

## LECTURE TWELVE

Model:

$$y_t = \beta_0 + \beta_1 \text{ReLU}(t-1) + \beta_2 \text{ReLU}(t-2) + \dots + \beta_{n-1} \text{ReLU}(t-(n-1)) + \epsilon_t$$

$t = 1, \dots, n$

→ high-dimension linear regression

$$Y = X\beta + \epsilon$$

$$X = \begin{bmatrix} 1 & 0 & 0 & & 0 \\ 1 & 1 & 0 & & 0 \\ 2 & & 1 & & \vdots \\ \vdots & & \vdots & & 0 \\ 1 & n-1 & n-2 & & 1 \end{bmatrix}$$

$n \times n$

① Least Squares → Overfitting  
(fitted values = data)

$$\hat{\beta} = (X^T X)^{-1} X^T Y$$

$$X \hat{\beta} = Y \rightarrow \text{overfitting}$$

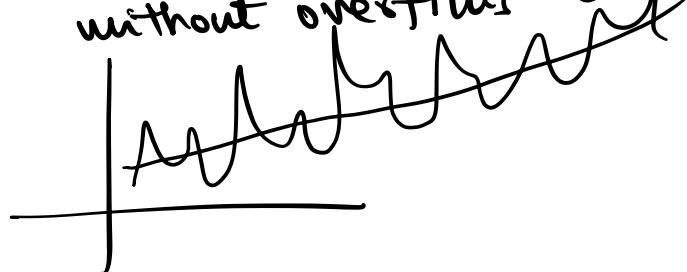
② Regularization

$$\text{Ridge: } \|Y - X\beta\|^2 + \lambda \left[ \sum_{j=2}^{n-1} \beta_j^2 \right]$$

$$\text{LASSO: } \|Y - X\beta\|^2 + \lambda \sum_{j=2}^{n-1} |\beta_j|$$

Why Regularize?

① Want smoother fits } without overfitting } Our preference



② Regularization leads to improved prediction.  
 (Underlies CV for  $\lambda$ -selection)  
 Cross Validation

---

Formula for the Ridge Estimator

$$\begin{aligned}
 & \|y - X\beta\|^2 + \lambda \sum_{j=2}^{n-1} \beta_j^2 \\
 & \nabla_{\beta} \left[ (\underbrace{y - X\beta}_{}^T (\underbrace{y - X\beta}_{})) - \lambda \underbrace{\beta^T X^T y}_{} + \underbrace{y^T y}_{} + \lambda \sum_{j=2}^{n-1} \beta_j^2 \right] \\
 & = \cancel{\lambda} \underbrace{\beta^T X^T X \beta}_{} - \cancel{\lambda} \underbrace{X^T y}_{} + \cancel{\lambda} \underbrace{2}_{=} \begin{pmatrix} 0 \\ 0 \\ \vdots \\ \beta_2 \\ \vdots \\ \beta_{n-1} \end{pmatrix} = 0 \\
 & \begin{pmatrix} 0 \\ 0 \\ \vdots \\ \beta_2 \\ \vdots \\ \beta_{n-1} \end{pmatrix} = J\beta \quad \text{where} \\
 & J = \begin{bmatrix} 0 & & & & 0 \\ & 0 & \ddots & & \\ & & \ddots & \ddots & \\ & & & 0 & \ddots \end{bmatrix}
 \end{aligned}$$

$$X^T X \beta - X^T y + \lambda J \beta = 0$$

$$(X^T X + \lambda J) \beta = X^T y$$

---

$$\hat{\beta}^{\text{Ridge}} = (X^T X + \lambda J)^{-1} X^T y$$

$\hat{\beta}^{\text{least square}} = (X^T X)^{-1} X^T y$

### Bayesian Regularization

$$y = X\beta + \varepsilon, \quad \varepsilon_i \stackrel{iid}{\sim} N(0, \sigma^2)$$

likelihood:  $\left(\frac{1}{\sqrt{2\pi}\sigma}\right)^n \exp\left[-\frac{\|y - X\beta\|^2}{2\sigma^2}\right]$

$$\beta (\beta_0, \beta_1, \dots, \beta_{n-1}) \propto \sigma$$

①  $\beta_0, \beta_1, \dots, \beta_{n-1} \stackrel{iid}{\sim} \text{Unif}(-C, C)$

$C \rightarrow \infty$

↓ when  $C \rightarrow \infty$

posterior  $\beta | \text{data}, \sigma \sim N((X^T X)^{-1} X^T y, \sigma^2 (X^T X)^{-1})$

