

# Variational Inference Masterclass: Variational Message Passing

Tui Nolan

8 November 2022



UNIVERSITY OF  
CAMBRIDGE



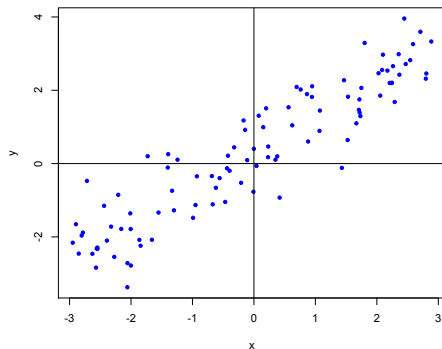
MRC  
Biostatistics  
Unit

## Part I

### Introduction and Motivation

# Linear Regression

Consider the following data:

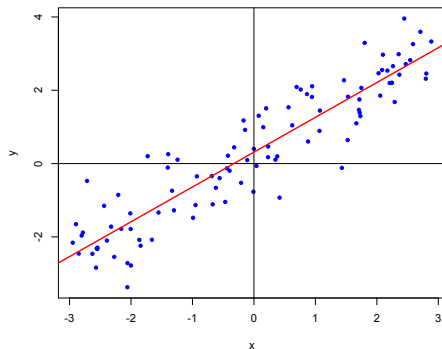


In a typical regression problem we solve using:

$$\hat{\beta} = \underset{\beta \in \mathbb{R}^2}{\operatorname{argmin}} \|\mathbf{y} - \mathbf{X}\beta\|_2$$

# Linear Regression

Consider the following data:



In a typical regression problem we solve using:

$$\hat{\beta} = \underset{\beta \in \mathbb{R}^2}{\operatorname{argmin}} \|\mathbf{y} - \mathbf{X}\beta\|_2$$

# Mean Field Variational Bayes

Tractable Bayesian inference is achieved by introducing a latent variable (Gelman, 2006; Huang and Wand, 2013):

$$\begin{aligned}y_i|x_i, \beta, \sigma^2 &\stackrel{\text{ind.}}{\sim} N(\beta_0 + \beta_1 x_i, \sigma^2) \quad i = 1, \dots, n \\ \beta &\sim N(\mathbf{0}, \Sigma_0) \\ \sigma^2|a &\sim \text{Inverse} - \chi^2(1, 1/a) \\ a &\sim \text{Inverse} - \chi^2(1, 1/A^2)\end{aligned}$$

Note that if

$$\sigma^2|a \sim \text{Inverse} - \chi^2(1, 1/a) \quad \text{and} \quad a \sim \text{Inverse} - \chi^2(1, 1/A^2)$$

We use the mean field assumption (Menictas and Wand, 2013):

$$\{\beta, a\} \perp\!\!\!\perp \sigma^2$$

That is,

$$p(\beta, \sigma^2, a|\mathbf{y}) \approx p(\beta, a|\mathbf{y})p(\sigma^2|\mathbf{y})$$

# Mean Field Variational Bayes

Tractable Bayesian inference is achieved by introducing a latent variable (Gelman, 2006; Huang and Wand, 2013):

$$\begin{aligned}y_i|x_i, \beta, \sigma^2 &\stackrel{\text{ind.}}{\sim} N(\beta_0 + \beta_1 x_i, \sigma^2) \quad i = 1, \dots, n \\ \beta &\sim N(\mathbf{0}, \Sigma_0) \\ \sigma^2|a &\sim \text{Inverse} - \chi^2(1, 1/a) \\ a &\sim \text{Inverse} - \chi^2(1, 1/A^2)\end{aligned}$$

Note that if

$$\sigma^2|a \sim \text{Inverse} - \chi^2(1, 1/a) \quad \text{and} \quad a \sim \text{Inverse} - \chi^2(1, 1/A^2)$$

We use the mean field assumption (Menictas and Wand, 2013):

$$\{\beta, a\} \perp\!\!\!\perp \sigma^2$$

That is,

$$p(\beta, \sigma^2, a|\mathbf{y}) \approx p(\beta, a|\mathbf{y})p(\sigma^2|\mathbf{y})$$

# Mean Field Variational Bayes

Tractable Bayesian inference is achieved by introducing a latent variable (Gelman, 2006; Huang and Wand, 2013):

$$\begin{aligned}y_i|x_i, \beta, \sigma^2 &\stackrel{\text{ind.}}{\sim} N(\beta_0 + \beta_1 x_i, \sigma^2) \quad i = 1, \dots, n \\ \beta &\sim N(\mathbf{0}, \Sigma_0) \\ \sigma^2|a &\sim \text{Inverse} - \chi^2(1, 1/a) \\ a &\sim \text{Inverse} - \chi^2(1, 1/A^2)\end{aligned}$$

Note that if

$$\sigma^2|a \sim \text{Inverse} - \chi^2(1, 1/a) \quad \text{and} \quad a \sim \text{Inverse} - \chi^2(1, 1/A^2)$$

We use the mean field assumption (Menictas and Wand, 2013):

$$\{\beta, a\} \perp\!\!\!\perp \sigma^2$$

That is,

$$p(\beta, \sigma^2, a|\mathbf{y}) \approx p(\beta, a|\mathbf{y})p(\sigma^2|\mathbf{y})$$

In variational inference, we approximate the posterior density function  $p(\beta, \sigma^2, a | \mathbf{y})$  by another density function  $q(\beta, \sigma^2, a)$ .

The “best” approximate density function  $q^*(\beta, \sigma^2, a)$  is the one that is “closest” to the true posterior density function  $p(\beta, \sigma^2, a | \mathbf{y})$ :

$$q^*(\beta, \sigma^2, a) = \underset{q \in \mathcal{Q}}{\operatorname{argmin}} D_{\text{KL}} \left\{ q(\beta, \sigma^2, a) \| p(\beta, \sigma^2, a | \mathbf{y}) \right\}$$

Our mean field assumption combined with results from Graph Theory, see Chapter 8 of Bishop (2006), allow us to enforce the product form:

$$p(\beta, \sigma^2, a | \mathbf{y}) \approx q(\beta, \sigma^2, a) = q(\beta, a)q(\sigma^2) = q(\beta)q(\sigma^2)q(a)$$

For  $\beta$ , the optimisation problem has the following solution:

$$\begin{aligned} q^*(\beta) &= C_1 \exp \left[ \mathbb{E}_{-q(\beta)} \log \{ p(\mathbf{y}, \beta, \sigma^2, a) \} \right] \\ &= C_2 \exp \left[ \mathbb{E}_{-q(\beta)} \log \{ p(\mathbf{y} | \beta, \sigma^2, a) p(\beta) \} \right] \end{aligned}$$

where

$$\mathbb{E}_{-q(\beta)} f(\mathbf{y}, \beta, \sigma^2, a) \equiv \int \int f(\mathbf{y}, \beta, \sigma^2, a) q(\sigma^2) q(a) da d\sigma^2$$



# Mean Field Variational Bayes

In variational inference, we approximate the posterior density function  $p(\beta, \sigma^2, a | \mathbf{y})$  by another density function  $q(\beta, \sigma^2, a)$ .

The “best” approximate density function  $q^*(\beta, \sigma^2, a)$  is the one that is “closest” to the true posterior density function  $p(\beta, \sigma^2, a | \mathbf{y})$ :

$$q^*(\beta, \sigma^2, a) = \underset{q \in \mathcal{Q}}{\operatorname{argmin}} D_{\text{KL}} \{ q(\beta, \sigma^2, a) \| p(\beta, \sigma^2, a | \mathbf{y}) \}$$

Our mean field assumption combined with results from Graph Theory, see Chapter 8 of Bishop (2006), allow us to enforce the product form:

$$p(\beta, \sigma^2, a | \mathbf{y}) \approx q(\beta, \sigma^2, a) = q(\beta, a)q(\sigma^2) = q(\beta)q(\sigma^2)q(a)$$

For  $\beta$ , the optimisation problem has the following solution:

$$\begin{aligned} q^*(\beta) &= C_1 \exp \left[ \mathbb{E}_{-q(\beta)} \log \{ p(\mathbf{y}, \beta, \sigma^2, a) \} \right] \\ &= C_2 \exp \left[ \mathbb{E}_{-q(\beta)} \log \{ p(\mathbf{y} | \beta, \sigma^2, a) p(\beta) \} \right] \end{aligned}$$

where

$$\mathbb{E}_{-q(\beta)} f(\mathbf{y}, \beta, \sigma^2, a) \equiv \int \int f(\mathbf{y}, \beta, \sigma^2, a) q(\sigma^2) q(a) da d\sigma^2$$

In variational inference, we approximate the posterior density function  $p(\beta, \sigma^2, a | \mathbf{y})$  by another density function  $q(\beta, \sigma^2, a)$ .

The “best” approximate density function  $q^*(\beta, \sigma^2, a)$  is the one that is “closest” to the true posterior density function  $p(\beta, \sigma^2, a | \mathbf{y})$ :

$$q^*(\beta, \sigma^2, a) = \underset{q \in \mathcal{Q}}{\operatorname{argmin}} D_{\text{KL}} \{ q(\beta, \sigma^2, a) \| p(\beta, \sigma^2, a | \mathbf{y}) \}$$

Our mean field assumption combined with results from Graph Theory, see Chapter 8 of Bishop (2006), allow us to enforce the product form:

$$p(\beta, \sigma^2, a | \mathbf{y}) \approx q(\beta, \sigma^2, a) = q(\beta, a)q(\sigma^2) = q(\beta)q(\sigma^2)q(a)$$

For  $\beta$ , the optimisation problem has the following solution:

$$\begin{aligned} q^*(\beta) &= C_1 \exp \left[ \mathbb{E}_{-q(\beta)} \log \{ p(\mathbf{y}, \beta, \sigma^2, a) \} \right] \\ &= C_2 \exp \left[ \mathbb{E}_{-q(\beta)} \log \{ p(\mathbf{y} | \beta, \sigma^2, a) p(\beta) \} \right] \end{aligned}$$

where

$$\mathbb{E}_{-q(\beta)} f(\mathbf{y}, \beta, \sigma^2, a) \equiv \int \int f(\mathbf{y}, \beta, \sigma^2, a) q(\sigma^2) q(a) da d\sigma^2$$

# Mean Field Variational Bayes

In variational inference, we approximate the posterior density function  $p(\beta, \sigma^2, a | \mathbf{y})$  by another density function  $q(\beta, \sigma^2, a)$ .

The “best” approximate density function  $q^*(\beta, \sigma^2, a)$  is the one that is “closest” to the true posterior density function  $p(\beta, \sigma^2, a | \mathbf{y})$ :

$$q^*(\beta, \sigma^2, a) = \underset{q \in \mathcal{Q}}{\operatorname{argmin}} D_{\text{KL}} \{q(\beta, \sigma^2, a) \| p(\beta, \sigma^2, a | \mathbf{y})\}$$

Our mean field assumption combined with results from Graph Theory, see Chapter 8 of Bishop (2006), allow us to enforce the product form:

$$p(\beta, \sigma^2, a | \mathbf{y}) \approx q(\beta, \sigma^2, a) = q(\beta, a)q(\sigma^2) = q(\beta)q(\sigma^2)q(a)$$

For  $\beta$ , the optimisation problem has the following solution:

$$\begin{aligned} q^*(\beta) &= C_1 \exp \left[ \mathbb{E}_{-q(\beta)} \log \{p(\mathbf{y}, \beta, \sigma^2, a)\} \right] \\ &= C_2 \exp \left[ \mathbb{E}_{-q(\beta)} \log \{p(\mathbf{y} | \beta, \sigma^2, a)p(\beta)\} \right] \end{aligned}$$

where

$$\mathbb{E}_{-q(\beta)} f(\mathbf{y}, \beta, \sigma^2, a) \equiv \int \int f(\mathbf{y}, \beta, \sigma^2, a) q(\sigma^2) q(a) da d\sigma^2$$

$$\begin{aligned}\mathbb{E}_{-q(\beta)} \log \{p(\mathbf{y}|\beta, \sigma^2)p(\beta)\} \\&= \mathbb{E}_{-q(\beta)} \log p(\mathbf{y}|\beta, \sigma^2) + \log p(\beta) \\&= \mathbb{E}_{-q(\beta)} \left\{ -\frac{1}{2\sigma^2} (\mathbf{y} - \mathbf{X}\beta)^T (\mathbf{y} - \mathbf{X}\beta) \right\} - \frac{1}{2} \beta^T \Sigma_0^{-1} \beta + \text{const.} \\&= -\frac{1}{2} \mathbb{E}_q \left( \frac{1}{\sigma^2} \right) (\mathbf{y} - \mathbf{X}\beta)^T (\mathbf{y} - \mathbf{X}\beta) - \frac{1}{2} \beta^T \Sigma_0^{-1} \beta + \text{const.} \\&= -\frac{1}{2} \beta^T \left\{ \mathbb{E}_q \left( \frac{1}{\sigma^2} \right) \mathbf{X}^T \mathbf{X} + \Sigma_0^{-1} \right\} \beta + \mathbb{E}_q \left( \frac{1}{\sigma^2} \right) \beta^T \mathbf{X}^T \mathbf{y} + \text{const.}\end{aligned}$$

Remember that  $q^*(\beta) \propto \exp \left[ \mathbb{E}_{-q(\beta)} \log \{p(\mathbf{y}|\beta, \sigma^2)p(\beta)\} \right]$

That is,  $q^*(\beta)$  is a normal density function.

Similar arguments can be used to show that  $q^*(\sigma^2)$  and  $q^*(a)$  are inverse chi-squared density functions.

$$\begin{aligned}\mathbb{E}_{-q(\beta)} \log \{p(\mathbf{y}|\beta, \sigma^2)p(\beta)\} \\&= \mathbb{E}_{-q(\beta)} \log p(\mathbf{y}|\beta, \sigma^2) + \log p(\beta) \\&= \mathbb{E}_{-q(\beta)} \left\{ -\frac{1}{2\sigma^2} (\mathbf{y} - \mathbf{X}\beta)^T (\mathbf{y} - \mathbf{X}\beta) \right\} - \frac{1}{2} \beta^T \Sigma_0^{-1} \beta + \text{const.} \\&= -\frac{1}{2} \mathbb{E}_q \left( \frac{1}{\sigma^2} \right) (\mathbf{y} - \mathbf{X}\beta)^T (\mathbf{y} - \mathbf{X}\beta) - \frac{1}{2} \beta^T \Sigma_0^{-1} \beta + \text{const.} \\&= -\frac{1}{2} \beta^T \left\{ \mathbb{E}_q \left( \frac{1}{\sigma^2} \right) \mathbf{X}^T \mathbf{X} + \Sigma_0^{-1} \right\} \beta + \mathbb{E}_q \left( \frac{1}{\sigma^2} \right) \beta^T \mathbf{X}^T \mathbf{y} + \text{const.}\end{aligned}$$

Remember that  $q^*(\beta) \propto \exp \left[ \mathbb{E}_{-q(\beta)} \log \{p(\mathbf{y}|\beta, \sigma^2)p(\beta)\} \right]$

That is,  $q^*(\beta)$  is a normal density function.

Similar arguments can be used to show that  $q^*(\sigma^2)$  and  $q^*(a)$  are inverse chi-squared density functions.

$$\begin{aligned}\mathbb{E}_{-q(\beta)} \log \{p(\mathbf{y}|\beta, \sigma^2)p(\beta)\} \\&= \mathbb{E}_{-q(\beta)} \log p(\mathbf{y}|\beta, \sigma^2) + \log p(\beta) \\&= \mathbb{E}_{-q(\beta)} \left\{ -\frac{1}{2\sigma^2} (\mathbf{y} - \mathbf{X}\beta)^T (\mathbf{y} - \mathbf{X}\beta) \right\} - \frac{1}{2} \beta^T \Sigma_0^{-1} \beta + \text{const.} \\&= -\frac{1}{2} \mathbb{E}_q \left( \frac{1}{\sigma^2} \right) (\mathbf{y} - \mathbf{X}\beta)^T (\mathbf{y} - \mathbf{X}\beta) - \frac{1}{2} \beta^T \Sigma_0^{-1} \beta + \text{const.} \\&= -\frac{1}{2} \beta^T \left\{ \mathbb{E}_q \left( \frac{1}{\sigma^2} \right) \mathbf{X}^T \mathbf{X} + \Sigma_0^{-1} \right\} \beta + \mathbb{E}_q \left( \frac{1}{\sigma^2} \right) \beta^T \mathbf{X}^T \mathbf{y} + \text{const.}\end{aligned}$$

Remember that  $q^*(\beta) \propto \exp \left[ \mathbb{E}_{-q(\beta)} \log \{p(\mathbf{y}|\beta, \sigma^2)p(\beta)\} \right]$

That is,  $q^*(\beta)$  is a normal density function.

Similar arguments can be used to show that  $q^*(\sigma^2)$  and  $q^*(a)$  are inverse chi-squared density functions.

$$\begin{aligned}\mathbb{E}_{-q(\beta)} \log \{p(\mathbf{y}|\beta, \sigma^2)p(\beta)\} \\&= \mathbb{E}_{-q(\beta)} \log p(\mathbf{y}|\beta, \sigma^2) + \log p(\beta) \\&= \mathbb{E}_{-q(\beta)} \left\{ -\frac{1}{2\sigma^2} (\mathbf{y} - \mathbf{X}\beta)^T (\mathbf{y} - \mathbf{X}\beta) \right\} - \frac{1}{2} \beta^T \Sigma_0^{-1} \beta + \text{const.} \\&= -\frac{1}{2} \mathbb{E}_q \left( \frac{1}{\sigma^2} \right) (\mathbf{y} - \mathbf{X}\beta)^T (\mathbf{y} - \mathbf{X}\beta) - \frac{1}{2} \beta^T \Sigma_0^{-1} \beta + \text{const.} \\&= -\frac{1}{2} \beta^T \left\{ \mathbb{E}_q \left( \frac{1}{\sigma^2} \right) \mathbf{X}^T \mathbf{X} + \Sigma_0^{-1} \right\} \beta + \mathbb{E}_q \left( \frac{1}{\sigma^2} \right) \beta^T \mathbf{X}^T \mathbf{y} + \text{const.}\end{aligned}$$

Remember that  $q^*(\beta) \propto \exp \left[ \mathbb{E}_{-q(\beta)} \log \{p(\mathbf{y}|\beta, \sigma^2)p(\beta)\} \right]$

That is,  $q^*(\beta)$  is a normal density function.

Similar arguments can be used to show that  $q^*(\sigma^2)$  and  $q^*(a)$  are inverse chi-squared density functions.

$q^*(\beta)$  is the  $N(\mu_{q(\beta)}, \Sigma_{q(\beta)})$  density function, where

$$\Sigma_{q(\beta)} \leftarrow \left\{ \left( \frac{N+1}{\lambda_{q(\sigma^2)}} \right) \mathbf{X}^T \mathbf{X} + \Sigma_0^{-1} \right\}^{-1}$$

and  $\mu_{q(\beta)} \leftarrow \left( \frac{N+1}{\lambda_{q(\sigma^2)}} \right) \Sigma_{q(\beta)} \mathbf{X}^T \mathbf{y}$

$q^*(\sigma^2)$  is the Inverse  $-\chi^2(N+1, \lambda_{q(\sigma^2)})$  density function, where

$$\lambda_{q(\sigma^2)} \leftarrow \mathbf{y}^T \mathbf{y} - 2\mu_{q(\beta)}^T \mathbf{X}^T \mathbf{y} + \text{tr} \left\{ \mathbf{X}^T \mathbf{X} (\Sigma_{q(\beta)} + \mu_{q(\beta)} \mu_{q(\beta)}^T) \right\} + \frac{2}{\lambda_{q(a)}}$$

$q^*(a)$  is the Inverse  $-\chi^2(2, \lambda_{q(a)})$  density function, where

$$\lambda_{q(a)} \leftarrow 2 \left( \frac{N+1}{\lambda_{q(\sigma^2)}} \right) + \frac{1}{A^2}$$



$q^*(\beta)$  is the  $N(\mu_{q(\beta)}, \Sigma_{q(\beta)})$  density function, where

$$\Sigma_{q(\beta)} \leftarrow \left\{ \left( \frac{N+1}{\lambda_{q(\sigma^2)}} \right) \mathbf{X}^T \mathbf{X} + \Sigma_0^{-1} \right\}^{-1}$$

and  $\mu_{q(\beta)} \leftarrow \left( \frac{N+1}{\lambda_{q(\sigma^2)}} \right) \Sigma_{q(\beta)} \mathbf{X}^T \mathbf{y}$

$q^*(\sigma^2)$  is the Inverse  $-\chi^2(N+1, \lambda_{q(\sigma^2)})$  density function, where

$$\lambda_{q(\sigma^2)} \leftarrow \mathbf{y}^T \mathbf{y} - 2\mu_{q(\beta)}^T \mathbf{X}^T \mathbf{y} + \text{tr} \left\{ \mathbf{X}^T \mathbf{X} (\Sigma_{q(\beta)} + \mu_{q(\beta)} \mu_{q(\beta)}^T) \right\} + \frac{2}{\lambda_{q(a)}}$$

$q^*(a)$  is the Inverse  $-\chi^2(2, \lambda_{q(a)})$  density function, where

$$\lambda_{q(a)} \leftarrow 2 \left( \frac{N+1}{\lambda_{q(\sigma^2)}} \right) + \frac{1}{A^2}$$

$q^*(\beta)$  is the  $N(\mu_{q(\beta)}, \Sigma_{q(\beta)})$  density function, where

$$\Sigma_{q(\beta)} \leftarrow \left\{ \left( \frac{N+1}{\lambda_{q(\sigma^2)}} \right) \mathbf{X}^T \mathbf{X} + \Sigma_0^{-1} \right\}^{-1}$$
$$\text{and } \mu_{q(\beta)} \leftarrow \left( \frac{N+1}{\lambda_{q(\sigma^2)}} \right) \Sigma_{q(\beta)} \mathbf{X}^T \mathbf{y}$$

$q^*(\sigma^2)$  is the Inverse- $\chi^2(N+1, \lambda_{q(\sigma^2)})$  density function, where

$$\lambda_{q(\sigma^2)} \leftarrow \mathbf{y}^T \mathbf{y} - 2\mu_{q(\beta)} \mathbf{X}^T \mathbf{y} + \text{tr} \left\{ \mathbf{X}^T \mathbf{X} (\Sigma_{q(\beta)} + \mu_{q(\beta)} \mu_{q(\beta)}^T) \right\} + \frac{2}{\lambda_{q(a)}}$$

$q^*(a)$  is the Inverse- $\chi^2(2, \lambda_{q(a)})$  density function, where

$$\lambda_{q(a)} \leftarrow 2 \left( \frac{N+1}{\lambda_{q(\sigma^2)}} \right) + \frac{1}{A^2}$$

## MFVB for the Bayesian linear regression model

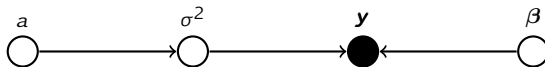
1. Initialise all optimal posterior density functions
2. Cycle:

$$q^*(a) \propto \exp\left\{\mathbb{E}_{-q^*(a)} \log p(\mathbf{y}, \beta, \sigma^2, a)\right\}$$

$$q^*(\sigma^2) \propto \exp\left\{\mathbb{E}_{-q^*(\sigma^2)} \log p(\mathbf{y}, \beta, \sigma^2, a)\right\}$$

$$q^*(\beta) \propto \exp\left\{\mathbb{E}_{-q^*(\beta)} \log p(\mathbf{y}, \beta, \sigma^2, a)\right\}$$

3. Stop:  $D_{\text{KL}}\{q(\beta, \sigma^2, a) \| p(\beta, \sigma^2, a | \mathbf{y})\}$  converges.



## MFVB for the Bayesian linear regression model

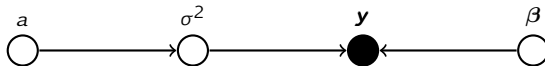
1. Initialise all optimal posterior density functions
2. Cycle:

$$q^*(a) \propto \exp\left\{\mathbb{E}_{-q^*(a)} \log p(\mathbf{y}, \beta, \sigma^2, a)\right\}$$

$$q^*(\sigma^2) \propto \exp\left\{\mathbb{E}_{-q^*(\sigma^2)} \log p(\mathbf{y}, \beta, \sigma^2, a)\right\}$$

$$q^*(\beta) \propto \exp\left\{\mathbb{E}_{-q^*(\beta)} \log p(\mathbf{y}, \beta, \sigma^2, a)\right\}$$

3. Stop:  $D_{\text{KL}}\{q(\beta, \sigma^2, a) \| p(\beta, \sigma^2, a | \mathbf{y})\}$  converges.



## MFVB for the Bayesian linear regression model

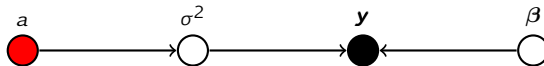
1. Initialise all optimal posterior density functions
2. Cycle:

$$q^*(a) \propto \exp\left\{\mathbb{E}_{-q^*(a)} \log p(\mathbf{y}, \beta, \sigma^2, a)\right\}$$

$$q^*(\sigma^2) \propto \exp\left\{\mathbb{E}_{-q^*(\sigma^2)} \log p(\mathbf{y}, \beta, \sigma^2, a)\right\}$$

$$q^*(\beta) \propto \exp\left\{\mathbb{E}_{-q^*(\beta)} \log p(\mathbf{y}, \beta, \sigma^2, a)\right\}$$

3. Stop:  $D_{\text{KL}}\{q(\beta, \sigma^2, a) \| p(\beta, \sigma^2, a | \mathbf{y})\}$  converges.



## MFVB for the Bayesian linear regression model

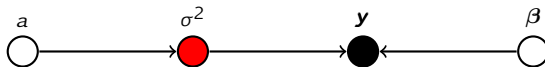
1. Initialise all optimal posterior density functions
2. Cycle:

$$q^*(a) \propto \exp\left\{\mathbb{E}_{-q^*(a)} \log p(\mathbf{y}, \beta, \sigma^2, a)\right\}$$

$$q^*(\sigma^2) \propto \exp\left\{\mathbb{E}_{-q^*(\sigma^2)} \log p(\mathbf{y}, \beta, \sigma^2, a)\right\}$$

$$q^*(\beta) \propto \exp\left\{\mathbb{E}_{-q^*(\beta)} \log p(\mathbf{y}, \beta, \sigma^2, a)\right\}$$

3. Stop:  $D_{\text{KL}}\{q(\beta, \sigma^2, a) \| p(\beta, \sigma^2, a | \mathbf{y})\}$  converges.



## MFVB for the Bayesian linear regression model

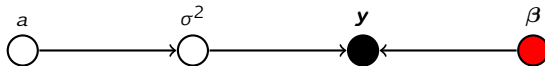
1. Initialise all optimal posterior density functions
2. Cycle:

$$q^*(a) \propto \exp\left\{\mathbb{E}_{-q^*(a)} \log p(\mathbf{y}, \beta, \sigma^2, a)\right\}$$

$$q^*(\sigma^2) \propto \exp\left\{\mathbb{E}_{-q^*(\sigma^2)} \log p(\mathbf{y}, \beta, \sigma^2, a)\right\}$$

$$q^*(\beta) \propto \exp\left\{\mathbb{E}_{-q^*(\beta)} \log p(\mathbf{y}, \beta, \sigma^2, a)\right\}$$

3. Stop:  $D_{\text{KL}}\{q(\beta, \sigma^2, a) \| p(\beta, \sigma^2, a | \mathbf{y})\}$  converges.



## MFVB for the Bayesian linear regression model

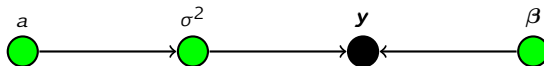
1. Initialise all optimal posterior density functions
2. Cycle:

$$q^*(a) \propto \exp\left\{\mathbb{E}_{-q^*(a)} \log p(\mathbf{y}, \boldsymbol{\beta}, \sigma^2, a)\right\}$$

$$q^*(\sigma^2) \propto \exp\left\{\mathbb{E}_{-q^*(\sigma^2)} \log p(\mathbf{y}, \boldsymbol{\beta}, \sigma^2, a)\right\}$$

$$q^*(\boldsymbol{\beta}) \propto \exp\left\{\mathbb{E}_{-q^*(\boldsymbol{\beta})} \log p(\mathbf{y}, \boldsymbol{\beta}, \sigma^2, a)\right\}$$

3. Stop:  $D_{\text{KL}}\{q(\boldsymbol{\beta}, \sigma^2, a) \| p(\boldsymbol{\beta}, \sigma^2, a | \mathbf{y})\}$  converges.





The advantage of MFVB over Monte Carlo alternatives is that it is a much faster algorithm

However, for Bayesian inference on more complex models, we have to re-do all the derivations.

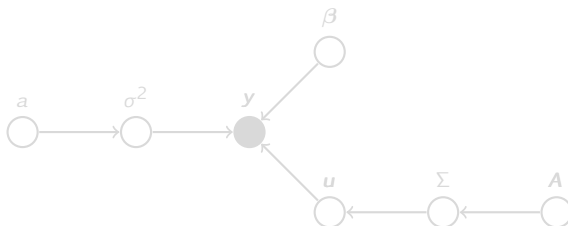
Take for example, the Bayesian Gaussian-response linear mixed model:

$$y_i | \beta, u_i, \sigma^2 \sim N(X_i \beta + X_i u_i, \sigma^2 I), \quad u_i | \Sigma \sim N(0, \Sigma), \quad \text{for } 1 \leq i \leq m$$

$$\beta \sim N(0, \sigma_\beta^2 I); \quad \sigma^2 | a \sim \text{Inverse} - \chi^2(1, 1/a)$$

$$a \sim \text{Inverse} - \chi^2(1, 1/2), \quad \Sigma | \mathbf{A} \sim \text{Inverse G-Wishart}(C_{\text{full}}, 2q, \mathbf{A}^{-1})$$

$$\mathbf{A} \sim \text{Inverse G-Wishart}(C_{\text{diag}}, 1, \frac{1}{2}I)$$



# Limitations of MFVB

The advantage of MFVB over Monte Carlo alternatives is that it is a much faster algorithm

However, for Bayesian inference on more complex models, we have to re-do all the derivations.

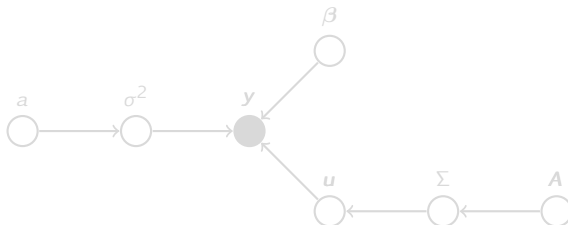
Take for example, the Bayesian Gaussian-response linear mixed model:

$$y_i | \beta, u_i, \sigma^2 \sim N(X_i \beta + X_i u_i, \sigma^2 I), \quad u_i | \Sigma \sim N(0, \Sigma), \quad \text{for } 1 \leq i \leq m$$

$$\beta \sim N(0, \sigma_\beta^2 I); \quad \sigma^2 | a \sim \text{Inverse-}\chi^2(1, 1/a)$$

$$a \sim \text{Inverse-}\chi^2(1, 1/2), \quad \Sigma | \mathbf{A} \sim \text{Inverse G-Wishart}(C_{\text{full}}, 2q, \mathbf{A}^{-1})$$

$$\mathbf{A} \sim \text{Inverse G-Wishart}(C_{\text{diag}}, 1, \frac{1}{2}I)$$



The advantage of MFVB over Monte Carlo alternatives is that it is a much faster algorithm

However, for Bayesian inference on more complex models, we have to re-do all the derivations.

