

# STAT 576 Bayesian Analysis

## Lecture 12: Bayesian Regression Models

Chencheng Cai

Washington State University

# Conditional Modeling

- ▶ Traditional regression models are based on the conditional distribution of the response variable given the covariates.

$$\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\epsilon}$$

# Conditional Modeling

- ▶ Traditional regression models are based on the conditional distribution of the response variable given the covariates.

$$\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\epsilon}$$

where

- ▶  $\mathbf{y}$  is the response variable ( $n \times 1$ ),
- ▶  $\mathbf{X}$  is the design matrix ( $n \times p$ ),
- ▶  $\boldsymbol{\beta}$  is the regression coefficients ( $p \times 1$ ),
- ▶  $\boldsymbol{\epsilon}$  is the error term ( $n \times 1$ ).

# Conditional Modeling

- ▶ Traditional regression models are based on the conditional distribution of the response variable given the covariates.

$$\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\epsilon}$$

where

- ▶  $\mathbf{y}$  is the response variable ( $n \times 1$ ),
  - ▶  $\mathbf{X}$  is the design matrix ( $n \times p$ ),
  - ▶  $\boldsymbol{\beta}$  is the regression coefficients ( $p \times 1$ ),
  - ▶  $\boldsymbol{\epsilon}$  is the error term ( $n \times 1$ ).
- ▶ It is often assumed that
    - ▶  $\mathbf{X}$  is fixed and known,
    - ▶  $\boldsymbol{\epsilon} \sim N(\mathbf{0}, \sigma^2 \mathbf{I})$ .

## Conditional Modeling

- ▶ The inference is based on the conditional distribution of  $\mathbf{y}$  given  $\mathbf{X}$ ,  $\boldsymbol{\beta}$  and  $\sigma^2$  .:

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

- ▶ Frequentists maximize the log-likelihood function:

$$\ell(\boldsymbol{\beta}, \sigma^2; \mathbf{y}) = -\frac{1}{2\sigma^2} \|\mathbf{y} - \mathbf{X}\boldsymbol{\beta}\|^2 - \frac{n}{2} \log \sigma^2$$

## Conditional Modeling

- ▶ The inference is based on the conditional distribution of  $\mathbf{y}$  given  $\mathbf{X}$ ,  $\boldsymbol{\beta}$  and  $\sigma^2$  .:

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

- ▶ Frequentists maximize the log-likelihood function:

$$\ell(\boldsymbol{\beta}, \sigma^2; \mathbf{y}) = -\frac{1}{2\sigma^2} \|\mathbf{y} - \mathbf{X}\boldsymbol{\beta}\|^2 - \frac{n}{2} \log \sigma^2$$

- ▶ The MLE therefore is given by

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

- ▶ However, it is **not** a full probabilistic model.

# Bayesian Linear Regression

- ▶ In Bayesian regression, we treat  $\beta$  and  $\sigma^2$  as random variables.
- ▶ We put priors on  $\beta$  and  $\sigma^2$ :

$$\beta \sim \pi(\beta),$$
$$\sigma^2 \sim \pi(\sigma^2).$$

# Bayesian Linear Regression

- ▶ In Bayesian regression, we treat  $\beta$  and  $\sigma^2$  as random variables.
- ▶ We put priors on  $\beta$  and  $\sigma^2$ :

$$\begin{aligned}\beta &\sim \pi(\beta), \\ \sigma^2 &\sim \pi(\sigma^2).\end{aligned}$$

- ▶ The joint distribution of  $\mathbf{y}$ ,  $\beta$  and  $\sigma^2$  is given by

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

- ▶ The posterior distribution of  $\beta$  and  $\sigma^2$  is given by

$$p(\beta, \sigma^2|\mathbf{y}) \propto p(\mathbf{y}|\beta, \sigma^2)p(\beta)p(\sigma^2).$$



# Bayesian Linear Regression

- ▶ The noninformative prior for  $\beta$  and  $\sigma^2$  is often taken as

$$\begin{aligned}\pi(\beta) &\propto 1, \\ \pi(\sigma^2) &\propto \frac{1}{\sigma^2}.\end{aligned}$$

Derivation: (1) Jeffreys prior (2) results for location-scale families.

# Bayesian Linear Regression

- ▶ The noninformative prior for  $\boldsymbol{\beta}$  and  $\sigma^2$  is often taken as

$$\begin{aligned}\pi(\boldsymbol{\beta}) &\propto 1, \\ \pi(\sigma^2) &\propto \frac{1}{\sigma^2}.\end{aligned}$$

Derivation: (1) Jeffreys prior (2) results for location-scale families.

- ▶ We can derive the posterior distribution of  $\boldsymbol{\beta}$  and  $\sigma^2$  by

$$\begin{aligned}p(\boldsymbol{\beta}, \sigma^2 | \mathbf{y}) &\propto p(\mathbf{y} | \boldsymbol{\beta}, \sigma^2) p(\boldsymbol{\beta}) p(\sigma^2) \\ &\propto \sigma^{-n} \exp\left(-\frac{1}{2\sigma^2} \|\mathbf{y} - \mathbf{X}\boldsymbol{\beta}\|^2\right) \times 1 \times \frac{1}{\sigma^2} \\ &\propto \sigma^{-n} \exp\left(-\frac{1}{2\sigma^2} \|\mathbf{y} - \mathbf{X}\boldsymbol{\beta}\|^2\right) \times \frac{1}{\sigma^2}.\end{aligned}$$

# Bayesian Linear Regression

- ▶ Notice that:

$$\|\mathbf{y} - \mathbf{X}\boldsymbol{\beta}\|^2 = (\boldsymbol{\beta} - \hat{\boldsymbol{\beta}})^T \mathbf{X}^T \mathbf{X} (\boldsymbol{\beta} - \hat{\boldsymbol{\beta}}) + \|\mathbf{y}\|^2 - \|\mathbf{X}\hat{\boldsymbol{\beta}}\|^2$$

