

# Bayesian Linear Regression

Yuki Shiraito

POLSCI 798 Advanced Topics in Quantitative Methodology  
University of Michigan

Winter 2019

# Linear Regression with Gaussian Errors

- **Regression:** Predicting conditional expectation of  $Y_i$  given  $X_i$
- Linear regression:  $\mathbb{E}[Y_i | X_i] = X_i^\top \beta$
- **Ordinary linear regression:**

$$Y_i = X_i^\top \beta + \varepsilon_i, \quad \varepsilon_i \stackrel{\text{i.i.d.}}{\sim} \mathcal{N}(0, \sigma^2), \quad i = 1, \dots, N$$

$$\Leftrightarrow Y | \beta, \sigma, \mathbf{X} \sim \mathcal{N}(\mathbf{X}\beta, \sigma^2 \mathbf{I}_N)$$

- Unknown parameters:  $\beta$  and  $\sigma$
- $\mathbf{X}$  may be data, but inference is conditional on it
- Likelihood–multivariate Gaussian density:

$$\begin{aligned} & p(Y | \beta, \sigma^2, \mathbf{X}) \\ &= \frac{1}{\sqrt{\det(2\pi\sigma^2 \mathbf{I})}} \exp\left(-\frac{1}{2} (Y - \mathbf{X}\beta)^\top (\sigma^2 \mathbf{I}_N)^{-1} (Y - \mathbf{X}\beta)\right) \\ &= \frac{1}{(2\pi\sigma^2)^{N/2}} \exp\left(-\frac{1}{2\sigma^2} (Y - \mathbf{X}\beta)^\top (Y - \mathbf{X}\beta)\right) \end{aligned}$$

# Variance Parameter

- Parameter  $\sigma^2$  in the regression model:
    - Not the variance of  $Y_i$ , but the variance of  $\varepsilon_i$
    - Cannot be estimated by the sample variance of  $Y_i$
    - Need to be estimated simultaneously with  $\beta$
- ↪ Not assuming  $\sigma^2$  is known, unlike randomized experiment

- Prior distribution  $p(\beta, \sigma^2)$ : Joint distribution on  $(\beta, \sigma^2)$

- Conjugate prior on  $(\beta, \sigma^2)$ :

$$\sigma^2 \sim \text{Inv-}\chi^2(v_0, \sigma_0^2), \quad \beta \mid \sigma^2 \sim \mathcal{N}(\beta_0, \sigma^2 \Sigma_\beta)$$

- Fully conjugate: Easier derivation of the posterior
  - Unclear interpretation: What is  $\Sigma_\beta$ ?
- Semi-conjugate (conditionally conjugate) prior on  $(\beta, \sigma^2)$ :

$$\sigma^2 \sim \text{Inv-}\chi^2(v_0, \sigma_0^2), \quad \beta \sim \mathcal{N}(\beta_0, \Sigma_\beta)$$

- Not fully conjugate: Trickier derivation
  - Clearer interpretation:  $\Sigma_\beta$  is the prior variance of  $\beta$

# Posterior with Conjugate Prior

- Prior density:

$$p(\beta, \sigma^2) = p(\sigma^2) p(\beta | \sigma^2) \propto \frac{e^{-\frac{v_0 \sigma_0^2}{2\sigma^2}}}{(\sigma^2)^{1+v/2}} \frac{e^{-\frac{1}{2\sigma^2}(\beta - \beta_0)^T \Sigma_\beta^{-1}(\beta - \beta_0)}}{(\sigma^2)^{K/2}}$$

- Posterior density:

$$p(\beta, \sigma^2 | Y, \mathbf{X}) \propto \frac{e^{-\frac{v_0 \sigma_0^2 + Ns^2}{2\sigma^2}}}{(\sigma^2)^{1+(v_0+N)/2}} \frac{e^{-\frac{1}{2\sigma^2}(\beta - \hat{\beta}_C)^T \hat{\Sigma}_C^{-1}(\beta - \hat{\beta}_C)}}{(\sigma^2)^{K/2}},$$

$$\hat{\Sigma}_C \equiv (\mathbf{X}^T \mathbf{X} + \Sigma_\beta^{-1})^{-1}, \quad \hat{\beta}_C \equiv \hat{\Sigma}_C (\mathbf{X}^T \mathbf{Y} + \Sigma_\beta^{-1} \beta_0),$$

$$s^2 \equiv \frac{1}{N-K} (\mathbf{Y} - \mathbf{X} \hat{\beta}_C)^T (\mathbf{Y} - \mathbf{X} \hat{\beta}_C)$$