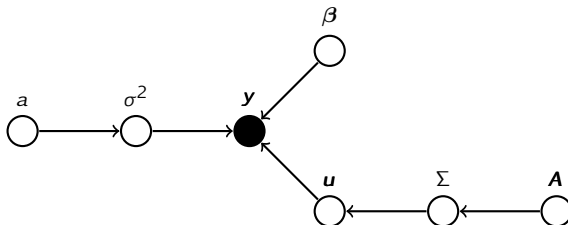
Take for example, the Bayesian Gaussian-response linear mixed model:

$$\mathbf{y}_i | \boldsymbol{\beta}, \mathbf{u}_i, \sigma^2 \sim N(\mathbf{X}_i \boldsymbol{\beta} + \mathbf{X}_i \mathbf{u}_i, \sigma^2 \mathbf{I}), \quad \mathbf{u}_i | \boldsymbol{\Sigma} \sim N(\mathbf{0}, \boldsymbol{\Sigma}), \quad \text{for } 1 \leq i \leq m$$

$$\boldsymbol{\beta} \sim N(\mathbf{0}, \sigma_{\beta}^2 \mathbf{I}); \quad \sigma^2 | a \sim \text{Inverse-}\chi^2(1, 1/a)$$

$$a \sim \text{Inverse-}\chi^2(1, 1/2), \quad \boldsymbol{\Sigma} | \mathbf{A} \sim \text{Inverse G-Wishart}(C_{\text{full}}, 2q, \mathbf{A}^{-1})$$

$$\mathbf{A} \sim \text{Inverse G-Wishart}(C_{\text{diag}}, 1, \frac{1}{2} \mathbf{I})$$



# Variational Message Passing

Prof. Matt Wand

Distinguished professor of statistics, University of Technology Sydney

*Fast Approximate Inference for Arbitrarily Large Semiparametric Regression Models via Message Passing, Journal of the American Statistical Association, 2017, VOL. 112, NO. 517, 137–168, Theory and Methods*



# Steps Towards Message Passing Algorithms

Let's reconsider the Bayesian linear regression model:

$$y_i|x_i, \beta, \sigma^2 \stackrel{\text{ind.}}{\sim} N(\beta_0 + \beta_1 x_i, \sigma^2) \quad i = 1, \dots, n, \quad \beta \sim N(\mathbf{0}, \Sigma_0) \\ \sigma^2|a \sim \text{Inverse} - \chi^2(1, 1/a) \quad a \sim \text{Inverse} - \chi^2(1, 1/A^2)$$

We showed that the optimal posterior density function for  $\beta$  is

$$q^*(\beta) = C \exp \left[ -\frac{1}{2} \beta^T \left\{ \mathbb{E}_q(1/\sigma^2) \mathbf{X}^T \mathbf{X} + \Sigma_0^{-1} \right\} \beta + \mathbb{E}_q(1/\sigma^2) \beta^T \mathbf{X}^T \mathbf{y} \right] \\ = C \exp \left\{ -\frac{1}{2} \mathbb{E}_q(1/\sigma^2) \beta^T \mathbf{X}^T \mathbf{X} \beta + \mathbb{E}_q(1/\sigma^2) \beta^T \mathbf{X}^T \mathbf{y} \right\} \exp \left\{ -\frac{1}{2} \beta^T \Sigma_0^{-1} \beta \right\}$$

That is, the optimal posterior density function can be factorised into a product of a Gaussian contribution from the likelihood and a Gaussian contribution from the prior distribution.

# Steps Towards Message Passing Algorithms

Let's reconsider the Bayesian linear regression model:

$$y_i | x_i, \beta, \sigma^2 \stackrel{\text{ind.}}{\sim} N(\beta_0 + \beta_1 x_i, \sigma^2) \quad i = 1, \dots, n, \quad \beta \sim N(\mathbf{0}, \Sigma_0) \\ \sigma^2 | a \sim \text{Inverse} - \chi^2(1, 1/a) \quad a \sim \text{Inverse} - \chi^2(1, 1/A^2)$$

We showed that the optimal posterior density function for  $\beta$  is

$$q^*(\beta) = C \exp \left[ -\frac{1}{2} \beta^T \left\{ \mathbb{E}_q(1/\sigma^2) \mathbf{X}^T \mathbf{X} + \Sigma_0^{-1} \right\} \beta + \mathbb{E}_q(1/\sigma^2) \beta^T \mathbf{X}^T \mathbf{y} \right] \\ = C \exp \left\{ -\frac{1}{2} \mathbb{E}_q(1/\sigma^2) \beta^T \mathbf{X}^T \mathbf{X} \beta + \mathbb{E}_q(1/\sigma^2) \beta^T \mathbf{X}^T \mathbf{y} \right\} \exp \left\{ -\frac{1}{2} \beta^T \Sigma_0^{-1} \beta \right\}$$

That is, the optimal posterior density function can be factorised into a product of a Gaussian contribution from the likelihood and a Gaussian contribution from the prior distribution.

# Steps Towards Message Passing Algorithms

Let's reconsider the Bayesian linear regression model:

$$y_i | x_i, \beta, \sigma^2 \stackrel{\text{ind.}}{\sim} N(\beta_0 + \beta_1 x_i, \sigma^2) \quad i = 1, \dots, n, \quad \beta \sim N(\mathbf{0}, \Sigma_0) \\ \sigma^2 | a \sim \text{Inverse} - \chi^2(1, 1/a) \quad a \sim \text{Inverse} - \chi^2(1, 1/A^2)$$

We showed that the optimal posterior density function for  $\beta$  is

$$q^*(\beta) = C \exp \left[ -\frac{1}{2} \beta^T \left\{ \mathbb{E}_q(1/\sigma^2) \mathbf{X}^T \mathbf{X} + \Sigma_0^{-1} \right\} \beta + \mathbb{E}_q(1/\sigma^2) \beta^T \mathbf{X}^T \mathbf{y} \right] \\ = C \exp \left\{ -\frac{1}{2} \mathbb{E}_q(1/\sigma^2) \beta^T \mathbf{X}^T \mathbf{X} \beta + \mathbb{E}_q(1/\sigma^2) \beta^T \mathbf{X}^T \mathbf{y} \right\} \exp \left\{ -\frac{1}{2} \beta^T \Sigma_0^{-1} \beta \right\}$$

That is, the optimal posterior density function can be factorised into a product of a Gaussian contribution from the likelihood and a Gaussian contribution from the prior distribution.

# Steps Towards Message Passing Algorithms

Let's reconsider the Bayesian linear regression model:

$$y_i | x_i, \beta, \sigma^2 \stackrel{\text{ind.}}{\sim} N(\beta_0 + \beta_1 x_i, \sigma^2) \quad i = 1, \dots, n, \quad \beta \sim N(\mathbf{0}, \Sigma_0) \\ \sigma^2 | a \sim \text{Inverse} - \chi^2(1, 1/a) \quad a \sim \text{Inverse} - \chi^2(1, 1/A^2)$$

We showed that the optimal posterior density function for  $\beta$  is

$$q^*(\beta) = C \exp \left[ -\frac{1}{2} \beta^T \left\{ \mathbb{E}_q(1/\sigma^2) \mathbf{X}^T \mathbf{X} + \Sigma_0^{-1} \right\} \beta + \mathbb{E}_q(1/\sigma^2) \beta^T \mathbf{X}^T \mathbf{y} \right] \\ = C \exp \left\{ -\frac{1}{2} \mathbb{E}_q(1/\sigma^2) \beta^T \mathbf{X}^T \mathbf{X} \beta + \mathbb{E}_q(1/\sigma^2) \beta^T \mathbf{X}^T \mathbf{y} \right\} \exp \left\{ -\frac{1}{2} \beta^T \Sigma_0^{-1} \beta \right\}$$

That is, the optimal posterior density function can be factorised into a product of a Gaussian contribution from the likelihood and a Gaussian contribution from the prior distribution.



# Steps Towards Message Passing Algorithms

Let's reconsider the Bayesian linear regression model:

$$y_i | x_i, \beta, \sigma^2 \stackrel{\text{ind.}}{\sim} N(\beta_0 + \beta_1 x_i, \sigma^2) \quad i = 1, \dots, n, \quad \beta \sim N(\mathbf{0}, \Sigma_0) \\ \sigma^2 | a \sim \text{Inverse} - \chi^2(1, 1/a) \quad a \sim \text{Inverse} - \chi^2(1, 1/A^2)$$

We showed that the optimal posterior density function for  $\beta$  is

$$q^*(\beta) = C \exp \left[ -\frac{1}{2} \beta^T \left\{ \mathbb{E}_q(1/\sigma^2) \mathbf{X}^T \mathbf{X} + \Sigma_0^{-1} \right\} \beta + \mathbb{E}_q(1/\sigma^2) \beta^T \mathbf{X}^T \mathbf{y} \right] \\ = C \exp \left\{ -\frac{1}{2} \mathbb{E}_q(1/\sigma^2) \beta^T \mathbf{X}^T \mathbf{X} \beta + \mathbb{E}_q(1/\sigma^2) \beta^T \mathbf{X}^T \mathbf{y} \right\} \exp \left\{ -\frac{1}{2} \beta^T \Sigma_0^{-1} \beta \right\}$$

That is, the optimal posterior density function can be factorised into a product of a Gaussian contribution from the **likelihood** and a Gaussian contribution from the prior distribution.

# Steps Towards Message Passing Algorithms

Let's reconsider the Bayesian linear regression model:

$$y_i | x_i, \beta, \sigma^2 \stackrel{\text{ind.}}{\sim} N(\beta_0 + \beta_1 x_i, \sigma^2) \quad i = 1, \dots, n, \quad \beta \sim N(\mathbf{0}, \Sigma_0) \\ \sigma^2 | a \sim \text{Inverse} - \chi^2(1, 1/a) \quad a \sim \text{Inverse} - \chi^2(1, 1/A^2)$$

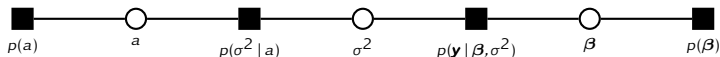
We showed that the optimal posterior density function for  $\beta$  is

$$q^*(\beta) = C \exp \left[ -\frac{1}{2} \beta^T \left\{ \mathbb{E}_q(1/\sigma^2) \mathbf{X}^T \mathbf{X} + \Sigma_0^{-1} \right\} \beta + \mathbb{E}_q(1/\sigma^2) \beta^T \mathbf{X}^T \mathbf{y} \right] \\ = C \exp \left\{ -\frac{1}{2} \mathbb{E}_q(1/\sigma^2) \beta^T \mathbf{X}^T \mathbf{X} \beta + \mathbb{E}_q(1/\sigma^2) \beta^T \mathbf{X}^T \mathbf{y} \right\} \exp \left\{ -\frac{1}{2} \beta^T \Sigma_0^{-1} \beta \right\}$$

That is, the optimal posterior density function can be factorised into a product of a Gaussian contribution from the likelihood and a Gaussian contribution from the **prior distribution**.

# Steps Towards Message Passing Algorithms

This factorised approach is key to the message passing infrastructure:



$$\mathbf{y}|\beta, \sigma^2 \sim N(\mathbf{X}\beta, \sigma^2 \mathbf{I}_n), \quad \beta \sim N(\mathbf{0}, \Sigma_0)$$

$$\sigma^2|a \sim \text{Inverse} - \chi^2(1, 1/a) \quad a \sim \text{Inverse} - \chi^2(1, 1/A^2)$$

We showed that the optimal posterior density function for  $\beta$  is

$$q^*(\beta) = C \exp \left\{ -\frac{1}{2} \mathbb{E}_q(1/\sigma^2) \beta^T \mathbf{X}^T \mathbf{X} \beta + \mathbb{E}_q(1/\sigma^2) \beta^T \mathbf{X}^T \mathbf{y} \right\} \exp \left\{ -\frac{1}{2} \beta^T \Sigma_0^{-1} \beta \right\}$$

We define the messages as

$$m_{p(\mathbf{y}|\beta, \sigma^2) \rightarrow \beta}(\beta) \equiv \exp \left\{ -\frac{1}{2} \mathbb{E}_q(1/\sigma^2) \beta^T \mathbf{X}^T \mathbf{X} \beta + \mathbb{E}_q(1/\sigma^2) \beta^T \mathbf{X}^T \mathbf{y} \right\}$$

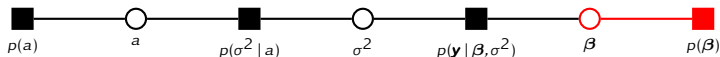
$$m_{p(\beta) \rightarrow \beta}(\beta) \equiv \exp \left\{ -\frac{1}{2} \beta^T \Sigma_0^{-1} \beta \right\}$$

The representation via a directed acyclic graph does not facilitate this factorised form:



# Steps Towards Message Passing Algorithms

This factorised approach is key to the message passing infrastructure:



$$\mathbf{y}|\beta, \sigma^2 \sim N(\mathbf{X}\beta, \sigma^2 \mathbf{I}_n), \quad \beta \sim N(\mathbf{0}, \Sigma_0)$$

$$\sigma^2|a \sim \text{Inverse} - \chi^2(1, 1/a) \quad a \sim \text{Inverse} - \chi^2(1, 1/A^2)$$

We showed that the optimal posterior density function for  $\beta$  is

$$q^*(\beta) = C \exp \left\{ -\frac{1}{2} \mathbb{E}_q(1/\sigma^2) \beta^T \mathbf{X}^T \mathbf{X} \beta + \mathbb{E}_q(1/\sigma^2) \beta^T \mathbf{X}^T \mathbf{y} \right\} \exp \left\{ -\frac{1}{2} \beta^T \Sigma_0^{-1} \beta \right\}$$

We define the messages as

$$m_{p(\mathbf{y}|\beta, \sigma^2) \rightarrow \beta}(\beta) \equiv \exp \left\{ -\frac{1}{2} \mathbb{E}_q(1/\sigma^2) \beta^T \mathbf{X}^T \mathbf{X} \beta + \mathbb{E}_q(1/\sigma^2) \beta^T \mathbf{X}^T \mathbf{y} \right\}$$

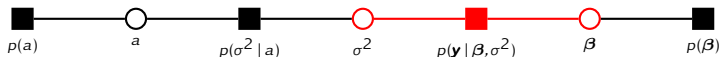
$$m_{p(\beta) \rightarrow \beta}(\beta) \equiv \exp \left\{ -\frac{1}{2} \beta^T \Sigma_0^{-1} \beta \right\}$$

The representation via a directed acyclic graph does not facilitate this factorised form:



# Steps Towards Message Passing Algorithms

This factorised approach is key to the message passing infrastructure:



$$y|\beta, \sigma^2 \sim N(\mathbf{X}\beta, \sigma^2 \mathbf{I}_n), \quad \beta \sim N(\mathbf{0}, \Sigma_0)$$

$$\sigma^2|a \sim \text{Inverse} - \chi^2(1, 1/a) \quad a \sim \text{Inverse} - \chi^2(1, 1/A^2)$$

We showed that the optimal posterior density function for  $\beta$  is

$$q^*(\beta) = C \exp \left\{ -\frac{1}{2} \mathbb{E}_q(1/\sigma^2) \beta^T \mathbf{X}^T \mathbf{X} \beta + \mathbb{E}_q(1/\sigma^2) \beta^T \mathbf{X}^T \mathbf{y} \right\} \exp \left\{ -\frac{1}{2} \beta^T \Sigma_0^{-1} \beta \right\}$$

We define the messages as

$$m_{p(y|\beta, \sigma^2) \rightarrow \beta}(\beta) \equiv \exp \left\{ -\frac{1}{2} \mathbb{E}_q(1/\sigma^2) \beta^T \mathbf{X}^T \mathbf{X} \beta + \mathbb{E}_q(1/\sigma^2) \beta^T \mathbf{X}^T \mathbf{y} \right\}$$

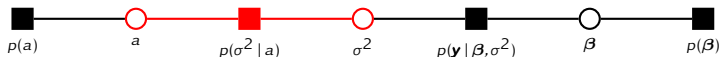
$$m_{p(\beta) \rightarrow \beta}(\beta) \equiv \exp \left\{ -\frac{1}{2} \beta^T \Sigma_0^{-1} \beta \right\}$$

The representation via a directed acyclic graph does not facilitate this factorised form:



# Steps Towards Message Passing Algorithms

This factorised approach is key to the message passing infrastructure:



$$\mathbf{y}|\beta, \sigma^2 \sim N(\mathbf{X}\beta, \sigma^2 \mathbf{I}_n), \quad \beta \sim N(\mathbf{0}, \Sigma_0)$$

$$\sigma^2|a \sim \text{Inverse} - \chi^2(1, 1/a) \quad a \sim \text{Inverse} - \chi^2(1, 1/A^2)$$

We showed that the optimal posterior density function for  $\beta$  is

$$q^*(\beta) = C \exp \left\{ -\frac{1}{2} \mathbb{E}_q(1/\sigma^2) \beta^T \mathbf{X}^T \mathbf{X} \beta + \mathbb{E}_q(1/\sigma^2) \beta^T \mathbf{X}^T \mathbf{y} \right\} \exp \left\{ -\frac{1}{2} \beta^T \Sigma_0^{-1} \beta \right\}$$

We define the messages as

$$m_{p(\mathbf{y}|\beta, \sigma^2) \rightarrow \beta}(\beta) \equiv \exp \left\{ -\frac{1}{2} \mathbb{E}_q(1/\sigma^2) \beta^T \mathbf{X}^T \mathbf{X} \beta + \mathbb{E}_q(1/\sigma^2) \beta^T \mathbf{X}^T \mathbf{y} \right\}$$

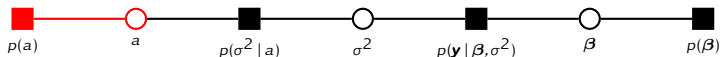
$$m_{p(\beta) \rightarrow \beta}(\beta) \equiv \exp \left\{ -\frac{1}{2} \beta^T \Sigma_0^{-1} \beta \right\}$$

The representation via a directed acyclic graph does not facilitate this factorised form:



# Steps Towards Message Passing Algorithms

This factorised approach is key to the message passing infrastructure:



$$\mathbf{y}|\beta, \sigma^2 \sim N(\mathbf{X}\beta, \sigma^2 \mathbf{I}_n), \quad \beta \sim N(\mathbf{0}, \Sigma_0)$$

$$\sigma^2|a \sim \text{Inverse} - \chi^2(1, 1/a) \quad a \sim \text{Inverse} - \chi^2(1, 1/A^2)$$

We showed that the optimal posterior density function for  $\beta$  is

$$q^*(\beta) = C \exp \left\{ -\frac{1}{2} \mathbb{E}_q(1/\sigma^2) \beta^T \mathbf{X}^T \mathbf{X} \beta + \mathbb{E}_q(1/\sigma^2) \beta^T \mathbf{X}^T \mathbf{y} \right\} \exp \left\{ -\frac{1}{2} \beta^T \Sigma_0^{-1} \beta \right\}$$

We define the messages as

$$m_{p(\mathbf{y}|\beta, \sigma^2) \rightarrow \beta}(\beta) \equiv \exp \left\{ -\frac{1}{2} \mathbb{E}_q(1/\sigma^2) \beta^T \mathbf{X}^T \mathbf{X} \beta + \mathbb{E}_q(1/\sigma^2) \beta^T \mathbf{X}^T \mathbf{y} \right\}$$

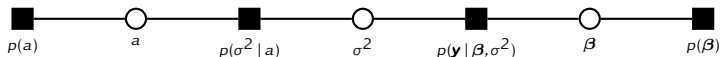
$$m_{p(\beta) \rightarrow \beta}(\beta) \equiv \exp \left\{ -\frac{1}{2} \beta^T \Sigma_0^{-1} \beta \right\}$$

The representation via a directed acyclic graph does not facilitate this factorised form:



# Steps Towards Message Passing Algorithms

This factorised approach is key to the message passing infrastructure:



$$y | \beta, \sigma^2 \sim N(\mathbf{X}\beta, \sigma^2 \mathbf{I}_n), \quad \beta \sim N(\mathbf{0}, \Sigma_0)$$

$$\sigma^2 | a \sim \text{Inverse} - \chi^2(1, 1/a) \quad a \sim \text{Inverse} - \chi^2(1, 1/A^2)$$

We showed that the optimal posterior density function for  $\beta$  is

$$q^*(\beta) = C \exp \left\{ -\frac{1}{2} \mathbb{E}_q(1/\sigma^2) \beta^T \mathbf{X}^T \mathbf{X} \beta + \mathbb{E}_q(1/\sigma^2) \beta^T \mathbf{X}^T \mathbf{y} \right\} \exp \left\{ -\frac{1}{2} \beta^T \Sigma_0^{-1} \beta \right\}$$

We define the messages as

$$m_{p(y|\beta, \sigma^2) \rightarrow \beta}(\beta) \equiv \exp \left\{ -\frac{1}{2} \mathbb{E}_q(1/\sigma^2) \beta^T \mathbf{X}^T \mathbf{X} \beta + \mathbb{E}_q(1/\sigma^2) \beta^T \mathbf{X}^T \mathbf{y} \right\}$$

$$m_{p(\beta) \rightarrow \beta}(\beta) \equiv \exp \left\{ -\frac{1}{2} \beta^T \Sigma_0^{-1} \beta \right\}$$

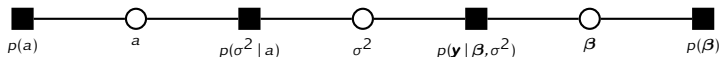
The representation via a directed acyclic graph does not facilitate this factorised form:





# Steps Towards Message Passing Algorithms

This factorised approach is key to the message passing infrastructure:



$$\mathbf{y}|\beta, \sigma^2 \sim N(\mathbf{X}\beta, \sigma^2 \mathbf{I}_n), \quad \beta \sim N(\mathbf{0}, \Sigma_0)$$

$$\sigma^2|a \sim \text{Inverse} - \chi^2(1, 1/a) \quad a \sim \text{Inverse} - \chi^2(1, 1/A^2)$$

We showed that the optimal posterior density function for  $\beta$  is

$$q^*(\beta) = C \exp \left\{ -\frac{1}{2} \mathbb{E}_q(1/\sigma^2) \beta^T \mathbf{X}^T \mathbf{X} \beta + \mathbb{E}_q(1/\sigma^2) \beta^T \mathbf{X}^T \mathbf{y} \right\} \exp \left\{ -\frac{1}{2} \beta^T \Sigma_0^{-1} \beta \right\}$$

We define the messages as

$$m_{p(\mathbf{y}|\beta, \sigma^2) \rightarrow \beta}(\beta) \equiv \exp \left\{ -\frac{1}{2} \mathbb{E}_q(1/\sigma^2) \beta^T \mathbf{X}^T \mathbf{X} \beta + \mathbb{E}_q(1/\sigma^2) \beta^T \mathbf{X}^T \mathbf{y} \right\}$$

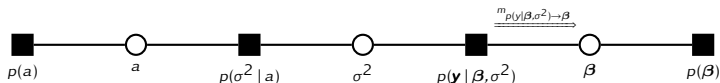
$$m_{p(\beta) \rightarrow \beta}(\beta) \equiv \exp \left\{ -\frac{1}{2} \beta^T \Sigma_0^{-1} \beta \right\}$$

The representation via a directed acyclic graph does not facilitate this factorised form:



# Steps Towards Message Passing Algorithms

This factorised approach is key to the message passing infrastructure:



$$\mathbf{y}|\beta, \sigma^2 \sim N(\mathbf{X}\beta, \sigma^2 \mathbf{I}_n), \quad \beta \sim N(\mathbf{0}, \Sigma_0)$$

$$\sigma^2|a \sim \text{Inverse} - \chi^2(1, 1/a) \quad a \sim \text{Inverse} - \chi^2(1, 1/A^2)$$

We showed that the optimal posterior density function for  $\beta$  is

$$q^*(\beta) = C \exp \left\{ -\frac{1}{2} \mathbb{E}_q(1/\sigma^2) \beta^T \mathbf{X}^T \mathbf{X} \beta + \mathbb{E}_q(1/\sigma^2) \beta^T \mathbf{X}^T \mathbf{y} \right\} \exp \left\{ -\frac{1}{2} \beta^T \Sigma_0^{-1} \beta \right\}$$

We define the messages as

$$m_{p(\mathbf{y}|\beta, \sigma^2) \rightarrow \beta}(\beta) \equiv \exp \left\{ -\frac{1}{2} \mathbb{E}_q(1/\sigma^2) \beta^T \mathbf{X}^T \mathbf{X} \beta + \mathbb{E}_q(1/\sigma^2) \beta^T \mathbf{X}^T \mathbf{y} \right\}$$

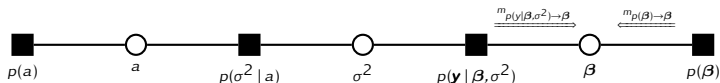
$$m_{p(\beta) \rightarrow \beta}(\beta) \equiv \exp \left\{ -\frac{1}{2} \beta^T \Sigma_0^{-1} \beta \right\}$$

The representation via a directed acyclic graph does not facilitate this factorised form:



# Steps Towards Message Passing Algorithms

This factorised approach is key to the message passing infrastructure:



$$\mathbf{y}|\beta, \sigma^2 \sim N(\mathbf{X}\beta, \sigma^2 \mathbf{I}_n), \quad \beta \sim N(\mathbf{0}, \Sigma_0)$$

$$\sigma^2|a \sim \text{Inverse} - \chi^2(1, 1/a) \quad a \sim \text{Inverse} - \chi^2(1, 1/A^2)$$

We showed that the optimal posterior density function for  $\beta$  is

$$q^*(\beta) = C \exp \left\{ -\frac{1}{2} \mathbb{E}_q(1/\sigma^2) \beta^T \mathbf{X}^T \mathbf{X} \beta + \mathbb{E}_q(1/\sigma^2) \beta^T \mathbf{X}^T \mathbf{y} \right\} \exp \left\{ -\frac{1}{2} \beta^T \Sigma_0^{-1} \beta \right\}$$

We define the messages as

$$m_{p(y|\beta, \sigma^2) \rightarrow \beta}(\beta) \equiv \exp \left\{ -\frac{1}{2} \mathbb{E}_q(1/\sigma^2) \beta^T \mathbf{X}^T \mathbf{X} \beta + \mathbb{E}_q(1/\sigma^2) \beta^T \mathbf{X}^T \mathbf{y} \right\}$$

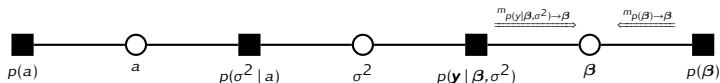
$$m_{p(\beta) \rightarrow \beta}(\beta) \equiv \exp \left\{ -\frac{1}{2} \beta^T \Sigma_0^{-1} \beta \right\}$$

The representation via a directed acyclic graph does not facilitate this factorised form:



# Steps Towards Message Passing Algorithms

This factorised approach is key to the message passing infrastructure:



$$\mathbf{y}|\beta, \sigma^2 \sim N(\mathbf{X}\beta, \sigma^2 \mathbf{I}_n), \quad \beta \sim N(\mathbf{0}, \Sigma_0)$$

$$\sigma^2|a \sim \text{Inverse} - \chi^2(1, 1/a) \quad a \sim \text{Inverse} - \chi^2(1, 1/A^2)$$

We showed that the optimal posterior density function for  $\beta$  is

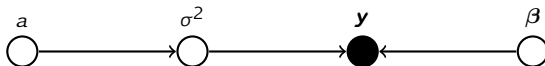
$$q^*(\beta) = C \exp \left\{ -\frac{1}{2} \mathbb{E}_q(1/\sigma^2) \beta^T \mathbf{X}^T \mathbf{X} \beta + \mathbb{E}_q(1/\sigma^2) \beta^T \mathbf{X}^T \mathbf{y} \right\} \exp \left\{ -\frac{1}{2} \beta^T \Sigma_0^{-1} \beta \right\}$$

We define the messages as

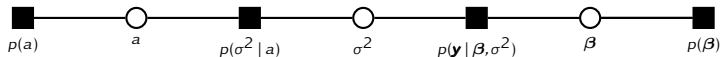
$$m_{p(\mathbf{y}|\beta, \sigma^2) \rightarrow \beta}(\beta) \equiv \exp \left\{ -\frac{1}{2} \mathbb{E}_q(1/\sigma^2) \beta^T \mathbf{X}^T \mathbf{X} \beta + \mathbb{E}_q(1/\sigma^2) \beta^T \mathbf{X}^T \mathbf{y} \right\}$$

$$m_{p(\beta) \rightarrow \beta}(\beta) \equiv \exp \left\{ -\frac{1}{2} \beta^T \Sigma_0^{-1} \beta \right\}$$

The representation via a directed acyclic graph does not facilitate this factorised form:



# General Message Passing Rules



The message passed from a factor  $f$  to its neighbouring parameter  $\theta$  is

$$m_{f \rightarrow \theta}(\theta) \propto \exp\{\mathbb{E}_{-q(\theta)}(\log f)\}$$

For example, the message from the Gaussian likelihood factor  $p(y|\beta, \sigma^2)$  to the parameter  $\beta$  is

$$m_{p(y|\beta, \sigma^2) \rightarrow \beta}(\beta) \propto \exp[\mathbb{E}_{-q(\beta)}\{\log p(y|\beta, \sigma^2)\}]$$

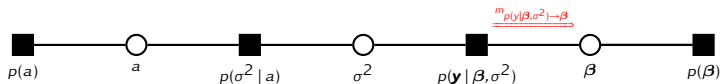
The message that  $\theta$  will pass on to another neighbouring factor is simply the message that it received from the factor  $f$ . For example, the message that  $\beta$  passes on to  $p(\beta)$  is

$$m_{\beta \rightarrow p(\beta)}(\beta) = m_{p(y|\beta, \sigma^2) \rightarrow \beta}(\beta)$$

The approximate  $q$ -density function for a parameter  $\theta$  is the product of all messages that it received. For example,  $q(\beta)$  is

$$q(\beta) = m_{p(y|\beta, \sigma^2) \rightarrow \beta}(\beta) \times m_{\beta \rightarrow p(\beta)}(\beta)$$

# General Message Passing Rules



The message passed from a factor  $f$  to its neighbouring parameter  $\theta$  is

$$m_{f \rightarrow \theta}(\theta) \propto \exp\left\{\mathbb{E}_{-q(\theta)}(\log f)\right\}$$

For example, the message from the Gaussian likelihood factor  $p(\mathbf{y}|\beta, \sigma^2)$  to the parameter  $\beta$  is

$$m_{p(\mathbf{y}|\beta, \sigma^2) \rightarrow \beta}(\beta) \propto \exp\left[\mathbb{E}_{-q(\beta)}\{\log p(\mathbf{y}|\beta, \sigma^2)\}\right]$$

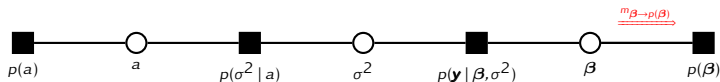
The message that  $\theta$  will pass on to another neighbouring factor is simply the message that it received from the factor  $f$ . For example, the message that  $\beta$  passes on to  $p(\beta)$  is

$$m_{\beta \rightarrow p(\beta)}(\beta) = m_{p(\mathbf{y}|\beta, \sigma^2) \rightarrow \beta}(\beta)$$

The approximate  $q$ -density function for a parameter  $\theta$  is the product of all messages that it received. For example,  $q(\beta)$  is

$$q(\beta) = m_{p(\mathbf{y}|\beta, \sigma^2) \rightarrow \beta}(\beta) \times m_{\beta \rightarrow p(\beta)}(\beta)$$

# General Message Passing Rules



The message passed from a factor  $f$  to its neighbouring parameter  $\theta$  is

$$m_{f \rightarrow \theta}(\theta) \propto \exp\{\mathbb{E}_{-q(\theta)}(\log f)\}$$

For example, the message from the Gaussian likelihood factor  $p(y|\beta, \sigma^2)$  to the parameter  $\beta$  is

$$m_{p(y|\beta, \sigma^2) \rightarrow \beta}(\beta) \propto \exp[\mathbb{E}_{-q(\beta)}\{\log p(y|\beta, \sigma^2)\}]$$

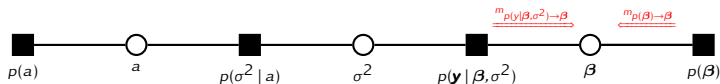
The message that  $\theta$  will pass on to another neighbouring factor is simply the message that it received from the factor  $f$ . For example, the message that  $\beta$  passes on to  $p(\beta)$  is

$$m_{\beta \rightarrow p(\beta)}(\beta) = m_{p(y|\beta, \sigma^2) \rightarrow \beta}(\beta)$$

The approximate  $q$ -density function for a parameter  $\theta$  is the product of all messages that it received. For example,  $q(\beta)$  is

$$q(\beta) = m_{p(y|\beta, \sigma^2) \rightarrow \beta}(\beta) \times m_{\beta \rightarrow p(\beta)}(\beta)$$

# General Message Passing Rules



The message passed from a factor  $f$  to its neighbouring parameter  $\theta$  is

$$m_{f \rightarrow \theta}(\theta) \propto \exp\left\{\mathbb{E}_{-q(\theta)}(\log f)\right\}$$

For example, the message from the Gaussian likelihood factor  $p(y|\beta, \sigma^2)$  to the parameter  $\beta$  is

$$m_{p(y|\beta, \sigma^2) \rightarrow \beta}(\beta) \propto \exp\left[\mathbb{E}_{-q(\beta)}\{\log p(y|\beta, \sigma^2)\}\right]$$

The message that  $\theta$  will pass on to another neighbouring factor is simply the message that it received from the factor  $f$ . For example, the message that  $\beta$  passes on to  $p(\beta)$  is

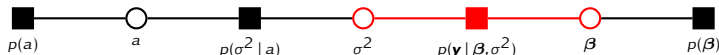
$$m_{\beta \rightarrow p(\beta)}(\beta) = m_{p(y|\beta, \sigma^2) \rightarrow \beta}(\beta)$$

The approximate  $q$ -density function for a parameter  $\theta$  is the product of all messages that it received. For example,  $q(\beta)$  is

$$q(\beta) = m_{p(y|\beta, \sigma^2) \rightarrow \beta}(\beta) \times m_{\beta \rightarrow p(\beta)}(\beta)$$



# The Gaussian Likelihood Fragment



The message from the likelihood specification to the parameter  $\beta$  is

$$m_{p(\mathbf{y}|\beta,\sigma^2)\rightarrow\beta}(\beta) \propto \exp \left\{ \left[ \begin{array}{c} \beta \\ \text{vec}(\beta\beta^T) \end{array} \right]^T \left[ \begin{array}{c} \mathbb{E}_q(1/\sigma^2)\mathbf{X}^T\mathbf{y} \\ -\frac{1}{2}\mathbb{E}_q(1/\sigma^2)\text{vec}(\mathbf{X}^T\mathbf{X}) \end{array} \right] \right\}$$

That is,  $m_{p(\mathbf{y}|\beta,\sigma^2)\rightarrow\beta}(\beta)$  is a multivariate normal density function in exponential family form:

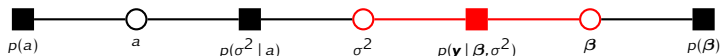
$$p(\mathbf{x}) = C \exp\{\mathbf{T}(\mathbf{x})^T \boldsymbol{\eta}\}.$$

The message to  $\sigma^2$  is

$$m_{p(\mathbf{y}|\beta,\sigma^2)\rightarrow\sigma^2}(\sigma^2) = \exp \left\{ \left[ \begin{array}{c} \log(\sigma^2) \\ 1/\sigma^2 \end{array} \right]^T \left[ \begin{array}{c} -n/2 \\ -\frac{1}{2}\mathbb{E}_q(\|\mathbf{y} - \mathbf{X}\beta\|^2) \end{array} \right] \right\},$$

which is an inverse chi-squared density function in exponential family form.

