

# Biostatistics 682: Applied Bayesian Inference

## Lecture 12: Bayesian Linear Regression

**Jian Kang**

Department of Biostatistics  
University of Michigan, Ann Arbor

# Bayesian linear regression

- Linear regression is the most common statistical model.
- The multiple linear regression model is

$$Y_i \sim \text{N} \left( \beta_0 + \sum_{j=1}^p X_{i,j} \beta_j, \sigma^2 \right),$$

for  $i = 1, \dots, n$ .  $Y_i$  are independently across the  $n$  observations.

- Bayesian and classical linear regression are similar if  $n \gg p$  and the priors are uninformative

# Review of least squares

- The least squares estimate of  $\beta = (\beta_0, \beta_1, \dots, \beta_p)^T$  is

$$\hat{\beta}_{\text{OLS}} = \arg \min_{\beta} \sum_{i=1}^n (Y_i - \mu_i)^2,$$

where  $\mu_i = \beta_0 + X_{i,1}\beta_1 + \dots + X_{i,p}\beta_p$ .

- $\hat{\beta}_{\text{OLS}}$  is unbiased even if the errors are non-Gaussian.
- If the errors are Gaussian then the likelihood is proportional to

$$\prod_{i=1}^n \exp \left\{ -\frac{(Y_i - \mu_i)^2}{2\sigma^2} \right\} = \exp \left\{ -\frac{\sum_{i=1}^n (Y_i - \mu_i)^2}{2\sigma^2} \right\}.$$

- Therefore, if the errors are Gaussian  $\hat{\beta}_{\text{OLS}}$  is also the MLE.

# Review of least squares

- Linear regression is often simpler to describe using linear algebra notation.
- Let  $\mathbf{Y} = (Y_1, \dots, Y_n)^T$  be the response vector and  $\mathbf{X}$  be the  $n \times (p + 1)$  matrix of covariates.
- Then the mean of  $\mathbf{Y}$  is  $\mathbf{X}\boldsymbol{\beta}$  and the least squares solution is

$$\boldsymbol{\beta}_{\text{OLS}} = \arg \min_{\boldsymbol{\beta}} (\mathbf{Y} - \mathbf{X}\boldsymbol{\beta})^T (\mathbf{Y} - \mathbf{X}\boldsymbol{\beta}) = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{Y}.$$

- If the errors are Gaussian then the sampling distribution is

$$\hat{\boldsymbol{\beta}}_{\text{OLS}} \sim N [\boldsymbol{\beta}, \sigma^2 (\mathbf{X}^T \mathbf{X})^{-1}].$$

- If the variance  $\sigma^2$  is estimated using the mean squared residual error then the sampling distribution is multivariate  $t$ .

- The likelihood remains

$$Y_i \mid \beta, \sigma^2 \sim N(\beta_0 + X_{i,1}\beta_1 + \dots + X_{i,p}\beta_p, \sigma^2)$$

independent for  $i = 1, \dots, n$  observations.

- A Bayesian analysis also requires priors for  $\beta$  and  $\sigma$ .
- We will focus on prior specification since this piece is uniquely Bayesian.

- For the purpose of setting priors, it is helpful to standardize both the response and each covariate to have mean zero and variance one.
- Many priors for  $\beta$  have been considered:
  - Improper priors
  - Gaussian priors
  - Double exponential priors
  - Many, many more ...

- The Jeffrey's prior is flat  $\pi(\boldsymbol{\beta}) \propto 1$ .
- This is improper, but the posterior is proper under the same conditions required by least squares.
- If  $\sigma$  is known then

$$\boldsymbol{\beta} \mid \mathbf{Y} \sim N \left[ \hat{\boldsymbol{\beta}}_{\text{OLS}}, \sigma^2 (\mathbf{X}^T \mathbf{X})^{-1} \right].$$

- How is this result different from the least squares?

- We rarely know  $\sigma^2$  in practice.
- The Jeffreys prior for  $(\beta, \sigma^2)$  is

$$\pi(\beta, \sigma^2) \propto 1/\sigma^2,$$

which is the limit case of an inverse gamma distribution with shape and rate parameters approaching zero.

