

# SDS 383C: Statistical Modeling I

## Fall 2022, Module IX

**Abhra Sarkar**

Department of Statistics and Data Sciences  
The University of Texas at Austin

"All models are wrong, but some are useful."- George E. P. Box

- Models linear in parameters

$$y_i = \mathbf{x}_i^T \boldsymbol{\beta} + \epsilon_i, \quad \epsilon_i \stackrel{iid}{\sim} f_\epsilon, \quad i = 1, \dots, n, \quad \text{with } \mathbb{E}_{f_\epsilon}(\epsilon) = 0.$$

- ▶  $y_i = \beta_0 + x_i \beta_1 + \epsilon_i, \quad \epsilon_i \stackrel{iid}{\sim} \text{Normal}(0, \sigma^2)$
- ▶  $y_i = \beta_0 + x_i \beta_1 + \epsilon_i, \quad \epsilon_i \stackrel{iid}{\sim} \text{Laplace}(0, b)$
- ▶  $y_i = \beta_0 + x_i \beta_1 + x_i^2 \beta_2 + x_i^3 \beta_3 + \epsilon_i, \quad \epsilon_i \stackrel{iid}{\sim} \text{Normal}(0, \sigma^2)$
- ▶  $y_i = \beta_0 + x_{i,1} \beta_1 + x_{i,2} \beta_2 + \dots + x_{i,p} \beta_p + \epsilon_i, \quad \epsilon_i \stackrel{iid}{\sim} \text{Normal}(0, \sigma^2)$
- Matrix-vector notation

$$\mathbf{y}^{n \times 1} = \mathbf{X}^{n \times p} \boldsymbol{\beta}^{p \times 1} + \mathbf{e}^{n \times 1}, \quad \mathbb{E}(\mathbf{e}) = \mathbf{0}, \quad \text{var}(\mathbf{e}) = \sigma^2 \mathbf{I}_n$$

$$\begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{bmatrix} = \begin{bmatrix} x_{1,1} & x_{1,2} & \dots & x_{1,p} \\ x_{2,1} & x_{2,2} & \dots & x_{2,p} \\ \vdots & \vdots & \ddots & \vdots \\ x_{n,1} & x_{n,2} & \dots & x_{n,p} \end{bmatrix} \begin{bmatrix} \beta_1 \\ \beta_2 \\ \vdots \\ \beta_p \end{bmatrix} + \begin{bmatrix} \epsilon_1 \\ \epsilon_2 \\ \vdots \\ \epsilon_n \end{bmatrix}$$

- Models linear in parameters

$$y_i = \mathbf{x}_i^T \boldsymbol{\beta} + \epsilon_i, \quad \epsilon_i \stackrel{iid}{\sim} f_\epsilon, \quad i = 1, \dots, n, \quad \text{with } \mathbb{E}_{f_\epsilon}(\epsilon) = 0.$$

►  $y_i = \beta_0 + x_i \beta_1 + \epsilon_i, \quad \epsilon_i \stackrel{iid}{\sim} \text{Normal}(0, \sigma^2)$

►  $y_i = \beta_0 + x_i \beta_1 + \epsilon_i, \quad \epsilon_i \stackrel{iid}{\sim} \text{Laplace}(0, b)$

►  $y_i = \beta_0 + x_i \beta_1 + x_i^2 \beta_2 + x_i^3 \beta_3 + \epsilon_i, \quad \epsilon_i \stackrel{iid}{\sim} \text{Normal}(0, \sigma^2)$

►  $y_i = \beta_0 + x_{i,1} \beta_1 + x_{i,2} \beta_2 + \dots + x_{i,p} \beta_p + \epsilon_i, \quad \epsilon_i \stackrel{iid}{\sim} \text{Normal}(0, \sigma^2)$

- Matrix-vector notation

$$\mathbf{y}^{n \times 1} = \mathbf{X}^{n \times p} \boldsymbol{\beta}^{p \times 1} + \mathbf{e}^{n \times 1}, \quad \mathbb{E}(\mathbf{e}) = \mathbf{0}, \quad \text{var}(\mathbf{e}) = \sigma^2 \mathbf{I}_n$$

$$\begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{bmatrix} = \begin{bmatrix} x_{1,1} & x_{1,2} & \dots & x_{1,p} \\ x_{2,1} & x_{2,2} & \dots & x_{2,p} \\ \vdots & \vdots & \ddots & \vdots \\ x_{n,1} & x_{n,2} & \dots & x_{n,p} \end{bmatrix} \begin{bmatrix} \beta_1 \\ \beta_2 \\ \vdots \\ \beta_p \end{bmatrix} + \begin{bmatrix} \epsilon_1 \\ \epsilon_2 \\ \vdots \\ \epsilon_n \end{bmatrix}$$

- Models linear in parameters

$$y_i = \mathbf{x}_i^T \boldsymbol{\beta} + \epsilon_i, \quad \epsilon_i \stackrel{iid}{\sim} f_\epsilon, \quad i = 1, \dots, n, \quad \text{with } \mathbb{E}_{f_\epsilon}(\epsilon) = 0.$$

►  $y_i = \beta_0 + x_i \beta_1 + \epsilon_i, \quad \epsilon_i \stackrel{iid}{\sim} \text{Normal}(0, \sigma^2)$

►  $y_i = \beta_0 + x_i \beta_1 + \epsilon_i, \quad \epsilon_i \stackrel{iid}{\sim} \text{Laplace}(0, b)$

►  $y_i = \beta_0 + x_i \beta_1 + x_i^2 \beta_2 + x_i^3 \beta_3 + \epsilon_i, \quad \epsilon_i \stackrel{iid}{\sim} \text{Normal}(0, \sigma^2)$

►  $y_i = \beta_0 + x_{i,1} \beta_1 + x_{i,2} \beta_2 + \dots + x_{i,p} \beta_p + \epsilon_i, \quad \epsilon_i \stackrel{iid}{\sim} \text{Normal}(0, \sigma^2)$

- Matrix-vector notation

$$\mathbf{y}^{n \times 1} = \mathbf{X}^{n \times p} \boldsymbol{\beta}^{p \times 1} + \mathbf{e}^{n \times 1}, \quad \mathbb{E}(\mathbf{e}) = \mathbf{0}, \quad \text{var}(\mathbf{e}) = \sigma^2 \mathbf{I}_n$$

$$\begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{bmatrix} = \begin{bmatrix} x_{1,1} & x_{1,2} & \dots & x_{1,p} \\ x_{2,1} & x_{2,2} & \dots & x_{2,p} \\ \vdots & \vdots & \ddots & \vdots \\ x_{n,1} & x_{n,2} & \dots & x_{n,p} \end{bmatrix} \begin{bmatrix} \beta_1 \\ \beta_2 \\ \vdots \\ \beta_p \end{bmatrix} + \begin{bmatrix} \epsilon_1 \\ \epsilon_2 \\ \vdots \\ \epsilon_n \end{bmatrix}$$

- Models linear in parameters

$$y_i = \mathbf{x}_i^T \boldsymbol{\beta} + \epsilon_i, \quad \epsilon_i \stackrel{iid}{\sim} f_\epsilon, \quad i = 1, \dots, n, \quad \text{with } \mathbb{E}_{f_\epsilon}(\epsilon) = 0.$$

- $y_i = \beta_0 + x_i \beta_1 + \epsilon_i, \quad \epsilon_i \stackrel{iid}{\sim} \text{Normal}(0, \sigma^2)$
- $y_i = \beta_0 + x_i \beta_1 + \epsilon_i, \quad \epsilon_i \stackrel{iid}{\sim} \text{Laplace}(0, b)$
- $y_i = \beta_0 + x_i \beta_1 + x_i^2 \beta_2 + x_i^3 \beta_3 + \epsilon_i, \quad \epsilon_i \stackrel{iid}{\sim} \text{Normal}(0, \sigma^2)$
- $y_i = \beta_0 + x_{i,1} \beta_1 + x_{i,2} \beta_2 + \dots + x_{i,p} \beta_p + \epsilon_i, \quad \epsilon_i \stackrel{iid}{\sim} \text{Normal}(0, \sigma^2)$
- Matrix-vector notation

$$\mathbf{y}^{n \times 1} = \mathbf{X}^{n \times p} \boldsymbol{\beta}^{p \times 1} + \mathbf{e}^{n \times 1}, \quad \mathbb{E}(\mathbf{e}) = \mathbf{0}, \quad \text{var}(\mathbf{e}) = \sigma^2 \mathbf{I}_n$$

$$\begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{bmatrix} = \begin{bmatrix} x_{1,1} & x_{1,2} & \dots & x_{1,p} \\ x_{2,1} & x_{2,2} & \dots & x_{2,p} \\ \vdots & \vdots & \ddots & \vdots \\ x_{n,1} & x_{n,2} & \dots & x_{n,p} \end{bmatrix} \begin{bmatrix} \beta_1 \\ \beta_2 \\ \vdots \\ \beta_p \end{bmatrix} + \begin{bmatrix} \epsilon_1 \\ \epsilon_2 \\ \vdots \\ \epsilon_n \end{bmatrix}$$

- Models linear in parameters

$$y_i = \mathbf{x}_i^T \boldsymbol{\beta} + \epsilon_i, \quad \epsilon_i \stackrel{iid}{\sim} f_\epsilon, \quad i = 1, \dots, n, \quad \text{with } \mathbb{E}_{f_\epsilon}(\epsilon) = 0.$$

- $y_i = \beta_0 + x_i \beta_1 + \epsilon_i, \quad \epsilon_i \stackrel{iid}{\sim} \text{Normal}(0, \sigma^2)$
- $y_i = \beta_0 + x_i \beta_1 + \epsilon_i, \quad \epsilon_i \stackrel{iid}{\sim} \text{Laplace}(0, b)$
- $y_i = \beta_0 + x_i \beta_1 + x_i^2 \beta_2 + x_i^3 \beta_3 + \epsilon_i, \quad \epsilon_i \stackrel{iid}{\sim} \text{Normal}(0, \sigma^2)$
- $y_i = \beta_0 + x_{i,1} \beta_1 + x_{i,2} \beta_2 + \dots + x_{i,p} \beta_p + \epsilon_i, \quad \epsilon_i \stackrel{iid}{\sim} \text{Normal}(0, \sigma^2)$

- Matrix-vector notation

$$\mathbf{y}^{n \times 1} = \mathbf{X}^{n \times p} \boldsymbol{\beta}^{p \times 1} + \mathbf{e}^{n \times 1}, \quad \mathbb{E}(\mathbf{e}) = \mathbf{0}, \quad \text{var}(\mathbf{e}) = \sigma^2 \mathbf{I}_n$$
$$\begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{bmatrix} = \begin{bmatrix} x_{1,1} & x_{1,2} & \dots & x_{1,p} \\ x_{2,1} & x_{2,2} & \dots & x_{2,p} \\ \vdots & \vdots & \ddots & \vdots \\ x_{n,1} & x_{n,2} & \dots & x_{n,p} \end{bmatrix} \begin{bmatrix} \beta_1 \\ \beta_2 \\ \vdots \\ \beta_p \end{bmatrix} + \begin{bmatrix} \epsilon_1 \\ \epsilon_2 \\ \vdots \\ \epsilon_n \end{bmatrix}$$

- Models linear in parameters

$$y_i = \mathbf{x}_i^T \boldsymbol{\beta} + \epsilon_i, \quad \epsilon_i \stackrel{iid}{\sim} f_\epsilon, \quad i = 1, \dots, n, \quad \text{with } \mathbb{E}_{f_\epsilon}(\epsilon) = 0.$$

- $y_i = \beta_0 + x_i \beta_1 + \epsilon_i, \quad \epsilon_i \stackrel{iid}{\sim} \text{Normal}(0, \sigma^2)$
- $y_i = \beta_0 + x_i \beta_1 + \epsilon_i, \quad \epsilon_i \stackrel{iid}{\sim} \text{Laplace}(0, b)$
- $y_i = \beta_0 + x_i \beta_1 + x_i^2 \beta_2 + x_i^3 \beta_3 + \epsilon_i, \quad \epsilon_i \stackrel{iid}{\sim} \text{Normal}(0, \sigma^2)$
- $y_i = \beta_0 + x_{i,1} \beta_1 + x_{i,2} \beta_2 + \dots + x_{i,p} \beta_p + \epsilon_i, \quad \epsilon_i \stackrel{iid}{\sim} \text{Normal}(0, \sigma^2)$
- Matrix-vector notation

$$\mathbf{y}^{n \times 1} = \mathbf{X}^{n \times p} \boldsymbol{\beta}^{n \times 1} + \mathbf{e}^{n \times 1}, \quad \mathbb{E}(\mathbf{e}) = \mathbf{0}, \quad \text{var}(\mathbf{e}) = \sigma^2 \mathbf{I}_n$$

$$\begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{bmatrix} = \begin{bmatrix} x_{1,1} & x_{1,2} & \dots & x_{1,p} \\ x_{2,1} & x_{2,2} & \dots & x_{2,p} \\ \vdots & \vdots & \ddots & \vdots \\ x_{n,1} & x_{n,2} & \dots & x_{n,p} \end{bmatrix} \begin{bmatrix} \beta_1 \\ \beta_2 \\ \vdots \\ \beta_p \end{bmatrix} + \begin{bmatrix} \epsilon_1 \\ \epsilon_2 \\ \vdots \\ \epsilon_n \end{bmatrix}$$

► **Laplace Distribution:**

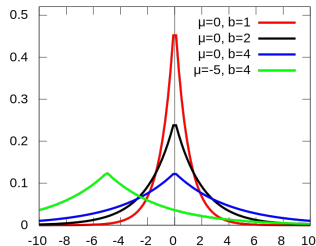
$$y \sim f(y \mid \mu, b) = \frac{1}{2b} \exp\left(-\frac{|y - \mu|}{b}\right).$$

►  $\mathbb{E}(y) = \mu$

►  $\text{var}(y) = 2b^2$

►  $\text{skewness}(y) = 0$

►  $\text{excess kurtosis}(y) = 3$





► **Laplace Distribution:**

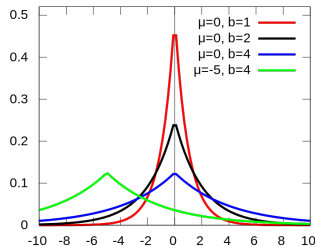
$$y \sim f(y \mid \mu, b) = \frac{1}{2b} \exp\left(-\frac{|y - \mu|}{b}\right).$$

►  $\mathbb{E}(y) = \mu$

►  $\text{var}(y) = 2b^2$

►  $\text{skewness}(y) = 0$

►  $\text{excess kurtosis}(y) = 3$



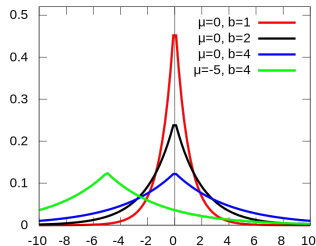
• Laplace as Normal scale mixture

$$\begin{aligned} \frac{a}{2} \exp(-a|z|) &= \int_0^\infty \frac{1}{\sqrt{2\pi s^2}} \exp\left(-\frac{z^2}{2s^2}\right) \frac{a^2}{2} \exp\left(-\frac{a^2 s^2}{2}\right) ds^2 \\ &= \int_0^\infty \text{Normal}(z \mid 0, s^2) \text{Exp}\left(s^2 \mid \frac{a^2}{2}\right) ds^2 \end{aligned}$$

► **Laplace Distribution:**

$$y \sim f(y \mid \mu, b) = \frac{1}{2b} \exp\left(-\frac{|y - \mu|}{b}\right).$$

- $\mathbb{E}(y) = \mu$
- $\text{var}(y) = 2b^2$
- $\text{skewness}(y) = 0$
- $\text{excess kurtosis}(y) = 3$



- Laplace as Normal scale mixture

$$\begin{aligned} \frac{a}{2} \exp(-a|z|) &= \int_0^\infty \frac{1}{\sqrt{2\pi s^2}} \exp\left(-\frac{z^2}{2s^2}\right) \frac{a^2}{2} \exp\left(-\frac{a^2 s^2}{2}\right) ds^2 \\ &= \int_0^\infty \text{Normal}(z \mid 0, s^2) \text{Exp}\left(s^2 \mid \frac{a^2}{2}\right) ds^2 \end{aligned}$$

→ Useful for Bayesian computation!

$$\mathbf{y}^{n \times 1} = \mathbf{X}^{n \times p} \boldsymbol{\beta}^{p \times 1} + \mathbf{e}^{n \times 1}, \quad \mathbb{E}(\mathbf{e}) = \mathbf{0}, \quad \text{cov}(\mathbf{e}) = \sigma^2 \mathbf{I}_n$$

- Minimize squared error loss:

$$\hat{\boldsymbol{\beta}} = \arg \min \sum_{i=1}^n (y_i - \mathbf{x}_i^T \boldsymbol{\beta})^2 = \arg \min (\mathbf{y} - \mathbf{X}\boldsymbol{\beta})^T (\mathbf{y} - \mathbf{X}\boldsymbol{\beta}) = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}$$

- $\mathbb{E}(\hat{\boldsymbol{\beta}}) = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{X} \boldsymbol{\beta} = \boldsymbol{\beta}$

- $\text{var}(\hat{\boldsymbol{\beta}}) = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \sigma^2 \mathbf{I}_n \mathbf{X} (\mathbf{X}^T \mathbf{X})^{-1} = \sigma^2 (\mathbf{X}^T \mathbf{X})^{-1}$

- $\text{var}(\hat{\boldsymbol{\beta}}) = \sigma^2 (\mathbf{X}^T \mathbf{X})^{-1} = \sigma^2 (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{X} (\mathbf{X}^T \mathbf{X})^{-1}$

- $\mathbf{H} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T = \mathbf{H}^T, \quad (\mathbf{I}_n - \mathbf{H})^T = (\mathbf{I}_n - \mathbf{H})$

- $\text{rank}(\mathbf{H}) = \text{rank}(\mathbf{X} (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T) = \text{rank}(\mathbf{X}^T \mathbf{X}) = \text{rank}(\mathbf{X}) = p$

- $\text{rank}(\mathbf{I}_n - \mathbf{H}) = \text{rank}(\mathbf{I}_n) - \text{rank}(\mathbf{H}) = (n - p)$

- $\mathbf{H} + \mathbf{I}_n - \mathbf{H} = \mathbf{I}_n$  and  $(\mathbf{I}_n - \mathbf{H}) + \mathbf{H} = (\mathbf{I}_n - \mathbf{H} + \mathbf{X} (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T) = \mathbf{I}_n$

- Fitted values:  $\hat{\mathbf{y}} = \mathbf{X}\hat{\boldsymbol{\beta}} = \mathbf{X}(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y} = \mathbf{H}\mathbf{y}$

- Mean-Variance:  $\mathbb{E}(\hat{\mathbf{y}}) = \mathbf{X}\mathbb{E}(\hat{\boldsymbol{\beta}}) = \mathbf{X}\boldsymbol{\beta}, \quad \text{var}(\hat{\mathbf{y}}) = \sigma^2 \mathbf{H}$

- Residuals:  $\hat{\mathbf{e}} = (\mathbf{y} - \hat{\mathbf{y}}) = \mathbf{y} - \mathbf{X}\hat{\boldsymbol{\beta}} = \{\mathbf{I}_n - \mathbf{X}(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T\} \mathbf{y} = (\mathbf{I}_n - \mathbf{H})\mathbf{y}$

- Mean-Variance:  $\mathbb{E}(\hat{\mathbf{e}}) = (\mathbf{I}_n - \mathbf{H})\mathbf{X}\boldsymbol{\beta} = \mathbf{0}, \quad \text{var}(\hat{\mathbf{e}}) = \mathbb{E}(\hat{\mathbf{e}}\hat{\mathbf{e}}^T) = \sigma^2 (\mathbf{I}_n - \mathbf{H})$

- $\text{cov}(\hat{\mathbf{y}}, \hat{\mathbf{e}}) = \text{cov}\{\mathbf{H}\mathbf{y}, (\mathbf{I}_n - \mathbf{H})\mathbf{y}\} = \sigma^2 \mathbf{H}(\mathbf{I}_n - \mathbf{H}) = \sigma^2 (\mathbf{H} - \mathbf{H}^2) = \mathbf{0}$

$$\mathbf{y}^{n \times 1} = \mathbf{X}^{n \times p} \boldsymbol{\beta}^{p \times 1} + \mathbf{e}^{n \times 1}, \quad \mathbb{E}(\mathbf{e}) = \mathbf{0}, \quad \text{cov}(\mathbf{e}) = \sigma^2 \mathbf{I}_n$$

- Minimize squared error loss:

$$\hat{\boldsymbol{\beta}} = \arg \min \sum_{i=1}^n (y_i - \mathbf{x}_i^T \boldsymbol{\beta})^2 = \arg \min (\mathbf{y} - \mathbf{X}\boldsymbol{\beta})^T (\mathbf{y} - \mathbf{X}\boldsymbol{\beta}) = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}$$

- $\mathbb{E}(\hat{\boldsymbol{\beta}}) = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{X} \boldsymbol{\beta} = \boldsymbol{\beta}$

- $\text{var}(\hat{\boldsymbol{\beta}}) = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \sigma^2 \mathbf{I}_n \mathbf{X} (\mathbf{X}^T \mathbf{X})^{-1} = \sigma^2 (\mathbf{X}^T \mathbf{X})^{-1}$

- Fitted values:  $\hat{\mathbf{y}} = \mathbf{X}\hat{\boldsymbol{\beta}} = \mathbf{X}(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y} = \mathbf{H}\mathbf{y}$

- Mean-Variance:  $\mathbb{E}(\hat{\mathbf{y}}) = \mathbf{X}\mathbb{E}(\hat{\boldsymbol{\beta}}) = \mathbf{X}\boldsymbol{\beta}$ ,  $\text{var}(\hat{\mathbf{y}}) = \sigma^2 \mathbf{H}$

- Residuals:  $\hat{\mathbf{e}} = (\mathbf{y} - \hat{\mathbf{y}}) = \mathbf{y} - \mathbf{X}\hat{\boldsymbol{\beta}} = \{\mathbf{I}_n - \mathbf{X}(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T\} \mathbf{y} = (\mathbf{I}_n - \mathbf{H})\mathbf{y}$

- Mean-Variance:  $\mathbb{E}(\hat{\mathbf{e}}) = (\mathbf{I}_n - \mathbf{H})\mathbf{X}\boldsymbol{\beta} = \mathbf{0}$ ,  $\text{var}(\hat{\mathbf{e}}) = \mathbb{E}(\hat{\mathbf{e}}\hat{\mathbf{e}}^T) = \sigma^2(\mathbf{I}_n - \mathbf{H})$

- $\text{cov}(\hat{\mathbf{y}}, \hat{\mathbf{e}}) = \text{cov}\{\mathbf{H}\mathbf{y}, (\mathbf{I}_n - \mathbf{H})\mathbf{y}\} = \sigma^2 \mathbf{H}(\mathbf{I}_n - \mathbf{H}) = \sigma^2(\mathbf{H} - \mathbf{H}^2) = \mathbf{0}$

$$\mathbf{y}^{n \times 1} = \mathbf{X}^{n \times p} \boldsymbol{\beta}^{p \times 1} + \mathbf{e}^{n \times 1}, \quad \mathbb{E}(\mathbf{e}) = \mathbf{0}, \quad \text{cov}(\mathbf{e}) = \sigma^2 \mathbf{I}_n$$

- Minimize squared error loss:

$$\hat{\boldsymbol{\beta}} = \arg \min \sum_{i=1}^n (y_i - \mathbf{x}_i^T \boldsymbol{\beta})^2 = \arg \min (\mathbf{y} - \mathbf{X}\boldsymbol{\beta})^T (\mathbf{y} - \mathbf{X}\boldsymbol{\beta}) = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}$$

- $\mathbb{E}(\hat{\boldsymbol{\beta}}) = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{X} \boldsymbol{\beta} = \boldsymbol{\beta}$
- $\text{var}(\hat{\boldsymbol{\beta}}) = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \sigma^2 \mathbf{I}_n \mathbf{X} (\mathbf{X}^T \mathbf{X})^{-1} = \sigma^2 (\mathbf{X}^T \mathbf{X})^{-1}$

- Hat matrix:  $\mathbf{H} = \mathbf{X}(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T$

- Idempotent:  $\mathbf{H}^2 = \mathbf{H}$ ,  $(\mathbf{I}_n - \mathbf{H})^2 = (\mathbf{I}_n - \mathbf{H})$

- $\text{trace}(\mathbf{H}) = \text{trace}\{\mathbf{X}(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T\} = \text{trace}\{\mathbf{X}^T \mathbf{X} (\mathbf{X}^T \mathbf{X})^{-1}\} = p$ ,

- $\text{trace}(\mathbf{I}_n - \mathbf{H}) = \text{trace}(\mathbf{I}_n) - \text{trace}(\mathbf{H}) = (n - p)$

- $\mathbf{H}\mathbf{X} = \mathbf{X}$  and  $(\mathbf{I}_n - \mathbf{H})\mathbf{X} = \{\mathbf{I}_n - \mathbf{X}(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T\}\mathbf{X} = \mathbf{0}$

- Fitted values:  $\hat{\mathbf{y}} = \mathbf{X}\hat{\boldsymbol{\beta}} = \mathbf{X}(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y} = \mathbf{H}\mathbf{y}$

- Mean-Variance:  $\mathbb{E}(\hat{\mathbf{y}}) = \mathbf{X}\mathbb{E}(\hat{\boldsymbol{\beta}}) = \mathbf{X}\boldsymbol{\beta}$ ,  $\text{var}(\hat{\mathbf{y}}) = \sigma^2 \mathbf{H}$

- Residuals:  $\hat{\mathbf{e}} = (\mathbf{y} - \hat{\mathbf{y}}) = \mathbf{y} - \mathbf{X}\hat{\boldsymbol{\beta}} = \{\mathbf{I}_n - \mathbf{X}(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T\}\mathbf{y} = (\mathbf{I}_n - \mathbf{H})\mathbf{y}$

- Mean-Variance:  $\mathbb{E}(\hat{\mathbf{e}}) = (\mathbf{I}_n - \mathbf{H})\mathbf{X}\boldsymbol{\beta} = \mathbf{0}$ ,  $\text{var}(\hat{\mathbf{e}}) = \mathbb{E}(\hat{\mathbf{e}}\hat{\mathbf{e}}^T) = \sigma^2 (\mathbf{I}_n - \mathbf{H})$

- $\text{cov}(\hat{\mathbf{y}}, \hat{\mathbf{e}}) = \text{cov}\{\mathbf{H}\mathbf{y}, (\mathbf{I}_n - \mathbf{H})\mathbf{y}\} = \sigma^2 \mathbf{H}(\mathbf{I}_n - \mathbf{H}) = \sigma^2 (\mathbf{H} - \mathbf{H}^2) = \mathbf{0}$

$$\mathbf{y}^{n \times 1} = \mathbf{X}^{n \times p} \boldsymbol{\beta}^{p \times 1} + \mathbf{e}^{n \times 1}, \quad \mathbb{E}(\mathbf{e}) = \mathbf{0}, \quad \text{cov}(\mathbf{e}) = \sigma^2 \mathbf{I}_n$$

- Minimize squared error loss:

$$\hat{\boldsymbol{\beta}} = \arg \min \sum_{i=1}^n (y_i - \mathbf{x}_i^T \boldsymbol{\beta})^2 = \arg \min (\mathbf{y} - \mathbf{X}\boldsymbol{\beta})^T (\mathbf{y} - \mathbf{X}\boldsymbol{\beta}) = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}$$

- $\mathbb{E}(\hat{\boldsymbol{\beta}}) = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{X} \boldsymbol{\beta} = \boldsymbol{\beta}$
- $\text{var}(\hat{\boldsymbol{\beta}}) = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \sigma^2 \mathbf{I}_n \mathbf{X} (\mathbf{X}^T \mathbf{X})^{-1} = \sigma^2 (\mathbf{X}^T \mathbf{X})^{-1}$

- **Hat matrix:  $\mathbf{H} = \mathbf{X}(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T$**

- Idempotent:  $\mathbf{H}^2 = \mathbf{H}$ ,  $(\mathbf{I}_n - \mathbf{H})^2 = (\mathbf{I}_n - \mathbf{H})$
    - $\text{trace}(\mathbf{H}) = \text{trace}\{\mathbf{X}(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T\} = \text{trace}\{\mathbf{X}^T \mathbf{X} (\mathbf{X}^T \mathbf{X})^{-1}\} = p$ ,
    - $\text{trace}(\mathbf{I}_n - \mathbf{H}) = \text{trace}(\mathbf{I}_n) - \text{trace}(\mathbf{H}) = (n - p)$
    - $\mathbf{H}\mathbf{X} = \mathbf{X}$  and  $(\mathbf{I}_n - \mathbf{H})\mathbf{X} = \{\mathbf{I}_n - \mathbf{X}(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T\}\mathbf{X} = \mathbf{0}$

- Fitted values:  $\hat{\mathbf{y}} = \mathbf{X}\hat{\boldsymbol{\beta}} = \mathbf{X}(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y} = \mathbf{H}\mathbf{y}$

- Mean-Variance:  $\mathbb{E}(\hat{\mathbf{y}}) = \mathbf{X}\mathbb{E}(\hat{\boldsymbol{\beta}}) = \mathbf{X}\boldsymbol{\beta}$ ,  $\text{var}(\hat{\mathbf{y}}) = \sigma^2 \mathbf{H}$

- Residuals:  $\hat{\mathbf{e}} = (\mathbf{y} - \hat{\mathbf{y}}) = \mathbf{y} - \mathbf{X}\hat{\boldsymbol{\beta}} = \{\mathbf{I}_n - \mathbf{X}(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T\}\mathbf{y} = (\mathbf{I}_n - \mathbf{H})\mathbf{y}$

- Mean-Variance:  $\mathbb{E}(\hat{\mathbf{e}}) = (\mathbf{I}_n - \mathbf{H})\mathbf{X}\boldsymbol{\beta} = \mathbf{0}$ ,  $\text{var}(\hat{\mathbf{e}}) = \mathbb{E}(\hat{\mathbf{e}}\hat{\mathbf{e}}^T) = \sigma^2 (\mathbf{I}_n - \mathbf{H})$

- $\text{cov}(\hat{\mathbf{y}}, \hat{\mathbf{e}}) = \text{cov}\{\mathbf{H}\mathbf{y}, (\mathbf{I}_n - \mathbf{H})\mathbf{y}\} = \sigma^2 \mathbf{H}(\mathbf{I}_n - \mathbf{H}) = \sigma^2 (\mathbf{H} - \mathbf{H}^2) = \mathbf{0}$

$$\mathbf{y}^{n \times 1} = \mathbf{X}^{n \times p} \boldsymbol{\beta}^{p \times 1} + \mathbf{e}^{n \times 1}, \quad \mathbb{E}(\mathbf{e}) = \mathbf{0}, \quad \text{cov}(\mathbf{e}) = \sigma^2 \mathbf{I}_n$$

- Minimize squared error loss:

$$\hat{\boldsymbol{\beta}} = \arg \min \sum_{i=1}^n (y_i - \mathbf{x}_i^T \boldsymbol{\beta})^2 = \arg \min (\mathbf{y} - \mathbf{X}\boldsymbol{\beta})^T (\mathbf{y} - \mathbf{X}\boldsymbol{\beta}) = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}$$

- $\mathbb{E}(\hat{\boldsymbol{\beta}}) = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{X} \boldsymbol{\beta} = \boldsymbol{\beta}$

- $\text{var}(\hat{\boldsymbol{\beta}}) = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \sigma^2 \mathbf{I}_n \mathbf{X} (\mathbf{X}^T \mathbf{X})^{-1} = \sigma^2 (\mathbf{X}^T \mathbf{X})^{-1}$

- Hat matrix:  $\mathbf{H} = \mathbf{X}(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T$

- Idempotent:  $\mathbf{H}^2 = \mathbf{H}, \quad (\mathbf{I}_n - \mathbf{H})^2 = (\mathbf{I}_n - \mathbf{H})$

- $\text{trace}(\mathbf{H}) = \text{trace}\{\mathbf{X}(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T\} = \text{trace}\{\mathbf{X}^T \mathbf{X} (\mathbf{X}^T \mathbf{X})^{-1}\} = p,$

- $\text{trace}(\mathbf{I}_n - \mathbf{H}) = \text{trace}(\mathbf{I}_n) - \text{trace}(\mathbf{H}) = (n - p)$

- $\mathbf{H}\mathbf{X} = \mathbf{X}$  and  $(\mathbf{I}_n - \mathbf{H})\mathbf{X} = \{\mathbf{I}_n - \mathbf{X}(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T\}\mathbf{X} = \mathbf{0}$

- Fitted values:  $\hat{\mathbf{y}} = \mathbf{X}\hat{\boldsymbol{\beta}} = \mathbf{X}(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y} = \mathbf{H}\mathbf{y}$

- Mean-Variance:  $\mathbb{E}(\hat{\mathbf{y}}) = \mathbf{X}\mathbb{E}(\hat{\boldsymbol{\beta}}) = \mathbf{X}\boldsymbol{\beta}, \quad \text{var}(\hat{\mathbf{y}}) = \sigma^2 \mathbf{H}$

- Residuals:  $\hat{\mathbf{e}} = (\mathbf{y} - \hat{\mathbf{y}}) = \mathbf{y} - \mathbf{X}\hat{\boldsymbol{\beta}} = \{\mathbf{I}_n - \mathbf{X}(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T\}\mathbf{y} = (\mathbf{I}_n - \mathbf{H})\mathbf{y}$

- Mean-Variance:  $\mathbb{E}(\hat{\mathbf{e}}) = (\mathbf{I}_n - \mathbf{H})\mathbf{X}\boldsymbol{\beta} = \mathbf{0}, \quad \text{var}(\hat{\mathbf{e}}) = \mathbb{E}(\hat{\mathbf{e}}\hat{\mathbf{e}}^T) = \sigma^2 (\mathbf{I}_n - \mathbf{H})$

- $\text{cov}(\hat{\mathbf{y}}, \hat{\mathbf{e}}) = \text{cov}\{\mathbf{H}\mathbf{y}, (\mathbf{I}_n - \mathbf{H})\mathbf{y}\} = \sigma^2 \mathbf{H}(\mathbf{I}_n - \mathbf{H}) = \sigma^2 (\mathbf{H} - \mathbf{H}^2) = \mathbf{0}$

$$\mathbf{y}^{n \times 1} = \mathbf{X}^{n \times p} \boldsymbol{\beta}^{p \times 1} + \mathbf{e}^{n \times 1}, \quad \mathbb{E}(\mathbf{e}) = \mathbf{0}, \quad \text{cov}(\mathbf{e}) = \sigma^2 \mathbf{I}_n$$

- Minimize squared error loss:

$$\hat{\boldsymbol{\beta}} = \arg \min \sum_{i=1}^n (y_i - \mathbf{x}_i^T \boldsymbol{\beta})^2 = \arg \min (\mathbf{y} - \mathbf{X}\boldsymbol{\beta})^T (\mathbf{y} - \mathbf{X}\boldsymbol{\beta}) = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}$$

- $\mathbb{E}(\hat{\boldsymbol{\beta}}) = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{X} \boldsymbol{\beta} = \boldsymbol{\beta}$
- $\text{var}(\hat{\boldsymbol{\beta}}) = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \sigma^2 \mathbf{I}_n \mathbf{X} (\mathbf{X}^T \mathbf{X})^{-1} = \sigma^2 (\mathbf{X}^T \mathbf{X})^{-1}$ 
  - Hat matrix:  $\mathbf{H} = \mathbf{X} (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T$
  - Idempotent:  $\mathbf{H}^2 = \mathbf{H}$ ,  $(\mathbf{I}_n - \mathbf{H})^2 = (\mathbf{I}_n - \mathbf{H})$
  - $\text{trace}(\mathbf{H}) = \text{trace}\{\mathbf{X} (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T\} = \text{trace}\{\mathbf{X}^T \mathbf{X} (\mathbf{X}^T \mathbf{X})^{-1}\} = p$ 
    - $\text{trace}(\mathbf{I}_n - \mathbf{H}) = \text{trace}(\mathbf{I}_n) - \text{trace}(\mathbf{H}) = (n - p)$
    - $\mathbf{H}\mathbf{X} = \mathbf{X}$  and  $(\mathbf{I}_n - \mathbf{H})\mathbf{X} = \{\mathbf{I}_n - \mathbf{X} (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T\} \mathbf{X} = \mathbf{0}$
- Fitted values:  $\hat{\mathbf{y}} = \mathbf{X} \hat{\boldsymbol{\beta}} = \mathbf{X} (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y} = \mathbf{H}\mathbf{y}$
- Mean-Variance:  $\mathbb{E}(\hat{\mathbf{y}}) = \mathbf{X} \mathbb{E}(\hat{\boldsymbol{\beta}}) = \mathbf{X} \boldsymbol{\beta}$ ,  $\text{var}(\hat{\mathbf{y}}) = \sigma^2 \mathbf{H}$
- Residuals:  $\hat{\mathbf{e}} = (\mathbf{y} - \hat{\mathbf{y}}) = \mathbf{y} - \mathbf{X} \hat{\boldsymbol{\beta}} = \{\mathbf{I}_n - \mathbf{X} (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T\} \mathbf{y} = (\mathbf{I}_n - \mathbf{H}) \mathbf{y}$
- Mean-Variance:  $\mathbb{E}(\hat{\mathbf{e}}) = (\mathbf{I}_n - \mathbf{H}) \mathbf{X} \boldsymbol{\beta} = \mathbf{0}$ ,  $\text{var}(\hat{\mathbf{e}}) = \mathbb{E}(\hat{\mathbf{e}} \hat{\mathbf{e}}^T) = \sigma^2 (\mathbf{I}_n - \mathbf{H})$
- $\text{cov}(\hat{\mathbf{y}}, \hat{\mathbf{e}}) = \text{cov}\{\mathbf{H}\mathbf{y}, (\mathbf{I}_n - \mathbf{H})\mathbf{y}\} = \sigma^2 \mathbf{H} (\mathbf{I}_n - \mathbf{H}) = \sigma^2 (\mathbf{H} - \mathbf{H}^2) = \mathbf{0}$



$$\mathbf{y}^{n \times 1} = \mathbf{X}^{n \times p} \boldsymbol{\beta}^{p \times 1} + \mathbf{e}^{n \times 1}, \quad \mathbb{E}(\mathbf{e}) = \mathbf{0}, \quad \text{cov}(\mathbf{e}) = \sigma^2 \mathbf{I}_n$$

- Minimize squared error loss:

$$\hat{\boldsymbol{\beta}} = \arg \min \sum_{i=1}^n (y_i - \mathbf{x}_i^T \boldsymbol{\beta})^2 = \arg \min (\mathbf{y} - \mathbf{X}\boldsymbol{\beta})^T (\mathbf{y} - \mathbf{X}\boldsymbol{\beta}) = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}$$

- $\mathbb{E}(\hat{\boldsymbol{\beta}}) = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{X} \boldsymbol{\beta} = \boldsymbol{\beta}$
- $\text{var}(\hat{\boldsymbol{\beta}}) = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \sigma^2 \mathbf{I}_n \mathbf{X} (\mathbf{X}^T \mathbf{X})^{-1} = \sigma^2 (\mathbf{X}^T \mathbf{X})^{-1}$ 
  - Hat matrix:  $\mathbf{H} = \mathbf{X}(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T$
  - Idempotent:  $\mathbf{H}^2 = \mathbf{H}$ ,  $(\mathbf{I}_n - \mathbf{H})^2 = (\mathbf{I}_n - \mathbf{H})$
  - $\text{trace}(\mathbf{H}) = \text{trace}\{\mathbf{X}(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T\} = \text{trace}\{\mathbf{X}^T \mathbf{X} (\mathbf{X}^T \mathbf{X})^{-1}\} = p$ ,
  - $\text{trace}(\mathbf{I}_n - \mathbf{H}) = \text{trace}(\mathbf{I}_n) - \text{trace}(\mathbf{H}) = (n - p)$
  - $\mathbf{H}\mathbf{X} = \mathbf{X}$  and  $(\mathbf{I}_n - \mathbf{H})\mathbf{X} = \{\mathbf{I}_n - \mathbf{X}(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T\}\mathbf{X} = \mathbf{0}$
- Fitted values:  $\hat{\mathbf{y}} = \mathbf{X}\hat{\boldsymbol{\beta}} = \mathbf{X}(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y} = \mathbf{H}\mathbf{y}$
- Mean-Variance:  $\mathbb{E}(\hat{\mathbf{y}}) = \mathbf{X}\mathbb{E}(\hat{\boldsymbol{\beta}}) = \mathbf{X}\boldsymbol{\beta}$ ,  $\text{var}(\hat{\mathbf{y}}) = \sigma^2 \mathbf{H}$
- Residuals:  $\hat{\mathbf{e}} = (\mathbf{y} - \hat{\mathbf{y}}) = \mathbf{y} - \mathbf{X}\hat{\boldsymbol{\beta}} = \{\mathbf{I}_n - \mathbf{X}(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T\}\mathbf{y} = (\mathbf{I}_n - \mathbf{H})\mathbf{y}$
- Mean-Variance:  $\mathbb{E}(\hat{\mathbf{e}}) = (\mathbf{I}_n - \mathbf{H})\mathbf{X}\boldsymbol{\beta} = \mathbf{0}$ ,  $\text{var}(\hat{\mathbf{e}}) = \mathbb{E}(\hat{\mathbf{e}}\hat{\mathbf{e}}^T) = \sigma^2 (\mathbf{I}_n - \mathbf{H})$
- $\text{cov}(\hat{\mathbf{y}}, \hat{\mathbf{e}}) = \text{cov}\{\mathbf{H}\mathbf{y}, (\mathbf{I}_n - \mathbf{H})\mathbf{y}\} = \sigma^2 \mathbf{H}(\mathbf{I}_n - \mathbf{H}) = \sigma^2 (\mathbf{H} - \mathbf{H}^2) = \mathbf{0}$

$$\mathbf{y}^{n \times 1} = \mathbf{X}^{n \times p} \boldsymbol{\beta}^{p \times 1} + \mathbf{e}^{n \times 1}, \quad \mathbb{E}(\mathbf{e}) = \mathbf{0}, \quad \text{cov}(\mathbf{e}) = \sigma^2 \mathbf{I}_n$$

- Minimize squared error loss:

$$\hat{\boldsymbol{\beta}} = \arg \min \sum_{i=1}^n (y_i - \mathbf{x}_i^T \boldsymbol{\beta})^2 = \arg \min (\mathbf{y} - \mathbf{X}\boldsymbol{\beta})^T (\mathbf{y} - \mathbf{X}\boldsymbol{\beta}) = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}$$

- $\mathbb{E}(\hat{\boldsymbol{\beta}}) = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{X} \boldsymbol{\beta} = \boldsymbol{\beta}$
- $\text{var}(\hat{\boldsymbol{\beta}}) = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \sigma^2 \mathbf{I}_n \mathbf{X} (\mathbf{X}^T \mathbf{X})^{-1} = \sigma^2 (\mathbf{X}^T \mathbf{X})^{-1}$ 
  - Hat matrix:  $\mathbf{H} = \mathbf{X}(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T$
  - Idempotent:  $\mathbf{H}^2 = \mathbf{H}$ ,  $(\mathbf{I}_n - \mathbf{H})^2 = (\mathbf{I}_n - \mathbf{H})$
  - $\text{trace}(\mathbf{H}) = \text{trace}\{\mathbf{X}(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T\} = \text{trace}\{\mathbf{X}^T \mathbf{X} (\mathbf{X}^T \mathbf{X})^{-1}\} = p$ ,
  - $\text{trace}(\mathbf{I}_n - \mathbf{H}) = \text{trace}(\mathbf{I}_n) - \text{trace}(\mathbf{H}) = (n - p)$
  - $\mathbf{H}\mathbf{X} = \mathbf{X}$  and  $(\mathbf{I}_n - \mathbf{H})\mathbf{X} = \{\mathbf{I}_n - \mathbf{X}(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T\}\mathbf{X} = \mathbf{0}$
- Fitted values:  $\hat{\mathbf{y}} = \mathbf{X}\hat{\boldsymbol{\beta}} = \mathbf{X}(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y} = \mathbf{H}\mathbf{y}$
- Mean-Variance:  $\mathbb{E}(\hat{\mathbf{y}}) = \mathbf{X}\mathbb{E}(\hat{\boldsymbol{\beta}}) = \mathbf{X}\boldsymbol{\beta}$ ,  $\text{var}(\hat{\mathbf{y}}) = \sigma^2 \mathbf{H}$
- Residuals:  $\hat{\mathbf{e}} = (\mathbf{y} - \hat{\mathbf{y}}) = \mathbf{y} - \mathbf{X}\hat{\boldsymbol{\beta}} = \{\mathbf{I}_n - \mathbf{X}(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T\}\mathbf{y} = (\mathbf{I}_n - \mathbf{H})\mathbf{y}$
- Mean-Variance:  $\mathbb{E}(\hat{\mathbf{e}}) = (\mathbf{I}_n - \mathbf{H})\mathbf{X}\boldsymbol{\beta} = \mathbf{0}$ ,  $\text{var}(\hat{\mathbf{e}}) = \mathbb{E}(\hat{\mathbf{e}}\hat{\mathbf{e}}^T) = \sigma^2 (\mathbf{I}_n - \mathbf{H})$
- $\text{cov}(\hat{\mathbf{y}}, \hat{\mathbf{e}}) = \text{cov}\{\mathbf{H}\mathbf{y}, (\mathbf{I}_n - \mathbf{H})\mathbf{y}\} = \sigma^2 \mathbf{H}(\mathbf{I}_n - \mathbf{H}) = \sigma^2 (\mathbf{H} - \mathbf{H}^2) = \mathbf{0}$

$$\mathbf{y}^{n \times 1} = \mathbf{X}^{n \times p} \boldsymbol{\beta}^{p \times 1} + \mathbf{e}^{n \times 1}, \quad \mathbb{E}(\mathbf{e}) = \mathbf{0}, \quad \text{cov}(\mathbf{e}) = \sigma^2 \mathbf{I}_n$$

- Minimize squared error loss:

$$\hat{\boldsymbol{\beta}} = \arg \min \sum_{i=1}^n (y_i - \mathbf{x}_i^T \boldsymbol{\beta})^2 = \arg \min (\mathbf{y} - \mathbf{X}\boldsymbol{\beta})^T (\mathbf{y} - \mathbf{X}\boldsymbol{\beta}) = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}$$

- $\mathbb{E}(\hat{\boldsymbol{\beta}}) = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{X} \boldsymbol{\beta} = \boldsymbol{\beta}$
- $\text{var}(\hat{\boldsymbol{\beta}}) = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \sigma^2 \mathbf{I}_n \mathbf{X} (\mathbf{X}^T \mathbf{X})^{-1} = \sigma^2 (\mathbf{X}^T \mathbf{X})^{-1}$ 
  - Hat matrix:  $\mathbf{H} = \mathbf{X} (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T$
  - Idempotent:  $\mathbf{H}^2 = \mathbf{H}$ ,  $(\mathbf{I}_n - \mathbf{H})^2 = (\mathbf{I}_n - \mathbf{H})$
  - $\text{trace}(\mathbf{H}) = \text{trace}\{\mathbf{X} (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T\} = \text{trace}\{\mathbf{X}^T \mathbf{X} (\mathbf{X}^T \mathbf{X})^{-1}\} = p$ ,
  - $\text{trace}(\mathbf{I}_n - \mathbf{H}) = \text{trace}(\mathbf{I}_n) - \text{trace}(\mathbf{H}) = (n - p)$
  - $\mathbf{H}\mathbf{X} = \mathbf{X}$  and  $(\mathbf{I}_n - \mathbf{H})\mathbf{X} = \{\mathbf{I}_n - \mathbf{X} (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T\} \mathbf{X} = \mathbf{0}$
- Fitted values:  $\hat{\mathbf{y}} = \mathbf{X}\hat{\boldsymbol{\beta}} = \mathbf{X} (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y} = \mathbf{H}\mathbf{y}$ 
  - Mean-Variance:  $\mathbb{E}(\hat{\mathbf{y}}) = \mathbf{X}\mathbb{E}(\hat{\boldsymbol{\beta}}) = \mathbf{X}\boldsymbol{\beta}$ ,  $\text{var}(\hat{\mathbf{y}}) = \sigma^2 \mathbf{H}$
  - Residuals:  $\hat{\mathbf{e}} = (\mathbf{y} - \hat{\mathbf{y}}) = \mathbf{y} - \mathbf{X}\hat{\boldsymbol{\beta}} = \{\mathbf{I}_n - \mathbf{X} (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T\} \mathbf{y} = (\mathbf{I}_n - \mathbf{H})\mathbf{y}$
  - Mean-Variance:  $\mathbb{E}(\hat{\mathbf{e}}) = (\mathbf{I}_n - \mathbf{H})\mathbf{X}\boldsymbol{\beta} = \mathbf{0}$ ,  $\text{var}(\hat{\mathbf{e}}) = \mathbb{E}(\hat{\mathbf{e}}\hat{\mathbf{e}}^T) = \sigma^2 (\mathbf{I}_n - \mathbf{H})$
  - $\text{cov}(\hat{\mathbf{y}}, \hat{\mathbf{e}}) = \text{cov}\{\mathbf{H}\mathbf{y}, (\mathbf{I}_n - \mathbf{H})\mathbf{y}\} = \sigma^2 \mathbf{H}(\mathbf{I}_n - \mathbf{H}) = \sigma^2 (\mathbf{H} - \mathbf{H}^2) = \mathbf{0}$

