

End-of-Semester Rambling #1  
Ridge Regression

# Ridge Regression

**Numeric Hack** Instead of

$$\hat{\beta} = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{y},$$

use

$$\hat{\beta}_\lambda = (\mathbf{X}^\top \mathbf{X} + \lambda \mathbf{I})^{-1} \mathbf{X}^\top \mathbf{y}.$$

**Question** Clearly makes sense when  $p > n$ , in which case  $\mathbf{X}^\top \mathbf{X}$  is not even invertible, but beneficial sometimes even when  $p < n$ .

Why?

**Remark** Really only makes sense if data are centered and scaled, i.e.,  $\mathbf{1}^\top \mathbf{y} = 0$ ,  $\|\mathbf{y}\| = 1$  and likewise for columns of  $\mathbf{X}$ . (Why?)

# Bias-Variance Analysis

## Bias

$$\mathbb{E}(\hat{\beta}) = \beta \quad [\text{unbiased}], \quad \mathbb{E}(\hat{\beta}_\lambda) \neq \beta \quad [\text{biased}]$$

## Variance

$$\begin{aligned} \text{Var}(\hat{\beta}) &= (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top [\text{Var}(\mathbf{y})] \mathbf{X} (\mathbf{X}^\top \mathbf{X})^{-1} \\ &= \sigma^2 (\mathbf{X}^\top \mathbf{X})^{-1} \end{aligned}$$

$$\begin{aligned} \text{Var}(\hat{\beta}_\lambda) &= (\mathbf{X}^\top \mathbf{X} + \lambda \mathbf{I})^{-1} \mathbf{X}^\top [\text{Var}(\mathbf{y})] \mathbf{X} (\mathbf{X}^\top \mathbf{X} + \lambda \mathbf{I})^{-1} \\ &= \sigma^2 [(\mathbf{X}^\top \mathbf{X} + \lambda \mathbf{I})^{-1} (\mathbf{X}^\top \mathbf{X}) (\mathbf{X}^\top \mathbf{X} + \lambda \mathbf{I})^{-1}] \end{aligned}$$

**Question** Clearly,  $\text{Bias}(\hat{\beta}_\lambda) \geq \text{Bias}(\hat{\beta})$ , but  $\text{Var}(\hat{\beta}_\lambda) \geq \text{Var}(\hat{\beta})$ ?

# Singular Value Decomposition (SVD)

**SVD** Assume  $\mathbf{X}$  is  $n \times p$  with  $n \geq p$ . There exists decomposition

$$\mathbf{X} = \underset{\substack{\uparrow \\ n \times p}}{\mathbf{U}} \overset{\substack{p \times p \\ \downarrow}}{\mathbf{D}} \underset{\substack{\uparrow \\ p \times p}}{\mathbf{V}^\top},$$

where  $\mathbf{U}$ ,  $\mathbf{V}$  orthonormal;  $\mathbf{D}$  diagonal with all  $\mathbf{D}_{ii} \geq 0$ .

## Remark

- (a) Assumption  $n \geq p$  not restrictive — if  $n \leq p$ , simply apply SVD to  $\mathbf{X}^\top$  and “transpose back”.
- (b) Probably the most important/useful result in matrix algebra.

# Bias-Variance Analysis

$$\begin{aligned} \mathbf{X} = \mathbf{U}\mathbf{D}\mathbf{V}^\top &\Rightarrow \mathbf{X}^\top \mathbf{X} = (\mathbf{V}\mathbf{D}\mathbf{U}^\top) \underbrace{(\mathbf{U}\mathbf{D}\mathbf{V}^\top)}_I = \mathbf{V}\mathbf{D}^2\mathbf{V}^\top \\ &\Rightarrow (\mathbf{X}^\top \mathbf{X})^{-1} = \mathbf{V}\mathbf{D}^{-2}\mathbf{V}^\top \\ &\Rightarrow (\mathbf{X}^\top \mathbf{X} + \lambda\mathbf{I})^{-1} = (\mathbf{V}\mathbf{D}^2\mathbf{V}^\top + \lambda\mathbf{V}\mathbf{V}^\top)^{-1} \\ &= \mathbf{V}(\mathbf{D}^2 + \lambda\mathbf{I})^{-1}\mathbf{V}^\top \end{aligned}$$

