

# Bayesian linear regression

Bayesian Methodology in Biostatistics (BST 249)

Jeffrey W. Miller

Department of Biostatistics  
Harvard T.H. Chan School of Public Health

# Outline

## Model, prior, and posterior

Conditionally conjugate prior on  $\beta$

Conditionally conjugate prior on  $\sigma^2$

## Setting the hyperparameters

Unit information prior

Zellner's  $g$ -prior

Computation with Zellner's  $g$ -prior

# Outline

## Model, prior, and posterior

Conditionally conjugate prior on  $\beta$

Conditionally conjugate prior on  $\sigma^2$

## Setting the hyperparameters

Unit information prior

Zellner's  $g$ -prior

Computation with Zellner's  $g$ -prior

## Linear regression model

- Data:  $(x_1, y_1), \dots, (x_n, y_n)$ , where  $x_i \in \mathbb{R}^p$  and  $y_i \in \mathbb{R}$ .
  - ▶  $x_i = (x_{i1}, \dots, x_{ip})^T =$  vector of covariates
  - ▶  $y_i =$  outcome
- Model:  $Y_i \sim \mathcal{N}(x_i^T \beta, \sigma^2)$  independently for  $i = 1, \dots, n$ .
- Throughout, we will treat  $x_{1:n}$  as fixed and known.
- Equivalently,  $Y \sim \mathcal{N}(X\beta, \sigma^2 I)$  where

$$Y = \begin{bmatrix} Y_1 \\ \vdots \\ Y_n \end{bmatrix} \in \mathbb{R}^n \quad X = \begin{bmatrix} x_1^T \\ \vdots \\ x_n^T \end{bmatrix} \in \mathbb{R}^{n \times p}$$

and  $I$  is the  $p \times p$  identity matrix.

## Conditionally conjugate prior on $\beta$

- Assume a multivariate normal prior on  $\beta$ :

$$\beta \sim \mathcal{N}(m_0, L_0^{-1}).$$

- For any fixed  $\sigma^2$ , this is a conjugate prior.
- The resulting posterior, with  $\sigma^2$  fixed, is:

$$\beta|y_{1:n} \sim \mathcal{N}(m_n, L_n^{-1})$$

where

$$\begin{aligned}L_n &= L_0 + X^T X / \sigma^2, \\m_n &= L_n^{-1} (L_0 m_0 + X^T y / \sigma^2),\end{aligned}$$

and  $y = (y_1, \dots, y_n)^T \in \mathbb{R}^n$ .

## Connections with the MLE and ridge regression

- The *maximum a posteriori* (MAP) estimate is the mode of the posterior.
- Since the posterior is normal, the mode is equal to the mean:

$$\hat{\beta}_{\text{MAP}} = m_n = (L_0 + X^T X / \sigma^2)^{-1} (L_0 m_0 + X^T y / \sigma^2).$$

- In the limit as  $L_0 \rightarrow 0$ , we have  $L_n \rightarrow X^T X / \sigma^2$  and

$$\hat{\beta}_{\text{MAP}} = m_n \rightarrow (X^T X)^{-1} X^T y = \hat{\beta}_{\text{MLE}}.$$

- Thus, in the limit of having “no prior information” about  $\beta$ ,

$$\beta | y_{1:n} \sim \mathcal{N}(\hat{\beta}_{\text{MLE}}, \sigma^2 (X^T X)^{-1}).$$

- Note that the frequentist sampling distribution of the MLE is

$$\hat{\beta}_{\text{MLE}} \sim \mathcal{N}(\beta, \sigma^2 (X^T X)^{-1})$$

when  $\beta$  is the true parameter, which exhibits perfect symmetry with the posterior above when we take  $L_0 \rightarrow 0$ .