# Ordinary Least Squares

$$\mathbf{y}^{n \times 1} = \mathbf{X}^{n \times p} \boldsymbol{\beta}^{p \times 1} + \mathbf{e}^{n \times 1}, \quad \mathbb{E}(\mathbf{e}) = \mathbf{0}, \quad \text{cov}(\mathbf{e}) = \sigma^2 \mathbf{I}_n$$

- Minimize squared error loss:

$$\hat{\boldsymbol{\beta}} = \arg \min \sum_{i=1}^n (y_i - \mathbf{x}_i^T \boldsymbol{\beta})^2 = \arg \min (\mathbf{y} - \mathbf{X}\boldsymbol{\beta})^T (\mathbf{y} - \mathbf{X}\boldsymbol{\beta}) = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}$$

- $\mathbb{E}(\hat{\boldsymbol{\beta}}) = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{X} \boldsymbol{\beta} = \boldsymbol{\beta}$
- $\text{var}(\hat{\boldsymbol{\beta}}) = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \sigma^2 \mathbf{I}_n \mathbf{X} (\mathbf{X}^T \mathbf{X})^{-1} = \sigma^2 (\mathbf{X}^T \mathbf{X})^{-1}$ 
  - Hat matrix:  $\mathbf{H} = \mathbf{X}(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T$
  - Idempotent:  $\mathbf{H}^2 = \mathbf{H}$ ,  $(\mathbf{I}_n - \mathbf{H})^2 = (\mathbf{I}_n - \mathbf{H})$
  - $\text{trace}(\mathbf{H}) = \text{trace}\{\mathbf{X}(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T\} = \text{trace}\{\mathbf{X}^T \mathbf{X} (\mathbf{X}^T \mathbf{X})^{-1}\} = p$ ,
  - $\text{trace}(\mathbf{I}_n - \mathbf{H}) = \text{trace}(\mathbf{I}_n) - \text{trace}(\mathbf{H}) = (n - p)$
  - $\mathbf{H}\mathbf{X} = \mathbf{X}$  and  $(\mathbf{I}_n - \mathbf{H})\mathbf{X} = \{\mathbf{I}_n - \mathbf{X}(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T\}\mathbf{X} = \mathbf{0}$
- Fitted values:  $\hat{\mathbf{y}} = \mathbf{X}\hat{\boldsymbol{\beta}} = \mathbf{X}(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y} = \mathbf{H}\mathbf{y}$
- Mean-Variance:  $\mathbb{E}(\hat{\mathbf{y}}) = \mathbf{X}\mathbb{E}(\hat{\boldsymbol{\beta}}) = \mathbf{X}\boldsymbol{\beta}$ ,  $\text{var}(\hat{\mathbf{y}}) = \sigma^2 \mathbf{H}$ 
  - Residuals:  $\hat{\mathbf{e}} = (\mathbf{y} - \hat{\mathbf{y}}) = \mathbf{y} - \mathbf{X}\hat{\boldsymbol{\beta}} = \{\mathbf{I}_n - \mathbf{X}(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T\}\mathbf{y} = (\mathbf{I}_n - \mathbf{H})\mathbf{y}$
  - Mean-Variance:  $\mathbb{E}(\hat{\mathbf{e}}) = (\mathbf{I}_n - \mathbf{H})\mathbf{X}\boldsymbol{\beta} = \mathbf{0}$ ,  $\text{var}(\hat{\mathbf{e}}) = \mathbb{E}(\hat{\mathbf{e}}\hat{\mathbf{e}}^T) = \sigma^2 (\mathbf{I}_n - \mathbf{H})$
  - $\text{cov}(\hat{\mathbf{y}}, \hat{\mathbf{e}}) = \text{cov}\{\mathbf{H}\mathbf{y}, (\mathbf{I}_n - \mathbf{H})\mathbf{y}\} = \sigma^2 \mathbf{H}(\mathbf{I}_n - \mathbf{H}) = \sigma^2 (\mathbf{H} - \mathbf{H}^2) = \mathbf{0}$

# Ordinary Least Squares

$$\mathbf{y}^{n \times 1} = \mathbf{X}^{n \times p} \boldsymbol{\beta}^{p \times 1} + \mathbf{e}^{n \times 1}, \quad \mathbb{E}(\mathbf{e}) = \mathbf{0}, \quad \text{cov}(\mathbf{e}) = \sigma^2 \mathbf{I}_n$$

- Minimize squared error loss:

$$\hat{\boldsymbol{\beta}} = \arg \min \sum_{i=1}^n (y_i - \mathbf{x}_i^T \boldsymbol{\beta})^2 = \arg \min (\mathbf{y} - \mathbf{X}\boldsymbol{\beta})^T (\mathbf{y} - \mathbf{X}\boldsymbol{\beta}) = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}$$

- $\mathbb{E}(\hat{\boldsymbol{\beta}}) = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{X} \boldsymbol{\beta} = \boldsymbol{\beta}$
- $\text{var}(\hat{\boldsymbol{\beta}}) = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \sigma^2 \mathbf{I}_n \mathbf{X} (\mathbf{X}^T \mathbf{X})^{-1} = \sigma^2 (\mathbf{X}^T \mathbf{X})^{-1}$ 
  - Hat matrix:  $\mathbf{H} = \mathbf{X} (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T$
  - Idempotent:  $\mathbf{H}^2 = \mathbf{H}$ ,  $(\mathbf{I}_n - \mathbf{H})^2 = (\mathbf{I}_n - \mathbf{H})$
  - $\text{trace}(\mathbf{H}) = \text{trace}\{\mathbf{X} (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T\} = \text{trace}\{\mathbf{X}^T \mathbf{X} (\mathbf{X}^T \mathbf{X})^{-1}\} = p$ ,
  - $\text{trace}(\mathbf{I}_n - \mathbf{H}) = \text{trace}(\mathbf{I}_n) - \text{trace}(\mathbf{H}) = (n - p)$
  - $\mathbf{H}\mathbf{X} = \mathbf{X}$  and  $(\mathbf{I}_n - \mathbf{H})\mathbf{X} = \{\mathbf{I}_n - \mathbf{X} (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T\} \mathbf{X} = \mathbf{0}$
- Fitted values:  $\hat{\mathbf{y}} = \mathbf{X}\hat{\boldsymbol{\beta}} = \mathbf{X} (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y} = \mathbf{H}\mathbf{y}$
- Mean-Variance:  $\mathbb{E}(\hat{\mathbf{y}}) = \mathbf{X}\mathbb{E}(\hat{\boldsymbol{\beta}}) = \mathbf{X}\boldsymbol{\beta}$ ,  $\text{var}(\hat{\mathbf{y}}) = \sigma^2 \mathbf{H}$
- Residuals:  $\hat{\mathbf{e}} = (\mathbf{y} - \hat{\mathbf{y}}) = \mathbf{y} - \mathbf{X}\hat{\boldsymbol{\beta}} = \{\mathbf{I}_n - \mathbf{X} (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T\} \mathbf{y} = (\mathbf{I}_n - \mathbf{H})\mathbf{y}$ 
  - Mean-Variance:  $\mathbb{E}(\hat{\mathbf{e}}) = (\mathbf{I}_n - \mathbf{H})\mathbf{X}\boldsymbol{\beta} = \mathbf{0}$ ,  $\text{var}(\hat{\mathbf{e}}) = \mathbb{E}(\hat{\mathbf{e}}\hat{\mathbf{e}}^T) = \sigma^2 (\mathbf{I}_n - \mathbf{H})$
  - $\text{cov}(\hat{\mathbf{y}}, \hat{\mathbf{e}}) = \text{cov}\{\mathbf{H}\mathbf{y}, (\mathbf{I}_n - \mathbf{H})\mathbf{y}\} = \sigma^2 \mathbf{H}(\mathbf{I}_n - \mathbf{H}) = \sigma^2 (\mathbf{H} - \mathbf{H}^2) = \mathbf{0}$

# Ordinary Least Squares

$$\mathbf{y}^{n \times 1} = \mathbf{X}^{n \times p} \boldsymbol{\beta}^{p \times 1} + \mathbf{e}^{n \times 1}, \quad \mathbb{E}(\mathbf{e}) = \mathbf{0}, \quad \text{cov}(\mathbf{e}) = \sigma^2 \mathbf{I}_n$$

- Minimize squared error loss:

$$\hat{\boldsymbol{\beta}} = \arg \min \sum_{i=1}^n (y_i - \mathbf{x}_i^T \boldsymbol{\beta})^2 = \arg \min (\mathbf{y} - \mathbf{X}\boldsymbol{\beta})^T (\mathbf{y} - \mathbf{X}\boldsymbol{\beta}) = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}$$

- $\mathbb{E}(\hat{\boldsymbol{\beta}}) = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{X} \boldsymbol{\beta} = \boldsymbol{\beta}$

- $\text{var}(\hat{\boldsymbol{\beta}}) = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \sigma^2 \mathbf{I}_n \mathbf{X} (\mathbf{X}^T \mathbf{X})^{-1} = \sigma^2 (\mathbf{X}^T \mathbf{X})^{-1}$

- Hat matrix:  $\mathbf{H} = \mathbf{X}(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T$

- Idempotent:  $\mathbf{H}^2 = \mathbf{H}$ ,  $(\mathbf{I}_n - \mathbf{H})^2 = (\mathbf{I}_n - \mathbf{H})$

- $\text{trace}(\mathbf{H}) = \text{trace}\{\mathbf{X}(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T\} = \text{trace}\{\mathbf{X}^T \mathbf{X} (\mathbf{X}^T \mathbf{X})^{-1}\} = p$ ,

- $\text{trace}(\mathbf{I}_n - \mathbf{H}) = \text{trace}(\mathbf{I}_n) - \text{trace}(\mathbf{H}) = (n - p)$

- $\mathbf{H}\mathbf{X} = \mathbf{X}$  and  $(\mathbf{I}_n - \mathbf{H})\mathbf{X} = \{\mathbf{I}_n - \mathbf{X}(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T\}\mathbf{X} = \mathbf{0}$

- Fitted values:  $\hat{\mathbf{y}} = \mathbf{X}\hat{\boldsymbol{\beta}} = \mathbf{X}(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y} = \mathbf{H}\mathbf{y}$

- Mean-Variance:  $\mathbb{E}(\hat{\mathbf{y}}) = \mathbf{X}\mathbb{E}(\hat{\boldsymbol{\beta}}) = \mathbf{X}\boldsymbol{\beta}$ ,  $\text{var}(\hat{\mathbf{y}}) = \sigma^2 \mathbf{H}$

- Residuals:  $\hat{\mathbf{e}} = (\mathbf{y} - \hat{\mathbf{y}}) = \mathbf{y} - \mathbf{X}\hat{\boldsymbol{\beta}} = \{\mathbf{I}_n - \mathbf{X}(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T\}\mathbf{y} = (\mathbf{I}_n - \mathbf{H})\mathbf{y}$

- Mean-Variance:  $\mathbb{E}(\hat{\mathbf{e}}) = (\mathbf{I}_n - \mathbf{H})\mathbf{X}\boldsymbol{\beta} = \mathbf{0}$ ,  $\text{var}(\hat{\mathbf{e}}) = \mathbb{E}(\hat{\mathbf{e}}\hat{\mathbf{e}}^T) = \sigma^2 (\mathbf{I}_n - \mathbf{H})$

- $\text{cov}(\hat{\mathbf{y}}, \hat{\mathbf{e}}) = \text{cov}\{\mathbf{H}\mathbf{y}, (\mathbf{I}_n - \mathbf{H})\mathbf{y}\} = \sigma^2 \mathbf{H}(\mathbf{I}_n - \mathbf{H}) = \sigma^2 (\mathbf{H} - \mathbf{H}^2) = \mathbf{0}$

# Ordinary Least Squares

$$\mathbf{y}^{n \times 1} = \mathbf{X}^{n \times p} \boldsymbol{\beta}^{p \times 1} + \mathbf{e}^{n \times 1}, \quad \mathbb{E}(\mathbf{e}) = \mathbf{0}, \quad \text{cov}(\mathbf{e}) = \sigma^2 \mathbf{I}_n$$

- Minimize squared error loss:

$$\hat{\boldsymbol{\beta}} = \arg \min \sum_{i=1}^n (y_i - \mathbf{x}_i^T \boldsymbol{\beta})^2 = \arg \min (\mathbf{y} - \mathbf{X}\boldsymbol{\beta})^T (\mathbf{y} - \mathbf{X}\boldsymbol{\beta}) = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}$$

- $\mathbb{E}(\hat{\boldsymbol{\beta}}) = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{X} \boldsymbol{\beta} = \boldsymbol{\beta}$
- $\text{var}(\hat{\boldsymbol{\beta}}) = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \sigma^2 \mathbf{I}_n \mathbf{X} (\mathbf{X}^T \mathbf{X})^{-1} = \sigma^2 (\mathbf{X}^T \mathbf{X})^{-1}$ 
  - Hat matrix:  $\mathbf{H} = \mathbf{X}(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T$
  - Idempotent:  $\mathbf{H}^2 = \mathbf{H}$ ,  $(\mathbf{I}_n - \mathbf{H})^2 = (\mathbf{I}_n - \mathbf{H})$
  - $\text{trace}(\mathbf{H}) = \text{trace}\{\mathbf{X}(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T\} = \text{trace}\{\mathbf{X}^T \mathbf{X} (\mathbf{X}^T \mathbf{X})^{-1}\} = p$ ,
  - $\text{trace}(\mathbf{I}_n - \mathbf{H}) = \text{trace}(\mathbf{I}_n) - \text{trace}(\mathbf{H}) = (n - p)$
  - $\mathbf{H}\mathbf{X} = \mathbf{X}$  and  $(\mathbf{I}_n - \mathbf{H})\mathbf{X} = \{\mathbf{I}_n - \mathbf{X}(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T\}\mathbf{X} = \mathbf{0}$
- Fitted values:  $\hat{\mathbf{y}} = \mathbf{X}\hat{\boldsymbol{\beta}} = \mathbf{X}(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y} = \mathbf{H}\mathbf{y}$
- Mean-Variance:  $\mathbb{E}(\hat{\mathbf{y}}) = \mathbf{X}\mathbb{E}(\hat{\boldsymbol{\beta}}) = \mathbf{X}\boldsymbol{\beta}$ ,  $\text{var}(\hat{\mathbf{y}}) = \sigma^2 \mathbf{H}$
- Residuals:  $\hat{\mathbf{e}} = (\mathbf{y} - \hat{\mathbf{y}}) = \mathbf{y} - \mathbf{X}\hat{\boldsymbol{\beta}} = \{\mathbf{I}_n - \mathbf{X}(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T\}\mathbf{y} = (\mathbf{I}_n - \mathbf{H})\mathbf{y}$
- Mean-Variance:  $\mathbb{E}(\hat{\mathbf{e}}) = (\mathbf{I}_n - \mathbf{H})\mathbf{X}\boldsymbol{\beta} = \mathbf{0}$ ,  $\text{var}(\hat{\mathbf{e}}) = \mathbb{E}(\hat{\mathbf{e}}\hat{\mathbf{e}}^T) = \sigma^2 (\mathbf{I}_n - \mathbf{H})$
- $\text{cov}(\hat{\mathbf{y}}, \hat{\mathbf{e}}) = \text{cov}\{\mathbf{H}\mathbf{y}, (\mathbf{I}_n - \mathbf{H})\mathbf{y}\} = \sigma^2 \mathbf{H}(\mathbf{I}_n - \mathbf{H}) = \sigma^2 (\mathbf{H} - \mathbf{H}^2) = \mathbf{0}$

$$\mathbf{y}^{n \times 1} = \mathbf{X}^{n \times p} \boldsymbol{\beta}^{p \times 1} + \mathbf{e}^{n \times 1}, \quad \mathbb{E}(\mathbf{e}) = \mathbf{0}, \quad \text{cov}(\mathbf{e}) = \sigma^2 \mathbf{I}_n$$

- Minimize squared error loss:

$$\hat{\boldsymbol{\beta}} = \arg \min \sum_{i=1}^n (y_i - \mathbf{x}_i^T \boldsymbol{\beta})^2 = \arg \min (\mathbf{y} - \mathbf{X}\boldsymbol{\beta})^T (\mathbf{y} - \mathbf{X}\boldsymbol{\beta}) = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}$$

- Fitted values:  $\hat{\mathbf{y}} = \mathbf{X}\hat{\boldsymbol{\beta}} = \mathbf{X}(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y} = \mathbf{H}\mathbf{y}$
- Mean-Variance:  $\mathbb{E}(\hat{\mathbf{y}}) = \mathbf{H}\mathbb{E}(\mathbf{y}) = \mathbf{H}\mathbf{X}\boldsymbol{\beta} = \mathbf{X}\boldsymbol{\beta}, \quad \text{var}(\hat{\mathbf{y}}) = \sigma^2 \mathbf{H}$
- Residuals:  $\hat{\mathbf{e}} = (\mathbf{y} - \hat{\mathbf{y}}) = \mathbf{y} - \mathbf{X}\hat{\boldsymbol{\beta}} = \{\mathbf{I}_n - \mathbf{X}(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T\} \mathbf{y} = (\mathbf{I}_n - \mathbf{H})\mathbf{y}$
- Mean-Variance:  $\mathbb{E}(\hat{\mathbf{e}}) = (\mathbf{I}_n - \mathbf{H})\mathbf{X}\boldsymbol{\beta} = \mathbf{0}, \quad \text{var}(\hat{\mathbf{e}}) = \mathbb{E}(\hat{\mathbf{e}}\hat{\mathbf{e}}^T) = \sigma^2 (\mathbf{I}_n - \mathbf{H})$
- $\text{cov}(\hat{\mathbf{y}}, \hat{\mathbf{e}}) = \text{cov}\{\mathbf{H}\mathbf{y}, (\mathbf{I}_n - \mathbf{H})\mathbf{y}\} = \sigma^2 \mathbf{H}(\mathbf{I}_n - \mathbf{H}) = \sigma^2 (\mathbf{H} - \mathbf{H}^2) = \mathbf{0}$



# Ordinary Least Squares

$$\mathbf{y}^{n \times 1} = \mathbf{X}^{n \times p} \boldsymbol{\beta}^{p \times 1} + \mathbf{e}^{n \times 1}, \quad \mathbb{E}(\mathbf{e}) = \mathbf{0}, \quad \text{cov}(\mathbf{e}) = \sigma^2 \mathbf{I}_n$$

- Minimize squared error loss:

$$\hat{\boldsymbol{\beta}} = \arg \min \sum_{i=1}^n (y_i - \mathbf{x}_i^T \boldsymbol{\beta})^2 = \arg \min (\mathbf{y} - \mathbf{X}\boldsymbol{\beta})^T (\mathbf{y} - \mathbf{X}\boldsymbol{\beta}) = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}$$

- Fitted values:  $\hat{\mathbf{y}} = \mathbf{X}\hat{\boldsymbol{\beta}} = \mathbf{X}(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y} = \mathbf{H}\mathbf{y}$
- Mean-Variance:  $\mathbb{E}(\hat{\mathbf{y}}) = \mathbf{H}\mathbb{E}(\mathbf{y}) = \mathbf{H}\mathbf{X}\boldsymbol{\beta} = \mathbf{X}\boldsymbol{\beta}, \quad \text{var}(\hat{\mathbf{y}}) = \sigma^2 \mathbf{H}$
- Residuals:  $\hat{\mathbf{e}} = (\mathbf{y} - \hat{\mathbf{y}}) = \mathbf{y} - \mathbf{X}\hat{\boldsymbol{\beta}} = \{\mathbf{I}_n - \mathbf{X}(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T\} \mathbf{y} = (\mathbf{I}_n - \mathbf{H})\mathbf{y}$
- Mean-Variance:  $\mathbb{E}(\hat{\mathbf{e}}) = (\mathbf{I}_n - \mathbf{H})\mathbf{X}\boldsymbol{\beta} = \mathbf{0}, \quad \text{var}(\hat{\mathbf{e}}) = \mathbb{E}(\hat{\mathbf{e}}\hat{\mathbf{e}}^T) = \sigma^2(\mathbf{I}_n - \mathbf{H})$
- $\text{cov}(\hat{\mathbf{y}}, \hat{\mathbf{e}}) = \text{cov}\{\mathbf{H}\mathbf{y}, (\mathbf{I}_n - \mathbf{H})\mathbf{y}\} = \sigma^2 \mathbf{H}(\mathbf{I}_n - \mathbf{H}) = \sigma^2(\mathbf{H} - \mathbf{H}^2) = \mathbf{0}$
- Let  $\hat{\sigma}^2 = (\mathbf{y} - \mathbf{X}\hat{\boldsymbol{\beta}})^T (\mathbf{y} - \mathbf{X}\hat{\boldsymbol{\beta}}) / n$ 
  - $\mathbb{E}(\hat{\sigma}^2) = \mathbb{E}\left\{(\mathbf{y} - \mathbf{X}\hat{\boldsymbol{\beta}})^T (\mathbf{y} - \mathbf{X}\hat{\boldsymbol{\beta}}) / n\right\} = \mathbb{E}(\hat{\mathbf{e}}^T \hat{\mathbf{e}}) / n = \mathbb{E}\{\text{trace}(\hat{\mathbf{e}}^T \hat{\mathbf{e}})\} / n$   
 $= \mathbb{E}\{\text{trace}(\hat{\mathbf{e}}\hat{\mathbf{e}}^T)\} / n = \text{trace}\{\mathbb{E}(\hat{\mathbf{e}}\hat{\mathbf{e}}^T)\} / n = \text{trace}\{\sigma^2(\mathbf{I}_n - \mathbf{H})\} / n = \sigma^2(n - p) / n$
  - $\mathbb{E}(s^2) = \mathbb{E}\{n\hat{\sigma}^2 / (n - p)\} = \mathbb{E}\left\{(\mathbf{y} - \mathbf{X}\hat{\boldsymbol{\beta}})^T (\mathbf{y} - \mathbf{X}\hat{\boldsymbol{\beta}}) / (n - p)\right\} = \sigma^2$

$$\mathbf{y}^{n \times 1} = \mathbf{X}^{n \times p} \boldsymbol{\beta}^{p \times 1} + \mathbf{e}^{n \times 1}, \quad \mathbb{E}(\mathbf{e}) = \mathbf{0}, \quad \text{cov}(\mathbf{e}) = \sigma^2 \mathbf{I}_n$$

- Minimize squared error loss:

$$\hat{\boldsymbol{\beta}} = \arg \min \sum_{i=1}^n (y_i - \mathbf{x}_i^T \boldsymbol{\beta})^2 = \arg \min (\mathbf{y} - \mathbf{X}\boldsymbol{\beta})^T (\mathbf{y} - \mathbf{X}\boldsymbol{\beta}) = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}$$

- Fitted values:  $\hat{\mathbf{y}} = \mathbf{X}\hat{\boldsymbol{\beta}} = \mathbf{X}(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y} = \mathbf{H}\mathbf{y}$
- Mean-Variance:  $\mathbb{E}(\hat{\mathbf{y}}) = \mathbf{H}\mathbb{E}(\mathbf{y}) = \mathbf{H}\mathbf{X}\boldsymbol{\beta} = \mathbf{X}\boldsymbol{\beta}, \quad \text{var}(\hat{\mathbf{y}}) = \sigma^2 \mathbf{H}$
- Residuals:  $\hat{\mathbf{e}} = (\mathbf{y} - \hat{\mathbf{y}}) = \mathbf{y} - \mathbf{X}\hat{\boldsymbol{\beta}} = \{\mathbf{I}_n - \mathbf{X}(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T\} \mathbf{y} = (\mathbf{I}_n - \mathbf{H})\mathbf{y}$
- Mean-Variance:  $\mathbb{E}(\hat{\mathbf{e}}) = (\mathbf{I}_n - \mathbf{H})\mathbf{X}\boldsymbol{\beta} = \mathbf{0}, \quad \text{var}(\hat{\mathbf{e}}) = \mathbb{E}(\hat{\mathbf{e}}\hat{\mathbf{e}}^T) = \sigma^2(\mathbf{I}_n - \mathbf{H})$
- $\text{cov}(\hat{\mathbf{y}}, \hat{\mathbf{e}}) = \text{cov}\{\mathbf{H}\mathbf{y}, (\mathbf{I}_n - \mathbf{H})\mathbf{y}\} = \sigma^2 \mathbf{H}(\mathbf{I}_n - \mathbf{H}) = \sigma^2(\mathbf{H} - \mathbf{H}^2) = \mathbf{0}$
- Let  $\hat{\sigma}^2 = (\mathbf{y} - \mathbf{X}\hat{\boldsymbol{\beta}})^T (\mathbf{y} - \mathbf{X}\hat{\boldsymbol{\beta}}) / n$
- $\mathbb{E}(\hat{\sigma}^2) = \mathbb{E}\left\{(\mathbf{y} - \mathbf{X}\hat{\boldsymbol{\beta}})^T (\mathbf{y} - \mathbf{X}\hat{\boldsymbol{\beta}}) / n\right\} = \mathbb{E}(\hat{\mathbf{e}}^T \hat{\mathbf{e}}) / n = \mathbb{E}\{\text{trace}(\hat{\mathbf{e}}^T \hat{\mathbf{e}})\} / n$   
 $= \mathbb{E}\{\text{trace}(\hat{\mathbf{e}}\hat{\mathbf{e}}^T)\} / n = \text{trace}\{\mathbb{E}(\hat{\mathbf{e}}\hat{\mathbf{e}}^T)\} / n = \text{trace}\{\sigma^2(\mathbf{I}_n - \mathbf{H})\} / n = \sigma^2(n - p) / n$
- $\mathbb{E}(s^2) = \mathbb{E}\{n\hat{\sigma}^2 / (n - p)\} = \mathbb{E}\left\{(\mathbf{y} - \mathbf{X}\hat{\boldsymbol{\beta}})^T (\mathbf{y} - \mathbf{X}\hat{\boldsymbol{\beta}}) / (n - p)\right\} = \sigma^2$

$$\mathbf{y}^{n \times 1} = \mathbf{X}^{n \times p} \boldsymbol{\beta}^{p \times 1} + \mathbf{e}^{n \times 1}, \quad \mathbb{E}(\mathbf{e}) = \mathbf{0}, \quad \text{cov}(\mathbf{e}) = \sigma^2 \mathbf{I}_n$$

- Minimize squared error loss:

$$\hat{\boldsymbol{\beta}} = \arg \min \sum_{i=1}^n (y_i - \mathbf{x}_i^T \boldsymbol{\beta})^2 = \arg \min (\mathbf{y} - \mathbf{X}\boldsymbol{\beta})^T (\mathbf{y} - \mathbf{X}\boldsymbol{\beta}) = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}$$

- Fitted values:  $\hat{\mathbf{y}} = \mathbf{X}\hat{\boldsymbol{\beta}} = \mathbf{X}(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y} = \mathbf{H}\mathbf{y}$
- Mean-Variance:  $\mathbb{E}(\hat{\mathbf{y}}) = \mathbf{H}\mathbb{E}(\mathbf{y}) = \mathbf{H}\mathbf{X}\boldsymbol{\beta} = \mathbf{X}\boldsymbol{\beta}, \quad \text{var}(\hat{\mathbf{y}}) = \sigma^2 \mathbf{H}$
- Residuals:  $\hat{\mathbf{e}} = (\mathbf{y} - \hat{\mathbf{y}}) = \mathbf{y} - \mathbf{X}\hat{\boldsymbol{\beta}} = \{\mathbf{I}_n - \mathbf{X}(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T\} \mathbf{y} = (\mathbf{I}_n - \mathbf{H})\mathbf{y}$
- Mean-Variance:  $\mathbb{E}(\hat{\mathbf{e}}) = (\mathbf{I}_n - \mathbf{H})\mathbf{X}\boldsymbol{\beta} = \mathbf{0}, \quad \text{var}(\hat{\mathbf{e}}) = \mathbb{E}(\hat{\mathbf{e}}\hat{\mathbf{e}}^T) = \sigma^2(\mathbf{I}_n - \mathbf{H})$
- $\text{cov}(\hat{\mathbf{y}}, \hat{\mathbf{e}}) = \text{cov}\{\mathbf{H}\mathbf{y}, (\mathbf{I}_n - \mathbf{H})\mathbf{y}\} = \sigma^2 \mathbf{H}(\mathbf{I}_n - \mathbf{H}) = \sigma^2(\mathbf{H} - \mathbf{H}^2) = \mathbf{0}$
- Let  $\hat{\sigma}^2 = (\mathbf{y} - \mathbf{X}\hat{\boldsymbol{\beta}})^T (\mathbf{y} - \mathbf{X}\hat{\boldsymbol{\beta}}) / n$
- $\mathbb{E}(\hat{\sigma}^2) = \mathbb{E}\left\{(\mathbf{y} - \mathbf{X}\hat{\boldsymbol{\beta}})^T (\mathbf{y} - \mathbf{X}\hat{\boldsymbol{\beta}}) / n\right\} = \mathbb{E}(\hat{\mathbf{e}}^T \hat{\mathbf{e}}) / n = \mathbb{E}\{\text{trace}(\hat{\mathbf{e}}^T \hat{\mathbf{e}})\} / n$   
 $= \mathbb{E}\{\text{trace}(\hat{\mathbf{e}}\hat{\mathbf{e}}^T)\} / n = \text{trace}\{\mathbb{E}(\hat{\mathbf{e}}\hat{\mathbf{e}}^T)\} / n = \text{trace}\{\sigma^2(\mathbf{I}_n - \mathbf{H})\} / n = \sigma^2(n - p) / n$
- $\mathbb{E}(s^2) = \mathbb{E}\{n\hat{\sigma}^2 / (n - p)\} = \mathbb{E}\left\{(\mathbf{y} - \mathbf{X}\hat{\boldsymbol{\beta}})^T (\mathbf{y} - \mathbf{X}\hat{\boldsymbol{\beta}}) / (n - p)\right\} = \sigma^2$

$$\mathbf{y}^{n \times 1} = \mathbf{X}^{n \times p} \boldsymbol{\beta}^{p \times 1} + \mathbf{e}^{n \times 1}, \quad \mathbb{E}(\mathbf{e}) = \mathbf{0}, \quad \text{cov}(\mathbf{e}) = \sigma^2 \mathbf{I}_n$$

$$\text{OLS: } \hat{\boldsymbol{\beta}} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}$$

- Fitted values:  $\hat{\mathbf{y}} = \mathbf{H}\mathbf{y}$ , and residuals:  $\hat{\mathbf{e}} = (\mathbf{y} - \hat{\mathbf{y}}) = (\mathbf{I}_n - \mathbf{H})\mathbf{y}$

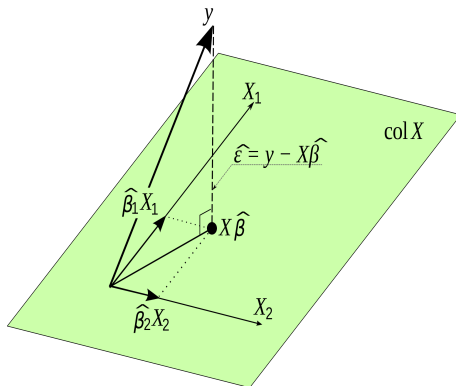
$$y = \beta_1 x_1 + \beta_2 x_2 + \epsilon$$

# Ordinary Least Squares

$$\mathbf{y}^{n \times 1} = \mathbf{X}^{n \times p} \boldsymbol{\beta}^{p \times 1} + \mathbf{e}^{n \times 1}, \quad \mathbb{E}(\mathbf{e}) = \mathbf{0}, \quad \text{cov}(\mathbf{e}) = \sigma^2 \mathbf{I}_n$$

$$\text{OLS: } \hat{\boldsymbol{\beta}} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}$$

- Fitted values:  $\hat{\mathbf{y}} = \mathbf{H}\mathbf{y}$ , and residuals:  $\hat{\mathbf{e}} = (\mathbf{y} - \hat{\mathbf{y}}) = (\mathbf{I}_n - \mathbf{H})\mathbf{y}$



$$y = \beta_1 x_1 + \beta_2 x_2 + \epsilon$$

$$\mathbf{y}^{n \times 1} = \mathbf{X}^{n \times p} \boldsymbol{\beta}^{p \times 1} + \mathbf{e}^{n \times 1}$$

- Minimize squared error loss:

$$\hat{\boldsymbol{\beta}}_{OLS} = \arg \min \sum_{i=1}^n (y_i - \mathbf{x}_i^T \boldsymbol{\beta})^2 = \arg \min (\mathbf{y} - \mathbf{X}\boldsymbol{\beta})^T (\mathbf{y} - \mathbf{X}\boldsymbol{\beta})$$

- $$\min_{\boldsymbol{\beta}} (\mathbf{y} - \mathbf{X}\boldsymbol{\beta})^T (\mathbf{y} - \mathbf{X}\boldsymbol{\beta}) \equiv \min_{\boldsymbol{\beta}} \left\{ \frac{1}{2\sigma^2} (\mathbf{y} - \mathbf{X}\boldsymbol{\beta})^T (\mathbf{y} - \mathbf{X}\boldsymbol{\beta}) \right\}$$
$$\equiv \max_{\boldsymbol{\beta}} \left\{ -\frac{n}{2} \log \sigma^2 - \frac{1}{2\sigma^2} (\mathbf{y} - \mathbf{X}\boldsymbol{\beta})^T (\mathbf{y} - \mathbf{X}\boldsymbol{\beta}) \right\}$$

- Maximize a normal log-likelihood function:

$$\max_{\boldsymbol{\beta}, \sigma^2} \mathcal{L}(\boldsymbol{\beta}, \sigma^2) \equiv \min_{\boldsymbol{\beta}, \sigma^2} \left\{ \frac{n}{2} \log \sigma^2 + \frac{1}{2\sigma^2} (\mathbf{y} - \mathbf{X}\boldsymbol{\beta})^T (\mathbf{y} - \mathbf{X}\boldsymbol{\beta}) \right\}$$

- MLE: 
$$\frac{\partial \mathcal{L}(\boldsymbol{\beta}, \sigma^2)}{\partial \boldsymbol{\beta}} = \frac{1}{\sigma^2} (\mathbf{X}^T \mathbf{y} - \mathbf{X}^T \mathbf{X} \boldsymbol{\beta}) = \mathbf{0}$$

$$\text{and } \frac{\partial \mathcal{L}(\boldsymbol{\beta}, \sigma^2)}{\partial \sigma^2} = -\frac{n}{2\sigma^2} + \frac{1}{2\sigma^4} (\mathbf{y} - \mathbf{X}\boldsymbol{\beta})^T (\mathbf{y} - \mathbf{X}\boldsymbol{\beta}) = 0$$

$$\Rightarrow \hat{\boldsymbol{\beta}}_{OLS} = \hat{\boldsymbol{\beta}}_{MLE} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}, \quad \hat{\sigma}_{MLE}^2 = (\mathbf{y} - \mathbf{X}\hat{\boldsymbol{\beta}})^T (\mathbf{y} - \mathbf{X}\hat{\boldsymbol{\beta}}) / n$$

$$\mathbf{y}^{n \times 1} = \mathbf{X}^{n \times p} \boldsymbol{\beta}^{p \times 1} + \mathbf{e}^{n \times 1}, \quad \mathbf{e} \sim \text{MVN}_n(\mathbf{0}, \sigma^2 \mathbf{I}_n)$$

- Minimize squared error loss:

$$\hat{\boldsymbol{\beta}}_{OLS} = \arg \min \sum_{i=1}^n (y_i - \mathbf{x}_i^T \boldsymbol{\beta})^2 = \arg \min (\mathbf{y} - \mathbf{X}\boldsymbol{\beta})^T (\mathbf{y} - \mathbf{X}\boldsymbol{\beta})$$

- $$\min_{\boldsymbol{\beta}} (\mathbf{y} - \mathbf{X}\boldsymbol{\beta})^T (\mathbf{y} - \mathbf{X}\boldsymbol{\beta}) \equiv \min_{\boldsymbol{\beta}} \left\{ \frac{1}{2\sigma^2} (\mathbf{y} - \mathbf{X}\boldsymbol{\beta})^T (\mathbf{y} - \mathbf{X}\boldsymbol{\beta}) \right\}$$
$$\equiv \max_{\boldsymbol{\beta}} \left\{ -\frac{n}{2} \log \sigma^2 - \frac{1}{2\sigma^2} (\mathbf{y} - \mathbf{X}\boldsymbol{\beta})^T (\mathbf{y} - \mathbf{X}\boldsymbol{\beta}) \right\}$$
- Maximize a normal log-likelihood function:

$$\max_{\boldsymbol{\beta}, \sigma^2} \mathcal{L}(\boldsymbol{\beta}, \sigma^2) \equiv \min_{\boldsymbol{\beta}, \sigma^2} \left\{ \frac{n}{2} \log \sigma^2 + \frac{1}{2\sigma^2} (\mathbf{y} - \mathbf{X}\boldsymbol{\beta})^T (\mathbf{y} - \mathbf{X}\boldsymbol{\beta}) \right\}$$

- MLE: 
$$\frac{\partial \mathcal{L}(\boldsymbol{\beta}, \sigma^2)}{\partial \boldsymbol{\beta}} = \frac{1}{\sigma^2} (\mathbf{X}^T \mathbf{y} - \mathbf{X}^T \mathbf{X} \boldsymbol{\beta}) = \mathbf{0}$$
and 
$$\frac{\partial \mathcal{L}(\boldsymbol{\beta}, \sigma^2)}{\partial \sigma^2} = -\frac{n}{2\sigma^2} + \frac{1}{2\sigma^4} (\mathbf{y} - \mathbf{X}\boldsymbol{\beta})^T (\mathbf{y} - \mathbf{X}\boldsymbol{\beta}) = 0$$
$$\Rightarrow \hat{\boldsymbol{\beta}}_{OLS} = \hat{\boldsymbol{\beta}}_{MLE} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}, \quad \hat{\sigma}_{MLE}^2 = (\mathbf{y} - \mathbf{X}\hat{\boldsymbol{\beta}})^T (\mathbf{y} - \mathbf{X}\hat{\boldsymbol{\beta}}) / n$$

$$\mathbf{y}^{n \times 1} = \mathbf{X}^{n \times p} \boldsymbol{\beta}^{p \times 1} + \mathbf{e}^{n \times 1}, \quad \mathbf{e} \sim \text{MVN}_n(\mathbf{0}, \sigma^2 \mathbf{I}_n)$$

- Minimize squared error loss:

$$\hat{\boldsymbol{\beta}}_{OLS} = \arg \min \sum_{i=1}^n (y_i - \mathbf{x}_i^T \boldsymbol{\beta})^2 = \arg \min (\mathbf{y} - \mathbf{X}\boldsymbol{\beta})^T (\mathbf{y} - \mathbf{X}\boldsymbol{\beta})$$

- $$\min_{\boldsymbol{\beta}} (\mathbf{y} - \mathbf{X}\boldsymbol{\beta})^T (\mathbf{y} - \mathbf{X}\boldsymbol{\beta}) \equiv \min_{\boldsymbol{\beta}} \left\{ \frac{1}{2\sigma^2} (\mathbf{y} - \mathbf{X}\boldsymbol{\beta})^T (\mathbf{y} - \mathbf{X}\boldsymbol{\beta}) \right\}$$
$$\equiv \max_{\boldsymbol{\beta}} \left\{ -\frac{n}{2} \log \sigma^2 - \frac{1}{2\sigma^2} (\mathbf{y} - \mathbf{X}\boldsymbol{\beta})^T (\mathbf{y} - \mathbf{X}\boldsymbol{\beta}) \right\}$$

- Maximize a normal log-likelihood function:

$$\max_{\boldsymbol{\beta}, \sigma^2} \mathcal{L}(\boldsymbol{\beta}, \sigma^2) \equiv \min_{\boldsymbol{\beta}, \sigma^2} \left\{ \frac{n}{2} \log \sigma^2 + \frac{1}{2\sigma^2} (\mathbf{y} - \mathbf{X}\boldsymbol{\beta})^T (\mathbf{y} - \mathbf{X}\boldsymbol{\beta}) \right\}$$

- MLE: 
$$\frac{\partial \mathcal{L}(\boldsymbol{\beta}, \sigma^2)}{\partial \boldsymbol{\beta}} = \frac{1}{\sigma^2} (\mathbf{X}^T \mathbf{y} - \mathbf{X}^T \mathbf{X} \boldsymbol{\beta}) = \mathbf{0}$$

$$\text{and } \frac{\partial \mathcal{L}(\boldsymbol{\beta}, \sigma^2)}{\partial \sigma^2} = -\frac{n}{2\sigma^2} + \frac{1}{2\sigma^4} (\mathbf{y} - \mathbf{X}\boldsymbol{\beta})^T (\mathbf{y} - \mathbf{X}\boldsymbol{\beta}) = 0$$

$$\Rightarrow \hat{\boldsymbol{\beta}}_{OLS} = \hat{\boldsymbol{\beta}}_{MLE} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}, \quad \hat{\sigma}_{MLE}^2 = (\mathbf{y} - \mathbf{X}\hat{\boldsymbol{\beta}})^T (\mathbf{y} - \mathbf{X}\hat{\boldsymbol{\beta}}) / n$$



$$\mathbf{y}^{n \times 1} = \mathbf{X}^{n \times p} \boldsymbol{\beta}^{p \times 1} + \mathbf{e}^{n \times 1}, \quad \mathbf{e} \sim \text{MVN}_n(\mathbf{0}, \sigma^2 \mathbf{I}_n)$$

$$\text{MLE: } \hat{\boldsymbol{\beta}} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}, \quad \hat{\sigma}^2 = \frac{1}{n} (\mathbf{y} - \mathbf{X} \hat{\boldsymbol{\beta}})^T (\mathbf{y} - \mathbf{X} \hat{\boldsymbol{\beta}})$$

- Log-likelihood:  $\mathcal{L}(\boldsymbol{\beta}, \sigma^2) = \left\{ -\frac{n}{2} \log \sigma^2 - \frac{1}{2\sigma^2} (\mathbf{y} - \mathbf{X}\boldsymbol{\beta})^T (\mathbf{y} - \mathbf{X}\boldsymbol{\beta}) \right\}$

- Gradients:  $\frac{\partial \mathcal{L}(\boldsymbol{\beta}, \sigma^2)}{\partial \boldsymbol{\beta}} = \frac{1}{\sigma^2} (\mathbf{X}^T \mathbf{y} - \mathbf{X}^T \mathbf{X} \boldsymbol{\beta})$

$$\frac{\partial \mathcal{L}(\boldsymbol{\beta}, \sigma^2)}{\partial \sigma^2} = -\frac{n}{2\sigma^2} + \frac{1}{2\sigma^4} (\mathbf{y} - \mathbf{X}\boldsymbol{\beta})^T (\mathbf{y} - \mathbf{X}\boldsymbol{\beta})$$

$$\mathbf{y}^{n \times 1} = \mathbf{X}^{n \times p} \boldsymbol{\beta}^{p \times 1} + \mathbf{e}^{n \times 1}, \quad \mathbf{e} \sim \text{MVN}_n(\mathbf{0}, \sigma^2 \mathbf{I}_n)$$

$$\text{MLE: } \hat{\boldsymbol{\beta}} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}, \quad \hat{\sigma}^2 = \frac{1}{n} (\mathbf{y} - \mathbf{X} \hat{\boldsymbol{\beta}})^T (\mathbf{y} - \mathbf{X} \hat{\boldsymbol{\beta}})$$

- Log-likelihood:  $\mathcal{L}(\boldsymbol{\beta}, \sigma^2) = \left\{ -\frac{n}{2} \log \sigma^2 - \frac{1}{2\sigma^2} (\mathbf{y} - \mathbf{X}\boldsymbol{\beta})^T (\mathbf{y} - \mathbf{X}\boldsymbol{\beta}) \right\}$

- Gradients:  $\frac{\partial \mathcal{L}(\boldsymbol{\beta}, \sigma^2)}{\partial \boldsymbol{\beta}} = \frac{1}{\sigma^2} (\mathbf{X}^T \mathbf{y} - \mathbf{X}^T \mathbf{X} \boldsymbol{\beta})$

$$\frac{\partial \mathcal{L}(\boldsymbol{\beta}, \sigma^2)}{\partial \sigma^2} = -\frac{n}{2\sigma^2} + \frac{1}{2\sigma^4} (\mathbf{y} - \mathbf{X}\boldsymbol{\beta})^T (\mathbf{y} - \mathbf{X}\boldsymbol{\beta})$$

- Hessian: 
$$\begin{aligned} -\frac{\partial^2 \mathcal{L}(\boldsymbol{\beta}, \sigma^2)}{\partial \boldsymbol{\beta} \partial \boldsymbol{\beta}^T} &= \frac{1}{\sigma^2} \mathbf{X}^T \mathbf{X} \\ -\frac{\partial \mathcal{L}(\boldsymbol{\beta}, \sigma^2)}{\partial (\sigma^2)^2} &= -\frac{n}{2\sigma^4} + \frac{1}{\sigma^6} (\mathbf{y} - \mathbf{X}\boldsymbol{\beta})^T (\mathbf{y} - \mathbf{X}\boldsymbol{\beta}) \\ -\frac{\partial \mathcal{L}(\boldsymbol{\beta}, \sigma^2)}{\partial \boldsymbol{\beta} \partial \sigma^2} &= \frac{-1}{\sigma^4} (\mathbf{X}^T \mathbf{y} - \mathbf{X}^T \mathbf{X} \boldsymbol{\beta}) \end{aligned}$$

- Fisher Information:  $\mathbf{I}(\boldsymbol{\beta}, \sigma^2) = \begin{bmatrix} \frac{1}{\sigma^2} \mathbf{X}^T \mathbf{X} & 0 \\ 0 & \frac{n}{2\sigma^4} \end{bmatrix}$

$$\mathbf{y}^{n \times 1} = \mathbf{X}^{n \times p} \boldsymbol{\beta}^{p \times 1} + \mathbf{e}^{n \times 1}, \quad \mathbf{e} \sim \text{MVN}_n(\mathbf{0}, \sigma^2 \mathbf{I}_n)$$

$$\text{MLE: } \hat{\boldsymbol{\beta}} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}, \quad \hat{\sigma}^2 = \frac{1}{n} (\mathbf{y} - \mathbf{X} \hat{\boldsymbol{\beta}})^T (\mathbf{y} - \mathbf{X} \hat{\boldsymbol{\beta}})$$

- Log-likelihood:  $\mathcal{L}(\boldsymbol{\beta}, \sigma^2) = \left\{ -\frac{n}{2} \log \sigma^2 - \frac{1}{2\sigma^2} (\mathbf{y} - \mathbf{X}\boldsymbol{\beta})^T (\mathbf{y} - \mathbf{X}\boldsymbol{\beta}) \right\}$

- Gradients:  $\frac{\partial \mathcal{L}(\boldsymbol{\beta}, \sigma^2)}{\partial \boldsymbol{\beta}} = \frac{1}{\sigma^2} (\mathbf{X}^T \mathbf{y} - \mathbf{X}^T \mathbf{X} \boldsymbol{\beta})$

$$\frac{\partial \mathcal{L}(\boldsymbol{\beta}, \sigma^2)}{\partial \sigma^2} = -\frac{n}{2\sigma^2} + \frac{1}{2\sigma^4} (\mathbf{y} - \mathbf{X}\boldsymbol{\beta})^T (\mathbf{y} - \mathbf{X}\boldsymbol{\beta})$$

- Hessian:  $-\frac{\partial^2 \mathcal{L}(\boldsymbol{\beta}, \sigma^2)}{\partial \boldsymbol{\beta} \partial \boldsymbol{\beta}^T} = \frac{1}{\sigma^2} \mathbf{X}^T \mathbf{X}$   
 $-\frac{\partial \mathcal{L}(\boldsymbol{\beta}, \sigma^2)}{\partial (\sigma^2)^2} = -\frac{n}{2\sigma^4} + \frac{1}{\sigma^6} (\mathbf{y} - \mathbf{X}\boldsymbol{\beta})^T (\mathbf{y} - \mathbf{X}\boldsymbol{\beta})$   
 $-\frac{\partial \mathcal{L}(\boldsymbol{\beta}, \sigma^2)}{\partial \boldsymbol{\beta} \partial \sigma^2} = \frac{-1}{\sigma^4} (\mathbf{X}^T \mathbf{y} - \mathbf{X}^T \mathbf{X} \boldsymbol{\beta})$

- Fisher Information:  $\mathbf{I}(\boldsymbol{\beta}, \sigma^2) = \begin{bmatrix} \frac{1}{\sigma^2} \mathbf{X}^T \mathbf{X} & 0 \\ 0 & \frac{n}{2\sigma^4} \end{bmatrix}$

$$\mathbf{y}^{n \times 1} = \mathbf{X}^{n \times p} \boldsymbol{\beta}^{p \times 1} + \mathbf{e}^{n \times 1}, \quad \mathbf{e} \sim \text{MVN}_n(\mathbf{0}, \sigma^2 \mathbf{I}_n)$$

$$\text{MLE: } \hat{\boldsymbol{\beta}} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}, \quad \hat{\sigma}^2 = \frac{1}{n} (\mathbf{y} - \mathbf{X} \hat{\boldsymbol{\beta}})^T (\mathbf{y} - \mathbf{X} \hat{\boldsymbol{\beta}})$$

- Log-likelihood:  $\mathcal{L}(\boldsymbol{\beta}, \sigma^2) = \left\{ -\frac{n}{2} \log \sigma^2 - \frac{1}{2\sigma^2} (\mathbf{y} - \mathbf{X}\boldsymbol{\beta})^T (\mathbf{y} - \mathbf{X}\boldsymbol{\beta}) \right\}$
- Fisher information:  $\mathbf{I}(\boldsymbol{\beta}, \sigma^2) = \begin{bmatrix} \frac{1}{\sigma^2} \mathbf{X}^T \mathbf{X} & 0 \\ 0 & \frac{n}{2\sigma^4} \end{bmatrix}$

$$\mathbf{y}^{n \times 1} = \mathbf{X}^{n \times p} \boldsymbol{\beta}^{p \times 1} + \mathbf{e}^{n \times 1}, \quad \mathbf{e} \sim \text{MVN}_n(\mathbf{0}, \sigma^2 \mathbf{I}_n)$$

$$\text{MLE: } \hat{\boldsymbol{\beta}} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}, \quad \hat{\sigma}^2 = \frac{1}{n} (\mathbf{y} - \mathbf{X} \hat{\boldsymbol{\beta}})^T (\mathbf{y} - \mathbf{X} \hat{\boldsymbol{\beta}})$$

- Log-likelihood:  $\mathcal{L}(\boldsymbol{\beta}, \sigma^2) = \left\{ -\frac{n}{2} \log \sigma^2 - \frac{1}{2\sigma^2} (\mathbf{y} - \mathbf{X}\boldsymbol{\beta})^T (\mathbf{y} - \mathbf{X}\boldsymbol{\beta}) \right\}$
- Fisher information:  $\mathbf{I}(\boldsymbol{\beta}, \sigma^2) = \begin{bmatrix} \frac{1}{\sigma^2} \mathbf{X}^T \mathbf{X} & 0 \\ 0 & \frac{n}{2\sigma^4} \end{bmatrix}$
- Asymptotic distribution:  
$$\begin{pmatrix} \hat{\boldsymbol{\beta}} \\ \hat{\sigma}^2 \end{pmatrix} \approx \text{MVN}_{p+1} \left[ \begin{pmatrix} \boldsymbol{\beta} \\ \sigma^2 \end{pmatrix}, \mathbf{I}(\boldsymbol{\beta}, \sigma^2)^{-1} = \begin{pmatrix} \sigma^2 (\mathbf{X}^T \mathbf{X})^{-1} & 0 \\ 0 & \frac{2\sigma^4}{n} \end{pmatrix} \right]$$

$$\mathbf{y}^{n \times 1} = \mathbf{X}^{n \times p} \boldsymbol{\beta}^{p \times 1} + \mathbf{e}^{n \times 1}, \quad \mathbf{e} \sim \text{MVN}_n(\mathbf{0}, \sigma^2 \mathbf{I}_n)$$

$$\text{OLS/MLE: } \hat{\boldsymbol{\beta}} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}$$

- What if  $\mathbf{X}^T \mathbf{X}$  is ill conditioned?