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

# Bayesian Linear Regression

- ▶ Notice that:

$$\|\mathbf{y} - \mathbf{X}\boldsymbol{\beta}\|^2 = (\boldsymbol{\beta} - \hat{\boldsymbol{\beta}})^T \mathbf{X}^T \mathbf{X} (\boldsymbol{\beta} - \hat{\boldsymbol{\beta}}) + \|\mathbf{y}\|^2 - \|\mathbf{X}\hat{\boldsymbol{\beta}}\|^2$$

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

- ▶ Therefore, the posterior distribution of  $\boldsymbol{\beta}$  and  $\sigma^2$  is given by

$$p(\boldsymbol{\beta}, \sigma^2 | \mathbf{y}) \propto \sigma^{-n-2} \exp\left(-\frac{(\boldsymbol{\beta} - \hat{\boldsymbol{\beta}})^T \mathbf{X}^T \mathbf{X} (\boldsymbol{\beta} - \hat{\boldsymbol{\beta}}) + \|\mathbf{y}\|^2 - \|\mathbf{X}\hat{\boldsymbol{\beta}}\|^2}{2\sigma^2}\right).$$

# Bayesian Linear Regression

- ▶ Notice that:

$$\|\mathbf{y} - \mathbf{X}\boldsymbol{\beta}\|^2 = (\boldsymbol{\beta} - \hat{\boldsymbol{\beta}})^T \mathbf{X}^T \mathbf{X} (\boldsymbol{\beta} - \hat{\boldsymbol{\beta}}) + \|\mathbf{y}\|^2 - \|\mathbf{X}\hat{\boldsymbol{\beta}}\|^2$$

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

- ▶ Therefore, the posterior distribution of  $\boldsymbol{\beta}$  and  $\sigma^2$  is given by

$$p(\boldsymbol{\beta}, \sigma^2 | \mathbf{y}) \propto \sigma^{-n-2} \exp\left(-\frac{(\boldsymbol{\beta} - \hat{\boldsymbol{\beta}})^T \mathbf{X}^T \mathbf{X} (\boldsymbol{\beta} - \hat{\boldsymbol{\beta}}) + \|\mathbf{y}\|^2 - \|\mathbf{X}\hat{\boldsymbol{\beta}}\|^2}{2\sigma^2}\right).$$

- ▶ Compared to Normal-Inverse-Gamma distribution, the normal component is replaced with a multivariate normal distribution.
- ▶ Compared to Normal-Inverse-Wishart distribution, the covariance component is replaced with  $\sigma^2(\mathbf{X}^T \mathbf{X})^{-1}$ .

## Bayesian Linear Regression

$$p(\boldsymbol{\beta}, \sigma^2 | \mathbf{y}) \propto \sigma^{-n-2} \exp \left( -\frac{(\boldsymbol{\beta} - \hat{\boldsymbol{\beta}})^T \mathbf{X}^T \mathbf{X} (\boldsymbol{\beta} - \hat{\boldsymbol{\beta}}) + \|\mathbf{y}\|^2 - \|\mathbf{X} \hat{\boldsymbol{\beta}}\|^2}{2\sigma^2} \right)$$

## Bayesian Linear Regression

$$p(\boldsymbol{\beta}, \sigma^2 | \mathbf{y}) \propto \sigma^{-n-2} \exp \left( -\frac{(\boldsymbol{\beta} - \hat{\boldsymbol{\beta}})^T \mathbf{X}^T \mathbf{X} (\boldsymbol{\beta} - \hat{\boldsymbol{\beta}}) + \|\mathbf{y}\|^2 - \|\mathbf{X} \hat{\boldsymbol{\beta}}\|^2}{2\sigma^2} \right)$$

- ▶ The conditional posterior of  $\boldsymbol{\beta}$  given  $\sigma^2$  and  $\mathbf{y}$  is given by

$$\boldsymbol{\beta} | \sigma^2, \mathbf{y} \sim \mathcal{N} \left( \hat{\boldsymbol{\beta}}, \sigma^2 (\mathbf{X}^T \mathbf{X})^{-1} \right)$$

## Bayesian Linear Regression

$$p(\boldsymbol{\beta}, \sigma^2 | \mathbf{y}) \propto \sigma^{-n-2} \exp \left( -\frac{(\boldsymbol{\beta} - \hat{\boldsymbol{\beta}})^T \mathbf{X}^T \mathbf{X} (\boldsymbol{\beta} - \hat{\boldsymbol{\beta}}) + \|\mathbf{y}\|^2 - \|\mathbf{X}\hat{\boldsymbol{\beta}}\|^2}{2\sigma^2} \right)$$

- ▶ The conditional posterior of  $\boldsymbol{\beta}$  given  $\sigma^2$  and  $\mathbf{y}$  is given by

$$\boldsymbol{\beta} | \sigma^2, \mathbf{y} \sim \mathcal{N} \left( \hat{\boldsymbol{\beta}}, \sigma^2 (\mathbf{X}^T \mathbf{X})^{-1} \right)$$

- ▶ The conditional posterior of  $\sigma^2$  given  $\boldsymbol{\beta}$  and  $\mathbf{y}$  is given by

$$\sigma^2 | \boldsymbol{\beta}, \mathbf{y} \sim \text{InvGamma} \left( \frac{n}{2}, \frac{1}{2} \|\mathbf{y} - \mathbf{X}\boldsymbol{\beta}\|^2 \right)$$



## Bayesian Linear Regression

$$p(\boldsymbol{\beta}, \sigma^2 | \mathbf{y}) \propto \sigma^{-n-2} \exp \left( -\frac{(\boldsymbol{\beta} - \hat{\boldsymbol{\beta}})^T \mathbf{X}^T \mathbf{X} (\boldsymbol{\beta} - \hat{\boldsymbol{\beta}}) + \|\mathbf{y}\|^2 - \|\mathbf{X}\hat{\boldsymbol{\beta}}\|^2}{2\sigma^2} \right)$$

- ▶ The conditional posterior of  $\boldsymbol{\beta}$  given  $\sigma^2$  and  $\mathbf{y}$  is given by

$$\boldsymbol{\beta} | \sigma^2, \mathbf{y} \sim \mathcal{N} \left( \hat{\boldsymbol{\beta}}, \sigma^2 (\mathbf{X}^T \mathbf{X})^{-1} \right)$$