# The Gaussian Likelihood Fragment



The message from the likelihood specification to the parameter  $\beta$  is

$$m_{p(\mathbf{y}|\beta,\sigma^2)\rightarrow\beta}(\beta) \propto \exp \left\{ \begin{bmatrix} \beta \\ \text{vec}(\beta\beta^T) \end{bmatrix}^T \begin{bmatrix} \mathbb{E}_q(1/\sigma^2)\mathbf{X}^T\mathbf{y} \\ -\frac{1}{2}\mathbb{E}_q(1/\sigma^2)\text{vec}(\mathbf{X}^T\mathbf{X}) \end{bmatrix} \right\}$$

That is,  $m_{p(\mathbf{y}|\beta,\sigma^2)\rightarrow\beta}(\beta)$  is a multivariate normal density function in exponential family form:

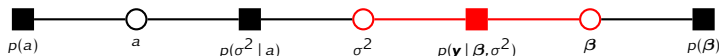
$$p(\mathbf{x}) = C \exp\{\mathbf{T}(\mathbf{x})^T \boldsymbol{\eta}\}.$$

The message to  $\sigma^2$  is

$$m_{p(\mathbf{y}|\beta,\sigma^2)\rightarrow\sigma^2}(\sigma^2) = \exp \left\{ \begin{bmatrix} \log(\sigma^2) \\ 1/\sigma^2 \end{bmatrix}^T \begin{bmatrix} -n/2 \\ -\frac{1}{2}\mathbb{E}_q(\|\mathbf{y} - \mathbf{X}\beta\|^2) \end{bmatrix} \right\},$$

which is an inverse chi-squared density function in exponential family form.

# The Gaussian Likelihood Fragment



The message from the likelihood specification to the parameter  $\beta$  is

$$m_{p(\mathbf{y}|\beta,\sigma^2)\rightarrow\beta}(\beta) \propto \exp \left\{ \begin{bmatrix} \beta \\ \text{vec}(\beta\beta^T) \end{bmatrix}^T \begin{bmatrix} \mathbb{E}_q(1/\sigma^2)\mathbf{X}^T\mathbf{y} \\ -\frac{1}{2}\mathbb{E}_q(1/\sigma^2)\text{vec}(\mathbf{X}^T\mathbf{X}) \end{bmatrix} \right\}$$

That is,  $m_{p(\mathbf{y}|\beta,\sigma^2)\rightarrow\beta}(\beta)$  is a multivariate normal density function in exponential family form:

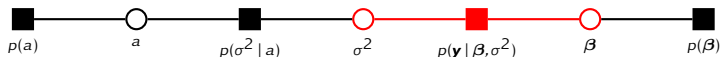
$$p(\mathbf{x}) = C \exp\{\mathbf{T}(\mathbf{x})^T \boldsymbol{\eta}\}.$$

The message to  $\sigma^2$  is

$$m_{p(\mathbf{y}|\beta,\sigma^2)\rightarrow\sigma^2}(\sigma^2) = \exp \left\{ \begin{bmatrix} \log(\sigma^2) \\ 1/\sigma^2 \end{bmatrix}^T \begin{bmatrix} -n/2 \\ -\frac{1}{2}\mathbb{E}_q(\|\mathbf{y} - \mathbf{X}\beta\|^2) \end{bmatrix} \right\},$$

which is an inverse chi-squared density function in exponential family form.

# The Gaussian Likelihood Fragment



By restricting the form of the messages passing updates to distributions in the exponential family form, we can completely characterise these message updates by their natural parameter vectors:

$$\eta_{p(\mathbf{y}|\beta,\sigma^2)\rightarrow\beta} = \begin{bmatrix} \mathbb{E}_q(1/\sigma^2)\mathbf{X}^T\mathbf{y} \\ -\frac{1}{2}\mathbb{E}_q(1/\sigma^2)\text{vec}(\mathbf{X}^T\mathbf{X}) \end{bmatrix}$$
$$\eta_{p(\mathbf{y}|\beta,\sigma^2)\rightarrow\sigma^2} = \begin{bmatrix} -n/2 \\ -\frac{1}{2}\mathbb{E}_q(\|\mathbf{y}-\mathbf{X}\beta\|^2) \end{bmatrix}$$

# Gaussian Density Function in Exponential Family Form

$$\mathbf{x}|\boldsymbol{\mu}, \Sigma \sim \mathcal{N}(\boldsymbol{\mu}, \Sigma)$$

The density function in exponential family form is:

$$p(\mathbf{x}|\boldsymbol{\mu}, \Sigma) = C \exp \left\{ \left[ \begin{array}{c} \boldsymbol{\beta} \\ \text{vec}(\boldsymbol{\beta}\boldsymbol{\beta}^\top) \end{array} \right]^\top \boldsymbol{\eta} \right\}, \quad \boldsymbol{\eta} = \begin{bmatrix} \boldsymbol{\eta}_1 \\ \boldsymbol{\eta}_2 \end{bmatrix} = \begin{bmatrix} \Sigma^{-1} \boldsymbol{\mu} \\ -\frac{1}{2} \text{vec}(\Sigma^{-1}) \end{bmatrix}$$

The transformation back to the common parameters is

$$\Sigma = -\frac{1}{2} \{\text{vec}^{-1}(\boldsymbol{\eta}_2)\}^{-1}, \quad \boldsymbol{\mu} = \Sigma \boldsymbol{\eta}_1$$

# Inverse Wishart Density Function in Exponential Family Form

$$\Sigma|\kappa, \mathbf{\Lambda} \sim \text{Inverse Wishart}(\kappa, \mathbf{\Lambda})$$

The density function in exponential family form is:

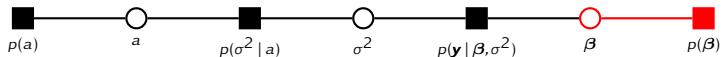
$$p(\Sigma|\kappa, \mathbf{\Lambda}) = C \exp \left\{ \begin{bmatrix} \log|\Sigma| \\ \text{vec}(\Sigma^{-1}) \end{bmatrix}^T \boldsymbol{\eta} \right\}, \quad \boldsymbol{\eta} = \begin{bmatrix} \eta_1 \\ \eta_2 \end{bmatrix} = \begin{bmatrix} -\frac{1}{2}(\kappa + 2) \\ -\frac{1}{2} \text{vec}(\mathbf{\Lambda}) \end{bmatrix}$$

The transformation back to the common parameters is

$$\kappa = -2\eta_1 - 2, \quad \mathbf{\Lambda} = -2\text{vec}^{-1}(\eta_2)$$

For variational inference, we require:

$$\mathbb{E}(\Sigma^{-1}) = \begin{cases} \{\eta_1 + \frac{1}{2}(d+1)\}\{\text{vec}^{-1}(\eta_2)\}^{-1}, & \text{if } \Sigma \text{ is a complete matrix} \\ \{\eta_1 + 1\}\{\text{vec}^{-1}(\eta_2)\}^{-1}, & \text{if } \Sigma \text{ is a diagonal; matrix} \end{cases}$$

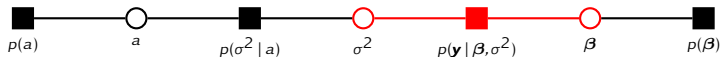


Gaussian prior fragment (Wand, 2017)

Gaussian likelihood fragment (Wand, 2017)

Iterated inverse Wishart fragment (Maestrini and Wand, 2020)

Inverse Wishart prior fragment (Maestrini and Wand, 2020)



Gaussian prior fragment (Wand, 2017)

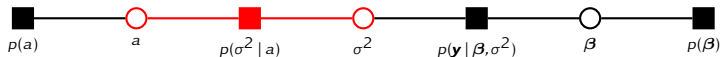
Gaussian likelihood fragment (Wand, 2017)

Iterated inverse Wishart fragment (Maestrini and Wand, 2020)

Inverse Wishart prior fragment (Maestrini and Wand, 2020)



# Fragments



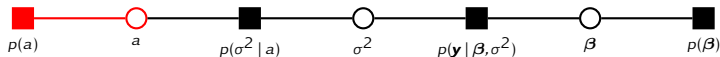
Gaussian prior fragment (Wand, 2017)

Gaussian likelihood fragment (Wand, 2017)

Iterated inverse Wishart fragment (Maestrini and Wand, 2020)

Inverse Wishart prior fragment (Maestrini and Wand, 2020)

# Fragments



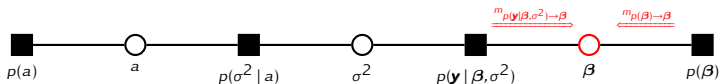
Gaussian prior fragment (Wand, 2017)

Gaussian likelihood fragment (Wand, 2017)

Iterated inverse Wishart fragment (Maestrini and Wand, 2020)

Inverse Wishart prior fragment (Maestrini and Wand, 2020)

## $q$ -Density Functions



Recall that the approximate  $q$ -density function for each parameter is the product of all the messages that it received.

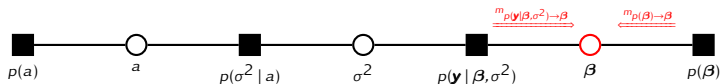
The computation for  $q(\beta)$  is

$$\begin{aligned} q(\beta) &= m_{p(y|\beta, \sigma^2) \rightarrow \beta}(\beta) \times m_{p(\beta) \rightarrow \beta}(\beta) \\ &= \exp \left\{ \begin{bmatrix} \beta \\ \text{vec}(\beta \beta^T) \end{bmatrix}^T (\eta_{p(y|\beta, \sigma^2) \rightarrow \beta} + \eta_{p(\beta) \rightarrow \beta}) \right\} \end{aligned}$$

So the  $q$ -density function is also in the exponential family of density functions with natural parameter vector:

$$\eta_{q^*}(\beta) = \eta_{p(y|\beta, \sigma^2) \rightarrow \beta} + \eta_{p(\beta) \rightarrow \beta}$$

# $q$ -Density Functions



Recall that the approximate  $q$ -density function for each parameter is the product of all the messages that it received.

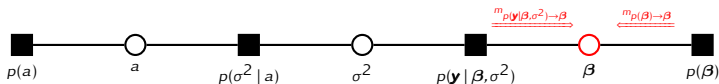
The computation for  $q(\beta)$  is

$$\begin{aligned} q(\beta) &= m_{p(\mathbf{y}|\beta,\sigma^2)\rightarrow\beta}(\beta) \times m_{p(\beta)\rightarrow\beta}(\beta) \\ &= \exp \left\{ \begin{bmatrix} \beta \\ \text{vec}(\beta\beta^T) \end{bmatrix}^T (\eta_{p(\mathbf{y}|\beta,\sigma^2)\rightarrow\beta} + \eta_{p(\beta)\rightarrow\beta}) \right\} \end{aligned}$$

So the  $q$ -density function is also in the exponential family of density functions with natural parameter vector:

$$\eta_{q^*}(\beta) = \eta_{p(\mathbf{y}|\beta,\sigma^2)\rightarrow\beta} + \eta_{p(\beta)\rightarrow\beta}$$

# $q$ -Density Functions



Recall that the approximate  $q$ -density function for each parameter is the product of all the messages that it received.

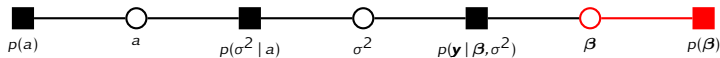
The computation for  $q(\beta)$  is

$$\begin{aligned} q(\beta) &= m_{p(\mathbf{y}|\beta,\sigma^2)\rightarrow\beta}(\beta) \times m_{p(\beta)\rightarrow\beta}(\beta) \\ &= \exp \left\{ \begin{bmatrix} \beta \\ \text{vec}(\beta\beta^T) \end{bmatrix}^T (\eta_{p(\mathbf{y}|\beta,\sigma^2)\rightarrow\beta} + \eta_{p(\beta)\rightarrow\beta}) \right\} \end{aligned}$$

So the  $q$ -density function is also in the exponential family of density functions with natural parameter vector:

$$\eta_{q^*}(\beta) = \eta_{p(\mathbf{y}|\beta,\sigma^2)\rightarrow\beta} + \eta_{p(\beta)\rightarrow\beta}$$

# The Gaussian Prior Fragment



$$\beta \sim N(\mathbf{0}, \sigma_\beta^2 \mathbf{I})$$

Inputs:

$$\sigma_\beta^2$$

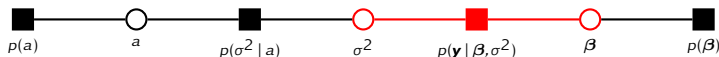
Updates:

$$\eta_{p(\beta) \rightarrow \beta} = \begin{bmatrix} \mathbf{0} \\ -\frac{1}{2\sigma_\beta^2} \text{vec}(\mathbf{I}) \end{bmatrix}$$

Outputs:

$$\eta_{p(\beta) \rightarrow \beta}$$

# The Gaussian Likelihood Fragment



$$\mathbf{y}|\beta, \sigma^2 \sim N(\mathbf{X}\beta, \sigma^2 \mathbf{I})$$

Inputs:

$$\boldsymbol{\eta}_{q^*}(\beta), \quad \boldsymbol{\eta}_{q^*}(\sigma^2)$$

Updates:

$$\text{Cov}_q(\beta) = -\frac{1}{2} \{\text{vec}^{-1}(\boldsymbol{\eta}_{q^*}(\beta))_2\}^{-1}, \quad \mathbb{E}_q(\beta) = \text{Cov}_q(\beta)(\boldsymbol{\eta}_{q^*}(\beta))_1,$$

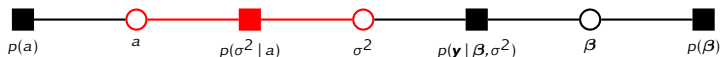
$$\mathbb{E}_q(1/\sigma^2) = \frac{(\boldsymbol{\eta}_{q^*}(\sigma^2))_1 + 1}{(\boldsymbol{\eta}_{q^*}(\sigma^2))_2}$$

$$\boldsymbol{\eta}_{p(\mathbf{y}|\beta, \sigma^2) \rightarrow \beta} = \begin{bmatrix} \mathbb{E}_q(1/\sigma^2) \mathbf{X}^T \mathbf{y} \\ -\frac{1}{2} \mathbb{E}_q(1/\sigma^2) \text{vec}(\mathbf{X}^T \mathbf{X}) \end{bmatrix}, \quad \boldsymbol{\eta}_{p(\mathbf{y}|\beta, \sigma^2) \rightarrow \sigma^2} = \begin{bmatrix} -n/2 \\ -\frac{1}{2} \mathbb{E}_q(\|\mathbf{y} - \mathbf{X}\beta\|^2) \end{bmatrix}$$

Outputs:

$$\boldsymbol{\eta}_{p(\mathbf{y}|\beta, \sigma^2) \rightarrow \beta}, \quad \boldsymbol{\eta}_{p(\mathbf{y}|\beta, \sigma^2) \rightarrow \sigma^2}$$

# Iterated Inverse Wishart Fragment



$$\Sigma | \mathbf{A} \sim \text{Inverse-Wishart}(\kappa, \mathbf{A}^{-1})$$

Inputs:

$$\boldsymbol{\eta}_{q^*}(\Sigma), \quad \boldsymbol{\eta}_{q^*}(\mathbf{A})$$

Updates:

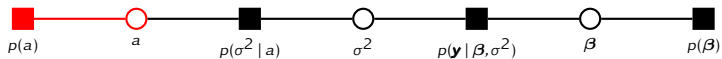
$$\begin{aligned} \mathbb{E}_q(\Sigma) &= \{(\boldsymbol{\eta}_{q^*}(\Sigma))_1 + 1\} \{\text{vec}^{-1}(\boldsymbol{\eta}_{q^*}(\Sigma))_2\}^{-1}, \\ \mathbb{E}_q(\mathbf{A}) &= \{(\boldsymbol{\eta}_{q^*}(\mathbf{A}))_1 + 1\} \{\text{vec}^{-1}(\boldsymbol{\eta}_{q^*}(\mathbf{A}))_2\}^{-1}, \\ \boldsymbol{\eta}_{p(\Sigma|\mathbf{A}) \rightarrow \Sigma} &= \begin{bmatrix} -\frac{1}{2}(\kappa + d + 1) \\ -\frac{1}{2}\{\text{vec } \mathbb{E}_q(\mathbf{A}^{-1})\} \end{bmatrix}, \quad \boldsymbol{\eta}_{p(\Sigma|\mathbf{A}) \rightarrow \mathbf{A}} = \begin{bmatrix} -\kappa/2 \\ -\frac{1}{2}\{\text{vec } \mathbb{E}_q(\Sigma^{-1})\} \end{bmatrix} \end{aligned}$$

Outputs:

$$\boldsymbol{\eta}_{p(\Sigma|\mathbf{A}) \rightarrow \Sigma}, \quad \boldsymbol{\eta}_{p(\Sigma|\mathbf{A}) \rightarrow \mathbf{A}}$$



# Inverse Wishart Prior Fragment



$$\mathbf{A} \sim \text{Inverse-Wishart}(\kappa, \mathbf{\Lambda})$$

Inputs:

$$\kappa, \mathbf{\Lambda}$$

Updates:

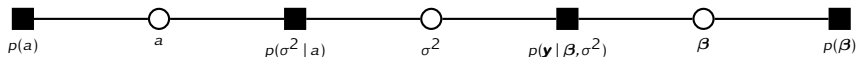
$$\eta_{p(\mathbf{A}) \rightarrow \mathbf{A}} = \begin{bmatrix} -\frac{1}{2}(\kappa + d + 1) \\ -\frac{1}{2} \text{vec}(\mathbf{\Lambda}) \end{bmatrix}$$

Outputs:

$$\eta_{p(\mathbf{A}) \rightarrow \mathbf{A}}$$

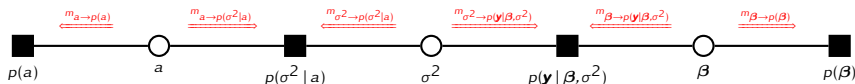
## VMP for the Bayesian linear regression model

1. Initialise all messages from stochastic nodes to factors
2. Cycle:
  - (i) Update all messages from factors to stochastic nodes
  - (ii) Update all messages from stochastic nodes to factors
  - (iii) Update all optimal posterior density functions
3. Stop:  $D_{\text{KL}} \{q(\beta, \sigma^2, a) || p(\beta, \sigma^2, a | y)\}$  converges.



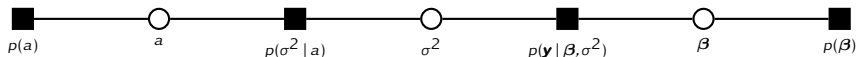
## VMP for the Bayesian linear regression model

1. Initialise all messages from stochastic nodes to factors
2. Cycle:
  - (i) Update all messages from factors to stochastic nodes
  - (ii) Update all messages from stochastic nodes to factors
  - (iii) Update all optimal posterior density functions
3. Stop:  $D_{\text{KL}} \{q(\beta, \sigma^2, a) \| p(\beta, \sigma^2, a | \mathbf{y})\}$  converges.



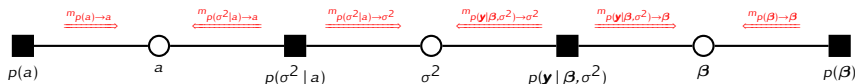
## VMP for the Bayesian linear regression model

1. Initialise all messages from stochastic nodes to factors
2. Cycle:
  - (i) Update all messages from factors to stochastic nodes
  - (ii) Update all messages from stochastic nodes to factors
  - (iii) Update all optimal posterior density functions
3. Stop:  $D_{\text{KL}} \{q(\beta, \sigma^2, a) \| p(\beta, \sigma^2, a | \mathbf{y})\}$  converges.



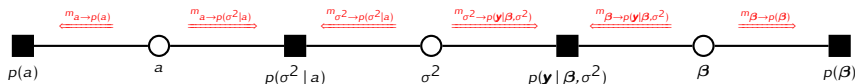
## VMP for the Bayesian linear regression model

1. Initialise all messages from stochastic nodes to factors
2. Cycle:
  - (i) Update all messages from factors to stochastic nodes
  - (ii) Update all messages from stochastic nodes to factors
  - (iii) Update all optimal posterior density functions
3. Stop:  $D_{\text{KL}} \{q(\beta, \sigma^2, a) \| p(\beta, \sigma^2, a | y)\}$  converges.



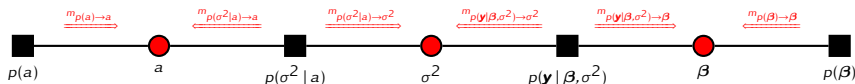
## VMP for the Bayesian linear regression model

1. Initialise all messages from stochastic nodes to factors
2. Cycle:
  - (i) Update all messages from factors to stochastic nodes
  - (ii) Update all messages from stochastic nodes to factors
  - (iii) Update all optimal posterior density functions
3. Stop:  $D_{\text{KL}} \{q(\beta, \sigma^2, a) \| p(\beta, \sigma^2, a | y)\}$  converges.



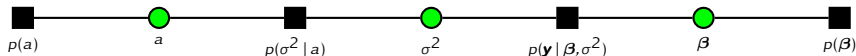
## VMP for the Bayesian linear regression model

1. Initialise all messages from stochastic nodes to factors
2. Cycle:
  - (i) Update all messages from factors to stochastic nodes
  - (ii) Update all messages from stochastic nodes to factors
  - (iii) Update all optimal posterior density functions
3. Stop:  $D_{\text{KL}} \{q(\beta, \sigma^2, a) \| p(\beta, \sigma^2, a | y)\}$  converges.



## VMP for the Bayesian linear regression model

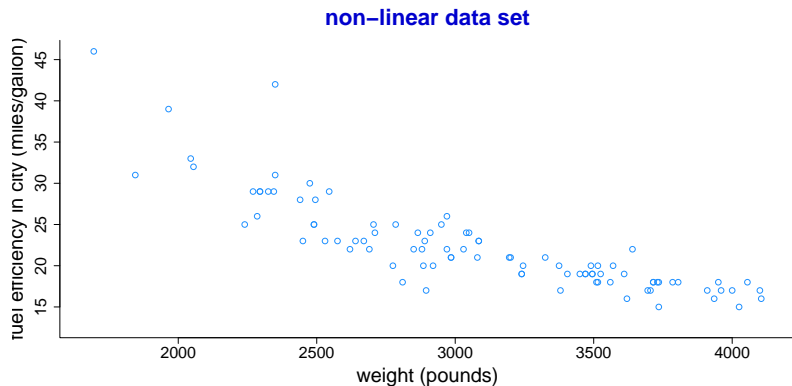
1. Initialise all messages from stochastic nodes to factors
2. Cycle:
  - (i) Update all messages from factors to stochastic nodes
  - (ii) Update all messages from stochastic nodes to factors
  - (iii) Update all optimal posterior density functions
3. Stop:  $D_{\text{KL}} \{q(\beta, \sigma^2, a) \| p(\beta, \sigma^2, a | \mathbf{y})\}$  converges.





## Part II

### Building Model Complexity



# Bayesian Semiparametric Regression

We can address nonlinear data sets using semiparametric regression techniques

Let  $x_i$  be the  $i$ th vehicle weight and  $y_i$  be its corresponding fuel efficiency score. The corresponding vectors are  $\mathbf{x}$  and  $\mathbf{y}$ .

We construct a fixed effects matrix

$$\mathbf{X} = \begin{bmatrix} 1 & x_1 \\ \vdots & \vdots \\ 1 & x_n \end{bmatrix}$$

Nonlinear effects are incorporated through a random effects matrix

$$\mathbf{Z} = \begin{bmatrix} z_1(x_1) & \dots & z_K(x_1) \\ \vdots & \ddots & \vdots \\ z_1(x_N) & \dots & z_K(x_N) \end{bmatrix},$$

where  $\{z_k\}_{k=1}^K$  is a suitable spline basis – see Wand and Ormerod (2008) for a description of canonical cubic O’Sullivan penalised splines.

# Bayesian Semiparametric Regression

We can address nonlinear data sets using semiparametric regression techniques

Let  $x_i$  be the  $i$ th vehicle weight and  $y_i$  be its corresponding fuel efficiency score. The corresponding vectors are  $\mathbf{x}$  and  $\mathbf{y}$ .

We construct a fixed effects matrix

$$\mathbf{X} = \begin{bmatrix} 1 & x_1 \\ \vdots & \vdots \\ 1 & x_n \end{bmatrix}$$

Nonlinear effects are incorporated through a random effects matrix

$$\mathbf{Z} = \begin{bmatrix} z_1(x_1) & \dots & z_K(x_1) \\ \vdots & \ddots & \vdots \\ z_1(x_N) & \dots & z_K(x_N) \end{bmatrix},$$

where  $\{z_k\}_{k=1}^K$  is a suitable spline basis – see Wand and Ormerod (2008) for a description of canonical cubic O’Sullivan penalised splines.

# Bayesian Semiparametric Regression

We can address nonlinear data sets using semiparametric regression techniques

Let  $x_i$  be the  $i$ th vehicle weight and  $y_i$  be its corresponding fuel efficiency score. The corresponding vectors are  $\mathbf{x}$  and  $\mathbf{y}$ .

We construct a fixed effects matrix

$$\mathbf{X} = \begin{bmatrix} 1 & x_1 \\ \vdots & \vdots \\ 1 & x_n \end{bmatrix}$$

Nonlinear effects are incorporated through a random effects matrix

$$\mathbf{Z} = \begin{bmatrix} z_1(x_1) & \dots & z_K(x_1) \\ \vdots & \ddots & \vdots \\ z_1(x_N) & \dots & z_K(x_N) \end{bmatrix},$$

where  $\{z_k\}_{k=1}^K$  is a suitable spline basis – see Wand and Ormerod (2008) for a description of canonical cubic O’Sullivan penalised splines.

# Bayesian Semiparametric Regression

We can address nonlinear data sets using semiparametric regression techniques

Let  $x_i$  be the  $i$ th vehicle weight and  $y_i$  be its corresponding fuel efficiency score. The corresponding vectors are  $\mathbf{x}$  and  $\mathbf{y}$ .

We construct a fixed effects matrix

$$\mathbf{X} = \begin{bmatrix} 1 & x_1 \\ \vdots & \vdots \\ 1 & x_n \end{bmatrix}$$

Nonlinear effects are incorporated through a random effects matrix

$$\mathbf{Z} = \begin{bmatrix} z_1(x_1) & \dots & z_K(x_1) \\ \vdots & \ddots & \vdots \\ z_1(x_N) & \dots & z_K(x_N) \end{bmatrix},$$

where  $\{z_k\}_{k=1}^K$  is a suitable spline basis – see Wand and Ormerod (2008) for a description of canonical cubic O’Sullivan penalised splines.