$$\left. \begin{aligned} \text{Var}(\hat{\boldsymbol{\beta}}) &= \sigma^2 \mathbf{V} \text{diag} \left[ \frac{1}{D_{ii}^2} \right] \mathbf{V}^\top \\ \text{Var}(\hat{\boldsymbol{\beta}}_\lambda) &= \sigma^2 \mathbf{V} \text{diag} \left[ \frac{D_{ii}^2}{(D_{ii}^2 + \lambda)^2} \right] \mathbf{V}^\top \end{aligned} \right\} \frac{D_{ii}^2}{(D_{ii}^2 + \lambda)^2} \leq \frac{1}{D_{ii}^2}$$

$$\text{Var}(\hat{\boldsymbol{\beta}}_\lambda) \leq \text{Var}(\hat{\boldsymbol{\beta}})$$

ridge regression most useful in “high-variance scenarios”

# Multicollinearity

If there exists  $\ell \in \mathbb{R}^p$  such that  $\|\ell\| = 1$  and  $\|\mathbf{X}\ell\| = \varepsilon$ , then

$$\|\mathbf{X}\ell\|^2 = \ell^\top \underbrace{\mathbf{X}^\top \mathbf{X}}_S \ell = \ell^\top \mathbf{S} \ell = \|\mathbf{S}^{1/2} \ell\|^2 = \varepsilon^2$$

and

$$\begin{aligned} \text{Var}(\ell^\top \hat{\beta}) &= \ell^\top \text{Var}(\hat{\beta}) \ell = \sigma^2 \ell^\top (\mathbf{X}^\top \mathbf{X})^{-1} \ell \\ &= \sigma^2 \ell^\top \mathbf{S}^{-1} \ell = \sigma^2 \|\mathbf{S}^{-1/2} \ell\|^2. \end{aligned}$$

But the [Cauchy-Schwarz inequality](#) implies

$$\|\mathbf{S}^{1/2} \ell\|^2 \|\mathbf{S}^{-1/2} \ell\|^2 \geq (\ell^\top \mathbf{S}^{1/2} \mathbf{S}^{-1/2} \ell)^2 = \|\ell\|^4 = 1,$$

so

$$\text{Var}(\ell^\top \hat{\beta}) \geq \frac{\sigma^2}{\varepsilon^2}.$$

End-of-Semester Rambling #2  
Generalized Linear Models



# Normal

$$f(y; \theta, \phi) = \exp \left[ \frac{y\theta - \gamma(\theta)}{\phi} + c(y; \phi) \right]$$

$$f(y) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp \left[ -\frac{(y - \mu)^2}{2\sigma^2} \right]$$

$$= \exp \left[ \underbrace{\frac{y\mu - \mu^2/2}{\sigma^2}}_{\gamma(\cdot)} + \underbrace{\log \frac{1}{\sqrt{2\pi\sigma^2}} - \frac{y^2/2}{\sigma^2}}_{c(y, \cdot)} \right]$$

$\theta$	$\phi$	$\gamma(\cdot)$	$c(y, \cdot)$
$\mu$	$\sigma^2$	$\mu^2/2$	$\log \frac{1}{\sqrt{2\pi\sigma^2}} - \frac{y^2/2}{\sigma^2}$

# Bernoulli

$$f(y; \theta, \phi) = \exp \left[ \frac{y\theta - \gamma(\theta)}{\phi} + c(y; \phi) \right]$$

$$f(y) = p^y(1-p)^{1-y} = \exp [y \log(p) + (1-y) \log(1-p)]$$

$$= \exp \left[ y \log \frac{p}{1-p} + \underbrace{\log(1-p)}_{-\gamma(\cdot)} \right]$$

$\theta$	$\phi$	$\gamma(\cdot)$	$c(y, \cdot)$
$\log[p/(1-p)]$	1	$-\log(1-p)$	0

$$\theta = \log \frac{p}{1-p} \quad \Rightarrow \quad p = \frac{e^\theta}{1+e^\theta} \quad \Rightarrow \quad \gamma(\theta) = \log(1+e^\theta)$$

# Logit vs Sigmoid Transforms

## Logit

$$\theta = \log \frac{p}{1-p} = \begin{cases} +\infty, & p = 1 \\ 0, & p = 1/2 \\ -\infty, & p = 0 \end{cases}$$

## Sigmoid

$$p = \frac{e^\theta}{1 + e^\theta} = \begin{cases} 1, & \theta \rightarrow +\infty \\ 1/2, & \theta = 0 \\ 0, & \theta \rightarrow -\infty \end{cases}$$

