

Ridge regression Setup

Setup: y, x_1, \dots, x_p with $p \gg 0$ or even $p > n$

Model: $y = \beta_0 + \beta_1 x_1 + \dots + \beta_p x_p + \varepsilon$ as "usual"

$$\hat{\beta}_0 = \bar{y}, \quad \beta = \begin{pmatrix} \beta_1 \\ \beta_2 \\ \vdots \\ \beta_p \end{pmatrix}_{p \times 1}$$

$$y_i \xrightarrow{\text{new } y} y_i - \bar{y} \xrightarrow{\text{old } y} \left[= y_i - \hat{\beta}_0 \right]$$

$$x_i \xrightarrow{\text{old } x_i} \frac{x_i - \bar{x}}{\sqrt{\frac{1}{n} \sum_{j=1}^n (x_j - \bar{x})^2}} \xrightarrow{\text{new } x_i}$$

$$\Rightarrow \begin{aligned} & \text{new } x_i \\ & x_i = 0 \\ & \sum_{i=1}^n x_i^2 = n \end{aligned}$$

X = design matrix of new features without 1st column of 1s, so is $n \times p$.

Ridge regression

" " estimate of β minimizes

$$\sum_{i=1}^n (y_i - (\beta_1 x_{i1} + \dots + \beta_p x_{ip}))^2 + \lambda \sum_{j=1}^p \beta_j^2, \quad \lambda \geq 0$$

$$= (y - X\beta)^T (y - X\beta) + \lambda \|\beta\|_2^2$$

$$= \dots + \lambda \beta^T \beta$$

same

Ex $p=1$, $y = \beta x_i + \epsilon$. OLS minimizes

$$\sum (y_i - \beta x_i)^2 \Rightarrow \hat{\beta}_{OLS} = \frac{\sum x_i y_i}{\sum x_i^2} = \frac{\sum_{i=1}^n x_i y_i}{n}$$

ridge minimizes

$$\sum (y_i - \beta x_i)^2 + \lambda \beta^2$$

$$\hat{\beta}_{OLS}$$

$$\frac{\partial}{\partial \beta} \Big| \Rightarrow -2 \sum_{i=1}^n x_i (y_i - \beta x_i) + 2\lambda \beta$$

$$= -2 \sum x_i y_i + 2\beta \underbrace{\sum x_i^2}_n + 2\lambda \beta$$

$$= -2 \sum x_i y_i + 2\beta (n + \lambda) \stackrel{\text{Set}}{=} 0$$

$$\Rightarrow \hat{\beta}_{\text{ridge}} = \frac{\sum x_i y_i}{n + \lambda}$$

$\Rightarrow \hat{\beta}_{\text{ridge}}$ is shrinking $\hat{\beta}_{OLS}$ toward 0

Note Ridge regression is a shrinkage estimator in that the estimated β s tend to be shrunk toward 0.

- $\lambda = 0 \Rightarrow \hat{\beta}_{\text{ridge}} = \hat{\beta}_{OLS}$

- $\lambda \rightarrow \infty \Rightarrow \hat{\beta}_{\text{ridge}} = 0 \Rightarrow \text{model reduces to } Y = \beta_0 + \varepsilon$

Solution

$\hat{\beta}_{\text{ridge}}$ minimizes $(y - X\beta)^T (y - X\beta) + \lambda \beta^T \beta$

$$\frac{\partial}{\partial \beta} \mid \Rightarrow -2X^T(y - X\beta) + 2\lambda\beta$$

$$= -2X^T y + 2X^T X\beta + 2\lambda\beta$$

$$= -2X^T y + 2(X^T X + \lambda I)\beta \stackrel{!}{=} 0$$

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

$$\Rightarrow \hat{\beta}_{\text{ridge}} = (X^T X + \lambda I)^{-1} X^T y$$

Aside

$(X^T X)$ is a nonnegative definite matrix (Semi positive definite).

choose any \underline{a}

0 \leq Show $\underline{a}^T (X^T X) \underline{a} = (\underline{a}^T X^T) (X \underline{a}) = (X \underline{a})^T (X \underline{a})$

$[X \underline{a} = \text{some vector} = \underline{b}]$
 $\rightarrow = \underline{b}^T \underline{b} = \sum_{i=1}^p b_i^2 \geq 0. \checkmark$

That implies that $(X^T X)$ has nonnegative eigenvalues.

But if $(X^T X)$ is not full rank, it has some eigenvalues that are exactly 0. $\Rightarrow (X^T X)$ is not

always invertible.

$$\hat{\beta}_{OLS} = \cancel{(X^T X)^{-1}} X^T y$$

The effect of ridge's $+ \lambda I$? Bumps up all eigenvalues of $X^T X$ by λ .