- Minimize penalized squared error loss:

$$\hat{\boldsymbol{\beta}}_{RR} = \arg \min \{ (\mathbf{y} - \mathbf{X}\boldsymbol{\beta})^T (\mathbf{y} - \mathbf{X}\boldsymbol{\beta}) + \lambda \boldsymbol{\beta}^T \boldsymbol{\beta} \}$$

- Ridge:  $\hat{\boldsymbol{\beta}} = (\mathbf{X}^T \mathbf{X} + \lambda \mathbf{I}_p)^{-1} \mathbf{X}^T \mathbf{y}$ ,  $\hat{\sigma}^2 = (\mathbf{y} - \mathbf{X}\hat{\boldsymbol{\beta}})^T (\mathbf{y} - \mathbf{X}\hat{\boldsymbol{\beta}}) / n$

- Mean-variance:  $\mathbb{E}(\hat{\boldsymbol{\beta}}) = (\mathbf{X}^T \mathbf{X} + \lambda \mathbf{I}_p)^{-1} \mathbf{X}^T \mathbf{X} \boldsymbol{\beta} \neq \boldsymbol{\beta}$ ,

$$\text{var}(\hat{\boldsymbol{\beta}}) = \sigma^2 (\mathbf{X}^T \mathbf{X} + \lambda \mathbf{I}_p)^{-1} \mathbf{X}^T \mathbf{X} (\mathbf{X}^T \mathbf{X} + \lambda \mathbf{I}_p)^{-1}$$

- Typically, cross-validation is used to determine  $\lambda$ .

$$\mathbf{y}^{n \times 1} = \mathbf{X}^{n \times p} \boldsymbol{\beta}^{p \times 1} + \mathbf{e}^{n \times 1}, \quad \mathbf{e} \sim \text{MVN}_n(\mathbf{0}, \sigma^2 \mathbf{I}_n)$$

$$\text{OLS/MLE: } \hat{\boldsymbol{\beta}} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}$$

- What if  $\mathbf{X}^T \mathbf{X}$  is ill conditioned?
- Minimize penalized squared error loss:

$$\hat{\boldsymbol{\beta}}_{RR} = \arg \min \{ (\mathbf{y} - \mathbf{X}\boldsymbol{\beta})^T (\mathbf{y} - \mathbf{X}\boldsymbol{\beta}) + \lambda \boldsymbol{\beta}^T \boldsymbol{\beta} \}$$

- Ridge:  $\hat{\boldsymbol{\beta}} = (\mathbf{X}^T \mathbf{X} + \lambda \mathbf{I}_p)^{-1} \mathbf{X}^T \mathbf{y}$ ,  $\hat{\sigma}^2 = (\mathbf{y} - \mathbf{X}\hat{\boldsymbol{\beta}})^T (\mathbf{y} - \mathbf{X}\hat{\boldsymbol{\beta}}) / n$
- Mean-variance:  $\mathbb{E}(\hat{\boldsymbol{\beta}}) = (\mathbf{X}^T \mathbf{X} + \lambda \mathbf{I}_p)^{-1} \mathbf{X}^T \mathbf{X} \boldsymbol{\beta} \neq \boldsymbol{\beta}$ ,  
 $\text{var}(\hat{\boldsymbol{\beta}}) = \sigma^2 (\mathbf{X}^T \mathbf{X} + \lambda \mathbf{I}_p)^{-1} \mathbf{X}^T \mathbf{X} (\mathbf{X}^T \mathbf{X} + \lambda \mathbf{I}_p)^{-1}$
- Typically, cross-validation is used to determine  $\lambda$ .

$$\mathbf{y}^{n \times 1} = \mathbf{X}^{n \times p} \boldsymbol{\beta}^{p \times 1} + \mathbf{e}^{n \times 1}, \quad \mathbf{e} \sim \text{MVN}_n(\mathbf{0}, \sigma^2 \mathbf{I}_n)$$

$$\text{OLS/MLE: } \hat{\boldsymbol{\beta}} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}$$

- What if  $\mathbf{X}^T \mathbf{X}$  is ill conditioned?
- Minimize penalized squared error loss:

$$\hat{\boldsymbol{\beta}}_{RR} = \arg \min \{ (\mathbf{y} - \mathbf{X}\boldsymbol{\beta})^T (\mathbf{y} - \mathbf{X}\boldsymbol{\beta}) + \lambda \boldsymbol{\beta}^T \boldsymbol{\beta} \}$$

- Ridge:  $\hat{\boldsymbol{\beta}} = (\mathbf{X}^T \mathbf{X} + \lambda \mathbf{I}_p)^{-1} \mathbf{X}^T \mathbf{y}$ ,  $\hat{\sigma}^2 = (\mathbf{y} - \mathbf{X}\hat{\boldsymbol{\beta}})^T (\mathbf{y} - \mathbf{X}\hat{\boldsymbol{\beta}}) / n$
- Mean-variance:  $\mathbb{E}(\hat{\boldsymbol{\beta}}) = (\mathbf{X}^T \mathbf{X} + \lambda \mathbf{I}_p)^{-1} \mathbf{X}^T \mathbf{X} \boldsymbol{\beta} \neq \boldsymbol{\beta}$ ,  
 $\text{var}(\hat{\boldsymbol{\beta}}) = \sigma^2 (\mathbf{X}^T \mathbf{X} + \lambda \mathbf{I}_p)^{-1} \mathbf{X}^T \mathbf{X} (\mathbf{X}^T \mathbf{X} + \lambda \mathbf{I}_p)^{-1}$
- Typically, cross-validation is used to determine  $\lambda$ .



$$\mathbf{y}^{n \times 1} = \mathbf{X}^{n \times p} \boldsymbol{\beta}^{p \times 1} + \mathbf{e}^{n \times 1}, \quad \mathbf{e} \sim \text{MVN}_n(\mathbf{0}, \sigma^2 \mathbf{I}_n)$$

$$\text{OLS/MLE: } \hat{\boldsymbol{\beta}} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}$$

- What if  $\mathbf{X}^T \mathbf{X}$  is ill conditioned?
- Minimize penalized squared error loss:

$$\hat{\boldsymbol{\beta}}_{RR} = \arg \min \{ (\mathbf{y} - \mathbf{X}\boldsymbol{\beta})^T (\mathbf{y} - \mathbf{X}\boldsymbol{\beta}) + \lambda \boldsymbol{\beta}^T \boldsymbol{\beta} \}$$

- Ridge:  $\hat{\boldsymbol{\beta}} = (\mathbf{X}^T \mathbf{X} + \lambda \mathbf{I}_p)^{-1} \mathbf{X}^T \mathbf{y}$ ,  $\hat{\sigma}^2 = (\mathbf{y} - \mathbf{X}\hat{\boldsymbol{\beta}})^T (\mathbf{y} - \mathbf{X}\hat{\boldsymbol{\beta}}) / n$
- Mean-variance:  $\mathbb{E}(\hat{\boldsymbol{\beta}}) = (\mathbf{X}^T \mathbf{X} + \lambda \mathbf{I}_p)^{-1} \mathbf{X}^T \mathbf{X} \boldsymbol{\beta} \neq \boldsymbol{\beta}$ ,  
 $\text{var}(\hat{\boldsymbol{\beta}}) = \sigma^2 (\mathbf{X}^T \mathbf{X} + \lambda \mathbf{I}_p)^{-1} \mathbf{X}^T \mathbf{X} (\mathbf{X}^T \mathbf{X} + \lambda \mathbf{I}_p)^{-1}$

- Typically, cross-validation is used to determine  $\lambda$ .

$$\mathbf{y}^{n \times 1} = \mathbf{X}^{n \times p} \boldsymbol{\beta}^{p \times 1} + \mathbf{e}^{n \times 1}, \quad \mathbf{e} \sim \text{MVN}_n(\mathbf{0}, \sigma^2 \mathbf{I}_n)$$

$$\text{OLS/MLE: } \hat{\boldsymbol{\beta}} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}$$

- What if  $\mathbf{X}^T \mathbf{X}$  is ill conditioned?
- Minimize penalized squared error loss:

$$\hat{\boldsymbol{\beta}}_{RR} = \arg \min \{ (\mathbf{y} - \mathbf{X}\boldsymbol{\beta})^T (\mathbf{y} - \mathbf{X}\boldsymbol{\beta}) + \lambda \boldsymbol{\beta}^T \boldsymbol{\beta} \}$$

- Ridge:  $\hat{\boldsymbol{\beta}} = (\mathbf{X}^T \mathbf{X} + \lambda \mathbf{I}_p)^{-1} \mathbf{X}^T \mathbf{y}$ ,  $\hat{\sigma}^2 = (\mathbf{y} - \mathbf{X}\hat{\boldsymbol{\beta}})^T (\mathbf{y} - \mathbf{X}\hat{\boldsymbol{\beta}}) / n$
- Mean-variance:  $\mathbb{E}(\hat{\boldsymbol{\beta}}) = (\mathbf{X}^T \mathbf{X} + \lambda \mathbf{I}_p)^{-1} \mathbf{X}^T \mathbf{X} \boldsymbol{\beta} \neq \boldsymbol{\beta}$ ,

$$\text{var}(\hat{\boldsymbol{\beta}}) = \sigma^2 (\mathbf{X}^T \mathbf{X} + \lambda \mathbf{I}_p)^{-1} \mathbf{X}^T \mathbf{X} (\mathbf{X}^T \mathbf{X} + \lambda \mathbf{I}_p)^{-1}$$

- Typically, cross-validation is used to determine  $\lambda$ .

## Least Absolute Shrinkage and Selection Operator (LASSO)

$$\mathbf{y}^{n \times 1} = \mathbf{X}^{n \times p} \boldsymbol{\beta}^{p \times 1} + \mathbf{e}^{n \times 1}, \quad \mathbf{e} \sim \text{MVN}_n(\mathbf{0}, \sigma^2 \mathbf{I}_n)$$

$$\text{OLS/MLE: } \hat{\boldsymbol{\beta}} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}$$

$$\text{Ridge: } \hat{\boldsymbol{\beta}} = (\mathbf{X}^T \mathbf{X} + \lambda \mathbf{I})^{-1} \mathbf{X}^T \mathbf{y}$$

- What if  $\boldsymbol{\beta}$  is sparse?

- Minimize penalized squared error loss:

$$\hat{\boldsymbol{\beta}}_{LASSO} = \arg \min \{ (\mathbf{y} - \mathbf{X}\boldsymbol{\beta})^T (\mathbf{y} - \mathbf{X}\boldsymbol{\beta}) + \lambda \|\boldsymbol{\beta}\|_1 \}$$

- Can be solved via the coordinate-wise gradient descent algorithm.
- Typically cross-validation is used to determine  $\lambda$ .

## Least Absolute Shrinkage and Selection Operator (LASSO)

$$\mathbf{y}^{n \times 1} = \mathbf{X}^{n \times p} \boldsymbol{\beta}^{p \times 1} + \mathbf{e}^{n \times 1}, \quad \mathbf{e} \sim \text{MVN}_n(\mathbf{0}, \sigma^2 \mathbf{I}_n)$$

$$\text{OLS/MLE: } \hat{\boldsymbol{\beta}} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}$$

$$\text{Ridge: } \hat{\boldsymbol{\beta}} = (\mathbf{X}^T \mathbf{X} + \lambda \mathbf{I})^{-1} \mathbf{X}^T \mathbf{y}$$

- What if  $\boldsymbol{\beta}$  is sparse?
- Minimize penalized squared error loss:

$$\hat{\boldsymbol{\beta}}_{LASSO} = \arg \min \{ (\mathbf{y} - \mathbf{X}\boldsymbol{\beta})^T (\mathbf{y} - \mathbf{X}\boldsymbol{\beta}) + \lambda \|\boldsymbol{\beta}\|_1 \}$$

- Can be solved via the coordinate-wise gradient descent algorithm.
- Typically cross-validation is used to determine  $\lambda$ .

## Least Absolute Shrinkage and Selection Operator (LASSO)

$$\mathbf{y}^{n \times 1} = \mathbf{X}^{n \times p} \boldsymbol{\beta}^{p \times 1} + \mathbf{e}^{n \times 1}, \quad \mathbf{e} \sim \text{MVN}_n(\mathbf{0}, \sigma^2 \mathbf{I}_n)$$

$$\text{OLS/MLE: } \hat{\boldsymbol{\beta}} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}$$

$$\text{Ridge: } \hat{\boldsymbol{\beta}} = (\mathbf{X}^T \mathbf{X} + \lambda \mathbf{I})^{-1} \mathbf{X}^T \mathbf{y}$$

- What if  $\boldsymbol{\beta}$  is sparse?
- Minimize penalized squared error loss:

$$\hat{\boldsymbol{\beta}}_{LASSO} = \arg \min \{ (\mathbf{y} - \mathbf{X}\boldsymbol{\beta})^T (\mathbf{y} - \mathbf{X}\boldsymbol{\beta}) + \lambda \|\boldsymbol{\beta}\|_1 \}$$

- Can be solved via the coordinate-wise gradient descent algorithm.
- Typically cross-validation is used to determine  $\lambda$ .

## Least Absolute Shrinkage and Selection Operator (LASSO)

$$\mathbf{y}^{n \times 1} = \mathbf{X}^{n \times p} \boldsymbol{\beta}^{p \times 1} + \mathbf{e}^{n \times 1}, \quad \mathbf{e} \sim \text{MVN}_n(\mathbf{0}, \sigma^2 \mathbf{I}_n)$$

$$\text{OLS/MLE: } \hat{\boldsymbol{\beta}} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}$$

$$\text{Ridge: } \hat{\boldsymbol{\beta}} = (\mathbf{X}^T \mathbf{X} + \lambda \mathbf{I})^{-1} \mathbf{X}^T \mathbf{y}$$

- What if  $\boldsymbol{\beta}$  is sparse?
- Minimize penalized squared error loss:

$$\hat{\boldsymbol{\beta}}_{LASSO} = \arg \min \{ (\mathbf{y} - \mathbf{X}\boldsymbol{\beta})^T (\mathbf{y} - \mathbf{X}\boldsymbol{\beta}) + \lambda \|\boldsymbol{\beta}\|_1 \}$$

- Can be solved via the coordinate-wise gradient descent algorithm.
- Typically cross-validation is used to determine  $\lambda$ .

$$\mathbf{y}^{n \times 1} = \mathbf{X}^{n \times p} \boldsymbol{\beta}^{p \times 1} + \mathbf{e}^{n \times 1}, \quad \mathbf{e} \sim \text{MVN}_n(\mathbf{0}, \sigma^2 \mathbf{I}_n)$$

$$\text{OLS/MLE: } \hat{\boldsymbol{\beta}}_{OLS} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}$$

$$\text{Ridge: } \hat{\boldsymbol{\beta}}_{RR} = (\mathbf{X}^T \mathbf{X} + \lambda \mathbf{I})^{-1} \mathbf{X}^T \mathbf{y}$$

$$\text{LASSO: } \hat{\boldsymbol{\beta}}_{LASSO} = \arg \min \{ (\mathbf{y} - \mathbf{X}\boldsymbol{\beta})^T (\mathbf{y} - \mathbf{X}\boldsymbol{\beta}) + \lambda \|\boldsymbol{\beta}\|_1 \}$$

- How to get the best of both worlds?

- Minimize  $L_1$  and  $L_2$  penalized squared error loss:

$$\hat{\boldsymbol{\beta}}_{ENET} = \arg \min [ (\mathbf{y} - \mathbf{X}\boldsymbol{\beta})^T (\mathbf{y} - \mathbf{X}\boldsymbol{\beta}) + \lambda \{ (1 - \alpha) \|\boldsymbol{\beta}\|_2^2 + 2\alpha \|\boldsymbol{\beta}\|_1 \} ]$$

- Typically cross-validation is used to determine  $\lambda$  for fixed values of  $\alpha$ .

$$\mathbf{y}^{n \times 1} = \mathbf{X}^{n \times p} \boldsymbol{\beta}^{p \times 1} + \mathbf{e}^{n \times 1}, \quad \mathbf{e} \sim \text{MVN}_n(\mathbf{0}, \sigma^2 \mathbf{I}_n)$$

$$\text{OLS/MLE: } \hat{\boldsymbol{\beta}}_{OLS} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}$$

$$\text{Ridge: } \hat{\boldsymbol{\beta}}_{RR} = (\mathbf{X}^T \mathbf{X} + \lambda \mathbf{I})^{-1} \mathbf{X}^T \mathbf{y}$$

$$\text{LASSO: } \hat{\boldsymbol{\beta}}_{LASSO} = \arg \min \{ (\mathbf{y} - \mathbf{X}\boldsymbol{\beta})^T (\mathbf{y} - \mathbf{X}\boldsymbol{\beta}) + \lambda \|\boldsymbol{\beta}\|_1 \}$$

- How to get the best of both worlds?
- Minimize  $L_1$  and  $L_2$  penalized squared error loss:

$$\hat{\boldsymbol{\beta}}_{ENET} = \arg \min [ (\mathbf{y} - \mathbf{X}\boldsymbol{\beta})^T (\mathbf{y} - \mathbf{X}\boldsymbol{\beta}) + \lambda \{ (1 - \alpha) \|\boldsymbol{\beta}\|_2^2 + 2\alpha \|\boldsymbol{\beta}\|_1 \} ]$$

- Typically cross-validation is used to determine  $\lambda$  for fixed values of  $\alpha$ .



$$\mathbf{y}^{n \times 1} = \mathbf{X}^{n \times p} \boldsymbol{\beta}^{p \times 1} + \mathbf{e}^{n \times 1}, \quad \mathbf{e} \sim \text{MVN}_n(\mathbf{0}, \sigma^2 \mathbf{I}_n)$$

$$\text{OLS/MLE: } \hat{\boldsymbol{\beta}}_{OLS} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}$$

$$\text{Ridge: } \hat{\boldsymbol{\beta}}_{RR} = (\mathbf{X}^T \mathbf{X} + \lambda \mathbf{I})^{-1} \mathbf{X}^T \mathbf{y}$$

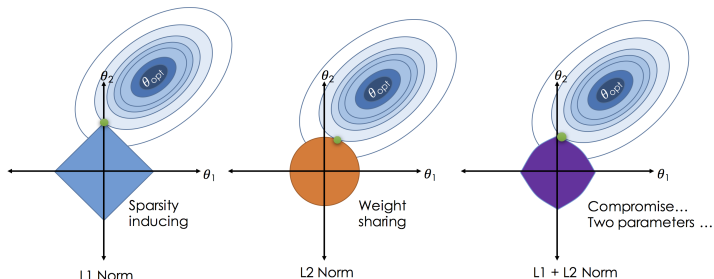
$$\text{LASSO: } \hat{\boldsymbol{\beta}}_{LASSO} = \arg \min \{ (\mathbf{y} - \mathbf{X}\boldsymbol{\beta})^T (\mathbf{y} - \mathbf{X}\boldsymbol{\beta}) + \lambda \|\boldsymbol{\beta}\|_1 \}$$

- How to get the best of both worlds?
- Minimize  $L_1$  and  $L_2$  penalized squared error loss:

$$\hat{\boldsymbol{\beta}}_{ENET} = \arg \min [ (\mathbf{y} - \mathbf{X}\boldsymbol{\beta})^T (\mathbf{y} - \mathbf{X}\boldsymbol{\beta}) + \lambda \{ (1 - \alpha) \|\boldsymbol{\beta}\|_2^2 + 2\alpha \|\boldsymbol{\beta}\|_1 \} ]$$

- Typically cross-validation is used to determine  $\lambda$  for fixed values of  $\alpha$ .

# Ridge vs LASSO vs Elastic Net

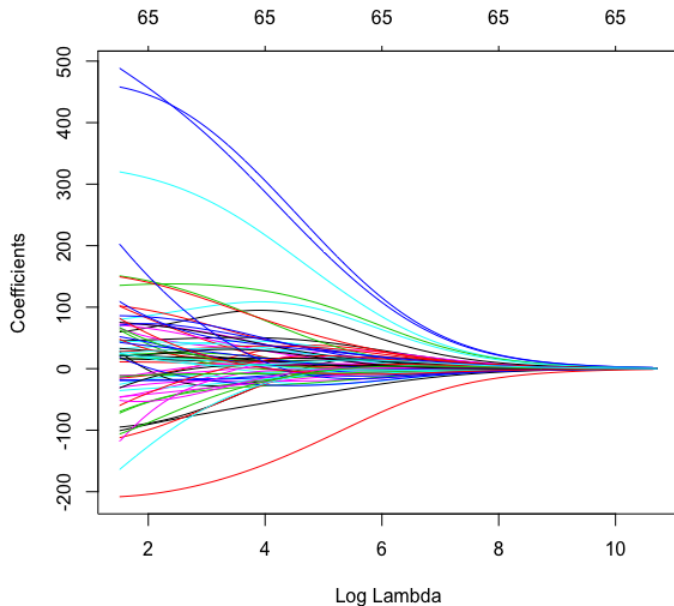


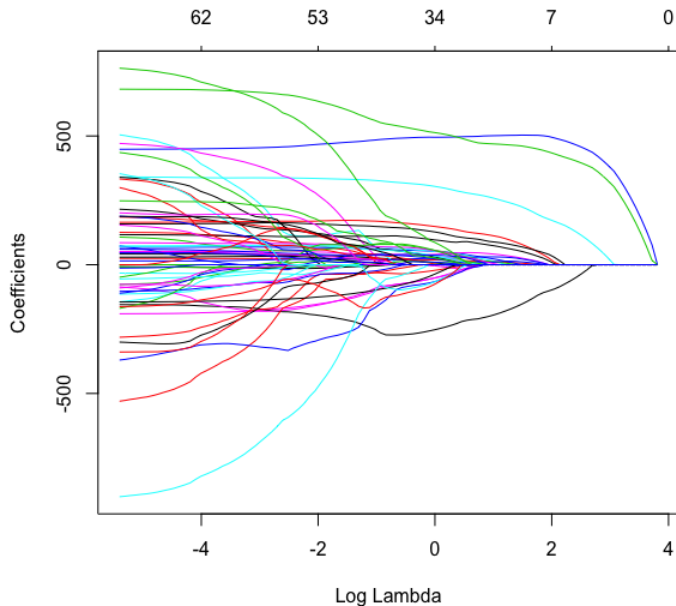
Contours of the error and constraint functions: the ellipses show the contours of the residual sum of squares (RSS), the regions around the origin show constraints corresponding to some given budget. The green dot is the smallest RSS that meets the budget.

- Data set with  $n = 442$  individuals and  $p = 64$  covariates.
- The response and the associated covariate values for the first six individuals and the first seven covariates are listed below.

	y	age	sex	bmi	map	tc	ldl	hdl
1	151	0.038075906	0.05068012	0.06169621	0.021872355	-0.044223498	-0.03482076	-0.043400846
2	75	-0.001882017	-0.04464164	-0.05147406	-0.026327835	-0.008448724	-0.01916334	0.074411564
3	141	0.085298906	0.05068012	0.04445121	-0.005670611	-0.045599451	-0.03419447	-0.032355932
4	206	-0.089062939	-0.04464164	-0.01159501	-0.036656447	0.012190569	0.02499059	-0.036037570
5	135	0.005383060	-0.04464164	-0.03638469	0.021872355	0.003934852	0.01559614	0.008142084
6	97	-0.092695478	-0.04464164	-0.04069594	-0.019442093	-0.068990650	-0.07928784	0.041276824

## Diabetes Data - Ridge





$$\mathbf{y}^{n \times 1} = \mathbf{X}^{n \times p} \boldsymbol{\beta}^{p \times 1} + \mathbf{e}^{n \times 1}, \quad \mathbf{e} \sim \text{MVN}_n(\mathbf{0}, \sigma^2 \mathbf{I}_n)$$

- Likelihood:  $L(\boldsymbol{\beta}, \sigma^2) = \frac{1}{(2\pi\sigma^2)^{n/2}} \exp \left\{ -\frac{1}{2\sigma^2} (\mathbf{y} - \mathbf{X}\boldsymbol{\beta})^T (\mathbf{y} - \mathbf{X}\boldsymbol{\beta}) \right\}$   
 $\propto \frac{1}{(\sigma^2)^{n/2}} \exp \left\{ -\frac{1}{2\sigma^2} (\boldsymbol{\beta}^T \mathbf{X}^T \mathbf{X} \boldsymbol{\beta} - 2\boldsymbol{\beta}^T \mathbf{X}^T \mathbf{y} + \mathbf{y}^T \mathbf{y}) \right\}$
- Semi-conjugate priors:  $p(\boldsymbol{\beta}, \sigma^2) = \text{MVN}(\boldsymbol{\mu}_\beta, \boldsymbol{\Sigma}_\beta) \cdot \text{Inv-Ga}(a_\sigma, b_\sigma)$   
 $\propto \frac{1}{(\sigma^2)^{a_\sigma+1}} \exp \left( -\frac{b_\sigma}{\sigma^2} \right) \exp \left\{ -\frac{1}{2} (\boldsymbol{\beta}^T \boldsymbol{\Sigma}_\beta^{-1} \boldsymbol{\beta} - 2\boldsymbol{\beta}^T \boldsymbol{\Sigma}_\beta^{-1} \boldsymbol{\mu}_\beta + \boldsymbol{\mu}_\beta^T \boldsymbol{\Sigma}_\beta^{-1} \boldsymbol{\mu}_\beta) \right\}$
- Posterior:  $p(\boldsymbol{\beta}, \sigma^2 \mid \mathbf{y}_{1:n}) \propto L(\boldsymbol{\beta}, \sigma^2) \times p(\boldsymbol{\beta}, \sigma^2)$
- Posterior full conditionals for block-Gibbs sampler:

$$p(\boldsymbol{\beta} \mid \sigma^2, \mathbf{y}) \propto \exp \left\{ -\frac{1}{2\sigma^2} (\boldsymbol{\beta}^T \mathbf{X}^T \mathbf{X} \boldsymbol{\beta} - 2\boldsymbol{\beta}^T \mathbf{X}^T \mathbf{y} + \mathbf{y}^T \mathbf{y}) \right\} \cdot \exp \left\{ -\frac{1}{2} (\boldsymbol{\beta}^T \boldsymbol{\Sigma}_\beta^{-1} \boldsymbol{\beta} - 2\boldsymbol{\beta}^T \boldsymbol{\Sigma}_\beta^{-1} \boldsymbol{\mu}_\beta + \boldsymbol{\mu}_\beta^T \boldsymbol{\Sigma}_\beta^{-1} \boldsymbol{\mu}_\beta) \right\}$$
$$= \text{MVN}(\boldsymbol{\beta}, \boldsymbol{\Sigma})$$
$$\boldsymbol{\Sigma}^{-1} = \mathbf{X}^T \mathbf{X} + \boldsymbol{\Sigma}_\beta^{-1}$$
$$\boldsymbol{\Sigma}^{-1} \boldsymbol{\mu} = \mathbf{X}^T \mathbf{y} + \boldsymbol{\Sigma}_\beta^{-1} \boldsymbol{\mu}_\beta$$

$$\mathbf{y}^{n \times 1} = \mathbf{X}^{n \times p} \boldsymbol{\beta}^{p \times 1} + \mathbf{e}^{n \times 1}, \quad \mathbf{e} \sim \text{MVN}_n(\mathbf{0}, \sigma^2 \mathbf{I}_n)$$

- Likelihood:  $L(\boldsymbol{\beta}, \sigma^2) = \frac{1}{(2\pi\sigma^2)^{n/2}} \exp \left\{ -\frac{1}{2\sigma^2} (\mathbf{y} - \mathbf{X}\boldsymbol{\beta})^T (\mathbf{y} - \mathbf{X}\boldsymbol{\beta}) \right\}$   
 $\propto \frac{1}{(\sigma^2)^{n/2}} \exp \left\{ -\frac{1}{2\sigma^2} (\boldsymbol{\beta}^T \mathbf{X}^T \mathbf{X} \boldsymbol{\beta} - 2\boldsymbol{\beta}^T \mathbf{X}^T \mathbf{y} + \mathbf{y}^T \mathbf{y}) \right\}$
- Semi-conjugate priors:  $p(\boldsymbol{\beta}, \sigma^2) = \text{MVN}(\boldsymbol{\mu}_\beta, \boldsymbol{\Sigma}_\beta) \cdot \text{Inv-Ga}(a_\sigma, b_\sigma)$   
 $\propto \exp \left[ -\frac{1}{2} \{ (\boldsymbol{\beta} - \boldsymbol{\mu}_\beta)^T \boldsymbol{\Sigma}_\beta^{-1} (\boldsymbol{\beta} - \boldsymbol{\mu}_\beta) \} \right] \frac{1}{(\sigma^2)^{a_\sigma+1}} \exp \left( -\frac{b_\sigma}{\sigma^2} \right)$   
 $\propto \frac{1}{(\sigma^2)^{a_\sigma+1}} \exp \left( -\frac{b_\sigma}{\sigma^2} \right) \exp \left\{ -\frac{1}{2} (\boldsymbol{\beta}^T \boldsymbol{\Sigma}_\beta^{-1} \boldsymbol{\beta} - 2\boldsymbol{\beta}^T \boldsymbol{\Sigma}_\beta^{-1} \boldsymbol{\mu}_\beta + \boldsymbol{\mu}_\beta^T \boldsymbol{\Sigma}_\beta^{-1} \boldsymbol{\mu}_\beta) \right\}$
- Posterior:  $p(\boldsymbol{\beta}, \sigma^2 \mid \mathbf{y}_{1:n}) \propto L(\boldsymbol{\beta}, \sigma^2) \times p(\boldsymbol{\beta}, \sigma^2)$
- Posterior full conditionals for block-Gibbs sampler:

$$\mathbf{y}^{n \times 1} = \mathbf{X}^{n \times p} \boldsymbol{\beta}^{p \times 1} + \mathbf{e}^{n \times 1}, \quad \mathbf{e} \sim \text{MVN}_n(\mathbf{0}, \sigma^2 \mathbf{I}_n)$$

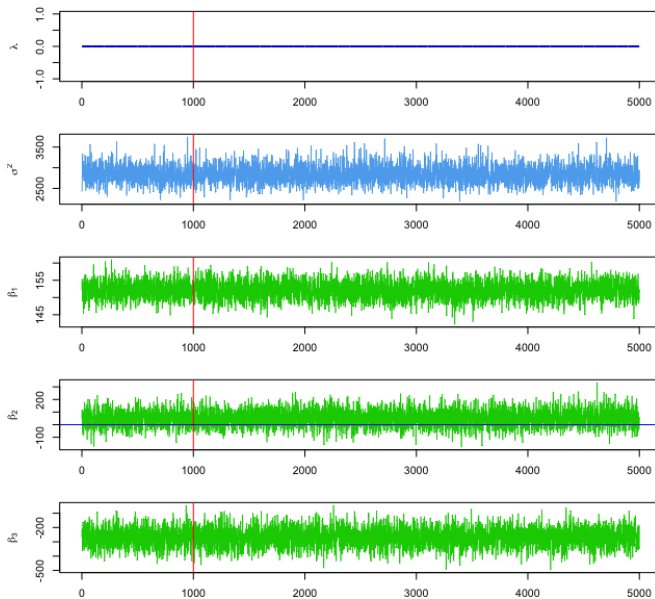
- Likelihood:  $L(\boldsymbol{\beta}, \sigma^2) = \frac{1}{(2\pi\sigma^2)^{n/2}} \exp \left\{ -\frac{1}{2\sigma^2} (\mathbf{y} - \mathbf{X}\boldsymbol{\beta})^T (\mathbf{y} - \mathbf{X}\boldsymbol{\beta}) \right\}$   
 $\propto \frac{1}{(\sigma^2)^{n/2}} \exp \left\{ -\frac{1}{2\sigma^2} (\boldsymbol{\beta}^T \mathbf{X}^T \mathbf{X} \boldsymbol{\beta} - 2\boldsymbol{\beta}^T \mathbf{X}^T \mathbf{y} + \mathbf{y}^T \mathbf{y}) \right\}$
- Semi-conjugate priors:  $p(\boldsymbol{\beta}, \sigma^2) = \text{MVN}(\boldsymbol{\mu}_\beta, \boldsymbol{\Sigma}_\beta) \cdot \text{Inv-Ga}(a_\sigma, b_\sigma)$   
 $\propto \frac{1}{(\sigma^2)^{a_\sigma+1}} \exp \left( -\frac{b_\sigma}{\sigma^2} \right) \exp \left\{ -\frac{1}{2} (\boldsymbol{\beta}^T \boldsymbol{\Sigma}_\beta^{-1} \boldsymbol{\beta} - 2\boldsymbol{\beta}^T \boldsymbol{\Sigma}_\beta^{-1} \boldsymbol{\mu}_\beta + \boldsymbol{\mu}_\beta^T \boldsymbol{\Sigma}_\beta^{-1} \boldsymbol{\mu}_\beta) \right\}$
- Posterior:  $p(\boldsymbol{\beta}, \sigma^2 \mid \mathbf{y}_{1:n}) \propto L(\boldsymbol{\beta}, \sigma^2) \times p(\boldsymbol{\beta}, \sigma^2)$
- Posterior full conditionals for block-Gibbs sampler:



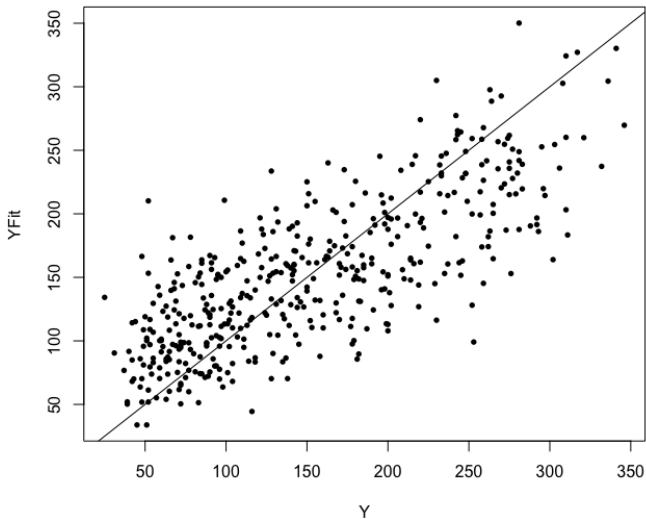
$$\mathbf{y}^{n \times 1} = \mathbf{X}^{n \times p} \boldsymbol{\beta}^{p \times 1} + \mathbf{e}^{n \times 1}, \quad \mathbf{e} \sim \text{MVN}_n(\mathbf{0}, \sigma^2 \mathbf{I}_n)$$

- Likelihood:  $L(\boldsymbol{\beta}, \sigma^2) = \frac{1}{(2\pi\sigma^2)^{n/2}} \exp \left\{ -\frac{1}{2\sigma^2} (\mathbf{y} - \mathbf{X}\boldsymbol{\beta})^T (\mathbf{y} - \mathbf{X}\boldsymbol{\beta}) \right\}$   
 $\propto \frac{1}{(\sigma^2)^{n/2}} \exp \left\{ -\frac{1}{2\sigma^2} (\boldsymbol{\beta}^T \mathbf{X}^T \mathbf{X} \boldsymbol{\beta} - 2\boldsymbol{\beta}^T \mathbf{X}^T \mathbf{y} + \mathbf{y}^T \mathbf{y}) \right\}$
- Semi-conjugate priors:  $p(\boldsymbol{\beta}, \sigma^2) = \text{MVN}(\boldsymbol{\mu}_\beta, \boldsymbol{\Sigma}_\beta) \cdot \text{Inv-Ga}(a_\sigma, b_\sigma)$   
 $\propto \frac{1}{(\sigma^2)^{a_\sigma+1}} \exp \left( -\frac{b_\sigma}{\sigma^2} \right) \exp \left\{ -\frac{1}{2} (\boldsymbol{\beta}^T \boldsymbol{\Sigma}_\beta^{-1} \boldsymbol{\beta} - 2\boldsymbol{\beta}^T \boldsymbol{\Sigma}_\beta^{-1} \boldsymbol{\mu}_\beta + \boldsymbol{\mu}_\beta^T \boldsymbol{\Sigma}_\beta^{-1} \boldsymbol{\mu}_\beta) \right\}$
- Posterior:  $p(\boldsymbol{\beta}, \sigma^2 \mid \mathbf{y}_{1:n}) \propto L(\boldsymbol{\beta}, \sigma^2) \times p(\boldsymbol{\beta}, \sigma^2)$
- Posterior full conditionals for block-Gibbs sampler:
  - $p(\boldsymbol{\beta} \mid -) \propto \exp \left[ -\frac{1}{2} \left\{ \boldsymbol{\beta}^T (\boldsymbol{\Sigma}_\beta^{-1} + \sigma^{-2} \mathbf{X}^T \mathbf{X}) \boldsymbol{\beta} - 2\boldsymbol{\beta}^T (\boldsymbol{\Sigma}_\beta^{-1} \boldsymbol{\mu}_\beta + \sigma^{-2} \mathbf{X}^T \mathbf{y}) \right\} \right]$   
 $\equiv \text{MVN}(\boldsymbol{\mu}_{\beta,n}, \boldsymbol{\Sigma}_{\beta,n}),$   
 $\boldsymbol{\Sigma}_{\beta,n} = (\boldsymbol{\Sigma}_\beta^{-1} + \sigma^{-2} \mathbf{X}^T \mathbf{X})^{-1}, \quad \boldsymbol{\mu}_{\beta,n} = \boldsymbol{\Sigma}_{\beta,n} (\boldsymbol{\Sigma}_\beta^{-1} \boldsymbol{\mu}_\beta + \sigma^{-2} \mathbf{X}^T \mathbf{y})$
  - $p(\sigma^2 \mid -) \propto \frac{1}{(\sigma^2)^{a_\sigma + \frac{n}{2} + 1}} \exp \left[ -\frac{1}{\sigma^2} \left\{ b_\sigma + \frac{1}{2} (\mathbf{y} - \mathbf{X}\boldsymbol{\beta})^T (\mathbf{y} - \mathbf{X}\boldsymbol{\beta}) \right\} \right]$   
 $\equiv \text{Inv-Ga} \left\{ a_\sigma + \frac{n}{2}, b_\sigma + \frac{1}{2} (\mathbf{y} - \mathbf{X}\boldsymbol{\beta})^T (\mathbf{y} - \mathbf{X}\boldsymbol{\beta}) \right\}$

# Diabetes Data - Bayesian Linear Models



## Diabetes Data - Bayesian Linear Models



$$\mathbf{y}^{n \times 1} = \mathbf{X}^{n \times p} \boldsymbol{\beta}^{p \times 1} + \mathbf{e}^{n \times 1}, \quad \mathbf{e} \sim \text{MVN}_n(\mathbf{0}, \sigma^2 \mathbf{I}_n)$$

- Ridge Regression:

- $\hat{\boldsymbol{\beta}}_{RR} = \arg \min \{(\mathbf{y} - \mathbf{X}\boldsymbol{\beta})^T (\mathbf{y} - \mathbf{X}\boldsymbol{\beta}) + \lambda \boldsymbol{\beta}^T \boldsymbol{\beta}\} = (\mathbf{X}^T \mathbf{X} + \lambda \mathbf{I}_p)^{-1} \mathbf{X}^T \mathbf{y}$

- An alternative view:

- 1.  $\hat{\boldsymbol{\beta}}_{RR} = \arg \min_{\boldsymbol{\beta}} \|\mathbf{y} - \mathbf{X}\boldsymbol{\beta}\|_2^2 + \lambda \|\boldsymbol{\beta}\|_2^2$
- 2.  $\hat{\boldsymbol{\beta}}_{RR} = \arg \min_{\boldsymbol{\beta}} \|\mathbf{y} - \mathbf{X}\boldsymbol{\beta}\|_2^2 + \lambda \|\boldsymbol{\beta}\|_2^2$
- 3.  $\hat{\boldsymbol{\beta}}_{RR} = \arg \min_{\boldsymbol{\beta}} \|\mathbf{y} - \mathbf{X}\boldsymbol{\beta}\|_2^2 + \lambda \|\boldsymbol{\beta}\|_2^2$

- 4.  $\hat{\boldsymbol{\beta}}_{RR} = \arg \min_{\boldsymbol{\beta}} \|\mathbf{y} - \mathbf{X}\boldsymbol{\beta}\|_2^2 + \lambda \|\boldsymbol{\beta}\|_2^2$
- 5.  $\hat{\boldsymbol{\beta}}_{RR} = \arg \min_{\boldsymbol{\beta}} \|\mathbf{y} - \mathbf{X}\boldsymbol{\beta}\|_2^2 + \lambda \|\boldsymbol{\beta}\|_2^2$

$$\mathbf{y}^{n \times 1} = \mathbf{X}^{n \times p} \boldsymbol{\beta}^{p \times 1} + \mathbf{e}^{n \times 1}, \quad \mathbf{e} \sim \text{MVN}_n(\mathbf{0}, \sigma^2 \mathbf{I}_n)$$

- Ridge Regression:

- $$\hat{\boldsymbol{\beta}}_{RR} = \arg \min \{ (\mathbf{y} - \mathbf{X}\boldsymbol{\beta})^T (\mathbf{y} - \mathbf{X}\boldsymbol{\beta}) + \lambda \boldsymbol{\beta}^T \boldsymbol{\beta} \} = (\mathbf{X}^T \mathbf{X} + \lambda \mathbf{I}_p)^{-1} \mathbf{X}^T \mathbf{y}$$

- An alternative view:

- $$\begin{aligned} \hat{\boldsymbol{\beta}}_{RR} &= \arg \min \left\{ \frac{1}{2\sigma^2} (\mathbf{y} - \mathbf{X}\boldsymbol{\beta})^T (\mathbf{y} - \mathbf{X}\boldsymbol{\beta}) + \frac{\lambda}{2\sigma^2} \boldsymbol{\beta}^T \boldsymbol{\beta} \right\} \\ &= \arg \max \left\{ -\frac{1}{2} (\mathbf{y} - \mathbf{X}\boldsymbol{\beta})^T (\sigma^2 \mathbf{I}_n)^{-1} (\mathbf{y} - \mathbf{X}\boldsymbol{\beta}) - \frac{1}{2} \boldsymbol{\beta}^T \left( \frac{\sigma^2}{\lambda} \mathbf{I}_p \right)^{-1} \boldsymbol{\beta} \right\} \\ &= \arg \max \{ \mathcal{L}(\boldsymbol{\beta} \mid \sigma^2) + \log p(\boldsymbol{\beta} \mid \lambda, \sigma^2) \} \end{aligned}$$

where  $p(\boldsymbol{\beta} \mid \lambda, \sigma^2) = \text{MVN} \left( \mathbf{0}, \frac{\sigma^2}{\lambda} \mathbf{I}_p \right)$

- $$\hat{\boldsymbol{\beta}}_{RR} = \hat{\boldsymbol{\beta}}_{MAP} = (\mathbf{X}^T \mathbf{X} + \lambda \mathbf{I}_p)^{-1} \mathbf{X}^T \mathbf{y} \quad \text{with} \quad p(\boldsymbol{\beta} \mid \lambda, \sigma^2) = \text{MVN} \left( \mathbf{0}, \frac{\sigma^2}{\lambda} \mathbf{I}_p \right).$$
- Full conditionals:

- $$\lambda \text{ can also be assigned } \text{Ga}(a_\lambda, b_\lambda) \text{ hyper-prior and sampled from its full conditional } p(\lambda \mid -) = \text{Ga} \left\{ a_\lambda + p/2, b_\lambda + \boldsymbol{\beta}^T \boldsymbol{\beta} / (2\sigma^2) \right\}.$$

$$\mathbf{y}^{n \times 1} = \mathbf{X}^{n \times p} \boldsymbol{\beta}^{p \times 1} + \mathbf{e}^{n \times 1}, \quad \mathbf{e} \sim \text{MVN}_n(\mathbf{0}, \sigma^2 \mathbf{I}_n)$$

- Ridge Regression:

- $\hat{\boldsymbol{\beta}}_{RR} = \arg \min \{ (\mathbf{y} - \mathbf{X}\boldsymbol{\beta})^T (\mathbf{y} - \mathbf{X}\boldsymbol{\beta}) + \lambda \boldsymbol{\beta}^T \boldsymbol{\beta} \} = (\mathbf{X}^T \mathbf{X} + \lambda \mathbf{I}_p)^{-1} \mathbf{X}^T \mathbf{y}$

- An alternative view:

- $\hat{\boldsymbol{\beta}}_{RR} = \arg \min \left\{ \frac{1}{2\sigma^2} (\mathbf{y} - \mathbf{X}\boldsymbol{\beta})^T (\mathbf{y} - \mathbf{X}\boldsymbol{\beta}) + \frac{\lambda}{2\sigma^2} \boldsymbol{\beta}^T \boldsymbol{\beta} \right\}$
  - $\hat{\boldsymbol{\beta}}_{RR} = \hat{\boldsymbol{\beta}}_{MAP} = (\mathbf{X}^T \mathbf{X} + \lambda \mathbf{I}_p)^{-1} \mathbf{X}^T \mathbf{y} \quad \text{with} \quad p(\boldsymbol{\beta} \mid \lambda, \sigma^2) = \text{MVN} \left( \mathbf{0}, \frac{\sigma^2}{\lambda} \mathbf{I}_p \right).$

- Full conditionals:

- $\lambda$  can also be assigned  $\text{Ga}(a_\lambda, b_\lambda)$  hyper-prior and sampled from its full conditional  $p(\lambda \mid -) = \text{Ga} \left\{ a_\lambda + p/2, b_\lambda + \boldsymbol{\beta}^T \boldsymbol{\beta} / (2\sigma^2) \right\}.$

$$\mathbf{y}^{n \times 1} = \mathbf{X}^{n \times p} \boldsymbol{\beta}^{p \times 1} + \mathbf{e}^{n \times 1}, \quad \mathbf{e} \sim \text{MVN}_n(\mathbf{0}, \sigma^2 \mathbf{I}_n)$$

- Ridge Regression:

- $\hat{\boldsymbol{\beta}}_{RR} = \arg \min \{(\mathbf{y} - \mathbf{X}\boldsymbol{\beta})^T (\mathbf{y} - \mathbf{X}\boldsymbol{\beta}) + \lambda \boldsymbol{\beta}^T \boldsymbol{\beta}\} = (\mathbf{X}^T \mathbf{X} + \lambda \mathbf{I}_p)^{-1} \mathbf{X}^T \mathbf{y}$

- An alternative view:

- $\hat{\boldsymbol{\beta}}_{RR} = \arg \min \left\{ \frac{1}{2\sigma^2} (\mathbf{y} - \mathbf{X}\boldsymbol{\beta})^T (\mathbf{y} - \mathbf{X}\boldsymbol{\beta}) + \frac{\lambda}{2\sigma^2} \boldsymbol{\beta}^T \boldsymbol{\beta} \right\}$
  - $\hat{\boldsymbol{\beta}}_{RR} = \hat{\boldsymbol{\beta}}_{MAP} = (\mathbf{X}^T \mathbf{X} + \lambda \mathbf{I}_p)^{-1} \mathbf{X}^T \mathbf{y} \quad \text{with} \quad p(\boldsymbol{\beta} \mid \lambda, \sigma^2) = \text{MVN} \left( \mathbf{0}, \frac{\sigma^2}{\lambda} \mathbf{I}_p \right).$