# Bayesian Semiparametric Regression

A fundamental ingredient, which facilitates the incorporation of nonlinear predictor effects, is that of mixed model-based penalized splines:

$$f(x) = \beta_0 + \beta_1 x + \sum_{k=1}^K u_k z_k(x), \quad u_k | \sigma_u^2 \stackrel{\text{ind.}}{\sim} N(0, \sigma_u^2), \quad k = 1, \dots, K.$$

The Bayesian semiparametric regression model is

$$\mathbf{y} | \boldsymbol{\beta}, \sigma^2 \sim N(\mathbf{X}\boldsymbol{\beta} + \mathbf{Z}\mathbf{u}, \sigma^2 \mathbf{I}), \quad \boldsymbol{\beta} \sim N(\mathbf{0}, \sigma_\beta^2 \mathbf{I}_2)$$

$$\mathbf{u} | \sigma_u^2 \sim N(\mathbf{0}, \sigma_u^2 \mathbf{I}_K)$$

$$\sigma_u^2 | a_u \sim \text{Inverse-}\chi^2(1, 1/a_u), \quad a_u \sim \text{inverse-}\chi^2(1, 1/A^2)$$

$$\sigma^2 | a \sim \text{Inverse-}\chi^2(1, 1/a), \quad a \sim \text{inverse-}\chi^2(1, 1/A^2)$$

# Bayesian Semiparametric Regression

A fundamental ingredient, which facilitates the incorporation of nonlinear predictor effects, is that of mixed model-based penalized splines:

$$f(x) = \beta_0 + \beta_1 x + \sum_{k=1}^K u_k z_k(x), \quad u_k | \sigma_u^2 \stackrel{\text{ind.}}{\sim} N(0, \sigma_u^2), \quad k = 1, \dots, K.$$

The Bayesian semiparametric regression model is

$$\mathbf{y} | \boldsymbol{\beta}, \sigma^2 \sim N(\mathbf{X}\boldsymbol{\beta} + \mathbf{Z}\mathbf{u}, \sigma^2 \mathbf{I}), \quad \boldsymbol{\beta} \sim N(\mathbf{0}, \sigma_\beta^2 \mathbf{I}_2)$$

$$\mathbf{u} | \sigma_u^2 \sim N(\mathbf{0}, \sigma_u^2 \mathbf{I}_K)$$

$$\sigma_u^2 | a_u \sim \text{Inverse-}\chi^2(1, 1/a_u), \quad a_u \sim \text{inverse-}\chi^2(1, 1/A^2)$$

$$\sigma^2 | a \sim \text{Inverse-}\chi^2(1, 1/a), \quad a \sim \text{inverse-}\chi^2(1, 1/A^2)$$

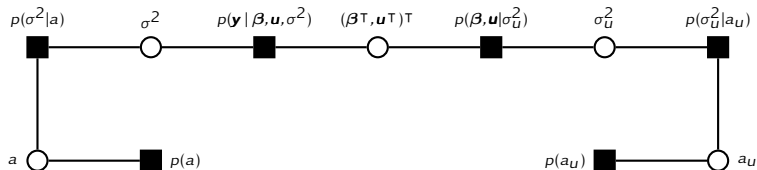


# Bayesian Semiparametric Regression

$$y|\beta, \mathbf{u}, \sigma^2 \sim N(\mathbf{X}\beta + \mathbf{Z}\mathbf{u}, \sigma^2 \mathbf{I}), \quad \begin{bmatrix} \beta \\ \mathbf{u} \end{bmatrix} \bigg| \sigma_u^2 \sim N\left(\begin{bmatrix} \mathbf{0} \\ \mathbf{0} \end{bmatrix}, \begin{bmatrix} \sigma_\beta^2 \mathbf{I}_2 & \mathbf{0}^T \\ \mathbf{0} & \sigma_u^2 \mathbf{I}_m \end{bmatrix}\right)$$

$$\sigma_u^2 | a_u \sim \text{Inverse-}\chi^2(1, 1/a_u), \quad a_u \sim \text{Inverse-}\chi^2(1, 1/A^2)$$

$$\sigma^2 | a \sim \text{Inverse-}\chi^2(1, 1/a), \quad a \sim \text{Inverse-}\chi^2(1, 1/A^2)$$



Gaussian likelihood fragment

Gaussian penalization fragment (Wand, 2017)

Iterated inverse Wishart fragment

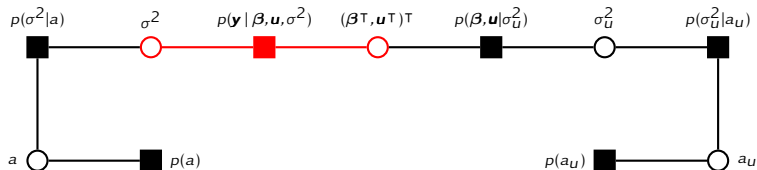
Inverse Wishart prior fragment

# Bayesian Semiparametric Regression

$$y|\beta, \mathbf{u}, \sigma^2 \sim N(\mathbf{X}\beta + \mathbf{Z}\mathbf{u}, \sigma^2 \mathbf{I}), \quad \begin{bmatrix} \beta \\ \mathbf{u} \end{bmatrix} \Big| \sigma_u^2 \sim N\left(\begin{bmatrix} \mathbf{0} \\ \mathbf{0} \end{bmatrix}, \begin{bmatrix} \sigma_\beta^2 \mathbf{I}_2 & \mathbf{0}^T \\ \mathbf{0} & \sigma_u^2 \mathbf{I}_m \end{bmatrix}\right)$$

$$\sigma_u^2 | a_u \sim \text{Inverse-}\chi^2(1, 1/a_u), \quad a_u \sim \text{Inverse-}\chi^2(1, 1/A^2)$$

$$\sigma^2 | a \sim \text{Inverse-}\chi^2(1, 1/a), \quad a \sim \text{Inverse-}\chi^2(1, 1/A^2)$$



Gaussian likelihood fragment

Gaussian penalization fragment (Wand, 2017)

Iterated inverse Wishart fragment

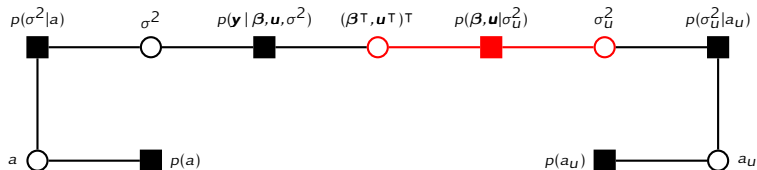
Inverse Wishart prior fragment

# Bayesian Semiparametric Regression

$$y|\beta, \mathbf{u}, \sigma^2 \sim N(\mathbf{X}\beta + \mathbf{Z}\mathbf{u}, \sigma^2 \mathbf{I}), \quad \begin{bmatrix} \beta \\ \mathbf{u} \end{bmatrix} \sim N\left(\begin{bmatrix} \mathbf{0} \\ \mathbf{0} \end{bmatrix}, \begin{bmatrix} \sigma_\beta^2 \mathbf{I}_2 & \mathbf{0}^T \\ \mathbf{0} & \sigma_u^2 \mathbf{I}_m \end{bmatrix}\right)$$

$$\sigma_u^2 | a_u \sim \text{Inverse-}\chi^2(1, 1/a_u), \quad a_u \sim \text{Inverse-}\chi^2(1, 1/A^2)$$

$$\sigma^2 | a \sim \text{Inverse-}\chi^2(1, 1/a), \quad a \sim \text{Inverse-}\chi^2(1, 1/A^2)$$



Gaussian likelihood fragment

Gaussian penalization fragment (Wand, 2017)

Iterated inverse Wishart fragment

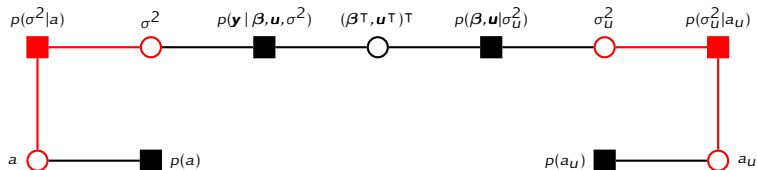
Inverse Wishart prior fragment

# Bayesian Semiparametric Regression

$$y|\beta, \mathbf{u}, \sigma^2 \sim N(\mathbf{X}\beta + \mathbf{Z}\mathbf{u}, \sigma^2 \mathbf{I}), \quad \begin{bmatrix} \beta \\ \mathbf{u} \end{bmatrix} \Big| \sigma_u^2 \sim N\left(\begin{bmatrix} \mathbf{0} \\ \mathbf{0} \end{bmatrix}, \begin{bmatrix} \sigma_\beta^2 \mathbf{I}_2 & \mathbf{0}^T \\ \mathbf{0} & \sigma_u^2 \mathbf{I}_m \end{bmatrix}\right)$$

$$\sigma_u^2 | a_u \sim \text{Inverse-}\chi^2(1, 1/a_u), \quad a_u \sim \text{Inverse-}\chi^2(1, 1/A^2)$$

$$\sigma^2 | a \sim \text{Inverse-}\chi^2(1, 1/a), \quad a \sim \text{Inverse-}\chi^2(1, 1/A^2)$$



Gaussian likelihood fragment

Gaussian penalization fragment (Wand, 2017)

Iterated inverse Wishart fragment

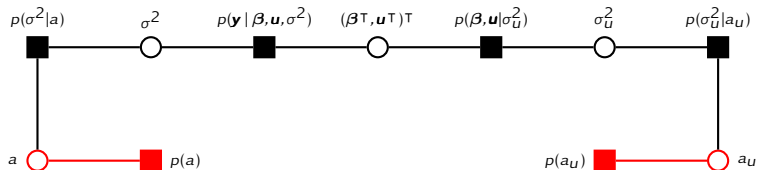
Inverse Wishart prior fragment

# Bayesian Semiparametric Regression

$$y|\beta, \mathbf{u}, \sigma^2 \sim N(\mathbf{X}\beta + \mathbf{Z}\mathbf{u}, \sigma^2 \mathbf{I}), \quad \begin{bmatrix} \beta \\ \mathbf{u} \end{bmatrix} \Big| \sigma_u^2 \sim N\left(\begin{bmatrix} \mathbf{0} \\ \mathbf{0} \end{bmatrix}, \begin{bmatrix} \sigma_\beta^2 \mathbf{I}_2 & \mathbf{0}^T \\ \mathbf{0} & \sigma_u^2 \mathbf{I}_m \end{bmatrix}\right)$$

$$\sigma_u^2 | a_u \sim \text{Inverse-}\chi^2(1, 1/a_u), \quad a_u \sim \text{Inverse-}\chi^2(1, 1/A^2)$$

$$\sigma^2 | a \sim \text{Inverse-}\chi^2(1, 1/a), \quad a \sim \text{Inverse-}\chi^2(1, 1/A^2)$$



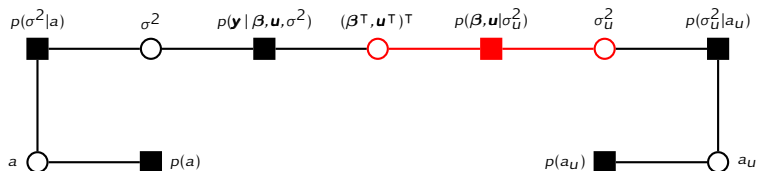
Gaussian likelihood fragment

Gaussian penalization fragment (Wand, 2017)

Iterated inverse Wishart fragment

Inverse Wishart prior fragment

# Gaussian Penalization Fragment



$$\begin{bmatrix} \beta \\ \mathbf{u} \end{bmatrix} \Big| \Sigma_U \sim \mathcal{N} \left( \begin{bmatrix} 0 \\ 0 \end{bmatrix}, \begin{bmatrix} \sigma_\beta^2 \mathbf{I}_2 & \mathbf{O}^T \\ \mathbf{O} & \Sigma_U \end{bmatrix} \right)$$

Inputs:

$$\eta_{q^*}(\beta, \mathbf{u}), \quad \eta_{q^*}(\Sigma_U)$$

Updates:

$$\text{Cov}_q\{(\beta^\top, \mathbf{u}^\top)^\top\} = -\frac{1}{2} \{\text{vec}^{-1}(\eta_{q^*}(\beta, \mathbf{u}))_2\}^{-1}, \quad \mathbb{E}_q\{(\beta^\top, \mathbf{u}^\top)^\top\} = \text{Cov}_q\{(\beta^\top, \mathbf{u}^\top)^\top\}(\eta_{q^*}(\beta, \mathbf{u}))_1,$$

$$\mathbb{E}_q(\Sigma_U^{-1}) = \{(\eta_{q^*}(\Sigma_U))_1 + 1\} \left\{ \text{vec}^{-1}(\eta_{q^*}(\sigma_U^2))_2 \right\}^{-1},$$

$$\eta_{p(\beta, \mathbf{u} | \Sigma_U) \rightarrow (\beta, \mathbf{u})} = \begin{bmatrix} \mathbf{0}_{K+2} \\ -\frac{1}{2} \text{vec} \left( \begin{bmatrix} (1/\sigma_\beta^2) \mathbf{I}_2 & \mathbf{O} \\ \mathbf{O} & \mathbb{E}_q(\Sigma_U^{-1}) \end{bmatrix} \right) \end{bmatrix}, \quad \eta_{p(\beta, \mathbf{u} | \Sigma_U) \rightarrow \Sigma_U} = \begin{bmatrix} -K/2 \\ -\frac{1}{2} \text{vec}(\mathbb{E}_q(\mathbf{u} \mathbf{u}^\top)) \end{bmatrix}$$

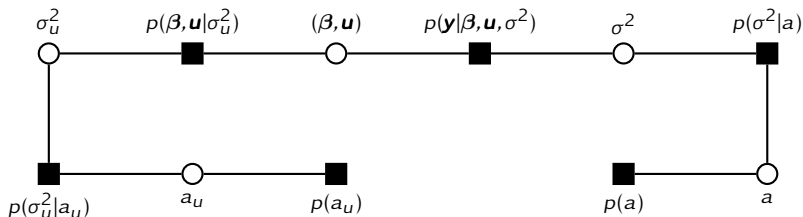
Outputs:

$$\eta_{p(\beta, \mathbf{u} | \Sigma_U) \rightarrow (\beta, \mathbf{u})}, \quad \eta_{p(\beta, \mathbf{u} | \Sigma_U) \rightarrow \Sigma_U}$$

# Bayesian Semiparametric Regression

## VMP for the Bayesian semiparametric regression model

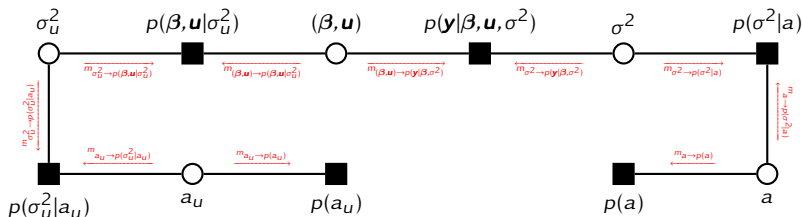
1. Initialise all messages from stochastic nodes to factors
2. Cycle:
  - (i) Update all messages from factors to stochastic nodes
  - (ii) Update all messages from stochastic nodes to factors
  - (iii) Update all optimal posterior density functions
3. Stop:  $D_{KL} \left\{ q(\beta, \mathbf{u}, \sigma^2, \sigma_u^2, a, a_u) \parallel p(\beta, \mathbf{u}, \sigma^2, \sigma_u^2, a, a_u | \mathbf{y}) \right\}$  converges.



# Bayesian Semiparametric Regression

## VMP for the Bayesian semiparametric regression model

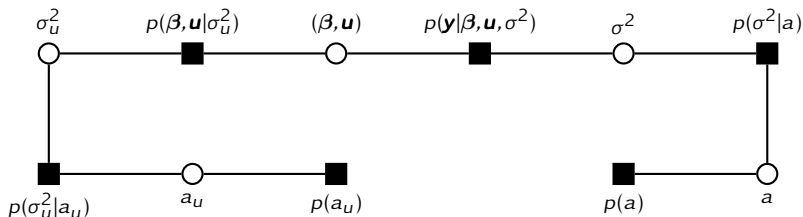
1. Initialise all messages from stochastic nodes to factors
2. Cycle:
  - (i) Update all messages from factors to stochastic nodes
  - (ii) Update all messages from stochastic nodes to factors
  - (iii) Update all optimal posterior density functions
3. Stop:  $D_{KL} \left\{ q(\beta, \mathbf{u}, \sigma^2, \sigma_u^2, a, a_u) \| p(\beta, \mathbf{u}, \sigma^2, \sigma_u^2, a, a_u | \mathbf{y}) \right\}$  converges.





## VMP for the Bayesian semiparametric regression model

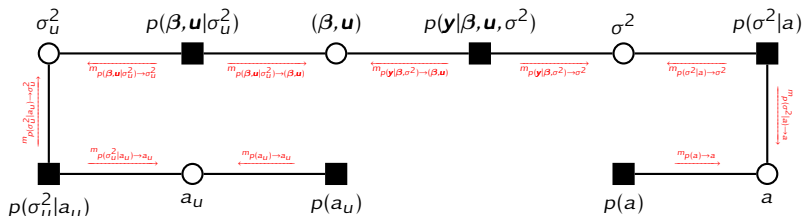
1. Initialise all messages from stochastic nodes to factors
2. Cycle:
  - (i) Update all messages from factors to stochastic nodes
  - (ii) Update all messages from stochastic nodes to factors
  - (iii) Update all optimal posterior density functions
3. Stop:  $D_{KL} \left\{ q(\beta, \mathbf{u}, \sigma^2, \sigma_u^2, a, a_u) \| p(\beta, \mathbf{u}, \sigma^2, \sigma_u^2, a, a_u | \mathbf{y}) \right\}$  converges.



# Bayesian Semiparametric Regression

## VMP for the Bayesian semiparametric regression model

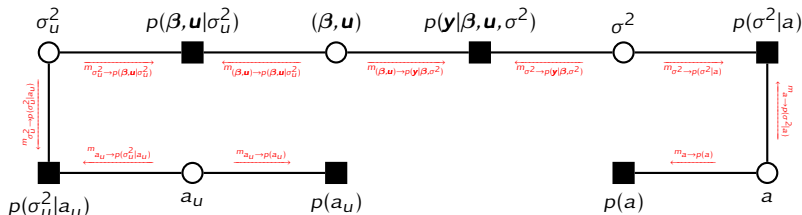
1. Initialise all messages from stochastic nodes to factors
2. Cycle:
  - (i) Update all messages from factors to stochastic nodes
  - (ii) Update all messages from stochastic nodes to factors
  - (iii) Update all optimal posterior density functions
3. Stop:  $D_{KL} \{q(\beta, \mathbf{u}, \sigma^2, \sigma_u^2, a, a_u) \| p(\beta, \mathbf{u}, \sigma^2, \sigma_u^2, a, a_u | \mathbf{y})\}$  converges.



# Bayesian Semiparametric Regression

## VMP for the Bayesian semiparametric regression model

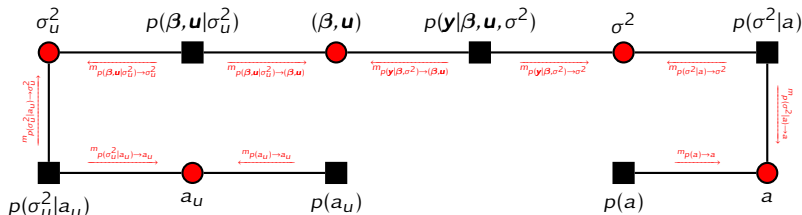
1. Initialise all messages from stochastic nodes to factors
2. Cycle:
  - (i) Update all messages from factors to stochastic nodes
  - (ii) Update all messages from stochastic nodes to factors
  - (iii) Update all optimal posterior density functions
3. Stop:  $D_{KL} \{ q(\beta, \mathbf{u}, \sigma^2, \sigma_u^2, a, a_u) \| p(\beta, \mathbf{u}, \sigma^2, \sigma_u^2, a, a_u | \mathbf{y}) \}$  converges.



# Bayesian Semiparametric Regression

## VMP for the Bayesian semiparametric regression model

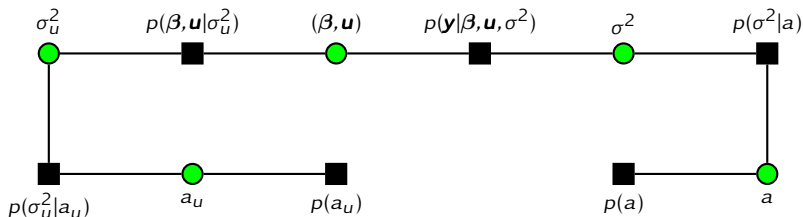
1. Initialise all messages from stochastic nodes to factors
2. Cycle:
  - (i) Update all messages from factors to stochastic nodes
  - (ii) Update all messages from stochastic nodes to factors
  - (iii) Update all optimal posterior density functions
3. Stop:  $D_{KL} \{q(\beta, \mathbf{u}, \sigma^2, \sigma_u^2, a, a_u) \| p(\beta, \mathbf{u}, \sigma^2, \sigma_u^2, a, a_u | \mathbf{y})\}$  converges.

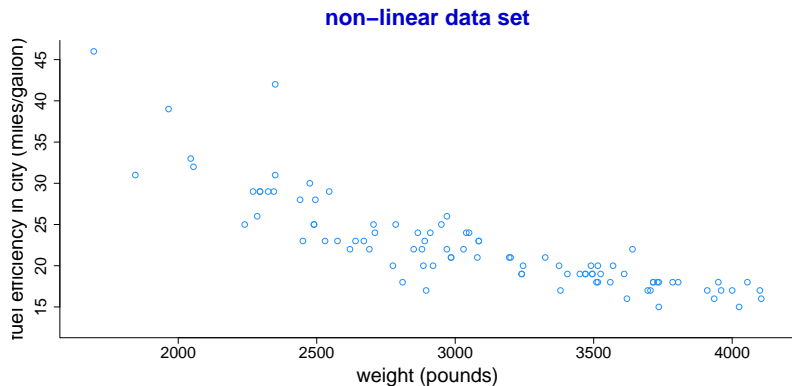


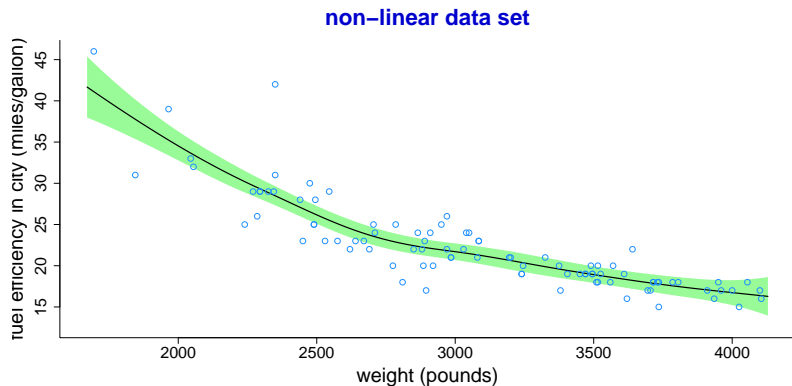
# Bayesian Semiparametric Regression

## VMP for the Bayesian semiparametric regression model

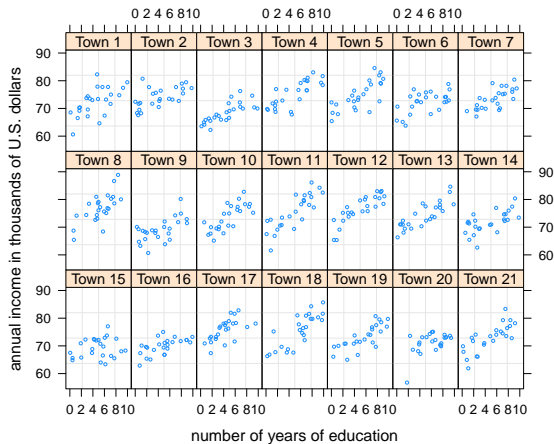
1. Initialise all messages from stochastic nodes to factors
2. Cycle:
  - (i) Update all messages from factors to stochastic nodes
  - (ii) Update all messages from stochastic nodes to factors
  - (iii) Update all optimal posterior density functions
3. Stop:  $D_{\text{KL}} \left\{ q(\beta, \mathbf{u}, \sigma^2, \sigma_u^2, a, a_u) \parallel p(\beta, \mathbf{u}, \sigma^2, \sigma_u^2, a, a_u | \mathbf{y}) \right\}$  converges.







# Multilevel Data





Suppose we have  $m$  groups and  $n_i$  subjects in the  $i$ th group.

Let  $x_{ij}$  be the predictor for the  $j$ th subject in the  $i$ th, and let  $y_{ij}$  be the corresponding observation.

A Gaussian response linear mixed model for such data consists of

$$y_{ij} = \beta_0 + \beta_1 x_{ij} + u_{i1} + u_{i2} x_{ij} + \epsilon_{ij}, \quad \begin{bmatrix} u_{i1} \\ u_{i2} \end{bmatrix} \stackrel{\text{ind.}}{\sim} N(\mathbf{0}_2, \Sigma_u),$$
$$\epsilon_{ij} \stackrel{\text{ind.}}{\sim} N(0, \sigma^2), \quad j = 1, \dots, n_i, \quad i = 1, \dots, m.$$

Next, set

$$\mathbf{y}_i \equiv \begin{bmatrix} y_{i1} \\ \vdots \\ y_{in_i} \end{bmatrix}, \quad \mathbf{x}_i \equiv \begin{bmatrix} x_{i1} \\ \vdots \\ x_{in_i} \end{bmatrix}, \quad \mathbf{X}_i \equiv \begin{bmatrix} \mathbf{1}_{n_i} & \mathbf{x}_i \end{bmatrix} \quad \text{and} \quad \mathbf{u}_i \equiv \begin{bmatrix} u_{i1} \\ u_{i2} \end{bmatrix}.$$

Suppose we have  $m$  groups and  $n_i$  subjects in the  $i$ th group.

Let  $x_{ij}$  be the predictor for the  $j$ th subject in the  $i$ th, and let  $y_{ij}$  be the corresponding observation.

A Gaussian response linear mixed model for such data consists of

$$y_{ij} = \beta_0 + \beta_1 x_{ij} + u_{i1} + u_{i2} x_{ij} + \epsilon_{ij}, \quad \begin{bmatrix} u_{i1} \\ u_{i2} \end{bmatrix} \stackrel{\text{ind.}}{\sim} N(\mathbf{0}_2, \Sigma_u),$$
$$\epsilon_{ij} \stackrel{\text{ind.}}{\sim} N(0, \sigma^2), \quad j = 1, \dots, n_i, \quad i = 1, \dots, m.$$

Next, set

$$\mathbf{y}_i \equiv \begin{bmatrix} y_{i1} \\ \vdots \\ y_{in_i} \end{bmatrix}, \quad \mathbf{x}_i \equiv \begin{bmatrix} x_{i1} \\ \vdots \\ x_{in_i} \end{bmatrix}, \quad \mathbf{X}_i \equiv \begin{bmatrix} \mathbf{1}_{n_i} & \mathbf{x}_i \end{bmatrix} \quad \text{and} \quad \mathbf{u}_i \equiv \begin{bmatrix} u_{i1} \\ u_{i2} \end{bmatrix}.$$

Suppose we have  $m$  groups and  $n_i$  subjects in the  $i$ th group.

Let  $x_{ij}$  be the predictor for the  $j$ th subject in the  $i$ th, and let  $y_{ij}$  be the corresponding observation.

A Gaussian response linear mixed model for such data consists of

$$y_{ij} = \beta_0 + \beta_1 x_{ij} + u_{i1} + u_{i2} x_{ij} + \epsilon_{ij}, \quad \begin{bmatrix} u_{i1} \\ u_{i2} \end{bmatrix} \stackrel{\text{ind.}}{\sim} N(\mathbf{0}_2, \Sigma_u),$$
$$\epsilon_{ij} \stackrel{\text{ind.}}{\sim} N(0, \sigma^2), \quad j = 1, \dots, n_i, \quad i = 1, \dots, m.$$

Next, set

$$\mathbf{y}_i \equiv \begin{bmatrix} y_{i1} \\ \vdots \\ y_{in_i} \end{bmatrix}, \quad \mathbf{x}_i \equiv \begin{bmatrix} x_{i1} \\ \vdots \\ x_{in_i} \end{bmatrix}, \quad \mathbf{X}_i \equiv \begin{bmatrix} \mathbf{1}_{n_i} & \mathbf{x}_i \end{bmatrix} \quad \text{and} \quad \mathbf{u}_i \equiv \begin{bmatrix} u_{i1} \\ u_{i2} \end{bmatrix}.$$

Suppose we have  $m$  groups and  $n_i$  subjects in the  $i$ th group.

Let  $x_{ij}$  be the predictor for the  $j$ th subject in the  $i$ th, and let  $y_{ij}$  be the corresponding observation.