- Factorization of the posterior:

① Conditional posterior of  $\beta$  given  $\sigma^2$ :  $\beta | \sigma^2, Y, \mathbf{X} \sim \mathcal{N}(\hat{\beta}_C, \sigma^2 \hat{\Sigma}_C)$

② Marginal posterior of  $\sigma^2$ :  $\sigma^2 | Y, \mathbf{X} \sim \text{Inv-}\chi^2(\hat{v}_C, \hat{\sigma}_C^2)$ , where

$$\hat{v}_C \equiv v_0 + N, \quad \hat{\sigma}_C^2 \equiv (v_0 \sigma_0^2 + Ns^2) / \hat{v}_C$$

# Posterior with Semi-conjugate Prior

- Joint posterior density

$$\begin{aligned}
 p(\beta, \sigma^2 \mid Y, \mathbf{X}) &\propto p(\sigma^2) p(\beta) p(Y \mid \beta, \sigma^2, \mathbf{X}) \\
 &\propto \frac{e^{-\frac{v_0 \sigma_0^2}{2\sigma^2}}}{(\sigma^2)^{1+v_0/2}} e^{-\frac{1}{2}(\beta - \beta_0)^\top \Sigma_\beta^{-1} (\beta - \beta_0)} \frac{1}{(\sigma^2)^{N/2}} e^{-\frac{1}{2\sigma^2} (Y - \mathbf{X}\beta)^\top (Y - \mathbf{X}\beta)}
 \end{aligned}$$

- Full conditional posterior:

- Conditional posterior of  $\beta$  given  $\sigma^2$ :

$$\beta \mid \sigma^2, Y, \mathbf{X} \sim \mathcal{N}(\hat{\beta}_{SC}, \hat{\Sigma}_{SC})$$

$$\hat{\Sigma}_{SC} \equiv \left( \frac{1}{\sigma^2} \mathbf{X}^\top \mathbf{X} + \Sigma_\beta^{-1} \right)^{-1}, \quad \hat{\beta}_{SC} \equiv \hat{\Sigma}_{SC} \left( \frac{1}{\sigma^2} \mathbf{X}^\top Y + \Sigma_\beta^{-1} \beta_0 \right)$$

- Conditional posterior of  $\sigma^2$  given  $\beta$ :

$$\sigma^2 \mid \beta, Y, \mathbf{X} \sim \text{Inv-}\chi^2(\hat{v}_C, \hat{\sigma}_{SC}^2)$$

$$S^2 \equiv (Y - \mathbf{X}\beta)^\top (Y - \mathbf{X}\beta), \quad \hat{\sigma}_{SC}^2 \equiv (v_0 \sigma_0^2 + S^2) / \hat{v}_C$$

# Gibbs Sampling

- Joint distributions:
  - Often difficult to summarize analytically
  - Need for Monte Carlo methods—simulations from a joint posterior
- Conjugate prior case:
  - 1 Draw  $\sigma^2$  from the marginal posterior  
`sigma2.samp <- rinvchisq(S, nu.c.hat, sigma2.c.hat)`
  - 2 Draw  $\beta$  from the conditional posterior given  $\sigma^2$   

```
for (s in 1:S) {
  beta.samp[s, ] <- mvrnorm(1, beta.c.hat, sigma2.samp[s] * Sigma.c.hat)
}
```
- Semi-conjugate prior case:
  - Hard to draw a sample directly from the joint posterior
  - Need a method to draw from the joint using the full conditionals
- **Gibbs sampling**: Iterative sampling from full conditionals  
 $s = 0$       Set an arbitrary initial value of  $\sigma^2$   
 $s = 1, 2, \dots$       Repeat
  - 1 Draw  $\beta^{(s)}$  from the conditional posterior given  $\sigma^{2(s-1)}$
  - 2 Draw  $\sigma^{2(s)}$  from the conditional posterior given  $\beta^{(s)}$

# Markov Chain

## ● Markov chain

- Discrete-time stochastic process:  $Z^{(1)}, Z^{(2)}, \dots$
- Markov property:

$$\begin{aligned} p\left(Z^{(s)} \mid Z^{(1)}, \dots, Z^{(s-1)}\right) &= p\left(Z^{(s)} \mid Z^{(s-1)}\right) \\ \Leftrightarrow \left(Z^{(1)}, \dots, Z^{(s-1)}\right) &\perp\!\!\!\perp \left(Z^{(s+1)}, \dots\right) \mid Z^{(s)} \end{aligned}$$