- Full conditionals:

- $p(\boldsymbol{\beta} \mid -) \equiv \text{MVN}(\boldsymbol{\mu}_{\boldsymbol{\beta},n}, \boldsymbol{\Sigma}_{\boldsymbol{\beta},n}),$

$$\boldsymbol{\Sigma}_{\boldsymbol{\beta},n} = \sigma^2 (\mathbf{X}^T \mathbf{X} + \lambda \mathbf{I}_p)^{-1}, \quad \boldsymbol{\mu}_{\boldsymbol{\beta},n} = \boldsymbol{\Sigma}_{\boldsymbol{\beta},n} (\sigma^{-2} \mathbf{X}^T \mathbf{y}) = (\mathbf{X}^T \mathbf{X} + \lambda \mathbf{I}_p)^{-1} \mathbf{X}^T \mathbf{y}$$

- $p(\sigma^2 \mid -) \equiv \text{Inv-Ga} \left\{ a_\sigma + \frac{n}{2}, b_\sigma + \frac{1}{2} (\mathbf{y} - \mathbf{X}\boldsymbol{\beta})^T (\mathbf{y} - \mathbf{X}\boldsymbol{\beta}) + \frac{\lambda}{2} \boldsymbol{\beta}^T \boldsymbol{\beta} \right\}$

- $\lambda$  can also be assigned  $\text{Ga}(a_\lambda, b_\lambda)$  hyper-prior and sampled from its full conditional

$$p(\lambda \mid -) = \text{Ga} \left\{ a_\lambda + p/2, b_\lambda + \boldsymbol{\beta}^T \boldsymbol{\beta} / (2\sigma^2) \right\}.$$

$$\mathbf{y}^{n \times 1} = \mathbf{X}^{n \times p} \boldsymbol{\beta}^{p \times 1} + \mathbf{e}^{n \times 1}, \quad \mathbf{e} \sim \text{MVN}_n(\mathbf{0}, \sigma^2 \mathbf{I}_n)$$

- Ridge Regression:

- $\hat{\boldsymbol{\beta}}_{RR} = \arg \min \{(\mathbf{y} - \mathbf{X}\boldsymbol{\beta})^T (\mathbf{y} - \mathbf{X}\boldsymbol{\beta}) + \lambda \boldsymbol{\beta}^T \boldsymbol{\beta}\} = (\mathbf{X}^T \mathbf{X} + \lambda \mathbf{I}_p)^{-1} \mathbf{X}^T \mathbf{y}$

- An alternative view:

- $\hat{\boldsymbol{\beta}}_{RR} = \arg \min \left\{ \frac{1}{2\sigma^2} (\mathbf{y} - \mathbf{X}\boldsymbol{\beta})^T (\mathbf{y} - \mathbf{X}\boldsymbol{\beta}) + \frac{\lambda}{2\sigma^2} \boldsymbol{\beta}^T \boldsymbol{\beta} \right\}$
  - $\hat{\boldsymbol{\beta}}_{RR} = \hat{\boldsymbol{\beta}}_{MAP} = (\mathbf{X}^T \mathbf{X} + \lambda \mathbf{I}_p)^{-1} \mathbf{X}^T \mathbf{y} \quad \text{with} \quad p(\boldsymbol{\beta} \mid \lambda, \sigma^2) = \text{MVN} \left( \mathbf{0}, \frac{\sigma^2}{\lambda} \mathbf{I}_p \right).$

- Full conditionals:

- $p(\boldsymbol{\beta} \mid -) \equiv \text{MVN}(\boldsymbol{\mu}_{\boldsymbol{\beta},n}, \boldsymbol{\Sigma}_{\boldsymbol{\beta},n}),$

$$\boldsymbol{\Sigma}_{\boldsymbol{\beta},n} = \sigma^2 (\mathbf{X}^T \mathbf{X} + \lambda \mathbf{I}_p)^{-1}, \quad \boldsymbol{\mu}_{\boldsymbol{\beta},n} = \boldsymbol{\Sigma}_{\boldsymbol{\beta},n} (\sigma^{-2} \mathbf{X}^T \mathbf{y}) = (\mathbf{X}^T \mathbf{X} + \lambda \mathbf{I}_p)^{-1} \mathbf{X}^T \mathbf{y}$$

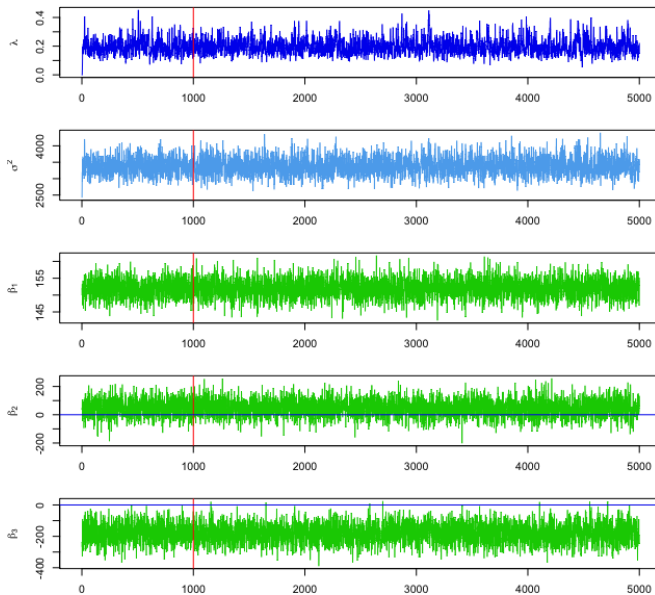
- $p(\sigma^2 \mid -) \equiv \text{Inv-Ga} \left\{ a_\sigma + \frac{n}{2}, b_\sigma + \frac{1}{2} (\mathbf{y} - \mathbf{X}\boldsymbol{\beta})^T (\mathbf{y} - \mathbf{X}\boldsymbol{\beta}) + \frac{\lambda}{2} \boldsymbol{\beta}^T \boldsymbol{\beta} \right\}$

- $\lambda$  can also be assigned  $\text{Ga}(a_\lambda, b_\lambda)$  hyper-prior and sampled from its full conditional

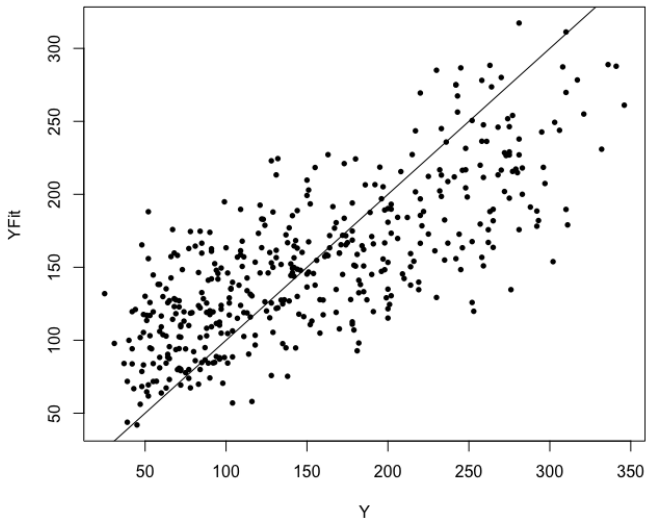
$$p(\lambda \mid -) = \text{Ga} \left\{ a_\lambda + p/2, b_\lambda + \boldsymbol{\beta}^T \boldsymbol{\beta} / (2\sigma^2) \right\}.$$



# Diabetes Data - Bayesian Ridge



## Diabetes Data - Bayesian Ridge



$$\mathbf{y}^{n \times 1} = \mathbf{X}^{n \times p} \boldsymbol{\beta}^{p \times 1} + \mathbf{e}^{n \times 1}, \quad \mathbf{e} \sim \text{MVN}_n(\mathbf{0}, \sigma^2 \mathbf{I}_n)$$

- LASSO Regression:

- $\hat{\boldsymbol{\beta}}_{LASSO} = \arg \min \{(\mathbf{y} - \mathbf{X}\boldsymbol{\beta})^T(\mathbf{y} - \mathbf{X}\boldsymbol{\beta}) + \lambda \|\boldsymbol{\beta}\|_1\}$

- An alternative view:

$$\hat{\boldsymbol{\beta}}_{LASSO} = \arg \min_{\boldsymbol{\beta}} \left\{ \frac{1}{2} (\mathbf{y} - \mathbf{X}\boldsymbol{\beta})^T (\mathbf{y} - \mathbf{X}\boldsymbol{\beta}) + \lambda \sum_{j=1}^p |\beta_j| \right\}$$
$$= \hat{\boldsymbol{\beta}}_{OLS} + \arg \min_{\boldsymbol{\beta}} \left\{ \frac{1}{2} (\mathbf{y} - \mathbf{X}\boldsymbol{\beta})^T (\mathbf{y} - \mathbf{X}\boldsymbol{\beta}) + \lambda \sum_{j=1}^p |\beta_j| \right\}$$

$$\mathbf{y}^{n \times 1} = \mathbf{X}^{n \times p} \boldsymbol{\beta}^{p \times 1} + \mathbf{e}^{n \times 1}, \quad \mathbf{e} \sim \text{MVN}_n(\mathbf{0}, \sigma^2 \mathbf{I}_n)$$

- LASSO Regression:

- $\hat{\boldsymbol{\beta}}_{LASSO} = \arg \min \{(\mathbf{y} - \mathbf{X}\boldsymbol{\beta})^T(\mathbf{y} - \mathbf{X}\boldsymbol{\beta}) + \lambda \|\boldsymbol{\beta}\|_1\}$

- An alternative view:

- $$\begin{aligned} \hat{\boldsymbol{\beta}}_{LASSO} &= \arg \min \left\{ \frac{1}{2\sigma^2} (\mathbf{y} - \mathbf{X}\boldsymbol{\beta})^T (\mathbf{y} - \mathbf{X}\boldsymbol{\beta}) + \frac{\lambda}{2\sigma^2} \sum_{j=1}^p |\beta_j| \right\} \\ &= \arg \max \left\{ -\frac{1}{2} (\mathbf{y} - \mathbf{X}\boldsymbol{\beta})^T (\sigma^2 \mathbf{I}_n)^{-1} (\mathbf{y} - \mathbf{X}\boldsymbol{\beta}) - \frac{\lambda}{2\sigma^2} \sum_{j=1}^p |\beta_j| \right\} \\ &= \arg \max \{ \mathcal{L}(\boldsymbol{\beta} \mid \sigma^2) + \log p(\boldsymbol{\beta} \mid \lambda, \sigma^2) \} \end{aligned}$$

where  $p(\boldsymbol{\beta} \mid \lambda, \sigma^2) = \prod_{j=1}^p \text{Laplace}(\beta_j \mid 0, \frac{2\sigma^2}{\lambda})$

- $\hat{\boldsymbol{\beta}}_{LASSO} = \hat{\boldsymbol{\beta}}_{MAP}$  with  $p(\boldsymbol{\beta} \mid \lambda, \sigma^2) = \prod_{j=1}^p \text{Laplace}(\beta_j \mid 0, \frac{2\sigma^2}{\lambda})$ .

$$\mathbf{y}^{n \times 1} = \mathbf{X}^{n \times p} \boldsymbol{\beta}^{p \times 1} + \mathbf{e}^{n \times 1}, \quad \mathbf{e} \sim \text{MVN}_n(\mathbf{0}, \sigma^2 \mathbf{I}_n)$$

- LASSO Regression:

- $\hat{\boldsymbol{\beta}}_{LASSO} = \arg \min \{ (\mathbf{y} - \mathbf{X}\boldsymbol{\beta})^T (\mathbf{y} - \mathbf{X}\boldsymbol{\beta}) + \lambda \|\boldsymbol{\beta}\|_1 \}$

- An alternative view:

- $\hat{\boldsymbol{\beta}}_{LASSO} = \arg \min \left\{ \frac{1}{2\sigma^2} (\mathbf{y} - \mathbf{X}\boldsymbol{\beta})^T (\mathbf{y} - \mathbf{X}\boldsymbol{\beta}) + \frac{\lambda}{2\sigma^2} \sum_{j=1}^p |\beta_j| \right\}$

- $\hat{\boldsymbol{\beta}}_{LASSO} = \hat{\boldsymbol{\beta}}_{MAP} \quad \text{with} \quad p(\boldsymbol{\beta} \mid \lambda, \sigma^2) = \prod_{j=1}^p \text{Laplace} \left( \beta_j \mid 0, \frac{2\sigma^2}{\lambda} \right).$

$$\mathbf{y}^{n \times 1} = \mathbf{X}^{n \times p} \boldsymbol{\beta}^{p \times 1} + \mathbf{e}^{n \times 1}, \quad \mathbf{e} \sim \text{MVN}_n(\mathbf{0}, \sigma^2 \mathbf{I}_n)$$

- LASSO Regression:

- $\hat{\boldsymbol{\beta}}_{LASSO} = \arg \min \{ (\mathbf{y} - \mathbf{X}\boldsymbol{\beta})^T (\mathbf{y} - \mathbf{X}\boldsymbol{\beta}) + \lambda \|\boldsymbol{\beta}\|_1 \}$

- Bayesian LASSO is motivated by the conditional Laplace prior

$$\begin{aligned} p(\boldsymbol{\beta} \mid \lambda, \sigma^2) &= \prod_{j=1}^p \text{Laplace}(\beta_j \mid 0, \frac{\sigma}{\lambda}) = \prod_{j=1}^p \frac{\lambda}{2\sqrt{\sigma^2}} \exp\left(-\lambda |\beta_j| / \sqrt{\sigma^2}\right) \\ &= \prod_{j=1}^p \int_0^\infty \frac{1}{\sqrt{2\pi\tau_j^2}} \exp\left(-\frac{\beta_j^2}{2\tau_j^2}\right) \frac{\lambda^2}{2\sigma^2} \exp\left(-\frac{\lambda^2 \tau_j^2}{2\sigma^2}\right) d\tau_j^2 \\ &= \prod_{j=1}^p \int_0^\infty \text{Normal}(\beta_j \mid 0, \tau_j^2) \text{Exp}\left(\tau_j^2 \mid \frac{\lambda^2}{2\sigma^2}\right) d\tau_j^2 \end{aligned}$$

$$\mathbf{y}^{n \times 1} = \mathbf{X}^{n \times p} \boldsymbol{\beta}^{p \times 1} + \mathbf{e}^{n \times 1}, \quad \mathbf{e} \sim \text{MVN}_n(\mathbf{0}, \sigma^2 \mathbf{I}_n)$$

- LASSO Regression:

- $\hat{\boldsymbol{\beta}}_{LASSO} = \arg \min \{(\mathbf{y} - \mathbf{X}\boldsymbol{\beta})^T(\mathbf{y} - \mathbf{X}\boldsymbol{\beta}) + \lambda \|\boldsymbol{\beta}\|_1\}$

- Bayesian LASSO is motivated by the conditional Laplace prior

$$p(\boldsymbol{\beta} \mid \lambda, \sigma^2) = \prod_{j=1}^p \frac{\lambda}{2\sqrt{\sigma^2}} \exp\left(-\frac{\lambda|\beta_j|}{\sqrt{\sigma^2}}\right) = \prod_{j=1}^p \int_0^\infty \text{Normal}(\beta_j \mid 0, \tau_j^2) \text{Exp}\left(\tau_j^2 \mid \frac{\lambda^2}{2\sigma^2}\right) d\tau_j^2$$

- Bayesian LASSO as a global-local prior

$$p(\boldsymbol{\beta} \mid \boldsymbol{\tau}^2) = \underbrace{\prod_{j=1}^p \text{Normal}(\beta_j \mid 0, \tau_j^2)}_{\text{Local components}}, \quad p(\boldsymbol{\tau}^2 \mid \lambda, \sigma^2) = \underbrace{\prod_{j=1}^p \text{Exp}\left(\tau_j^2 \mid \frac{\lambda^2}{2\sigma^2}\right)}_{\text{Global components}}$$

$$\mathbf{y}^{n \times 1} = \mathbf{X}^{n \times p} \boldsymbol{\beta}^{p \times 1} + \mathbf{e}^{n \times 1}, \quad \mathbf{e} \sim \text{MVN}_n(\mathbf{0}, \sigma^2 \mathbf{I}_n)$$

- LASSO Regression:

- $\hat{\boldsymbol{\beta}}_{LASSO} = \arg \min \{(\mathbf{y} - \mathbf{X}\boldsymbol{\beta})^T(\mathbf{y} - \mathbf{X}\boldsymbol{\beta}) + \lambda \|\boldsymbol{\beta}\|_1\}$

- Bayesian LASSO is motivated by the conditional Laplace prior

$$p(\boldsymbol{\beta} \mid \lambda, \sigma^2) = \prod_{j=1}^p \frac{\lambda}{2\sqrt{\sigma^2}} \exp\left(-\frac{\lambda|\beta_j|}{\sqrt{\sigma^2}}\right) = \prod_{j=1}^p \int_0^\infty \text{Normal}(\beta_j \mid 0, \tau_j^2) \text{Exp}\left(\tau_j^2 \mid \frac{\lambda^2}{2\sigma^2}\right) d\tau_j^2$$

- Posterior full conditionals for block-Gibbs sampler:

- $p(\boldsymbol{\beta} \mid -) = \text{MVN}(\boldsymbol{\mu}_{\boldsymbol{\beta},n}, \boldsymbol{\Sigma}_{\boldsymbol{\beta},n}), \quad \mathbf{D}_\tau = \text{diag}(\tau_1^2, \dots, \tau_p^2),$

$$\boldsymbol{\Sigma}_{\boldsymbol{\beta},n} = \sigma^2 (\mathbf{X}^T \mathbf{X} + \mathbf{D}_\tau^{-1})^{-1}, \quad \boldsymbol{\mu}_{\boldsymbol{\beta},n} = \boldsymbol{\Sigma}_{\boldsymbol{\beta},n} (\sigma^{-2} \mathbf{X}^T \mathbf{y}) = (\mathbf{X}^T \mathbf{X} + \mathbf{D}_\tau^{-1})^{-1} \mathbf{X}^T \mathbf{y}$$

- $p(\sigma^2 \mid -) = \text{Inv-Ga} \left\{ a_\sigma + \frac{n+p}{2}, b_\sigma + \frac{(\mathbf{y} - \mathbf{X}\boldsymbol{\beta})^T(\mathbf{y} - \mathbf{X}\boldsymbol{\beta})}{2} + \frac{\boldsymbol{\beta}^T \mathbf{D}_\tau^{-1} \boldsymbol{\beta}}{2} \right\}$

- $p(\tau_j^2 \mid -) \propto \frac{1}{\tau_j} \exp\left(-\frac{\beta_j^2}{2\tau_j^2}\right) \exp\left(-\frac{\lambda^2 \tau_j^2}{2\sigma^2}\right), \quad \tau_j^2 \rightarrow w_j = \tau_j^{-2}$

$$p(w_j \mid -) \propto w_j^{-3/2} \exp\left(-\frac{\beta_j^2 w_j}{2}\right) \exp\left(-\frac{\lambda^2}{2\sigma^2 w_j}\right) \equiv \text{Inv-Gs}(\mu', \lambda'), \quad \mu' = \frac{\lambda}{\sigma|\beta_j|}, \quad \lambda' = \frac{\lambda^2}{\sigma^2}.$$

- $\lambda^2$  can also be assigned  $\text{Ga}(a_\lambda, b_\lambda)$  hyper-prior and sampled from its full conditional

$$\begin{aligned} p(\lambda^2 \mid -) &\propto (\lambda^2)^{a_\lambda - 1} \exp(-\lambda^2 b_\lambda) \prod_{j=1}^p \left\{ \frac{\lambda^2}{2\sigma^2} \exp\left(-\frac{\lambda^2 \tau_j^2}{2\sigma^2}\right) \right\} \\ &\equiv \text{Ga} \left\{ a_\lambda + p, b_\lambda + \sum_{j=1}^p \tau_j^2 / (2\sigma^2) \right\}. \end{aligned}$$



$$\mathbf{y}^{n \times 1} = \mathbf{X}^{n \times p} \boldsymbol{\beta}^{p \times 1} + \mathbf{e}^{n \times 1}, \quad \mathbf{e} \sim \text{MVN}_n(\mathbf{0}, \sigma^2 \mathbf{I}_n)$$

- LASSO Regression:

- $\hat{\boldsymbol{\beta}}_{LASSO} = \arg \min \{(\mathbf{y} - \mathbf{X}\boldsymbol{\beta})^T(\mathbf{y} - \mathbf{X}\boldsymbol{\beta}) + \lambda \|\boldsymbol{\beta}\|_1\}$

- Bayesian LASSO is motivated by the conditional Laplace prior

$$p(\boldsymbol{\beta} \mid \lambda, \sigma^2) = \prod_{j=1}^p \frac{\lambda}{2\sqrt{\sigma^2}} \exp\left(-\frac{\lambda|\beta_j|}{\sqrt{\sigma^2}}\right) = \prod_{j=1}^p \int_0^\infty \text{Normal}(\beta_j \mid 0, \tau_j^2) \text{Exp}\left(\tau_j^2 \mid \frac{\lambda^2}{2\sigma^2}\right) d\tau_j^2$$

- Posterior full conditionals for block-Gibbs sampler:

- $p(\boldsymbol{\beta} \mid -) = \text{MVN}(\boldsymbol{\mu}_{\boldsymbol{\beta},n}, \boldsymbol{\Sigma}_{\boldsymbol{\beta},n}), \quad \mathbf{D}_\tau = \text{diag}(\tau_1^2, \dots, \tau_p^2),$

$$\boldsymbol{\Sigma}_{\boldsymbol{\beta},n} = \sigma^2 (\mathbf{X}^T \mathbf{X} + \mathbf{D}_\tau^{-1})^{-1}, \quad \boldsymbol{\mu}_{\boldsymbol{\beta},n} = \boldsymbol{\Sigma}_{\boldsymbol{\beta},n} (\sigma^{-2} \mathbf{X}^T \mathbf{y}) = (\mathbf{X}^T \mathbf{X} + \mathbf{D}_\tau^{-1})^{-1} \mathbf{X}^T \mathbf{y}$$

- $p(\sigma^2 \mid -) = \text{Inv-Ga} \left\{ a_\sigma + \frac{n+p}{2}, b_\sigma + \frac{(\mathbf{y} - \mathbf{X}\boldsymbol{\beta})^T(\mathbf{y} - \mathbf{X}\boldsymbol{\beta})}{2} + \frac{\boldsymbol{\beta}^T \mathbf{D}_\tau^{-1} \boldsymbol{\beta}}{2} \right\}$

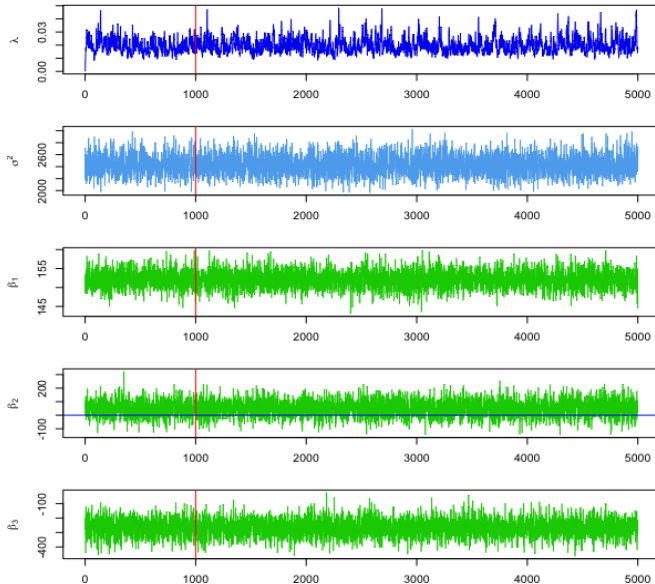
- $p(\tau_j^2 \mid -) \propto \frac{1}{\tau_j} \exp\left(-\frac{\beta_j^2}{2\tau_j^2}\right) \exp\left(-\frac{\lambda^2 \tau_j^2}{2\sigma^2}\right), \quad \tau_j^2 \rightarrow w_j = \tau_j^{-2}$

$$p(w_j \mid -) \propto w_j^{-3/2} \exp\left(-\frac{\beta_j^2 w_j}{2}\right) \exp\left(-\frac{\lambda^2}{2\sigma^2 w_j}\right) \equiv \text{Inv-Gs}(\mu', \lambda'), \quad \mu' = \frac{\lambda}{\sigma|\beta_j|}, \quad \lambda' = \frac{\lambda^2}{\sigma^2}.$$

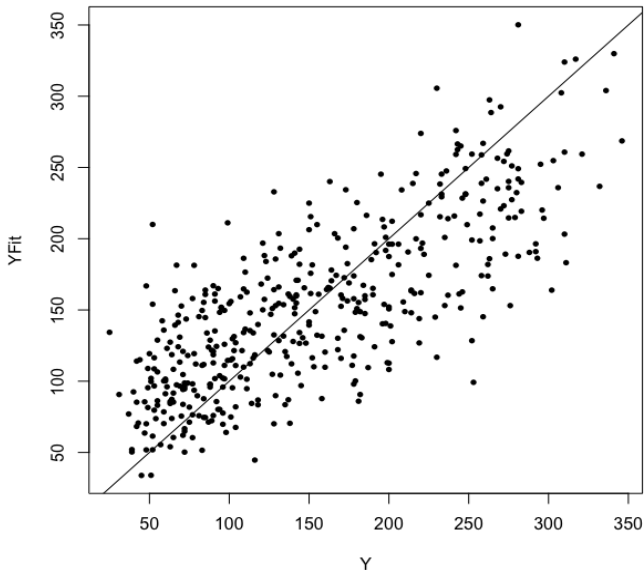
- $\lambda^2$  can also be assigned  $\text{Ga}(a_\lambda, b_\lambda)$  hyper-prior and sampled from its full conditional

$$\begin{aligned} p(\lambda^2 \mid -) &\propto (\lambda^2)^{a_\lambda - 1} \exp(-\lambda^2 b_\lambda) \prod_{j=1}^p \left\{ \frac{\lambda^2}{2\sigma^2} \exp\left(-\frac{\lambda^2 \tau_j^2}{2\sigma^2}\right) \right\} \\ &\equiv \text{Ga} \left\{ a_\lambda + p, b_\lambda + \sum_{j=1}^p \tau_j^2 / (2\sigma^2) \right\}. \end{aligned}$$

# Diabetes Data - Bayesian LASSO



## Diabetes Data - Bayesian LASSO



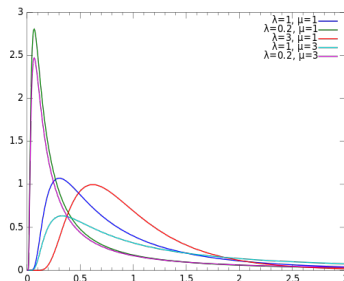
► **Inverse Gaussian Distribution:**

$$y \sim f(y \mid \mu, \lambda)$$

$$= \sqrt{\frac{\lambda}{2\pi y^3}} \exp\left(-\frac{\lambda(y - \mu)^2}{2\mu^2 y}\right)$$
$$y, \mu, \lambda \in \mathbb{R}^+$$

►  $\mathbb{E}(y) = \mu$

►  $\text{var}(y) = \frac{\mu^3}{\lambda}$



$$\mathbf{y}^{n \times 1} = \mathbf{X}^{n \times p} \boldsymbol{\beta}^{p \times 1} + \mathbf{e}^{n \times 1}, \quad \mathbb{E}(\mathbf{e}) = \mathbf{0}, \quad \text{cov}(\mathbf{e}) = \sigma^2 \mathbf{V}, \quad \mathbf{V} \text{ known.}$$

- Minimize squared error loss:

$$\hat{\boldsymbol{\beta}} = \arg \min (\mathbf{y} - \mathbf{X}\boldsymbol{\beta})^T \mathbf{V}^{-1} (\mathbf{y} - \mathbf{X}\boldsymbol{\beta}) = (\mathbf{X}^T \mathbf{V}^{-1} \mathbf{X})^{-1} \mathbf{X}^T \mathbf{V}^{-1} \mathbf{y}$$

- $\mathbb{E}(\hat{\boldsymbol{\beta}}) = (\mathbf{X}^T \mathbf{V}^{-1} \mathbf{X})^{-1} \mathbf{X}^T \mathbf{V}^{-1} \mathbf{X} \boldsymbol{\beta} = \boldsymbol{\beta}$
- $\text{var}(\hat{\boldsymbol{\beta}}) = (\mathbf{X}^T \mathbf{V}^{-1} \mathbf{X})^{-1} \mathbf{X}^T \mathbf{V}^{-1} \sigma^2 \mathbf{V} \mathbf{V}^{-1} \mathbf{X} (\mathbf{X}^T \mathbf{V}^{-1} \mathbf{X})^{-1} = \sigma^2 (\mathbf{X}^T \mathbf{V}^{-1} \mathbf{X})^{-1}$   
 $\Rightarrow \text{var}(\hat{\beta}_j) = \sigma^2 [(\mathbf{X}^T \mathbf{V}^{-1} \mathbf{X})^{-1}]_{jj} = \sigma^2 \frac{1}{\sum_{i=1}^n \frac{1}{V_{ii}} x_{ij}^2}$   
 $\Rightarrow \text{var}(\hat{\beta}_j) = \sigma^2 \frac{1}{\sum_{i=1}^n \frac{1}{V_{ii}} x_{ij}^2} = \sigma^2 \frac{1}{\sum_{i=1}^n \frac{1}{V_{ii}} x_{ij}^2} = \sigma^2 \frac{1}{\sum_{i=1}^n \frac{1}{V_{ii}} x_{ij}^2}$   
 $\Rightarrow \text{var}(\hat{\beta}_j) = \sigma^2 \frac{1}{\sum_{i=1}^n \frac{1}{V_{ii}} x_{ij}^2} = \sigma^2 \frac{1}{\sum_{i=1}^n \frac{1}{V_{ii}} x_{ij}^2} = \sigma^2 \frac{1}{\sum_{i=1}^n \frac{1}{V_{ii}} x_{ij}^2}$   
 $\Rightarrow \text{var}(\hat{\beta}_j) = \sigma^2 \frac{1}{\sum_{i=1}^n \frac{1}{V_{ii}} x_{ij}^2} = \sigma^2 \frac{1}{\sum_{i=1}^n \frac{1}{V_{ii}} x_{ij}^2} = \sigma^2 \frac{1}{\sum_{i=1}^n \frac{1}{V_{ii}} x_{ij}^2}$   
 $\Rightarrow \text{var}(\hat{\beta}_j) = \sigma^2 \frac{1}{\sum_{i=1}^n \frac{1}{V_{ii}} x_{ij}^2} = \sigma^2 \frac{1}{\sum_{i=1}^n \frac{1}{V_{ii}} x_{ij}^2} = \sigma^2 \frac{1}{\sum_{i=1}^n \frac{1}{V_{ii}} x_{ij}^2}$
- Fitted values:  $\hat{\mathbf{y}} = \mathbf{X}\hat{\boldsymbol{\beta}} = \mathbf{X}(\mathbf{X}^T \mathbf{V}^{-1} \mathbf{X})^{-1} \mathbf{X}^T \mathbf{V}^{-1} \mathbf{y} = \mathbf{H}_V \mathbf{y}$
- Mean-Variance:  $\mathbb{E}(\hat{\mathbf{y}}) = \mathbf{X}\mathbb{E}(\hat{\boldsymbol{\beta}}) = \mathbf{X}\boldsymbol{\beta}$ ,  $\text{var}(\hat{\mathbf{y}}) = \sigma^2 \mathbf{H}_V$
- Residuals:  $\hat{\mathbf{e}} = (\mathbf{y} - \hat{\mathbf{y}}) = \mathbf{y} - \mathbf{X}\hat{\boldsymbol{\beta}} = (\mathbf{I}_n - \mathbf{H}_V) \mathbf{y}$
- Mean-Variance:  $\mathbb{E}(\hat{\mathbf{e}}) = (\mathbf{I}_n - \mathbf{H}_V) \mathbf{X} \boldsymbol{\beta} = \mathbf{0}$ ,  $\text{var}(\hat{\mathbf{e}}) = \mathbb{E}(\hat{\mathbf{e}} \hat{\mathbf{e}}^T) = \sigma^2 (\mathbf{I}_n - \mathbf{H}_V)$
- $\text{cov}(\hat{\mathbf{y}}, \hat{\mathbf{e}}) = \text{cov}\{\mathbf{H}_V \mathbf{y}, (\mathbf{I}_n - \mathbf{H}_V) \mathbf{y}\} = \sigma^2 \mathbf{H}_V (\mathbf{I}_n - \mathbf{H}_V) = \sigma^2 (\mathbf{H}_V - \mathbf{H}_V^2) = \mathbf{0}$

$$\mathbf{y}^{n \times 1} = \mathbf{X}^{n \times p} \boldsymbol{\beta}^{p \times 1} + \mathbf{e}^{n \times 1}, \quad \mathbb{E}(\mathbf{e}) = \mathbf{0}, \quad \text{cov}(\mathbf{e}) = \sigma^2 \mathbf{V}, \quad \mathbf{V} \text{ known.}$$

- Minimize squared error loss:

$$\hat{\boldsymbol{\beta}} = \arg \min (\mathbf{y} - \mathbf{X}\boldsymbol{\beta})^T \mathbf{V}^{-1} (\mathbf{y} - \mathbf{X}\boldsymbol{\beta}) = (\mathbf{X}^T \mathbf{V}^{-1} \mathbf{X})^{-1} \mathbf{X}^T \mathbf{V}^{-1} \mathbf{y}$$

- $\mathbb{E}(\hat{\boldsymbol{\beta}}) = (\mathbf{X}^T \mathbf{V}^{-1} \mathbf{X})^{-1} \mathbf{X}^T \mathbf{V}^{-1} \mathbf{X} \boldsymbol{\beta} = \boldsymbol{\beta}$

- $\text{var}(\hat{\boldsymbol{\beta}}) = (\mathbf{X}^T \mathbf{V}^{-1} \mathbf{X})^{-1} \mathbf{X}^T \mathbf{V}^{-1} \sigma^2 \mathbf{V} \mathbf{V}^{-1} \mathbf{X} (\mathbf{X}^T \mathbf{V}^{-1} \mathbf{X})^{-1} = \sigma^2 (\mathbf{X}^T \mathbf{V}^{-1} \mathbf{X})^{-1}$

- Fitted values:  $\hat{\mathbf{y}} = \mathbf{X}\hat{\boldsymbol{\beta}} = \mathbf{X}(\mathbf{X}^T \mathbf{V}^{-1} \mathbf{X})^{-1} \mathbf{X}^T \mathbf{V}^{-1} \mathbf{y} = \mathbf{H}_V \mathbf{y}$

- Mean-Variance:  $\mathbb{E}(\hat{\mathbf{y}}) = \mathbf{X}\mathbb{E}(\hat{\boldsymbol{\beta}}) = \mathbf{X}\boldsymbol{\beta}$ ,  $\text{var}(\hat{\mathbf{y}}) = \sigma^2 \mathbf{H}_V$

- Residuals:  $\hat{\mathbf{e}} = (\mathbf{y} - \hat{\mathbf{y}}) = \mathbf{y} - \mathbf{X}\hat{\boldsymbol{\beta}} = (\mathbf{I}_n - \mathbf{H}_V) \mathbf{y}$

- Mean-Variance:  $\mathbb{E}(\hat{\mathbf{e}}) = (\mathbf{I}_n - \mathbf{H}_V) \mathbf{X} \boldsymbol{\beta} = \mathbf{0}$ ,  $\text{var}(\hat{\mathbf{e}}) = \mathbb{E}(\hat{\mathbf{e}} \hat{\mathbf{e}}^T) = \sigma^2 (\mathbf{I}_n - \mathbf{H}_V)$

- $\text{cov}(\hat{\mathbf{y}}, \hat{\mathbf{e}}) = \text{cov}\{\mathbf{H}_V \mathbf{y}, (\mathbf{I}_n - \mathbf{H}_V) \mathbf{y}\} = \sigma^2 \mathbf{H}_V (\mathbf{I}_n - \mathbf{H}_V) = \sigma^2 (\mathbf{H}_V - \mathbf{H}_V^2) = \mathbf{0}$

## Generalized / Weighted Least Squares

$$\mathbf{y}^{n \times 1} = \mathbf{X}^{n \times p} \boldsymbol{\beta}^{p \times 1} + \mathbf{e}^{n \times 1}, \quad \mathbb{E}(\mathbf{e}) = \mathbf{0}, \quad \text{cov}(\mathbf{e}) = \sigma^2 \mathbf{V}, \quad \mathbf{V} \text{ known.}$$

- Minimize squared error loss:

$$\hat{\boldsymbol{\beta}} = \arg \min (\mathbf{y} - \mathbf{X}\boldsymbol{\beta})^T \mathbf{V}^{-1} (\mathbf{y} - \mathbf{X}\boldsymbol{\beta}) = (\mathbf{X}^T \mathbf{V}^{-1} \mathbf{X})^{-1} \mathbf{X}^T \mathbf{V}^{-1} \mathbf{y}$$

- $\mathbb{E}(\hat{\boldsymbol{\beta}}) = (\mathbf{X}^T \mathbf{V}^{-1} \mathbf{X})^{-1} \mathbf{X}^T \mathbf{V}^{-1} \mathbf{X} \boldsymbol{\beta} = \boldsymbol{\beta}$
- $\text{var}(\hat{\boldsymbol{\beta}}) = (\mathbf{X}^T \mathbf{V}^{-1} \mathbf{X})^{-1} \mathbf{X}^T \mathbf{V}^{-1} \sigma^2 \mathbf{V} \mathbf{V}^{-1} \mathbf{X} (\mathbf{X}^T \mathbf{V}^{-1} \mathbf{X})^{-1} = \sigma^2 (\mathbf{X}^T \mathbf{V}^{-1} \mathbf{X})^{-1}$ 
  - Hat matrix:  $\mathbf{H}_V = \mathbf{X} (\mathbf{X}^T \mathbf{V}^{-1} \mathbf{X})^{-1} \mathbf{X}^T \mathbf{V}^{-1}$
  - Idempotent:  $\mathbf{H}_V^2 = \mathbf{H}_V, \quad (\mathbf{I}_n - \mathbf{H}_V)^2 = (\mathbf{I}_n - \mathbf{H}_V)$
  - $\text{trace}(\mathbf{H}_V) = \text{trace}\{\mathbf{X} (\mathbf{X}^T \mathbf{V}^{-1} \mathbf{X})^{-1} \mathbf{X}^T \mathbf{V}^{-1}\} = \text{trace}\{\mathbf{X}^T \mathbf{V}^{-1} \mathbf{X} (\mathbf{X}^T \mathbf{V}^{-1} \mathbf{X})^{-1}\} = p$
  - $\text{trace}(\mathbf{I}_n - \mathbf{H}_V) = \text{trace}(\mathbf{I}_n) - \text{trace}(\mathbf{H}_V) = (n - p)$
  - $\mathbf{H}_V \mathbf{X} = \mathbf{X}$  and  $(\mathbf{I}_n - \mathbf{H}_V) \mathbf{X} = \{\mathbf{I}_n - \mathbf{X} (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T\} \mathbf{X} = \mathbf{0}$
- Fitted values:  $\hat{\mathbf{y}} = \mathbf{X} \hat{\boldsymbol{\beta}} = \mathbf{X} (\mathbf{X}^T \mathbf{V}^{-1} \mathbf{X})^{-1} \mathbf{X}^T \mathbf{V}^{-1} \mathbf{y} = \mathbf{H}_V \mathbf{y}$
- Mean-Variance:  $\mathbb{E}(\hat{\mathbf{y}}) = \mathbf{X} \mathbb{E}(\hat{\boldsymbol{\beta}}) = \mathbf{X} \boldsymbol{\beta}, \quad \text{var}(\hat{\mathbf{y}}) = \sigma^2 \mathbf{H}_V$
- Residuals:  $\hat{\mathbf{e}} = (\mathbf{y} - \hat{\mathbf{y}}) = \mathbf{y} - \mathbf{X} \hat{\boldsymbol{\beta}} = (\mathbf{I}_n - \mathbf{H}_V) \mathbf{y}$
- Mean-Variance:  $\mathbb{E}(\hat{\mathbf{e}}) = (\mathbf{I}_n - \mathbf{H}_V) \mathbf{X} \boldsymbol{\beta} = \mathbf{0}, \quad \text{var}(\hat{\mathbf{e}}) = \mathbb{E}(\hat{\mathbf{e}} \hat{\mathbf{e}}^T) = \sigma^2 (\mathbf{I}_n - \mathbf{H}_V)$
- $\text{cov}(\hat{\mathbf{y}}, \hat{\mathbf{e}}) = \text{cov}\{\mathbf{H}_V \mathbf{y}, (\mathbf{I}_n - \mathbf{H}_V) \mathbf{y}\} = \sigma^2 \mathbf{H}_V (\mathbf{I}_n - \mathbf{H}_V) = \sigma^2 (\mathbf{H}_V - \mathbf{H}_V^2) = \mathbf{0}$

## Generalized / Weighted Least Squares

$$\mathbf{y}^{n \times 1} = \mathbf{X}^{n \times p} \boldsymbol{\beta}^{p \times 1} + \mathbf{e}^{n \times 1}, \quad \mathbb{E}(\mathbf{e}) = \mathbf{0}, \quad \text{cov}(\mathbf{e}) = \sigma^2 \mathbf{V}, \quad \mathbf{V} \text{ known.}$$

- Minimize squared error loss:

$$\hat{\boldsymbol{\beta}} = \arg \min (\mathbf{y} - \mathbf{X}\boldsymbol{\beta})^T \mathbf{V}^{-1} (\mathbf{y} - \mathbf{X}\boldsymbol{\beta}) = (\mathbf{X}^T \mathbf{V}^{-1} \mathbf{X})^{-1} \mathbf{X}^T \mathbf{V}^{-1} \mathbf{y}$$

- $\mathbb{E}(\hat{\boldsymbol{\beta}}) = (\mathbf{X}^T \mathbf{V}^{-1} \mathbf{X})^{-1} \mathbf{X}^T \mathbf{V}^{-1} \mathbf{X} \boldsymbol{\beta} = \boldsymbol{\beta}$
- $\text{var}(\hat{\boldsymbol{\beta}}) = (\mathbf{X}^T \mathbf{V}^{-1} \mathbf{X})^{-1} \mathbf{X}^T \mathbf{V}^{-1} \sigma^2 \mathbf{V} \mathbf{V}^{-1} \mathbf{X} (\mathbf{X}^T \mathbf{V}^{-1} \mathbf{X})^{-1} = \sigma^2 (\mathbf{X}^T \mathbf{V}^{-1} \mathbf{X})^{-1}$ 
  - Hat matrix:  $\mathbf{H}_V = \mathbf{X}(\mathbf{X}^T \mathbf{V}^{-1} \mathbf{X})^{-1} \mathbf{X}^T \mathbf{V}^{-1}$** 
    - Idempotent:  $\mathbf{H}_V^2 = \mathbf{H}_V, \quad (\mathbf{I}_n - \mathbf{H}_V)^2 = (\mathbf{I}_n - \mathbf{H}_V)$
    - $\text{trace}(\mathbf{H}_V) = \text{trace}\{\mathbf{X}(\mathbf{X}^T \mathbf{V}^{-1} \mathbf{X})^{-1} \mathbf{X}^T \mathbf{V}^{-1}\} = \text{trace}\{\mathbf{X}^T \mathbf{V}^{-1} \mathbf{X}(\mathbf{X}^T \mathbf{V}^{-1} \mathbf{X})^{-1}\} = p$
    - $\text{trace}(\mathbf{I}_n - \mathbf{H}_V) = \text{trace}(\mathbf{I}_n) - \text{trace}(\mathbf{H}_V) = (n - p)$
    - $\mathbf{H}_V \mathbf{X} = \mathbf{X}$  and  $(\mathbf{I}_n - \mathbf{H}_V) \mathbf{X} = \{\mathbf{I}_n - \mathbf{X}(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T\} \mathbf{X} = \mathbf{0}$
- Fitted values:  $\hat{\mathbf{y}} = \mathbf{X}\hat{\boldsymbol{\beta}} = \mathbf{X}(\mathbf{X}^T \mathbf{V}^{-1} \mathbf{X})^{-1} \mathbf{X}^T \mathbf{V}^{-1} \mathbf{y} = \mathbf{H}_V \mathbf{y}$
- Mean-Variance:  $\mathbb{E}(\hat{\mathbf{y}}) = \mathbf{X}\mathbb{E}(\hat{\boldsymbol{\beta}}) = \mathbf{X}\boldsymbol{\beta}, \quad \text{var}(\hat{\mathbf{y}}) = \sigma^2 \mathbf{H}_V$
- Residuals:  $\hat{\mathbf{e}} = (\mathbf{y} - \hat{\mathbf{y}}) = \mathbf{y} - \mathbf{X}\hat{\boldsymbol{\beta}} = (\mathbf{I}_n - \mathbf{H}_V) \mathbf{y}$
- Mean-Variance:  $\mathbb{E}(\hat{\mathbf{e}}) = (\mathbf{I}_n - \mathbf{H}_V) \mathbf{X} \boldsymbol{\beta} = \mathbf{0}, \quad \text{var}(\hat{\mathbf{e}}) = \mathbb{E}(\hat{\mathbf{e}} \hat{\mathbf{e}}^T) = \sigma^2 (\mathbf{I}_n - \mathbf{H}_V)$
- $\text{cov}(\hat{\mathbf{y}}, \hat{\mathbf{e}}) = \text{cov}\{\mathbf{H}_V \mathbf{y}, (\mathbf{I}_n - \mathbf{H}_V) \mathbf{y}\} = \sigma^2 \mathbf{H}_V (\mathbf{I}_n - \mathbf{H}_V) = \sigma^2 (\mathbf{H}_V - \mathbf{H}_V^2) = \mathbf{0}$



## Generalized / Weighted Least Squares

$$\mathbf{y}^{n \times 1} = \mathbf{X}^{n \times p} \boldsymbol{\beta}^{p \times 1} + \mathbf{e}^{n \times 1}, \quad \mathbb{E}(\mathbf{e}) = \mathbf{0}, \quad \text{cov}(\mathbf{e}) = \sigma^2 \mathbf{V}, \quad \mathbf{V} \text{ known.}$$

- Minimize squared error loss:

$$\hat{\boldsymbol{\beta}} = \arg \min (\mathbf{y} - \mathbf{X}\boldsymbol{\beta})^T \mathbf{V}^{-1} (\mathbf{y} - \mathbf{X}\boldsymbol{\beta}) = (\mathbf{X}^T \mathbf{V}^{-1} \mathbf{X})^{-1} \mathbf{X}^T \mathbf{V}^{-1} \mathbf{y}$$

- $\mathbb{E}(\hat{\boldsymbol{\beta}}) = (\mathbf{X}^T \mathbf{V}^{-1} \mathbf{X})^{-1} \mathbf{X}^T \mathbf{V}^{-1} \mathbf{X} \boldsymbol{\beta} = \boldsymbol{\beta}$
- $\text{var}(\hat{\boldsymbol{\beta}}) = (\mathbf{X}^T \mathbf{V}^{-1} \mathbf{X})^{-1} \mathbf{X}^T \mathbf{V}^{-1} \sigma^2 \mathbf{V} \mathbf{V}^{-1} \mathbf{X} (\mathbf{X}^T \mathbf{V}^{-1} \mathbf{X})^{-1} = \sigma^2 (\mathbf{X}^T \mathbf{V}^{-1} \mathbf{X})^{-1}$ 
  - Hat matrix:  $\mathbf{H}_V = \mathbf{X} (\mathbf{X}^T \mathbf{V}^{-1} \mathbf{X})^{-1} \mathbf{X}^T \mathbf{V}^{-1}$
  - Idempotent:  $\mathbf{H}_V^2 = \mathbf{H}_V, \quad (\mathbf{I}_n - \mathbf{H}_V)^2 = (\mathbf{I}_n - \mathbf{H}_V)$ 
    - $\text{trace}(\mathbf{H}_V) = \text{trace}\{\mathbf{X} (\mathbf{X}^T \mathbf{V}^{-1} \mathbf{X})^{-1} \mathbf{X}^T \mathbf{V}^{-1}\} = \text{trace}\{\mathbf{X}^T \mathbf{V}^{-1} \mathbf{X} (\mathbf{X}^T \mathbf{V}^{-1} \mathbf{X})^{-1}\} = p$
    - $\text{trace}(\mathbf{I}_n - \mathbf{H}_V) = \text{trace}(\mathbf{I}_n) - \text{trace}(\mathbf{H}_V) = (n - p)$
    - $\mathbf{H}_V \mathbf{X} = \mathbf{X}$  and  $(\mathbf{I}_n - \mathbf{H}_V) \mathbf{X} = \{\mathbf{I}_n - \mathbf{X} (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T\} \mathbf{X} = \mathbf{0}$
- Fitted values:  $\hat{\mathbf{y}} = \mathbf{X} \hat{\boldsymbol{\beta}} = \mathbf{X} (\mathbf{X}^T \mathbf{V}^{-1} \mathbf{X})^{-1} \mathbf{X}^T \mathbf{V}^{-1} \mathbf{y} = \mathbf{H}_V \mathbf{y}$
- Mean-Variance:  $\mathbb{E}(\hat{\mathbf{y}}) = \mathbf{X} \mathbb{E}(\hat{\boldsymbol{\beta}}) = \mathbf{X} \boldsymbol{\beta}, \quad \text{var}(\hat{\mathbf{y}}) = \sigma^2 \mathbf{H}_V$
- Residuals:  $\hat{\mathbf{e}} = (\mathbf{y} - \hat{\mathbf{y}}) = \mathbf{y} - \mathbf{X} \hat{\boldsymbol{\beta}} = (\mathbf{I}_n - \mathbf{H}_V) \mathbf{y}$
- Mean-Variance:  $\mathbb{E}(\hat{\mathbf{e}}) = (\mathbf{I}_n - \mathbf{H}_V) \mathbf{X} \boldsymbol{\beta} = \mathbf{0}, \quad \text{var}(\hat{\mathbf{e}}) = \mathbb{E}(\hat{\mathbf{e}} \hat{\mathbf{e}}^T) = \sigma^2 (\mathbf{I}_n - \mathbf{H}_V)$
- $\text{cov}(\hat{\mathbf{y}}, \hat{\mathbf{e}}) = \text{cov}\{\mathbf{H}_V \mathbf{y}, (\mathbf{I}_n - \mathbf{H}_V) \mathbf{y}\} = \sigma^2 \mathbf{H}_V (\mathbf{I}_n - \mathbf{H}_V) = \sigma^2 (\mathbf{H}_V - \mathbf{H}_V^2) = \mathbf{0}$