A Gaussian response linear mixed model for such data consists of

$$y_{ij} = \beta_0 + \beta_1 x_{ij} + u_{i1} + u_{i2} x_{ij} + \epsilon_{ij}, \quad \begin{bmatrix} u_{i1} \\ u_{i2} \end{bmatrix} \stackrel{\text{ind.}}{\sim} N(\mathbf{0}_2, \Sigma_u),$$
$$\epsilon_{ij} \stackrel{\text{ind.}}{\sim} N(0, \sigma^2), \quad j = 1, \dots, n_i, \quad i = 1, \dots, m.$$

Next, set

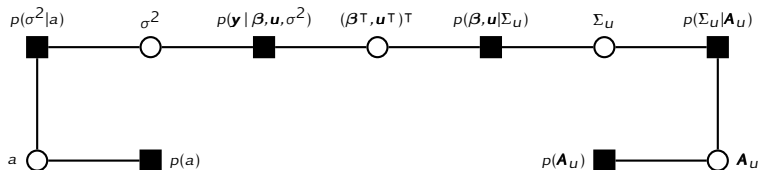
$$\mathbf{y}_i \equiv \begin{bmatrix} y_{i1} \\ \vdots \\ y_{in_i} \end{bmatrix}, \quad \mathbf{x}_i \equiv \begin{bmatrix} x_{i1} \\ \vdots \\ x_{in_i} \end{bmatrix}, \quad \mathbf{X}_i \equiv \begin{bmatrix} \mathbf{1}_{n_i} & \mathbf{x}_i \end{bmatrix} \quad \text{and} \quad \mathbf{u}_i \equiv \begin{bmatrix} u_{i1} \\ u_{i2} \end{bmatrix}.$$

# Bayesian Multilevel Data Model

$$\mathbf{y}_i | \boldsymbol{\beta}, \mathbf{u}_i, \sigma^2 \sim \mathcal{N}(\mathbf{X}_i \boldsymbol{\beta} + \mathbf{X}_i \mathbf{u}_i, \sigma^2 \mathbf{I}), \quad \begin{bmatrix} \boldsymbol{\beta} \\ \mathbf{u}_i \end{bmatrix} \Big| \Sigma_u \stackrel{\text{ind.}}{\sim} \mathcal{N} \left( \begin{bmatrix} \mathbf{0} \\ \mathbf{0} \end{bmatrix}, \begin{bmatrix} \sigma_\beta^2 \mathbf{I}_2 & \mathbf{O}^T \\ \mathbf{O} & \Sigma_u \end{bmatrix} \right)$$

$$\Sigma_u | \mathbf{A}_u \sim \text{Inverse Wishart}(1, \mathbf{A}_u^{-1}), \quad \mathbf{A}_u \sim \text{Inverse Wishart}\left(1, \frac{1}{A^2} \mathbf{I}\right)$$

$$\sigma^2 | a \sim \text{Inverse-}\chi^2(1, 1/a), \quad a \sim \text{Inverse-}\chi^2(1, 1/A^2)$$



Gaussian likelihood fragment

Gaussian penalization fragment

Iterated inverse Wishart fragment

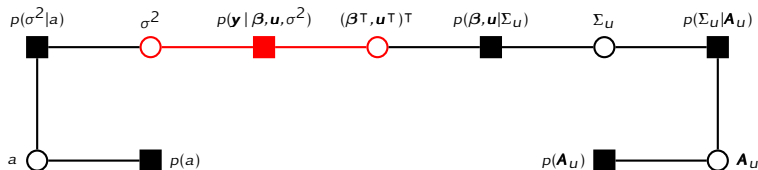
Inverse Wishart prior fragment

# Bayesian Multilevel Data Model

$$\mathbf{y}_i | \boldsymbol{\beta}, \mathbf{u}_i, \sigma^2 \sim \mathcal{N}(\mathbf{X}_i \boldsymbol{\beta} + \mathbf{X}_i \mathbf{u}_i, \sigma^2 \mathbf{I}), \quad \begin{bmatrix} \boldsymbol{\beta} \\ \mathbf{u}_i \end{bmatrix} \Big| \Sigma_u \stackrel{\text{ind.}}{\sim} \mathcal{N} \left( \begin{bmatrix} \mathbf{0} \\ \mathbf{0} \end{bmatrix}, \begin{bmatrix} \sigma_\beta^2 \mathbf{I}_2 & \mathbf{O}^T \\ \mathbf{O} & \Sigma_u \end{bmatrix} \right)$$

$$\Sigma_u | \mathbf{A}_u \sim \text{Inverse Wishart}(1, \mathbf{A}_u^{-1}), \quad \mathbf{A}_u \sim \text{Inverse Wishart}\left(1, \frac{1}{A^2} \mathbf{I}\right)$$

$$\sigma^2 | a \sim \text{Inverse-}\chi^2(1, 1/a), \quad a \sim \text{Inverse-}\chi^2(1, 1/A^2)$$



Gaussian likelihood fragment

Gaussian penalization fragment

Iterated inverse Wishart fragment

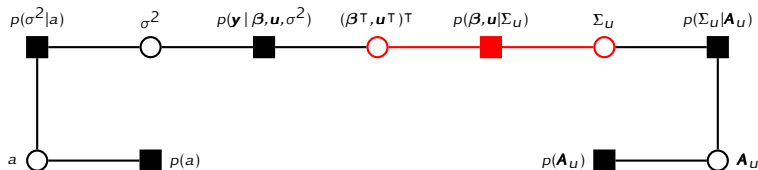
Inverse Wishart prior fragment

# Bayesian Multilevel Data Model

$$\mathbf{y}_i | \boldsymbol{\beta}, \mathbf{u}_i, \sigma^2 \sim \mathcal{N}(\mathbf{X}_i \boldsymbol{\beta} + \mathbf{X}_i \mathbf{u}_i, \sigma^2 \mathbf{I}), \quad \begin{bmatrix} \boldsymbol{\beta} \\ \mathbf{u}_i \end{bmatrix} \left| \Sigma_u \right. \stackrel{\text{ind.}}{\sim} \mathcal{N} \left( \begin{bmatrix} \mathbf{0} \\ \mathbf{0} \end{bmatrix}, \begin{bmatrix} \sigma_\beta^2 \mathbf{I}_2 & \mathbf{O}^T \\ \mathbf{O} & \Sigma_u \end{bmatrix} \right)$$

$$\Sigma_u | \mathbf{A}_u \sim \text{Inverse Wishart}(1, \mathbf{A}_u^{-1}), \quad \mathbf{A}_u \sim \text{Inverse Wishart}\left(1, \frac{1}{A^2} \mathbf{I}\right)$$

$$\sigma^2 | a \sim \text{Inverse-}\chi^2(1, 1/a), \quad a \sim \text{Inverse-}\chi^2(1, 1/A^2)$$



Gaussian likelihood fragment

Gaussian penalization fragment

Iterated inverse Wishart fragment

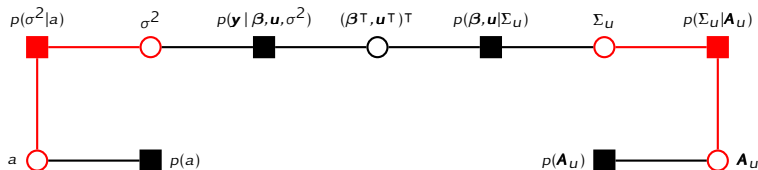
Inverse Wishart prior fragment

# Bayesian Multilevel Data Model

$$\mathbf{y}_i | \boldsymbol{\beta}, \mathbf{u}_i, \sigma^2 \sim \mathcal{N}(\mathbf{X}_i \boldsymbol{\beta} + \mathbf{X}_i \mathbf{u}_i, \sigma^2 \mathbf{I}), \quad \begin{bmatrix} \boldsymbol{\beta} \\ \mathbf{u}_i \end{bmatrix} \Big| \Sigma_u \stackrel{\text{ind.}}{\sim} \mathcal{N} \left( \begin{bmatrix} \mathbf{0} \\ \mathbf{0} \end{bmatrix}, \begin{bmatrix} \sigma_\beta^2 \mathbf{I}_2 & \mathbf{O}^T \\ \mathbf{O} & \Sigma_u \end{bmatrix} \right)$$

$$\Sigma_u | \mathbf{A}_u \sim \text{Inverse Wishart}(1, \mathbf{A}_u^{-1}), \quad \mathbf{A}_u \sim \text{Inverse Wishart}\left(1, \frac{1}{A^2} \mathbf{I}\right)$$

$$\sigma^2 | a \sim \text{Inverse-}\chi^2(1, 1/a), \quad a \sim \text{Inverse-}\chi^2(1, 1/A^2)$$



Gaussian likelihood fragment

Gaussian penalization fragment

Iterated inverse Wishart fragment

Inverse Wishart prior fragment

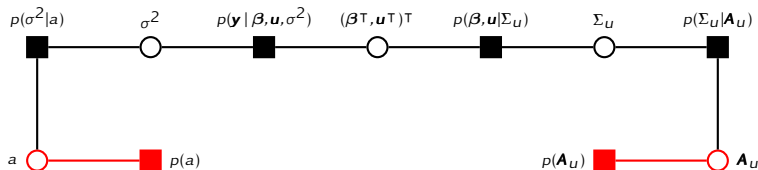


# Bayesian Multilevel Data Model

$$\mathbf{y}_i | \boldsymbol{\beta}, \mathbf{u}_i, \sigma^2 \sim \mathcal{N}(\mathbf{X}_i \boldsymbol{\beta} + \mathbf{X}_i \mathbf{u}_i, \sigma^2 \mathbf{I}), \quad \begin{bmatrix} \boldsymbol{\beta} \\ \mathbf{u}_i \end{bmatrix} \bigg| \Sigma_u \stackrel{\text{ind.}}{\sim} \mathcal{N} \left( \begin{bmatrix} \mathbf{0} \\ \mathbf{0} \end{bmatrix}, \begin{bmatrix} \sigma_\beta^2 \mathbf{I}_2 & \mathbf{O}^T \\ \mathbf{O} & \Sigma_u \end{bmatrix} \right)$$

$$\Sigma_u | \mathbf{A}_u \sim \text{Inverse Wishart}(1, \mathbf{A}_u^{-1}), \quad \mathbf{A}_u \sim \text{Inverse Wishart}\left(1, \frac{1}{A^2} \mathbf{I}\right)$$

$$\sigma^2 | a \sim \text{Inverse-}\chi^2(1, 1/a), \quad a \sim \text{Inverse-}\chi^2(1, 1/A^2)$$



Gaussian likelihood fragment

Gaussian penalization fragment

Iterated inverse Wishart fragment

Inverse Wishart prior fragment

The parameters for the posterior of  $\boldsymbol{\nu} \equiv (\boldsymbol{\beta}^\top, \boldsymbol{u}^\top)^\top$  take the form - Lee and Wand (2016):

$$\mathbb{E}_q(\boldsymbol{\nu}) \equiv \mathbf{A}^{-1} \mathbf{a}, \quad \text{Cov}_q(\boldsymbol{\nu}) \equiv \mathbf{A}^{-1}$$

where  $\mathbf{A}$  and  $\mathbf{a}$  are known

The posterior of  $\sigma^2$  also depends on norms and determinants involving  $\mathbf{A}$ .

What is the issue with  $\mathbf{A}$ ?

Consider our fictitious model of income against number of years of education for 21 towns.

The parameters for the posterior of  $\boldsymbol{\nu} \equiv (\boldsymbol{\beta}^\top, \boldsymbol{u}^\top)^\top$  take the form - Lee and Wand (2016):

$$\mathbb{E}_q(\boldsymbol{\nu}) \equiv \mathbf{A}^{-1} \mathbf{a}, \quad \text{Cov}_q(\boldsymbol{\nu}) \equiv \mathbf{A}^{-1}$$

where  $\mathbf{A}$  and  $\mathbf{a}$  are known

The posterior of  $\sigma^2$  also depends on norms and determinants involving  $\mathbf{A}$ .

What is the issue with  $\mathbf{A}$ ?

Consider our fictitious model of income against number of years of education for 21 towns.

The parameters for the posterior of  $\boldsymbol{\nu} \equiv (\boldsymbol{\beta}^\top, \boldsymbol{u}^\top)^\top$  take the form - Lee and Wand (2016):

$$\mathbb{E}_q(\boldsymbol{\nu}) \equiv \mathbf{A}^{-1} \mathbf{a}, \quad \text{Cov}_q(\boldsymbol{\nu}) \equiv \mathbf{A}^{-1}$$

where  $\mathbf{A}$  and  $\mathbf{a}$  are known

The posterior of  $\sigma^2$  also depends on norms and determinants involving  $\mathbf{A}$ .

What is the issue with  $\mathbf{A}$ ?

Consider our fictitious model of income against number of years of education for 21 towns.

The parameters for the posterior of  $\boldsymbol{\nu} \equiv (\boldsymbol{\beta}^\top, \boldsymbol{u}^\top)^\top$  take the form - Lee and Wand (2016):

$$\mathbb{E}_q(\boldsymbol{\nu}) \equiv \mathbf{A}^{-1} \mathbf{a}, \quad \text{Cov}_q(\boldsymbol{\nu}) \equiv \mathbf{A}^{-1}$$

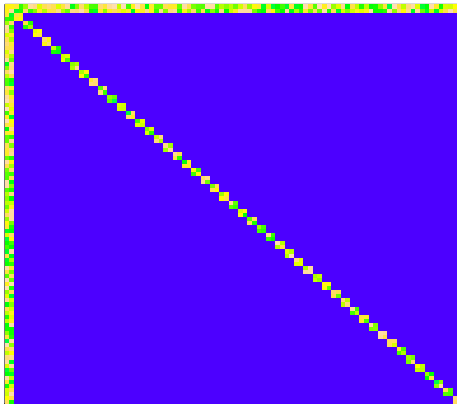
where  $\mathbf{A}$  and  $\mathbf{a}$  are known

The posterior of  $\sigma^2$  also depends on norms and determinants involving  $\mathbf{A}$ .

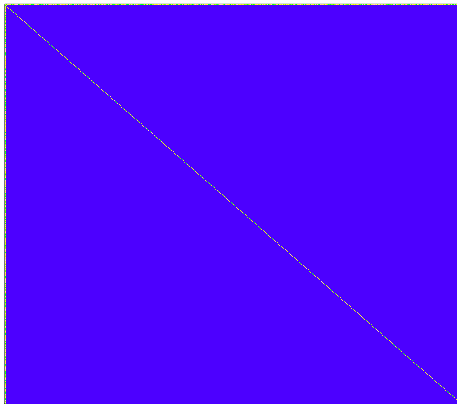
What is the issue with  $\mathbf{A}$ ?

Consider our fictitious model of income against number of years of education for 21 towns.

$A =$



$A =$



The solution comes from a simple observation:

$$\text{Cov}_q(\boldsymbol{\nu}) = \mathbf{A}^{-1} \implies \mathbf{A} \text{Cov}_q(\boldsymbol{\nu}) = \mathbf{I}$$

$$\mathbb{E}_q(\boldsymbol{\nu}) = \mathbf{A}^{-1} \mathbf{a} \implies \mathbf{A} \mathbb{E}_q(\boldsymbol{\nu}) = \mathbf{a}$$

where  $\mathbf{A}$  and  $\mathbf{a}$  are known.

Matrices and vectors are partitioned as:

$$\mathbf{A} \equiv \begin{bmatrix} \mathbf{A}_{11} & \mathbf{A}_{12,1} & \mathbf{A}_{12,2} & \mathbf{A}_{12,3} \\ \mathbf{A}_{12,1}^T & \mathbf{A}_{22,1} & \mathbf{0} & \mathbf{0} \\ \mathbf{A}_{12,2}^T & \mathbf{0} & \mathbf{A}_{22,2} & \mathbf{0} \\ \mathbf{A}_{12,3}^T & \mathbf{0} & \mathbf{0} & \mathbf{A}_{22,3} \end{bmatrix} \quad \text{Cov}_q(\boldsymbol{\nu}) \equiv \begin{bmatrix} \text{Cov}_q(\beta) & \text{Cov}_q(\beta, \mathbf{u}_1) & \text{Cov}_q(\beta, \mathbf{u}_2) & \text{Cov}_q(\beta, \mathbf{u}_3) \\ \text{Cov}_q(\beta, \mathbf{u}_1)^T & \text{Cov}_q(\mathbf{u}_1) & \times & \times \\ \text{Cov}_q(\beta, \mathbf{u}_2)^T & \times & \text{Cov}_q(\mathbf{u}_2) & \times \\ \text{Cov}_q(\beta, \mathbf{u}_3)^T & \times & \times & \text{Cov}_q(\mathbf{u}_3) \end{bmatrix}$$

$$\mathbb{E}_q(\boldsymbol{\nu}) \equiv \begin{bmatrix} \mathbb{E}_q(\beta) \\ \mathbb{E}_q(\mathbf{u}_1) \\ \mathbb{E}_q(\mathbf{u}_2) \\ \mathbb{E}_q(\mathbf{u}_3) \end{bmatrix} \quad \mathbf{a} \equiv \begin{bmatrix} \mathbf{a}_1 \\ \mathbf{a}_{2,1} \\ \mathbf{a}_{2,2} \\ \mathbf{a}_{2,3} \end{bmatrix}$$



The solution comes from a simple observation:

$$\text{Cov}_q(\boldsymbol{\nu}) = \mathbf{A}^{-1} \implies \mathbf{A} \text{Cov}_q(\boldsymbol{\nu}) = \mathbf{I}$$

$$\mathbb{E}_q(\boldsymbol{\nu}) = \mathbf{A}^{-1} \mathbf{a} \implies \mathbf{A} \mathbb{E}_q(\boldsymbol{\nu}) = \mathbf{a}$$

where  $\mathbf{A}$  and  $\mathbf{a}$  are known.

Matrices and vectors are partitioned as:

$$\mathbf{A} \equiv \begin{bmatrix} \mathbf{A}_{11} & \mathbf{A}_{12,1} & \mathbf{A}_{12,2} & \mathbf{A}_{12,3} \\ \mathbf{A}_{12,1}^T & \mathbf{A}_{22,1} & \mathbf{0} & \mathbf{0} \\ \mathbf{A}_{12,2}^T & \mathbf{0} & \mathbf{A}_{22,2} & \mathbf{0} \\ \mathbf{A}_{12,3}^T & \mathbf{0} & \mathbf{0} & \mathbf{A}_{22,3} \end{bmatrix} \quad \text{Cov}_q(\boldsymbol{\nu}) \equiv \begin{bmatrix} \text{Cov}_q(\beta) & \text{Cov}_q(\beta, \mathbf{u}_1) & \text{Cov}_q(\beta, \mathbf{u}_2) & \text{Cov}_q(\beta, \mathbf{u}_3) \\ \text{Cov}_q(\beta, \mathbf{u}_1)^T & \text{Cov}_q(\mathbf{u}_1) & \times & \times \\ \text{Cov}_q(\beta, \mathbf{u}_2)^T & \times & \text{Cov}_q(\mathbf{u}_2) & \times \\ \text{Cov}_q(\beta, \mathbf{u}_3)^T & \times & \times & \text{Cov}_q(\mathbf{u}_3) \end{bmatrix}$$

$$\mathbb{E}_q(\boldsymbol{\nu}) \equiv \begin{bmatrix} \mathbb{E}_q(\beta) \\ \mathbb{E}_q(\mathbf{u}_1) \\ \mathbb{E}_q(\mathbf{u}_2) \\ \mathbb{E}_q(\mathbf{u}_3) \end{bmatrix} \quad \mathbf{a} \equiv \begin{bmatrix} \mathbf{a}_1 \\ \mathbf{a}_{2,1} \\ \mathbf{a}_{2,2} \\ \mathbf{a}_{2,3} \end{bmatrix}$$

The solution comes from a simple observation:

$$\text{Cov}_q(\boldsymbol{\nu}) = \mathbf{A}^{-1} \Rightarrow \mathbf{A} \text{Cov}_q(\boldsymbol{\nu}) = \mathbf{I}$$

$$\mathbb{E}_q(\boldsymbol{\nu}) = \mathbf{A}^{-1} \mathbf{a} \Rightarrow \mathbf{A} \mathbb{E}_q(\boldsymbol{\nu}) = \mathbf{a}$$

where  $\mathbf{A}$  and  $\mathbf{a}$  are known.

Matrices and vectors are partitioned as:

$$\mathbf{A} \equiv \begin{bmatrix} \mathbf{A}_{11} & \mathbf{A}_{12,1} & \mathbf{A}_{12,2} & \mathbf{A}_{12,3} \\ \mathbf{A}_{12,1}^T & \mathbf{A}_{22,1} & \mathbf{0} & \mathbf{0} \\ \mathbf{A}_{12,2}^T & \mathbf{0} & \mathbf{A}_{22,2} & \mathbf{0} \\ \mathbf{A}_{12,3}^T & \mathbf{0} & \mathbf{0} & \mathbf{A}_{22,3} \end{bmatrix} \quad \text{Cov}_q(\boldsymbol{\nu}) \equiv \begin{bmatrix} \text{Cov}_q(\beta) & \text{Cov}_q(\beta, \mathbf{u}_1) & \text{Cov}_q(\beta, \mathbf{u}_2) & \text{Cov}_q(\beta, \mathbf{u}_3) \\ \text{Cov}_q(\beta, \mathbf{u}_1)^T & \text{Cov}_q(\mathbf{u}_1) & \times & \times \\ \text{Cov}_q(\beta, \mathbf{u}_2)^T & \times & \text{Cov}_q(\mathbf{u}_2) & \times \\ \text{Cov}_q(\beta, \mathbf{u}_3)^T & \times & \times & \text{Cov}_q(\mathbf{u}_3) \end{bmatrix}$$

$$\mathbb{E}_q(\boldsymbol{\nu}) \equiv \begin{bmatrix} \mathbb{E}_q(\beta) \\ \mathbb{E}_q(\mathbf{u}_1) \\ \mathbb{E}_q(\mathbf{u}_2) \\ \mathbb{E}_q(\mathbf{u}_3) \end{bmatrix} \quad \mathbf{a} \equiv \begin{bmatrix} \mathbf{a}_1 \\ \mathbf{a}_{2,1} \\ \mathbf{a}_{2,2} \\ \mathbf{a}_{2,3} \end{bmatrix}$$

The solution comes from a simple observation:

$$\text{Cov}_q(\boldsymbol{\nu}) = \mathbf{A}^{-1} \Rightarrow \mathbf{A} \text{Cov}_q(\boldsymbol{\nu}) = \mathbf{I}$$

$$\mathbb{E}_q(\boldsymbol{\nu}) = \mathbf{A}^{-1} \mathbf{a} \Rightarrow \mathbf{A} \mathbb{E}_q(\boldsymbol{\nu}) = \mathbf{a}$$

where  $\mathbf{A}$  and  $\mathbf{a}$  are known.

Matrices and vectors are partitioned as:

$$\mathbf{A} \equiv \begin{bmatrix} \mathbf{A}_{11} & \mathbf{A}_{12,1} & \mathbf{A}_{12,2} & \mathbf{A}_{12,3} \\ \mathbf{A}_{12,1}^T & \mathbf{A}_{22,1} & \mathbf{0} & \mathbf{0} \\ \mathbf{A}_{12,2}^T & \mathbf{0} & \mathbf{A}_{22,2} & \mathbf{0} \\ \mathbf{A}_{12,3}^T & \mathbf{0} & \mathbf{0} & \mathbf{A}_{22,3} \end{bmatrix} \quad \text{Cov}_q(\boldsymbol{\nu}) \equiv \begin{bmatrix} \text{Cov}_q(\boldsymbol{\beta}) & \text{Cov}_q(\boldsymbol{\beta}, \mathbf{u}_1) & \text{Cov}_q(\boldsymbol{\beta}, \mathbf{u}_2) & \text{Cov}_q(\boldsymbol{\beta}, \mathbf{u}_3) \\ \text{Cov}_q(\boldsymbol{\beta}, \mathbf{u}_1)^T & \text{Cov}_q(\mathbf{u}_1) & \times & \times \\ \text{Cov}_q(\boldsymbol{\beta}, \mathbf{u}_2)^T & \times & \text{Cov}_q(\mathbf{u}_2) & \times \\ \text{Cov}_q(\boldsymbol{\beta}, \mathbf{u}_3)^T & \times & \times & \text{Cov}_q(\mathbf{u}_3) \end{bmatrix}$$

$$\mathbb{E}_q(\boldsymbol{\nu}) \equiv \begin{bmatrix} \mathbb{E}_q(\boldsymbol{\beta}) \\ \mathbb{E}_q(\mathbf{u}_1) \\ \mathbb{E}_q(\mathbf{u}_2) \\ \mathbb{E}_q(\mathbf{u}_3) \end{bmatrix} \quad \mathbf{a} \equiv \begin{bmatrix} \mathbf{a}_1 \\ \mathbf{a}_{2,1} \\ \mathbf{a}_{2,2} \\ \mathbf{a}_{2,3} \end{bmatrix}$$

The solution comes from a simple observation:

$$\text{Cov}_q(\boldsymbol{\nu}) = \mathbf{A}^{-1} \implies \mathbf{A} \text{Cov}_q(\boldsymbol{\nu}) = \mathbf{I}$$

$$\mathbb{E}_q(\boldsymbol{\nu}) = \mathbf{A}^{-1} \mathbf{a} \implies \mathbf{A} \mathbb{E}_q(\boldsymbol{\nu}) = \mathbf{a}$$

where  $\mathbf{A}$  and  $\mathbf{a}$  are known.

Matrices and vectors are partitioned as:

$$\mathbf{A} \equiv \begin{bmatrix} \mathbf{A}_{11} & \mathbf{A}_{12,1} & \mathbf{A}_{12,2} & \mathbf{A}_{12,3} \\ \mathbf{A}_{12,1}^T & \mathbf{A}_{22,1} & \mathbf{0} & \mathbf{0} \\ \mathbf{A}_{12,2}^T & \mathbf{0} & \mathbf{A}_{22,2} & \mathbf{0} \\ \mathbf{A}_{12,3}^T & \mathbf{0} & \mathbf{0} & \mathbf{A}_{22,3} \end{bmatrix} \quad \text{Cov}_q(\boldsymbol{\nu}) \equiv \begin{bmatrix} \text{Cov}_q(\beta) & \text{Cov}_q(\beta, \mathbf{u}_1) & \text{Cov}_q(\beta, \mathbf{u}_2) & \text{Cov}_q(\beta, \mathbf{u}_3) \\ \text{Cov}_q(\beta, \mathbf{u}_1)^T & \text{Cov}_q(\mathbf{u}_1) & \times & \times \\ \text{Cov}_q(\beta, \mathbf{u}_2)^T & \times & \text{Cov}_q(\mathbf{u}_2) & \times \\ \text{Cov}_q(\beta, \mathbf{u}_3)^T & \times & \times & \text{Cov}_q(\mathbf{u}_3) \end{bmatrix}$$

$$\mathbb{E}_q(\boldsymbol{\nu}) \equiv \begin{bmatrix} \mathbb{E}_q(\beta) \\ \mathbb{E}_q(\mathbf{u}_1) \\ \mathbb{E}_q(\mathbf{u}_2) \\ \mathbb{E}_q(\mathbf{u}_3) \end{bmatrix} \quad \mathbf{a} \equiv \begin{bmatrix} \mathbf{a}_1 \\ \mathbf{a}_{2,1} \\ \mathbf{a}_{2,2} \\ \mathbf{a}_{2,3} \end{bmatrix}$$

The solution comes from a simple observation:

$$\text{Cov}_q(\boldsymbol{\nu}) = \mathbf{A}^{-1} \implies \mathbf{A} \text{Cov}_q(\boldsymbol{\nu}) = \mathbf{I}$$

$$\mathbb{E}_q(\boldsymbol{\nu}) = \mathbf{A}^{-1} \mathbf{a} \implies \mathbf{A} \mathbb{E}_q(\boldsymbol{\nu}) = \mathbf{a}$$

where  $\mathbf{A}$  and  $\mathbf{a}$  are known.

Matrices and vectors are partitioned as:

$$\mathbf{A} \equiv \begin{bmatrix} \mathbf{A}_{11} & \mathbf{A}_{12,1} & \mathbf{A}_{12,2} & \mathbf{A}_{12,3} \\ \mathbf{A}_{12,1}^T & \mathbf{A}_{22,1} & \mathbf{0} & \mathbf{0} \\ \mathbf{A}_{12,2}^T & \mathbf{0} & \mathbf{A}_{22,2} & \mathbf{0} \\ \mathbf{A}_{12,3}^T & \mathbf{0} & \mathbf{0} & \mathbf{A}_{22,3} \end{bmatrix} \quad \text{Cov}_q(\boldsymbol{\nu}) \equiv \begin{bmatrix} \text{Cov}_q(\beta) & \text{Cov}_q(\beta, \mathbf{u}_1) & \text{Cov}_q(\beta, \mathbf{u}_2) & \text{Cov}_q(\beta, \mathbf{u}_3) \\ \text{Cov}_q(\beta, \mathbf{u}_1)^T & \text{Cov}_q(\mathbf{u}_1) & \times & \times \\ \text{Cov}_q(\beta, \mathbf{u}_2)^T & \times & \text{Cov}_q(\mathbf{u}_2) & \times \\ \text{Cov}_q(\beta, \mathbf{u}_3)^T & \times & \times & \text{Cov}_q(\mathbf{u}_3) \end{bmatrix}$$

$$\mathbb{E}_q(\boldsymbol{\nu}) \equiv \begin{bmatrix} \mathbb{E}_q(\beta) \\ \mathbb{E}_q(\mathbf{u}_1) \\ \mathbb{E}_q(\mathbf{u}_2) \\ \mathbb{E}_q(\mathbf{u}_3) \end{bmatrix} \quad \mathbf{a} \equiv \begin{bmatrix} \mathbf{a}_1 \\ \mathbf{a}_{2,1} \\ \mathbf{a}_{2,2} \\ \mathbf{a}_{2,3} \end{bmatrix}$$

The solution comes from a simple observation:

$$\text{Cov}_q(\boldsymbol{\nu}) = \mathbf{A}^{-1} \Rightarrow \mathbf{A} \text{Cov}_q(\boldsymbol{\nu}) = \mathbf{I}$$

$$\mathbb{E}_q(\boldsymbol{\nu}) = \mathbf{A}^{-1} \mathbf{a} \Rightarrow \mathbf{A} \mathbb{E}_q(\boldsymbol{\nu}) = \mathbf{a}$$

where  $\mathbf{A}$  and  $\mathbf{a}$  are known.

Matrices and vectors are partitioned as:

$$\mathbf{A} \equiv \begin{bmatrix} \mathbf{A}_{11} & \mathbf{A}_{12,1} & \mathbf{A}_{12,2} & \mathbf{A}_{12,3} \\ \mathbf{A}_{12,1}^T & \mathbf{A}_{22,1} & \mathbf{0} & \mathbf{0} \\ \mathbf{A}_{12,2}^T & \mathbf{0} & \mathbf{A}_{22,2} & \mathbf{0} \\ \mathbf{A}_{12,3}^T & \mathbf{0} & \mathbf{0} & \mathbf{A}_{22,3} \end{bmatrix} \quad \text{Cov}_q(\boldsymbol{\nu}) \equiv \begin{bmatrix} \text{Cov}_q(\beta) & \text{Cov}_q(\beta, \mathbf{u}_1) & \text{Cov}_q(\beta, \mathbf{u}_2) & \text{Cov}_q(\beta, \mathbf{u}_3) \\ \text{Cov}_q(\beta, \mathbf{u}_1)^T & \text{Cov}_q(\mathbf{u}_1) & \times & \times \\ \text{Cov}_q(\beta, \mathbf{u}_2)^T & \times & \text{Cov}_q(\mathbf{u}_2) & \times \\ \text{Cov}_q(\beta, \mathbf{u}_3)^T & \times & \times & \text{Cov}_q(\mathbf{u}_3) \end{bmatrix}$$

$$\mathbb{E}_q(\boldsymbol{\nu}) \equiv \begin{bmatrix} \mathbb{E}_q(\beta) \\ \mathbb{E}_q(\mathbf{u}_1) \\ \mathbb{E}_q(\mathbf{u}_2) \\ \mathbb{E}_q(\mathbf{u}_3) \end{bmatrix} \quad \mathbf{a} \equiv \begin{bmatrix} \mathbf{a}_1 \\ \mathbf{a}_{2,1} \\ \mathbf{a}_{2,2} \\ \mathbf{a}_{2,3} \end{bmatrix}$$

# Multilevel Model ( $\mathbf{A}\text{Cov}_q(\boldsymbol{\nu}) = \mathbf{I}$ )

$$\mathbf{A}\text{Cov}_q(\boldsymbol{\nu}) = \begin{bmatrix} \mathbf{A}_{11} & \mathbf{A}_{12,1} & \mathbf{A}_{12,2} & \mathbf{A}_{12,3} \\ \mathbf{A}_{12,1}^T & \mathbf{A}_{22,1} & \mathbf{O} & \mathbf{O} \\ \mathbf{A}_{12,2}^T & \mathbf{O} & \mathbf{A}_{22,2} & \mathbf{O} \\ \mathbf{A}_{12,3}^T & \mathbf{O} & \mathbf{O} & \mathbf{A}_{22,3} \end{bmatrix} \begin{bmatrix} \text{Cov}_q(\boldsymbol{\beta}) & \text{Cov}_q(\boldsymbol{\beta}, \mathbf{u}_1) & \text{Cov}_q(\boldsymbol{\beta}, \mathbf{u}_2) & \text{Cov}_q(\boldsymbol{\beta}, \mathbf{u}_3) \\ \text{Cov}_q(\boldsymbol{\beta}, \mathbf{u}_1)^T & \text{Cov}_q(\mathbf{u}_1) & \times & \times \\ \text{Cov}_q(\boldsymbol{\beta}, \mathbf{u}_2)^T & \times & \text{Cov}_q(\mathbf{u}_2) & \times \\ \text{Cov}_q(\boldsymbol{\beta}, \mathbf{u}_3)^T & \times & \times & \text{Cov}_q(\mathbf{u}_3) \end{bmatrix}$$

$$= \begin{bmatrix} \mathbf{I} & \mathbf{O} & \mathbf{O} & \mathbf{O} \\ \mathbf{O} & \mathbf{I} & \mathbf{O} & \mathbf{O} \\ \mathbf{O} & \mathbf{O} & \mathbf{I} & \mathbf{O} \\ \mathbf{O} & \mathbf{O} & \mathbf{O} & \mathbf{I} \end{bmatrix}$$

$$\mathbf{A}_{11} \text{Cov}_q(\boldsymbol{\beta}) + \sum_{i=1}^3 \mathbf{A}_{12,i} \text{Cov}_q(\boldsymbol{\beta}, \mathbf{u}_i) = \mathbf{I}$$

$$\mathbf{A}_{12,i}^T \text{Cov}_q(\boldsymbol{\beta}, \mathbf{u}_i) + \mathbf{A}_{22,i} \text{Cov}_q(\mathbf{u}_i) = \mathbf{I}, \quad i = 1, 2, 3$$

$$\mathbf{A}_{12,i}^T \text{Cov}_q(\boldsymbol{\beta}) + \mathbf{A}_{22,i} \text{Cov}_q(\boldsymbol{\beta}, \mathbf{u}_i) = \mathbf{O}, \quad i = 1, 2, 3$$