- Then the posterior of  $\beta$  follows a multivariate  $t$  centered on  $\hat{\beta}_{\text{OLS}}$ .



# Multivariate normal prior

- Another common prior for is Zellner's g-prior

$$\boldsymbol{\beta} \mid \sigma^2 \sim N \left[ 0, \frac{\sigma^2}{g} (\mathbf{X}^T \mathbf{X})^{-1} \right].$$

- This prior is proper assuming  $\mathbf{X}$  is full rank.
- Then

$$\boldsymbol{\beta} \mid \sigma^2, \mathbf{Y}, \mathbf{X} \sim N \left[ \frac{1}{1+g} \hat{\boldsymbol{\beta}}_{\text{OLS}}, \frac{1}{1+g} \sigma^2 (\mathbf{X}^T \mathbf{X})^{-1} \right].$$

- This shrinks the least squares estimate towards zero.
- $g$  controls the amount of shrinkage.
- $g = 1/n$  is common, and called the unit information prior.

# Univariate Gaussian priors

- If there are many covariates or the covariates are collinear, then  $\hat{\beta}_{OLS}$  is unstable.
- Independent priors can counteract collinearity

$$\beta_j \sim N(0, \sigma^2/g)$$

independent over  $j$ .

- The posterior mode is

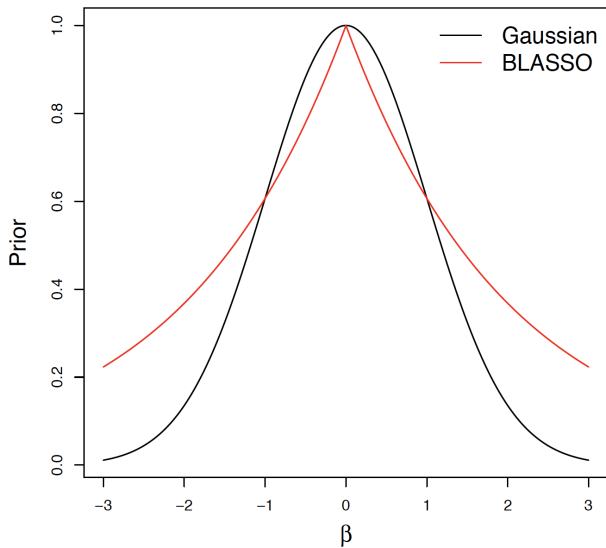
$$\arg \min_{\beta} \sum_{i=1}^n (Y_i - \mu_i)^2 + g \sum_{j=1}^p \beta_j^2.$$

- In classical statistics, this is known as the ridge regression solution and is used to stabilize the least squares solution.

- An increasingly-popular prior is the double exponential or Bayesian LASSO prior
- The prior is  $\beta_j \sim \text{DE}(\tau^2)$  which has the probability density function

$$\pi(\beta_j) \propto \exp\left(-\frac{|\beta_j|}{\tau^2}\right).$$

- The square in the Gaussian prior is replaced with an absolute value
- The shape of the PDF is thus more peaked at zero
- The BLASSO prior favors settings where there are many  $\beta_j$  near zero and a few large  $\beta_j$ .
- That is,  $p$  is large but most of the covariates are noise.



- The posterior model is

$$\arg \min_{\boldsymbol{\beta}} \sum_{i=1}^n (Y_i - \mu_i)^2 + \tau^{-2} \sum_{j=1}^p |\beta_j|.$$

- In classical statistics, this is known as the LASSO solution
- It is popular because it adds stability by shrinking estimates towards zero, and also sets some coefficient to zero
- Covariates with coefficients set to zero and can be excluded from the model.
- LASSO performs variable selection and estimation simultaneously.