- E.g.: "Drunkard's walk"
- Gibbs sampling generates a Markov chain:

$$\begin{aligned} p\left(\beta^{(s)}, \sigma^{2(s)} \mid \beta^{(1)}, \sigma^{2(1)}, \dots, \beta^{(s-1)}, \sigma^{2(s-1)}, Y, \mathbf{X}\right) \\ = p\left(\sigma^{2(s)} \mid \beta^{(s)}, Y, \mathbf{X}\right) p\left(\beta^{(s)} \mid \sigma^{2(s-1)}, Y, \mathbf{X}\right) \\ \underbrace{\hspace{15em}}_{=p(\beta^{(s)}, \sigma^{2(s)} \mid \sigma^{2(s-1)}, Y, \mathbf{X})} \end{aligned}$$

## ● Stationary distribution, $p^*(Z)$

- Invariance:  $\forall z', p^*(Z = z') = \int_z p(Z^{(s+1)} = z' \mid Z^{(s)} = z) p^*(Z = z) dz$
- Under some conditions, the Markov chain:
  - Has a unique stationary distribution
  - Starts drawing samples from the stationary distribution ("converge")

# Stationary Distribution of the Gibbs Sampler

- Gibbs sampler for Bayesian linear regression:

$$\begin{aligned} & \int \int p\left(\sigma^{2(s)} \mid \beta^{(s)}, Y, \mathbf{X}\right) p\left(\beta^{(s)} \mid \sigma^{2(s-1)} = \sigma^2, Y, \mathbf{X}\right) \\ & \quad \times p\left(\beta, \sigma^2 \mid Y, \mathbf{X}\right) d\beta d\sigma^2 \\ & = p\left(\sigma^{2(s)} \mid \beta^{(s)}, Y, \mathbf{X}\right) p\left(\beta^{(s)} \mid Y, \mathbf{X}\right) = p\left(\beta^{(s)}, \sigma^{2(s)} \mid Y, \mathbf{X}\right) \end{aligned}$$

- Generally, the Gibbs sampler for parameters  $\theta \equiv (\theta_1, \dots, \theta_K)$ 
  - Iterations of alternating draws:

- $\theta_1^{(s)}$  from  $p\left(\theta_1 \mid \theta_2^{(s-1)}, \dots, \theta_K^{(s-1)}, \text{Data}\right)$

- $\theta_2^{(s)}$  from  $p\left(\theta_2 \mid \theta_1^{(s)}, \theta_3^{(s-1)}, \dots, \theta_K^{(s-1)}, \text{Data}\right)$

⋮

- $\theta_K^{(s)}$  from  $p\left(\theta_K \mid \theta_1^{(s)}, \dots, \theta_{K-1}^{(s)}, \text{Data}\right)$

- Stationary distribution:

$$\int p\left(\theta^{(s)} \mid \theta^{(s-1)} = \theta, \text{Data}\right) p(\theta \mid \text{Data}) d\theta = p\left(\theta^{(s)} \mid \text{Data}\right)$$

- The Gibbs sampler converges to the joint posterior distribution



# Unequal Variances

- Pros of the Gibbs sampler:
  - Easy even for many parameters
  - **Readily adapted to extended models with added parameters**
- Linear regression with unequal variances (heteroscedasticity):

$$Y_i = \mathbf{X}_i^\top \boldsymbol{\beta} + \varepsilon_i, \quad \varepsilon_i \sim \mathcal{N}\left(0, \sigma_{j[i]}^2\right)$$

- $J$  groups of observations, and obs.  $i$  belongs to group  $j[i]$
- Group-specific variance,  $\sigma_{j[i]}^2$
- Randomized experiment as a special case
- Joint posterior density of  $\boldsymbol{\beta}, \sigma_1^2, \dots, \sigma_J^2$ :

$$p\left(\boldsymbol{\beta}, \sigma^2 \mid \mathbf{Y}, \mathbf{X}\right) \propto \prod_{j=1}^J \frac{e^{-\frac{v_0 \sigma_0^2}{2\sigma_j^2}}}{\left(\sigma_j^2\right)^{1+v_0/2}} e^{-\frac{1}{2}(\boldsymbol{\beta}-\boldsymbol{\beta}_0)^\top \boldsymbol{\Sigma}_\beta^{-1}(\boldsymbol{\beta}-\boldsymbol{\beta}_0)}$$

$$\times \frac{1}{\left(\sigma_j^2\right)^{N_j/2}} e^{-\frac{1}{2\sigma_j^2}(\mathbf{Y}_j - \mathbf{X}_j \boldsymbol{\beta})^\top (\mathbf{Y}_j - \mathbf{X}_j \boldsymbol{\beta})}$$

where  $\mathbf{Y}_j, \mathbf{X}_j$  are the data of group  $j$

# Modifying the Gibbs Sampler

- Conditional posterior density of:

$\sigma_j^2$  – Equivalent to using the data of group  $j$  only

$$\sigma_j^2 \mid \beta, \mathbf{Y}, \mathbf{X} \sim \text{Inv-}\chi^2 \left( \hat{v}_{UV,j}, \hat{\sigma}_{UV,j}^2 \right)$$

$$\hat{v}_{UV,j} \equiv v_0 + N_j, \quad S_j^2 \equiv (\mathbf{Y}_j - \mathbf{X}_j \beta)^\top (\mathbf{Y}_j - \mathbf{X}_j \beta), \quad \hat{\sigma}_{UV,j}^2 \equiv \frac{(v_0 \sigma_0^2 + S_j^2)}{\hat{v}_{UV,j}}$$

$\beta$  – Equivalent to adding the data of each group sequentially

$$\beta \mid \sigma_{1:j}^2, \mathbf{Y}, \mathbf{X} \sim \mathcal{N} \left( \hat{\beta}_{UV}, \hat{\Sigma}_{UV} \right)$$

$$\hat{\Sigma}_{UV} \equiv \left( \sum_{j=1}^J \frac{1}{\sigma_j^2} \mathbf{X}_j^\top \mathbf{X}_j + \Sigma_\beta^{-1} \right)^{-1}, \quad \hat{\beta}_{UV} \equiv \hat{\Sigma}_{UV} \left( \sum_{j=1}^J \frac{1}{\sigma_j^2} \mathbf{X}_j^\top \mathbf{Y}_j + \Sigma_\beta^{-1} \beta_0 \right)$$

- Gibbs sampler:

① Draw  $\beta^{(s)}$  given  $\sigma_1^{2(s-1)}, \dots, \sigma_J^{2(s-1)}$

② Draw  $\sigma_1^{2(s)}, \dots, \sigma_J^{2(s)}$  independently conditional on  $\beta^{(s)}$

- Conditioning on all the other parameters

↪ Can ignore changes of the other parameters' sampling

# Prior as Additional "Data Points"

- Prior:
  - Extra source of information
  - "Additional data points" (c.f. Beta-Binomial model)
- "Additional data points" for linear regression:
  - $\beta$  –  $K$  observations with  $\mathbf{Y} = \beta_0$ ,  $\mathbf{X} = \mathbf{I}_K$ , and known variance  $\Sigma_\beta$

E.g., If  $\Sigma_\beta = \sigma_\beta^2 \mathbf{I}_K$ ,

$$\hat{\Sigma}_{SC} = \left( \frac{1}{\sigma^2} \mathbf{X}^T \mathbf{X} + \frac{1}{\sigma_\beta^2} \mathbf{I}_K^T \mathbf{I}_K \right)^{-1}, \quad \hat{\beta}_{SC} = \hat{\Sigma}_{SC} \left( \frac{1}{\sigma^2} \mathbf{X}^T \mathbf{Y} + \frac{1}{\sigma_\beta^2} \mathbf{I}_K^T \beta_0 \right)$$

$\sigma^2$  –  $v_0$  observations with sample variance  $\sigma_0^2$

$$\hat{v}_C = \underbrace{v_0}_{\text{prior "N"}} + N,$$

$$\hat{v}_C \hat{\sigma}_{SC}^2 = \underbrace{v_0 \sigma_0^2}_{\text{prior "N" } \times \text{ prior sample variance}} + N \times \underbrace{S^2}_{\text{mean squared residuals}}$$

- Same principle applies to the conjugate case as well

# Noninformative Prior and MLE

- **Noninformative prior:** What does “noninformative” mean?
  - “Zero additional data points”:

$$p(\sigma^2) p(\beta) \propto \underbrace{\frac{e^{-\frac{v_0 \sigma_0^2}{2\sigma^2}}}{1 + v_0/2}}_{(\sigma^2) \rightarrow 1/\sigma^2} \underbrace{e^{-\frac{1}{2}(\beta - \beta_0)^\top \Sigma_\beta^{-1} (\beta - \beta_0)}}_{\rightarrow 1}$$

$$\Rightarrow p(\sigma^2, \beta) \propto \frac{1}{\sigma^2} \Leftrightarrow p(\log \sigma, \beta) \propto 1$$