## Generalized / Weighted Least Squares

$$\mathbf{y}^{n \times 1} = \mathbf{X}^{n \times p} \boldsymbol{\beta}^{p \times 1} + \mathbf{e}^{n \times 1}, \quad \mathbb{E}(\mathbf{e}) = \mathbf{0}, \quad \text{cov}(\mathbf{e}) = \sigma^2 \mathbf{V}, \quad \mathbf{V} \text{ known.}$$

- Minimize squared error loss:

$$\hat{\boldsymbol{\beta}} = \arg \min (\mathbf{y} - \mathbf{X}\boldsymbol{\beta})^T \mathbf{V}^{-1} (\mathbf{y} - \mathbf{X}\boldsymbol{\beta}) = (\mathbf{X}^T \mathbf{V}^{-1} \mathbf{X})^{-1} \mathbf{X}^T \mathbf{V}^{-1} \mathbf{y}$$

- $\mathbb{E}(\hat{\boldsymbol{\beta}}) = (\mathbf{X}^T \mathbf{V}^{-1} \mathbf{X})^{-1} \mathbf{X}^T \mathbf{V}^{-1} \mathbf{X} \boldsymbol{\beta} = \boldsymbol{\beta}$
- $\text{var}(\hat{\boldsymbol{\beta}}) = (\mathbf{X}^T \mathbf{V}^{-1} \mathbf{X})^{-1} \mathbf{X}^T \mathbf{V}^{-1} \sigma^2 \mathbf{V} \mathbf{V}^{-1} \mathbf{X} (\mathbf{X}^T \mathbf{V}^{-1} \mathbf{X})^{-1} = \sigma^2 (\mathbf{X}^T \mathbf{V}^{-1} \mathbf{X})^{-1}$ 
  - Hat matrix:  $\mathbf{H}_V = \mathbf{X} (\mathbf{X}^T \mathbf{V}^{-1} \mathbf{X})^{-1} \mathbf{X}^T \mathbf{V}^{-1}$
  - Idempotent:  $\mathbf{H}_V^2 = \mathbf{H}_V, \quad (\mathbf{I}_n - \mathbf{H}_V)^2 = (\mathbf{I}_n - \mathbf{H}_V)$
  - $\text{trace}(\mathbf{H}_V) = \text{trace}\{\mathbf{X} (\mathbf{X}^T \mathbf{V}^{-1} \mathbf{X})^{-1} \mathbf{X}^T \mathbf{V}^{-1}\} = \text{trace}\{\mathbf{X}^T \mathbf{V}^{-1} \mathbf{X} (\mathbf{X}^T \mathbf{V}^{-1} \mathbf{X})^{-1}\} = p$ 
    - $\text{trace}(\mathbf{I}_n - \mathbf{H}_V) = \text{trace}(\mathbf{I}_n) - \text{trace}(\mathbf{H}_V) = (n - p)$
    - $\mathbf{H}_V \mathbf{X} = \mathbf{X}$  and  $(\mathbf{I}_n - \mathbf{H}_V) \mathbf{X} = \{\mathbf{I}_n - \mathbf{X} (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T\} \mathbf{X} = \mathbf{0}$
- Fitted values:  $\hat{\mathbf{y}} = \mathbf{X} \hat{\boldsymbol{\beta}} = \mathbf{X} (\mathbf{X}^T \mathbf{V}^{-1} \mathbf{X})^{-1} \mathbf{X}^T \mathbf{V}^{-1} \mathbf{y} = \mathbf{H}_V \mathbf{y}$
- Mean-Variance:  $\mathbb{E}(\hat{\mathbf{y}}) = \mathbf{X} \mathbb{E}(\hat{\boldsymbol{\beta}}) = \mathbf{X} \boldsymbol{\beta}, \quad \text{var}(\hat{\mathbf{y}}) = \sigma^2 \mathbf{H}_V$
- Residuals:  $\hat{\mathbf{e}} = (\mathbf{y} - \hat{\mathbf{y}}) = \mathbf{y} - \mathbf{X} \hat{\boldsymbol{\beta}} = (\mathbf{I}_n - \mathbf{H}_V) \mathbf{y}$
- Mean-Variance:  $\mathbb{E}(\hat{\mathbf{e}}) = (\mathbf{I}_n - \mathbf{H}_V) \mathbf{X} \boldsymbol{\beta} = \mathbf{0}, \quad \text{var}(\hat{\mathbf{e}}) = \mathbb{E}(\hat{\mathbf{e}} \hat{\mathbf{e}}^T) = \sigma^2 (\mathbf{I}_n - \mathbf{H}_V)$
- $\text{cov}(\hat{\mathbf{y}}, \hat{\mathbf{e}}) = \text{cov}\{\mathbf{H}_V \mathbf{y}, (\mathbf{I}_n - \mathbf{H}_V) \mathbf{y}\} = \sigma^2 \mathbf{H}_V (\mathbf{I}_n - \mathbf{H}_V) = \sigma^2 (\mathbf{H}_V - \mathbf{H}_V^2) = \mathbf{0}$

## Generalized / Weighted Least Squares

$$\mathbf{y}^{n \times 1} = \mathbf{X}^{n \times p} \boldsymbol{\beta}^{p \times 1} + \mathbf{e}^{n \times 1}, \quad \mathbb{E}(\mathbf{e}) = \mathbf{0}, \quad \text{cov}(\mathbf{e}) = \sigma^2 \mathbf{V}, \quad \mathbf{V} \text{ known.}$$

- Minimize squared error loss:

$$\hat{\boldsymbol{\beta}} = \arg \min (\mathbf{y} - \mathbf{X}\boldsymbol{\beta})^T \mathbf{V}^{-1} (\mathbf{y} - \mathbf{X}\boldsymbol{\beta}) = (\mathbf{X}^T \mathbf{V}^{-1} \mathbf{X})^{-1} \mathbf{X}^T \mathbf{V}^{-1} \mathbf{y}$$

- $\mathbb{E}(\hat{\boldsymbol{\beta}}) = (\mathbf{X}^T \mathbf{V}^{-1} \mathbf{X})^{-1} \mathbf{X}^T \mathbf{V}^{-1} \mathbf{X} \boldsymbol{\beta} = \boldsymbol{\beta}$
- $\text{var}(\hat{\boldsymbol{\beta}}) = (\mathbf{X}^T \mathbf{V}^{-1} \mathbf{X})^{-1} \mathbf{X}^T \mathbf{V}^{-1} \sigma^2 \mathbf{V} \mathbf{V}^{-1} \mathbf{X} (\mathbf{X}^T \mathbf{V}^{-1} \mathbf{X})^{-1} = \sigma^2 (\mathbf{X}^T \mathbf{V}^{-1} \mathbf{X})^{-1}$ 
  - Hat matrix:  $\mathbf{H}_V = \mathbf{X} (\mathbf{X}^T \mathbf{V}^{-1} \mathbf{X})^{-1} \mathbf{X}^T \mathbf{V}^{-1}$
  - Idempotent:  $\mathbf{H}_V^2 = \mathbf{H}_V, \quad (\mathbf{I}_n - \mathbf{H}_V)^2 = (\mathbf{I}_n - \mathbf{H}_V)$
  - $\text{trace}(\mathbf{H}_V) = \text{trace}\{\mathbf{X} (\mathbf{X}^T \mathbf{V}^{-1} \mathbf{X})^{-1} \mathbf{X}^T \mathbf{V}^{-1}\} = \text{trace}\{\mathbf{X}^T \mathbf{V}^{-1} \mathbf{X} (\mathbf{X}^T \mathbf{V}^{-1} \mathbf{X})^{-1}\} = p$
  - $\text{trace}(\mathbf{I}_n - \mathbf{H}_V) = \text{trace}(\mathbf{I}_n) - \text{trace}(\mathbf{H}_V) = (n - p)$
  - $\mathbf{H}_V \mathbf{X} = \mathbf{X}$  and  $(\mathbf{I}_n - \mathbf{H}_V) \mathbf{X} = \{\mathbf{I}_n - \mathbf{X} (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T\} \mathbf{X} = \mathbf{0}$
- Fitted values:  $\hat{\mathbf{y}} = \mathbf{X} \hat{\boldsymbol{\beta}} = \mathbf{X} (\mathbf{X}^T \mathbf{V}^{-1} \mathbf{X})^{-1} \mathbf{X}^T \mathbf{V}^{-1} \mathbf{y} = \mathbf{H}_V \mathbf{y}$
- Mean-Variance:  $\mathbb{E}(\hat{\mathbf{y}}) = \mathbf{X} \mathbb{E}(\hat{\boldsymbol{\beta}}) = \mathbf{X} \boldsymbol{\beta}, \quad \text{var}(\hat{\mathbf{y}}) = \sigma^2 \mathbf{H}_V$
- Residuals:  $\hat{\mathbf{e}} = (\mathbf{y} - \hat{\mathbf{y}}) = \mathbf{y} - \mathbf{X} \hat{\boldsymbol{\beta}} = (\mathbf{I}_n - \mathbf{H}_V) \mathbf{y}$
- Mean-Variance:  $\mathbb{E}(\hat{\mathbf{e}}) = (\mathbf{I}_n - \mathbf{H}_V) \mathbf{X} \boldsymbol{\beta} = \mathbf{0}, \quad \text{var}(\hat{\mathbf{e}}) = \mathbb{E}(\hat{\mathbf{e}} \hat{\mathbf{e}}^T) = \sigma^2 (\mathbf{I}_n - \mathbf{H}_V)$
- $\text{cov}(\hat{\mathbf{y}}, \hat{\mathbf{e}}) = \text{cov}\{\mathbf{H}_V \mathbf{y}, (\mathbf{I}_n - \mathbf{H}_V) \mathbf{y}\} = \sigma^2 \mathbf{H}_V (\mathbf{I}_n - \mathbf{H}_V) = \sigma^2 (\mathbf{H}_V - \mathbf{H}_V^2) = \mathbf{0}$

## Generalized / Weighted Least Squares

$$\mathbf{y}^{n \times 1} = \mathbf{X}^{n \times p} \boldsymbol{\beta}^{p \times 1} + \mathbf{e}^{n \times 1}, \quad \mathbb{E}(\mathbf{e}) = \mathbf{0}, \quad \text{cov}(\mathbf{e}) = \sigma^2 \mathbf{V}, \quad \mathbf{V} \text{ known.}$$

- Minimize squared error loss:

$$\hat{\boldsymbol{\beta}} = \arg \min (\mathbf{y} - \mathbf{X}\boldsymbol{\beta})^T \mathbf{V}^{-1} (\mathbf{y} - \mathbf{X}\boldsymbol{\beta}) = (\mathbf{X}^T \mathbf{V}^{-1} \mathbf{X})^{-1} \mathbf{X}^T \mathbf{V}^{-1} \mathbf{y}$$

- $\mathbb{E}(\hat{\boldsymbol{\beta}}) = (\mathbf{X}^T \mathbf{V}^{-1} \mathbf{X})^{-1} \mathbf{X}^T \mathbf{V}^{-1} \mathbf{X} \boldsymbol{\beta} = \boldsymbol{\beta}$
- $\text{var}(\hat{\boldsymbol{\beta}}) = (\mathbf{X}^T \mathbf{V}^{-1} \mathbf{X})^{-1} \mathbf{X}^T \mathbf{V}^{-1} \sigma^2 \mathbf{V} \mathbf{V}^{-1} \mathbf{X} (\mathbf{X}^T \mathbf{V}^{-1} \mathbf{X})^{-1} = \sigma^2 (\mathbf{X}^T \mathbf{V}^{-1} \mathbf{X})^{-1}$ 
  - Hat matrix:  $\mathbf{H}_V = \mathbf{X} (\mathbf{X}^T \mathbf{V}^{-1} \mathbf{X})^{-1} \mathbf{X}^T \mathbf{V}^{-1}$
  - Idempotent:  $\mathbf{H}_V^2 = \mathbf{H}_V, \quad (\mathbf{I}_n - \mathbf{H}_V)^2 = (\mathbf{I}_n - \mathbf{H}_V)$
  - $\text{trace}(\mathbf{H}_V) = \text{trace}\{\mathbf{X} (\mathbf{X}^T \mathbf{V}^{-1} \mathbf{X})^{-1} \mathbf{X}^T \mathbf{V}^{-1}\} = \text{trace}\{\mathbf{X}^T \mathbf{V}^{-1} \mathbf{X} (\mathbf{X}^T \mathbf{V}^{-1} \mathbf{X})^{-1}\} = p$
  - $\text{trace}(\mathbf{I}_n - \mathbf{H}_V) = \text{trace}(\mathbf{I}_n) - \text{trace}(\mathbf{H}_V) = (n - p)$
  - $\mathbf{H}_V \mathbf{X} = \mathbf{X}$  and  $(\mathbf{I}_n - \mathbf{H}_V) \mathbf{X} = \{\mathbf{I}_n - \mathbf{X} (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T\} \mathbf{X} = \mathbf{0}$
- Fitted values:  $\hat{\mathbf{y}} = \mathbf{X} \hat{\boldsymbol{\beta}} = \mathbf{X} (\mathbf{X}^T \mathbf{V}^{-1} \mathbf{X})^{-1} \mathbf{X}^T \mathbf{V}^{-1} \mathbf{y} = \mathbf{H}_V \mathbf{y}$
- Mean-Variance:  $\mathbb{E}(\hat{\mathbf{y}}) = \mathbf{X} \mathbb{E}(\hat{\boldsymbol{\beta}}) = \mathbf{X} \boldsymbol{\beta}, \quad \text{var}(\hat{\mathbf{y}}) = \sigma^2 \mathbf{H}_V$
- Residuals:  $\hat{\mathbf{e}} = (\mathbf{y} - \hat{\mathbf{y}}) = \mathbf{y} - \mathbf{X} \hat{\boldsymbol{\beta}} = (\mathbf{I}_n - \mathbf{H}_V) \mathbf{y}$
- Mean-Variance:  $\mathbb{E}(\hat{\mathbf{e}}) = (\mathbf{I}_n - \mathbf{H}_V) \mathbf{X} \boldsymbol{\beta} = \mathbf{0}, \quad \text{var}(\hat{\mathbf{e}}) = \mathbb{E}(\hat{\mathbf{e}} \hat{\mathbf{e}}^T) = \sigma^2 (\mathbf{I}_n - \mathbf{H}_V)$
- $\text{cov}(\hat{\mathbf{y}}, \hat{\mathbf{e}}) = \text{cov}\{\mathbf{H}_V \mathbf{y}, (\mathbf{I}_n - \mathbf{H}_V) \mathbf{y}\} = \sigma^2 \mathbf{H}_V (\mathbf{I}_n - \mathbf{H}_V) = \sigma^2 (\mathbf{H}_V - \mathbf{H}_V^2) = \mathbf{0}$

## Generalized / Weighted Least Squares

$$\mathbf{y}^{n \times 1} = \mathbf{X}^{n \times p} \boldsymbol{\beta}^{p \times 1} + \mathbf{e}^{n \times 1}, \quad \mathbb{E}(\mathbf{e}) = \mathbf{0}, \quad \text{cov}(\mathbf{e}) = \sigma^2 \mathbf{V}, \quad \mathbf{V} \text{ known.}$$

- Minimize squared error loss:

$$\hat{\boldsymbol{\beta}} = \arg \min (\mathbf{y} - \mathbf{X}\boldsymbol{\beta})^T \mathbf{V}^{-1} (\mathbf{y} - \mathbf{X}\boldsymbol{\beta}) = (\mathbf{X}^T \mathbf{V}^{-1} \mathbf{X})^{-1} \mathbf{X}^T \mathbf{V}^{-1} \mathbf{y}$$

- $\mathbb{E}(\hat{\boldsymbol{\beta}}) = (\mathbf{X}^T \mathbf{V}^{-1} \mathbf{X})^{-1} \mathbf{X}^T \mathbf{V}^{-1} \mathbf{X} \boldsymbol{\beta} = \boldsymbol{\beta}$
- $\text{var}(\hat{\boldsymbol{\beta}}) = (\mathbf{X}^T \mathbf{V}^{-1} \mathbf{X})^{-1} \mathbf{X}^T \mathbf{V}^{-1} \sigma^2 \mathbf{V} \mathbf{V}^{-1} \mathbf{X} (\mathbf{X}^T \mathbf{V}^{-1} \mathbf{X})^{-1} = \sigma^2 (\mathbf{X}^T \mathbf{V}^{-1} \mathbf{X})^{-1}$ 
  - Hat matrix:  $\mathbf{H}_V = \mathbf{X} (\mathbf{X}^T \mathbf{V}^{-1} \mathbf{X})^{-1} \mathbf{X}^T \mathbf{V}^{-1}$
  - Idempotent:  $\mathbf{H}_V^2 = \mathbf{H}_V, \quad (\mathbf{I}_n - \mathbf{H}_V)^2 = (\mathbf{I}_n - \mathbf{H}_V)$
  - $\text{trace}(\mathbf{H}_V) = \text{trace}\{\mathbf{X} (\mathbf{X}^T \mathbf{V}^{-1} \mathbf{X})^{-1} \mathbf{X}^T \mathbf{V}^{-1}\} = \text{trace}\{\mathbf{X}^T \mathbf{V}^{-1} \mathbf{X} (\mathbf{X}^T \mathbf{V}^{-1} \mathbf{X})^{-1}\} = p$
  - $\text{trace}(\mathbf{I}_n - \mathbf{H}_V) = \text{trace}(\mathbf{I}_n) - \text{trace}(\mathbf{H}_V) = (n - p)$
  - $\mathbf{H}_V \mathbf{X} = \mathbf{X}$  and  $(\mathbf{I}_n - \mathbf{H}_V) \mathbf{X} = \{\mathbf{I}_n - \mathbf{X} (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T\} \mathbf{X} = \mathbf{0}$
- Fitted values:  $\hat{\mathbf{y}} = \mathbf{X} \hat{\boldsymbol{\beta}} = \mathbf{X} (\mathbf{X}^T \mathbf{V}^{-1} \mathbf{X})^{-1} \mathbf{X}^T \mathbf{V}^{-1} \mathbf{y} = \mathbf{H}_V \mathbf{y}$ 
  - Mean-Variance:  $\mathbb{E}(\hat{\mathbf{y}}) = \mathbf{X} \mathbb{E}(\hat{\boldsymbol{\beta}}) = \mathbf{X} \boldsymbol{\beta}, \quad \text{var}(\hat{\mathbf{y}}) = \sigma^2 \mathbf{H}_V$
  - Residuals:  $\hat{\mathbf{e}} = (\mathbf{y} - \hat{\mathbf{y}}) = \mathbf{y} - \mathbf{X} \hat{\boldsymbol{\beta}} = (\mathbf{I}_n - \mathbf{H}_V) \mathbf{y}$
  - Mean-Variance:  $\mathbb{E}(\hat{\mathbf{e}}) = (\mathbf{I}_n - \mathbf{H}_V) \mathbf{X} \boldsymbol{\beta} = \mathbf{0}, \quad \text{var}(\hat{\mathbf{e}}) = \mathbb{E}(\hat{\mathbf{e}} \hat{\mathbf{e}}^T) = \sigma^2 (\mathbf{I}_n - \mathbf{H}_V)$
  - $\text{cov}(\hat{\mathbf{y}}, \hat{\mathbf{e}}) = \text{cov}\{\mathbf{H}_V \mathbf{y}, (\mathbf{I}_n - \mathbf{H}_V) \mathbf{y}\} = \sigma^2 \mathbf{H}_V (\mathbf{I}_n - \mathbf{H}_V) = \sigma^2 (\mathbf{H}_V - \mathbf{H}_V^2) = \mathbf{0}$

## Generalized / Weighted Least Squares

$$\mathbf{y}^{n \times 1} = \mathbf{X}^{n \times p} \boldsymbol{\beta}^{p \times 1} + \mathbf{e}^{n \times 1}, \quad \mathbb{E}(\mathbf{e}) = \mathbf{0}, \quad \text{cov}(\mathbf{e}) = \sigma^2 \mathbf{V}, \quad \mathbf{V} \text{ known.}$$

- Minimize squared error loss:

$$\hat{\boldsymbol{\beta}} = \arg \min (\mathbf{y} - \mathbf{X}\boldsymbol{\beta})^T \mathbf{V}^{-1} (\mathbf{y} - \mathbf{X}\boldsymbol{\beta}) = (\mathbf{X}^T \mathbf{V}^{-1} \mathbf{X})^{-1} \mathbf{X}^T \mathbf{V}^{-1} \mathbf{y}$$

- $\mathbb{E}(\hat{\boldsymbol{\beta}}) = (\mathbf{X}^T \mathbf{V}^{-1} \mathbf{X})^{-1} \mathbf{X}^T \mathbf{V}^{-1} \mathbf{X} \boldsymbol{\beta} = \boldsymbol{\beta}$
- $\text{var}(\hat{\boldsymbol{\beta}}) = (\mathbf{X}^T \mathbf{V}^{-1} \mathbf{X})^{-1} \mathbf{X}^T \mathbf{V}^{-1} \sigma^2 \mathbf{V} \mathbf{V}^{-1} \mathbf{X} (\mathbf{X}^T \mathbf{V}^{-1} \mathbf{X})^{-1} = \sigma^2 (\mathbf{X}^T \mathbf{V}^{-1} \mathbf{X})^{-1}$ 
  - Hat matrix:  $\mathbf{H}_V = \mathbf{X} (\mathbf{X}^T \mathbf{V}^{-1} \mathbf{X})^{-1} \mathbf{X}^T \mathbf{V}^{-1}$
  - Idempotent:  $\mathbf{H}_V^2 = \mathbf{H}_V, \quad (\mathbf{I}_n - \mathbf{H}_V)^2 = (\mathbf{I}_n - \mathbf{H}_V)$
  - $\text{trace}(\mathbf{H}_V) = \text{trace}\{\mathbf{X} (\mathbf{X}^T \mathbf{V}^{-1} \mathbf{X})^{-1} \mathbf{X}^T \mathbf{V}^{-1}\} = \text{trace}\{\mathbf{X}^T \mathbf{V}^{-1} \mathbf{X} (\mathbf{X}^T \mathbf{V}^{-1} \mathbf{X})^{-1}\} = p$
  - $\text{trace}(\mathbf{I}_n - \mathbf{H}_V) = \text{trace}(\mathbf{I}_n) - \text{trace}(\mathbf{H}_V) = (n - p)$
  - $\mathbf{H}_V \mathbf{X} = \mathbf{X}$  and  $(\mathbf{I}_n - \mathbf{H}_V) \mathbf{X} = \{\mathbf{I}_n - \mathbf{X} (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T\} \mathbf{X} = \mathbf{0}$
- Fitted values:  $\hat{\mathbf{y}} = \mathbf{X} \hat{\boldsymbol{\beta}} = \mathbf{X} (\mathbf{X}^T \mathbf{V}^{-1} \mathbf{X})^{-1} \mathbf{X}^T \mathbf{V}^{-1} \mathbf{y} = \mathbf{H}_V \mathbf{y}$
- Mean-Variance:  $\mathbb{E}(\hat{\mathbf{y}}) = \mathbf{X} \mathbb{E}(\hat{\boldsymbol{\beta}}) = \mathbf{X} \boldsymbol{\beta}, \quad \text{var}(\hat{\mathbf{y}}) = \sigma^2 \mathbf{H}_V$
- Residuals:  $\hat{\mathbf{e}} = (\mathbf{y} - \hat{\mathbf{y}}) = \mathbf{y} - \mathbf{X} \hat{\boldsymbol{\beta}} = (\mathbf{I}_n - \mathbf{H}_V) \mathbf{y}$
- Mean-Variance:  $\mathbb{E}(\hat{\mathbf{e}}) = (\mathbf{I}_n - \mathbf{H}_V) \mathbf{X} \boldsymbol{\beta} = \mathbf{0}, \quad \text{var}(\hat{\mathbf{e}}) = \mathbb{E}(\hat{\mathbf{e}} \hat{\mathbf{e}}^T) = \sigma^2 (\mathbf{I}_n - \mathbf{H}_V)$
- $\text{cov}(\hat{\mathbf{y}}, \hat{\mathbf{e}}) = \text{cov}\{\mathbf{H}_V \mathbf{y}, (\mathbf{I}_n - \mathbf{H}_V) \mathbf{y}\} = \sigma^2 \mathbf{H}_V (\mathbf{I}_n - \mathbf{H}_V) = \sigma^2 (\mathbf{H}_V - \mathbf{H}_V^2) = \mathbf{0}$

## Generalized / Weighted Least Squares

$$\mathbf{y}^{n \times 1} = \mathbf{X}^{n \times p} \boldsymbol{\beta}^{p \times 1} + \mathbf{e}^{n \times 1}, \quad \mathbb{E}(\mathbf{e}) = \mathbf{0}, \quad \text{cov}(\mathbf{e}) = \sigma^2 \mathbf{V}, \quad \mathbf{V} \text{ known.}$$

- Minimize squared error loss:

$$\hat{\boldsymbol{\beta}} = \arg \min (\mathbf{y} - \mathbf{X}\boldsymbol{\beta})^T \mathbf{V}^{-1} (\mathbf{y} - \mathbf{X}\boldsymbol{\beta}) = (\mathbf{X}^T \mathbf{V}^{-1} \mathbf{X})^{-1} \mathbf{X}^T \mathbf{V}^{-1} \mathbf{y}$$

- $\mathbb{E}(\hat{\boldsymbol{\beta}}) = (\mathbf{X}^T \mathbf{V}^{-1} \mathbf{X})^{-1} \mathbf{X}^T \mathbf{V}^{-1} \mathbf{X} \boldsymbol{\beta} = \boldsymbol{\beta}$
- $\text{var}(\hat{\boldsymbol{\beta}}) = (\mathbf{X}^T \mathbf{V}^{-1} \mathbf{X})^{-1} \mathbf{X}^T \mathbf{V}^{-1} \sigma^2 \mathbf{V} \mathbf{V}^{-1} \mathbf{X} (\mathbf{X}^T \mathbf{V}^{-1} \mathbf{X})^{-1} = \sigma^2 (\mathbf{X}^T \mathbf{V}^{-1} \mathbf{X})^{-1}$ 
  - Hat matrix:  $\mathbf{H}_V = \mathbf{X} (\mathbf{X}^T \mathbf{V}^{-1} \mathbf{X})^{-1} \mathbf{X}^T \mathbf{V}^{-1}$
  - Idempotent:  $\mathbf{H}_V^2 = \mathbf{H}_V, \quad (\mathbf{I}_n - \mathbf{H}_V)^2 = (\mathbf{I}_n - \mathbf{H}_V)$
  - $\text{trace}(\mathbf{H}_V) = \text{trace}\{\mathbf{X} (\mathbf{X}^T \mathbf{V}^{-1} \mathbf{X})^{-1} \mathbf{X}^T \mathbf{V}^{-1}\} = \text{trace}\{\mathbf{X}^T \mathbf{V}^{-1} \mathbf{X} (\mathbf{X}^T \mathbf{V}^{-1} \mathbf{X})^{-1}\} = p$
  - $\text{trace}(\mathbf{I}_n - \mathbf{H}_V) = \text{trace}(\mathbf{I}_n) - \text{trace}(\mathbf{H}_V) = (n - p)$
  - $\mathbf{H}_V \mathbf{X} = \mathbf{X}$  and  $(\mathbf{I}_n - \mathbf{H}_V) \mathbf{X} = \{\mathbf{I}_n - \mathbf{X} (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T\} \mathbf{X} = \mathbf{0}$
- Fitted values:  $\hat{\mathbf{y}} = \mathbf{X} \hat{\boldsymbol{\beta}} = \mathbf{X} (\mathbf{X}^T \mathbf{V}^{-1} \mathbf{X})^{-1} \mathbf{X}^T \mathbf{V}^{-1} \mathbf{y} = \mathbf{H}_V \mathbf{y}$
- Mean-Variance:  $\mathbb{E}(\hat{\mathbf{y}}) = \mathbf{X} \mathbb{E}(\hat{\boldsymbol{\beta}}) = \mathbf{X} \boldsymbol{\beta}, \quad \text{var}(\hat{\mathbf{y}}) = \sigma^2 \mathbf{H}_V$
- Residuals:  $\hat{\mathbf{e}} = (\mathbf{y} - \hat{\mathbf{y}}) = \mathbf{y} - \mathbf{X} \hat{\boldsymbol{\beta}} = (\mathbf{I}_n - \mathbf{H}_V) \mathbf{y}$ 
  - Mean-Variance:  $\mathbb{E}(\hat{\mathbf{e}}) = (\mathbf{I}_n - \mathbf{H}_V) \mathbf{X} \boldsymbol{\beta} = \mathbf{0}, \quad \text{var}(\hat{\mathbf{e}}) = \mathbb{E}(\hat{\mathbf{e}} \hat{\mathbf{e}}^T) = \sigma^2 (\mathbf{I}_n - \mathbf{H}_V)$
  - $\text{cov}(\hat{\mathbf{y}}, \hat{\mathbf{e}}) = \text{cov}\{\mathbf{H}_V \mathbf{y}, (\mathbf{I}_n - \mathbf{H}_V) \mathbf{y}\} = \sigma^2 \mathbf{H}_V (\mathbf{I}_n - \mathbf{H}_V) = \sigma^2 (\mathbf{H}_V - \mathbf{H}_V^2) = \mathbf{0}$

## Generalized / Weighted Least Squares

$$\mathbf{y}^{n \times 1} = \mathbf{X}^{n \times p} \boldsymbol{\beta}^{p \times 1} + \mathbf{e}^{n \times 1}, \quad \mathbb{E}(\mathbf{e}) = \mathbf{0}, \quad \text{cov}(\mathbf{e}) = \sigma^2 \mathbf{V}, \quad \mathbf{V} \text{ known.}$$

- Minimize squared error loss:

$$\hat{\boldsymbol{\beta}} = \arg \min (\mathbf{y} - \mathbf{X}\boldsymbol{\beta})^T \mathbf{V}^{-1} (\mathbf{y} - \mathbf{X}\boldsymbol{\beta}) = (\mathbf{X}^T \mathbf{V}^{-1} \mathbf{X})^{-1} \mathbf{X}^T \mathbf{V}^{-1} \mathbf{y}$$

- $\mathbb{E}(\hat{\boldsymbol{\beta}}) = (\mathbf{X}^T \mathbf{V}^{-1} \mathbf{X})^{-1} \mathbf{X}^T \mathbf{V}^{-1} \mathbf{X} \boldsymbol{\beta} = \boldsymbol{\beta}$
- $\text{var}(\hat{\boldsymbol{\beta}}) = (\mathbf{X}^T \mathbf{V}^{-1} \mathbf{X})^{-1} \mathbf{X}^T \mathbf{V}^{-1} \sigma^2 \mathbf{V} \mathbf{V}^{-1} \mathbf{X} (\mathbf{X}^T \mathbf{V}^{-1} \mathbf{X})^{-1} = \sigma^2 (\mathbf{X}^T \mathbf{V}^{-1} \mathbf{X})^{-1}$ 
  - Hat matrix:  $\mathbf{H}_V = \mathbf{X} (\mathbf{X}^T \mathbf{V}^{-1} \mathbf{X})^{-1} \mathbf{X}^T \mathbf{V}^{-1}$
  - Idempotent:  $\mathbf{H}_V^2 = \mathbf{H}_V, \quad (\mathbf{I}_n - \mathbf{H}_V)^2 = (\mathbf{I}_n - \mathbf{H}_V)$
  - $\text{trace}(\mathbf{H}_V) = \text{trace}\{\mathbf{X} (\mathbf{X}^T \mathbf{V}^{-1} \mathbf{X})^{-1} \mathbf{X}^T \mathbf{V}^{-1}\} = \text{trace}\{\mathbf{X}^T \mathbf{V}^{-1} \mathbf{X} (\mathbf{X}^T \mathbf{V}^{-1} \mathbf{X})^{-1}\} = p$
  - $\text{trace}(\mathbf{I}_n - \mathbf{H}_V) = \text{trace}(\mathbf{I}_n) - \text{trace}(\mathbf{H}_V) = (n - p)$
  - $\mathbf{H}_V \mathbf{X} = \mathbf{X}$  and  $(\mathbf{I}_n - \mathbf{H}_V) \mathbf{X} = \{\mathbf{I}_n - \mathbf{X} (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T\} \mathbf{X} = \mathbf{0}$
- Fitted values:  $\hat{\mathbf{y}} = \mathbf{X} \hat{\boldsymbol{\beta}} = \mathbf{X} (\mathbf{X}^T \mathbf{V}^{-1} \mathbf{X})^{-1} \mathbf{X}^T \mathbf{V}^{-1} \mathbf{y} = \mathbf{H}_V \mathbf{y}$
- Mean-Variance:  $\mathbb{E}(\hat{\mathbf{y}}) = \mathbf{X} \mathbb{E}(\hat{\boldsymbol{\beta}}) = \mathbf{X} \boldsymbol{\beta}, \quad \text{var}(\hat{\mathbf{y}}) = \sigma^2 \mathbf{H}_V$
- Residuals:  $\hat{\mathbf{e}} = (\mathbf{y} - \hat{\mathbf{y}}) = \mathbf{y} - \mathbf{X} \hat{\boldsymbol{\beta}} = (\mathbf{I}_n - \mathbf{H}_V) \mathbf{y}$
- Mean-Variance:  $\mathbb{E}(\hat{\mathbf{e}}) = (\mathbf{I}_n - \mathbf{H}_V) \mathbf{X} \boldsymbol{\beta} = \mathbf{0}, \quad \text{var}(\hat{\mathbf{e}}) = \mathbb{E}(\hat{\mathbf{e}} \hat{\mathbf{e}}^T) = \sigma^2 (\mathbf{I}_n - \mathbf{H}_V)$
- $\text{cov}(\hat{\mathbf{y}}, \hat{\mathbf{e}}) = \text{cov}\{\mathbf{H}_V \mathbf{y}, (\mathbf{I}_n - \mathbf{H}_V) \mathbf{y}\} = \sigma^2 \mathbf{H}_V (\mathbf{I}_n - \mathbf{H}_V) = \sigma^2 (\mathbf{H}_V - \mathbf{H}_V^2) = \mathbf{0}$



$$\mathbf{y}^{n \times 1} = \mathbf{X}^{n \times p} \boldsymbol{\beta}^{p \times 1} + \mathbf{e}^{n \times 1}, \quad \mathbb{E}(\mathbf{e}) = \mathbf{0}, \quad \text{cov}(\mathbf{e}) = \sigma^2 \mathbf{V}, \quad \mathbf{V} \text{ known.}$$

- Minimize squared error loss:

$$\hat{\boldsymbol{\beta}} = \arg \min (\mathbf{y} - \mathbf{X}\boldsymbol{\beta})^T \mathbf{V}^{-1} (\mathbf{y} - \mathbf{X}\boldsymbol{\beta}) = (\mathbf{X}^T \mathbf{V}^{-1} \mathbf{X})^{-1} \mathbf{X}^T \mathbf{V}^{-1} \mathbf{y}$$

- $\mathbb{E}(\hat{\boldsymbol{\beta}}) = (\mathbf{X}^T \mathbf{V}^{-1} \mathbf{X})^{-1} \mathbf{X}^T \mathbf{V}^{-1} \mathbf{X} \boldsymbol{\beta} = \boldsymbol{\beta}$
- $\text{var}(\hat{\boldsymbol{\beta}}) = (\mathbf{X}^T \mathbf{V}^{-1} \mathbf{X})^{-1} \mathbf{X}^T \mathbf{V}^{-1} \sigma^2 \mathbf{V} \mathbf{V}^{-1} \mathbf{X} (\mathbf{X}^T \mathbf{V}^{-1} \mathbf{X})^{-1} = \sigma^2 (\mathbf{X}^T \mathbf{V}^{-1} \mathbf{X})^{-1}$ 
  - Hat matrix:  $\mathbf{H}_V = \mathbf{X} (\mathbf{X}^T \mathbf{V}^{-1} \mathbf{X})^{-1} \mathbf{X}^T \mathbf{V}^{-1}$
  - Idempotent:  $\mathbf{H}_V^2 = \mathbf{H}_V, \quad (\mathbf{I}_n - \mathbf{H}_V)^2 = (\mathbf{I}_n - \mathbf{H}_V)$
  - $\text{trace}(\mathbf{H}_V) = \text{trace}\{\mathbf{X} (\mathbf{X}^T \mathbf{V}^{-1} \mathbf{X})^{-1} \mathbf{X}^T \mathbf{V}^{-1}\} = \text{trace}\{\mathbf{X}^T \mathbf{V}^{-1} \mathbf{X} (\mathbf{X}^T \mathbf{V}^{-1} \mathbf{X})^{-1}\} = p$
  - $\text{trace}(\mathbf{I}_n - \mathbf{H}_V) = \text{trace}(\mathbf{I}_n) - \text{trace}(\mathbf{H}_V) = (n - p)$
  - $\mathbf{H}_V \mathbf{X} = \mathbf{X}$  and  $(\mathbf{I}_n - \mathbf{H}_V) \mathbf{X} = \{\mathbf{I}_n - \mathbf{X} (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T\} \mathbf{X} = \mathbf{0}$
- Fitted values:  $\hat{\mathbf{y}} = \mathbf{X} \hat{\boldsymbol{\beta}} = \mathbf{X} (\mathbf{X}^T \mathbf{V}^{-1} \mathbf{X})^{-1} \mathbf{X}^T \mathbf{V}^{-1} \mathbf{y} = \mathbf{H}_V \mathbf{y}$
- Mean-Variance:  $\mathbb{E}(\hat{\mathbf{y}}) = \mathbf{X} \mathbb{E}(\hat{\boldsymbol{\beta}}) = \mathbf{X} \boldsymbol{\beta}, \quad \text{var}(\hat{\mathbf{y}}) = \sigma^2 \mathbf{H}_V$
- Residuals:  $\hat{\mathbf{e}} = (\mathbf{y} - \hat{\mathbf{y}}) = \mathbf{y} - \mathbf{X} \hat{\boldsymbol{\beta}} = (\mathbf{I}_n - \mathbf{H}_V) \mathbf{y}$
- Mean-Variance:  $\mathbb{E}(\hat{\mathbf{e}}) = (\mathbf{I}_n - \mathbf{H}_V) \mathbf{X} \boldsymbol{\beta} = \mathbf{0}, \quad \text{var}(\hat{\mathbf{e}}) = \mathbb{E}(\hat{\mathbf{e}} \hat{\mathbf{e}}^T) = \sigma^2 (\mathbf{I}_n - \mathbf{H}_V)$
- $\text{cov}(\hat{\mathbf{y}}, \hat{\mathbf{e}}) = \text{cov}\{\mathbf{H}_V \mathbf{y}, (\mathbf{I}_n - \mathbf{H}_V) \mathbf{y}\} = \sigma^2 \mathbf{H}_V (\mathbf{I}_n - \mathbf{H}_V) = \sigma^2 (\mathbf{H}_V - \mathbf{H}_V^2) = \mathbf{0}$

$$\mathbf{y}^{n \times 1} = \mathbf{X}^{n \times p} \boldsymbol{\beta}^{p \times 1} + \mathbf{e}^{n \times 1}, \quad \mathbb{E}(\mathbf{e}) = \mathbf{0}, \quad \text{cov}(\mathbf{e}) = \sigma^2 \mathbf{V}, \quad \mathbf{V} \text{ known.}$$

- Likelihood:  $L(\boldsymbol{\beta}, \sigma^2) = \frac{1}{(2\pi\sigma^2)^{n/2}} \exp \left\{ -\frac{1}{2\sigma^2} (\mathbf{y} - \mathbf{X}\boldsymbol{\beta})^T \mathbf{V}^{-1} (\mathbf{y} - \mathbf{X}\boldsymbol{\beta}) \right\}$   

$$\propto \frac{1}{(\sigma^2)^{n/2}} \exp \left\{ -\frac{1}{2\sigma^2} (\boldsymbol{\beta}^T \mathbf{X}^T \mathbf{V}^{-1} \mathbf{X} \boldsymbol{\beta} - 2\boldsymbol{\beta}^T \mathbf{X}^T \mathbf{V}^{-1} \mathbf{y} + \mathbf{y}^T \mathbf{y}) \right\}$$
- Semi-conjugate priors:  $p(\boldsymbol{\beta}, \sigma^2) = \text{MVN}(\boldsymbol{\mu}_\beta, \boldsymbol{\Sigma}_\beta) \cdot \text{Inv-Ga}(a_\sigma, b_\sigma)$   

$$\propto \frac{1}{(\sigma^2)^{a_\sigma+1}} \exp \left( -\frac{b_\sigma}{\sigma^2} \right) \exp \left\{ -\frac{1}{2} (\boldsymbol{\beta}^T \boldsymbol{\Sigma}_\beta^{-1} \boldsymbol{\beta} - 2\boldsymbol{\beta}^T \boldsymbol{\Sigma}_\beta^{-1} \boldsymbol{\mu}_\beta + \boldsymbol{\mu}_\beta^T \boldsymbol{\Sigma}_\beta^{-1} \boldsymbol{\mu}_\beta) \right\}$$
- Posterior:  $p(\boldsymbol{\beta}, \sigma^2 \mid \mathbf{y}_{1:n}) \propto L(\boldsymbol{\beta}, \sigma^2) \times p(\boldsymbol{\beta}, \sigma^2)$
- Posterior full conditionals for block-Gibbs sampler:

$$p(\boldsymbol{\beta} \mid \sigma^2, \mathbf{y}) \propto \exp \left\{ -\frac{1}{2\sigma^2} (\boldsymbol{\beta}^T \mathbf{X}^T \mathbf{V}^{-1} \mathbf{X} \boldsymbol{\beta} - 2\boldsymbol{\beta}^T \mathbf{X}^T \mathbf{V}^{-1} \mathbf{y} + \mathbf{y}^T \mathbf{y}) - \frac{1}{2} (\boldsymbol{\beta}^T \boldsymbol{\Sigma}_\beta^{-1} \boldsymbol{\beta} - 2\boldsymbol{\beta}^T \boldsymbol{\Sigma}_\beta^{-1} \boldsymbol{\mu}_\beta + \boldsymbol{\mu}_\beta^T \boldsymbol{\Sigma}_\beta^{-1} \boldsymbol{\mu}_\beta) \right\}$$

$$p(\sigma^2 \mid \boldsymbol{\beta}, \mathbf{y}) \propto \frac{1}{(\sigma^2)^{a_\sigma+1}} \exp \left( -\frac{b_\sigma}{\sigma^2} \right) \exp \left\{ -\frac{1}{2\sigma^2} (\mathbf{y} - \mathbf{X}\boldsymbol{\beta})^T \mathbf{V}^{-1} (\mathbf{y} - \mathbf{X}\boldsymbol{\beta}) \right\}$$

$$p(\boldsymbol{\beta}_j \mid \sigma^2, \boldsymbol{\beta}_{-j}, \mathbf{y}) \propto \exp \left\{ -\frac{1}{2\sigma^2} (\boldsymbol{\beta}_j^T \mathbf{X}_j^T \mathbf{V}^{-1} \mathbf{X}_j \boldsymbol{\beta}_j - 2\boldsymbol{\beta}_j^T \mathbf{X}_j^T \mathbf{V}^{-1} \mathbf{y} + \mathbf{y}^T \mathbf{y}) - \frac{1}{2} (\boldsymbol{\beta}_j^T \boldsymbol{\Sigma}_{\beta_j}^{-1} \boldsymbol{\beta}_j - 2\boldsymbol{\beta}_j^T \boldsymbol{\Sigma}_{\beta_j}^{-1} \boldsymbol{\mu}_{\beta_j} + \boldsymbol{\mu}_{\beta_j}^T \boldsymbol{\Sigma}_{\beta_j}^{-1} \boldsymbol{\mu}_{\beta_j}) \right\}$$

$$p(\sigma^2 \mid \boldsymbol{\beta}_j, \boldsymbol{\beta}_{-j}, \mathbf{y}) \propto \frac{1}{(\sigma^2)^{a_\sigma+1}} \exp \left( -\frac{b_\sigma}{\sigma^2} \right) \exp \left\{ -\frac{1}{2\sigma^2} (\mathbf{y} - \mathbf{X}\boldsymbol{\beta}_j - \mathbf{X}_{-j}\boldsymbol{\beta}_{-j})^T \mathbf{V}^{-1} (\mathbf{y} - \mathbf{X}\boldsymbol{\beta}_j - \mathbf{X}_{-j}\boldsymbol{\beta}_{-j}) \right\}$$

$$p(\boldsymbol{\beta}_j \mid \sigma^2, \boldsymbol{\beta}_{-j}, \mathbf{y}) \propto \exp \left\{ -\frac{1}{2\sigma^2} (\boldsymbol{\beta}_j^T \mathbf{X}_j^T \mathbf{V}^{-1} \mathbf{X}_j \boldsymbol{\beta}_j - 2\boldsymbol{\beta}_j^T \mathbf{X}_j^T \mathbf{V}^{-1} \mathbf{y} + \mathbf{y}^T \mathbf{y}) - \frac{1}{2} (\boldsymbol{\beta}_j^T \boldsymbol{\Sigma}_{\beta_j}^{-1} \boldsymbol{\beta}_j - 2\boldsymbol{\beta}_j^T \boldsymbol{\Sigma}_{\beta_j}^{-1} \boldsymbol{\mu}_{\beta_j} + \boldsymbol{\mu}_{\beta_j}^T \boldsymbol{\Sigma}_{\beta_j}^{-1} \boldsymbol{\mu}_{\beta_j}) \right\}$$

- $\mathbf{y}^{n \times 1} = \mathbf{X}^{n \times p} \boldsymbol{\beta}^{p \times 1} + \mathbf{e}^{n \times 1}, \quad \mathbb{E}(\mathbf{e}) = \mathbf{0}, \quad \text{cov}(\mathbf{e}) = \sigma^2 \mathbf{V}, \quad \mathbf{V} \text{ known},$
- Likelihood:  $L(\boldsymbol{\beta}, \sigma^2) = \frac{1}{(2\pi\sigma^2)^{n/2}} \exp \left\{ -\frac{1}{2\sigma^2} (\mathbf{y} - \mathbf{X}\boldsymbol{\beta})^T \mathbf{V}^{-1} (\mathbf{y} - \mathbf{X}\boldsymbol{\beta}) \right\}$   
 $\propto \frac{1}{(\sigma^2)^{n/2}} \exp \left\{ -\frac{1}{2\sigma^2} (\boldsymbol{\beta}^T \mathbf{X}^T \mathbf{V}^{-1} \mathbf{X} \boldsymbol{\beta} - 2\boldsymbol{\beta}^T \mathbf{X}^T \mathbf{V}^{-1} \mathbf{y} + \mathbf{y}^T \mathbf{y}) \right\}$
- Semi-conjugate priors:  $p(\boldsymbol{\beta}, \sigma^2) = \text{MVN}(\boldsymbol{\mu}_\beta, \boldsymbol{\Sigma}_\beta) \cdot \text{Inv-Ga}(a_\sigma, b_\sigma)$   
 $\propto \frac{1}{(\sigma^2)^{a_\sigma+1}} \exp \left( -\frac{b_\sigma}{\sigma^2} \right) \exp \left\{ -\frac{1}{2} (\boldsymbol{\beta}^T \boldsymbol{\Sigma}_\beta^{-1} \boldsymbol{\beta} - 2\boldsymbol{\beta}^T \boldsymbol{\Sigma}_\beta^{-1} \boldsymbol{\mu}_\beta + \boldsymbol{\mu}_\beta^T \boldsymbol{\Sigma}_\beta^{-1} \boldsymbol{\mu}_\beta) \right\}$
- Posterior:  $p(\boldsymbol{\beta}, \sigma^2 \mid \mathbf{y}_{1:n}) \propto L(\boldsymbol{\beta}, \sigma^2) \times p(\boldsymbol{\beta}, \sigma^2)$
- Posterior full conditionals for block-Gibbs sampler:

- $\mathbf{y}^{n \times 1} = \mathbf{X}^{n \times p} \boldsymbol{\beta}^{p \times 1} + \mathbf{e}^{n \times 1}, \quad \mathbb{E}(\mathbf{e}) = \mathbf{0}, \quad \text{cov}(\mathbf{e}) = \sigma^2 \mathbf{V}, \quad \mathbf{V} \text{ known.}$
- Likelihood:  $L(\boldsymbol{\beta}, \sigma^2) = \frac{1}{(2\pi\sigma^2)^{n/2}} \exp \left\{ -\frac{1}{2\sigma^2} (\mathbf{y} - \mathbf{X}\boldsymbol{\beta})^T \mathbf{V}^{-1} (\mathbf{y} - \mathbf{X}\boldsymbol{\beta}) \right\}$   
 $\propto \frac{1}{(\sigma^2)^{n/2}} \exp \left\{ -\frac{1}{2\sigma^2} (\boldsymbol{\beta}^T \mathbf{X}^T \mathbf{V}^{-1} \mathbf{X} \boldsymbol{\beta} - 2\boldsymbol{\beta}^T \mathbf{X}^T \mathbf{V}^{-1} \mathbf{y} + \mathbf{y}^T \mathbf{y}) \right\}$
- Semi-conjugate priors:  $p(\boldsymbol{\beta}, \sigma^2) = \text{MVN}(\boldsymbol{\mu}_\beta, \boldsymbol{\Sigma}_\beta) \cdot \text{Inv-Ga}(a_\sigma, b_\sigma)$   
 $\propto \frac{1}{(\sigma^2)^{a_\sigma+1}} \exp \left( -\frac{b_\sigma}{\sigma^2} \right) \exp \left\{ -\frac{1}{2} (\boldsymbol{\beta}^T \boldsymbol{\Sigma}_\beta^{-1} \boldsymbol{\beta} - 2\boldsymbol{\beta}^T \boldsymbol{\Sigma}_\beta^{-1} \boldsymbol{\mu}_\beta + \boldsymbol{\mu}_\beta^T \boldsymbol{\Sigma}_\beta^{-1} \boldsymbol{\mu}_\beta) \right\}$
- Posterior:  $p(\boldsymbol{\beta}, \sigma^2 \mid \mathbf{y}_{1:n}) \propto L(\boldsymbol{\beta}, \sigma^2) \times p(\boldsymbol{\beta}, \sigma^2)$
- Posterior full conditionals for block-Gibbs sampler:

- $\mathbf{y}^{n \times 1} = \mathbf{X}^{n \times p} \boldsymbol{\beta}^{p \times 1} + \mathbf{e}^{n \times 1}, \quad \mathbb{E}(\mathbf{e}) = \mathbf{0}, \quad \text{cov}(\mathbf{e}) = \sigma^2 \mathbf{V}, \quad \mathbf{V} \text{ known},$
- Likelihood:  $L(\boldsymbol{\beta}, \sigma^2) = \frac{1}{(2\pi\sigma^2)^{n/2}} \exp \left\{ -\frac{1}{2\sigma^2} (\mathbf{y} - \mathbf{X}\boldsymbol{\beta})^T \mathbf{V}^{-1} (\mathbf{y} - \mathbf{X}\boldsymbol{\beta}) \right\}$   
 $\propto \frac{1}{(\sigma^2)^{n/2}} \exp \left\{ -\frac{1}{2\sigma^2} (\boldsymbol{\beta}^T \mathbf{X}^T \mathbf{V}^{-1} \mathbf{X} \boldsymbol{\beta} - 2\boldsymbol{\beta}^T \mathbf{X}^T \mathbf{V}^{-1} \mathbf{y} + \mathbf{y}^T \mathbf{y}) \right\}$
- Semi-conjugate priors:  $p(\boldsymbol{\beta}, \sigma^2) = \text{MVN}(\boldsymbol{\mu}_\beta, \boldsymbol{\Sigma}_\beta) \cdot \text{Inv-Ga}(a_\sigma, b_\sigma)$   
 $\propto \frac{1}{(\sigma^2)^{a_\sigma+1}} \exp \left( -\frac{b_\sigma}{\sigma^2} \right) \exp \left\{ -\frac{1}{2} (\boldsymbol{\beta}^T \boldsymbol{\Sigma}_\beta^{-1} \boldsymbol{\beta} - 2\boldsymbol{\beta}^T \boldsymbol{\Sigma}_\beta^{-1} \boldsymbol{\mu}_\beta + \boldsymbol{\mu}_\beta^T \boldsymbol{\Sigma}_\beta^{-1} \boldsymbol{\mu}_\beta) \right\}$
- Posterior:  $p(\boldsymbol{\beta}, \sigma^2 \mid \mathbf{y}_{1:n}) \propto L(\boldsymbol{\beta}, \sigma^2) \times p(\boldsymbol{\beta}, \sigma^2)$
- Posterior full conditionals for block-Gibbs sampler:
  - $p(\boldsymbol{\beta} \mid -) \propto \exp \left[ -\frac{1}{2} \left\{ \boldsymbol{\beta}^T (\boldsymbol{\Sigma}_\beta^{-1} + \sigma^{-2} \mathbf{X}^T \mathbf{V}^{-1} \mathbf{X}) \boldsymbol{\beta} - 2\boldsymbol{\beta}^T (\boldsymbol{\Sigma}_\beta^{-1} \boldsymbol{\mu}_\beta + \sigma^{-2} \mathbf{X}^T \mathbf{V}^{-1} \mathbf{y}) \right\} \right]$   
 $\equiv \text{MVN}(\boldsymbol{\mu}_{\beta,n}, \boldsymbol{\Sigma}_{\beta,n}),$   
 $\boldsymbol{\Sigma}_{\beta,n} = (\boldsymbol{\Sigma}_\beta^{-1} + \sigma^{-2} \mathbf{X}^T \mathbf{V}^{-1} \mathbf{X})^{-1}, \quad \boldsymbol{\mu}_{\beta,n} = \boldsymbol{\Sigma}_{\beta,n} (\boldsymbol{\Sigma}_\beta^{-1} \boldsymbol{\mu}_\beta + \sigma^{-2} \mathbf{X}^T \mathbf{V}^{-1} \mathbf{y})$
  - $p(\sigma^2 \mid -) \propto \frac{1}{(\sigma^2)^{a_\sigma + \frac{n}{2} + 1}} \exp \left[ -\frac{1}{\sigma^2} \left\{ b_\sigma + \frac{1}{2} (\mathbf{y} - \mathbf{X}\boldsymbol{\beta})^T \mathbf{V}^{-1} (\mathbf{y} - \mathbf{X}\boldsymbol{\beta}) \right\} \right]$   
 $\equiv \text{Inv-Ga} \left\{ a_\sigma + \frac{n}{2}, b_\sigma + \frac{1}{2} (\mathbf{y} - \mathbf{X}\boldsymbol{\beta})^T \mathbf{V}^{-1} (\mathbf{y} - \mathbf{X}\boldsymbol{\beta}) \right\}$

- $\mathbf{y}^{n \times 1} = \mathbf{X}^{n \times p} \boldsymbol{\beta}^{p \times 1} + \mathbf{e}^{n \times 1}, \quad \mathbb{E}(\mathbf{e}) = \mathbf{0}, \quad \text{cov}(\mathbf{e}) = \sigma^2 \mathbf{V}, \quad \mathbf{V} \text{ known},$   
Likelihood:  $L(\boldsymbol{\beta}, \sigma^2) = \frac{1}{(2\pi\sigma^2)^{n/2}} \exp \left\{ -\frac{1}{2\sigma^2} (\mathbf{y} - \mathbf{X}\boldsymbol{\beta})^T \mathbf{V}^{-1} (\mathbf{y} - \mathbf{X}\boldsymbol{\beta}) \right\}$   
 $\propto \frac{1}{(\sigma^2)^{n/2}} \exp \left\{ -\frac{1}{2\sigma^2} (\boldsymbol{\beta}^T \mathbf{X}^T \mathbf{V}^{-1} \mathbf{X} \boldsymbol{\beta} - 2\boldsymbol{\beta}^T \mathbf{X}^T \mathbf{V}^{-1} \mathbf{y} + \mathbf{y}^T \mathbf{y}) \right\}$
- Semi-conjugate priors:  $p(\boldsymbol{\beta}, \sigma^2) = \text{MVN}(\boldsymbol{\mu}_\beta, \boldsymbol{\Sigma}_\beta) \cdot \text{Inv-Ga}(a_\sigma, b_\sigma)$
- Posterior:  $p(\boldsymbol{\beta}, \sigma^2 \mid \mathbf{y}_{1:n}) \propto L(\boldsymbol{\beta}, \sigma^2) \times p(\boldsymbol{\beta}, \sigma^2)$
- Posterior full conditionals for block-Gibbs sampler:
  - $p(\boldsymbol{\beta} \mid -) = \text{MVN}(\boldsymbol{\mu}_{\beta,n}, \boldsymbol{\Sigma}_{\beta,n}),$   
 $\boldsymbol{\Sigma}_{\beta,n} = (\boldsymbol{\Sigma}_\beta^{-1} + \sigma^{-2} \mathbf{X}^T \mathbf{V}^{-1} \mathbf{X})^{-1}, \quad \boldsymbol{\mu}_{\beta,n} = \boldsymbol{\Sigma}_{\beta,n} \left( \boldsymbol{\Sigma}_\beta^{-1} \boldsymbol{\mu}_\beta + \sigma^{-2} \mathbf{X}^T \mathbf{V}^{-1} \mathbf{y} \right)$
  - $p(\sigma^2 \mid -) = \text{Inv-Ga} \left\{ a_\sigma + \frac{n}{2}, b_\sigma + \frac{1}{2} (\mathbf{y} - \mathbf{X}\boldsymbol{\beta})^T \mathbf{V}^{-1} (\mathbf{y} - \mathbf{X}\boldsymbol{\beta}) \right\}$
- Bayesian weighted Ridge and LASSO linear models can be similarly developed.

$$y_{i,j} = \mathbf{x}_i^T \boldsymbol{\beta} + \mathbf{z}_i^T \mathbf{u}_i + \epsilon_{i,j}, \quad \mathbf{u}_i \stackrel{iid}{\sim} f_{\mathbf{u}}, \quad \epsilon_{i,j} \stackrel{iid}{\sim} f_{\epsilon},$$

$$i = 1, \dots, n, \quad j = 1, \dots, m_i$$

with  $\mathbb{E}_{f_{\mathbf{u}}}(\mathbf{u}) = \mathbf{0}$  and  $\mathbb{E}_{f_{\epsilon}}(\epsilon) = 0$ .

- ▶  $y_{i,j} = \beta_0 + x_i \beta_1 + u_i + \epsilon_{i,j}, \quad u_i \stackrel{iid}{\sim} \text{Normal}(0, \sigma_u^2), \quad \epsilon_{i,j} \stackrel{iid}{\sim} \text{Normal}(0, \sigma_{\epsilon}^2)$
- ▶  $y_{i,j} = \beta_0 + x_i \beta_1 + x_i^2 \beta_2 + u_i + \epsilon_{i,j}, \quad u_i \stackrel{iid}{\sim} \text{Normal}(0, \sigma_u^2), \quad \epsilon_{i,j} \stackrel{iid}{\sim} \text{Normal}(0, \sigma_{\epsilon}^2)$
- ▶  $y_{i,j} = \beta_0 + x_{i,1} \beta_1 + x_{i,2} \beta_2 + \dots + x_{i,p} \beta_p + u_i + \epsilon_{i,j},$   
 $u_i \stackrel{iid}{\sim} \text{Normal}(0, \sigma_u^2), \quad \epsilon_{i,j} \stackrel{iid}{\sim} \text{Normal}(0, \sigma_{\epsilon}^2)$
- ▶  $y_{i,j} = \beta_0 + x_i \beta_1 + u_{0,i} + x_i u_{1,i} + \epsilon_{i,j}, \quad \mathbf{u}_i \stackrel{iid}{\sim} \text{MVN}(\mathbf{0}, \boldsymbol{\Sigma}_u), \quad \epsilon_{i,j} \stackrel{iid}{\sim} \text{Normal}(0, \sigma_{\epsilon}^2)$
- ▶  $y_{i,j} = \beta_0 + x_i \beta_1 + x_i^2 \beta_2 + u_{0,i} + x_i u_{1,i} + \epsilon_{i,j},$   
 $\mathbf{u}_i \stackrel{iid}{\sim} \text{MVN}(\mathbf{0}, \boldsymbol{\Sigma}_u), \quad \epsilon_{i,j} \stackrel{iid}{\sim} \text{Normal}(0, \sigma_{\epsilon}^2)$
- ▶  $y_{i,j} = \beta_0 + x_{i,1} \beta_1 + x_{i,2} \beta_2 + u_{0,i} + x_{1,i} u_{1,i} + \epsilon_{i,j},$   
 $\mathbf{u}_i \stackrel{iid}{\sim} \text{MVN}(\mathbf{0}, \boldsymbol{\Sigma}_u), \quad \epsilon_{i,j} \stackrel{iid}{\sim} \text{Normal}(0, \sigma_{\epsilon}^2)$
- ▶  $y_{i,j} = \beta_0 + x_i \beta_1 + u_{0,i} + z_i u_{1,i} + \epsilon_{i,j},$   
 $\mathbf{u}_i \stackrel{iid}{\sim} \text{MVN}(\mathbf{0}, \boldsymbol{\Sigma}_u), \quad \epsilon_{i,j} \stackrel{iid}{\sim} \text{Normal}(0, \sigma_{\epsilon}^2)$

$$y_{i,j} = \mathbf{x}_i^T \boldsymbol{\beta} + \mathbf{z}_i^T \mathbf{u}_i + \epsilon_{i,j}, \quad \mathbf{u}_i \stackrel{iid}{\sim} f_{\mathbf{u}}, \quad \epsilon_{i,j} \stackrel{iid}{\sim} f_{\epsilon},$$

$$i = 1, \dots, n, \quad j = 1, \dots, m_i$$

with  $\mathbb{E}_{f_{\mathbf{u}}}(\mathbf{u}) = \mathbf{0}$  and  $\mathbb{E}_{f_{\epsilon}}(\epsilon) = 0$ .

- ▶  $y_{i,j} = \beta_0 + x_i \beta_1 + u_i + \epsilon_{i,j}, \quad u_i \stackrel{iid}{\sim} \text{Normal}(0, \sigma_u^2), \quad \epsilon_{i,j} \stackrel{iid}{\sim} \text{Normal}(0, \sigma_{\epsilon}^2)$
- ▶  $y_{i,j} = \beta_0 + x_i \beta_1 + x_i^2 \beta_2 + u_i + \epsilon_{i,j}, \quad u_i \stackrel{iid}{\sim} \text{Normal}(0, \sigma_u^2), \quad \epsilon_{i,j} \stackrel{iid}{\sim} \text{Normal}(0, \sigma_{\epsilon}^2)$
- ▶  $y_{i,j} = \beta_0 + x_{i,1} \beta_1 + x_{i,2} \beta_2 + \dots + x_{i,p} \beta_p + u_i + \epsilon_{i,j},$   
 $u_i \stackrel{iid}{\sim} \text{Normal}(0, \sigma_u^2), \quad \epsilon_{i,j} \stackrel{iid}{\sim} \text{Normal}(0, \sigma_{\epsilon}^2)$
- ▶  $y_{i,j} = \beta_0 + x_i \beta_1 + u_{0,i} + x_i u_{1,i} + \epsilon_{i,j}, \quad \mathbf{u}_i \stackrel{iid}{\sim} \text{MVN}(\mathbf{0}, \boldsymbol{\Sigma}_u), \quad \epsilon_{i,j} \stackrel{iid}{\sim} \text{Normal}(0, \sigma_{\epsilon}^2)$
- ▶  $y_{i,j} = \beta_0 + x_i \beta_1 + x_i^2 \beta_2 + u_{0,i} + x_i u_{1,i} + \epsilon_{i,j},$   
 $\mathbf{u}_i \stackrel{iid}{\sim} \text{MVN}(\mathbf{0}, \boldsymbol{\Sigma}_u), \quad \epsilon_{i,j} \stackrel{iid}{\sim} \text{Normal}(0, \sigma_{\epsilon}^2)$
- ▶  $y_{i,j} = \beta_0 + x_{i,1} \beta_1 + x_{i,2} \beta_2 + u_{0,i} + x_{1,i} u_{1,i} + \epsilon_{i,j},$   
 $\mathbf{u}_i \stackrel{iid}{\sim} \text{MVN}(\mathbf{0}, \boldsymbol{\Sigma}_u), \quad \epsilon_{i,j} \stackrel{iid}{\sim} \text{Normal}(0, \sigma_{\epsilon}^2)$
- ▶  $y_{i,j} = \beta_0 + x_i \beta_1 + u_{0,i} + z_i u_{1,i} + \epsilon_{i,j},$   
 $\mathbf{u}_i \stackrel{iid}{\sim} \text{MVN}(\mathbf{0}, \boldsymbol{\Sigma}_u), \quad \epsilon_{i,j} \stackrel{iid}{\sim} \text{Normal}(0, \sigma_{\epsilon}^2)$



$$y_{i,j} = \mathbf{x}_i^T \boldsymbol{\beta} + \mathbf{z}_i^T \mathbf{u}_i + \epsilon_{i,j}, \quad \mathbf{u}_i \stackrel{iid}{\sim} f_{\mathbf{u}}, \quad \epsilon_{i,j} \stackrel{iid}{\sim} f_{\epsilon},$$

$$i = 1, \dots, n, \quad j = 1, \dots, m_i$$

with  $\mathbb{E}_{f_{\mathbf{u}}}(\mathbf{u}) = \mathbf{0}$  and  $\mathbb{E}_{f_{\epsilon}}(\epsilon) = 0$ .

- ▶  $y_{i,j} = \beta_0 + x_i \beta_1 + u_i + \epsilon_{i,j}, \quad u_i \stackrel{iid}{\sim} \text{Normal}(0, \sigma_u^2), \quad \epsilon_{i,j} \stackrel{iid}{\sim} \text{Normal}(0, \sigma_{\epsilon}^2)$
- ▶  $y_{i,j} = \beta_0 + x_i \beta_1 + x_i^2 \beta_2 + u_i + \epsilon_{i,j}, \quad u_i \stackrel{iid}{\sim} \text{Normal}(0, \sigma_u^2), \quad \epsilon_{i,j} \stackrel{iid}{\sim} \text{Normal}(0, \sigma_{\epsilon}^2)$
- ▶  $y_{i,j} = \beta_0 + x_{i,1} \beta_1 + x_{i,2} \beta_2 + \dots + x_{i,p} \beta_p + u_i + \epsilon_{i,j},$   
 $u_i \stackrel{iid}{\sim} \text{Normal}(0, \sigma_u^2), \quad \epsilon_{i,j} \stackrel{iid}{\sim} \text{Normal}(0, \sigma_{\epsilon}^2)$
- ▶  $y_{i,j} = \beta_0 + x_i \beta_1 + u_{0,i} + x_i u_{1,i} + \epsilon_{i,j}, \quad \mathbf{u}_i \stackrel{iid}{\sim} \text{MVN}(\mathbf{0}, \Sigma_u), \quad \epsilon_{i,j} \stackrel{iid}{\sim} \text{Normal}(0, \sigma_{\epsilon}^2)$
- ▶  $y_{i,j} = \beta_0 + x_i \beta_1 + x_i^2 \beta_2 + u_{0,i} + x_i u_{1,i} + \epsilon_{i,j},$   
 $\mathbf{u}_i \stackrel{iid}{\sim} \text{MVN}(\mathbf{0}, \Sigma_u), \quad \epsilon_{i,j} \stackrel{iid}{\sim} \text{Normal}(0, \sigma_{\epsilon}^2)$
- ▶  $y_{i,j} = \beta_0 + x_{i,1} \beta_1 + x_{i,2} \beta_2 + u_{0,i} + x_{1,i} u_{1,i} + \epsilon_{i,j},$   
 $\mathbf{u}_i \stackrel{iid}{\sim} \text{MVN}(\mathbf{0}, \Sigma_u), \quad \epsilon_{i,j} \stackrel{iid}{\sim} \text{Normal}(0, \sigma_{\epsilon}^2)$
- ▶  $y_{i,j} = \beta_0 + x_i \beta_1 + u_{0,i} + z_i u_{1,i} + \epsilon_{i,j},$   
 $\mathbf{u}_i \stackrel{iid}{\sim} \text{MVN}(\mathbf{0}, \Sigma_u), \quad \epsilon_{i,j} \stackrel{iid}{\sim} \text{Normal}(0, \sigma_{\epsilon}^2)$

$$y_{i,j} = \mathbf{x}_i^T \boldsymbol{\beta} + \mathbf{z}_i^T \mathbf{u}_i + \epsilon_{i,j}, \quad \mathbf{u}_i \stackrel{iid}{\sim} f_{\mathbf{u}}, \quad \epsilon_{i,j} \stackrel{iid}{\sim} f_{\epsilon},$$

$$i = 1, \dots, n, \quad j = 1, \dots, m_i$$

with  $\mathbb{E}_{f_{\mathbf{u}}}(\mathbf{u}) = \mathbf{0}$  and  $\mathbb{E}_{f_{\epsilon}}(\epsilon) = 0$ .

- ▶  $y_{i,j} = \beta_0 + x_i \beta_1 + u_i + \epsilon_{i,j}, \quad u_i \stackrel{iid}{\sim} \text{Normal}(0, \sigma_u^2), \quad \epsilon_{i,j} \stackrel{iid}{\sim} \text{Normal}(0, \sigma_{\epsilon}^2)$
- ▶  $y_{i,j} = \beta_0 + x_i \beta_1 + x_i^2 \beta_2 + u_i + \epsilon_{i,j}, \quad u_i \stackrel{iid}{\sim} \text{Normal}(0, \sigma_u^2), \quad \epsilon_{i,j} \stackrel{iid}{\sim} \text{Normal}(0, \sigma_{\epsilon}^2)$
- ▶  $y_{i,j} = \beta_0 + x_{i,1} \beta_1 + x_{i,2} \beta_2 + \dots + x_{i,p} \beta_p + u_i + \epsilon_{i,j},$   
 $u_i \stackrel{iid}{\sim} \text{Normal}(0, \sigma_u^2), \quad \epsilon_{i,j} \stackrel{iid}{\sim} \text{Normal}(0, \sigma_{\epsilon}^2)$
- ▶  $y_{i,j} = \beta_0 + x_i \beta_1 + u_{0,i} + x_i u_{1,i} + \epsilon_{i,j}, \quad \mathbf{u}_i \stackrel{iid}{\sim} \text{MVN}(\mathbf{0}, \boldsymbol{\Sigma}_u), \quad \epsilon_{i,j} \stackrel{iid}{\sim} \text{Normal}(0, \sigma_{\epsilon}^2)$
- ▶  $y_{i,j} = \beta_0 + x_i \beta_1 + x_i^2 \beta_2 + u_{0,i} + x_i u_{1,i} + \epsilon_{i,j},$   
 $\mathbf{u}_i \stackrel{iid}{\sim} \text{MVN}(\mathbf{0}, \boldsymbol{\Sigma}_u), \quad \epsilon_{i,j} \stackrel{iid}{\sim} \text{Normal}(0, \sigma_{\epsilon}^2)$
- ▶  $y_{i,j} = \beta_0 + x_{i,1} \beta_1 + x_{i,2} \beta_2 + u_{0,i} + x_{1,i} u_{1,i} + \epsilon_{i,j},$   
 $\mathbf{u}_i \stackrel{iid}{\sim} \text{MVN}(\mathbf{0}, \boldsymbol{\Sigma}_u), \quad \epsilon_{i,j} \stackrel{iid}{\sim} \text{Normal}(0, \sigma_{\epsilon}^2)$
- ▶  $y_{i,j} = \beta_0 + x_i \beta_1 + u_{0,i} + z_i u_{1,i} + \epsilon_{i,j},$   
 $\mathbf{u}_i \stackrel{iid}{\sim} \text{MVN}(\mathbf{0}, \boldsymbol{\Sigma}_u), \quad \epsilon_{i,j} \stackrel{iid}{\sim} \text{Normal}(0, \sigma_{\epsilon}^2)$

$$y_{i,j} = \mathbf{x}_i^T \boldsymbol{\beta} + \mathbf{z}_i^T \mathbf{u}_i + \epsilon_{i,j}, \quad \mathbf{u}_i \stackrel{iid}{\sim} f_{\mathbf{u}}, \quad \epsilon_{i,j} \stackrel{iid}{\sim} f_{\epsilon},$$

$$i = 1, \dots, n, \quad j = 1, \dots, m_i$$

with  $\mathbb{E}_{f_{\mathbf{u}}}(\mathbf{u}) = \mathbf{0}$  and  $\mathbb{E}_{f_{\epsilon}}(\epsilon) = 0$ .

- ▶  $y_{i,j} = \beta_0 + x_i \beta_1 + u_i + \epsilon_{i,j}, \quad u_i \stackrel{iid}{\sim} \text{Normal}(0, \sigma_u^2), \quad \epsilon_{i,j} \stackrel{iid}{\sim} \text{Normal}(0, \sigma_{\epsilon}^2)$
- ▶  $y_{i,j} = \beta_0 + x_i \beta_1 + x_i^2 \beta_2 + u_i + \epsilon_{i,j}, \quad u_i \stackrel{iid}{\sim} \text{Normal}(0, \sigma_u^2), \quad \epsilon_{i,j} \stackrel{iid}{\sim} \text{Normal}(0, \sigma_{\epsilon}^2)$
- ▶  $y_{i,j} = \beta_0 + x_{i,1} \beta_1 + x_{i,2} \beta_2 + \dots + x_{i,p} \beta_p + u_i + \epsilon_{i,j},$   
 $u_i \stackrel{iid}{\sim} \text{Normal}(0, \sigma_u^2), \quad \epsilon_{i,j} \stackrel{iid}{\sim} \text{Normal}(0, \sigma_{\epsilon}^2)$
- ▶  $y_{i,j} = \beta_0 + x_i \beta_1 + u_{0,i} + x_i u_{1,i} + \epsilon_{i,j}, \quad \mathbf{u}_i \stackrel{iid}{\sim} \text{MVN}(\mathbf{0}, \boldsymbol{\Sigma}_u), \quad \epsilon_{i,j} \stackrel{iid}{\sim} \text{Normal}(0, \sigma_{\epsilon}^2)$
- ▶  $y_{i,j} = \beta_0 + x_i \beta_1 + x_i^2 \beta_2 + u_{0,i} + x_i u_{1,i} + \epsilon_{i,j},$   
 $\mathbf{u}_i \stackrel{iid}{\sim} \text{MVN}(\mathbf{0}, \boldsymbol{\Sigma}_u), \quad \epsilon_{i,j} \stackrel{iid}{\sim} \text{Normal}(0, \sigma_{\epsilon}^2)$
- ▶  $y_{i,j} = \beta_0 + x_{i,1} \beta_1 + x_{i,2} \beta_2 + u_{0,i} + x_{1,i} u_{1,i} + \epsilon_{i,j},$   
 $\mathbf{u}_i \stackrel{iid}{\sim} \text{MVN}(\mathbf{0}, \boldsymbol{\Sigma}_u), \quad \epsilon_{i,j} \stackrel{iid}{\sim} \text{Normal}(0, \sigma_{\epsilon}^2)$
- ▶  $y_{i,j} = \beta_0 + x_i \beta_1 + u_{0,i} + z_i u_{1,i} + \epsilon_{i,j},$   
 $\mathbf{u}_i \stackrel{iid}{\sim} \text{MVN}(\mathbf{0}, \boldsymbol{\Sigma}_u), \quad \epsilon_{i,j} \stackrel{iid}{\sim} \text{Normal}(0, \sigma_{\epsilon}^2)$

$$y_{i,j} = \mathbf{x}_i^T \boldsymbol{\beta} + \mathbf{z}_i^T \mathbf{u}_i + \epsilon_{i,j}, \quad \mathbf{u}_i \stackrel{iid}{\sim} f_{\mathbf{u}}, \quad \epsilon_{i,j} \stackrel{iid}{\sim} f_{\epsilon},$$

$$i = 1, \dots, n, \quad j = 1, \dots, m_i$$

with  $\mathbb{E}_{f_{\mathbf{u}}}(\mathbf{u}) = \mathbf{0}$  and  $\mathbb{E}_{f_{\epsilon}}(\epsilon) = 0$ .

- ▶  $y_{i,j} = \beta_0 + x_i \beta_1 + u_i + \epsilon_{i,j}, \quad u_i \stackrel{iid}{\sim} \text{Normal}(0, \sigma_u^2), \quad \epsilon_{i,j} \stackrel{iid}{\sim} \text{Normal}(0, \sigma_{\epsilon}^2)$
- ▶  $y_{i,j} = \beta_0 + x_i \beta_1 + x_i^2 \beta_2 + u_i + \epsilon_{i,j}, \quad u_i \stackrel{iid}{\sim} \text{Normal}(0, \sigma_u^2), \quad \epsilon_{i,j} \stackrel{iid}{\sim} \text{Normal}(0, \sigma_{\epsilon}^2)$
- ▶  $y_{i,j} = \beta_0 + x_{i,1} \beta_1 + x_{i,2} \beta_2 + \dots + x_{i,p} \beta_p + u_i + \epsilon_{i,j},$   
 $u_i \stackrel{iid}{\sim} \text{Normal}(0, \sigma_u^2), \quad \epsilon_{i,j} \stackrel{iid}{\sim} \text{Normal}(0, \sigma_{\epsilon}^2)$
- ▶  $y_{i,j} = \beta_0 + x_i \beta_1 + u_{0,i} + x_i u_{1,i} + \epsilon_{i,j}, \quad \mathbf{u}_i \stackrel{iid}{\sim} \text{MVN}(\mathbf{0}, \boldsymbol{\Sigma}_u), \quad \epsilon_{i,j} \stackrel{iid}{\sim} \text{Normal}(0, \sigma_{\epsilon}^2)$
- ▶  $y_{i,j} = \beta_0 + x_i \beta_1 + x_i^2 \beta_2 + u_{0,i} + x_i u_{1,i} + \epsilon_{i,j},$   
 $\mathbf{u}_i \stackrel{iid}{\sim} \text{MVN}(\mathbf{0}, \boldsymbol{\Sigma}_u), \quad \epsilon_{i,j} \stackrel{iid}{\sim} \text{Normal}(0, \sigma_{\epsilon}^2)$
- ▶  $y_{i,j} = \beta_0 + x_{i,1} \beta_1 + x_{i,2} \beta_2 + u_{0,i} + x_{1,i} u_{1,i} + \epsilon_{i,j},$   
 $\mathbf{u}_i \stackrel{iid}{\sim} \text{MVN}(\mathbf{0}, \boldsymbol{\Sigma}_u), \quad \epsilon_{i,j} \stackrel{iid}{\sim} \text{Normal}(0, \sigma_{\epsilon}^2)$
- ▶  $y_{i,j} = \beta_0 + x_i \beta_1 + u_{0,i} + z_i u_{1,i} + \epsilon_{i,j},$   
 $\mathbf{u}_i \stackrel{iid}{\sim} \text{MVN}(\mathbf{0}, \boldsymbol{\Sigma}_u), \quad \epsilon_{i,j} \stackrel{iid}{\sim} \text{Normal}(0, \sigma_{\epsilon}^2)$

$$y_{i,j} = \mathbf{x}_i^T \boldsymbol{\beta} + \mathbf{z}_i^T \mathbf{u}_i + \epsilon_{i,j}, \quad \mathbf{u}_i \stackrel{iid}{\sim} f_{\mathbf{u}}, \quad \epsilon_{i,j} \stackrel{iid}{\sim} f_{\epsilon},$$

$$i = 1, \dots, n, \quad j = 1, \dots, m_i$$

with  $\mathbb{E}_{f_{\mathbf{u}}}(\mathbf{u}) = \mathbf{0}$  and  $\mathbb{E}_{f_{\epsilon}}(\epsilon) = 0$ .

- ▶  $y_{i,j} = \beta_0 + x_i \beta_1 + u_i + \epsilon_{i,j}, \quad u_i \stackrel{iid}{\sim} \text{Normal}(0, \sigma_u^2), \quad \epsilon_{i,j} \stackrel{iid}{\sim} \text{Normal}(0, \sigma_{\epsilon}^2)$
- ▶  $y_{i,j} = \beta_0 + x_i \beta_1 + x_i^2 \beta_2 + u_i + \epsilon_{i,j}, \quad u_i \stackrel{iid}{\sim} \text{Normal}(0, \sigma_u^2), \quad \epsilon_{i,j} \stackrel{iid}{\sim} \text{Normal}(0, \sigma_{\epsilon}^2)$
- ▶  $y_{i,j} = \beta_0 + x_{i,1} \beta_1 + x_{i,2} \beta_2 + \dots + x_{i,p} \beta_p + u_i + \epsilon_{i,j},$   
 $u_i \stackrel{iid}{\sim} \text{Normal}(0, \sigma_u^2), \quad \epsilon_{i,j} \stackrel{iid}{\sim} \text{Normal}(0, \sigma_{\epsilon}^2)$
- ▶  $y_{i,j} = \beta_0 + x_i \beta_1 + u_{0,i} + x_i u_{1,i} + \epsilon_{i,j}, \quad \mathbf{u}_i \stackrel{iid}{\sim} \text{MVN}(\mathbf{0}, \boldsymbol{\Sigma}_u), \quad \epsilon_{i,j} \stackrel{iid}{\sim} \text{Normal}(0, \sigma_{\epsilon}^2)$
- ▶  $y_{i,j} = \beta_0 + x_i \beta_1 + x_i^2 \beta_2 + u_{0,i} + x_i u_{1,i} + \epsilon_{i,j},$   
 $\mathbf{u}_i \stackrel{iid}{\sim} \text{MVN}(\mathbf{0}, \boldsymbol{\Sigma}_u), \quad \epsilon_{i,j} \stackrel{iid}{\sim} \text{Normal}(0, \sigma_{\epsilon}^2)$
- ▶  $y_{i,j} = \beta_0 + x_{i,1} \beta_1 + x_{i,2} \beta_2 + u_{0,i} + x_{1,i} u_{1,i} + \epsilon_{i,j},$   
 $\mathbf{u}_i \stackrel{iid}{\sim} \text{MVN}(\mathbf{0}, \boldsymbol{\Sigma}_u), \quad \epsilon_{i,j} \stackrel{iid}{\sim} \text{Normal}(0, \sigma_{\epsilon}^2)$
- ▶  $y_{i,j} = \beta_0 + x_i \beta_1 + u_{0,i} + z_i u_{1,i} + \epsilon_{i,j},$   
 $\mathbf{u}_i \stackrel{iid}{\sim} \text{MVN}(\mathbf{0}, \boldsymbol{\Sigma}_u), \quad \epsilon_{i,j} \stackrel{iid}{\sim} \text{Normal}(0, \sigma_{\epsilon}^2)$

Individual level model:  $\mathbb{E}_{\epsilon}(y_{i,j}) = \mathbf{x}_i^T \boldsymbol{\beta} + \mathbf{z}_i^T \mathbf{u}_i$ ,

Population level model:  $\mathbb{E}_{\epsilon, \mathbf{u}}(y_{i,j}) = \mathbf{x}_i^T \boldsymbol{\beta}$ .

- ▶  $y_{i,j} = \beta_0 + x_i \beta_1 + u_i + \epsilon_{i,j},$   $\mathbb{E}_{\epsilon, \mathbf{u}}(y_{i,j}) = \beta_0 + x_i \beta_1$
- ▶  $y_{i,j} = \beta_0 + x_i \beta_1 + x_i^2 \beta_2 + u_i + \epsilon_{i,j},$   $\mathbb{E}_{\epsilon, \mathbf{u}}(y_{i,j}) = \beta_0 + x_i \beta_1 + x_i^2 \beta_2$
- ▶  $y_{i,j} = \beta_0 + x_{i,1} \beta_1 + \cdots + x_{i,p} \beta_p + u_i + \epsilon_{i,j},$   $\mathbb{E}_{\epsilon, \mathbf{u}}(y_{i,j}) = \beta_0 + x_{i,1} \beta_1 + \cdots + x_{i,p} \beta_p$
- ▶  $y_{i,j} = \beta_0 + x_i \beta_1 + u_{0,i} + x_i u_{1,i} + \epsilon_{i,j},$   $\mathbb{E}_{\epsilon, \mathbf{u}}(y_{i,j}) = \beta_0 + x_i \beta_1$
- ▶  $y_{i,j} = \beta_0 + x_i \beta_1 + x_i^2 \beta_2 + u_{0,i} + x_i u_{1,i} + \epsilon_{i,j},$   $\mathbb{E}_{\epsilon, \mathbf{u}}(y_{i,j}) = \beta_0 + x_i \beta_1 + x_i^2 \beta_2$
- ▶  $y_{i,j} = \beta_0 + x_{i,1} \beta_1 + x_{i,2} \beta_2 + u_{0,i} + x_{i,1} u_{1,i} + \epsilon_{i,j},$   $\mathbb{E}_{\epsilon, \mathbf{u}}(y_{i,j}) = \beta_0 + x_{i,1} \beta_1 + x_{i,2} \beta_2$
- ▶  $y_{i,j} = \beta_0 + x_{i,1} \beta_1 + u_{0,i} + z_i u_{1,i} + \epsilon_{i,j},$   $\mathbb{E}_{\epsilon, \mathbf{u}}(y_{i,j}) = \beta_0 + x_{i,1} \beta_1$

$$\begin{aligned}\text{Likelihood: } \mathcal{L}(\boldsymbol{\theta} \mid \mathbf{y}) &= \log \prod_{i=1}^n \prod_{j=1}^{m_i} f_y(y_{i,j} \mid \boldsymbol{\theta}) \\ &= \log \prod_{i=1}^n \prod_{j=1}^{m_i} \int_{\mathbb{R}^q} f_y(y_{i,j}, \mathbf{u}_i \mid \boldsymbol{\theta}) f_{\mathbf{u}}(\mathbf{u}_i \mid \boldsymbol{\theta}) d\mathbf{u}_i \\ \text{MLE: } \hat{\boldsymbol{\theta}} &= \arg \max \mathcal{L}(\boldsymbol{\theta} \mid \mathbf{y})\end{aligned}$$

- We have  $\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \mathbf{Z}\mathbf{u} + \mathbf{e} \sim \text{MVN}(\mathbf{X}\boldsymbol{\beta}, \mathbf{Z}\boldsymbol{\Sigma}_u\mathbf{Z}^T + \sigma^2\mathbf{I}) = \text{MVN}(\mathbf{X}\boldsymbol{\beta}, \mathbf{V})$ .
- If the variance components  $\mathbf{V}$  were known, we could estimate  $\boldsymbol{\beta}$  by

$$\hat{\boldsymbol{\beta}} = (\mathbf{X}^T \mathbf{V}^{-1} \mathbf{X})^{-1} \mathbf{X}^T \mathbf{V}^{-1} \mathbf{y}.$$

- If there exists a  $\mathbf{B}$ , such that  $\mathbf{B}\mathbf{X} = \mathbf{0}$ , we have

$$\mathbf{B}\mathbf{y} = \mathbf{B}\mathbf{X}\boldsymbol{\beta} + \mathbf{B}\mathbf{Z}\mathbf{u} + \mathbf{B}\mathbf{e} = \mathbf{B}\mathbf{Z}\mathbf{u} + \mathbf{B}\mathbf{e} \sim \text{MVN}(\mathbf{0}, \mathbf{B}\mathbf{V}\mathbf{B}^T).$$

- $(\mathbf{I} - \mathbf{H})$ , with  $\mathbf{H} = \mathbf{X}(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T$ , is a candidate for  $\mathbf{B}$  since  $(\mathbf{I} - \mathbf{H})\mathbf{X} = \mathbf{0}$ .
- Use  $\mathbf{B}\mathbf{y}$  to estimate  $\mathbf{V}$  by some  $\hat{\mathbf{V}}_{REML}$ .
- Estimate  $\boldsymbol{\beta}$  then by  $\hat{\boldsymbol{\beta}}_{REML} = (\mathbf{X}^T \hat{\mathbf{V}}_{REML}^{-1} \mathbf{X})^{-1} \mathbf{X}^T \hat{\mathbf{V}}_{REML}^{-1} \mathbf{y}$ .

$$\begin{aligned}\text{Likelihood: } \mathcal{L}(\boldsymbol{\theta} \mid \mathbf{y}) &= \log \prod_{i=1}^n \prod_{j=1}^{m_i} f_y(y_{i,j} \mid \boldsymbol{\theta}) \\ &= \log \prod_{i=1}^n \prod_{j=1}^{m_i} \int_{\mathbb{R}^q} f_y(y_{i,j}, \mathbf{u}_i \mid \boldsymbol{\theta}) f_{\mathbf{u}}(\mathbf{u}_i \mid \boldsymbol{\theta}) d\mathbf{u}_i \\ \text{MLE: } \hat{\boldsymbol{\theta}} &= \arg \max \mathcal{L}(\boldsymbol{\theta} \mid \mathbf{y})\end{aligned}$$

- We have  $\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \mathbf{Z}\mathbf{u} + \mathbf{e} \sim \text{MVN}(\mathbf{X}\boldsymbol{\beta}, \mathbf{Z}\boldsymbol{\Sigma}_u\mathbf{Z}^T + \sigma^2\mathbf{I}) = \text{MVN}(\mathbf{X}\boldsymbol{\beta}, \mathbf{V})$ .
- If the variance components  $\mathbf{V}$  were known, we could estimate  $\boldsymbol{\beta}$  by

$$\hat{\boldsymbol{\beta}} = (\mathbf{X}^T \mathbf{V}^{-1} \mathbf{X})^{-1} \mathbf{X}^T \mathbf{V}^{-1} \mathbf{y}.$$

- If there exists a  $\mathbf{B}$ , such that  $\mathbf{B}\mathbf{X} = \mathbf{0}$ , we have

$$\mathbf{B}\mathbf{y} = \mathbf{B}\mathbf{X}\boldsymbol{\beta} + \mathbf{B}\mathbf{Z}\mathbf{u} + \mathbf{B}\mathbf{e} = \mathbf{B}\mathbf{Z}\mathbf{u} + \mathbf{B}\mathbf{e} \sim \text{MVN}(\mathbf{0}, \mathbf{B}\mathbf{V}\mathbf{B}^T).$$

- $(\mathbf{I} - \mathbf{H})$ , with  $\mathbf{H} = \mathbf{X}(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T$ , is a candidate for  $\mathbf{B}$  since  $(\mathbf{I} - \mathbf{H})\mathbf{X} = \mathbf{0}$ .
- Use  $\mathbf{B}\mathbf{y}$  to estimate  $\mathbf{V}$  by some  $\hat{\mathbf{V}}_{REML}$ .
- Estimate  $\boldsymbol{\beta}$  then by  $\hat{\boldsymbol{\beta}}_{REML} = (\mathbf{X}^T \hat{\mathbf{V}}_{REML}^{-1} \mathbf{X})^{-1} \mathbf{X}^T \hat{\mathbf{V}}_{REML}^{-1} \mathbf{y}$ .



$$\begin{aligned}\text{Likelihood: } \mathcal{L}(\boldsymbol{\theta} \mid \mathbf{y}) &= \log \prod_{i=1}^n \prod_{j=1}^{m_i} f_y(y_{i,j} \mid \boldsymbol{\theta}) \\ &= \log \prod_{i=1}^n \prod_{j=1}^{m_i} \int_{\mathbb{R}^q} f_y(y_{i,j}, \mathbf{u}_i \mid \boldsymbol{\theta}) f_{\mathbf{u}}(\mathbf{u}_i \mid \boldsymbol{\theta}) d\mathbf{u}_i \\ \text{MLE: } \hat{\boldsymbol{\theta}} &= \arg \max \mathcal{L}(\boldsymbol{\theta} \mid \mathbf{y})\end{aligned}$$

- We have  $\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \mathbf{Z}\mathbf{u} + \mathbf{e} \sim \text{MVN}(\mathbf{X}\boldsymbol{\beta}, \mathbf{Z}\boldsymbol{\Sigma}_u\mathbf{Z}^T + \sigma^2\mathbf{I}) = \text{MVN}(\mathbf{X}\boldsymbol{\beta}, \mathbf{V})$ .
- If the variance components  $\mathbf{V}$  were known, we could estimate  $\boldsymbol{\beta}$  by

$$\hat{\boldsymbol{\beta}} = (\mathbf{X}^T \mathbf{V}^{-1} \mathbf{X})^{-1} \mathbf{X}^T \mathbf{V}^{-1} \mathbf{y}.$$

- If there exists a  $\mathbf{B}$ , such that  $\mathbf{B}\mathbf{X} = \mathbf{0}$ , we have

$$\mathbf{B}\mathbf{y} = \mathbf{B}\mathbf{X}\boldsymbol{\beta} + \mathbf{B}\mathbf{Z}\mathbf{u} + \mathbf{B}\mathbf{e} = \mathbf{B}\mathbf{Z}\mathbf{u} + \mathbf{B}\mathbf{e} \sim \text{MVN}(\mathbf{0}, \mathbf{B}\mathbf{V}\mathbf{B}^T).$$

- $(\mathbf{I} - \mathbf{H})$ , with  $\mathbf{H} = \mathbf{X}(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T$ , is a candidate for  $\mathbf{B}$  since  $(\mathbf{I} - \mathbf{H})\mathbf{X} = \mathbf{0}$ .
- Use  $\mathbf{B}\mathbf{y}$  to estimate  $\mathbf{V}$  by some  $\hat{\mathbf{V}}_{REML}$ .
- Estimate  $\boldsymbol{\beta}$  then by  $\hat{\boldsymbol{\beta}}_{REML} = (\mathbf{X}^T \hat{\mathbf{V}}_{REML}^{-1} \mathbf{X})^{-1} \mathbf{X}^T \hat{\mathbf{V}}_{REML}^{-1} \mathbf{y}$ .

$$\begin{aligned}\text{Likelihood: } \mathcal{L}(\boldsymbol{\theta} \mid \mathbf{y}) &= \log \prod_{i=1}^n \prod_{j=1}^{m_i} f_y(y_{i,j} \mid \boldsymbol{\theta}) \\ &= \log \prod_{i=1}^n \prod_{j=1}^{m_i} \int_{\mathbb{R}^q} f_y(y_{i,j}, \mathbf{u}_i \mid \boldsymbol{\theta}) f_{\mathbf{u}}(\mathbf{u}_i \mid \boldsymbol{\theta}) d\mathbf{u}_i \\ \text{MLE: } \hat{\boldsymbol{\theta}} &= \arg \max \mathcal{L}(\boldsymbol{\theta} \mid \mathbf{y})\end{aligned}$$

- We have  $\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \mathbf{Z}\mathbf{u} + \mathbf{e} \sim \text{MVN}(\mathbf{X}\boldsymbol{\beta}, \mathbf{Z}\boldsymbol{\Sigma}_u\mathbf{Z}^T + \sigma^2\mathbf{I}) = \text{MVN}(\mathbf{X}\boldsymbol{\beta}, \mathbf{V})$ .
- If the variance components  $\mathbf{V}$  were known, we could estimate  $\boldsymbol{\beta}$  by

$$\hat{\boldsymbol{\beta}} = (\mathbf{X}^T \mathbf{V}^{-1} \mathbf{X})^{-1} \mathbf{X}^T \mathbf{V}^{-1} \mathbf{y}.$$

- If there exists a  $\mathbf{B}$ , such that  $\mathbf{B}\mathbf{X} = \mathbf{0}$ , we have

$$\mathbf{B}\mathbf{y} = \mathbf{B}\mathbf{X}\boldsymbol{\beta} + \mathbf{B}\mathbf{Z}\mathbf{u} + \mathbf{B}\mathbf{e} = \mathbf{B}\mathbf{Z}\mathbf{u} + \mathbf{B}\mathbf{e} \sim \text{MVN}(\mathbf{0}, \mathbf{B}\mathbf{V}\mathbf{B}^T).$$

- $(\mathbf{I} - \mathbf{H})$ , with  $\mathbf{H} = \mathbf{X}(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T$ , is a candidate for  $\mathbf{B}$  since  $(\mathbf{I} - \mathbf{H})\mathbf{X} = \mathbf{0}$ .
- Use  $\mathbf{B}\mathbf{y}$  to estimate  $\mathbf{V}$  by some  $\hat{\mathbf{V}}_{REML}$ .
- Estimate  $\boldsymbol{\beta}$  then by  $\hat{\boldsymbol{\beta}}_{REML} = (\mathbf{X}^T \hat{\mathbf{V}}_{REML}^{-1} \mathbf{X})^{-1} \mathbf{X}^T \hat{\mathbf{V}}_{REML}^{-1} \mathbf{y}$ .

$$\begin{aligned}\text{Likelihood: } \mathcal{L}(\boldsymbol{\theta} \mid \mathbf{y}) &= \log \prod_{i=1}^n \prod_{j=1}^{m_i} f_y(y_{i,j} \mid \boldsymbol{\theta}) \\ &= \log \prod_{i=1}^n \prod_{j=1}^{m_i} \int_{\mathbb{R}^q} f_y(y_{i,j}, \mathbf{u}_i \mid \boldsymbol{\theta}) f_{\mathbf{u}}(\mathbf{u}_i \mid \boldsymbol{\theta}) d\mathbf{u}_i \\ \text{MLE: } \hat{\boldsymbol{\theta}} &= \arg \max \mathcal{L}(\boldsymbol{\theta} \mid \mathbf{y})\end{aligned}$$

- We have  $\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \mathbf{Z}\mathbf{u} + \mathbf{e} \sim \text{MVN}(\mathbf{X}\boldsymbol{\beta}, \mathbf{Z}\boldsymbol{\Sigma}_u\mathbf{Z}^T + \sigma^2\mathbf{I}) = \text{MVN}(\mathbf{X}\boldsymbol{\beta}, \mathbf{V})$ .
- If the variance components  $\mathbf{V}$  were known, we could estimate  $\boldsymbol{\beta}$  by

$$\hat{\boldsymbol{\beta}} = (\mathbf{X}^T \mathbf{V}^{-1} \mathbf{X})^{-1} \mathbf{X}^T \mathbf{V}^{-1} \mathbf{y}.$$

- If there exists a  $\mathbf{B}$ , such that  $\mathbf{B}\mathbf{X} = \mathbf{0}$ , we have

$$\mathbf{B}\mathbf{y} = \mathbf{B}\mathbf{X}\boldsymbol{\beta} + \mathbf{B}\mathbf{Z}\mathbf{u} + \mathbf{B}\mathbf{e} = \mathbf{B}\mathbf{Z}\mathbf{u} + \mathbf{B}\mathbf{e} \sim \text{MVN}(\mathbf{0}, \mathbf{B}\mathbf{V}\mathbf{B}^T).$$

- $(\mathbf{I} - \mathbf{H})$ , with  $\mathbf{H} = \mathbf{X}(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T$ , is a candidate for  $\mathbf{B}$  since  $(\mathbf{I} - \mathbf{H})\mathbf{X} = \mathbf{0}$ .
- Use  $\mathbf{B}\mathbf{y}$  to estimate  $\mathbf{V}$  by some  $\hat{\mathbf{V}}_{REML}$ .
- Estimate  $\boldsymbol{\beta}$  then by  $\hat{\boldsymbol{\beta}}_{REML} = (\mathbf{X}^T \hat{\mathbf{V}}_{REML}^{-1} \mathbf{X})^{-1} \mathbf{X}^T \hat{\mathbf{V}}_{REML}^{-1} \mathbf{y}$ .

$$\begin{aligned}\text{Likelihood: } \mathcal{L}(\boldsymbol{\theta} \mid \mathbf{y}) &= \log \prod_{i=1}^n \prod_{j=1}^{m_i} f_y(y_{i,j} \mid \boldsymbol{\theta}) \\ &= \log \prod_{i=1}^n \prod_{j=1}^{m_i} \int_{\mathbb{R}^q} f_y(y_{i,j}, \mathbf{u}_i \mid \boldsymbol{\theta}) f_{\mathbf{u}}(\mathbf{u}_i \mid \boldsymbol{\theta}) d\mathbf{u}_i \\ \text{MLE: } \hat{\boldsymbol{\theta}} &= \arg \max \mathcal{L}(\boldsymbol{\theta} \mid \mathbf{y})\end{aligned}$$

- We have  $\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \mathbf{Z}\mathbf{u} + \mathbf{e} \sim \text{MVN}(\mathbf{X}\boldsymbol{\beta}, \mathbf{Z}\boldsymbol{\Sigma}_u\mathbf{Z}^T + \sigma^2\mathbf{I}) = \text{MVN}(\mathbf{X}\boldsymbol{\beta}, \mathbf{V})$ .
- If the variance components  $\mathbf{V}$  were known, we could estimate  $\boldsymbol{\beta}$  by

$$\hat{\boldsymbol{\beta}} = (\mathbf{X}^T \mathbf{V}^{-1} \mathbf{X})^{-1} \mathbf{X}^T \mathbf{V}^{-1} \mathbf{y}.$$

- If there exists a  $\mathbf{B}$ , such that  $\mathbf{B}\mathbf{X} = \mathbf{0}$ , we have

$$\mathbf{B}\mathbf{y} = \mathbf{B}\mathbf{X}\boldsymbol{\beta} + \mathbf{B}\mathbf{Z}\mathbf{u} + \mathbf{B}\mathbf{e} = \mathbf{B}\mathbf{Z}\mathbf{u} + \mathbf{B}\mathbf{e} \sim \text{MVN}(\mathbf{0}, \mathbf{B}\mathbf{V}\mathbf{B}^T).$$

- $(\mathbf{I} - \mathbf{H})$ , with  $\mathbf{H} = \mathbf{X}(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T$ , is a candidate for  $\mathbf{B}$  since  $(\mathbf{I} - \mathbf{H})\mathbf{X} = \mathbf{0}$ .
- Use  $\mathbf{B}\mathbf{y}$  to estimate  $\mathbf{V}$  by some  $\hat{\mathbf{V}}_{REML}$ .
- Estimate  $\boldsymbol{\beta}$  then by  $\hat{\boldsymbol{\beta}}_{REML} = (\mathbf{X}^T \hat{\mathbf{V}}_{REML}^{-1} \mathbf{X})^{-1} \mathbf{X}^T \hat{\mathbf{V}}_{REML}^{-1} \mathbf{y}$ .