## Connections with the MLE and ridge regression

- Taking  $L_0 \rightarrow 0$  is equivalent to using a *flat prior*  $p(\beta) \propto 1$  and formally applying Bayes' theorem:  $p(\beta|y_{1:n}) \propto p(y_{1:n}|\beta)$ .
- $p(\beta) \propto 1$  is called an *improper* prior since it cannot be normalized to a probability density on  $\mathbb{R}^p$ .
- Even though the flat prior on  $\beta$  is improper, defining  $p(\beta|y_{1:n}) \propto p(y_{1:n}|\beta)$  still results in a well-defined posterior distribution as long as  $n \geq 1$ .
  
- This is an example of a *noninformative* prior.
  
- The MAP estimate also generalizes ridge regression.
- Specifically, if  $m_0 = 0$  and  $L_0 = \alpha I / \sigma^2$  then

$$\hat{\beta}_{\text{MAP}} = (X^T X + \alpha I)^{-1} X^T y = \hat{\beta}_{\text{ridge}}.$$

# Outline

## Model, prior, and posterior

Conditionally conjugate prior on  $\beta$

Conditionally conjugate prior on  $\sigma^2$

## Setting the hyperparameters

Unit information prior

Zellner's  $g$ -prior

Computation with Zellner's  $g$ -prior



## Conjugate prior on $\sigma^2$

- Similar to the case of a univariate normal model, the Inverse Gamma distribution is a conjugate prior for  $\sigma^2$ .
- Suppose  $\beta$  is fixed, and define the prior on  $\sigma^2$  as:

$$\sigma^2 \sim \text{InvGamma}\left(\frac{1}{2}\nu_0, \frac{1}{2}\nu_0\sigma_0^2\right).$$

- This parametrization facilitates interpretation:
  - ▶  $\sigma_0^2 =$  prior guess of  $\sigma^2$  (since  $\sigma_0$  is between the mean and the mode of this prior on  $\sigma^2$ )
  - ▶  $\nu_0 =$  confidence in the choice of  $\sigma_0^2$ , in units of sample size — i.e., the strength of the prior is equivalent to  $\nu_0$  samples.
- The resulting posterior on  $\sigma^2$  is then

$$\sigma^2|y_{1:n} \sim ???$$

(Whiteboard activity)

## Conjugate prior on $\sigma^2$

- Similar to the case of a univariate normal model, the Inverse Gamma distribution is a conjugate prior for  $\sigma^2$ .
- Suppose  $\beta$  is fixed, and define the prior on  $\sigma^2$  as:

$$\sigma^2 \sim \text{InvGamma}\left(\frac{1}{2}\nu_0, \frac{1}{2}\nu_0\sigma_0^2\right).$$

- This parametrization facilitates interpretation:
  - ▶  $\sigma_0^2 =$  prior guess of  $\sigma^2$  (since  $\sigma_0$  is between the mean and the mode of this prior on  $\sigma^2$ )
  - ▶  $\nu_0 =$  confidence in the choice of  $\sigma_0^2$ , in units of sample size — i.e., the strength of the prior is equivalent to  $\nu_0$  samples.
- The resulting posterior on  $\sigma^2$  is then

$$\sigma^2|y_{1:n} \sim \text{InvGamma}\left(\frac{1}{2}(\nu_0 + n), \frac{1}{2}(\nu_0\sigma_0^2 + \text{SSR}(\beta))\right)$$

where  $\text{SSR}(\beta) = \sum_{i=1}^n (y_i - x_i^T\beta)^2$ .

## Inference for $\beta$ and $\sigma^2$ jointly

- We've seen conjugate priors for  $\beta$  and  $\sigma^2$ , given the other.
- Thus, these define a conditionally conjugate prior, combined:

$$\beta \sim \mathcal{N}(m_0, L_0^{-1}) \quad \sigma^2 \sim \text{InvGamma}(\frac{1}{2}\nu_0, \frac{1}{2}\nu_0\sigma_0^2)$$

independently.