# Multilevel Model ( $\mathbf{A}\text{Cov}_q(\boldsymbol{\nu}) = \mathbf{I}$ )

$$\mathbf{A}\text{Cov}_q(\boldsymbol{\nu}) = \begin{bmatrix} \mathbf{A}_{11} & \mathbf{A}_{12,1} & \mathbf{A}_{12,2} & \mathbf{A}_{12,3} \\ \mathbf{A}_{12,1}^T & \mathbf{A}_{22,1} & \mathbf{O} & \mathbf{O} \\ \mathbf{A}_{12,2}^T & \mathbf{O} & \mathbf{A}_{22,2} & \mathbf{O} \\ \mathbf{A}_{12,3}^T & \mathbf{O} & \mathbf{O} & \mathbf{A}_{22,3} \end{bmatrix} \begin{bmatrix} \text{Cov}_q(\boldsymbol{\beta}) & \text{Cov}_q(\boldsymbol{\beta}, \mathbf{u}_1) & \text{Cov}_q(\boldsymbol{\beta}, \mathbf{u}_2) & \text{Cov}_q(\boldsymbol{\beta}, \mathbf{u}_3) \\ \text{Cov}_q(\boldsymbol{\beta}, \mathbf{u}_1)^T & \text{Cov}_q(\mathbf{u}_1) & \times & \times \\ \text{Cov}_q(\boldsymbol{\beta}, \mathbf{u}_2)^T & \times & \text{Cov}_q(\mathbf{u}_2) & \times \\ \text{Cov}_q(\boldsymbol{\beta}, \mathbf{u}_3)^T & \times & \times & \text{Cov}_q(\mathbf{u}_3) \end{bmatrix}$$

$$= \begin{bmatrix} \mathbf{I} & \mathbf{O} & \mathbf{O} & \mathbf{O} \\ \mathbf{O} & \mathbf{I} & \mathbf{O} & \mathbf{O} \\ \mathbf{O} & \mathbf{O} & \mathbf{I} & \mathbf{O} \\ \mathbf{O} & \mathbf{O} & \mathbf{O} & \mathbf{I} \end{bmatrix}$$

$$\mathbf{A}_{11} \text{Cov}_q(\boldsymbol{\beta}) + \sum_{i=1}^3 \mathbf{A}_{12,i} \text{Cov}_q(\boldsymbol{\beta}, \mathbf{u}_i) = \mathbf{I}$$

$$\mathbf{A}_{12,i}^T \text{Cov}_q(\boldsymbol{\beta}, \mathbf{u}_i) + \mathbf{A}_{22,i} \text{Cov}_q(\mathbf{u}_i) = \mathbf{I}, \quad i = 1, 2, 3$$

$$\mathbf{A}_{12,i}^T \text{Cov}_q(\boldsymbol{\beta}) + \mathbf{A}_{22,i} \text{Cov}_q(\boldsymbol{\beta}, \mathbf{u}_i) = \mathbf{O}, \quad i = 1, 2, 3$$



# Multilevel Model ( $\mathbf{A}\text{Cov}_q(\boldsymbol{\nu}) = \mathbf{I}$ )

$$\mathbf{A}\text{Cov}_q(\boldsymbol{\nu}) = \begin{bmatrix} \mathbf{A}_{11} & \mathbf{A}_{12,1} & \mathbf{A}_{12,2} & \mathbf{A}_{12,3} \\ \mathbf{A}_{12,1}^T & \mathbf{A}_{22,1} & \mathbf{O} & \mathbf{O} \\ \mathbf{A}_{12,2}^T & \mathbf{O} & \mathbf{A}_{22,2} & \mathbf{O} \\ \mathbf{A}_{12,3}^T & \mathbf{O} & \mathbf{O} & \mathbf{A}_{22,3} \end{bmatrix} \begin{bmatrix} \text{Cov}_q(\boldsymbol{\beta}) & \text{Cov}_q(\boldsymbol{\beta}, \mathbf{u}_1) & \text{Cov}_q(\boldsymbol{\beta}, \mathbf{u}_2) & \text{Cov}_q(\boldsymbol{\beta}, \mathbf{u}_3) \\ \text{Cov}_q(\boldsymbol{\beta}, \mathbf{u}_1)^T & \text{Cov}_q(\mathbf{u}_1) & \times & \times \\ \text{Cov}_q(\boldsymbol{\beta}, \mathbf{u}_2)^T & \times & \text{Cov}_q(\mathbf{u}_2) & \times \\ \text{Cov}_q(\boldsymbol{\beta}, \mathbf{u}_3)^T & \times & \times & \text{Cov}_q(\mathbf{u}_3) \end{bmatrix}$$

$$= \begin{bmatrix} \mathbf{I} & \mathbf{O} & \mathbf{O} & \mathbf{O} \\ \mathbf{O} & \mathbf{I} & \mathbf{O} & \mathbf{O} \\ \mathbf{O} & \mathbf{O} & \mathbf{I} & \mathbf{O} \\ \mathbf{O} & \mathbf{O} & \mathbf{O} & \mathbf{I} \end{bmatrix}$$

$$\mathbf{A}_{11} \text{Cov}_q(\boldsymbol{\beta}) + \sum_{i=1}^3 \mathbf{A}_{12,i} \text{Cov}_q(\boldsymbol{\beta}, \mathbf{u}_i) = \mathbf{I}$$

$$\mathbf{A}_{12,i}^T \text{Cov}_q(\boldsymbol{\beta}, \mathbf{u}_i) + \mathbf{A}_{22,i} \text{Cov}_q(\mathbf{u}_i) = \mathbf{I}, \quad i = 1, 2, 3$$

$$\mathbf{A}_{12,i}^T \text{Cov}_q(\boldsymbol{\beta}) + \mathbf{A}_{22,i} \text{Cov}_q(\boldsymbol{\beta}, \mathbf{u}_i) = \mathbf{O}, \quad i = 1, 2, 3$$

# Multilevel Model ( $\mathbf{A}\text{Cov}_q(\boldsymbol{\nu}) = \mathbf{I}$ )

$$\mathbf{A}\text{Cov}_q(\boldsymbol{\nu}) = \begin{bmatrix} \mathbf{A}_{11} & \mathbf{A}_{12,1} & \mathbf{A}_{12,2} & \mathbf{A}_{12,3} \\ \mathbf{A}_{12,1}^T & \mathbf{A}_{22,1} & \mathbf{O} & \mathbf{O} \\ \mathbf{A}_{12,2}^T & \mathbf{O} & \mathbf{A}_{22,2} & \mathbf{O} \\ \mathbf{A}_{12,3}^T & \mathbf{O} & \mathbf{O} & \mathbf{A}_{22,3} \end{bmatrix} \begin{bmatrix} \text{Cov}_q(\boldsymbol{\beta}) & \text{Cov}_q(\boldsymbol{\beta}, \mathbf{u}_1) & \text{Cov}_q(\boldsymbol{\beta}, \mathbf{u}_2) & \text{Cov}_q(\boldsymbol{\beta}, \mathbf{u}_3) \\ \text{Cov}_q(\boldsymbol{\beta}, \mathbf{u}_1)^T & \text{Cov}_q(\mathbf{u}_1) & \times & \times \\ \text{Cov}_q(\boldsymbol{\beta}, \mathbf{u}_2)^T & \times & \text{Cov}_q(\mathbf{u}_2) & \times \\ \text{Cov}_q(\boldsymbol{\beta}, \mathbf{u}_3)^T & \times & \times & \text{Cov}_q(\mathbf{u}_3) \end{bmatrix}$$

$$= \begin{bmatrix} \mathbf{I} & \mathbf{O} & \mathbf{O} & \mathbf{O} \\ \mathbf{O} & \mathbf{I} & \mathbf{O} & \mathbf{O} \\ \mathbf{O} & \mathbf{O} & \mathbf{I} & \mathbf{O} \\ \mathbf{O} & \mathbf{O} & \mathbf{O} & \mathbf{I} \end{bmatrix}$$

$$\mathbf{A}_{11} \text{Cov}_q(\boldsymbol{\beta}) + \sum_{i=1}^3 \mathbf{A}_{12,i} \text{Cov}_q(\boldsymbol{\beta}, \mathbf{u}_i) = \mathbf{I}$$

$$\mathbf{A}_{12,i}^T \text{Cov}_q(\boldsymbol{\beta}, \mathbf{u}_i) + \mathbf{A}_{22,i} \text{Cov}_q(\mathbf{u}_i) = \mathbf{I}, \quad i = 1, 2, 3$$

$$\mathbf{A}_{12,i}^T \text{Cov}_q(\boldsymbol{\beta}) + \mathbf{A}_{22,i} \text{Cov}_q(\boldsymbol{\beta}, \mathbf{u}_i) = \mathbf{O}, \quad i = 1, 2, 3$$

# Multilevel Model ( $\mathbf{A} \text{Cov}_q(\boldsymbol{\nu}) = \mathbf{I}$ )

$$\mathbf{A} \text{Cov}_q(\boldsymbol{\nu}) = \begin{bmatrix} \mathbf{A}_{11} & \mathbf{A}_{12,1} & \mathbf{A}_{12,2} & \mathbf{A}_{12,3} \\ \mathbf{A}_{12,1}^T & \mathbf{A}_{22,1} & \mathbf{O} & \mathbf{O} \\ \mathbf{A}_{12,2}^T & \mathbf{O} & \mathbf{A}_{22,2} & \mathbf{O} \\ \mathbf{A}_{12,3}^T & \mathbf{O} & \mathbf{O} & \mathbf{A}_{22,3} \end{bmatrix} \begin{bmatrix} \text{Cov}_q(\boldsymbol{\beta}) & \text{Cov}_q(\boldsymbol{\beta}, \mathbf{u}_1) & \text{Cov}_q(\boldsymbol{\beta}, \mathbf{u}_2) & \text{Cov}_q(\boldsymbol{\beta}, \mathbf{u}_3) \\ \text{Cov}_q(\boldsymbol{\beta}, \mathbf{u}_1)^T & \text{Cov}_q(\mathbf{u}_1) & \times & \times \\ \text{Cov}_q(\boldsymbol{\beta}, \mathbf{u}_2)^T & \times & \text{Cov}_q(\mathbf{u}_2) & \times \\ \text{Cov}_q(\boldsymbol{\beta}, \mathbf{u}_3)^T & \times & \times & \text{Cov}_q(\mathbf{u}_3) \end{bmatrix}$$

$$= \begin{bmatrix} \mathbf{I} & \mathbf{O} & \mathbf{O} & \mathbf{O} \\ \mathbf{O} & \mathbf{I} & \mathbf{O} & \mathbf{O} \\ \mathbf{O} & \mathbf{O} & \mathbf{I} & \mathbf{O} \\ \mathbf{O} & \mathbf{O} & \mathbf{O} & \mathbf{I} \end{bmatrix}$$

$$\mathbf{A}_{11} \text{Cov}_q(\boldsymbol{\beta}) + \sum_{i=1}^3 \mathbf{A}_{12,i} \text{Cov}_q(\boldsymbol{\beta}, \mathbf{u}_i) = \mathbf{I}$$

$$\mathbf{A}_{12,i}^T \text{Cov}_q(\boldsymbol{\beta}, \mathbf{u}_i) + \mathbf{A}_{22,i} \text{Cov}_q(\mathbf{u}_i) = \mathbf{I}, \quad i = 1, 2, 3$$

$$\mathbf{A}_{12,i}^T \text{Cov}_q(\boldsymbol{\beta}) + \mathbf{A}_{22,i} \text{Cov}_q(\boldsymbol{\beta}, \mathbf{u}_i) = \mathbf{O}, \quad i = 1, 2, 3$$

# Multilevel Model ( $\mathbf{A}\text{Cov}_q(\boldsymbol{\nu}) = \mathbf{I}$ )

$$\mathbf{A}\text{Cov}_q(\boldsymbol{\nu}) = \begin{bmatrix} \mathbf{A}_{11} & \mathbf{A}_{12,1} & \mathbf{A}_{12,2} & \mathbf{A}_{12,3} \\ \mathbf{A}_{12,1}^T & \mathbf{A}_{22,1} & \mathbf{O} & \mathbf{O} \\ \mathbf{A}_{12,2}^T & \mathbf{O} & \mathbf{A}_{22,2} & \mathbf{O} \\ \mathbf{A}_{12,3}^T & \mathbf{O} & \mathbf{O} & \mathbf{A}_{22,3} \end{bmatrix} \begin{bmatrix} \text{Cov}_q(\boldsymbol{\beta}) & \text{Cov}_q(\boldsymbol{\beta}, \mathbf{u}_1) & \text{Cov}_q(\boldsymbol{\beta}, \mathbf{u}_2) & \text{Cov}_q(\boldsymbol{\beta}, \mathbf{u}_3) \\ \text{Cov}_q(\boldsymbol{\beta}, \mathbf{u}_1)^T & \text{Cov}_q(\mathbf{u}_1) & \times & \times \\ \text{Cov}_q(\boldsymbol{\beta}, \mathbf{u}_2)^T & \times & \text{Cov}_q(\mathbf{u}_2) & \times \\ \text{Cov}_q(\boldsymbol{\beta}, \mathbf{u}_3)^T & \times & \times & \text{Cov}_q(\mathbf{u}_3) \end{bmatrix}$$

$$= \begin{bmatrix} \mathbf{I} & \mathbf{O} & \mathbf{O} & \mathbf{O} \\ \mathbf{O} & \mathbf{I} & \mathbf{O} & \mathbf{O} \\ \mathbf{O} & \mathbf{O} & \mathbf{I} & \mathbf{O} \\ \mathbf{O} & \mathbf{O} & \mathbf{O} & \mathbf{I} \end{bmatrix}$$

$$\mathbf{A}_{11} \text{Cov}_q(\boldsymbol{\beta}) + \sum_{i=1}^3 \mathbf{A}_{12,i} \text{Cov}_q(\boldsymbol{\beta}, \mathbf{u}_i) = \mathbf{I}$$

$$\mathbf{A}_{12,i}^T \text{Cov}_q(\boldsymbol{\beta}, \mathbf{u}_i) + \mathbf{A}_{22,i} \text{Cov}_q(\mathbf{u}_i) = \mathbf{I}, \quad i = 1, 2, 3$$

$$\mathbf{A}_{12,i}^T \text{Cov}_q(\boldsymbol{\beta}) + \mathbf{A}_{22,i} \text{Cov}_q(\boldsymbol{\beta}, \mathbf{u}_i) = \mathbf{O}, \quad i = 1, 2, 3$$

# Multilevel Model ( $\mathbf{A}\text{Cov}_q(\boldsymbol{\nu}) = \mathbf{I}$ )

$$\mathbf{A}\text{Cov}_q(\boldsymbol{\nu}) = \begin{bmatrix} \mathbf{A}_{11} & \mathbf{A}_{12,1} & \mathbf{A}_{12,2} & \mathbf{A}_{12,3} \\ \mathbf{A}_{12,1}^T & \mathbf{A}_{22,1} & \mathbf{O} & \mathbf{O} \\ \mathbf{A}_{12,2}^T & \mathbf{O} & \mathbf{A}_{22,2} & \mathbf{O} \\ \mathbf{A}_{12,3}^T & \mathbf{O} & \mathbf{O} & \mathbf{A}_{22,3} \end{bmatrix} \begin{bmatrix} \text{Cov}_q(\boldsymbol{\beta}) & \text{Cov}_q(\boldsymbol{\beta}, \mathbf{u}_1) & \text{Cov}_q(\boldsymbol{\beta}, \mathbf{u}_2) & \text{Cov}_q(\boldsymbol{\beta}, \mathbf{u}_3) \\ \text{Cov}_q(\boldsymbol{\beta}, \mathbf{u}_1)^T & \text{Cov}_q(\mathbf{u}_1) & \times & \times \\ \text{Cov}_q(\boldsymbol{\beta}, \mathbf{u}_2)^T & \times & \text{Cov}_q(\mathbf{u}_2) & \times \\ \text{Cov}_q(\boldsymbol{\beta}, \mathbf{u}_3)^T & \times & \times & \text{Cov}_q(\mathbf{u}_3) \end{bmatrix}$$

$$= \begin{bmatrix} \mathbf{I} & \mathbf{O} & \mathbf{O} & \mathbf{O} \\ \mathbf{O} & \mathbf{I} & \mathbf{O} & \mathbf{O} \\ \mathbf{O} & \mathbf{O} & \mathbf{I} & \mathbf{O} \\ \mathbf{O} & \mathbf{O} & \mathbf{O} & \mathbf{I} \end{bmatrix}$$

$$\mathbf{A}_{11} \text{Cov}_q(\boldsymbol{\beta}) + \sum_{i=1}^3 \mathbf{A}_{12,i} \text{Cov}_q(\boldsymbol{\beta}, \mathbf{u}_i) = \mathbf{I}$$

$$\mathbf{A}_{12,i}^T \text{Cov}_q(\boldsymbol{\beta}, \mathbf{u}_i) + \mathbf{A}_{22,i} \text{Cov}_q(\mathbf{u}_i) = \mathbf{I}, \quad i = 1, 2, 3$$

$$\mathbf{A}_{12,i}^T \text{Cov}_q(\boldsymbol{\beta}) + \mathbf{A}_{22,i} \text{Cov}_q(\boldsymbol{\beta}, \mathbf{u}_i) = \mathbf{O}, \quad i = 1, 2, 3$$

# Multilevel Model ( $\mathbf{A}\text{Cov}_q(\boldsymbol{\nu}) = \mathbf{I}$ )

$$\mathbf{A}\text{Cov}_q(\boldsymbol{\nu}) = \begin{bmatrix} \mathbf{A}_{11} & \mathbf{A}_{12,1} & \mathbf{A}_{12,2} & \mathbf{A}_{12,3} \\ \mathbf{A}_{12,1}^T & \mathbf{A}_{22,1} & \mathbf{O} & \mathbf{O} \\ \mathbf{A}_{12,2}^T & \mathbf{O} & \mathbf{A}_{22,2} & \mathbf{O} \\ \mathbf{A}_{12,3}^T & \mathbf{O} & \mathbf{O} & \mathbf{A}_{22,3} \end{bmatrix} \begin{bmatrix} \text{Cov}_q(\boldsymbol{\beta}) & \text{Cov}_q(\boldsymbol{\beta}, \mathbf{u}_1) & \text{Cov}_q(\boldsymbol{\beta}, \mathbf{u}_2) & \text{Cov}_q(\boldsymbol{\beta}, \mathbf{u}_3) \\ \text{Cov}_q(\boldsymbol{\beta}, \mathbf{u}_1)^T & \text{Cov}_q(\mathbf{u}_1) & \times & \times \\ \text{Cov}_q(\boldsymbol{\beta}, \mathbf{u}_2)^T & \times & \text{Cov}_q(\mathbf{u}_2) & \times \\ \text{Cov}_q(\boldsymbol{\beta}, \mathbf{u}_3)^T & \times & \times & \text{Cov}_q(\mathbf{u}_3) \end{bmatrix}$$

$$= \begin{bmatrix} \mathbf{I} & \mathbf{O} & \mathbf{O} & \mathbf{O} \\ \mathbf{O} & \mathbf{I} & \mathbf{O} & \mathbf{O} \\ \mathbf{O} & \mathbf{O} & \mathbf{I} & \mathbf{O} \\ \mathbf{O} & \mathbf{O} & \mathbf{O} & \mathbf{I} \end{bmatrix}$$

$$\mathbf{A}_{11} \text{Cov}_q(\boldsymbol{\beta}) + \sum_{i=1}^3 \mathbf{A}_{12,i} \text{Cov}_q(\boldsymbol{\beta}, \mathbf{u}_i) = \mathbf{I}$$

$$\mathbf{A}_{12,i}^T \text{Cov}_q(\boldsymbol{\beta}, \mathbf{u}_i) + \mathbf{A}_{22,i} \text{Cov}_q(\mathbf{u}_i) = \mathbf{I}, \quad i = 1, 2, 3$$

$$\mathbf{A}_{12,i}^T \text{Cov}_q(\boldsymbol{\beta}) + \mathbf{A}_{22,i} \text{Cov}_q(\boldsymbol{\beta}, \mathbf{u}_i) = \mathbf{O}, \quad i = 1, 2, 3$$

# Multilevel Model ( $\mathbf{A}\text{Cov}_q(\boldsymbol{\nu}) = \mathbf{I}$ )

$$\mathbf{A}\text{Cov}_q(\boldsymbol{\nu}) = \begin{bmatrix} \mathbf{A}_{11} & \mathbf{A}_{12,1} & \mathbf{A}_{12,2} & \mathbf{A}_{12,3} \\ \mathbf{A}_{12,1}^T & \mathbf{A}_{22,1} & \mathbf{O} & \mathbf{O} \\ \mathbf{A}_{12,2}^T & \mathbf{O} & \mathbf{A}_{22,2} & \mathbf{O} \\ \mathbf{A}_{12,3}^T & \mathbf{O} & \mathbf{O} & \mathbf{A}_{22,3} \end{bmatrix} \begin{bmatrix} \text{Cov}_q(\boldsymbol{\beta}) & \text{Cov}_q(\boldsymbol{\beta}, \mathbf{u}_1) & \text{Cov}_q(\boldsymbol{\beta}, \mathbf{u}_2) & \text{Cov}_q(\boldsymbol{\beta}, \mathbf{u}_3) \\ \text{Cov}_q(\boldsymbol{\beta}, \mathbf{u}_1)^T & \text{Cov}_q(\mathbf{u}_1) & \times & \times \\ \text{Cov}_q(\boldsymbol{\beta}, \mathbf{u}_2)^T & \times & \text{Cov}_q(\mathbf{u}_2) & \times \\ \text{Cov}_q(\boldsymbol{\beta}, \mathbf{u}_3)^T & \times & \times & \text{Cov}_q(\mathbf{u}_3) \end{bmatrix}$$

$$= \begin{bmatrix} \mathbf{I} & \mathbf{O} & \mathbf{O} & \mathbf{O} \\ \mathbf{O} & \mathbf{I} & \mathbf{O} & \mathbf{O} \\ \mathbf{O} & \mathbf{O} & \mathbf{I} & \mathbf{O} \\ \mathbf{O} & \mathbf{O} & \mathbf{O} & \mathbf{I} \end{bmatrix}$$

$$\mathbf{A}_{11} \text{Cov}_q(\boldsymbol{\beta}) + \sum_{i=1}^3 \mathbf{A}_{12,i} \text{Cov}_q(\boldsymbol{\beta}, \mathbf{u}_i) = \mathbf{I}$$

$$\mathbf{A}_{12,i}^T \text{Cov}_q(\boldsymbol{\beta}, \mathbf{u}_i) + \mathbf{A}_{22,i} \text{Cov}_q(\mathbf{u}_i) = \mathbf{I}, \quad i = 1, 2, 3$$

$$\mathbf{A}_{12,i}^T \text{Cov}_q(\boldsymbol{\beta}) + \mathbf{A}_{22,i} \text{Cov}_q(\boldsymbol{\beta}, \mathbf{u}_i) = \mathbf{O}, \quad i = 1, 2, 3$$

# Multilevel Model ( $\mathbf{A}\text{Cov}_q(\boldsymbol{\nu}) = \mathbf{I}$ )

$$\mathbf{A}\text{Cov}_q(\boldsymbol{\nu}) = \begin{bmatrix} \mathbf{A}_{11} & \mathbf{A}_{12,1} & \mathbf{A}_{12,2} & \mathbf{A}_{12,3} \\ \mathbf{A}_{12,1}^T & \mathbf{A}_{22,1} & \mathbf{O} & \mathbf{O} \\ \mathbf{A}_{12,2}^T & \mathbf{O} & \mathbf{A}_{22,2} & \mathbf{O} \\ \mathbf{A}_{12,3}^T & \mathbf{O} & \mathbf{O} & \mathbf{A}_{22,3} \end{bmatrix} \begin{bmatrix} \text{Cov}_q(\boldsymbol{\beta}) & \text{Cov}_q(\boldsymbol{\beta}, \mathbf{u}_1) & \text{Cov}_q(\boldsymbol{\beta}, \mathbf{u}_2) & \text{Cov}_q(\boldsymbol{\beta}, \mathbf{u}_3) \\ \text{Cov}_q(\boldsymbol{\beta}, \mathbf{u}_1)^T & \text{Cov}_q(\mathbf{u}_1) & \times & \times \\ \text{Cov}_q(\boldsymbol{\beta}, \mathbf{u}_2)^T & \times & \text{Cov}_q(\mathbf{u}_2) & \times \\ \text{Cov}_q(\boldsymbol{\beta}, \mathbf{u}_3)^T & \times & \times & \text{Cov}_q(\mathbf{u}_3) \end{bmatrix}$$

$$= \begin{bmatrix} \mathbf{I} & \mathbf{O} & \mathbf{O} & \mathbf{O} \\ \mathbf{O} & \mathbf{I} & \mathbf{O} & \mathbf{O} \\ \mathbf{O} & \mathbf{O} & \mathbf{I} & \mathbf{O} \\ \mathbf{O} & \mathbf{O} & \mathbf{O} & \mathbf{I} \end{bmatrix}$$

$$\mathbf{A}_{11} \text{Cov}_q(\boldsymbol{\beta}) + \sum_{i=1}^3 \mathbf{A}_{12,i} \text{Cov}_q(\boldsymbol{\beta}, \mathbf{u}_i) = \mathbf{I}$$

$$\mathbf{A}_{12,i}^T \text{Cov}_q(\boldsymbol{\beta}, \mathbf{u}_i) + \mathbf{A}_{22,i} \text{Cov}_q(\mathbf{u}_i) = \mathbf{I}, \quad i = 1, 2, 3$$

$$\mathbf{A}_{12,i}^T \text{Cov}_q(\boldsymbol{\beta}) + \mathbf{A}_{22,i} \text{Cov}_q(\boldsymbol{\beta}, \mathbf{u}_i) = \mathbf{O}, \quad i = 1, 2, 3$$



# Multilevel Model ( $\mathbf{A} \mathbb{E}_q(\boldsymbol{\nu}) = \mathbf{a}$ )

$$\mathbf{A} \mathbb{E}_q(\boldsymbol{\nu}) = \begin{bmatrix} \mathbf{A}_{11} & \mathbf{A}_{12,1} & \mathbf{A}_{12,2} & \mathbf{A}_{12,3} \\ \mathbf{A}_{12,1}^T & \mathbf{A}_{22,1} & \mathbf{O} & \mathbf{O} \\ \mathbf{A}_{12,2}^T & \mathbf{O} & \mathbf{A}_{22,2} & \mathbf{O} \\ \mathbf{A}_{12,3}^T & \mathbf{O} & \mathbf{O} & \mathbf{A}_{22,3} \end{bmatrix} \begin{bmatrix} \mathbb{E}_q(\beta) \\ \mathbb{E}_q(\mathbf{u}_1) \\ \mathbb{E}_q(\mathbf{u}_2) \\ \mathbb{E}_q(\mathbf{u}_3) \end{bmatrix} = \begin{bmatrix} \mathbf{a}_1 \\ \mathbf{a}_{2,1} \\ \mathbf{a}_{2,2} \\ \mathbf{a}_{2,3} \end{bmatrix}$$

$$\mathbf{A}_{11} \mathbb{E}_q(\beta) + \sum_{i=1}^3 \mathbf{A}_{12,i} \mathbb{E}_q(\mathbf{u}_i) = \mathbf{a}_1$$

$$\mathbf{A}_{12,i}^T \mathbb{E}_q(\beta) + \mathbf{A}_{22,i} \mathbb{E}_q(\mathbf{u}_i) = \mathbf{I}, \quad i = 1, 2, 3$$

# Multilevel Model ( $\mathbf{A} \mathbb{E}_q(\boldsymbol{\nu}) = \mathbf{a}$ )

$$\mathbf{A} \mathbb{E}_q(\boldsymbol{\nu}) = \begin{bmatrix} \mathbf{A}_{11} & \mathbf{A}_{12,1} & \mathbf{A}_{12,2} & \mathbf{A}_{12,3} \\ \mathbf{A}_{12,1}^T & \mathbf{A}_{22,1} & \mathbf{O} & \mathbf{O} \\ \mathbf{A}_{12,2}^T & \mathbf{O} & \mathbf{A}_{22,2} & \mathbf{O} \\ \mathbf{A}_{12,3}^T & \mathbf{O} & \mathbf{O} & \mathbf{A}_{22,3} \end{bmatrix} \begin{bmatrix} \mathbb{E}_q(\beta) \\ \mathbb{E}_q(\mathbf{u}_1) \\ \mathbb{E}_q(\mathbf{u}_2) \\ \mathbb{E}_q(\mathbf{u}_3) \end{bmatrix} = \begin{bmatrix} \mathbf{a}_1 \\ \mathbf{a}_{2,1} \\ \mathbf{a}_{2,2} \\ \mathbf{a}_{2,3} \end{bmatrix}$$

$$\mathbf{A}_{11} \mathbb{E}_q(\beta) + \sum_{i=1}^3 \mathbf{A}_{12,i} \mathbb{E}_q(\mathbf{u}_i) = \mathbf{a}_1$$

$$\mathbf{A}_{12,i}^T \mathbb{E}_q(\beta) + \mathbf{A}_{22,i} \mathbb{E}_q(\mathbf{u}_i) = \mathbf{I}, \quad i = 1, 2, 3$$

# Multilevel Model ( $\mathbf{A} \mathbb{E}_q(\boldsymbol{\nu}) = \mathbf{a}$ )

$$\mathbf{A} \mathbb{E}_q(\boldsymbol{\nu}) = \begin{bmatrix} \mathbf{A}_{11} & \mathbf{A}_{12,1} & \mathbf{A}_{12,2} & \mathbf{A}_{12,3} \\ \mathbf{A}_{12,1}^T & \mathbf{A}_{22,1} & \mathbf{O} & \mathbf{O} \\ \mathbf{A}_{12,2}^T & \mathbf{O} & \mathbf{A}_{22,2} & \mathbf{O} \\ \mathbf{A}_{12,3}^T & \mathbf{O} & \mathbf{O} & \mathbf{A}_{22,3} \end{bmatrix} \begin{bmatrix} \mathbb{E}_q(\boldsymbol{\beta}) \\ \mathbb{E}_q(\mathbf{u}_1) \\ \mathbb{E}_q(\mathbf{u}_2) \\ \mathbb{E}_q(\mathbf{u}_3) \end{bmatrix} = \begin{bmatrix} \mathbf{a}_1 \\ \mathbf{a}_{2,1} \\ \mathbf{a}_{2,2} \\ \mathbf{a}_{2,3} \end{bmatrix}$$

$$\mathbf{A}_{11} \mathbb{E}_q(\boldsymbol{\beta}) + \sum_{i=1}^3 \mathbf{A}_{12,i} \mathbb{E}_q(\mathbf{u}_i) = \mathbf{a}_1$$

$$\mathbf{A}_{12,i}^T \mathbb{E}_q(\boldsymbol{\beta}) + \mathbf{A}_{22,i} \mathbb{E}_q(\mathbf{u}_i) = \mathbf{I}, \quad i = 1, 2, 3$$

# Multilevel Model ( $\mathbf{A} \mathbb{E}_q(\boldsymbol{\nu}) = \mathbf{a}$ )

$$\mathbf{A} \mathbb{E}_q(\boldsymbol{\nu}) = \begin{bmatrix} \mathbf{A}_{11} & \mathbf{A}_{12,1} & \mathbf{A}_{12,2} & \mathbf{A}_{12,3} \\ \mathbf{A}_{12,1}^T & \mathbf{A}_{22,1} & \mathbf{O} & \mathbf{O} \\ \mathbf{A}_{12,2}^T & \mathbf{O} & \mathbf{A}_{22,2} & \mathbf{O} \\ \mathbf{A}_{12,3}^T & \mathbf{O} & \mathbf{O} & \mathbf{A}_{22,3} \end{bmatrix} \begin{bmatrix} \mathbb{E}_q(\beta) \\ \mathbb{E}_q(\mathbf{u}_1) \\ \mathbb{E}_q(\mathbf{u}_2) \\ \mathbb{E}_q(\mathbf{u}_3) \end{bmatrix} = \begin{bmatrix} \mathbf{a}_1 \\ \mathbf{a}_{2,1} \\ \mathbf{a}_{2,2} \\ \mathbf{a}_{2,3} \end{bmatrix}$$

$$\mathbf{A}_{11} \mathbb{E}_q(\beta) + \sum_{i=1}^3 \mathbf{A}_{12,i} \mathbb{E}_q(\mathbf{u}_i) = \mathbf{a}_1$$

$$\mathbf{A}_{12,i}^T \mathbb{E}_q(\beta) + \mathbf{A}_{22,i} \mathbb{E}_q(\mathbf{u}_i) = \mathbf{I}, \quad i = 1, 2, 3$$