# Spike and Slab Priors

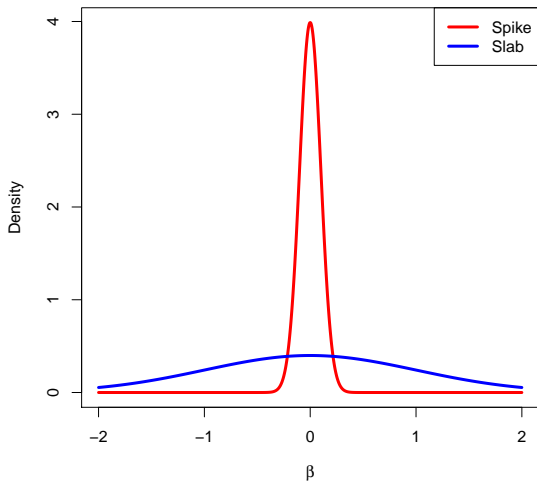
- Mixture Prior

$$\beta_j \mid \gamma_j \sim (1 - \gamma_j)\mathcal{N}(0, c_0^2) + \gamma_j\mathcal{N}(0, c_1^2).$$

$$\gamma_j \sim \text{Bernoulli}(\pi).$$

- The constant  $c_0^2$  is small, so that if  $\gamma_j = 0$ , “ $\beta_j$  could be safely estimated by 0”.
- The constant  $c_1^2$  is large, so that if  $\gamma_j = 1$ , “a non-zero estimate of  $\beta_j$  should probably be included in the final model”.
- It works well for computing the marginal inclusion probability of each covariate and for model averaging
- This model is computationally convenient and extremely flexible

# Spike and Slab Priors



# Posterior computation for Bayesian linear model

- With flat or Gaussian (with fixed prior variance) priors the posterior is available in closed-form and Monte Carlo sampling is not needed
- With Gaussian priors all full conditionals are Gaussian or inverse gamma, and so Gibbs sampling is simple and fast
- With the BLASSO prior the full conditionals are more complicated
  - There is a trick to make all full conditional conjugate so that Gibbs sampling can be used
  - Metropolis sampling works fine too
- With the Spike Slab prior the full conditionals are available
- JAGS can handle all of them



# Summarizing the results

- The standard summary is a table with marginal means and 95% intervals for each  $\beta_j$ .
- This becomes unwieldy for large  $p$
- Picking a subset of covariates is a crucial step in a linear regression analysis
- Common methods include cross-validation and information criteria.

- Say we have a new covariate vector  $\mathbf{X}_{\text{new}}$  and we would like to predict the corresponding response  $Y_{\text{new}}$ .
- A plug-in approach would fix  $\beta$  and  $\sigma$  at their posterior means  $\hat{\beta}$  and  $\hat{\sigma}$  to make predictions

$$Y_{\text{new}} \mid \beta, \sigma^2 \sim N(\mathbf{X}_{\text{new}}\hat{\beta}, \hat{\sigma}^2).$$

- However, this plug-in approach suppresses uncertainty about  $\beta$  and  $\sigma^2$ .
- Therefore these prediction intervals will be slightly too narrow leading to under coverage.

# Posterior predictive distribution (PPD)

- We should really account for all uncertainty when making predictions, including our uncertainty about  $\beta$  and  $\sigma^2$ .
- We really want to PPD

$$\begin{aligned}\pi(Y_{\text{new}} \mid \mathbf{Y}) &= \int \pi(Y_{\text{new}}, \beta, \sigma^2 \mid \mathbf{Y}) d\beta d\sigma^2 \\ &= \int \pi(Y_{\text{new}} \mid \beta, \sigma^2) \pi(\beta, \sigma^2 \mid \mathbf{Y}) d\beta d\sigma^2\end{aligned}$$

- Marginalizing over the model parameters accounts for their uncertainty

# Posterior predictive distribution (PPD)

- MCMC naturally gives draws from  $Y_{\text{new}}$ 's PPD
  - For MCMC iteration  $t$  we have  $\beta^{(t)}$  and  $\sigma^{2(t)}$ .
  - For MCMC iteration  $t$  we sample

$$Y_{\text{new}}^{(t)} \sim N\left(\mathbf{X}\beta^{(t)}, \sigma^{2(t)}\right).$$

- $Y_{\text{new}}^{(1)}, \dots, Y_{\text{new}}^{(S)}$  are samples from the PPD.
- Thus, “Bayesian methods” naturally quantify uncertainty.
- JAGS can handle it.