## Sleep Study Data

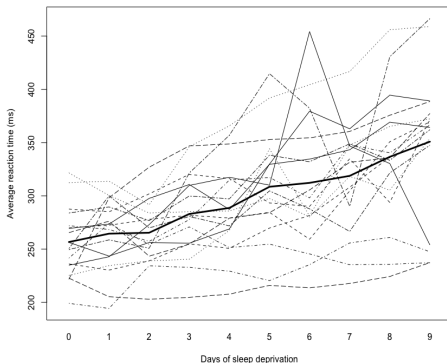
- $n = 18$  individuals, each measured  $m_i = 10$  times.
- On day 0 the subjects had their normal amount of sleep. Starting that night, they were restricted to 3 hours of sleep per night.

## Sleep Study Data

- $n = 18$  individuals, each measured  $m_i = 10$  times.
- On day 0 the subjects had their normal amount of sleep. Starting that night, they were restricted to 3 hours of sleep per night.

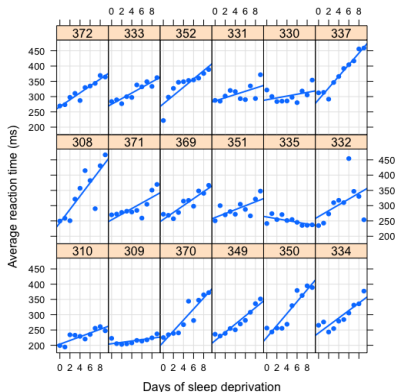
## Sleep Study Data

- $n = 18$  individuals, each measured  $m_i = 10$  times.
- On day 0 the subjects had their normal amount of sleep. Starting that night, they were restricted to 3 hours of sleep per night.
- The figure below shows average reaction time versus days of sleep deprivation for different subjects super-imposed over their averages in solid bold.



# Sleep Study Data

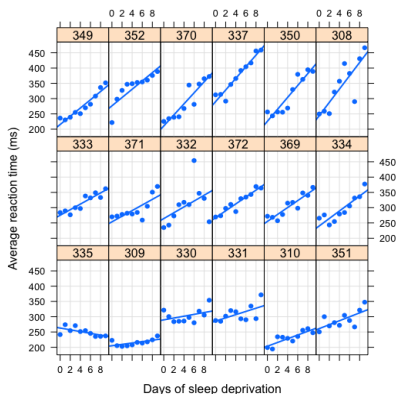
- $n = 18$  individuals, each measured  $m_i = 10$  times.
- On day 0 the subjects had their normal amount of sleep. Starting that night, they were restricted to 3 hours of sleep per night.
- The figure below shows average reaction time versus days of sleep deprivation by subject. Subjects ordered (from bottom-left to top-right) by increasing intercept of subject specific linear regressions.





## Sleep Study Data

- $n = 18$  individuals, each measured  $m_i = 10$  times.
- On day 0 the subjects had their normal amount of sleep. Starting that night, they were restricted to 3 hours of sleep per night.



- The figure above shows average reaction time versus days of sleep deprivation by subject. Subjects ordered (from bottom-left to top-right) by increasing slope of subject specific linear regressions.

$$y_{i,j} = \mathbf{x}_i^T \boldsymbol{\beta} + \mathbf{z}_i^T \mathbf{u}_i + \epsilon_{i,j}, \quad \mathbf{u}_i \stackrel{iid}{\sim} f_{\mathbf{u}}, \quad \epsilon_{i,j} \stackrel{iid}{\sim} f_{\epsilon}, \\ i = 1, \dots, n, \quad j = 1, \dots, m_i \\ \text{with } \mathbb{E}_{f_{\mathbf{u}}}(\mathbf{u}) = \mathbf{0} \quad \text{and} \quad \mathbb{E}_{f_{\epsilon}}(\epsilon) = 0.$$

$$\begin{aligned} \text{Likelihood: } \mathcal{L}(\boldsymbol{\theta} \mid \mathbf{y}) &= \log \prod_{i=1}^n \prod_{j=1}^{m_i} f_y(y_{i,j} \mid \boldsymbol{\theta}) \\ &= \log \prod_{i=1}^n \prod_{j=1}^{m_i} \int_{\mathbb{R}^q} f_y(y_{i,j}, \mathbf{u}_i \mid \boldsymbol{\theta}) f_{\mathbf{u}}(\mathbf{u}_i \mid \boldsymbol{\theta}) d\mathbf{u}_i \\ \text{Unrestricted MLE: } \hat{\boldsymbol{\theta}} &= \arg \max \mathcal{L}(\boldsymbol{\theta} \mid \mathbf{y}) \end{aligned}$$

- Penalized likelihood estimates:

$$\hat{\boldsymbol{\beta}}_{ENET} = \arg \min [\mathcal{L}(\boldsymbol{\theta} \mid \mathbf{y}) + \lambda \{ (1 - \alpha) \|\boldsymbol{\beta}\|_2^2 + 2\alpha \|\boldsymbol{\beta}\|_1 \}]$$

- Conceptually straightforward but computationally extremely challenging!

$$y_{i,j} = \mathbf{x}_i^T \boldsymbol{\beta} + \mathbf{z}_i^T \mathbf{u}_i + \epsilon_{i,j}, \quad \mathbf{u}_i \stackrel{iid}{\sim} f_{\mathbf{u}}, \quad \epsilon_{i,j} \stackrel{iid}{\sim} f_{\epsilon}, \\ i = 1, \dots, n, \quad j = 1, \dots, m_i \\ \text{with } \mathbb{E}_{f_{\mathbf{u}}}(\mathbf{u}) = \mathbf{0} \quad \text{and} \quad \mathbb{E}_{f_{\epsilon}}(\epsilon) = 0.$$

$$\begin{aligned} \text{Likelihood: } \mathcal{L}(\boldsymbol{\theta} \mid \mathbf{y}) &= \log \prod_{i=1}^n \prod_{j=1}^{m_i} f_y(y_{i,j} \mid \boldsymbol{\theta}) \\ &= \log \prod_{i=1}^n \prod_{j=1}^{m_i} \int_{\mathbb{R}^q} f_y(y_{i,j}, \mathbf{u}_i \mid \boldsymbol{\theta}) f_{\mathbf{u}}(\mathbf{u}_i \mid \boldsymbol{\theta}) d\mathbf{u}_i \\ \text{Unrestricted MLE: } \hat{\boldsymbol{\theta}} &= \arg \max \mathcal{L}(\boldsymbol{\theta} \mid \mathbf{y}) \end{aligned}$$

- Penalized likelihood estimates:

$$\hat{\boldsymbol{\beta}}_{ENET} = \arg \min [\mathcal{L}(\boldsymbol{\theta} \mid \mathbf{y}) + \lambda \{ (1 - \alpha) \|\boldsymbol{\beta}\|_2^2 + 2\alpha \|\boldsymbol{\beta}\|_1 \}]$$

- Conceptually straightforward but computationally extremely challenging!

$$y_{i,j} = \mathbf{x}_i^T \boldsymbol{\beta} + \mathbf{z}_i^T \mathbf{u}_i + \epsilon_{i,j}, \quad \mathbf{u}_i \stackrel{iid}{\sim} \text{MVN}(\mathbf{0}, \boldsymbol{\Sigma}_u), \quad \epsilon_{i,j} \stackrel{iid}{\sim} \text{Normal}(0, \sigma_\epsilon^2), \\ i = 1, \dots, n, \quad j = 1, \dots, m_i, \quad N = \sum_{i=1}^n m_i.$$

- Conditional likelihood:

$$\begin{aligned} L(\boldsymbol{\beta}, \sigma_\epsilon^2 \mid \mathbf{u}_{1:n}) &= \frac{1}{(2\pi\sigma_\epsilon^2)^{N/2}} \exp \left\{ -\frac{1}{2\sigma_\epsilon^2} (\mathbf{y} - \mathbf{X}\boldsymbol{\beta} - \mathbf{Z}\mathbf{u})^T (\mathbf{y} - \mathbf{X}\boldsymbol{\beta} - \mathbf{Z}\mathbf{u}) \right\} \\ &\propto \frac{1}{(\sigma_\epsilon^2)^{N/2}} \exp \left[ -\frac{1}{2\sigma_\epsilon^2} \left\{ \boldsymbol{\beta}^T \mathbf{X}^T \mathbf{X} \boldsymbol{\beta} - 2\boldsymbol{\beta}^T \mathbf{X}^T (\mathbf{y} - \mathbf{Z}\mathbf{u}) + (\mathbf{y} - \mathbf{Z}\mathbf{u})^T (\mathbf{y} - \mathbf{Z}\mathbf{u}) \right\} \right] \\ &\propto \frac{1}{(\sigma_\epsilon^2)^{N/2}} \exp \left\{ -\frac{1}{2\sigma_\epsilon^2} \left( \boldsymbol{\beta}^T \mathbf{X}^T \mathbf{X} \boldsymbol{\beta} - 2\boldsymbol{\beta}^T \mathbf{X}^T \tilde{\mathbf{y}} + \tilde{\mathbf{y}}^T \tilde{\mathbf{y}} \right) \right\} \quad \text{with } \tilde{\mathbf{y}} = \mathbf{y} - \mathbf{Z}\mathbf{u} \end{aligned}$$

- Semi-conjugate priors on  $(\boldsymbol{\beta}, \sigma_\epsilon^2)$ :  $p(\boldsymbol{\beta}, \sigma_\epsilon^2) = \text{MVN}(\boldsymbol{\mu}_\beta, \boldsymbol{\Sigma}_\beta) \cdot \text{Inv-Ga}(a_\sigma, b_\sigma)$

$$\propto \frac{1}{(\sigma_\epsilon^2)^{a_\sigma+1}} \exp \left( -\frac{b_\sigma}{\sigma_\epsilon^2} \right) \exp \left\{ -\frac{1}{2} \left( \boldsymbol{\beta}^T \boldsymbol{\Sigma}_\beta^{-1} \boldsymbol{\beta} - 2\boldsymbol{\beta}^T \boldsymbol{\Sigma}_\beta^{-1} \boldsymbol{\mu}_\beta + \boldsymbol{\mu}_\beta^T \boldsymbol{\Sigma}_\beta^{-1} \boldsymbol{\mu}_\beta \right) \right\}$$

- Posterior full conditionals of  $\boldsymbol{\beta}, \sigma_\epsilon^2$ :

- Conjugate prior on  $\boldsymbol{\Sigma}_u$ :

$$p(\boldsymbol{\Sigma}_u) = \text{IW}(\boldsymbol{\Sigma}_u \mid \nu_0, \boldsymbol{\Sigma}_0) \propto |\boldsymbol{\Sigma}_u|^{-\left(\frac{\nu_0+d+1}{2} + \frac{1}{2}\right)} \exp \left[ -\frac{1}{2} \left\{ \text{trace}(\boldsymbol{\Sigma}_u^{-1} \boldsymbol{\Sigma}_0) \right\} \right]$$

- Posterior full conditional of  $\boldsymbol{\Sigma}_u$ :

$$y_{i,j} = \mathbf{x}_i^T \boldsymbol{\beta} + \mathbf{z}_i^T \mathbf{u}_i + \epsilon_{i,j}, \quad \mathbf{u}_i \stackrel{iid}{\sim} \text{MVN}(\mathbf{0}, \boldsymbol{\Sigma}_u), \quad \epsilon_{i,j} \stackrel{iid}{\sim} \text{Normal}(0, \sigma_\epsilon^2), \\ i = 1, \dots, n, \quad j = 1, \dots, m_i, \quad N = \sum_{i=1}^n m_i.$$

- Conditional likelihood:

$$L(\boldsymbol{\beta}, \sigma_\epsilon^2 \mid \mathbf{u}_{1:n}) \propto \frac{1}{(\sigma_\epsilon^2)^{N/2}} \exp \left\{ -\frac{1}{2\sigma_\epsilon^2} \left( \boldsymbol{\beta}^T \mathbf{X}^T \mathbf{X} \boldsymbol{\beta} - 2\boldsymbol{\beta}^T \mathbf{X}^T \tilde{\mathbf{y}} + \tilde{\mathbf{y}}^T \tilde{\mathbf{y}} \right) \right\} \quad \text{with } \tilde{\mathbf{y}} = \mathbf{y} - \mathbf{Z}\mathbf{u}$$

- Semi-conjugate priors on  $(\boldsymbol{\beta}, \sigma_\epsilon^2)$ :  $p(\boldsymbol{\beta}, \sigma_\epsilon^2) = \text{MVN}(\boldsymbol{\mu}_\beta, \boldsymbol{\Sigma}_\beta) \cdot \text{Inv-Ga}(a_\sigma, b_\sigma)$

$$\propto \frac{1}{(\sigma_\epsilon^2)^{a_\sigma+1}} \exp \left( -\frac{b_\sigma}{\sigma_\epsilon^2} \right) \exp \left\{ -\frac{1}{2} \left( \boldsymbol{\beta}^T \boldsymbol{\Sigma}_\beta^{-1} \boldsymbol{\beta} - 2\boldsymbol{\beta}^T \boldsymbol{\Sigma}_\beta^{-1} \boldsymbol{\mu}_\beta + \boldsymbol{\mu}_\beta^T \boldsymbol{\Sigma}_\beta^{-1} \boldsymbol{\mu}_\beta \right) \right\}$$

- Posterior full conditionals of  $\boldsymbol{\beta}, \sigma_\epsilon^2$ :

$$p(\boldsymbol{\beta} \mid \sigma_\epsilon^2, \mathbf{u}_{1:n}) \propto \exp \left\{ -\frac{1}{2\sigma_\epsilon^2} \left( \boldsymbol{\beta}^T \mathbf{X}^T \mathbf{X} \boldsymbol{\beta} - 2\boldsymbol{\beta}^T \mathbf{X}^T \tilde{\mathbf{y}} + \tilde{\mathbf{y}}^T \tilde{\mathbf{y}} \right) \right\}$$

$$p(\sigma_\epsilon^2 \mid \boldsymbol{\beta}, \mathbf{u}_{1:n}) \propto \frac{1}{(\sigma_\epsilon^2)^{a_\sigma+1}} \exp \left( -\frac{b_\sigma}{\sigma_\epsilon^2} - \frac{1}{2\sigma_\epsilon^2} \left( \boldsymbol{\beta}^T \mathbf{X}^T \mathbf{X} \boldsymbol{\beta} - 2\boldsymbol{\beta}^T \mathbf{X}^T \tilde{\mathbf{y}} + \tilde{\mathbf{y}}^T \tilde{\mathbf{y}} \right) \right)$$

- Conjugate prior on  $\boldsymbol{\Sigma}_u$ :

$$p(\boldsymbol{\Sigma}_u) = \text{IW}(\boldsymbol{\Sigma}_u \mid \nu_0, \boldsymbol{\Sigma}_0) \propto |\boldsymbol{\Sigma}_u|^{-\left(\frac{\nu_0+d+1}{2} + \frac{1}{2}\right)} \exp \left[ -\frac{1}{2} \left\{ \text{trace}(\boldsymbol{\Sigma}_u^{-1} \boldsymbol{\Sigma}_0) \right\} \right]$$

- Posterior full conditional of  $\boldsymbol{\Sigma}_u$ :

$$p(\boldsymbol{\Sigma}_u \mid \boldsymbol{\beta}, \sigma_\epsilon^2, \mathbf{u}_{1:n}) \propto |\boldsymbol{\Sigma}_u|^{-\left(\frac{\nu_0+d+1}{2} + \frac{1}{2}\right)} \exp \left[ -\frac{1}{2} \left\{ \text{trace}(\boldsymbol{\Sigma}_u^{-1} \boldsymbol{\Sigma}_0) + \text{trace}(\boldsymbol{\Sigma}_u^{-1} \mathbf{S}) \right\} \right]$$

- Posterior full conditional of  $\mathbf{u}_i$ :

$$p(\mathbf{u}_i \mid \boldsymbol{\beta}, \sigma_\epsilon^2, \mathbf{u}_{-i}) \propto \exp \left\{ -\frac{1}{2\sigma_\epsilon^2} \left( \mathbf{z}_i^T \mathbf{u}_i + \mathbf{x}_i^T \boldsymbol{\beta} + \sum_{j \neq i} \mathbf{z}_i^T \mathbf{u}_i \right)^2 \right\}$$

$$y_{i,j} = \mathbf{x}_i^T \boldsymbol{\beta} + \mathbf{z}_i^T \mathbf{u}_i + \epsilon_{i,j}, \quad \mathbf{u}_i \stackrel{iid}{\sim} \text{MVN}(\mathbf{0}, \boldsymbol{\Sigma}_u), \quad \epsilon_{i,j} \stackrel{iid}{\sim} \text{Normal}(0, \sigma_\epsilon^2), \\ i = 1, \dots, n, \quad j = 1, \dots, m_i, \quad N = \sum_{i=1}^n m_i.$$

- Conditional likelihood:

$$L(\boldsymbol{\beta}, \sigma_\epsilon^2 \mid \mathbf{u}_{1:n}) \propto \frac{1}{(\sigma_\epsilon^2)^{N/2}} \exp \left\{ -\frac{1}{2\sigma_\epsilon^2} \left( \boldsymbol{\beta}^T \mathbf{X}^T \mathbf{X} \boldsymbol{\beta} - 2\boldsymbol{\beta}^T \mathbf{X}^T \tilde{\mathbf{y}} + \tilde{\mathbf{y}}^T \tilde{\mathbf{y}} \right) \right\} \quad \text{with } \tilde{\mathbf{y}} = \mathbf{y} - \mathbf{Z}\mathbf{u}$$

- Semi-conjugate priors on  $(\boldsymbol{\beta}, \sigma_\epsilon^2)$ :  $p(\boldsymbol{\beta}, \sigma_\epsilon^2) = \text{MVN}(\boldsymbol{\mu}_\beta, \boldsymbol{\Sigma}_\beta) \cdot \text{Inv-Ga}(a_\sigma, b_\sigma)$

$$\propto \frac{1}{(\sigma_\epsilon^2)^{a_\sigma+1}} \exp \left( -\frac{b_\sigma}{\sigma_\epsilon^2} \right) \exp \left\{ -\frac{1}{2} \left( \boldsymbol{\beta}^T \boldsymbol{\Sigma}_\beta^{-1} \boldsymbol{\beta} - 2\boldsymbol{\beta}^T \boldsymbol{\Sigma}_\beta^{-1} \boldsymbol{\mu}_\beta + \boldsymbol{\mu}_\beta^T \boldsymbol{\Sigma}_\beta^{-1} \boldsymbol{\mu}_\beta \right) \right\}$$

- Posterior full conditionals of  $\boldsymbol{\beta}, \sigma_\epsilon^2$ :

$$\bullet \quad p(\boldsymbol{\beta} \mid -) \propto \exp \left[ -\frac{1}{2} \left\{ \boldsymbol{\beta}^T (\boldsymbol{\Sigma}_\beta^{-1} + \sigma_\epsilon^{-2} \mathbf{X}^T \mathbf{X}) \boldsymbol{\beta} - 2\boldsymbol{\beta}^T (\boldsymbol{\Sigma}_\beta^{-1} \boldsymbol{\mu}_\beta + \sigma_\epsilon^{-2} \mathbf{X}^T \tilde{\mathbf{y}}) \right\} \right]$$

$$\equiv \text{MVN}(\boldsymbol{\mu}_{\beta,N}, \boldsymbol{\Sigma}_{\beta,N}),$$

$$\boldsymbol{\Sigma}_{\beta,N} = (\boldsymbol{\Sigma}_\beta^{-1} + \sigma_\epsilon^{-2} \mathbf{X}^T \mathbf{X})^{-1}, \quad \boldsymbol{\mu}_{\beta,N} = \boldsymbol{\Sigma}_{\beta,N} \left( \boldsymbol{\Sigma}_\beta^{-1} \boldsymbol{\mu}_\beta + \sigma_\epsilon^{-2} \mathbf{X}^T \tilde{\mathbf{y}} \right)$$

$$\bullet \quad p(\sigma_\epsilon^2 \mid -) \propto \frac{1}{(\sigma_\epsilon^2)^{a_\sigma + \frac{N}{2} + 1}} \exp \left[ -\frac{1}{\sigma_\epsilon^2} \left\{ b_\sigma + \frac{1}{2} (\tilde{\mathbf{y}} - \mathbf{X}\boldsymbol{\beta})^T (\tilde{\mathbf{y}} - \mathbf{X}\boldsymbol{\beta}) \right\} \right]$$

$$\equiv \text{Inv-Ga} \left\{ a_\sigma + \frac{N}{2}, b_\sigma + \frac{1}{2} (\tilde{\mathbf{y}} - \mathbf{X}\boldsymbol{\beta})^T (\tilde{\mathbf{y}} - \mathbf{X}\boldsymbol{\beta}) \right\}$$

- Conjugate prior on  $\boldsymbol{\Sigma}_u$ :

$$p(\boldsymbol{\Sigma}_u) = \text{IW}(\boldsymbol{\Sigma}_u \mid \nu_0, \boldsymbol{\Sigma}_0) \propto |\boldsymbol{\Sigma}_u|^{-\left(\frac{\nu_0+d+1}{2} + \frac{1}{2}\right)} \exp \left[ -\frac{1}{2} \left\{ \text{trace}(\boldsymbol{\Sigma}_u^{-1} \boldsymbol{\Sigma}_0) \right\} \right]$$

- Posterior full conditional of  $\boldsymbol{\Sigma}_u$ :

$$y_{i,j} = \mathbf{x}_i^T \boldsymbol{\beta} + \mathbf{z}_i^T \mathbf{u}_i + \epsilon_{i,j}, \quad \mathbf{u}_i \stackrel{iid}{\sim} \text{MVN}(\mathbf{0}, \boldsymbol{\Sigma}_u), \quad \epsilon_{i,j} \stackrel{iid}{\sim} \text{Normal}(0, \sigma_\epsilon^2), \\ i = 1, \dots, n, \quad j = 1, \dots, m_i, \quad N = \sum_{i=1}^n m_i.$$

- Conditional likelihood:

$$L(\boldsymbol{\beta}, \sigma_\epsilon^2 \mid \mathbf{u}_{1:n}) \propto \frac{1}{(\sigma_\epsilon^2)^{N/2}} \exp \left\{ -\frac{1}{2\sigma_\epsilon^2} \left( \boldsymbol{\beta}^T \mathbf{X}^T \mathbf{X} \boldsymbol{\beta} - 2\boldsymbol{\beta}^T \mathbf{X}^T \tilde{\mathbf{y}} + \tilde{\mathbf{y}}^T \tilde{\mathbf{y}} \right) \right\} \quad \text{with } \tilde{\mathbf{y}} = \mathbf{y} - \mathbf{Z}\mathbf{u}$$

- Semi-conjugate priors on  $(\boldsymbol{\beta}, \sigma_\epsilon^2)$ :  $p(\boldsymbol{\beta}, \sigma_\epsilon^2) = \text{MVN}(\boldsymbol{\mu}_\beta, \boldsymbol{\Sigma}_\beta) \cdot \text{Inv-Ga}(a_\sigma, b_\sigma)$

- Posterior full conditionals of  $\boldsymbol{\beta}, \sigma_\epsilon^2$ :

- $p(\boldsymbol{\beta} \mid -) = \text{MVN}(\boldsymbol{\mu}_{\beta,N}, \boldsymbol{\Sigma}_{\beta,N})$ ,

$$\boldsymbol{\Sigma}_{\beta,N} = (\boldsymbol{\Sigma}_\beta^{-1} + \sigma_\epsilon^{-2} \mathbf{X}^T \mathbf{X})^{-1}, \quad \boldsymbol{\mu}_{\beta,N} = \boldsymbol{\Sigma}_{\beta,N} \left( \boldsymbol{\Sigma}_\beta^{-1} \boldsymbol{\mu}_\beta + \sigma_\epsilon^{-2} \mathbf{X}^T \tilde{\mathbf{y}} \right)$$

- $p(\sigma_\epsilon^2 \mid -) = \text{Inv-Ga} \left\{ a_\sigma + \frac{N}{2}, b_\sigma + \frac{1}{2} (\tilde{\mathbf{y}} - \mathbf{X}\boldsymbol{\beta})^T (\tilde{\mathbf{y}} - \mathbf{X}\boldsymbol{\beta}) \right\}$

- Conjugate prior on  $\boldsymbol{\Sigma}_u$ :

$$p(\boldsymbol{\Sigma}_u) = \text{IW}(\boldsymbol{\Sigma}_u \mid \nu_0, \boldsymbol{\Sigma}_0) \propto |\boldsymbol{\Sigma}_u|^{-\left(\frac{\nu_0+d+1}{2} + \frac{1}{2}\right)} \exp \left[ -\frac{1}{2} \{ \text{trace}(\boldsymbol{\Sigma}_u^{-1} \boldsymbol{\Sigma}_0) \} \right]$$

- Posterior full conditional of  $\boldsymbol{\Sigma}_u$ :

- Posterior full conditional of  $\mathbf{u}_i$ :

$$y_{i,j} = \mathbf{x}_i^T \boldsymbol{\beta} + \mathbf{z}_i^T \mathbf{u}_i + \epsilon_{i,j}, \quad \mathbf{u}_i \stackrel{iid}{\sim} \text{MVN}(\mathbf{0}, \boldsymbol{\Sigma}_u), \quad \epsilon_{i,j} \stackrel{iid}{\sim} \text{Normal}(0, \sigma_\epsilon^2), \\ i = 1, \dots, n, \quad j = 1, \dots, m_i, \quad N = \sum_{i=1}^n m_i.$$

- Conditional likelihood:

$$L(\boldsymbol{\beta}, \sigma_\epsilon^2 \mid \mathbf{u}_{1:n}) \propto \frac{1}{(\sigma_\epsilon^2)^{N/2}} \exp \left\{ -\frac{1}{2\sigma_\epsilon^2} \left( \boldsymbol{\beta}^T \mathbf{X}^T \mathbf{X} \boldsymbol{\beta} - 2\boldsymbol{\beta}^T \mathbf{X}^T \tilde{\mathbf{y}} + \tilde{\mathbf{y}}^T \tilde{\mathbf{y}} \right) \right\} \quad \text{with } \tilde{\mathbf{y}} = \mathbf{y} - \mathbf{Z}\mathbf{u}$$

- Semi-conjugate priors on  $(\boldsymbol{\beta}, \sigma_\epsilon^2)$ :  $p(\boldsymbol{\beta}, \sigma_\epsilon^2) = \text{MVN}(\boldsymbol{\mu}_\beta, \boldsymbol{\Sigma}_\beta) \cdot \text{Inv-Ga}(a_\sigma, b_\sigma)$

- Posterior full conditionals of  $\boldsymbol{\beta}, \sigma_\epsilon^2$ :

- $p(\boldsymbol{\beta} \mid -) = \text{MVN}(\boldsymbol{\mu}_{\beta,N}, \boldsymbol{\Sigma}_{\beta,N}),$

$$\boldsymbol{\Sigma}_{\beta,N} = (\boldsymbol{\Sigma}_\beta^{-1} + \sigma_\epsilon^{-2} \mathbf{X}^T \mathbf{X})^{-1}, \quad \boldsymbol{\mu}_{\beta,N} = \boldsymbol{\Sigma}_{\beta,N} \left( \boldsymbol{\Sigma}_\beta^{-1} \boldsymbol{\mu}_\beta + \sigma_\epsilon^{-2} \mathbf{X}^T \tilde{\mathbf{y}} \right)$$

- $p(\sigma_\epsilon^2 \mid -) = \text{Inv-Ga} \left\{ a_\sigma + \frac{N}{2}, b_\sigma + \frac{1}{2} (\tilde{\mathbf{y}} - \mathbf{X}\boldsymbol{\beta})^T (\tilde{\mathbf{y}} - \mathbf{X}\boldsymbol{\beta}) \right\}$

- Conjugate prior on  $\boldsymbol{\Sigma}_u$ :

$$p(\boldsymbol{\Sigma}_u) = \text{IW}(\boldsymbol{\Sigma}_u \mid \nu_0, \boldsymbol{\Sigma}_0) \propto |\boldsymbol{\Sigma}_u|^{-\left(\frac{\nu_0+d+1}{2} + \frac{1}{2}\right)} \exp \left[ -\frac{1}{2} \left\{ \text{trace}(\boldsymbol{\Sigma}_u^{-1} \boldsymbol{\Sigma}_0) \right\} \right]$$

- Posterior full conditional of  $\boldsymbol{\Sigma}_u$ :

- $p(\boldsymbol{\Sigma}_u \mid -) = \text{IW}(n + \nu_0, \boldsymbol{\Sigma}_0 + \sum_{i=1}^n \mathbf{u}_i \mathbf{u}_i^T)$

- Posterior full conditional of  $\mathbf{u}_i$ :



$$y_{i,j} = \mathbf{x}_i^T \boldsymbol{\beta} + \mathbf{z}_i^T \mathbf{u}_i + \epsilon_{i,j}, \quad \mathbf{u}_i \stackrel{iid}{\sim} \text{MVN}(\mathbf{0}, \boldsymbol{\Sigma}_u), \quad \epsilon_{i,j} \stackrel{iid}{\sim} \text{Normal}(0, \sigma_\epsilon^2),$$

$$i = 1, \dots, n, \quad j = 1, \dots, m_i, \quad N = \sum_{i=1}^n m_i.$$

- Conditional likelihood:

$$L(\boldsymbol{\beta}, \sigma_\epsilon^2 \mid \mathbf{u}_{1:n}) \propto \frac{1}{(\sigma_\epsilon^2)^{N/2}} \exp \left\{ -\frac{1}{2\sigma_\epsilon^2} \left( \boldsymbol{\beta}^T \mathbf{X}^T \mathbf{X} \boldsymbol{\beta} - 2\boldsymbol{\beta}^T \mathbf{X}^T \tilde{\mathbf{y}} + \tilde{\mathbf{y}}^T \tilde{\mathbf{y}} \right) \right\} \quad \text{with } \tilde{\mathbf{y}} = \mathbf{y} - \mathbf{Z}\mathbf{u}$$

- Semi-conjugate priors on  $(\boldsymbol{\beta}, \sigma_\epsilon^2)$ :  $p(\boldsymbol{\beta}, \sigma_\epsilon^2) = \text{MVN}(\boldsymbol{\mu}_\beta, \boldsymbol{\Sigma}_\beta) \cdot \text{Inv-Ga}(a_\sigma, b_\sigma)$

- Posterior full conditionals of  $\boldsymbol{\beta}, \sigma_\epsilon^2$ :

- $p(\boldsymbol{\beta} \mid -) = \text{MVN}(\boldsymbol{\mu}_{\beta,N}, \boldsymbol{\Sigma}_{\beta,N}),$

$$\boldsymbol{\Sigma}_{\beta,N} = (\boldsymbol{\Sigma}_\beta^{-1} + \sigma_\epsilon^{-2} \mathbf{X}^T \mathbf{X})^{-1}, \quad \boldsymbol{\mu}_{\beta,N} = \boldsymbol{\Sigma}_{\beta,N} \left( \boldsymbol{\Sigma}_\beta^{-1} \boldsymbol{\mu}_\beta + \sigma_\epsilon^{-2} \mathbf{X}^T \tilde{\mathbf{y}} \right)$$

- $p(\sigma_\epsilon^2 \mid -) = \text{Inv-Ga} \left\{ a_\sigma + \frac{N}{2}, b_\sigma + \frac{1}{2} (\tilde{\mathbf{y}} - \mathbf{X}\boldsymbol{\beta})^T (\tilde{\mathbf{y}} - \mathbf{X}\boldsymbol{\beta}) \right\}$

- Conjugate prior on  $\boldsymbol{\Sigma}_u$ :

$$p(\boldsymbol{\Sigma}_u) = \text{IW}(\boldsymbol{\Sigma}_u \mid \nu_0, \boldsymbol{\Sigma}_0) \propto |\boldsymbol{\Sigma}_u|^{-\left(\frac{\nu_0+d+1}{2} + \frac{1}{2}\right)} \exp \left[ -\frac{1}{2} \left\{ \text{trace}(\boldsymbol{\Sigma}_u^{-1} \boldsymbol{\Sigma}_0) \right\} \right]$$

- Posterior full conditional of  $\boldsymbol{\Sigma}_u$ :

- $p(\boldsymbol{\Sigma}_u \mid -) = \text{IW}(n + \nu_0, \boldsymbol{\Sigma}_0 + \sum_{i=1}^n \mathbf{u}_i \mathbf{u}_i^T)$

- Posterior full conditional of  $\mathbf{u}_i$ :

- $p(\mathbf{u}_i \mid -) \propto \exp \left[ -\frac{1}{2} \left\{ \mathbf{u}_i^T \boldsymbol{\Sigma}_u^{-1} \mathbf{u}_i + \sigma_\epsilon^{-2} \sum_{j=1}^{m_i} (y_{i,j} - \mathbf{x}_i^T \boldsymbol{\beta} - \mathbf{z}_i^T \mathbf{u}_i)^2 \right\} \right]$   
 $\propto \exp \left[ -\frac{1}{2} \left\{ \mathbf{u}_i^T (\boldsymbol{\Sigma}_u^{-1} + \sigma_\epsilon^{-2} m_i \mathbf{z}_i \mathbf{z}_i^T) \mathbf{u}_i - 2\sigma_\epsilon^{-2} \sum_{j=1}^{m_i} (y_{i,j} - \mathbf{x}_i^T \boldsymbol{\beta}) \mathbf{z}_i^T \mathbf{u}_i \right\} \right]$

$$y_{i,j} = \mathbf{x}_i^T \boldsymbol{\beta} + \mathbf{z}_i^T \mathbf{u}_i + \epsilon_{i,j}, \quad \mathbf{u}_i \stackrel{iid}{\sim} \text{MVN}(\mathbf{0}, \boldsymbol{\Sigma}_u), \quad \epsilon_{i,j} \stackrel{iid}{\sim} \text{Normal}(0, \sigma_\epsilon^2),$$

$$i = 1, \dots, n, \quad j = 1, \dots, m_i, \quad N = \sum_{i=1}^n m_i.$$

- Conditional likelihood:

$$L(\boldsymbol{\beta}, \sigma_\epsilon^2 \mid \mathbf{u}_{1:n}) \propto \frac{1}{(\sigma_\epsilon^2)^{N/2}} \exp \left\{ -\frac{1}{2\sigma_\epsilon^2} \left( \boldsymbol{\beta}^T \mathbf{X}^T \mathbf{X} \boldsymbol{\beta} - 2\boldsymbol{\beta}^T \mathbf{X}^T \tilde{\mathbf{y}} + \tilde{\mathbf{y}}^T \tilde{\mathbf{y}} \right) \right\} \quad \text{with } \tilde{\mathbf{y}} = \mathbf{y} - \mathbf{Z}\mathbf{u}$$

- Semi-conjugate priors on  $(\boldsymbol{\beta}, \sigma_\epsilon^2)$ :  $p(\boldsymbol{\beta}, \sigma_\epsilon^2) = \text{MVN}(\boldsymbol{\mu}_\beta, \boldsymbol{\Sigma}_\beta) \cdot \text{Inv-Ga}(a_\sigma, b_\sigma)$

- Posterior full conditionals of  $\boldsymbol{\beta}, \sigma_\epsilon^2$ :

- $p(\boldsymbol{\beta} \mid -) = \text{MVN}(\boldsymbol{\mu}_{\beta,N}, \boldsymbol{\Sigma}_{\beta,N}),$

$$\boldsymbol{\Sigma}_{\beta,N} = (\boldsymbol{\Sigma}_\beta^{-1} + \sigma_\epsilon^{-2} \mathbf{X}^T \mathbf{X})^{-1}, \quad \boldsymbol{\mu}_{\beta,N} = \boldsymbol{\Sigma}_{\beta,N} \left( \boldsymbol{\Sigma}_\beta^{-1} \boldsymbol{\mu}_\beta + \sigma_\epsilon^{-2} \mathbf{X}^T \tilde{\mathbf{y}} \right)$$

- $p(\sigma_\epsilon^2 \mid -) = \text{Inv-Ga} \left\{ a_\sigma + \frac{N}{2}, b_\sigma + \frac{1}{2} (\tilde{\mathbf{y}} - \mathbf{X}\boldsymbol{\beta})^T (\tilde{\mathbf{y}} - \mathbf{X}\boldsymbol{\beta}) \right\}$

- Conjugate prior on  $\boldsymbol{\Sigma}_u$ :

$$p(\boldsymbol{\Sigma}_u) = \text{IW}(\boldsymbol{\Sigma}_u \mid \nu_0, \boldsymbol{\Sigma}_0) \propto |\boldsymbol{\Sigma}_u|^{-\left(\frac{\nu_0+d+1}{2} + \frac{1}{2}\right)} \exp \left[ -\frac{1}{2} \left\{ \text{trace}(\boldsymbol{\Sigma}_u^{-1} \boldsymbol{\Sigma}_0) \right\} \right]$$

- Posterior full conditional of  $\boldsymbol{\Sigma}_u$ :

- $p(\boldsymbol{\Sigma}_u \mid -) = \text{IW}(n + \nu_0, \boldsymbol{\Sigma}_0 + \sum_{i=1}^n \mathbf{u}_i \mathbf{u}_i^T)$

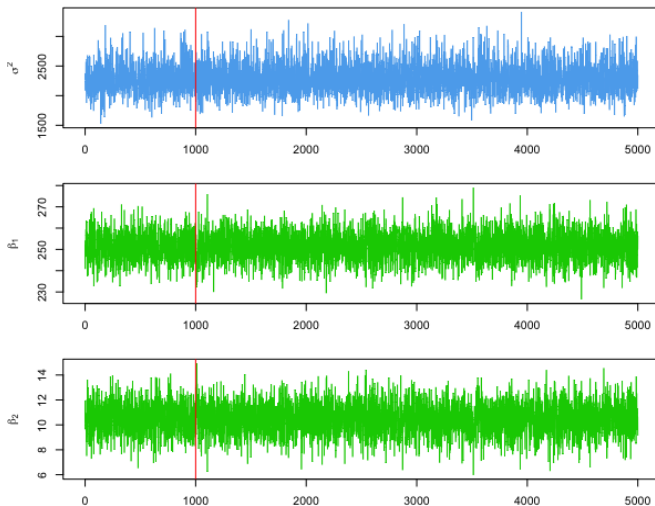
- Posterior full conditional of  $\mathbf{u}_i$ :

- $p(\mathbf{u}_i \mid -) = \text{MVN}(\boldsymbol{\mu}_{i,u,N}, \boldsymbol{\Sigma}_{i,u,N}),$

$$\boldsymbol{\Sigma}_{i,u,N} = (\boldsymbol{\Sigma}_u^{-1} + \sigma_\epsilon^{-2} m_i \mathbf{z}_i \mathbf{z}_i^T)^{-1}, \quad \boldsymbol{\mu}_{i,u,N} = \boldsymbol{\Sigma}_{i,u,N} \left\{ \sigma_\epsilon^{-2} \sum_{j=1}^{m_i} (y_{i,j} - \mathbf{x}_i^T \boldsymbol{\beta}) \mathbf{z}_i \right\}$$

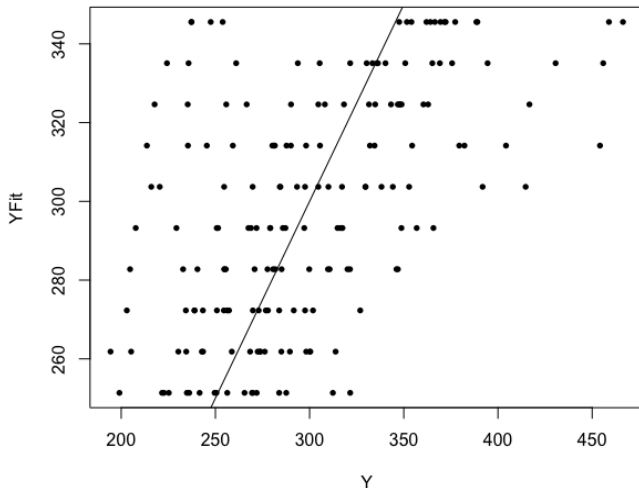
## Sleep Study Data - Bayesian LM - Fixed Effects Only

$$\text{ReactionTime}_{i,j} = \beta_0 + \beta_1 \text{Day}_j + \epsilon_{i,j}, \quad \epsilon_{i,j} \stackrel{iid}{\sim} \text{Normal}(0, \sigma_\epsilon^2),$$
$$i = 1, \dots, n = 18, \quad j = 1, \dots, m_i = 10, \quad N = \sum_{i=1}^n m_i = 180.$$



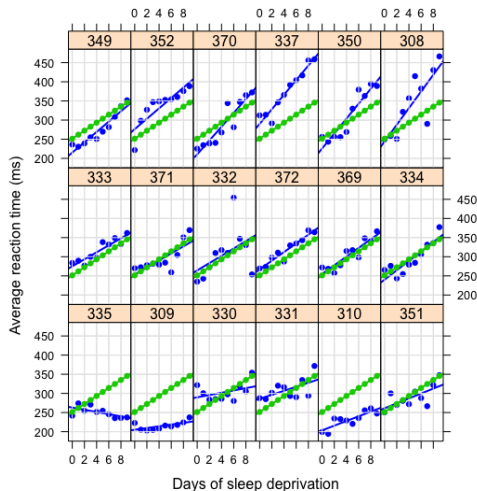
## Sleep Study Data - Bayesian LM - Fixed Effects Only

$$\text{ReactionTime}_{i,j} = \beta_0 + \beta_1 \text{Day}_j + \epsilon_{i,j}, \quad \epsilon_{i,j} \stackrel{iid}{\sim} \text{Normal}(0, \sigma_\epsilon^2),$$
$$i = 1, \dots, n = 18, \quad j = 1, \dots, m_i = 10, \quad N = \sum_{i=1}^n m_i = 180.$$



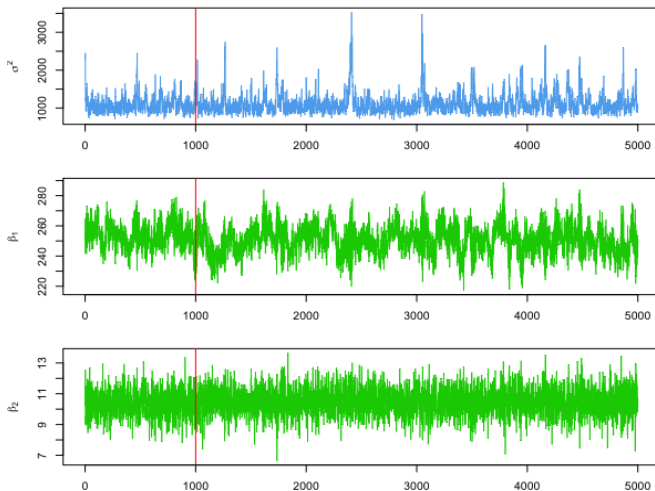
## Sleep Study Data - Bayesian LM - Fixed Effects Only

$$\text{ReactionTime}_{i,j} = \beta_0 + \beta_1 \text{Day}_j + \epsilon_{i,j}, \quad \epsilon_{i,j} \stackrel{iid}{\sim} \text{Normal}(0, \sigma_\epsilon^2),$$
$$i = 1, \dots, n = 18, \quad j = 1, \dots, m_i = 10, \quad N = \sum_{i=1}^n m_i = 180.$$



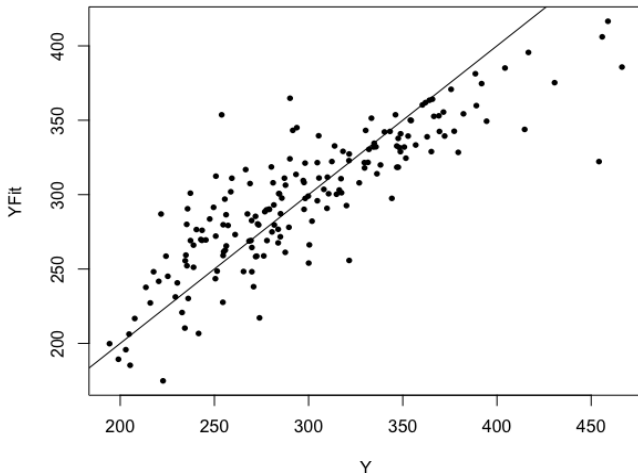
## Sleep Study Data - Bayesian LMM - Random Intercept Only

$$\text{ReactionTime}_{i,j} = \beta_0 + \beta_1 \text{Day}_j + u_i + \epsilon_{i,j}, \quad u_i \sim \text{Normal}(\sigma_u^2), \quad \epsilon_{i,j} \stackrel{iid}{\sim} \text{Normal}(0, \sigma_\epsilon^2), \\ i = 1, \dots, n = 18, \quad j = 1, \dots, m_i = 10, \quad N = \sum_{i=1}^n m_i = 180.$$