# Multilevel Model ( $\mathbf{A} \mathbb{E}_q(\boldsymbol{\nu}) = \mathbf{a}$ )

$$\mathbf{A} \mathbb{E}_q(\boldsymbol{\nu}) = \begin{bmatrix} \mathbf{A}_{11} & \mathbf{A}_{12,1} & \mathbf{A}_{12,2} & \mathbf{A}_{12,3} \\ \mathbf{A}_{12,1}^T & \mathbf{A}_{22,1} & \mathbf{O} & \mathbf{O} \\ \mathbf{A}_{12,2}^T & \mathbf{O} & \mathbf{A}_{22,2} & \mathbf{O} \\ \mathbf{A}_{12,3}^T & \mathbf{O} & \mathbf{O} & \mathbf{A}_{22,3} \end{bmatrix} \begin{bmatrix} \mathbb{E}_q(\beta) \\ \mathbb{E}_q(\mathbf{u}_1) \\ \mathbb{E}_q(\mathbf{u}_2) \\ \mathbb{E}_q(\mathbf{u}_3) \end{bmatrix} = \begin{bmatrix} \mathbf{a}_1 \\ \mathbf{a}_{2,1} \\ \mathbf{a}_{2,2} \\ \mathbf{a}_{2,3} \end{bmatrix}$$

$$\mathbf{A}_{11} \mathbb{E}_q(\beta) + \sum_{i=1}^3 \mathbf{A}_{12,i} \mathbb{E}_q(\mathbf{u}_i) = \mathbf{a}_1$$

$$\mathbf{A}_{12,i}^T \mathbb{E}_q(\beta) + \mathbf{A}_{22,i} \mathbb{E}_q(\mathbf{u}_i) = \mathbf{I}, \quad i = 1, 2, 3$$

# Multilevel Model ( $\mathbf{A} \mathbb{E}_q(\boldsymbol{\nu}) = \mathbf{a}$ )

$$\mathbf{A} \mathbb{E}_q(\boldsymbol{\nu}) = \begin{bmatrix} \mathbf{A}_{11} & \mathbf{A}_{12,1} & \mathbf{A}_{12,2} & \mathbf{A}_{12,3} \\ \mathbf{A}_{12,1}^T & \mathbf{A}_{22,1} & \mathbf{O} & \mathbf{O} \\ \mathbf{A}_{12,2}^T & \mathbf{O} & \mathbf{A}_{22,2} & \mathbf{O} \\ \mathbf{A}_{12,3}^T & \mathbf{O} & \mathbf{O} & \mathbf{A}_{22,3} \end{bmatrix} \begin{bmatrix} \mathbb{E}_q(\beta) \\ \mathbb{E}_q(\mathbf{u}_1) \\ \mathbb{E}_q(\mathbf{u}_2) \\ \mathbb{E}_q(\mathbf{u}_3) \end{bmatrix} = \begin{bmatrix} \mathbf{a}_1 \\ \mathbf{a}_{2,1} \\ \mathbf{a}_{2,2} \\ \mathbf{a}_{2,3} \end{bmatrix}$$

$$\mathbf{A}_{11} \mathbb{E}_q(\beta) + \sum_{i=1}^3 \mathbf{A}_{12,i} \mathbb{E}_q(\mathbf{u}_i) = \mathbf{a}_1$$

$$\mathbf{A}_{12,i}^T \mathbb{E}_q(\beta) + \mathbf{A}_{22,i} \mathbb{E}_q(\mathbf{u}_i) = \mathbf{I}, \quad i = 1, 2, 3$$

## Theorem: Nolan and Wand (2020)

For the two-level mean field variational Bayes model, the solutions for the required sub-blocks of  $\text{Cov}_q(\boldsymbol{\nu})$  are

$$\text{Cov}_q(\boldsymbol{\beta}) = \left( \mathbf{A}_{11} - \sum_{i=1}^m \mathbf{A}_{12,i} \mathbf{A}_{22,i}^{-1} \mathbf{A}_{12,i}^T \right)^{-1}$$

$$\text{Cov}_q(\boldsymbol{\beta}, \mathbf{u}_i) = -(\mathbf{A}_{22,i}^{-1} \mathbf{A}_{12,i}^T \text{Cov}_q(\boldsymbol{\beta}))^T, \quad \text{Cov}_q(\mathbf{u}_i) = \mathbf{A}_{22,i}^{-1} (\mathbf{I} - \mathbf{A}_{12,i}^T \text{Cov}_q(\boldsymbol{\beta}, \mathbf{u}_i)), \quad 1 \leq i \leq m.$$

The determinant of  $\text{Cov}_q(\boldsymbol{\nu})$  is

$$|\text{Cov}_q(\boldsymbol{\nu})| = |\text{Cov}_q(\boldsymbol{\beta})| \prod_{i=1}^m |\text{Cov}_q(\mathbf{u}_i)|.$$

The solutions for the sub-vectors of  $\boldsymbol{\mu}_q(\boldsymbol{\beta}, \mathbf{u})$  are

$$\mathbb{E}_q(\boldsymbol{\beta}) = \text{Cov}_q(\boldsymbol{\beta}) \left( \mathbf{a}_1 - \sum_{i=1}^m \mathbf{A}_{12,i} \mathbf{A}_{22,i}^{-1} \mathbf{a}_{2,i} \right) \quad \text{and} \quad \mathbb{E}_q(\mathbf{u}_i) = \mathbf{A}_{22,i}^{-1} (\mathbf{a}_{2,i} - \mathbf{A}_{12,i}^T \mathbb{E}_q(\boldsymbol{\beta})).$$

## Theorem: Nolan and Wand (2020)

For the two-level mean field variational Bayes model, the solutions for the required sub-blocks of  $\text{Cov}_q(\boldsymbol{\nu})$  are

$$\text{Cov}_q(\boldsymbol{\beta}) = \left( \mathbf{A}_{11} - \sum_{i=1}^m \mathbf{A}_{12,i} \mathbf{A}_{22,i}^{-1} \mathbf{A}_{12,i}^T \right)^{-1}$$

$$\text{Cov}_q(\boldsymbol{\beta}, \mathbf{u}_i) = -(\mathbf{A}_{22,i}^{-1} \mathbf{A}_{12,i}^T \text{Cov}_q(\boldsymbol{\beta}))^T, \quad \text{Cov}_q(\mathbf{u}_i) = \mathbf{A}_{22,i}^{-1} (\mathbf{I} - \mathbf{A}_{12,i}^T \text{Cov}_q(\boldsymbol{\beta}, \mathbf{u}_i)), \quad 1 \leq i \leq m.$$

The determinant of  $\text{Cov}_q(\boldsymbol{\nu})$  is

$$|\text{Cov}_q(\boldsymbol{\nu})| = |\text{Cov}_q(\boldsymbol{\beta})| \prod_{i=1}^m |\text{Cov}_q(\mathbf{u}_i)|.$$

The solutions for the sub-vectors of  $\boldsymbol{\mu}_q(\boldsymbol{\beta}, \mathbf{u})$  are

$$\mathbb{E}_q(\boldsymbol{\beta}) = \text{Cov}_q(\boldsymbol{\beta}) \left( \mathbf{a}_1 - \sum_{i=1}^m \mathbf{A}_{12,i} \mathbf{A}_{22,i}^{-1} \mathbf{a}_{2,i} \right) \quad \text{and} \quad \mathbb{E}_q(\mathbf{u}_i) = \mathbf{A}_{22,i}^{-1} (\mathbf{a}_{2,i} - \mathbf{A}_{12,i}^T \mathbb{E}_q(\boldsymbol{\beta})).$$



## Theorem: Nolan and Wand (2020)

For the two-level mean field variational Bayes model, the solutions for the required sub-blocks of  $\text{Cov}_q(\boldsymbol{\nu})$  are

$$\text{Cov}_q(\boldsymbol{\beta}) = \left( \mathbf{A}_{11} - \sum_{i=1}^m \mathbf{A}_{12,i} \mathbf{A}_{22,i}^{-1} \mathbf{A}_{12,i}^T \right)^{-1}$$

$$\text{Cov}_q(\boldsymbol{\beta}, \mathbf{u}_i) = -(\mathbf{A}_{22,i}^{-1} \mathbf{A}_{12,i}^T \text{Cov}_q(\boldsymbol{\beta}))^T, \quad \text{Cov}_q(\mathbf{u}_i) = \mathbf{A}_{22,i}^{-1} (\mathbf{I} - \mathbf{A}_{12,i}^T \text{Cov}_q(\boldsymbol{\beta}, \mathbf{u}_i)), \quad 1 \leq i \leq m.$$

The determinant of  $\text{Cov}_q(\boldsymbol{\nu})$  is

$$|\text{Cov}_q(\boldsymbol{\nu})| = |\text{Cov}_q(\boldsymbol{\beta})| \prod_{i=1}^m |\text{Cov}_q(\mathbf{u}_i)|.$$

The solutions for the sub-vectors of  $\boldsymbol{\mu}_q(\boldsymbol{\beta}, \mathbf{u})$  are

$$\mathbb{E}_q(\boldsymbol{\beta}) = \text{Cov}_q(\boldsymbol{\beta}) \left( \mathbf{a}_1 - \sum_{i=1}^m \mathbf{A}_{12,i} \mathbf{A}_{22,i}^{-1} \mathbf{a}_{2,i} \right) \quad \text{and} \quad \mathbb{E}_q(\mathbf{u}_i) = \mathbf{A}_{22,i}^{-1} (\mathbf{a}_{2,i} - \mathbf{A}_{12,i}^T \mathbb{E}_q(\boldsymbol{\beta})).$$

## Theorem: Nolan and Wand (2020)

For the two-level mean field variational Bayes model, the solutions for the required sub-blocks of  $\text{Cov}_q(\boldsymbol{\nu})$  are

$$\text{Cov}_q(\boldsymbol{\beta}) = \left( \mathbf{A}_{11} - \sum_{i=1}^m \mathbf{A}_{12,i} \mathbf{A}_{22,i}^{-1} \mathbf{A}_{12,i}^T \right)^{-1}$$

$$\text{Cov}_q(\boldsymbol{\beta}, \mathbf{u}_i) = -(\mathbf{A}_{22,i}^{-1} \mathbf{A}_{12,i}^T \text{Cov}_q(\boldsymbol{\beta}))^T, \quad \text{Cov}_q(\mathbf{u}_i) = \mathbf{A}_{22,i}^{-1} (\mathbf{I} - \mathbf{A}_{12,i}^T \text{Cov}_q(\boldsymbol{\beta}, \mathbf{u}_i)), \quad 1 \leq i \leq m.$$

The determinant of  $\text{Cov}_q(\boldsymbol{\nu})$  is

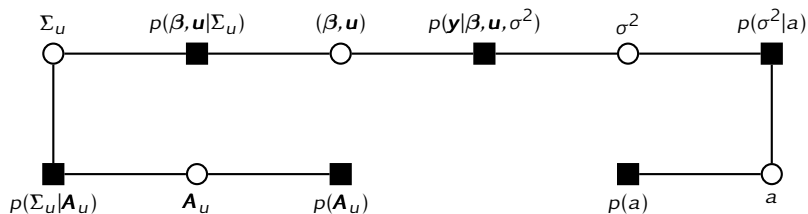
$$|\text{Cov}_q(\boldsymbol{\nu})| = |\text{Cov}_q(\boldsymbol{\beta})| \prod_{i=1}^m |\text{Cov}_q(\mathbf{u}_i)|.$$

The solutions for the sub-vectors of  $\boldsymbol{\mu}_q(\boldsymbol{\beta}, \mathbf{u})$  are

$$\mathbb{E}_q(\boldsymbol{\beta}) = \text{Cov}_q(\boldsymbol{\beta}) \left( \mathbf{a}_1 - \sum_{i=1}^m \mathbf{A}_{12,i} \mathbf{A}_{22,i}^{-1} \mathbf{a}_{2,i} \right) \quad \text{and} \quad \mathbb{E}_q(\mathbf{u}_i) = \mathbf{A}_{22,i}^{-1} (\mathbf{a}_{2,i} - \mathbf{A}_{12,i}^T \mathbb{E}_q(\boldsymbol{\beta})).$$

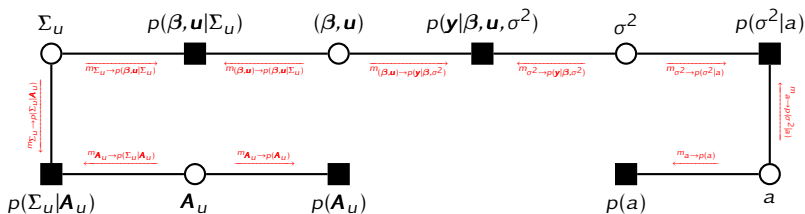
## VMP for the Bayesian multilevel data model

1. Initialise all messages from stochastic nodes to factors
2. Cycle:
  - (i) Update all messages from factors to stochastic nodes
  - (ii) Update all messages from stochastic nodes to factors
  - (iii) Update all optimal posterior density functions
3. Stop:  $D_{\text{KL}} \{q(\beta, \mathbf{u}, \sigma^2, \Sigma_u, a, \mathbf{A}_u) \| p(\beta, \mathbf{u}, \sigma^2, \Sigma_u, a, \mathbf{A}_u | \mathbf{y})\}$  converges.



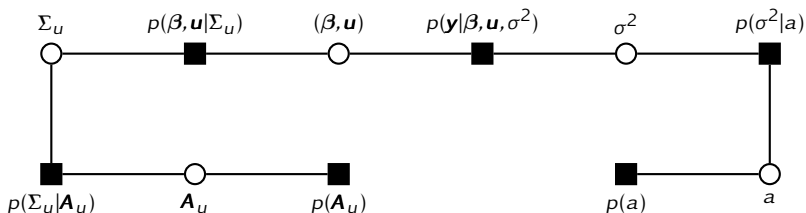
## VMP for the Bayesian multilevel data model

1. Initialise all messages from stochastic nodes to factors
2. Cycle:
  - (i) Update all messages from factors to stochastic nodes
  - (ii) Update all messages from stochastic nodes to factors
  - (iii) Update all optimal posterior density functions
3. Stop:  $D_{\text{KL}} \{q(\beta, \mathbf{u}, \sigma^2, \Sigma_u, a, \mathbf{A}_u) \| p(\beta, \mathbf{u}, \sigma^2, \Sigma_u, a, \mathbf{A}_u | \mathbf{y})\}$  converges.



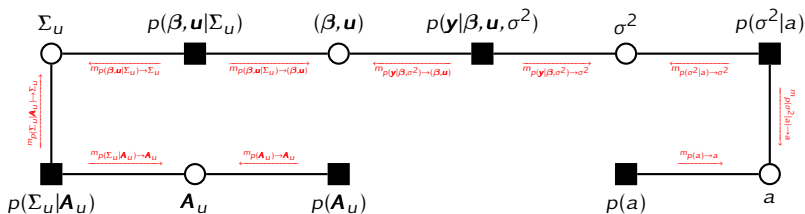
## VMP for the Bayesian multilevel data model

1. Initialise all messages from stochastic nodes to factors
2. Cycle:
  - (i) Update all messages from factors to stochastic nodes
  - (ii) Update all messages from stochastic nodes to factors
  - (iii) Update all optimal posterior density functions
3. Stop:  $D_{\text{KL}} \{q(\beta, \mathbf{u}, \sigma^2, \Sigma_U, a, \mathbf{A}_U) \| p(\beta, \mathbf{u}, \sigma^2, \Sigma_U, a, \mathbf{A}_U | \mathbf{y})\}$  converges.



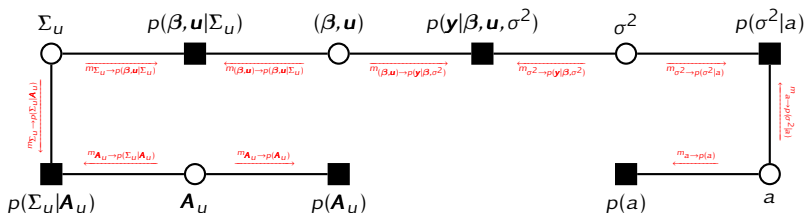
## VMP for the Bayesian multilevel data model

1. Initialise all messages from stochastic nodes to factors
2. Cycle:
  - (i) Update all messages from factors to stochastic nodes
  - (ii) Update all messages from stochastic nodes to factors
  - (iii) Update all optimal posterior density functions
3. Stop:  $D_{KL} \{q(\beta, \mathbf{u}, \sigma^2, \Sigma_U, a, \mathbf{A}_U) \| p(\beta, \mathbf{u}, \sigma^2, \Sigma_U, a, \mathbf{A}_U | \mathbf{y})\}$  converges.



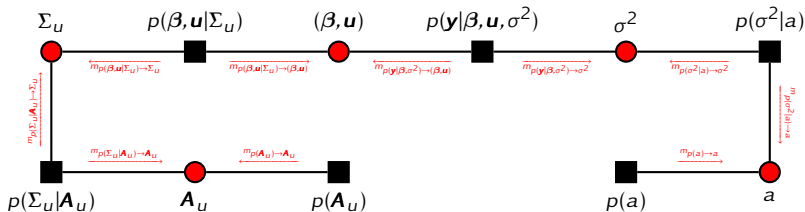
## VMP for the Bayesian multilevel data model

1. Initialise all messages from stochastic nodes to factors
2. Cycle:
  - (i) Update all messages from factors to stochastic nodes
  - (ii) Update all messages from stochastic nodes to factors
  - (iii) Update all optimal posterior density functions
3. Stop:  $D_{\text{KL}} \{q(\beta, \mathbf{u}, \sigma^2, \Sigma_U, a, \mathbf{A}_U) \| p(\beta, \mathbf{u}, \sigma^2, \Sigma_U, a, \mathbf{A}_U | \mathbf{y})\}$  converges.



## VMP for the Bayesian multilevel data model

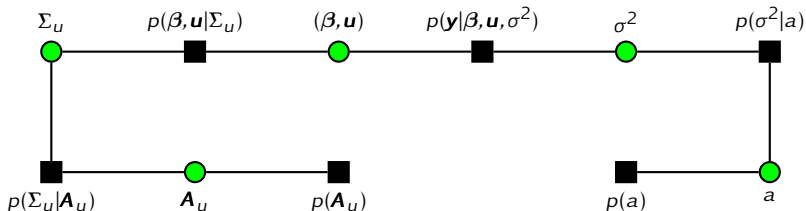
1. Initialise all messages from stochastic nodes to factors
2. Cycle:
  - (i) Update all messages from factors to stochastic nodes
  - (ii) Update all messages from stochastic nodes to factors
  - (iii) Update all optimal posterior density functions
3. Stop:  $D_{\text{KL}} \{q(\beta, \mathbf{u}, \sigma^2, \Sigma_U, a, \mathbf{A}_U) \| p(\beta, \mathbf{u}, \sigma^2, \Sigma_U, a, \mathbf{A}_U | \mathbf{y})\}$  converges.



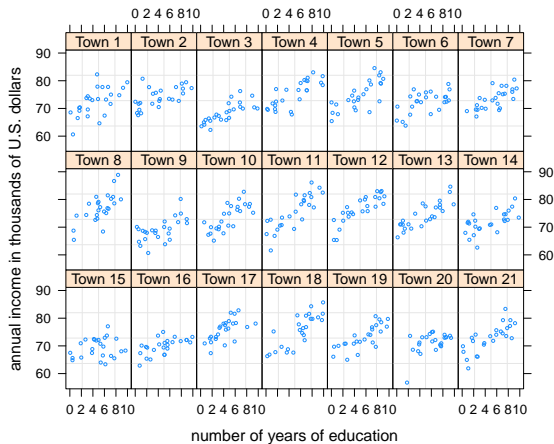


## VMP for the Bayesian multilevel data model

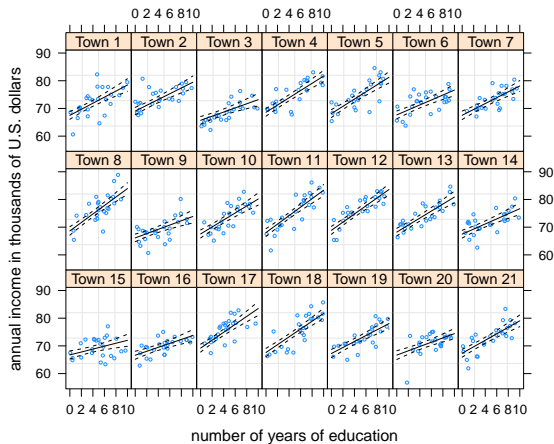
1. Initialise all messages from stochastic nodes to factors
2. Cycle:
  - (i) Update all messages from factors to stochastic nodes
  - (ii) Update all messages from stochastic nodes to factors
  - (iii) Update all optimal posterior density functions
3. Stop:  $D_{\text{KL}}\{q(\beta, \mathbf{u}, \sigma^2, \Sigma_u, a, \mathbf{A}_u) \| p(\beta, \mathbf{u}, \sigma^2, \Sigma_u, a, \mathbf{A}_u | \mathbf{y})\}$  converges.



# Multilevel Data



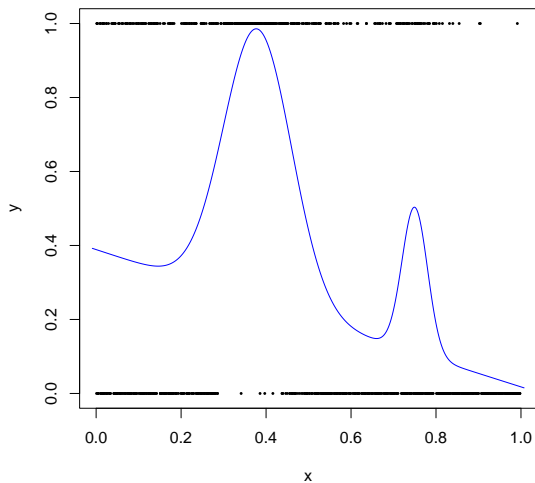
# Multilevel Data



## Part III

### Nonconjugate Models

## Binary Response Data



# Bayesian Logistic Regression Model

Consider the Bayesian logistic regression model:

$$\mathbf{y} \mid \boldsymbol{\beta} \sim \text{Bernoulli}\left\{\left[1 + \exp\{-\mathbf{X}\boldsymbol{\beta}\}\right]^{-1}\right\}$$
$$\boldsymbol{\beta} \sim N(\mathbf{0}, \sigma_{\boldsymbol{\beta}}^2 \mathbf{I})$$

The factor graph is



Gaussian Prior Fragment

Logistic Likelihood Fragment (Nolan and Wand, 2017)

The message takes the form:

$$m_{p(\mathbf{y}|\boldsymbol{\beta}) \rightarrow \boldsymbol{\beta}}(\boldsymbol{\beta}) \leftarrow \exp\left\{\mathbf{y}^T \mathbf{X} \boldsymbol{\beta} - \mathbf{1}^T \mathbb{E}_q \log[1 + \exp(\mathbf{X} \boldsymbol{\beta})]\right\}$$

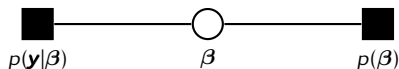
The messages that are typically passed to  $\boldsymbol{\beta}$  from other fragments are Gaussian density messages.

# Bayesian Logistic Regression Model

Consider the Bayesian logistic regression model:

$$\mathbf{y} \mid \boldsymbol{\beta} \sim \text{Bernoulli}\left\{\left[1 + \exp\{-\mathbf{X}\boldsymbol{\beta}\}\right]^{-1}\right\}$$
$$\boldsymbol{\beta} \sim N(\mathbf{0}, \sigma_{\boldsymbol{\beta}}^2 \mathbf{I})$$

The factor graph is



Gaussian Prior Fragment

Logistic Likelihood Fragment (Nolan and Wand, 2017)

The message takes the form:

$$m_{p(\mathbf{y}|\boldsymbol{\beta}) \rightarrow \boldsymbol{\beta}}(\boldsymbol{\beta}) \leftarrow \exp\left\{\mathbf{y}^T \mathbf{X} \boldsymbol{\beta} - \mathbf{1}^T \mathbb{E}_q \log[1 + \exp(\mathbf{X} \boldsymbol{\beta})]\right\}$$

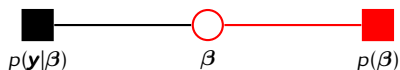
The messages that are typically passed to  $\boldsymbol{\beta}$  from other fragments are Gaussian density messages.

# Bayesian Logistic Regression Model

Consider the Bayesian logistic regression model:

$$\mathbf{y} \mid \boldsymbol{\beta} \sim \text{Bernoulli}\left\{\left[1 + \exp\{-\mathbf{X}\boldsymbol{\beta}\}\right]^{-1}\right\}$$
$$\boldsymbol{\beta} \sim N(\mathbf{0}, \sigma_{\boldsymbol{\beta}}^2 \mathbf{I})$$

The factor graph is



Gaussian Prior Fragment

Logistic Likelihood Fragment (Nolan and Wand, 2017)

The message takes the form:

$$m_{p(\mathbf{y}|\boldsymbol{\beta}) \rightarrow \boldsymbol{\beta}}(\boldsymbol{\beta}) \leftarrow \exp\left\{\mathbf{y}^T \mathbf{X} \boldsymbol{\beta} - \mathbf{1}^T \mathbb{E}_q \log[1 + \exp(\mathbf{X} \boldsymbol{\beta})]\right\}$$

The messages that are typically passed to  $\boldsymbol{\beta}$  from other fragments are Gaussian density messages.

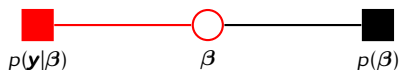


# Bayesian Logistic Regression Model

Consider the Bayesian logistic regression model:

$$\mathbf{y} \mid \boldsymbol{\beta} \sim \text{Bernoulli}\left\{\left[1 + \exp\{-\mathbf{X}\boldsymbol{\beta}\}\right]^{-1}\right\}$$
$$\boldsymbol{\beta} \sim N(\mathbf{0}, \sigma_{\boldsymbol{\beta}}^2 \mathbf{I})$$

The factor graph is



Gaussian Prior Fragment

Logistic Likelihood Fragment (Nolan and Wand, 2017)

The message takes the form:

$$m_{p(\mathbf{y}|\boldsymbol{\beta}) \rightarrow \boldsymbol{\beta}}(\boldsymbol{\beta}) \leftarrow \exp\left\{\mathbf{y}^T \mathbf{X} \boldsymbol{\beta} - \mathbf{1}^T \mathbb{E}_q \log[1 + \exp(\mathbf{X} \boldsymbol{\beta})]\right\}$$

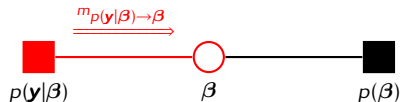
The messages that are typically passed to  $\boldsymbol{\beta}$  from other fragments are Gaussian density messages.

# Bayesian Logistic Regression Model

Consider the Bayesian logistic regression model:

$$\mathbf{y} \mid \beta \sim \text{Bernoulli}\left\{\left[1 + \exp\{-\mathbf{X}\beta\}\right]^{-1}\right\}$$
$$\beta \sim N(\mathbf{0}, \sigma_\beta^2 \mathbf{I})$$

The factor graph is



Gaussian Prior Fragment

Logistic Likelihood Fragment (Nolan and Wand, 2017)

The message takes the form:

$$m_{p(\mathbf{y}|\beta) \rightarrow \beta}(\beta) \leftarrow \exp\left\{\mathbf{y}^T \mathbf{X} \beta - \mathbf{1}^T \mathbb{E}_q \log[1 + \exp(\mathbf{X} \beta)]\right\}$$

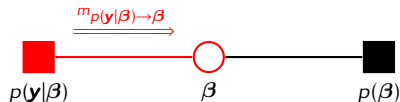
The messages that are typically passed to  $\beta$  from other fragments are Gaussian density messages.

# Bayesian Logistic Regression Model

Consider the Bayesian logistic regression model:

$$\mathbf{y} \mid \beta \sim \text{Bernoulli}\left\{\left[1 + \exp\{-\mathbf{X}\beta\}\right]^{-1}\right\}$$
$$\beta \sim N(\mathbf{0}, \sigma_\beta^2 \mathbf{I})$$

The factor graph is



Gaussian Prior Fragment

Logistic Likelihood Fragment (Nolan and Wand, 2017)

The message takes the form:

$$m_{p(\mathbf{y}|\beta) \rightarrow \beta}(\beta) \leftarrow \exp\left\{\mathbf{y}^T \mathbf{X} \beta - \mathbf{1}^T \mathbb{E}_q \log[1 + \exp(\mathbf{X} \beta)]\right\}$$

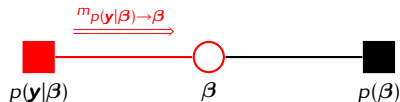
The messages that are typically passed to  $\beta$  from other fragments are Gaussian density messages.

# Bayesian Logistic Regression Model

Consider the Bayesian logistic regression model:

$$\mathbf{y} \mid \boldsymbol{\beta} \sim \text{Bernoulli}\left\{\left[1 + \exp\{-\mathbf{X}\boldsymbol{\beta}\}\right]^{-1}\right\}$$
$$\boldsymbol{\beta} \sim N(\mathbf{0}, \sigma_{\boldsymbol{\beta}}^2 \mathbf{I})$$

The factor graph is



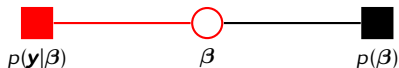
Gaussian Prior Fragment

Logistic Likelihood Fragment (Nolan and Wand, 2017)

The message takes the form:

$$m_{p(\mathbf{y}|\boldsymbol{\beta}) \rightarrow \boldsymbol{\beta}}(\boldsymbol{\beta}) \leftarrow \exp\left\{\mathbf{y}^T \mathbf{X} \boldsymbol{\beta} - \mathbf{1}^T \mathbb{E}_q \log[1 + \exp(\mathbf{X} \boldsymbol{\beta})]\right\}$$

The messages that are typically passed to  $\boldsymbol{\beta}$  from other fragments are Gaussian density messages.