# Canonical GLMs

- $y_1, y_2, \dots, y_n \in \mathbb{R} \Rightarrow$  “ordinary” regression

$$y_i \sim N(\mu_i, \sigma^2)$$

where

$$\mu_i = \mathbf{x}_i^\top \boldsymbol{\beta}$$

- $y_1, y_2, \dots, y_n \in \{0, 1\} \Rightarrow$  logistic regression

$$y_i \sim \text{Bernoulli}(p_i)$$

where

$$\log \frac{p_i}{1 - p_i} = \mathbf{x}_i^\top \boldsymbol{\beta} \quad \text{or} \quad p_i = \frac{\exp(\mathbf{x}_i^\top \boldsymbol{\beta})}{1 + \exp(\mathbf{x}_i^\top \boldsymbol{\beta})}$$

# Logistic Regression

- log-likelihood:

$$\sum_{i=1}^n y_i \log \frac{p_i}{1 - p_i} + \log(1 - p_i)$$

$$\log \frac{p_i}{1 - p_i} = \mathbf{x}_i^\top \boldsymbol{\beta} \quad \Rightarrow \quad \Downarrow \quad \Leftarrow \quad p_i = \frac{\exp(\mathbf{x}_i^\top \boldsymbol{\beta})}{1 + \exp(\mathbf{x}_i^\top \boldsymbol{\beta})}$$

$$\ell(\boldsymbol{\beta}; y_i, \mathbf{x}_i) = \sum_{i=1}^n y_i \mathbf{x}_i^\top \boldsymbol{\beta} - \log(1 + e^{\mathbf{x}_i^\top \boldsymbol{\beta}})$$

- MLE:

$$\max_{\boldsymbol{\beta}} \ell(\boldsymbol{\beta}; y_i, \mathbf{x}_i)$$

by [Newton-Raphson](#)

# Newton-Raphson

$$\ell(\boldsymbol{\beta}; y_i, \mathbf{x}_i) = \sum_{i=1}^n y_i \mathbf{x}_i^\top \boldsymbol{\beta} - \log(1 + e^{\mathbf{x}_i^\top \boldsymbol{\beta}})$$

↓

$$\ell'(\boldsymbol{\beta}) = \dots = \sum_{i=1}^n (y_i - p_i) \mathbf{x}_i$$

$$= \begin{bmatrix} \mathbf{x}_1 & \dots & \mathbf{x}_n \end{bmatrix} \begin{bmatrix} y_1 - p_1 \\ \vdots \\ y_n - p_n \end{bmatrix} = \mathbf{X}^\top (\mathbf{y} - \mathbf{p})$$

# Newton-Raphson

$$\begin{aligned} & \ell''(\boldsymbol{\beta}) \\ &= \vdots \\ &= - \begin{bmatrix} \mathbf{x}_1 & \dots & \mathbf{x}_n \end{bmatrix} \begin{bmatrix} p_1(1-p_1) & 0 & 0 \\ 0 & \ddots & 0 \\ 0 & 0 & p_n(1-p_n) \end{bmatrix} \begin{bmatrix} \mathbf{x}_1^\top \\ \vdots \\ \mathbf{x}_n^\top \end{bmatrix} \\ &= -\mathbf{X}^\top \mathbf{W} \mathbf{X} \end{aligned}$$

## Iteratively Reweighted LSQ

$$\begin{aligned}\beta_{new} &= \beta_{old} - [\ell''(\beta_{old})]^{-1} [\ell'(\beta_{old})] \\ &= \beta_{old} + [\mathbf{X}^\top \mathbf{W} \mathbf{X}]^{-1} [\mathbf{X}^\top (\mathbf{y} - \mathbf{p})] \\ &= \underbrace{[\mathbf{X}^\top \mathbf{W} \mathbf{X}]^{-1} \mathbf{X}^\top \mathbf{W}}_{\text{weighted LSQ operator}} \underbrace{[\mathbf{X} \beta_{old} + \mathbf{W}^{-1} (\mathbf{y} - \mathbf{p})]}_{\mathbf{z}; \text{ working response}}\end{aligned}$$

(notice that both  $\mathbf{W}$  and  $\mathbf{p}$  depend on  $\beta_{old}$  as well)

$\mathbf{W}_{ii} \equiv p_i(1 - p_i)$  **max** at  $p_i = 1/2$  and **min** at  $p_i = 0$  or  $p_i = 1$

⇓

estimates influenced mostly by points near decision boundary