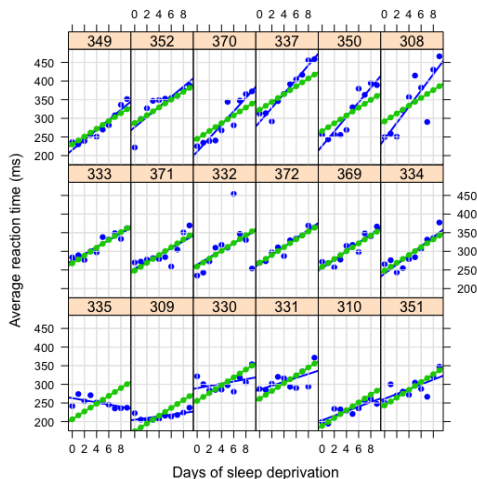
## Sleep Study Data - Bayesian LMM - Random Intercept Only

$$\text{ReactionTime}_{i,j} = \beta_0 + \beta_1 \text{Day}_j + u_i + \epsilon_{i,j}, \quad u_i \sim \text{Normal}(\sigma_u^2), \quad \epsilon_{i,j} \stackrel{iid}{\sim} \text{Normal}(0, \sigma_\epsilon^2), \\ i = 1, \dots, n = 18, \quad j = 1, \dots, m_i = 10, \quad N = \sum_{i=1}^n m_i = 180.$$



# Sleep Study Data - Bayesian LMM - Random Intercept Only

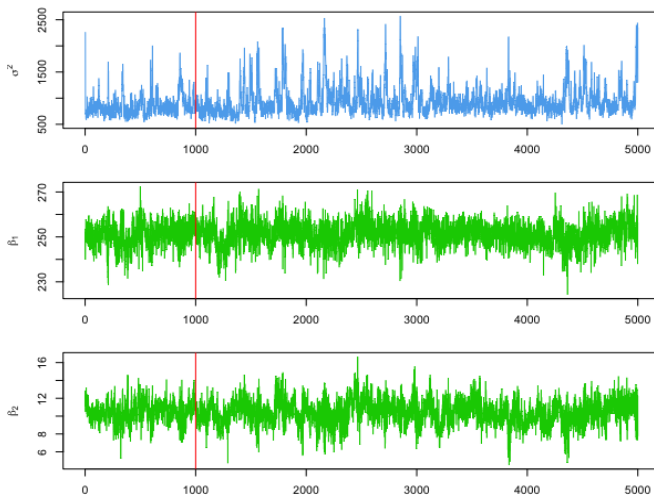
$$\text{ReactionTime}_{i,j} = \beta_0 + \beta_1 \text{Day}_j + u_i + \epsilon_{i,j}, \quad u_i \sim \text{Normal}(\sigma_u^2), \quad \epsilon_{i,j} \stackrel{iid}{\sim} \text{Normal}(0, \sigma_\epsilon^2),$$
$$i = 1, \dots, n = 18, \quad j = 1, \dots, m_i = 10, \quad N = \sum_{i=1}^n m_i = 180.$$





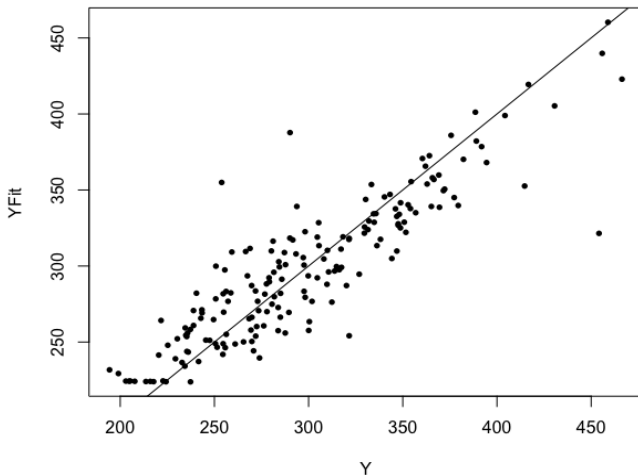
## Sleep Study Data - Bayesian LMM - Random Intercept and Slope

$$\text{ReactionTime}_{i,j} = \beta_0 + \beta_1 \text{Day}_j + u_{0,i} + u_{1,i} \text{Day}_j + \epsilon_{i,j}, \quad u_i \sim \text{Normal}(\sigma_u^2), \\ \epsilon_{i,j} \stackrel{iid}{\sim} \text{Normal}(0, \sigma_\epsilon^2), \quad i = 1, \dots, n = 18, \quad j = 1, \dots, m_i = 10, \quad N = \sum_{i=1}^n m_i = 180.$$



## Sleep Study Data - Bayesian LMM - Random Intercept and Slope

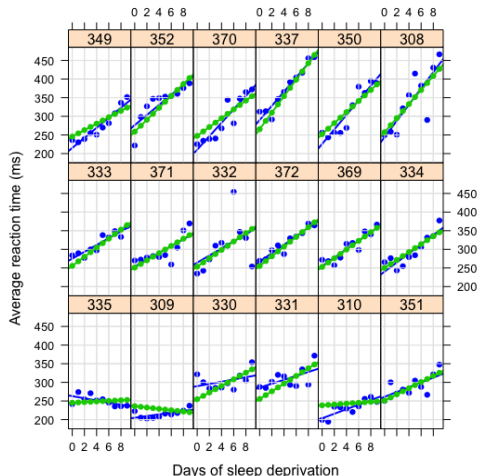
$$\text{ReactionTime}_{i,j} = \beta_0 + \beta_1 \text{Day}_j + u_{0,i} + u_{1,i} \text{Day}_j + \epsilon_{i,j}, \quad u_i \sim \text{Normal}(\sigma_u^2), \\ \epsilon_{i,j} \stackrel{iid}{\sim} \text{Normal}(0, \sigma_\epsilon^2), \quad i = 1, \dots, n = 18, \quad j = 1, \dots, m_i = 10, \quad N = \sum_{i=1}^n m_i = 180.$$



# Sleep Study Data - Bayesian LMM - Random Intercept and Slope

$$\text{ReactionTime}_{i,j} = \beta_0 + \beta_1 \text{Day}_j + u_{0,i} + u_{1,i} \text{Day}_j + \epsilon_{i,j}, \quad u_i \sim \text{Normal}(\sigma_u^2),$$

$$\epsilon_{i,j} \stackrel{iid}{\sim} \text{Normal}(0, \sigma_\epsilon^2), \quad i = 1, \dots, n = 18, \quad j = 1, \dots, m_i = 10, \quad N = \sum_{i=1}^n m_i = 180.$$



# Bayesian Ridge for Linear Mixed Models

$$y_{i,j} = \mathbf{x}_i^T \boldsymbol{\beta} + \mathbf{z}_i^T \mathbf{u}_i + \epsilon_{i,j}, \quad \mathbf{u}_i \stackrel{iid}{\sim} \text{MVN}(\mathbf{0}, \boldsymbol{\Sigma}_u), \quad \epsilon_{i,j} \stackrel{iid}{\sim} \text{Normal}(0, \sigma_\epsilon^2), \\ i = 1, \dots, n, \quad j = 1, \dots, m_i, \quad N = \sum_{i=1}^n m_i.$$

- Semi-conjugate priors on  $(\boldsymbol{\beta}, \sigma_\epsilon^2)$ :

$$p(\boldsymbol{\beta}, \sigma_\epsilon^2) = \text{MVN}(\boldsymbol{\mu}_\beta, \boldsymbol{\Sigma}_\beta) \cdot \text{Inv-Ga}(a_\sigma, b_\sigma) \quad \text{with } \boldsymbol{\Sigma}_\beta = \lambda^{-1} \sigma_\epsilon^2 \mathbf{I}_p$$

- Posterior full conditionals of  $\boldsymbol{\beta}, \sigma_\epsilon^2$ :

$$p(\boldsymbol{\beta} | \sigma_\epsilon^2, \mathbf{y}) \propto \exp\left\{-\frac{1}{2\sigma_\epsilon^2} \left(\mathbf{y} - \mathbf{X}\boldsymbol{\beta} - \mathbf{Z}\mathbf{U}\right)^T \left(\mathbf{y} - \mathbf{X}\boldsymbol{\beta} - \mathbf{Z}\mathbf{U}\right)\right\} \exp\left\{-\frac{1}{2\sigma_\epsilon^2} \boldsymbol{\beta}^T \mathbf{X}^T \mathbf{X} \boldsymbol{\beta}\right\} \exp\left\{-\frac{1}{2\sigma_\epsilon^2} \boldsymbol{\beta}^T \mathbf{X}^T \mathbf{y}\right\} \\ \exp\left\{-\frac{1}{2\sigma_\epsilon^2} \mathbf{y}^T \mathbf{y}\right\} \exp\left\{-\frac{1}{2\sigma_\epsilon^2} \mathbf{y}^T \mathbf{Z} \mathbf{U}\right\} \exp\left\{-\frac{1}{2\sigma_\epsilon^2} \mathbf{U}^T \mathbf{Z}^T \mathbf{y}\right\} \exp\left\{-\frac{1}{2\sigma_\epsilon^2} \mathbf{U}^T \mathbf{Z}^T \mathbf{Z} \mathbf{U}\right\}$$

- Conjugate prior on  $\lambda$ :  $p(\lambda) = \text{Ga}(a_\lambda, b_\lambda)$
- Posterior full conditional of  $\lambda$ :

$$p(\lambda | \boldsymbol{\beta}, \sigma_\epsilon^2, \mathbf{y}) \propto \lambda^{a_\lambda - 1} \exp\left\{-\frac{b_\lambda}{\lambda}\right\} \exp\left\{-\frac{1}{2\lambda} \boldsymbol{\beta}^T \mathbf{X}^T \mathbf{X} \boldsymbol{\beta}\right\}$$

- Conjugate prior on  $\boldsymbol{\Sigma}_u$ :  $p(\boldsymbol{\Sigma}_u) = \text{IW}(\boldsymbol{\Sigma}_u | \nu_0, \boldsymbol{\Sigma}_0)$
- Posterior full conditional of  $\boldsymbol{\Sigma}_u$ :

$$p(\boldsymbol{\Sigma}_u | \boldsymbol{\beta}, \sigma_\epsilon^2, \mathbf{y}) \propto |\boldsymbol{\Sigma}_u|^{-\frac{\nu_0 + n}{2}} \exp\left\{-\frac{1}{2} \text{tr}\left(\boldsymbol{\Sigma}_u^{-1} \left(\mathbf{U}^T \mathbf{U} + \boldsymbol{\Sigma}_0\right)\right)\right\}$$

- Posterior full conditional of  $\mathbf{u}_i$ :

$$p(\mathbf{u}_i | \boldsymbol{\beta}, \sigma_\epsilon^2, \mathbf{y}) \propto \exp\left\{-\frac{1}{2\sigma_\epsilon^2} \left(\mathbf{y}_i - \mathbf{x}_i^T \boldsymbol{\beta} - \mathbf{z}_i^T \mathbf{u}_i\right)^T \left(\mathbf{y}_i - \mathbf{x}_i^T \boldsymbol{\beta} - \mathbf{z}_i^T \mathbf{u}_i\right)\right\} \exp\left\{-\frac{1}{2} \mathbf{u}_i^T \boldsymbol{\Sigma}_u^{-1} \mathbf{u}_i\right\}$$

# Bayesian Ridge for Linear Mixed Models

$$y_{i,j} = \mathbf{x}_i^T \boldsymbol{\beta} + \mathbf{z}_i^T \mathbf{u}_i + \epsilon_{i,j}, \quad \mathbf{u}_i \stackrel{iid}{\sim} \text{MVN}(\mathbf{0}, \boldsymbol{\Sigma}_u), \quad \epsilon_{i,j} \stackrel{iid}{\sim} \text{Normal}(0, \sigma_\epsilon^2), \\ i = 1, \dots, n, \quad j = 1, \dots, m_i, \quad N = \sum_{i=1}^n m_i.$$

- Semi-conjugate priors on  $(\boldsymbol{\beta}, \sigma_\epsilon^2)$ :

$$p(\boldsymbol{\beta}, \sigma_\epsilon^2) = \text{MVN}(\boldsymbol{\mu}_\beta, \boldsymbol{\Sigma}_\beta) \cdot \text{Inv-Ga}(a_\sigma, b_\sigma) \quad \text{with} \quad \boldsymbol{\Sigma}_\beta = \lambda^{-1} \sigma_\epsilon^2 \mathbf{I}_p$$

- Posterior full conditionals of  $\boldsymbol{\beta}, \sigma_\epsilon^2$ :

- $p(\boldsymbol{\beta} | -) = \text{MVN}(\boldsymbol{\mu}_{\beta,N}, \boldsymbol{\Sigma}_{\beta,N}),$

$$\boldsymbol{\Sigma}_{\beta,N} = (\boldsymbol{\Sigma}_\beta^{-1} + \sigma_\epsilon^{-2} \mathbf{X}^T \mathbf{X})^{-1}, \quad \boldsymbol{\mu}_{\beta,N} = \boldsymbol{\Sigma}_{\beta,N} \left( \boldsymbol{\Sigma}_\beta^{-1} \boldsymbol{\mu}_\beta + \sigma_\epsilon^{-2} \mathbf{X}^T \tilde{\mathbf{y}} \right)$$

- $p(\sigma_\epsilon^2 | -) = \text{Inv-Ga} \left\{ a_\sigma + \frac{N}{2}, b_\sigma + \frac{1}{2} (\tilde{\mathbf{y}} - \mathbf{X}\boldsymbol{\beta})^T (\tilde{\mathbf{y}} - \mathbf{X}\boldsymbol{\beta}) + \frac{\lambda}{2} \boldsymbol{\beta}^T \boldsymbol{\beta} \right\}$

- Conjugate prior on  $\lambda$ :  $p(\lambda) = \text{Ga}(a_\lambda, b_\lambda)$

- Posterior full conditional of  $\lambda$ :

$$p(\lambda | -) = \text{Ga} \left( a_\lambda + \frac{N}{2}, b_\lambda + \frac{1}{2} (\tilde{\mathbf{y}} - \mathbf{X}\boldsymbol{\beta})^T (\tilde{\mathbf{y}} - \mathbf{X}\boldsymbol{\beta}) + \frac{\lambda}{2} \boldsymbol{\beta}^T \boldsymbol{\beta} \right)$$

- Conjugate prior on  $\boldsymbol{\Sigma}_u$ :  $p(\boldsymbol{\Sigma}_u) = \text{IW}(\boldsymbol{\Sigma}_u | \nu_0, \boldsymbol{\Sigma}_0)$

- Posterior full conditional of  $\boldsymbol{\Sigma}_u$ :

$$p(\boldsymbol{\Sigma}_u | -) = \text{IW} \left( \boldsymbol{\Sigma}_u \mid \nu_0 + n, \boldsymbol{\Sigma}_0 + \sum_{i=1}^n (\mathbf{z}_i \mathbf{z}_i^T + \tilde{\mathbf{u}}_i \tilde{\mathbf{u}}_i^T) \right)$$

- Posterior full conditional of  $\mathbf{u}_i$ :

$$p(\mathbf{u}_i | -) = \text{MVN} \left( \mathbf{u}_i \mid \boldsymbol{\Sigma}_u^{-1} (\mathbf{z}_i^T \tilde{\mathbf{y}} + \boldsymbol{\Sigma}_u^{-1} \mathbf{u}_i) \right)$$

# Bayesian Ridge for Linear Mixed Models

$$y_{i,j} = \mathbf{x}_i^T \boldsymbol{\beta} + \mathbf{z}_i^T \mathbf{u}_i + \epsilon_{i,j}, \quad \mathbf{u}_i \stackrel{iid}{\sim} \text{MVN}(\mathbf{0}, \boldsymbol{\Sigma}_u), \quad \epsilon_{i,j} \stackrel{iid}{\sim} \text{Normal}(0, \sigma_\epsilon^2), \\ i = 1, \dots, n, \quad j = 1, \dots, m_i, \quad N = \sum_{i=1}^n m_i.$$

- Semi-conjugate priors on  $(\boldsymbol{\beta}, \sigma_\epsilon^2)$ :

$$p(\boldsymbol{\beta}, \sigma_\epsilon^2) = \text{MVN}(\boldsymbol{\mu}_\beta, \boldsymbol{\Sigma}_\beta) \cdot \text{Inv-Ga}(a_\sigma, b_\sigma) \quad \text{with} \quad \boldsymbol{\Sigma}_\beta = \lambda^{-1} \sigma_\epsilon^2 \mathbf{I}_p$$

- Posterior full conditionals of  $\boldsymbol{\beta}, \sigma_\epsilon^2$ :

- $p(\boldsymbol{\beta} \mid -) = \text{MVN}(\boldsymbol{\mu}_{\beta,N}, \boldsymbol{\Sigma}_{\beta,N}),$

$$\boldsymbol{\Sigma}_{\beta,N} = (\boldsymbol{\Sigma}_\beta^{-1} + \sigma_\epsilon^{-2} \mathbf{X}^T \mathbf{X})^{-1}, \quad \boldsymbol{\mu}_{\beta,N} = \boldsymbol{\Sigma}_{\beta,N} \left( \boldsymbol{\Sigma}_\beta^{-1} \boldsymbol{\mu}_\beta + \sigma_\epsilon^{-2} \mathbf{X}^T \tilde{\mathbf{y}} \right)$$

- $p(\sigma_\epsilon^2 \mid -) = \text{Inv-Ga} \left\{ a_\sigma + \frac{N}{2}, b_\sigma + \frac{1}{2} (\tilde{\mathbf{y}} - \mathbf{X}\boldsymbol{\beta})^T (\tilde{\mathbf{y}} - \mathbf{X}\boldsymbol{\beta}) + \frac{\lambda}{2} \boldsymbol{\beta}^T \boldsymbol{\beta} \right\}$

- Conjugate prior on  $\lambda$ :  $p(\lambda) = \text{Ga}(a_\lambda, b_\lambda)$

- Posterior full conditional of  $\lambda$ :

- Conjugate prior on  $\boldsymbol{\Sigma}_u$ :  $p(\boldsymbol{\Sigma}_u) = \text{IW}(\boldsymbol{\Sigma}_u \mid \nu_0, \boldsymbol{\Sigma}_0)$

- Posterior full conditional of  $\boldsymbol{\Sigma}_u$ :

- Posterior full conditional of  $\mathbf{u}_i$ :

# Bayesian Ridge for Linear Mixed Models

$$y_{i,j} = \mathbf{x}_i^T \boldsymbol{\beta} + \mathbf{z}_i^T \mathbf{u}_i + \epsilon_{i,j}, \quad \mathbf{u}_i \stackrel{iid}{\sim} \text{MVN}(\mathbf{0}, \boldsymbol{\Sigma}_u), \quad \epsilon_{i,j} \stackrel{iid}{\sim} \text{Normal}(0, \sigma_\epsilon^2), \\ i = 1, \dots, n, \quad j = 1, \dots, m_i, \quad N = \sum_{i=1}^n m_i.$$

- Semi-conjugate priors on  $(\boldsymbol{\beta}, \sigma_\epsilon^2)$ :

$$p(\boldsymbol{\beta}, \sigma_\epsilon^2) = \text{MVN}(\boldsymbol{\mu}_\beta, \boldsymbol{\Sigma}_\beta) \cdot \text{Inv-Ga}(a_\sigma, b_\sigma) \quad \text{with} \quad \boldsymbol{\Sigma}_\beta = \lambda^{-1} \sigma_\epsilon^2 \mathbf{I}_p$$

- Posterior full conditionals of  $\boldsymbol{\beta}, \sigma_\epsilon^2$ :

- $p(\boldsymbol{\beta} \mid -) = \text{MVN}(\boldsymbol{\mu}_{\beta,N}, \boldsymbol{\Sigma}_{\beta,N}),$

$$\boldsymbol{\Sigma}_{\beta,N} = (\boldsymbol{\Sigma}_\beta^{-1} + \sigma_\epsilon^{-2} \mathbf{X}^T \mathbf{X})^{-1}, \quad \boldsymbol{\mu}_{\beta,N} = \boldsymbol{\Sigma}_{\beta,N} \left( \boldsymbol{\Sigma}_\beta^{-1} \boldsymbol{\mu}_\beta + \sigma_\epsilon^{-2} \mathbf{X}^T \tilde{\mathbf{y}} \right)$$

- $p(\sigma_\epsilon^2 \mid -) = \text{Inv-Ga} \left\{ a_\sigma + \frac{N}{2}, b_\sigma + \frac{1}{2} (\tilde{\mathbf{y}} - \mathbf{X}\boldsymbol{\beta})^T (\tilde{\mathbf{y}} - \mathbf{X}\boldsymbol{\beta}) + \frac{\lambda}{2} \boldsymbol{\beta}^T \boldsymbol{\beta} \right\}$

- Conjugate prior on  $\lambda$ :  $p(\lambda) = \text{Ga}(a_\lambda, b_\lambda)$

- Posterior full conditional of  $\lambda$ :

- Conjugate prior on  $\boldsymbol{\Sigma}_u$ :  $p(\boldsymbol{\Sigma}_u) = \text{IW}(\boldsymbol{\Sigma}_u \mid \nu_0, \boldsymbol{\Sigma}_0)$

- Posterior full conditional of  $\boldsymbol{\Sigma}_u$ :

- Posterior full conditional of  $\mathbf{u}_i$ :

# Bayesian Ridge for Linear Mixed Models

$$y_{i,j} = \mathbf{x}_i^T \boldsymbol{\beta} + \mathbf{z}_i^T \mathbf{u}_i + \epsilon_{i,j}, \quad \mathbf{u}_i \stackrel{iid}{\sim} \text{MVN}(\mathbf{0}, \boldsymbol{\Sigma}_u), \quad \epsilon_{i,j} \stackrel{iid}{\sim} \text{Normal}(0, \sigma_\epsilon^2), \\ i = 1, \dots, n, \quad j = 1, \dots, m_i, \quad N = \sum_{i=1}^n m_i.$$

- Semi-conjugate priors on  $(\boldsymbol{\beta}, \sigma_\epsilon^2)$ :

$$p(\boldsymbol{\beta}, \sigma_\epsilon^2) = \text{MVN}(\boldsymbol{\mu}_\beta, \boldsymbol{\Sigma}_\beta) \cdot \text{Inv-Ga}(a_\sigma, b_\sigma) \quad \text{with} \quad \boldsymbol{\Sigma}_\beta = \lambda^{-1} \sigma_\epsilon^2 \mathbf{I}_p$$

- Posterior full conditionals of  $\boldsymbol{\beta}, \sigma_\epsilon^2$ :

- $p(\boldsymbol{\beta} \mid -) = \text{MVN}(\boldsymbol{\mu}_{\beta,N}, \boldsymbol{\Sigma}_{\beta,N}),$

$$\boldsymbol{\Sigma}_{\beta,N} = (\boldsymbol{\Sigma}_\beta^{-1} + \sigma_\epsilon^{-2} \mathbf{X}^T \mathbf{X})^{-1}, \quad \boldsymbol{\mu}_{\beta,N} = \boldsymbol{\Sigma}_{\beta,N} \left( \boldsymbol{\Sigma}_\beta^{-1} \boldsymbol{\mu}_\beta + \sigma_\epsilon^{-2} \mathbf{X}^T \tilde{\mathbf{y}} \right)$$

- $p(\sigma_\epsilon^2 \mid -) = \text{Inv-Ga} \left\{ a_\sigma + \frac{N}{2}, b_\sigma + \frac{1}{2} (\tilde{\mathbf{y}} - \mathbf{X}\boldsymbol{\beta})^T (\tilde{\mathbf{y}} - \mathbf{X}\boldsymbol{\beta}) + \frac{\lambda}{2} \boldsymbol{\beta}^T \boldsymbol{\beta} \right\}$

- Conjugate prior on  $\lambda$ :  $p(\lambda) = \text{Ga}(a_\lambda, b_\lambda)$

- Posterior full conditional of  $\lambda$ :

- $p(\lambda \mid -) = \text{Ga} \left\{ a_\lambda + p/2, b_\lambda + \boldsymbol{\beta}^T \boldsymbol{\beta} / (2\sigma_\epsilon^2) \right\}$

- Conjugate prior on  $\boldsymbol{\Sigma}_u$ :  $p(\boldsymbol{\Sigma}_u) = \text{IW}(\boldsymbol{\Sigma}_u \mid \nu_0, \boldsymbol{\Sigma}_0)$

- Posterior full conditional of  $\boldsymbol{\Sigma}_u$ :

- Posterior full conditional of  $\mathbf{u}_i$ :



# Bayesian Ridge for Linear Mixed Models

$$y_{i,j} = \mathbf{x}_i^T \boldsymbol{\beta} + \mathbf{z}_i^T \mathbf{u}_i + \epsilon_{i,j}, \quad \mathbf{u}_i \stackrel{iid}{\sim} \text{MVN}(\mathbf{0}, \boldsymbol{\Sigma}_u), \quad \epsilon_{i,j} \stackrel{iid}{\sim} \text{Normal}(0, \sigma_\epsilon^2), \\ i = 1, \dots, n, \quad j = 1, \dots, m_i, \quad N = \sum_{i=1}^n m_i.$$

- Semi-conjugate priors on  $(\boldsymbol{\beta}, \sigma_\epsilon^2)$ :

$$p(\boldsymbol{\beta}, \sigma_\epsilon^2) = \text{MVN}(\boldsymbol{\mu}_\beta, \boldsymbol{\Sigma}_\beta) \cdot \text{Inv-Ga}(a_\sigma, b_\sigma) \quad \text{with} \quad \boldsymbol{\Sigma}_\beta = \lambda^{-1} \sigma_\epsilon^2 \mathbf{I}_p$$

- Posterior full conditionals of  $\boldsymbol{\beta}, \sigma_\epsilon^2$ :

- $p(\boldsymbol{\beta} \mid -) = \text{MVN}(\boldsymbol{\mu}_{\beta,N}, \boldsymbol{\Sigma}_{\beta,N}),$

$$\boldsymbol{\Sigma}_{\beta,N} = (\boldsymbol{\Sigma}_\beta^{-1} + \sigma_\epsilon^{-2} \mathbf{X}^T \mathbf{X})^{-1}, \quad \boldsymbol{\mu}_{\beta,N} = \boldsymbol{\Sigma}_{\beta,N} \left( \boldsymbol{\Sigma}_\beta^{-1} \boldsymbol{\mu}_\beta + \sigma_\epsilon^{-2} \mathbf{X}^T \tilde{\mathbf{y}} \right)$$

- $p(\sigma_\epsilon^2 \mid -) = \text{Inv-Ga} \left\{ a_\sigma + \frac{N}{2}, b_\sigma + \frac{1}{2} (\tilde{\mathbf{y}} - \mathbf{X}\boldsymbol{\beta})^T (\tilde{\mathbf{y}} - \mathbf{X}\boldsymbol{\beta}) + \frac{\lambda}{2} \boldsymbol{\beta}^T \boldsymbol{\beta} \right\}$

- Conjugate prior on  $\lambda$ :  $p(\lambda) = \text{Ga}(a_\lambda, b_\lambda)$

- Posterior full conditional of  $\lambda$ :

- $p(\lambda \mid -) = \text{Ga} \left\{ a_\lambda + p/2, b_\lambda + \boldsymbol{\beta}^T \boldsymbol{\beta} / (2\sigma_\epsilon^2) \right\}$

- Conjugate prior on  $\boldsymbol{\Sigma}_u$ :  $p(\boldsymbol{\Sigma}_u) = \text{IW}(\boldsymbol{\Sigma}_u \mid \nu_0, \boldsymbol{\Sigma}_0)$

- Posterior full conditional of  $\boldsymbol{\Sigma}_u$ :

- $p(\boldsymbol{\Sigma}_u \mid -) = \text{IW}(n + \nu_0, \boldsymbol{\Sigma}_0 + \sum_{i=1}^n \mathbf{u}_i \mathbf{u}_i^T)$

- Posterior full conditional of  $\mathbf{u}_i$ :

- $p(\mathbf{u}_i \mid -) = \text{MVN}(\boldsymbol{\mu}_{i,u,N}, \boldsymbol{\Sigma}_{i,u,N}),$

$$\boldsymbol{\Sigma}_{i,u,N} = (\boldsymbol{\Sigma}_u^{-1} + \sigma_\epsilon^{-2} m_i \mathbf{z}_i \mathbf{z}_i^T)^{-1}, \quad \boldsymbol{\mu}_{i,u,N} = \boldsymbol{\Sigma}_{i,u,N} \left\{ \sigma_\epsilon^{-2} \sum_{j=1}^{m_i} (y_{i,j} - \mathbf{x}_i^T \boldsymbol{\beta}) \mathbf{z}_i \right\}$$

$$y_{i,j} = \mathbf{x}_i^T \boldsymbol{\beta} + \mathbf{z}_i^T \mathbf{u}_i + \epsilon_{i,j}, \quad \mathbf{u}_i \stackrel{iid}{\sim} \text{MVN}(\mathbf{0}, \boldsymbol{\Sigma}_u), \quad \epsilon_{i,j} \stackrel{iid}{\sim} \text{Normal}(0, \sigma_\epsilon^2), \\ i = 1, \dots, n, \quad j = 1, \dots, m_i, \quad N = \sum_{i=1}^n m_i.$$

- Semi-conjugate LASSO priors on  $\boldsymbol{\beta}$ :

$$p(\boldsymbol{\beta} \mid \lambda, \sigma_\epsilon^2) = \prod_{j=1}^p \int_0^\infty \text{Normal}(\beta_j \mid 0, \tau_j^2) \text{Exp} \left( \tau_j^2 \mid \frac{\lambda^2}{2\sigma_\epsilon^2} \right) d\tau_j^2$$

- Posterior full conditionals for block-Gibbs sampler:

- Conjugate prior on  $\lambda^2$ :  $p(\lambda^2) = \text{Ga}(a_\lambda, b_\lambda)$

- Posterior full conditional of  $\lambda^2$ :

- Posterior full conditional of  $\boldsymbol{\Sigma}_u$ :

- Posterior full conditional of  $\mathbf{u}_i$ :

# Bayesian LASSO for Linear Mixed Models

$$y_{i,j} = \mathbf{x}_i^T \boldsymbol{\beta} + \mathbf{z}_i^T \mathbf{u}_i + \epsilon_{i,j}, \quad \mathbf{u}_i \stackrel{iid}{\sim} \text{MVN}(\mathbf{0}, \boldsymbol{\Sigma}_u), \quad \epsilon_{i,j} \stackrel{iid}{\sim} \text{Normal}(0, \sigma_\epsilon^2), \\ i = 1, \dots, n, \quad j = 1, \dots, m_i, \quad N = \sum_{i=1}^n m_i.$$

- Semi-conjugate LASSO priors on  $\boldsymbol{\beta}$ :

$$p(\boldsymbol{\beta} \mid \lambda, \sigma_\epsilon^2) = \prod_{j=1}^p \int_0^\infty \text{Normal}(\beta_j \mid 0, \tau_j^2) \text{Exp}\left(\tau_j^2 \mid \frac{\lambda^2}{2\sigma_\epsilon^2}\right) d\tau_j^2$$

- Posterior full conditionals for block-Gibbs sampler:

- $p(\boldsymbol{\beta} \mid -) = \text{MVN}(\boldsymbol{\mu}_{\boldsymbol{\beta}, N}, \boldsymbol{\Sigma}_{\boldsymbol{\beta}, N}), \quad \mathbf{D}_\tau = \text{diag}(\tau_1^2, \dots, \tau_p^2),$

$$\boldsymbol{\Sigma}_{\boldsymbol{\beta}, N} = \sigma_\epsilon^2 (\mathbf{X}^T \mathbf{X} + \mathbf{D}_\tau^{-1})^{-1}, \quad \boldsymbol{\mu}_{\boldsymbol{\beta}, N} = \boldsymbol{\Sigma}_{\boldsymbol{\beta}, N} (\sigma_\epsilon^{-2} \mathbf{X}^T \tilde{\mathbf{y}}) = (\mathbf{X}^T \mathbf{X} + \mathbf{D}_\tau^{-1})^{-1} \mathbf{X}^T \tilde{\mathbf{y}}$$

- $p(\sigma_\epsilon^2 \mid -) = \text{Inv-Ga} \left\{ a_\sigma + \frac{n+p}{2}, b_\sigma + \frac{(\tilde{\mathbf{y}} - \mathbf{X}\boldsymbol{\beta})^T (\tilde{\mathbf{y}} - \mathbf{X}\boldsymbol{\beta})}{2} + \frac{\boldsymbol{\beta}^T \mathbf{D}_\tau^{-1} \boldsymbol{\beta}}{2} \right\}$
- $\tau_j^2 \rightarrow w_j = \tau_j^{-2}, p(w_j \mid -) = \text{Inv-Gs}(\mu', \lambda'), \quad \mu' = \frac{\lambda}{\sigma_\epsilon |\beta_j|}, \quad \lambda' = \frac{\lambda^2}{\sigma_\epsilon^2}$

- Conjugate prior on  $\lambda^2$ :  $p(\lambda^2) = \text{Ga}(a_\lambda, b_\lambda)$

- Posterior full conditional of  $\lambda^2$ :

$$p(\lambda^2 \mid -) \propto \lambda^{a_\lambda - 1} \exp\left(-\frac{b_\lambda + \sum_{j=1}^p w_j}{\lambda^2}\right)$$

- Posterior full conditional of  $\boldsymbol{\Sigma}_u$ :

$$p(\boldsymbol{\Sigma}_u \mid -) \propto |\boldsymbol{\Sigma}_u|^{-\frac{n}{2}} \exp\left(-\frac{1}{2} \sum_{i=1}^n \mathbf{z}_i^T \boldsymbol{\Sigma}_u^{-1} \mathbf{z}_i\right)$$

- Posterior full conditional of  $\mathbf{u}_i$ :

$$p(\mathbf{u}_i \mid -) \propto \exp\left(-\frac{1}{2} \mathbf{z}_i^T \boldsymbol{\Sigma}_u^{-1} \mathbf{z}_i - \frac{1}{2} \mathbf{u}_i^T \boldsymbol{\Sigma}_u^{-1} \mathbf{u}_i\right)$$

$$y_{i,j} = \mathbf{x}_i^T \boldsymbol{\beta} + \mathbf{z}_i^T \mathbf{u}_i + \epsilon_{i,j}, \quad \mathbf{u}_i \stackrel{iid}{\sim} \text{MVN}(\mathbf{0}, \boldsymbol{\Sigma}_u), \quad \epsilon_{i,j} \stackrel{iid}{\sim} \text{Normal}(0, \sigma_\epsilon^2), \\ i = 1, \dots, n, \quad j = 1, \dots, m_i, \quad N = \sum_{i=1}^n m_i.$$

- Semi-conjugate LASSO priors on  $\boldsymbol{\beta}$ :

$$p(\boldsymbol{\beta} \mid \lambda, \sigma_\epsilon^2) = \prod_{j=1}^p \int_0^\infty \text{Normal}(\beta_j \mid 0, \tau_j^2) \text{Exp}\left(\tau_j^2 \mid \frac{\lambda^2}{2\sigma_\epsilon^2}\right) d\tau_j^2$$

- Posterior full conditionals for block-Gibbs sampler:

$$\bullet \quad p(\boldsymbol{\beta} \mid -) = \text{MVN}(\boldsymbol{\mu}_{\boldsymbol{\beta}, N}, \boldsymbol{\Sigma}_{\boldsymbol{\beta}, N}), \quad \mathbf{D}_\tau = \text{diag}(\tau_1^2, \dots, \tau_p^2),$$

$$\boldsymbol{\Sigma}_{\boldsymbol{\beta}, N} = \sigma_\epsilon^2 (\mathbf{X}^T \mathbf{X} + \mathbf{D}_\tau^{-1})^{-1}, \quad \boldsymbol{\mu}_{\boldsymbol{\beta}, N} = \boldsymbol{\Sigma}_{\boldsymbol{\beta}, N} (\sigma_\epsilon^{-2} \mathbf{X}^T \tilde{\mathbf{y}}) = (\mathbf{X}^T \mathbf{X} + \mathbf{D}_\tau^{-1})^{-1} \mathbf{X}^T \tilde{\mathbf{y}}$$

$$\bullet \quad p(\sigma_\epsilon^2 \mid -) = \text{Inv-Ga} \left\{ a_\sigma + \frac{n+p}{2}, b_\sigma + \frac{(\tilde{\mathbf{y}} - \mathbf{X}\boldsymbol{\beta})^T (\tilde{\mathbf{y}} - \mathbf{X}\boldsymbol{\beta})}{2} + \frac{\boldsymbol{\beta}^T \mathbf{D}_\tau^{-1} \boldsymbol{\beta}}{2} \right\}$$

$$\bullet \quad \tau_j^2 \rightarrow w_j = \tau_j^{-2}, \quad p(w_j \mid -) = \text{Inv-Gs}(\mu', \lambda'), \quad \mu' = \frac{\lambda}{\sigma_\epsilon |\beta_j|}, \quad \lambda' = \frac{\lambda^2}{\sigma_\epsilon^2}$$

- Conjugate prior on  $\lambda^2$ :  $p(\lambda^2) = \text{Ga}(a_\lambda, b_\lambda)$

- Posterior full conditional of  $\lambda^2$ :

$$p(\lambda^2 \mid -) = \text{Ga} \left( a_\lambda + \sum_{j=1}^p \frac{1}{w_j}, b_\lambda + \sum_{j=1}^p \frac{1}{w_j} \right)$$

- Posterior full conditional of  $\boldsymbol{\Sigma}_u$ :

$$p(\boldsymbol{\Sigma}_u \mid -) = \text{Inv-Wish} \left( \nu, \boldsymbol{\Sigma}_u^{-1} + \sum_{i=1}^n \mathbf{z}_i \mathbf{z}_i^T \right)$$

- Posterior full conditional of  $\mathbf{u}_i$ :

$$p(\mathbf{u}_i \mid -) = \text{MVN}(\boldsymbol{\mu}_i, \boldsymbol{\Sigma}_i)$$

# Bayesian LASSO for Linear Mixed Models

$$y_{i,j} = \mathbf{x}_i^T \boldsymbol{\beta} + \mathbf{z}_i^T \mathbf{u}_i + \epsilon_{i,j}, \quad \mathbf{u}_i \stackrel{iid}{\sim} \text{MVN}(\mathbf{0}, \boldsymbol{\Sigma}_u), \quad \epsilon_{i,j} \stackrel{iid}{\sim} \text{Normal}(0, \sigma_\epsilon^2),$$

$$i = 1, \dots, n, \quad j = 1, \dots, m_i, \quad N = \sum_{i=1}^n m_i.$$

- Semi-conjugate LASSO priors on  $\boldsymbol{\beta}$ :

$$p(\boldsymbol{\beta} \mid \lambda, \sigma_\epsilon^2) = \prod_{j=1}^p \int_0^\infty \text{Normal}(\beta_j \mid 0, \tau_j^2) \text{Exp} \left( \tau_j^2 \mid \frac{\lambda^2}{2\sigma_\epsilon^2} \right) d\tau_j^2$$

- Posterior full conditionals for block-Gibbs sampler:

- $p(\boldsymbol{\beta} \mid -) = \text{MVN}(\boldsymbol{\mu}_{\boldsymbol{\beta}, N}, \boldsymbol{\Sigma}_{\boldsymbol{\beta}, N}), \quad \mathbf{D}_\tau = \text{diag}(\tau_1^2, \dots, \tau_p^2),$

$$\boldsymbol{\Sigma}_{\boldsymbol{\beta}, N} = \sigma_\epsilon^2 (\mathbf{X}^T \mathbf{X} + \mathbf{D}_\tau^{-1})^{-1}, \quad \boldsymbol{\mu}_{\boldsymbol{\beta}, N} = \boldsymbol{\Sigma}_{\boldsymbol{\beta}, N} (\sigma_\epsilon^{-2} \mathbf{X}^T \tilde{\mathbf{y}}) = (\mathbf{X}^T \mathbf{X} + \mathbf{D}_\tau^{-1})^{-1} \mathbf{X}^T \tilde{\mathbf{y}}$$

- $p(\sigma_\epsilon^2 \mid -) = \text{Inv-Ga} \left\{ a_\sigma + \frac{n+p}{2}, b_\sigma + \frac{(\tilde{\mathbf{y}} - \mathbf{X}\boldsymbol{\beta})^T (\tilde{\mathbf{y}} - \mathbf{X}\boldsymbol{\beta})}{2} + \frac{\boldsymbol{\beta}^T \mathbf{D}_\tau^{-1} \boldsymbol{\beta}}{2} \right\}$
- $\tau_j^2 \rightarrow w_j = \tau_j^{-2}, p(w_j \mid -) = \text{Inv-Gs}(\mu', \lambda'), \quad \mu' = \frac{\lambda}{\sigma_\epsilon |\beta_j|}, \quad \lambda' = \frac{\lambda^2}{\sigma_\epsilon^2}$

- Conjugate prior on  $\lambda^2$ :  $p(\lambda^2) = \text{Ga}(a_\lambda, b_\lambda)$

- Posterior full conditional of  $\lambda^2$ :

- $p(\lambda^2 \mid -) = \text{Ga} \left\{ a_\lambda + p, b_\lambda + \sum_{j=1}^p \tau_j^2 / (2\sigma_\epsilon^2) \right\}.$

- Posterior full conditional of  $\boldsymbol{\Sigma}_u$ :

- $p(\boldsymbol{\Sigma}_u \mid -) = \text{IW}(n + \nu_0, \boldsymbol{\Sigma}_0 + \sum_{i=1}^n \mathbf{u}_i \mathbf{u}_i^T)$

- Posterior full conditional of  $\mathbf{u}_i$ :

- $p(\mathbf{u}_i \mid -) = \text{MVN}(\boldsymbol{\mu}_{i,u,N}, \boldsymbol{\Sigma}_{i,u,N}),$

$$\boldsymbol{\Sigma}_{i,u,N} = (\boldsymbol{\Sigma}_u^{-1} + \sigma_\epsilon^{-2} m_i \mathbf{z}_i \mathbf{z}_i^T)^{-1}, \quad \boldsymbol{\mu}_{i,u,N} = \boldsymbol{\Sigma}_{i,u,N} \left\{ \sigma_\epsilon^{-2} \sum_{j=1}^{m_i} (y_{i,j} - \mathbf{x}_i^T \boldsymbol{\beta}) \mathbf{z}_i \right\}$$

- Linear models are regression models linear in parameters.
  - Ordinary least squares minimizes the  $L_2$  distance between the observed and the model hypothesized response values.
  - MLE under a normal likelihood is naturally connected to OLS.
  - Ridge regression is useful when the model matrix is ill-conditioned.
  - LASSO is useful when the model is sparse.
  - Bayesian MCMC based inference is straightforward under (semi)conjugate priors.
  - Ridge estimates can be obtained as Bayesian MAP under a type of independent normal priors on the regression coefficients.
  - LASSO estimates can be obtained as Bayesian MAP under a type of independent Laplace priors on the regression coefficients.
  - Mixed models are useful for including population level fixed effects as well as individual level random effects.
  - Mixed models are generally computation intensive but can usually be relatively easily handled using Bayesian hierarchies.

- Linear models are regression models linear in parameters.
  - Ordinary least squares minimizes the  $L_2$  distance between the observed and the model hypothesized response values.
  - MLE under a normal likelihood is naturally connected to OLS.
  - Ridge regression is useful when the model matrix is ill-conditioned.
  - LASSO is useful when the model is sparse.
  - Bayesian MCMC based inference is straightforward under (semi)conjugate priors.
  - Ridge estimates can be obtained as Bayesian MAP under a type of independent normal priors on the regression coefficients.
  - LASSO estimates can be obtained as Bayesian MAP under a type of independent Laplace priors on the regression coefficients.
  - Mixed models are useful for including population level fixed effects as well as individual level random effects.
  - Mixed models are generally computation intensive but can usually be relatively easily handled using Bayesian hierarchies.

- Linear models are regression models linear in parameters.
  - Ordinary least squares minimizes the  $L_2$  distance between the observed and the model hypothesized response values.
  - MLE under a normal likelihood is naturally connected to OLS.
  - Ridge regression is useful when the model matrix is ill-conditioned.
  - LASSO is useful when the model is sparse.
  - Bayesian MCMC based inference is straightforward under (semi)conjugate priors.
  - Ridge estimates can be obtained as Bayesian MAP under a type of independent normal priors on the regression coefficients.
  - LASSO estimates can be obtained as Bayesian MAP under a type of independent Laplace priors on the regression coefficients.
  - Mixed models are useful for including population level fixed effects as well as individual level random effects.
  - Mixed models are generally computation intensive but can usually be relatively easily handled using Bayesian hierarchies.



- Linear models are regression models linear in parameters.
  - Ordinary least squares minimizes the  $L_2$  distance between the observed and the model hypothesized response values.
  - MLE under a normal likelihood is naturally connected to OLS.
  - Ridge regression is useful when the model matrix is ill-conditioned.
  - LASSO is useful when the model is sparse.
  - Bayesian MCMC based inference is straightforward under (semi)conjugate priors.
  - Ridge estimates can be obtained as Bayesian MAP under a type of independent normal priors on the regression coefficients.
  - LASSO estimates can be obtained as Bayesian MAP under a type of independent Laplace priors on the regression coefficients.
  - Mixed models are useful for including population level fixed effects as well as individual level random effects.
  - Mixed models are generally computation intensive but can usually be relatively easily handled using Bayesian hierarchies.

- Linear models are regression models linear in parameters.
  - Ordinary least squares minimizes the  $L_2$  distance between the observed and the model hypothesized response values.
  - MLE under a normal likelihood is naturally connected to OLS.
  - Ridge regression is useful when the model matrix is ill-conditioned.
  - **LASSO is useful when the model is sparse.**
  - Bayesian MCMC based inference is straightforward under (semi)conjugate priors.
  - Ridge estimates can be obtained as Bayesian MAP under a type of independent normal priors on the regression coefficients.
  - LASSO estimates can be obtained as Bayesian MAP under a type of independent Laplace priors on the regression coefficients.
  - Mixed models are useful for including population level fixed effects as well as individual level random effects.
  - Mixed models are generally computation intensive but can usually be relatively easily handled using Bayesian hierarchies.

- Linear models are regression models linear in parameters.
  - Ordinary least squares minimizes the  $L_2$  distance between the observed and the model hypothesized response values.
  - MLE under a normal likelihood is naturally connected to OLS.
  - Ridge regression is useful when the model matrix is ill-conditioned.
  - LASSO is useful when the model is sparse.
  - Bayesian MCMC based inference is straightforward under (semi)conjugate priors.
  - Ridge estimates can be obtained as Bayesian MAP under a type of independent normal priors on the regression coefficients.
  - LASSO estimates can be obtained as Bayesian MAP under a type of independent Laplace priors on the regression coefficients.
  - Mixed models are useful for including population level fixed effects as well as individual level random effects.
  - Mixed models are generally computation intensive but can usually be relatively easily handled using Bayesian hierarchies.

- Linear models are regression models linear in parameters.
  - Ordinary least squares minimizes the  $L_2$  distance between the observed and the model hypothesized response values.
  - MLE under a normal likelihood is naturally connected to OLS.
  - Ridge regression is useful when the model matrix is ill-conditioned.
  - LASSO is useful when the model is sparse.
  - Bayesian MCMC based inference is straightforward under (semi)conjugate priors.
  - Ridge estimates can be obtained as Bayesian MAP under a type of independent normal priors on the regression coefficients.
  - LASSO estimates can be obtained as Bayesian MAP under a type of independent Laplace priors on the regression coefficients.
  - Mixed models are useful for including population level fixed effects as well as individual level random effects.
  - Mixed models are generally computation intensive but can usually be relatively easily handled using Bayesian hierarchies.

- Linear models are regression models linear in parameters.
  - Ordinary least squares minimizes the  $L_2$  distance between the observed and the model hypothesized response values.
  - MLE under a normal likelihood is naturally connected to OLS.
  - Ridge regression is useful when the model matrix is ill-conditioned.
  - LASSO is useful when the model is sparse.
  - Bayesian MCMC based inference is straightforward under (semi)conjugate priors.
  - Ridge estimates can be obtained as Bayesian MAP under a type of independent normal priors on the regression coefficients.
  - LASSO estimates can be obtained as Bayesian MAP under a type of independent Laplace priors on the regression coefficients.
  - Mixed models are useful for including population level fixed effects as well as individual level random effects.
  - Mixed models are generally computation intensive but can usually be relatively easily handled using Bayesian hierarchies.

- Linear models are regression models linear in parameters.
  - Ordinary least squares minimizes the  $L_2$  distance between the observed and the model hypothesized response values.
  - MLE under a normal likelihood is naturally connected to OLS.
  - Ridge regression is useful when the model matrix is ill-conditioned.
  - LASSO is useful when the model is sparse.
  - Bayesian MCMC based inference is straightforward under (semi)conjugate priors.
  - Ridge estimates can be obtained as Bayesian MAP under a type of independent normal priors on the regression coefficients.
  - LASSO estimates can be obtained as Bayesian MAP under a type of independent Laplace priors on the regression coefficients.
  - Mixed models are useful for including population level fixed effects as well as individual level random effects.
  - Mixed models are generally computation intensive but can usually be relatively easily handled using Bayesian hierarchies.

- Linear models are regression models linear in parameters.
  - Ordinary least squares minimizes the  $L_2$  distance between the observed and the model hypothesized response values.
  - MLE under a normal likelihood is naturally connected to OLS.
  - Ridge regression is useful when the model matrix is ill-conditioned.
  - LASSO is useful when the model is sparse.
  - Bayesian MCMC based inference is straightforward under (semi)conjugate priors.
  - Ridge estimates can be obtained as Bayesian MAP under a type of independent normal priors on the regression coefficients.
  - LASSO estimates can be obtained as Bayesian MAP under a type of independent Laplace priors on the regression coefficients.
  - Mixed models are useful for including population level fixed effects as well as individual level random effects.
  - Mixed models are generally computation intensive but can usually be relatively easily handled using Bayesian hierarchies.