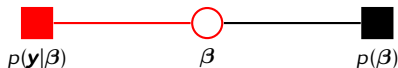


A fixed point iterative scheme for  $\eta_{p(\mathbf{y}|\beta) \rightarrow \beta}$  (Wand, 2014):

$$\begin{aligned} \boldsymbol{\mu} &= \mathbf{X} \mathbb{E}_q(\beta), \quad \boldsymbol{\sigma}^2 = \text{diagonal} \{ \mathbf{X} \text{Cov}_q(\beta) \mathbf{X}^T \} \\ \eta_{p(\mathbf{y}|\beta) \rightarrow \beta} &= \begin{bmatrix} \mathbf{X}^T \left\{ \mathbf{y} - \mathcal{B}_0(\boldsymbol{\mu}, \boldsymbol{\sigma}^2) + \mathcal{B}_1(\boldsymbol{\mu}, \boldsymbol{\sigma}^2) \odot \frac{\boldsymbol{\mu}}{\boldsymbol{\sigma}} \right\} \\ -\frac{1}{2} \text{vec} \{ \mathbf{X}^T \text{diag}(\omega_2) \mathbf{X} \} \end{bmatrix} \end{aligned}$$

where

$$\mathcal{B}_r(\boldsymbol{\mu}, \boldsymbol{\sigma}^2) \equiv \int_{-\infty}^{\infty} x^r \frac{d}{d\boldsymbol{\mu}} \log \{ 1 + \exp(\boldsymbol{\mu} + \boldsymbol{\sigma} x) \} \phi(x) dx, \quad \text{for } r = 0, 1$$



A fixed point iterative scheme for  $\eta_{p(\mathbf{y}|\beta) \rightarrow \beta}$  (Wand, 2014):

$$\begin{aligned} \boldsymbol{\mu} &= \mathbf{X} \mathbb{E}_q(\beta), \quad \boldsymbol{\sigma}^2 = \text{diagonal} \{ \mathbf{X} \text{Cov}_q(\beta) \mathbf{X}^T \} \\ \eta_{p(\mathbf{y}|\beta) \rightarrow \beta} &= \begin{bmatrix} \mathbf{X}^T \left\{ \mathbf{y} - \mathcal{B}_0(\boldsymbol{\mu}, \boldsymbol{\sigma}^2) + \mathcal{B}_1(\boldsymbol{\mu}, \boldsymbol{\sigma}^2) \odot \frac{\boldsymbol{\mu}}{\boldsymbol{\sigma}} \right\} \\ -\frac{1}{2} \text{vec} \{ \mathbf{X}^T \text{diag}(\omega_2) \mathbf{X} \} \end{bmatrix} \end{aligned}$$

where

$$\mathcal{B}_r(\boldsymbol{\mu}, \boldsymbol{\sigma}^2) \equiv \int_{-\infty}^{\infty} x^r \frac{d}{d\boldsymbol{\mu}} \log \{ 1 + \exp(\boldsymbol{\mu} + \boldsymbol{\sigma} x) \} \phi(x) dx, \quad \text{for } r = 0, 1$$

A fixed point iterative scheme for  $q(\beta)$  - Wand (2014):

However, we come across the following integrals:

$$B_r(\mu, \sigma) \equiv \int_{-\infty}^{\infty} x^r [1 + \exp(\mu + \sigma x)]^{-1} \phi(x) dx, \quad \text{for } r = 0, 1$$

↓ Monahan and Stefanski (1992)

$$B_r(\mu, \sigma) \approx \int_{-\infty}^{\infty} x^r \sum_{i=1}^k p_i \Phi\left(s_i \frac{x - \mu}{\sigma}\right) \phi(x) dx, \quad \text{for } r = 0, 1$$

A fixed point iterative scheme for  $q(\beta)$  - Wand (2014):

However, we come across the following integrals:

$$B_r(\mu, \sigma) \equiv \int_{-\infty}^{\infty} x^r \{1 + \exp(\mu + \sigma x)\}^{-1} \phi(x) dx, \quad \text{for } r = 0, 1$$

↓ Monahan and Stefanski (1992)

$$B_r(\mu, \sigma) \approx \int_{-\infty}^{\infty} x^r \sum_{i=1}^k p_i \Phi\left(s_i \frac{x - \mu}{\sigma}\right) \phi(x) dx, \quad \text{for } r = 0, 1$$



A fixed point iterative scheme for  $q(\beta)$  - Wand (2014):

However, we come across the following integrals:

$$B_r(\mu, \sigma) \equiv \int_{-\infty}^{\infty} x^r \{1 + \exp(\mu + \sigma x)\}^{-1} \phi(x) dx, \quad \text{for } r = 0, 1$$

⇓ Monahan and Stefanski (1992)

$$B_r(\mu, \sigma) \approx \int_{-\infty}^{\infty} x^r \sum_{i=1}^k p_i \Phi\left(s_i \frac{x - \mu}{\sigma}\right) \phi(x) dx, \quad \text{for } r = 0, 1$$

A fixed point iterative scheme for  $q(\beta)$  - Wand (2014):

However, we come across the following integrals:

$$B_r(\mu, \sigma) \equiv \int_{-\infty}^{\infty} x^r \{1 + \exp(\mu + \sigma x)\}^{-1} \phi(x) dx, \quad \text{for } r = 0, 1$$

↓ Monahan and Stefanski (1992)

$$B_r(\mu, \sigma) \approx \int_{-\infty}^{\infty} x^r \sum_{i=1}^k p_i \Phi\left(s_i \frac{x - \mu}{\sigma}\right) \phi(x) dx, \quad \text{for } r = 0, 1$$

# Normal Mixture Approximation

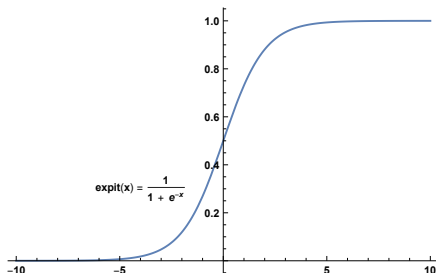
A fixed point iterative scheme for  $q(\beta)$  - Wand (2014):

However, we come across the following integrals:

$$B_r(\mu, \sigma) \equiv \int_{-\infty}^{\infty} x^r \{1 + \exp(\mu + \sigma x)\}^{-1} \phi(x) dx, \quad \text{for } r = 0, 1$$

⇓ Monahan and Stefanski (1992)

$$B_r(\mu, \sigma) \approx \int_{-\infty}^{\infty} x^r \sum_{i=1}^k p_i \Phi\left(s_i \frac{x - \mu}{\sigma}\right) \phi(x) dx, \quad \text{for } r = 0, 1$$



# Normal Mixture Approximation

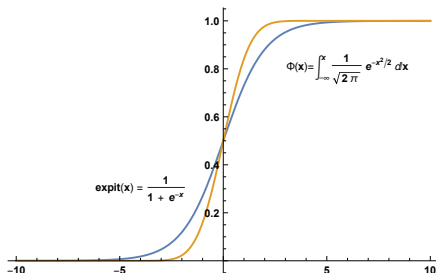
A fixed point iterative scheme for  $q(\beta)$  - Wand (2014):

However, we come across the following integrals:

$$B_r(\mu, \sigma) \equiv \int_{-\infty}^{\infty} x^r \{1 + \exp(\mu + \sigma x)\}^{-1} \phi(x) dx, \quad \text{for } r = 0, 1$$

⇓ Monahan and Stefanski (1992)

$$B_r(\mu, \sigma) \approx \int_{-\infty}^{\infty} x^r \sum_{i=1}^k p_i \Phi\left(s_i \frac{x - \mu}{\sigma}\right) \phi(x) dx, \quad \text{for } r = 0, 1$$



# Normal Mixture Approximation

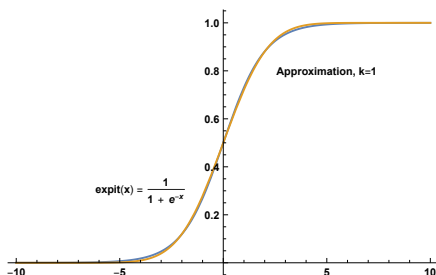
A fixed point iterative scheme for  $q(\beta)$  - Wand (2014):

However, we come across the following integrals:

$$B_r(\mu, \sigma) \equiv \int_{-\infty}^{\infty} x^r \{1 + \exp(\mu + \sigma x)\}^{-1} \phi(x) dx, \quad \text{for } r = 0, 1$$

⇓ Monahan and Stefanski (1992)

$$B_r(\mu, \sigma) \approx \int_{-\infty}^{\infty} x^r \sum_{i=1}^k p_i \Phi\left(s_i \frac{x-\mu}{\sigma}\right) \phi(x) dx, \quad \text{for } r = 0, 1$$



$p$	$s$
1.00000	0.58763

# Normal Mixture Approximation

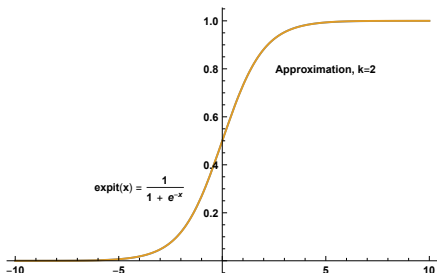
A fixed point iterative scheme for  $q(\beta)$  - Wand (2014):

However, we come across the following integrals:

$$B_r(\mu, \sigma) \equiv \int_{-\infty}^{\infty} x^r \{1 + \exp(\mu + \sigma x)\}^{-1} \phi(x) dx, \quad \text{for } r = 0, 1$$

⇓ Monahan and Stefanski (1992)

$$B_r(\mu, \sigma) \approx \int_{-\infty}^{\infty} x^r \sum_{i=1}^k p_i \Phi\left(s_i \frac{x - \mu}{\sigma}\right) \phi(x) dx, \quad \text{for } r = 0, 1$$



$p$	$s$
0.56442	0.76862
0.43557	0.43525

# Normal Mixture Approximation

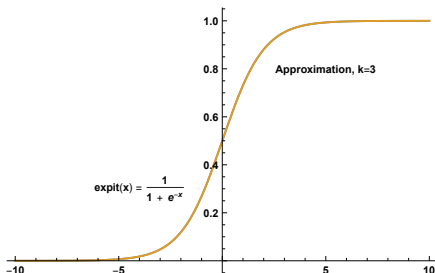
A fixed point iterative scheme for  $q(\beta)$  - Wand (2014):

However, we come across the following integrals:

$$B_r(\mu, \sigma) \equiv \int_{-\infty}^{\infty} x^r \{1 + \exp(\mu + \sigma x)\}^{-1} \phi(x) dx, \quad \text{for } r = 0, 1$$

⇓ Monahan and Stefanski (1992)

$$B_r(\mu, \sigma) \approx \int_{-\infty}^{\infty} x^r \sum_{i=1}^k p_i \Phi\left(s_i \frac{x-\mu}{\sigma}\right) \phi(x) dx, \quad \text{for } r = 0, 1$$



$p$	$s$
0.25220	0.90793
0.58522	0.57778
0.16257	0.36403

# Normal Mixture Approximation

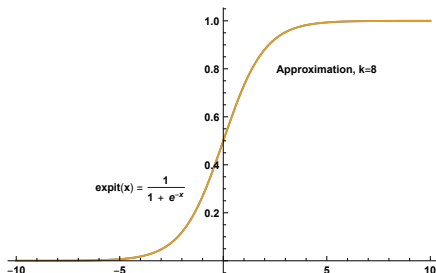
A fixed point iterative scheme for  $q(\beta)$  - Wand (2014):

However, we come across the following integrals:

$$B_r(\mu, \sigma) \equiv \int_{-\infty}^{\infty} x^r \{1 + \exp(\mu + \sigma x)\}^{-1} \phi(x) dx, \quad \text{for } r = 0, 1$$

⇓ Monahan and Stefanski (1992)

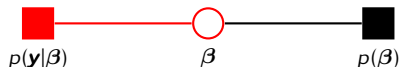
$$B_r(\mu, \sigma) \approx \int_{-\infty}^{\infty} x^r \sum_{i=1}^k p_i \Phi\left(s_i \frac{x - \mu}{\sigma}\right) \phi(x) dx, \quad \text{for } r = 0, 1$$



$p$	$s$
0.00324	1.36534
0.05151	1.05952
0.19507	0.83079
0.31556	0.65073
0.27414	0.50813
0.13107	0.39631
0.02791	0.30890
0.00144	0.23821



# Logistic Likelihood Fragment



$$\mathbf{y}|\beta \sim \text{Bernoulli}\left[\frac{1}{1 + \exp(-\mathbf{X}\beta)}\right]$$

Inputs:

$$\eta_{q^*}(\beta)$$

Updates:

$$\text{Cov}_q\{(\beta)\} = -\frac{1}{2}[\text{vec}^{-1}(\eta_{q^*}(\beta)_2)]^{-1}, \quad \mathbb{E}_q\{(\beta)\} = \text{Cov}_q\{(\beta)\}(\eta_{q^*}(\beta)_1),$$

$$\mu = \mathbf{X} \mathbb{E}_q(\beta), \quad \sigma^2 = \text{diagonal}\{\mathbf{X} \text{Cov}_q(\beta) \mathbf{X}^T\},$$

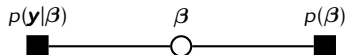
$$\eta_{p(\mathbf{y}|\beta) \rightarrow \beta} = \begin{bmatrix} \mathbf{X}^T \left\{ \mathbf{y} - \mathcal{B}_0(\mu, \sigma^2) + \mathcal{B}_1(\mu, \sigma^2) \odot \frac{\mu}{\sigma} \right\} \\ -\frac{1}{2} \text{vec}\left\{ \mathbf{X}^T \text{diag}(\omega_2) \mathbf{X} \right\} \end{bmatrix}$$

Outputs:

$$\eta_{p(\mathbf{y}|\beta) \rightarrow \beta}$$

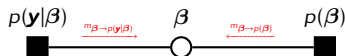
## VMP for the Bayesian logistic regression model

1. Initialise all messages from stochastic nodes to factors
2. Cycle:
  - (i) Update all messages from factors to stochastic nodes
  - (ii) Update all messages from stochastic nodes to factors
  - (iii) Update all optimal posterior density functions
3. Stop:  $D_{\text{KL}} \{q(\beta) || p(\beta|\mathbf{y})\}$  converges.



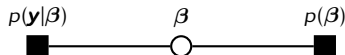
## VMP for the Bayesian logistic regression model

1. Initialise all messages from stochastic nodes to factors
2. Cycle:
  - (i) Update all messages from factors to stochastic nodes
  - (ii) Update all messages from stochastic nodes to factors
  - (iii) Update all optimal posterior density functions
3. Stop:  $D_{\text{KL}} \{q(\beta) || p(\beta|y)\}$  converges.



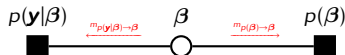
## VMP for the Bayesian logistic regression model

1. Initialise all messages from stochastic nodes to factors
2. Cycle:
  - (i) Update all messages from factors to stochastic nodes
  - (ii) Update all messages from stochastic nodes to factors
  - (iii) Update all optimal posterior density functions
3. Stop:  $D_{KL} \{q(\beta) || p(\beta|y)\}$  converges.



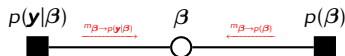
## VMP for the Bayesian logistic regression model

1. Initialise all messages from stochastic nodes to factors
2. Cycle:
  - (i) Update all messages from factors to stochastic nodes
  - (ii) Update all messages from stochastic nodes to factors
  - (iii) Update all optimal posterior density functions
3. Stop:  $D_{KL} \{q(\beta) || p(\beta|y)\}$  converges.



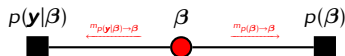
## VMP for the Bayesian logistic regression model

1. Initialise all messages from stochastic nodes to factors
2. Cycle:
  - (i) Update all messages from factors to stochastic nodes
  - (ii) Update all messages from stochastic nodes to factors
  - (iii) Update all optimal posterior density functions
3. Stop:  $D_{\text{KL}} \{q(\beta) \| p(\beta | \mathbf{y})\}$  converges.



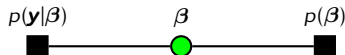
## VMP for the Bayesian logistic regression model

1. Initialise all messages from stochastic nodes to factors
2. Cycle:
  - (i) Update all messages from factors to stochastic nodes
  - (ii) Update all messages from stochastic nodes to factors
  - (iii) Update all optimal posterior density functions
3. Stop:  $D_{KL} \{q(\beta) || p(\beta|y)\}$  converges.



## VMP for the Bayesian logistic regression model

1. Initialise all messages from stochastic nodes to factors
2. Cycle:
  - (i) Update all messages from factors to stochastic nodes
  - (ii) Update all messages from stochastic nodes to factors
  - (iii) Update all optimal posterior density functions
3. Stop:  $D_{\text{KL}} \{q(\beta) || p(\beta|\mathbf{y})\}$  converges.

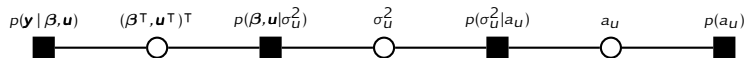




# Bayesian Logistic Semiparametric Regression

$$\mathbf{y} \mid \boldsymbol{\beta}, \mathbf{u} \sim \text{Bernoulli} \left\{ [1 + \exp \{ -(\mathbf{X}\boldsymbol{\beta} + \mathbf{X}\mathbf{u}) \}]^{-1} \right\}, \quad \begin{bmatrix} \boldsymbol{\beta} \\ \mathbf{u} \end{bmatrix} \mid \sigma_u^2 \sim \mathcal{N} \left( \begin{bmatrix} \mathbf{0} \\ \mathbf{0} \end{bmatrix}, \begin{bmatrix} \sigma_\beta^2 \mathbf{I}_2 & \mathbf{0}^T \\ \mathbf{0} & \sigma_u^2 \mathbf{I}_m \end{bmatrix} \right)$$

$$\sigma_u^2 \mid a_u \sim \text{Inverse-}\chi^2(1, 1/a_u), \quad a_u \sim \text{Inverse-}\chi^2(1, 1/A^2)$$



Logistic likelihood fragment

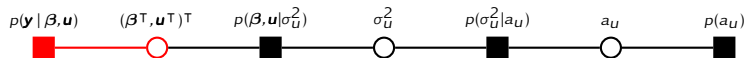
Gaussian penalization fragment (Wand, 2017)

Iterated inverse Wishart fragment

Inverse Wishart prior fragment

# Bayesian Logistic Semiparametric Regression

$$\mathbf{y} \mid \boldsymbol{\beta}, \mathbf{u} \sim \text{Bernoulli} \left\{ [1 + \exp \{-(\mathbf{X}\boldsymbol{\beta} + \mathbf{X}\mathbf{u})\}]^{-1} \right\}, \quad \begin{bmatrix} \boldsymbol{\beta} \\ \mathbf{u} \end{bmatrix} \mid \sigma_u^2 \sim \mathcal{N} \left( \begin{bmatrix} \mathbf{0} \\ \mathbf{0} \end{bmatrix}, \begin{bmatrix} \sigma_\beta^2 \mathbf{I}_2 & \mathbf{0}^T \\ \mathbf{0} & \sigma_u^2 \mathbf{I}_m \end{bmatrix} \right)$$
$$\sigma_u^2 \mid a_u \sim \text{Inverse-}\chi^2(1, 1/a_u), \quad a_u \sim \text{Inverse-}\chi^2(1, 1/A^2)$$



Logistic likelihood fragment

Gaussian penalization fragment (Wand, 2017)

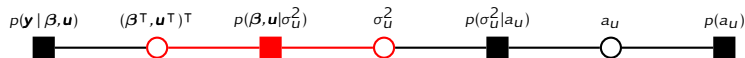
Iterated inverse Wishart fragment

Inverse Wishart prior fragment

# Bayesian Logistic Semiparametric Regression

$$\mathbf{y} \mid \boldsymbol{\beta}, \mathbf{u} \sim \text{Bernoulli} \left\{ \left[ 1 + \exp \{ -(\mathbf{X}\boldsymbol{\beta} + \mathbf{X}\mathbf{u}) \} \right]^{-1} \right\}, \quad \begin{bmatrix} \boldsymbol{\beta} \\ \mathbf{u} \end{bmatrix} \mid \sigma_u^2 \sim \text{N} \left( \begin{bmatrix} \mathbf{0} \\ \mathbf{0} \end{bmatrix}, \begin{bmatrix} \sigma_\beta^2 \mathbf{I}_2 & \mathbf{0}^T \\ \mathbf{0} & \sigma_u^2 \mathbf{I}_m \end{bmatrix} \right)$$

$$\sigma_u^2 \mid a_u \sim \text{Inverse-}\chi^2(1, 1/a_u), \quad a_u \sim \text{Inverse-}\chi^2(1, 1/A^2)$$



Logistic likelihood fragment

Gaussian penalization fragment (Wand, 2017)

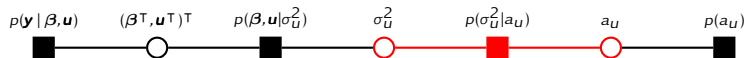
Iterated inverse Wishart fragment

Inverse Wishart prior fragment

# Bayesian Logistic Semiparametric Regression

$$\mathbf{y} \mid \boldsymbol{\beta}, \mathbf{u} \sim \text{Bernoulli} \left\{ [1 + \exp \{ -(\mathbf{X}\boldsymbol{\beta} + \mathbf{X}\mathbf{u}) \}]^{-1} \right\}, \quad \begin{bmatrix} \boldsymbol{\beta} \\ \mathbf{u} \end{bmatrix} \mid \sigma_u^2 \sim \mathcal{N} \left( \begin{bmatrix} \mathbf{0} \\ \mathbf{0} \end{bmatrix}, \begin{bmatrix} \sigma_\beta^2 \mathbf{I}_2 & \mathbf{0}^T \\ \mathbf{0} & \sigma_u^2 \mathbf{I}_m \end{bmatrix} \right)$$

$$\sigma_u^2 \mid a_u \sim \text{Inverse-}\chi^2(1, 1/a_u), \quad a_u \sim \text{Inverse-}\chi^2(1, 1/A^2)$$



Logistic likelihood fragment

Gaussian penalization fragment (Wand, 2017)

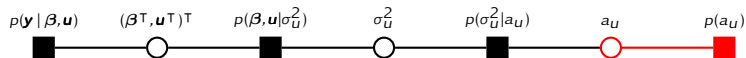
Iterated inverse Wishart fragment

Inverse Wishart prior fragment

# Bayesian Logistic Semiparametric Regression

$$\mathbf{y} \mid \boldsymbol{\beta}, \mathbf{u} \sim \text{Bernoulli} \left\{ [1 + \exp \{ -(\mathbf{X}\boldsymbol{\beta} + \mathbf{X}\mathbf{u}) \}]^{-1} \right\}, \quad \begin{bmatrix} \boldsymbol{\beta} \\ \mathbf{u} \end{bmatrix} \mid \sigma_u^2 \sim \mathcal{N} \left( \begin{bmatrix} \mathbf{0} \\ \mathbf{0} \end{bmatrix}, \begin{bmatrix} \sigma_\beta^2 \mathbf{I}_2 & \mathbf{0}^T \\ \mathbf{0} & \sigma_u^2 \mathbf{I}_m \end{bmatrix} \right)$$

$$\sigma_u^2 \mid a_u \sim \text{Inverse-}\chi^2(1, 1/a_u), \quad a_u \sim \text{Inverse-}\chi^2(1, 1/A^2)$$



Logistic likelihood fragment

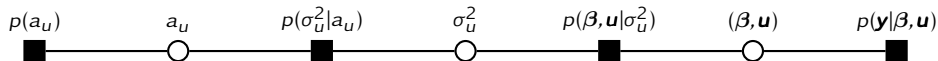
Gaussian penalization fragment (Wand, 2017)

Iterated inverse Wishart fragment

Inverse Wishart prior fragment

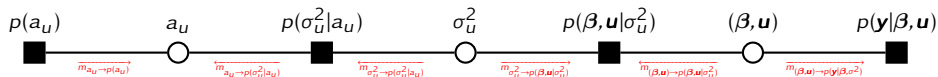
## VMP for the Bayesian logistic semiparametric regression model

1. Initialise all messages from stochastic nodes to factors
2. Cycle:
  - (i) Update all messages from factors to stochastic nodes
  - (ii) Update all messages from stochastic nodes to factors
  - (iii) Update all optimal posterior density functions
3. Stop:  $D_{\text{KL}} \{q(\beta, \mathbf{u}, \sigma_u^2, a_u) \| p(\beta, \mathbf{u}, \sigma_u^2, a_u | \mathbf{y})\}$  converges.



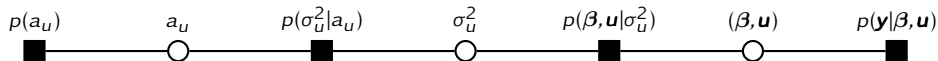
## VMP for the Bayesian logistic semiparametric regression model

1. Initialise all messages from stochastic nodes to factors
2. Cycle:
  - (i) Update all messages from factors to stochastic nodes
  - (ii) Update all messages from stochastic nodes to factors
  - (iii) Update all optimal posterior density functions
3. Stop:  $D_{KL} \{q(\beta, \mathbf{u}, \sigma_u^2, a_u) \| p(\beta, \mathbf{u}, \sigma_u^2, a_u | \mathbf{y})\}$  converges.



## VMP for the Bayesian logistic semiparametric regression model

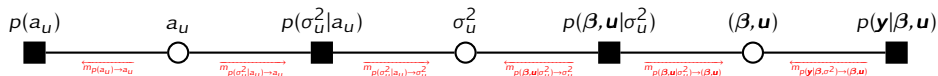
1. Initialise all messages from stochastic nodes to factors
2. Cycle:
  - (i) Update all messages from factors to stochastic nodes
  - (ii) Update all messages from stochastic nodes to factors
  - (iii) Update all optimal posterior density functions
3. Stop:  $D_{\text{KL}} \{q(\beta, \mathbf{u}, \sigma_u^2, a_u) \| p(\beta, \mathbf{u}, \sigma_u^2, a_u | \mathbf{y})\}$  converges.





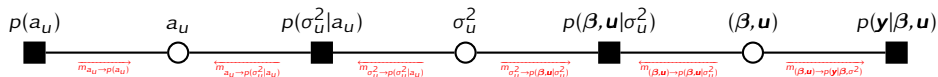
## VMP for the Bayesian logistic semiparametric regression model

1. Initialise all messages from stochastic nodes to factors
2. Cycle:
  - (i) Update all messages from factors to stochastic nodes
  - (ii) Update all messages from stochastic nodes to factors
  - (iii) Update all optimal posterior density functions
3. Stop:  $D_{\text{KL}} \{q(\beta, \mathbf{u}, \sigma_U^2, a_U) \| p(\beta, \mathbf{u}, \sigma_U^2, a_U | \mathbf{y})\}$  converges.



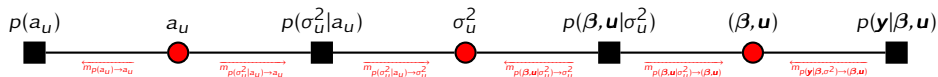
## VMP for the Bayesian logistic semiparametric regression model

1. Initialise all messages from stochastic nodes to factors
2. Cycle:
  - (i) Update all messages from factors to stochastic nodes
  - (ii) Update all messages from stochastic nodes to factors
  - (iii) Update all optimal posterior density functions
3. Stop:  $D_{KL} \{q(\beta, \mathbf{u}, \sigma_u^2, a_u) \| p(\beta, \mathbf{u}, \sigma_u^2, a_u | \mathbf{y})\}$  converges.



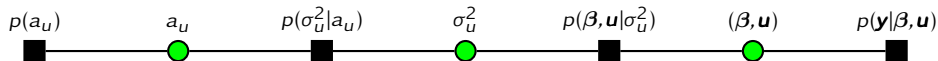
## VMP for the Bayesian logistic semiparametric regression model

1. Initialise all messages from stochastic nodes to factors
2. Cycle:
  - (i) Update all messages from factors to stochastic nodes
  - (ii) Update all messages from stochastic nodes to factors
  - (iii) Update all optimal posterior density functions
3. Stop:  $D_{\text{KL}} \{q(\beta, \mathbf{u}, \sigma_U^2, \mathbf{a}_U) \| p(\beta, \mathbf{u}, \sigma_U^2, \mathbf{a}_U | \mathbf{y})\}$  converges.

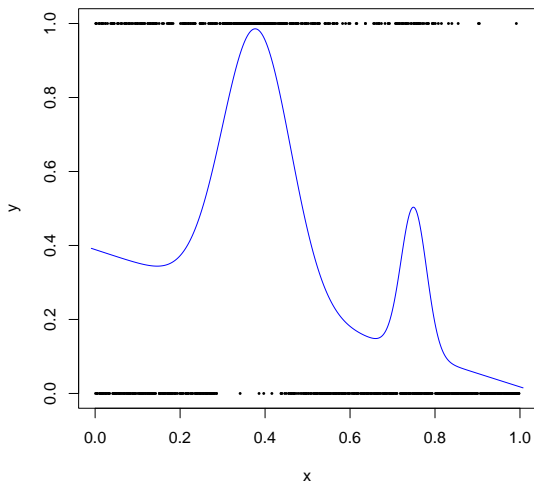


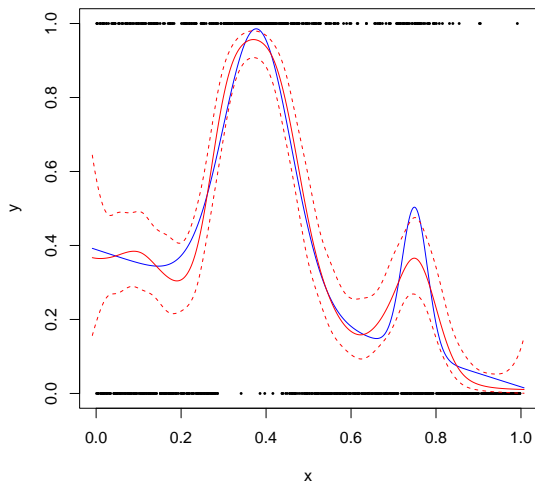
## VMP for the Bayesian logistic semiparametric regression model

1. Initialise all messages from stochastic nodes to factors
2. Cycle:
  - (i) Update all messages from factors to stochastic nodes
  - (ii) Update all messages from stochastic nodes to factors
  - (iii) Update all optimal posterior density functions
3. Stop:  $D_{\text{KL}}\{q(\boldsymbol{\beta}, \mathbf{u}, \sigma_u^2, a_u) \| p(\boldsymbol{\beta}, \mathbf{u}, \sigma_u^2, a_u | \mathbf{y})\}$  converges.



## Binary Response Data





- Bishop, C. M. (2006). *Pattern Recognition and Machine Learning*. Springer, New York.
- Chamberlain, S., Anderson, B., Salmon, M., Erickson, A., Potter, N., Stachelek, J., Simmons, A., Ram, K., Edmund, H., and rOpenSci (2021). 'NOAA' weather data from r. R package version 1.3.0.
- Gelman, A. (2006). Prior distributions for variance parameters in hierarchical models (comment on article by Browne and Draper). *Bayesian Analysis*, 1:515–534.
- Goldsmith, J., Zippunnikov, V., and Schrack, J. (2015). Generalized multilevel function-on-scalar regression and principal component analysis. *Biometrics*, 71:344–353.
- Happ, C. and Greven, S. (2018). Multivariate functional principal component analysis for data observed on different (dimensional) domains. *Journal of the American Statistical Association*, 113:649–659.
- Huang, A. and Wand, M. P. (2013). Simple marginally noninformative prior distributions for covariance matrices. *Bayesian Analysis*, 8:439–452.
- Lee, C. Y. Y. and Wand, M. P. (2016). Streamlined mean field variational Bayes for longitudinal and multilevel data analysis. *Biometrical Journal*, 58:868–895.
- Maestrini, L. and Wand, M. P. (2020). The Inverse G-Wishart distribution and variational message passing. *arXiv e-prints*, page arXiv:2005.09876.
- Menictas, M. and Wand, M. P. (2013). Variational inference for marginal longitudinal semiparametric regression. *Stat*, 2:61–71.
- Minka, T. (2005). Divergence measures and message passing. Technical report, Microsoft Research Ltd., Cambridge, UK.
- Monahan, J. F. and Stefanski, L. A. (1992). Normal scale mixture approximations to  $F^*(z)$  and computation of the logistic-normal integral. In *Handbook of the Logistic Distribution*, pages 529–540. CRC Press.
- Nolan, T. H., Goldsmith, J., and Ruppert, D. (2021). Bayesian functional principal components analysis via variational message passing. *arXiv e-prints*, page arXiv:2104.00645.
- Nolan, T. H. and Wand, M. P. (2017). Accurate logistic variational message passing: Algebraic and numerical details. *Stat*, 6:102–112.
- Nolan, T. H. and Wand, M. P. (2020). Streamlined solutions to multilevel sparse matrix problems. *ANZIAM Journal*, 62:18–41.
- Ramsay, J. O. and Silverman, B. W. (2005). *Functional Data Analysis*. Springer, New York.
- Wand, M. P. (2014). Fully simplified multivariate normal updates in non-conjugate variational message passing. *Journal of Machine Learning Research*, 15:1351–1369.
- Wand, M. P. (2017). Fast approximate inference for arbitrarily large semiparametric regression models via message passing (with discussion). *Journal of the American Statistical Association*, 112:137–168.
- Wand, M. P. and Ormerod, J. T. (2008). On semiparametric regression with O'Sullivan penalized splines. *Australian & New Zealand Journal of Statistics*, 50:179–198.
- Wang, J. L., Chiou, J. M., and Müller, H. G. (2016). Functional data analysis. *Annual Review of Statistics and Its Applications*, 3:257–295.