- *Not necessarily uniform*
- *Not necessarily invariant to transformation*
- Improper for unbounded parameters:  $\int_0^\infty \int_{-\infty}^\infty \frac{1}{\sigma^2} d\beta d\sigma^2 = \infty$
- Leads to proper posterior in linear regression, but not always
- Often recommended “weakly informative prior”:
  - $\sigma \sim \text{Unif}(0, r)$  for large  $r$
  - Allowing heavier tail while preventing diverge
  - Posterior mode  $\rightarrow$  MLE (c.f. Beta-Binomial model)

# Posterior Predictive Distributions

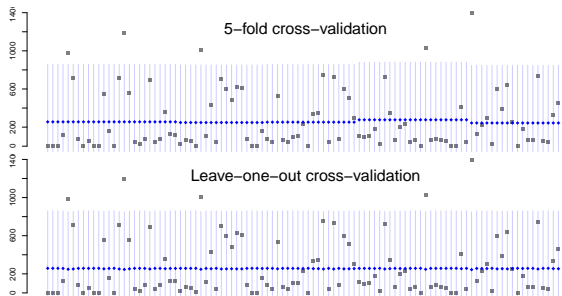
- Purpose of regression analysis: **Prediction**
  - Fit the model to a data set  $(\mathbf{X}, \mathbf{Y})$
  - If observe a new data point  $\mathbf{X}_{\text{new}}$ , what is  $Y_{\text{new}}$  given the model?
- **Posterior predictive distribution** of  $Y_{\text{new}}$ :

$$\begin{aligned}
 p(Y_{\text{new}} | \mathbf{X}_{\text{new}}, \mathbf{X}, \mathbf{Y}) &= \int \int p(Y_{\text{new}}, \beta, \sigma^2 | \mathbf{X}_{\text{new}}, \mathbf{X}, \mathbf{Y}) d\beta d\sigma^2 \\
 &= \int \int \underbrace{p(Y_{\text{new}} | \mathbf{X}_{\text{new}}, \beta, \sigma^2)}_{\text{Model}} \underbrace{p(\beta, \sigma^2 | \mathbf{X}, \mathbf{Y})}_{\text{Posterior}} d\beta d\sigma^2
 \end{aligned}$$

- Computation via Monte Carlo:
  - 1 Draw of  $(\beta^{(s)}, \sigma^{2(s)})$  from the posterior
  - 2 Draw  $Y_{\text{new}}^{(s)}$  from  $\mathcal{N}(\mathbf{X}_{\text{new}}^\top \beta^{(s)}, \sigma^{2(s)})$
- Goodness of fit: Out-of-sample prediction
  - How good  $p(Y_{\text{new}} | \mathbf{X}_{\text{new}}, \mathbf{X}, \mathbf{Y})$  is as a prediction of  $Y_{\text{new}}$
  - Need new data, but data collection is usually one-shot  
 $\rightsquigarrow$  cross-validation

# Cross-validation

- **Cross-validation:**
  - Approximation for out-of-sample prediction
  - $K$ -fold cross-validation
    - 1 Split the data into  $K$  subsets
    - 2 Hold a subset out as a test set
    - 3 Fit the model to the remaining  $K - 1$  subsets
    - 4 Check the prediction on the test set
    - 5 Repeat  $K$  times so each subset is held-out once
  - Leave-one-out cross-validation (LOO-CV):  $K = N$



# Summary

- Linear regression with Gaussian errors
  - Building block of many other models
  - Parameters:  $\beta$  and  $\sigma^2$
  - Conjugate and semi-conjugate priors
    - Conjugate: Prior dependence, mathematical tractability
    - Semi-conjugate: Prior independence, no explicit joint posterior
- Gibbs sampler
  - Alternating draws from full conditional posteriors
  - Distribution of draws converges to the joint posterior
  - Semi-conjugate model, extensions
- Posterior predictive distributions and cross-validation
- Readings for review:
  - 1 Noninformative prior: **BDA3** Sections 2.8–9
  - 2 Bayesian linear regression: **BDA3** Chapters 3 and 14
  - 3 Gibbs sampler: **BDA3** Sections 11.0-1
  - 4 (Optional) Asymptotic approximation: **BDA3** Chapter 4