- ▶ The conditional posterior of  $\sigma^2$  given  $\boldsymbol{\beta}$  and  $\mathbf{y}$  is given by

$$\sigma^2 | \boldsymbol{\beta}, \mathbf{y} \sim \text{InvGamma} \left( \frac{n}{2}, \frac{1}{2} \|\mathbf{y} - \mathbf{X}\boldsymbol{\beta}\|^2 \right)$$

- ▶ The marginal posterior of  $\sigma^2$  is given by

$$\sigma^2 | \mathbf{y} \sim \text{InvGamma} \left( \frac{n-p}{2}, \frac{\|\mathbf{y}\|^2 - \|\mathbf{X}\hat{\boldsymbol{\beta}}\|^2}{2} \right)$$

## Bayesian Linear Regression

$$p(\boldsymbol{\beta}, \sigma^2 | \mathbf{y}) \propto \sigma^{-n-2} \exp \left( -\frac{(\boldsymbol{\beta} - \hat{\boldsymbol{\beta}})^T \mathbf{X}^T \mathbf{X} (\boldsymbol{\beta} - \hat{\boldsymbol{\beta}}) + \|\mathbf{y}\|^2 - \|\mathbf{X}\hat{\boldsymbol{\beta}}\|^2}{2\sigma^2} \right)$$

- ▶ The marginal posterior of  $\boldsymbol{\beta}$  can be obtained by

$$\begin{aligned} p(\boldsymbol{\beta} | \mathbf{y}) &= \frac{p(\boldsymbol{\beta}, \sigma^2 | \mathbf{y})}{p(\sigma^2 | \boldsymbol{\beta}, \mathbf{y})} \propto \|\mathbf{y} - \mathbf{X}\boldsymbol{\beta}\|^{-n} \\ &\propto \left( (\boldsymbol{\beta} - \hat{\boldsymbol{\beta}})^T \mathbf{X}^T \mathbf{X} (\boldsymbol{\beta} - \hat{\boldsymbol{\beta}}) + \|\mathbf{y}\|^2 - \|\mathbf{X}\hat{\boldsymbol{\beta}}\|^2 \right)^{-n/2} \\ &\propto \left( 1 + \frac{(\boldsymbol{\beta} - \hat{\boldsymbol{\beta}})^T \mathbf{X}^T \mathbf{X} (\boldsymbol{\beta} - \hat{\boldsymbol{\beta}})}{\|\mathbf{y}\|^2 - \|\mathbf{X}\hat{\boldsymbol{\beta}}\|^2} \right)^{-n/2} \end{aligned}$$

## Bayesian Linear Regression

$$p(\boldsymbol{\beta}, \sigma^2 | \mathbf{y}) \propto \sigma^{-n-2} \exp \left( -\frac{(\boldsymbol{\beta} - \hat{\boldsymbol{\beta}})^T \mathbf{X}^T \mathbf{X} (\boldsymbol{\beta} - \hat{\boldsymbol{\beta}}) + \|\mathbf{y}\|^2 - \|\mathbf{X}\hat{\boldsymbol{\beta}}\|^2}{2\sigma^2} \right)$$

- ▶ The marginal posterior of  $\boldsymbol{\beta}$  can be obtained by

$$\begin{aligned} p(\boldsymbol{\beta} | \mathbf{y}) &= \frac{p(\boldsymbol{\beta}, \sigma^2 | \mathbf{y})}{p(\sigma^2 | \boldsymbol{\beta}, \mathbf{y})} \propto \|\mathbf{y} - \mathbf{X}\boldsymbol{\beta}\|^{-n} \\ &\propto \left( (\boldsymbol{\beta} - \hat{\boldsymbol{\beta}})^T \mathbf{X}^T \mathbf{X} (\boldsymbol{\beta} - \hat{\boldsymbol{\beta}}) + \|\mathbf{y}\|^2 - \|\mathbf{X}\hat{\boldsymbol{\beta}}\|^2 \right)^{-n/2} \\ &\propto \left( 1 + \frac{(\boldsymbol{\beta} - \hat{\boldsymbol{\beta}})^T \mathbf{X}^T \mathbf{X} (\boldsymbol{\beta} - \hat{\boldsymbol{\beta}})}{\|\mathbf{y}\|^2 - \|\mathbf{X}\hat{\boldsymbol{\beta}}\|^2} \right)^{-n/2} \end{aligned}$$

- ▶ This is a multivariate t distribution with degree  $n - p$ , mean  $\hat{\boldsymbol{\beta}}$  and covariance  $\frac{\|\mathbf{y}\|^2 - \|\mathbf{X}\hat{\boldsymbol{\beta}}\|^2}{n-p} (\mathbf{X}^T \mathbf{X})^{-1}$ .

# Sampling from the Posterior

- ▶ Easier way:

$$\sigma^2 \mid \mathbf{y} \sim \text{InvGamma} \left( \frac{n-p}{2}, \frac{\|\mathbf{y}\|^2 - \|\mathbf{X}\hat{\boldsymbol{\beta}}\|^2}{2} \right)$$
$$\boldsymbol{\beta} \mid \sigma^2, \mathbf{y} \sim \mathcal{N} \left( \hat{\boldsymbol{\beta}}, \sigma^2 (\mathbf{X}^T \mathbf{X})^{-1} \right)$$

# Sampling from the Posterior

- ▶ Easier way:

$$\sigma^2 \mid \mathbf{y} \sim \text{InvGamma} \left( \frac{n-p}{2}, \frac{\|\mathbf{y}\|^2 - \|\mathbf{X}\hat{\boldsymbol{\beta}}\|^2}{2} \right)$$
$$\boldsymbol{\beta} \mid \sigma^2, \mathbf{y} \sim \mathcal{N} \left( \hat{\boldsymbol{\beta}}, \sigma^2 (\mathbf{X}^T \mathbf{X})^{-1} \right)$$