- We know how to compute  $\beta|\sigma^2, y_{1:n}$  and  $\sigma^2|\beta, y_{1:n}$ . How can we infer  $\beta$  and  $\sigma^2$  jointly?
- One way: Initialize  $\beta, \sigma^2$  and iteratively repeat the following:
  1. update  $\beta$  by sampling from  $\beta|\sigma^2, y_{1:n}$
  2. update  $\sigma^2$  by sampling from  $\sigma^2|\beta, y_{1:n}$
- It turns out that this generates approximate samples from  $\beta, \sigma^2|y_{1:n}$ . This is an example of a *Gibbs sampler*, which is type of a Markov chain Monte Carlo algorithm.

## Group activity: Statistics trivia challenge!

Go to breakout rooms and work together to answer these questions:

<https://forms.gle/agm9P7PpLTWNUmDr5>

(Three people per room, randomly assigned. 5 minutes.)

# Outline

## Model, prior, and posterior

Conditionally conjugate prior on  $\beta$

Conditionally conjugate prior on  $\sigma^2$

## Setting the hyperparameters

Unit information prior

Zellner's  $g$ -prior

Computation with Zellner's  $g$ -prior

## Setting hyperparameters: Data-dependent, Unit info

- The subjective Bayesian approach is to set hyperparameters based solely on prior beliefs. However, sometimes we prefer to use *default priors* that don't require subjective input.
- Data-dependent priors are a useful way of creating default priors, even though, strictly speaking, they violate the principle of not using the data twice.

- Kass & Wasserman (1995) propose the following settings:

$$m_0 = \hat{\beta}_{\text{MLE}} \quad L_0 = X^T X / (n\sigma^2).$$

- Similarly, we can define a data-dependent prior on  $\sigma^2$ :

$$\sigma_0^2 = \hat{\sigma}_{\text{MLE}}^2 \quad \nu_0 = 1$$

where  $\hat{\sigma}_{\text{MLE}}^2$  is the maximum likelihood estimate of  $\sigma^2$ .

- These are both *unit information priors* — the strength of the prior is equivalent to one sample.

## Setting hyperparameters: Zellner's $g$ -prior

- A compelling property of  $\hat{\beta}_{\text{MLE}}$  is that the scale of the predictors doesn't matter, in the following sense:
  - ▶ Suppose one of the predictors is age in years.
  - ▶ If we change the units of age to months, then entry of  $\hat{\beta}_{\text{MLE}}$  corresponding to age is scaled accordingly by 1/12.
  - ▶ Consequently,  $X\hat{\beta}_{\text{MLE}}$  is invariant to the choice of units.
- Mathematically, changing the units means using the model  $Y \sim \mathcal{N}(\tilde{X}\beta, \sigma^2 I)$  with  $\tilde{X} = XH$  in place of  $X$ , where  $H$  is a diagonal matrix with  $H_{jj} > 0$  for all  $j$ . Observe that

$$\tilde{X}\tilde{\beta}_{\text{MLE}} = ??? = X\hat{\beta}_{\text{MLE}}.$$

(Whiteboard activity)

## Setting hyperparameters: Zellner's $g$ -prior

- A compelling property of  $\hat{\beta}_{\text{MLE}}$  is that the scale of the predictors doesn't matter, in the following sense:
  - ▶ Suppose one of the predictors is age in years.
  - ▶ If we change the units of age to months, then entry of  $\hat{\beta}_{\text{MLE}}$  corresponding to age is scaled accordingly by  $1/12$ .
  - ▶ Consequently,  $X\hat{\beta}_{\text{MLE}}$  is invariant to the choice of units.
- Mathematically, changing the units means using the model  $Y \sim \mathcal{N}(\tilde{X}\beta, \sigma^2 I)$  with  $\tilde{X} = XH$  in place of  $X$ , where  $H$  is a diagonal matrix with  $H_{jj} > 0$  for all  $j$ . Observe that

$$\begin{aligned}\tilde{X}\tilde{\beta}_{\text{MLE}} &= \tilde{X}(\tilde{X}^T\tilde{X})^{-1}\tilde{X}^T y \\ &= XH(H^T X^T XH)^{-1}H^T X^T y \\ &= XHH^{-1}(X^T X)^{-1}H^{-T}H^T X^T y \\ &= X(X^T X)^{-1}X^T y = X\hat{\beta}_{\text{MLE}}.\end{aligned}$$



