- Bernoulli (and Multinomial)

**Example:**

- $h(x) = 0$
- $\boldsymbol{\eta}(p) = \log(p/1 - p)$
- $\mathbf{t}(x) = x$
- $A(p) = -\log(1 - p)$

$$\begin{aligned} f_X(x | \theta) &= \exp(h(x) + \boldsymbol{\eta}(\theta)^\top \mathbf{t}(x) - A(\theta)) \\ &= \exp\{\log[p/(1 - p)] \cdot x + \log(1 - p)\} \\ &= [p/(1 - p)]^x \cdot (1 - p) \\ &= p^x (1 - p)^{1-x} \end{aligned}$$

**The Bernoulli is in the exponential family**

# Definition – univariate exponential family model

Consider a univariate distribution  $f_X(x | \theta)$ , written as:

$$f_X(x | \theta) = \exp \left( h(x) + \boldsymbol{\eta}(\theta)^\top \mathbf{t}(x) - A(\theta) \right)$$

Here we define:

- $h(x)$ : log base measure
- $\boldsymbol{\eta}(\theta)$ : the natural parameters
- $\mathbf{t}(x)$ : the sufficient statistics
- $A(\theta)$ : the log partition function

**Members:**

- Bernoulli (and Multinomial)
- Gaussians

**Example:**

$$\begin{aligned} \bullet \quad h(x) &= -\frac{1}{2} \log(2\pi) \\ \bullet \quad \boldsymbol{\eta}(\mu, \tau) &= \left[ \tau\mu, -\frac{\tau}{2} \right]^\top \\ \bullet \quad \mathbf{t}(x) &= [x, x^2]^\top \\ \bullet \quad A(\theta) &= \frac{\tau\mu^2}{2} - \frac{1}{2} \log |\tau|. \end{aligned} \quad \begin{aligned} f_X(x | \theta) &= \exp \left\{ \left[ \tau\mu, -\frac{\tau}{2} \right] \begin{bmatrix} x \\ x^2 \end{bmatrix} \right. \\ &\quad \left. - \frac{1}{2} \log(2\pi) - (\tau\mu^2/2 - \frac{1}{2} \log \tau) \right\} \\ &= \sqrt{\frac{\tau}{2\pi}} \cdot \exp \left( -\tau(x - \mu)^2/2 \right) \end{aligned}$$

**The Gaussian is in the exponential family**

Consider a univariate distribution  $f_X(x | \theta)$ , written as:

$$f_X(x | \theta) = \exp \left( h(x) + \boldsymbol{\eta}(\theta)^\top \mathbf{t}(x) - A(\theta) \right)$$

Here we define:

- $h(x)$ : log base measure
- $\boldsymbol{\eta}(\theta)$ : the natural parameters
- $\mathbf{t}(x)$ : the sufficient statistics
- $A(\theta)$ : the log partition function

**Members:**

- Bernoulli (and Multinomial)
- Gaussians
- Gamma + Inverse Gamma
- ...

**Example:**

- $h(x) = 0$
- $\boldsymbol{\eta}(\alpha, \beta) = [-\beta, (\alpha - 1)]^\top$
- $\mathbf{t}(x) = [x, \log(x)]^\top$
- $A(\theta) = \log(\Gamma(\alpha)) - \alpha \log(\beta)$

$$\begin{aligned} f_X(x | \theta) &= \exp \left( 0 + [-\beta, (\alpha - 1)]^\top [x, \log(x)]^\top \right. \\ &\quad \left. + \alpha \log(\beta) - \log(\Gamma(\alpha)) \right) \\ &= \frac{\beta^\alpha}{\Gamma(\alpha)} x^{(\alpha-1)} \exp(-\beta x) \end{aligned}$$

**The Gamma is in the exponential family**



## Multivariate distributions

Consider a multi-variate distribution  $f_{\mathbf{x}}(\mathbf{x} | \boldsymbol{\theta})$ , and assume it can be written as:

$$f_{\mathbf{x}}(\mathbf{x} | \boldsymbol{\theta}) = \exp \left( h(\mathbf{x}) + \boldsymbol{\eta}(\boldsymbol{\theta})^\top \mathbf{t}(\mathbf{x}) - A(\boldsymbol{\theta}) \right)$$

We define  $h(\mathbf{x})$  (the log base measure) and  $\mathbf{t}(\mathbf{x})$  (the sufficient statistics) as taking vector-inputs, but otherwise the definition is identical to the univariate case.