- ▶ Harder way:

$$\boldsymbol{\beta} \mid \mathbf{y} \sim t_{n-p} \left( \hat{\boldsymbol{\beta}}, \frac{\|\mathbf{y}\|^2 - \|\mathbf{X}\hat{\boldsymbol{\beta}}\|^2}{n-p} (\mathbf{X}^T \mathbf{X})^{-1} \right)$$
$$\sigma^2 \mid \boldsymbol{\beta}, \mathbf{y} \sim \text{InvGamma} \left( \frac{n}{2}, \frac{1}{2} \|\mathbf{y} - \mathbf{X}\boldsymbol{\beta}\|^2 \right)$$

## Sampling from the Posterior

$$\sigma^2 \mid \mathbf{y} \sim \text{InvGamma} \left( \frac{n-p}{2}, \frac{\|\mathbf{y}\|^2 - \|\mathbf{X}\hat{\boldsymbol{\beta}}\|^2}{2} \right)$$

$$\boldsymbol{\beta} \mid \sigma^2, \mathbf{y} \sim \mathcal{N} \left( \hat{\boldsymbol{\beta}}, \sigma^2 (\mathbf{X}^T \mathbf{X})^{-1} \right)$$

## Sampling from the Posterior

$$\sigma^2 \mid \mathbf{y} \sim \text{InvGamma} \left( \frac{n-p}{2}, \frac{\|\mathbf{y}\|^2 - \|\mathbf{X}\hat{\boldsymbol{\beta}}\|^2}{2} \right)$$

$$\boldsymbol{\beta} \mid \sigma^2, \mathbf{y} \sim \mathcal{N} \left( \hat{\boldsymbol{\beta}}, \sigma^2 (\mathbf{X}^T \mathbf{X})^{-1} \right)$$

- ▶ Sampling from  $\text{InvGamma}(\alpha, \beta)$ :
  - ▶ Generate  $x \sim \chi_{2\alpha}^2$ ,
  - ▶ Then  $y = \frac{\beta}{2x}$ .

## Sampling from the Posterior

$$\sigma^2 \mid \mathbf{y} \sim \text{InvGamma} \left( \frac{n-p}{2}, \frac{\|\mathbf{y}\|^2 - \|\mathbf{X}\hat{\boldsymbol{\beta}}\|^2}{2} \right)$$

$$\boldsymbol{\beta} \mid \sigma^2, \mathbf{y} \sim \mathcal{N} \left( \hat{\boldsymbol{\beta}}, \sigma^2 (\mathbf{X}^T \mathbf{X})^{-1} \right)$$

- ▶ Sampling from  $\text{InvGamma}(\alpha, \beta)$ :
  - ▶ Generate  $x \sim \chi_{2\alpha}^2$ ,
  - ▶ Then  $y = \frac{\beta}{2x}$ .
- ▶ Sampling from  $\mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma})$ :
  - ▶ Cholesky decomposition:  $\boldsymbol{\Sigma} = \mathbf{L}\mathbf{L}^T$ , where  $\mathbf{L}$  is lower triangular,
  - ▶ Generate  $\mathbf{z} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$ ,
  - ▶ Then  $\mathbf{x} = \boldsymbol{\mu} + \mathbf{L}\mathbf{z}$ .



## Predictive Distribution

Suppose  $\sigma^2$  is known.

- ▶ The distribution for new observation  $\tilde{y}$  given new covariate  $\tilde{\mathbf{X}}$  is given by

$$\tilde{y} | \mathbf{y}, \sigma^2 \sim \mathcal{N}(\tilde{\mathbf{X}} \hat{\boldsymbol{\beta}}, \sigma^2 \mathbf{I} + \sigma^2 \tilde{\mathbf{X}} (\mathbf{X}^T \mathbf{X})^{-1} \tilde{\mathbf{X}}^T).$$

- ▶ The mean is  $\tilde{\mathbf{X}} \hat{\boldsymbol{\beta}}$ ,
- ▶ The variance is  $\sigma^2 (\mathbf{I} + \tilde{\mathbf{X}} (\mathbf{X}^T \mathbf{X})^{-1} \tilde{\mathbf{X}}^T)$ .

## Predictive Distribution

Suppose  $\sigma^2$  is known.

- ▶ The distribution for new observation  $\tilde{y}$  given new covariate  $\tilde{\mathbf{X}}$  is given by

$$\tilde{y}|\mathbf{y}, \sigma^2 \sim \mathcal{N}(\tilde{\mathbf{X}}\hat{\boldsymbol{\beta}}, \sigma^2\mathbf{I} + \sigma^2\tilde{\mathbf{X}}(\mathbf{X}^T\mathbf{X})^{-1}\tilde{\mathbf{X}}^T).$$

- ▶ The mean is  $\tilde{\mathbf{X}}\hat{\boldsymbol{\beta}}$ ,
- ▶ The variance is  $\sigma^2\left(\mathbf{I} + \tilde{\mathbf{X}}(\mathbf{X}^T\mathbf{X})^{-1}\tilde{\mathbf{X}}^T\right)$ .

Suppose  $\sigma^2$  is unknown.

- ▶ The distribution for new observation  $\tilde{y}$  given new covariate  $\tilde{\mathbf{X}}$  is a linear transformation of a multivariate t distribution plus a Gaussian noise.
- ▶ The mean is  $\tilde{\mathbf{X}}\hat{\boldsymbol{\beta}}$ ,
- ▶ The variance is  $\frac{\|\mathbf{y}\|^2 - \|\mathbf{X}\hat{\boldsymbol{\beta}}\|^2}{n-p-2}\tilde{\mathbf{X}}(\mathbf{X}^T\mathbf{X})^{-1}\tilde{\mathbf{X}}^T + \sigma^2\mathbf{I}$

## Example

- ▶ Example from textbook Sec. 14.3.
- ▶ The data contains the election data for the U.S. House of Representatives in the past century (1900 – 2000).
- ▶ We would like to study the relationship between the percentage of votes for the incumbent party and the decision whether the incumbent officeholder runs for reelection.
- ▶ Goal: check if there is an advantage for the incumbent officeholder to reelect.
- ▶ Some facts of the data:
  - ▶ Election every two years.
  - ▶ The incumbent party is the party that won the previous election.
  - ▶ 435 districts in the U.S. House of Representatives.
  - ▶ Roughly 100 - 150 districts are uncontested.