## Setting hyperparameters: Zellner's $g$ -prior

- This is nice since, intuitively, the units shouldn't matter.
- For the Bayesian model, the invariance property we seek is:
  - ▶ Suppose  $\tilde{X} = XH$  for some invertible  $H \in \mathbb{R}^{p \times p}$ .
  - ▶ Let  $\beta$  be distributed according to the posterior using  $X$ .
  - ▶ Let  $\tilde{\beta}$  be distributed according to the posterior using  $\tilde{X}$ .
  - ▶ We would like  $X\beta$  and  $\tilde{X}\tilde{\beta}$  to have the same distribution.
- Unfortunately, for most choices of prior, the posterior doesn't have this invariance property.
- However, if we can make it so that  $\beta$  and  $H\tilde{\beta}$  have the same distribution, then this will work, since then  $X\beta$  has the same distribution as  $XH\tilde{\beta} = \tilde{X}\tilde{\beta}$ .

## Setting hyperparameters: Zellner's $g$ -prior

- This invariance property is satisfied by Zellner's  $g$ -prior:

$$\beta | \sigma^2 \sim \mathcal{N}(m_0, L_0^{-1}) \text{ where } m_0 = 0 \text{ and } L_0 = X^T X / (g\sigma^2).$$

- Here,  $g > 0$  is a free parameter.
- Given  $\sigma^2$ , the posterior on  $\beta$  under a  $g$ -prior simplifies to:

$$\beta | \sigma^2, y_{1:n} \sim \mathcal{N}(m_n, L_n^{-1})$$

where

$$L_n = L_0 + X^T X / \sigma^2 = \frac{g+1}{g\sigma^2} X^T X,$$

$$m_n = L_n^{-1} (L_0 m_0 + X^T y / \sigma^2) = \frac{g}{g+1} (X^T X)^{-1} X^T y.$$

## Setting hyperparameters: Zellner's $g$ -prior

- To check that the invariance property is satisfied, suppose:

$$\tilde{m}_0 = 0 \quad \tilde{L}_0 = (XH)^T(XH)/(g\sigma^2).$$

- Then the posterior is  $\tilde{\beta} \mid \sigma^2, y_{1:n} \sim \mathcal{N}(\tilde{m}_n, \tilde{L}_n^{-1})$  where

$$\tilde{L}_n = ???$$

$$\tilde{m}_n = ???.$$

- Since  $H\tilde{\beta} \mid \sigma^2, y_{1:n} \sim \mathcal{N}(???, ???)$ , the invariance property requires that

$$H\tilde{m}_n = m_n \quad \text{and} \quad H\tilde{L}_n^{-1}H^T = L_n^{-1}.$$

- **Group activity: Go to breakout rooms and work together to (1) fill in the blanks above and (2) check that the equations above in red hold. (15 minutes)**

# Outline

## Model, prior, and posterior

Conditionally conjugate prior on  $\beta$

Conditionally conjugate prior on  $\sigma^2$

## Setting the hyperparameters

Unit information prior

Zellner's  $g$ -prior

Computation with Zellner's  $g$ -prior

## Posterior computation with Zellner's $g$ -prior

- The  $g$ -prior also simplifies posterior inference for  $(\beta, \sigma^2)$ .
- Assume a  $g$ -prior on  $\beta|\sigma^2$  and an InvGamma prior on  $\sigma^2$ :

$$\begin{aligned}\beta|\sigma^2 &\sim \mathcal{N}(0, g\sigma^2(X^T X)^{-1}) \\ \sigma^2 &\sim \text{InvGamma}\left(\frac{1}{2}\nu_0, \frac{1}{2}\nu_0\sigma_0^2\right).\end{aligned}$$

- It turns out that we can generate i.i.d. samples from the posterior on  $(\beta, \sigma^2)$ , so we don't need to use MCMC.
- These samples can be used for Monte Carlo approximation of posterior expectations, e.g., for any integrable function  $h$ ,

$$\mathbb{E}(h(\beta, \sigma^2) | y_{1:n}) \approx \frac{1}{T} \sum_{t=1}^T h(\beta_t, \sigma_t^2).$$

- Monte Carlo has various advantages over MCMC: simplicity, efficiency, and reliable quantification of approximation error.