### Example:

Consider a  $d$ -dimensional  $\mathbf{x}$  with expectation  $\boldsymbol{\mu}$  and inverse covariance  $\mathbf{Q} = \boldsymbol{\Sigma}^{-1}$ .

Define

- $h(\mathbf{x}) = -\frac{d}{2} \log(2\pi)$
- $\boldsymbol{\eta}(\boldsymbol{\mu}, \mathbf{Q}) = [\mathbf{Q}\boldsymbol{\mu}, -\frac{1}{2}\mathbf{Q}]^\top$
- $\mathbf{t}(\mathbf{x}) = [\mathbf{x}, \mathbf{x}\mathbf{x}^\top]^\top$
- $A(\boldsymbol{\theta}) = \frac{1}{2}\boldsymbol{\mu}^\top \mathbf{Q}\boldsymbol{\mu} - \frac{1}{2} \log |\mathbf{Q}|.$

This gives us the multivariate Gaussian distribution;  $f_{\mathbf{x}}(\mathbf{x}) = \mathcal{N}(\mathbf{x} | \boldsymbol{\mu}, \mathbf{Q}^{-1})$ .

## Multivariate distributions

Consider a multi-variate distribution  $f_{\mathbf{X}}(\mathbf{x} \mid \boldsymbol{\theta})$ , and assume it can be written as:

$$f_{\mathbf{X}}(\mathbf{x} \mid \boldsymbol{\theta}) = \exp \left( h(\mathbf{x}) + \boldsymbol{\eta}(\boldsymbol{\theta})^{\top} \mathbf{t}(\mathbf{x}) - A(\boldsymbol{\theta}) \right)$$

We define  $h(\mathbf{x})$  (the log base measure) and  $\mathbf{t}(\mathbf{x})$  (the sufficient statistics) as taking vector-inputs, but otherwise the definition is identical to the univariate case.

### Notational simplification:

Notice how  $\boldsymbol{\eta}$  take the role of the parameters  $\boldsymbol{\theta}$ , and – for the model we consider – a one-to-one mapping between  $\boldsymbol{\eta}$  and  $\boldsymbol{\theta}$ .

For instance,  $\boldsymbol{\eta} = [\boldsymbol{\eta}_1, \boldsymbol{\eta}_2]^{\top} = [\boldsymbol{\Sigma}^{-1} \boldsymbol{\mu}, -\frac{1}{2} \boldsymbol{\Sigma}^{-1}]$  in the Gaussian distribution, meaning

$$\boldsymbol{\theta} = [\boldsymbol{\mu}, \boldsymbol{\Sigma}]^{\top} = \left[ -\frac{1}{2} \boldsymbol{\eta}_2^{-1} \boldsymbol{\eta}_1, -\frac{1}{2} \boldsymbol{\eta}_2^{-1} \right]^{\top}.$$

Given that  $\boldsymbol{\theta}$  and  $\boldsymbol{\eta}$  are interchangeable, we will simplify notation, and use this form:

$$f_{\mathbf{X}}(\mathbf{x} \mid \boldsymbol{\eta}) = \exp \left( h(\mathbf{x}) + \boldsymbol{\eta}^{\top} \mathbf{t}(\mathbf{x}) - A(\boldsymbol{\eta}) \right)$$

Remember the definition or the Exponential Family:

$$f_X(x | \boldsymbol{\eta}) = \exp(h(\mathbf{x}) + \boldsymbol{\eta}^\top \mathbf{t}(\mathbf{x}) - A(\boldsymbol{\eta})) ,$$