## Example

We formulate the problem as a simple linear regression model.

$$y_i = \alpha + \beta R_i + \epsilon_i$$

- ▶  $y_i$ : the percentage of votes for the **incumbent party** in district  $i$ .
- ▶  $R_i$ : a binary variable indicating whether the **incumbent officeholder** runs for reelection.
- ▶  $\alpha$ : the expected percentage of votes for the incumbent party when they incumbent officeholder **does not** run for reelection.
- ▶  $\alpha + \beta$ : the expected percentage of votes for the incumbent party when the incumbent officeholder **does** run for reelection.
- ▶  $\beta$ : **incumbency advantage**.

## Example

- ▶ The current model may have selection bias in the dataset.
- ▶ I.e. some variables may affect both the decision of reelection and the percentage of votes.
- ▶ We should include those variables in the model as well.

$$y_i = \alpha + \beta R_i + \gamma z_i + \delta P_i + \epsilon_i$$

- ▶  $z_i$ : the percentage of votes for the incumbent party in the **previous election**.
- ▶  $P_i$ : the indicator for Democratic party (1) or Republican party (0) controlling the seat.

## Example

With noninformative priors, the posterior inferences for the year 1988 are displayed below.

Variable	Posterior quantiles				
	2.5%	25%	median	75%	97.5%
Incumbency	0.084	0.103	0.114	0.124	0.144
Vote proportion in 1986	0.576	0.627	0.654	0.680	0.731
Incumbent party	-0.014	-0.009	-0.007	-0.004	0.001
Constant term	0.066	0.106	0.127	0.148	0.188
$\sigma$ (residual sd)	0.061	0.064	0.066	0.068	0.071

- ▶ The incumbency advantage is estimated to be 11.4% and is significant.
- ▶ It shows a strong autoregressive effect in the percentage of votes for the incumbent party.
- ▶ Party difference is not significant.

# Generalizations

We consider the following generalizations of the linear regression model in the subsequent slides.

- ▶ **Diverse Covariance Structures:** We may consider different covariance structures for the error term.
- ▶ **Regularization:** Sometimes we would like to choose a prior that encourages sparsity in the regression coefficients to prevent overfitting.
- ▶ **Hierarchical Linear Models:** We assume the regression coefficients are drawn from a common distribution for different subsets of data.

# Covariance Structure

In the general case, we may consider the following covariance structures for the error term:

$$\epsilon \sim \mathcal{N}(\mathbf{0}, \Sigma)$$

where  $\Sigma$  is a positive definite matrix, that allows for different variances and correlations between the errors.



# Covariance Structure

In the general case, we may consider the following covariance structures for the error term:

$$\boldsymbol{\epsilon} \sim \mathcal{N}(\mathbf{0}, \boldsymbol{\Sigma})$$

where  $\boldsymbol{\Sigma}$  is a positive definite matrix, that allows for different variances and correlations between the errors.

In this case, the model is given by

$$\mathbf{y} \sim \mathcal{N}(\mathbf{X}\boldsymbol{\beta}, \boldsymbol{\Sigma})$$

## Covariance Structure — Known Covariance

If  $\Sigma$  is known, the posterior distribution of  $\beta$  is given by

$$\begin{aligned} p(\beta|\mathbf{y}, \Sigma) &\propto p(\mathbf{y}|\beta, \Sigma)p(\beta) \\ &\propto \exp\left(-\frac{1}{2}(\mathbf{y} - \mathbf{X}\beta)^T \Sigma^{-1}(\mathbf{y} - \mathbf{X}\beta)\right) \times 1 \\ &\propto \exp\left(-\frac{1}{2}(\beta - \hat{\beta})^T \mathbf{X}^T \Sigma^{-1} \mathbf{X}(\beta - \hat{\beta})\right) \\ &\sim \mathcal{N}\left(\hat{\beta}, (\mathbf{X}^T \Sigma^{-1} \mathbf{X})^{-1}\right) \end{aligned}$$

with

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

## Covariance Structure — Unknown Covariance

If  $\Sigma$  is unknown, we may put a prior on  $\Sigma$  as well.

$$\begin{aligned} p(\Sigma | \mathbf{y}, \beta) &\propto \frac{p(\beta, \Sigma | \mathbf{y})}{p(\beta | \mathbf{y}, \Sigma)} \\ &\propto p(\Sigma) |\Sigma|^{-1/2} |\mathbf{X}^T \Sigma^{-1} \mathbf{X}|^{1/2} \exp\left(-\frac{1}{2}(\mathbf{y} - \mathbf{X}\hat{\beta})^T \Sigma^{-1}(\mathbf{y} - \mathbf{X}\hat{\beta})\right) \end{aligned}$$

## Covariance Structure — Unknown Covariance

If  $\Sigma$  is unknown, we may put a prior on  $\Sigma$  as well.

$$p(\Sigma | \mathbf{y}, \beta) \propto \frac{p(\beta, \Sigma | \mathbf{y})}{p(\beta | \mathbf{y}, \Sigma)}$$
$$\propto p(\Sigma) |\Sigma|^{-1/2} |\mathbf{X}^T \Sigma^{-1} \mathbf{X}|^{1/2} \exp\left(-\frac{1}{2}(\mathbf{y} - \mathbf{X}\hat{\beta})^T \Sigma^{-1}(\mathbf{y} - \mathbf{X}\hat{\beta})\right)$$

- ▶ It is difficult to set up a prior for  $\Sigma$ .
- ▶ It is difficult to draw from this posterior distribution.
- ▶ Therefore, we often need some further simplification on  $\Sigma$ .

## Covariance Structure — Simplified Covariance

If the covariance matrix  $\Sigma$  is proportional to a known matrix  $Q$ , that is

$$\Sigma = \sigma^2 Q.$$

## Covariance Structure — Simplified Covariance

If the covariance matrix  $\Sigma$  is proportional to a known matrix  $Q$ , that is

$$\Sigma = \sigma^2 Q.$$

Then the posterior distribution of  $\beta$  is multivariate  $t$  and the posterior distribution of  $\sigma^2$  is inverse gamma.