## Posterior computation with Zellner's $g$ -prior

- Under the prior on the previous slide, it can be shown that the posterior is (Hoff, 2009):

$$\sigma^2 \mid y_{1:n} \sim \text{InvGamma} \left( \frac{1}{2}(\nu_0 + n), \frac{1}{2}(\nu_0 \sigma_0^2 + \text{SSR}_g) \right)$$
$$\beta \mid \sigma^2, y_{1:n} \sim \mathcal{N} \left( \frac{g}{g+1} \hat{\beta}_{\text{MLE}}, \frac{g}{g+1} \sigma^2 (X^T X)^{-1} \right)$$

where

$$\text{SSR}_g = y^T y - \frac{g}{g+1} y^T X (X^T X)^{-1} X^T y.$$

- Thus, we can generate i.i.d. samples from the posterior by:
  1. sampling  $\sigma^2 \mid y_{1:n}$ ,
  2. sampling  $\beta \mid \sigma^2, y_{1:n}$ .

## References and supplements

- Hoff, P. D. (2009). A First Course in Bayesian Statistical Methods. Springer, New York.
- Kass, R. E., & Wasserman, L. (1995). A reference Bayesian test for nested hypotheses and its relationship to the Schwarz criterion. *Journal of the American Statistical Association*, 90(431), 928-934.
- Zellner, A. (1986). On assessing prior distributions and Bayesian regression analysis with g-prior distributions. In: *Bayesian inference and Decision Techniques*, vol 6, North-Holland, Amsterdam, 233-243.

## Setting hyperparameters: Zellner's $g$ -prior (SOLUTIONS)

- To check that the invariance property is satisfied, suppose:

$$\tilde{m}_0 = 0 \quad \tilde{L}_0 = (XH)^T(XH)/(g\sigma^2).$$

- Then the posterior is  $\tilde{\beta} \mid \sigma^2, y_{1:n} \sim \mathcal{N}(\tilde{m}_n, \tilde{L}_n^{-1})$  where

$$\begin{aligned}\tilde{L}_n &= \frac{g+1}{g\sigma^2} H^T X^T X H \\ \tilde{m}_n &= \frac{g}{g+1} (H^T X^T X H)^{-1} H^T X^T y.\end{aligned}$$

- Since  $H\tilde{\beta} \mid \sigma^2, y_{1:n} \sim \mathcal{N}(H\tilde{m}_n, H\tilde{L}_n^{-1}H^T)$ , the invariance property requires that

$$H\tilde{m}_n = m_n \quad \text{and} \quad H\tilde{L}_n^{-1}H^T = L_n^{-1}.$$



## Setting hyperparameters: Zellner's $g$ -prior (SOLUTIONS)

These equations hold since:

$$\begin{aligned} H\tilde{m}_n &= \frac{g}{g+1} H(H^T X^T X H)^{-1} H^T X^T y \\ &= \frac{g}{g+1} H H^{-1} (X^T X)^{-1} H^{-T} H^T X^T y \\ &= \frac{g}{g+1} (X^T X)^{-1} X^T y = m_n \end{aligned}$$

and

$$\begin{aligned} H\tilde{L}_n^{-1}H^T &= H\left(\frac{g+1}{g\sigma^2} H^T X^T X H\right)^{-1} H^T \\ &= H H^{-1} \left(\frac{g+1}{g\sigma^2} X^T X\right)^{-1} H^{-T} H^T \\ &= \left(\frac{g+1}{g\sigma^2} X^T X\right)^{-1} = L_n^{-1}. \end{aligned}$$