Derivatives of the log normalizer:

$\nabla^k A(\boldsymbol{\eta})$  has an interesting form:

- $\nabla A(\boldsymbol{\eta}) = \frac{dA(\boldsymbol{\eta})}{d\boldsymbol{\eta}} = \mathbb{E}[\mathbf{t}(\mathbf{X})]$

**Proof:**

$$\begin{aligned} \exp(A(\boldsymbol{\eta})) &= \int_{\mathbf{x}} \exp(h(\mathbf{x}) + \boldsymbol{\eta}^\top \mathbf{t}(\mathbf{x})) d\mathbf{x} \\ \frac{dA(\boldsymbol{\eta})}{d\boldsymbol{\eta}} &= \frac{d}{d\boldsymbol{\eta}} \log \int_{\mathbf{x}} \exp(h(\mathbf{x}) + \boldsymbol{\eta}^\top \mathbf{t}(\mathbf{x})) d\mathbf{x} \\ &= \frac{\frac{d}{d\boldsymbol{\eta}} \int_{\mathbf{x}} \exp(h(\mathbf{x}) + \boldsymbol{\eta}^\top \mathbf{t}(\mathbf{x})) d\mathbf{x}}{\int_{\mathbf{x}} \exp(h(\mathbf{x}) + \boldsymbol{\eta}^\top \mathbf{t}(\mathbf{x})) d\mathbf{x}} \\ &= \frac{\int_{\mathbf{x}} \mathbf{t}(\mathbf{x}) \exp(h(\mathbf{x}) + \boldsymbol{\eta}^\top \mathbf{t}(\mathbf{x})) d\mathbf{x}}{\exp(A(\boldsymbol{\eta}))} \\ &= \int_{\mathbf{x}} \mathbf{t}(\mathbf{x}) \exp(h(\mathbf{x}) + \boldsymbol{\eta}^\top \mathbf{t}(\mathbf{x}) - A(\boldsymbol{\eta})) d\mathbf{x} = \mathbb{E}[\mathbf{t}(\mathbf{X})] \end{aligned}$$

Remember the definition or the Exponential Family:

$$f_X(x | \boldsymbol{\eta}) = \exp(h(\mathbf{x}) + \boldsymbol{\eta}^\top \mathbf{t}(\mathbf{x}) - A(\boldsymbol{\eta})) ,$$

Derivatives of the log normalizer:

$\nabla^k A(\boldsymbol{\eta})$  has an interesting form:

- $\nabla A(\boldsymbol{\eta}) = \frac{dA(\boldsymbol{\eta})}{d\boldsymbol{\eta}} = \mathbb{E}[\mathbf{t}(\mathbf{X})]$
- $\nabla^2 A(\boldsymbol{\eta}) = \mathbb{E}[(\mathbf{t}(\mathbf{X}) - \mathbb{E}[\mathbf{t}(\mathbf{X})])^2] = \mathbb{V}[\mathbf{t}(\mathbf{X})]$ 
  - ... which also shows that  $A(\boldsymbol{\eta})$  is *convex*.
- $\nabla^k A(\boldsymbol{\eta}) = \mathbb{E}[(\mathbf{t}(\mathbf{X}) - \mathbb{E}[\mathbf{t}(\mathbf{X})])^k]$

Proofs for  $k > 1$  are more of the same manipulations (left out for simplicity).

## Code Task: Translation between moment-based and ExpFam representations

In this task you will translate between moment-based representations and the corresponding ExpFam representations. The task is fairly straight-forward, but is intended to give you “hands-on” experience with ExpFam representation.

Start from

```
students_translator.ipynb
```

You will see that there are some supporting functions there, and a translation for univariate Gaussians (to show how the supporting functions work).

Your task is to implement the same for the Gamma distribution. If time, other ExpFam distributions can of course be given the same treatment.