## Covariance Structure — Simplified Covariance

If the covariance matrix  $\Sigma$  is proportional to a known matrix  $Q$ , that is

$$\Sigma = \sigma^2 Q.$$

Then the posterior distribution of  $\beta$  is multivariate  $t$  and the posterior distribution of  $\sigma^2$  is inverse gamma.

- ▶ One can derive it from the posterior distribution of  $\beta$  and  $\sigma^2$  on the previous few slides.

## Covariance Structure — Simplified Covariance

If the covariance matrix  $\Sigma$  is proportional to a known matrix  $Q$ , that is

$$\Sigma = \sigma^2 Q.$$

Then the posterior distribution of  $\beta$  is multivariate  $t$  and the posterior distribution of  $\sigma^2$  is inverse gamma.

- ▶ One can derive it from the posterior distribution of  $\beta$  and  $\sigma^2$  on the previous few slides.
- ▶ Or, it can be seen from the following transformation of data:

$$\begin{aligned} \mathbf{y}^* &= Q^{-1/2} \mathbf{y}, \\ \mathbf{X}^* &= Q^{-1/2} \mathbf{X}. \end{aligned}$$

$Q^{-1/2}$  is any matrix such that  $(Q^{-1/2})^T Q Q^{-1/2} = I$ .

Then the linear regression problem becomes regress  $\mathbf{y}^*$  on  $\mathbf{X}^*$  with i.i.d. noise.

All previous results apply.



## Covariance Structure — Simplified Covariance

In a weighted regression model, we may consider the following covariance structure for the error term:

$$\Sigma_{ii} = \sigma^2/w_i$$

where  $w_i$  is the weight for the  $i$ th observation, and  $\Sigma_{ii}$  is the  $i$ th diagonal element of  $\Sigma$ .

## Covariance Structure — Simplified Covariance

In a weighted regression model, we may consider the following covariance structure for the error term:

$$\Sigma_{ii} = \sigma^2/w_i$$

where  $w_i$  is the weight for the  $i$ th observation, and  $\Sigma_{ii}$  is the  $i$ th diagonal element of  $\Sigma$ .

- ▶ The model is the same as the previous one, with

$$Q = \text{diag}(w_1, \dots, w_n)$$

- ▶ All previous results apply.

## Covariance Structure — Simplified Covariance

The unequal weights can be generalized to a more general setting by introducing the unequalness parameter  $\phi$  such that

$$\Sigma_{ii} = \sigma^2 v(w_i, \phi)$$

where  $\phi \in [0, 1]$  controls the unequalness.

## Covariance Structure — Simplified Covariance

The unequal weights can be generalized to a more general setting by introducing the unequalness parameter  $\phi$  such that

$$\Sigma_{ii} = \sigma^2 v(w_i, \phi)$$

where  $\phi \in [0, 1]$  controls the unequalness.

- ▶ Example:  $v(w_i, \phi) = w_i^{-\phi}$ .  $\phi = 0$  is the equal weight case and  $\phi = 1$  is the inverse weight case.
- ▶ Example:  $v(w_i, \phi) = 1 - \phi + \phi/w_i$ .  $\phi = 0$  is the equal weight case and  $\phi = 1$  is the inverse weight case.

## Covariance Structure — Simplified Covariance

The unequal weights can be generalized to a more general setting by introducing the unequalness parameter  $\phi$  such that

$$\Sigma_{ii} = \sigma^2 v(w_i, \phi)$$

where  $\phi \in [0, 1]$  controls the unequalness.

- ▶ Example:  $v(w_i, \phi) = w_i^{-\phi}$ .  $\phi = 0$  is the equal weight case and  $\phi = 1$  is the inverse weight case.
- ▶ Example:  $v(w_i, \phi) = 1 - \phi + \phi/w_i$ .  $\phi = 0$  is the equal weight case and  $\phi = 1$  is the inverse weight case.
- ▶ A natural noninformative prior for  $\phi$  is the uniform distribution on  $[0, 1]$ .
- ▶ For the posterior and its sampling, please check textbook Eq. (14.21) and (14.22).

# Regularization

In linear regression problem, the regularized least squares minimize the following objective function:

$$\min_{\boldsymbol{\beta}} \|\mathbf{y} - \mathbf{X}\boldsymbol{\beta}\|^2 + \lambda R(\boldsymbol{\beta}),$$

where  $R(\boldsymbol{\beta})$  is a penalty term that penalizes the complexity of the model.

# Regularization

In linear regression problem, the regularized least squares minimize the following objective function:

$$\min_{\boldsymbol{\beta}} \|\mathbf{y} - \mathbf{X}\boldsymbol{\beta}\|^2 + \lambda R(\boldsymbol{\beta}),$$

where  $R(\boldsymbol{\beta})$  is a penalty term that penalizes the complexity of the model.

- ▶ Ridge regression:  $R(\boldsymbol{\beta}) = \|\boldsymbol{\beta}\|^2$ .
- ▶ Lasso regression:  $R(\boldsymbol{\beta}) = \|\boldsymbol{\beta}\|_1$ .
- ▶ Elastic net:  $R(\boldsymbol{\beta}) = \alpha\|\boldsymbol{\beta}\|_1 + (1 - \alpha)\|\boldsymbol{\beta}\|^2$ .

# Regularization

In linear regression problem, the regularized least squares minimize the following objective function:

$$\min_{\beta} \|\mathbf{y} - \mathbf{X}\beta\|^2 + \lambda R(\beta),$$

where  $R(\beta)$  is a penalty term that penalizes the complexity of the model.

- ▶ Ridge regression:  $R(\beta) = \|\beta\|^2$ .
- ▶ Lasso regression:  $R(\beta) = \|\beta\|_1$ .
- ▶ Elastic net:  $R(\beta) = \alpha\|\beta\|_1 + (1 - \alpha)\|\beta\|^2$ .
- ▶ Notice that the sum of squared errors is equivalent to the negative log-likelihood function.
- ▶ The regularized least squares is equivalent to the maximum a posteriori estimation with a prior on  $\beta$  that corresponds to the exponential of the negative penalty.