Side: If \underline{v} is eigenvector of $(X^T X)$ with eigenvalue α . Then

$$\begin{aligned}(X^T X + \lambda I) \underline{v} &= (X^T X) \underline{v} + \lambda \underline{v} = \alpha \underline{v} + \lambda \underline{v} \\ &= (\alpha + \lambda) \underline{v}\end{aligned}$$

$\Rightarrow \underline{v}$ is still an eigenvector of $(X^T X + \lambda I)$ but
it has eigenvalue $\alpha + \lambda$

\Rightarrow All zero eigenvalues of $(X^T X)$ are now > 0

$\Rightarrow (X^T X + \lambda I)$ is invertible

\Rightarrow unique solution + allows us to fit
models with more parameters than data.

Properties

$$\begin{aligned}
 \bullet \quad E \hat{\beta}_{\text{ridge}} &= E \left((X^T X + \lambda I)^{-1} X^T y \right) \\
 &= (X^T X + \lambda I)^{-1} X^T E y \\
 &= (X^T X + \lambda I)^{-1} X^T X \beta \neq \beta
 \end{aligned}$$

$\Rightarrow \hat{\beta}_{\text{ridge}}$ is biased

$$\bullet \quad \text{Var } \hat{\beta}_{\text{ridge}} = \sigma^2 (X^T X + \lambda I)^{-1} (X^T X) (X^T X + \lambda I)^{-1}$$

$$\text{Var } \hat{\beta}_{\text{OLS}} = \sigma^2 (X^T X)^{-1}$$

nono

$$\text{Var } \hat{\beta}_{\text{ridge}}$$

$$\text{Var } \hat{\beta}_{\text{OLS}}$$

$$\begin{aligned}
 (a + \lambda)^{-1} a (a + \lambda)^{-1} &= \frac{a}{(a + \lambda)^2} \\
 \frac{1}{a}
 \end{aligned}$$

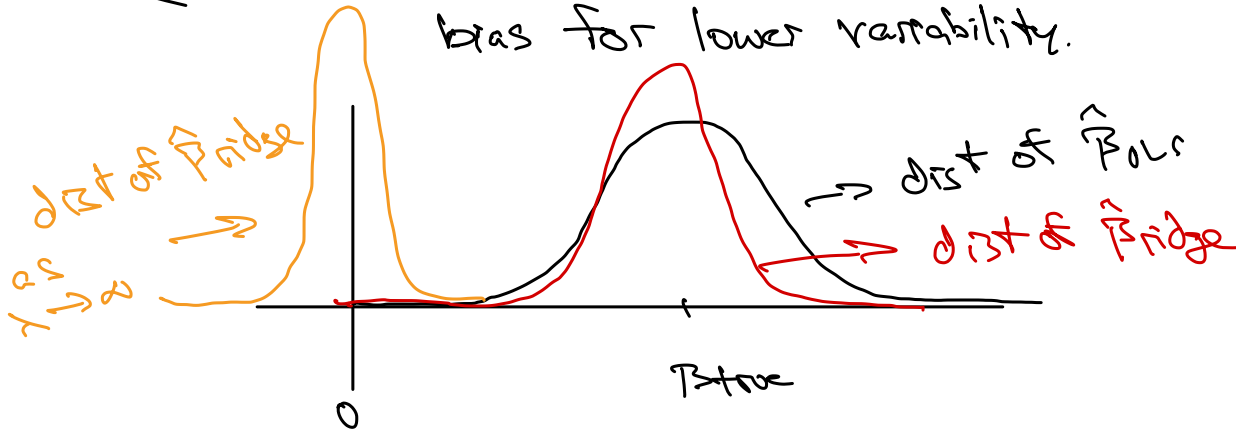
Is Var Bridge \leftarrow Var OLS? \leftarrow

$$\frac{\sigma^2}{(n+1)^2} < \frac{1}{n} \quad -$$

$$\frac{\sigma^2}{(n+1)^2} < 1 \quad \checkmark$$

Take home point

Ridge regression sacrifices some bias for lower variability.



Recall The mean squared error (MSE) is

$$E \left[\sum_{j=1}^p (\beta_j - \hat{\beta}_j)^2 \right] = E \left[(\beta - \hat{\beta})^T (\beta - \hat{\beta}) \right]$$

For ridge:

$$MSE(\lambda) = \beta^T (I - M(\lambda) X)^T (I - M(\lambda) X) \beta + \sigma^2 \{ \text{trace} [M(\lambda)^T M(\lambda)] \}$$

$$\text{where } M(\lambda) = (X^T X + \lambda I)^{-1} X^T$$

For certain choices of λ , $MSE_{\text{ridge}} < MSE_{\text{OLS}}$

$$MSE = (\text{squared bias}) + (\text{var})$$