↓  
n-variate normal

Homework,

②  $\beta_0, \beta_1, \dots, \beta_{n-1} \stackrel{iid}{\sim} N(0, C)$

$(C \text{ large})$

$\beta | \text{data}, \sigma$   $\downarrow$  posterior

$\frac{1}{\sqrt{2\pi} C} \exp\left[-\frac{\beta_i^2}{2C}\right]$ 

behaves like a constant

$$\sim N\left(\left(\frac{X^T X}{\sigma^2} + \frac{I}{C}\right)^{-1} \frac{X^T y}{\sigma^2}, \left(\frac{X^T X}{\sigma^2} + \frac{I}{C}\right)^{-1}\right)$$

Valid for every  $C$   
(not just too large  $C$ )

for most  $\beta_j$   
very similar  
to  $\text{Unif}(-C, C)$

$$\boxed{\frac{I}{C} \{ -C < \beta_j < C \}} \\ 2C$$

$$\text{If } C \rightarrow \infty : N\left(\left(\frac{X^T X}{\sigma^2}\right)^{-1} \frac{X^T y}{\sigma^2}, \left(\frac{X^T X}{\sigma^2}\right)^{-1}\right) \\ = N\left(\left(\frac{X^T X}{\sigma^2}\right)^{-1} X^T y, \sigma^2 \left(\frac{X^T X}{\sigma^2}\right)^{-1}\right)$$

Coincides with posterior for  
 $\text{Unif}(-\infty, \infty)$  prior.

$$\frac{1}{\sqrt{2\pi} C} \exp\left(-\frac{\beta_j^2}{2C}\right)$$

$$C = 10^{10}$$

$$\beta_j \in (-10, 10)$$

### Summarize

- ① If you are using  $C$  large (either in  $\text{Unif}(-C, C)$  or  $N(0, C)$ )  
the posterior mean = least squares  
will lead to overfitting.

- ③  $\beta_0, \beta_1, \dots, \beta_{n-1} \stackrel{iid}{\sim} N(0, C)$ ,  $\beta_n \sim N(0, \tau^2)$   
 $C$ : large for some small  $\tau > 0$

This will lead to smooth fits without overfitting (depending on the chosen value of  $\tau$ )

The formula for the posterior:

$$\beta \sim N_n \left( \begin{pmatrix} 0 \\ 0 \\ \vdots \\ 0 \end{pmatrix}, \underbrace{\begin{pmatrix} C & & & 0 \\ & C\tau^2 & & \\ & & \ddots & \\ & & & C\tau^2 \end{pmatrix}}_Q \right)$$

$\beta \sim N(0, Q)$

Check: posterior now becomes

$$N \left( \left( \frac{X^T X}{\sigma^2} + Q^{-1} \right)^{-1} \frac{X Y}{\sigma^2}, \left( \frac{X^T X}{\sigma^2} + Q^{-1} \right)^{-1} \right)$$

$\beta | \text{data}$  posterior mean:

$$\left( \frac{X^T X}{\sigma^2} + Q^{-1} \right)^{-1} \frac{X Y}{\sigma^2}$$

$$Q^{-1} = \begin{pmatrix} 1/C & & & 0 \\ & 1/C & & \\ & & 1/\tau^2 & \\ 0 & & & 1/\tau^2 \end{pmatrix}$$

$$\approx \begin{pmatrix} 0 & 0 & 0 \\ 0 & 0 & 0 \\ 0 & 0 & 1/\tau^2 \end{pmatrix} = \frac{1}{\tau^2} J$$

$$\text{posterior mean} : \left( \frac{\mathbf{X}^T \mathbf{X}}{\sigma^2} + \frac{1}{\tau^2} \mathbf{J} \right)^{-1} \frac{\mathbf{X}^T \mathbf{y}}{\sigma^2} = \left( \mathbf{X}^T \mathbf{X} + \frac{\sigma^2}{\tau^2} \mathbf{J} \right)^{-1} \mathbf{X}^T \mathbf{y}$$

$$\text{Ridge} : \left( \mathbf{X}^T \mathbf{X} + \lambda \mathbf{J} \right)^{-1} \mathbf{X}^T \mathbf{y}$$

Coincide if

$$\lambda = \frac{\sigma^2}{\tau^2}$$

$$\beta_j \stackrel{iid}{\sim} N(0, C) \quad (\text{or } \text{Unif}(-C, C))$$

$C \text{ large}$

Bayesian = Frequentist

$$(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}$$

$$\beta_0, \beta_1 \stackrel{iid}{\sim} N(0, C), \quad \beta_j \stackrel{iid}{\sim} N(0, \tau^2) \quad j=2, \dots, n-1$$

Frequentist  
with Ridge  
Regularization = Bayesian

Prior 1  $\rightarrow \beta_j \stackrel{iid}{\sim} N(0, C) \rightarrow$  Least squares

Prior 2  $\rightarrow \beta_0, \beta_1 \stackrel{iid}{\sim} N(0, C), \quad \beta_j \stackrel{iid}{\sim} N(0, \tau^2)$   $\rightarrow$  Ridge  
Regularization  
if  $\tau^2$  is small

Prior 2 is restrictive.