## Regularization — Ridge

In Ridge regression, we put a Gaussian prior on  $\beta$ :

$$p(\beta) \propto \exp\left(-\frac{\lambda}{2\sigma^2}\|\beta\|^2\right)$$

This is a multivariate normal distribution with mean  $\mathbf{0}$  and covariance  $\frac{\sigma^2}{\lambda}\mathbf{I}$ .

## Regularization — Ridge

In Ridge regression, we put a Gaussian prior on  $\beta$ :

$$p(\beta) \propto \exp\left(-\frac{\lambda}{2\sigma^2}\|\beta\|^2\right)$$

This is a multivariate normal distribution with mean  $\mathbf{0}$  and covariance  $\frac{\sigma^2}{\lambda}\mathbf{I}$ .

The posterior is (under noninformative prior for  $\sigma^2$ )

$$\begin{aligned} p(\beta, \sigma^2 | \mathbf{y}) &\propto \sigma^{-n-2} \exp\left(-\frac{1}{2\sigma^2}\|\mathbf{y} - \mathbf{X}\beta\|^2 - \lambda\|\beta\|^2\right) \\ &\propto \sigma^{-n-2} \exp\left(-\frac{1}{2\sigma^2}(\beta - \hat{\beta})^T(\mathbf{X}^T\mathbf{X} + \lambda\mathbf{I})(\beta - \hat{\beta})\right) \\ &\quad \times \exp\left(-\frac{1}{2\sigma^2}(\|\mathbf{y}\|^2 - \mathbf{y}^T\mathbf{X}(\mathbf{X}^T\mathbf{X} + \lambda\mathbf{I})^{-1}\mathbf{X}^T\mathbf{y})\right) \end{aligned}$$

with  $\hat{\beta} = (\mathbf{X}^T\mathbf{X} + \lambda\mathbf{I})^{-1}\mathbf{X}^T\mathbf{y}$ . The conditional/marginal posteriors are the similar as before except that  $\mathbf{X}^T\mathbf{X}$  is replaced with  $\mathbf{X}^T\mathbf{X} + \lambda\mathbf{I}$ .

## Regularization — LASSO

In LASSO (Least Absolute Shrinkage and Selection Operator) regression, we put a Laplace prior on  $\beta$ :

$$p(\beta) \propto \exp\left(-\frac{\lambda}{2\sigma^2}\|\beta\|_1\right)$$

## Regularization — LASSO

In LASSO (Least Absolute Shrinkage and Selection Operator) regression, we put a Laplace prior on  $\beta$ :

$$p(\beta) \propto \exp\left(-\frac{\lambda}{2\sigma^2}\|\beta\|_1\right)$$

- ▶ The posterior distribution is not a standard distribution.
- ▶ We usually do not have a closed form for the posterior mode.

## Regularization — LASSO

In LASSO (Least Absolute Shrinkage and Selection Operator) regression, we put a Laplace prior on  $\beta$ :

$$p(\beta) \propto \exp\left(-\frac{\lambda}{2\sigma^2}\|\beta\|_1\right)$$

- ▶ The posterior distribution is not a standard distribution.
- ▶ We usually do not have a closed form for the posterior mode.
- ▶ The posterior mode can force some coefficients to be exactly zero, resulting in a sparse model.
- ▶ The sparsity is due to the non-differentiability of the prior at 0.
- ▶ Or, the sub-derivative of the prior at 0 contains a neighborhood of 0.

## Regularization — Spike-and-Slab

Besides the Ridge and LASSO, which “encourage” coefficients to be small through the prior, we may also consider the Spike-and-Slab prior that directly set a probability for the coefficient to be zero.

## Regularization — Spike-and-Slab

Besides the Ridge and LASSO, which “encourage” coefficients to be small through the prior, we may also consider the Spike-and-Slab prior that directly set a probability for the coefficient to be zero.

Specifically, for each coefficient  $\beta_j$ , we set a prior as

$$p(\beta_j) = \theta \underbrace{\delta(\beta_j)}_{\text{spike}} + (1 - \theta) \underbrace{p_{\text{slab}}(\beta_j)}_{\text{slab}},$$

- ▶ The prior is a mixture of a point mass at 0 and a continuous distribution.
- ▶  $\delta(\beta_j)$  is the Dirac delta function at 0 corresponding to the “spike” component.
- ▶  $p_{\text{slab}}(\beta_j)$  is the continuous distribution corresponding to the “slab” component.  
 $p_{\text{slab}}$  can be chosen as uniform, Gaussian, etc..
- ▶  $\theta$  is the probability of sparsity that controls the mixture rate between the two components.

## Regularization — Spike-and-Slab

In practice, several modifications can be used to make inference with the Spike-and-Slab prior:



## Regularization — Spike-and-Slab

In practice, several modifications can be used to make inference with the Spike-and-Slab prior:

- ▶ It is often more convenient to introduce a binary variable  $z_j$  such that

$$z_j \sim \text{Bernoulli}(\theta),$$

$$\beta_j \mid z_j = 1 \sim \delta_0,$$

$$\beta_j \mid z_j = 0 \sim p_{slab}.$$

## Regularization — Spike-and-Slab

In practice, several modifications can be used to make inference with the Spike-and-Slab prior:

- ▶ It is often more convenient to introduce a binary variable  $z_j$  such that

$$\begin{aligned}z_j &\sim \text{Bernoulli}(\theta), \\ \beta_j \mid z_j = 1 &\sim \delta_0, \\ \beta_j \mid z_j = 0 &\sim p_{slab}.\end{aligned}$$

- ▶ It is often more convenient to set the spike component as a Gaussian distribution with a very small variance, and the slab component as a Gaussian distribution with a larger variance.

## Regularization — Spike-and-Slab

In practice, several modifications can be used to make inference with the Spike-and-Slab prior:

- ▶ It is often more convenient to introduce a binary variable  $z_j$  such that

$$\begin{aligned}z_j &\sim \text{Bernoulli}(\theta), \\ \beta_j \mid z_j = 1 &\sim \delta_0, \\ \beta_j \mid z_j = 0 &\sim p_{slab}.\end{aligned}$$

- ▶ It is often more convenient to set the spike component as a Gaussian distribution with a very small variance, and the slab component as a Gaussian distribution with a larger variance.
- ▶ Sampling from the posterior distribution is often done by Gibbs sampling for  $(\beta, z)$ .

## Hierarchical Linear Models

If we have linear regression models for different subsets of data, we may assume that the regression coefficients are drawn from a common distribution.

$$\mathbf{y}_i = \mathbf{X}_i \boldsymbol{\beta}_i + \boldsymbol{\epsilon}_i$$

with

$$\boldsymbol{\beta}_i \sim P, i.i.d.$$

where  $P$  is common distribution for the linear regression coefficients.

# Hierarchical Linear Models

If we have linear regression models for different subsets of data, we may assume that the regression coefficients are drawn from a common distribution.

$$\mathbf{y}_i = \mathbf{X}_i\boldsymbol{\beta}_i + \boldsymbol{\epsilon}_i$$

with

$$\boldsymbol{\beta}_i \sim P, i.i.d.$$

where  $P$  is common distribution for the linear regression coefficients.

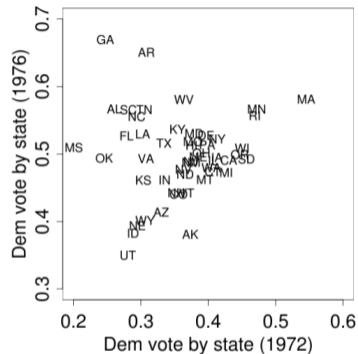
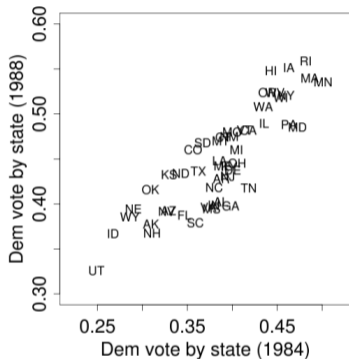
- ▶ When  $P$  is Gaussian, the model is also called a random effects model.
- ▶ Sometimes, only part of the  $\boldsymbol{\beta}_i$  are random effects, and the rest are fixed effects (same for all groups).
- ▶ If the random effects in above are normal, the model is also called a mixed effects model.

## Example: U.S. presidential elections

The data contains results from the U.S. presidential elections for all states from 1948 to 1988.

- ▶ 511 records by removing the District of Columbia and all third-party victories.
- ▶ The response variable is the percentage of votes for the Democratic party.

## Example: U.S. presidential elections



- ▶ Previous election results have a strong effect on the current election results.
- ▶ Some outliers from the southern states. (Upper left on the second graph)

# Example: U.S. presidential elections

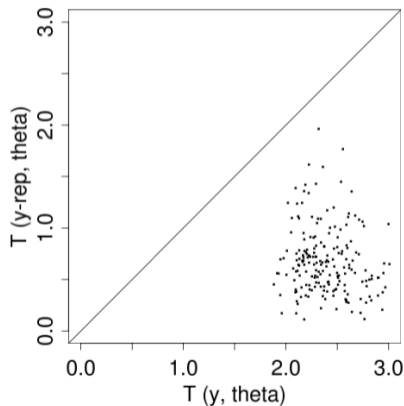
All covariates used for linear regression:

Description of variable	Sample quantiles		
	min	median	max
<b>Nationwide variables:</b>			
Support for Dem. candidate in Sept. poll	0.37	0.46	0.69
(Presidential approval in July poll) $\times$ Inc	-0.69	-0.47	0.74
(Presidential approval in July poll) $\times$ Presinc	-0.69	0	0.74
(2nd quarter GNP growth) $\times$ Inc	-0.024	-0.005	0.018
<b>Statewide variables:</b>			
Dem. share of state vote in last election	-0.23	-0.02	0.41
Dem. share of state vote two elections ago	-0.48	-0.02	0.41
Home states of presidential candidates	-1	0	1
Home states of vice-presidential candidates	-1	0	1
Democratic majority in the state legislature	-0.49	0.07	0.50
(State economic growth in past year) $\times$ Inc	-0.22	-0.00	0.26
Measure of state ideology	-0.78	-0.02	0.69
Ideological compatibility with candidates	-0.32	-0.05	0.32
Proportion Catholic in 1960 (compared to U.S. avg.)	-0.21	0	0.38
<b>Regional/subregional variables:</b>			
South	0	0	1
(South in 1964) $\times$ (-1)	-1	0	0
(Deep South in 1964) $\times$ (-1)	-1	0	0
New England in 1964	0	0	1
North Central in 1972	0	0	1
(West in 1976) $\times$ (-1)	-1	0	0



## Example: U.S. presidential elections

We compare the values of the test variable  $T(\mathbf{y}, \boldsymbol{\theta})$  from the posterior simulations of  $\beta$  to the hypothetical replicated values under the model,  $T(\mathbf{y}^{(rep)}, \boldsymbol{\theta})$ .



The performance is not satisfactory.

## Example: U.S. presidential elections

Now we consider a hierarchical model for the data.

$$y_{st} \sim \mathcal{N}(X_{st}\boldsymbol{\beta} + \gamma_{r(s)t} + \delta_t, \sigma^2),$$
$$\gamma_{rt} \sim \begin{cases} \mathcal{N}(0, \tau_{\gamma_1}^2) & \text{for } r = 1, 2, 3 \text{ (non-south)} \\ \mathcal{N}(0, \tau_{\gamma_2}^2) & \text{for } r = 4 \text{ (south)} \end{cases}$$
$$\delta_t \sim \mathcal{N}(0, \tau_{\delta}^2)$$

- ▶  $\gamma_{rt}$ : different intercepts for different regions.
- ▶  $\delta_t$ : different intercepts for different years.
- ▶  $\boldsymbol{\beta}$  dependence on other covariates is assumed to be the same for all regions and years.
- ▶ Hyperprior for the hyperparameters are set to uniform.

## Example: U.S. presidential elections

We conduct the Bayesian predictive checks for the hierarchical model.