Prior 3:  $\beta_0, \beta_i \sim N(0, C)$ ,  $\beta_j \sim N(0, \tau^2)$   
 $\log \tau \sim \text{Unif}(-C, C)$ ,  $\log \sigma \sim \text{Unif}(0, \infty)$

More General

prior:

$$f(\beta, \sigma, \tau) \propto \frac{1}{\sqrt{\det Q}} \exp\left(-\frac{\beta^T Q^{-1} \beta}{2}\right) I\{\tau < \log \sigma\}$$

$$\beta \sim N\left(\begin{pmatrix} 0 \\ \vdots \\ 0 \end{pmatrix}, \begin{pmatrix} C & & \\ & \ddots & \\ & & \tau^2 \end{pmatrix}\right)$$

$$N(0, Q)$$

$$= \frac{\tau^{-1} \sigma^{-1}}{\sqrt{\det Q}} \exp\left(-\frac{\beta^T Q^{-1} \beta}{2}\right)$$

Likelihood:

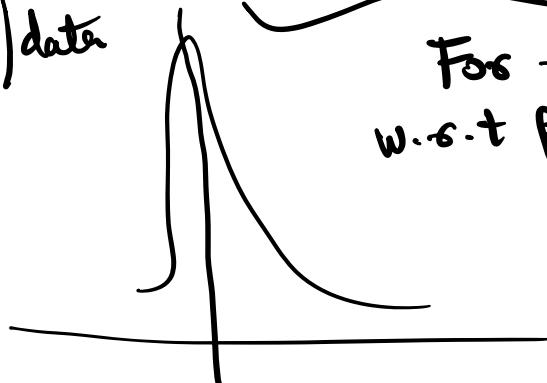
$$\sigma^{-n} \exp\left(-\frac{\|y - X\beta\|^2}{2\sigma^2}\right)$$

Posterior:

$$f(\beta, \tau, \sigma) \propto \frac{\sigma^{-n-1} \tau^{-1}}{\sqrt{\det Q}} \exp\left[-\frac{1}{2} \left( \frac{\|y - X\beta\|^2}{\sigma^2} + \beta^T Q^{-1} \beta \right)\right]$$

$$f(\beta, \tau, \sigma) | \text{data}$$

For fixed  $\sigma, \tau$ , integration w.r.t  $\beta$  is tractable.



$$\frac{\|y - X\beta\|^2}{\sigma^2} + \beta^T Q^{-1} \beta$$

$$x^2 - 4x + 6 = (x-2)^2 + 2$$

$$= (\beta - b)^T \left( \frac{X^T X}{\sigma^2} + Q^{-1} \right) (\beta - b) + \frac{y^T y}{\sigma^2} - b^T \left( \frac{X^T X}{\sigma^2} + Q^{-1} \right) b$$

$$b = \left( \frac{X^T X}{\sigma^2} + Q^{-1} \right)^{-1} \frac{X^T y}{\sigma^2}$$

$f(\sigma, \tau)$  & integral of  
 $\sigma, \tau | \text{data}$

$f(\beta, \sigma, \tau)$   
 $\beta, \sigma, \tau | \text{data}$   
 w.r.t  $\beta$

$$\propto \frac{\sigma^{-n-1} \tau^{-1}}{\sqrt{\det Q}} \sqrt{\det \left[ \left( \frac{X^T X}{\sigma^2} + Q^{-1} \right)^{-1} \right]} \exp \left( -\frac{y^T y}{2\sigma^2} \right)$$

$$\exp \left( \frac{y^T X \left( \frac{X^T X}{\sigma^2} + Q^{-1} \right)^{-1} \frac{X^T y}{\sigma^2}}{2\sigma^2} \right)$$

- ① Take a grid of values of  $\sigma$  &  $\tau$
- ② Compute  $f_{\sigma, \tau | \text{data}}$  over the grid
- ③  $\underbrace{\text{posterior mean/mode of } \sigma \& \tau}_{\text{Sample values of } \sigma \& \tau}$
- ④ Sample  $\beta$  from  $\beta | \sigma, \tau$  : normal posterior

## Why no overfitting?

$f_{\text{data} | \beta, \sigma}$  : likelihood  
maximization leads to  
 $\hat{\beta} = \text{unregularized } \hat{\beta} = 0$   
least squares

$$f_{\text{data} | \tau, \sigma} = \int f_{\text{data} | \beta, \sigma} d\beta$$

$f_{\text{data} | \beta, \sigma}$  (date)  $f_{\beta | \tau}$

$f_{\tau, \sigma | \text{data}}$

$$f_{\tau, \sigma | \text{data}} \propto f_{\text{data} | \tau, \sigma} f_{\tau, \sigma}$$

$f_{\text{data} | \tau, \sigma}$

$f_{\text{data} | \beta, \sigma}$