

Ridge Regression: $\hat{\beta}_\lambda^R = \underset{\tilde{\beta} \in \mathbb{R}^p}{\operatorname{argmin}} \left\{ \underbrace{\sum_{i=1}^n (y_i - X_i^T \tilde{\beta})^2}_{\text{Squared Error}} + \underbrace{\lambda \sum_{j=1}^p \tilde{\beta}_j^2}_{\text{penalty on using "big" values for } \tilde{\beta}} \right\}$

$$\hat{\beta}_\lambda^R = (X^T X + \lambda I_p)^{-1} X^T y$$

λI_p is "stabilizing" the inverse of $X^T X$

- if we have Multicollinearity then $(X^T X)^{-1}$ is unstable leading to a large variance in $\hat{\beta}$
- if $p > n$ then $(X^T X)^{-1}$ Does not exist.

However $(X^T X + \lambda I_p)^{-1}$ always exists for any $\lambda > 0$.

\Rightarrow Reduce variance and stop us from overfitting

Theorem: there exists a $\lambda > 0$, small enough, s.t.

$$\text{MSE}(\hat{\beta}_\lambda^R) < \text{MSE}(\hat{\beta})$$

Ridge < Least Squares

Unbiased!

Proof: • $\text{MSE}(\hat{\beta}) = \text{Var}(\hat{\beta})$

$$= \sigma^2 (X^T X)^{-1}$$

For simplicity, we assume this exists.

$$\begin{aligned} \text{MSE}(\hat{\beta}_\lambda^R) &:= E[(\hat{\beta}_\lambda^R - \beta)(\hat{\beta}_\lambda^R - \beta)^T] \\ &= \text{Var}(\hat{\beta}_\lambda^R) + \text{bias}(\hat{\beta}_\lambda^R) \text{bias}(\hat{\beta}_\lambda^R)^T \end{aligned}$$

$$\begin{aligned} \text{bias}(\hat{\beta}_\lambda^R) &= E\hat{\beta}_\lambda^R - \beta \\ &= E[(X^T X + \lambda I_p)^{-1} X^T y] - \beta \\ &= (X^T X + \lambda I_p)^{-1} X^T E y - \beta \\ &= (X^T X + \lambda I_p)^{-1} X^T X \beta - \beta \\ &= \left[(X^T X + \lambda I_p)^{-1} X^T X - I \right] \beta \\ &= \left[(X^T X + \lambda \underbrace{I_p}_{\text{I}})^{-1} X^T X - (X^T X)^{-1} (X^T X) \right] \beta \\ &= \left[\underbrace{(X^T X + \lambda I_p)^{-1}} - \underbrace{(X^T X)^{-1}} \right] (X^T X) \beta \end{aligned}$$

Here we use an outer product $\hat{\beta} \hat{\beta}^T$ a matrix instead of inner product $\hat{\beta}^T \hat{\beta} \in \mathbb{R}$

Note: $\text{tr}(\hat{\beta} \hat{\beta}^T) = \hat{\beta}^T \hat{\beta}$

Point: Bias is negative wrt β \therefore if $\beta_j > 0$

$$= \left[(\underline{X^T X + \lambda I_p}) - (\underline{X^T X}) \right] (X^T X)^{-1} \beta$$

Point: Bias is negative wrt β . i.e. if $\beta_i > 0$ then the bias $(\hat{\beta}_{\lambda i}^R) < 0$ and if $\beta_i < 0$ then the bias $(\hat{\beta}_{\lambda i}^R) > 0$

\Rightarrow Shrinkage estimator as the bias is always point towards zero.

$$\begin{aligned} \text{Var}(\hat{\beta}_{\lambda}^R) &= \text{Var}\left(\underline{(X^T X + \lambda I_p)^{-1} X^T Y}\right) \\ &= \left[(X^T X + \lambda I_p)^{-1} X^T\right] \text{Var}(Y) \left[(X^T X + \lambda I_p)^{-1} X^T\right]^T \\ &= \left[(X^T X + \lambda I_p)^{-1} X^T\right] (\sigma^2 I_n) \left[\underbrace{X (X^T X + \lambda I_p)^{-1}}_{\text{Symmetric}}\right] \\ &= \sigma^2 (X^T X + \lambda I_p)^{-1} (X^T X) (X^T X + \lambda I_p)^{-1} \end{aligned}$$

$$\begin{aligned} \text{Var}(\hat{\beta}) &= \sigma^2 (\underline{X^T X})^{-1} (\cancel{X^T X}) (\cancel{X^T X})^{-1} \quad \left\{ \begin{array}{l} \text{Multiplied by} \\ (X^T X)^{-1} (X^T X) \text{ and} \\ \text{grouped terms} \\ \text{together} \end{array} \right. \\ &= \sigma^2 (\underline{X^T X})^{-1} \left[(\underline{X^T X}) (\underline{X^T X + \lambda I_p})^{-1} \right]^2 \end{aligned}$$

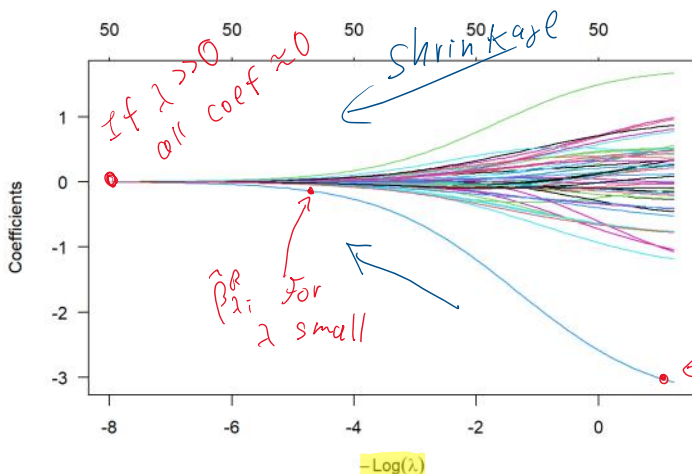
\hookrightarrow Claim: $\lambda < I_p$ then we have reduced the variance.

Simple case: $X = \begin{pmatrix} x_1 \\ \vdots \\ x_n \end{pmatrix}$

$$\text{Now, } (X^T X) (X^T X + \lambda I_p)^{-1} = \frac{\sum x_i^2}{\sum x_i^2 + \lambda} < 1$$

```
# Run a ridge regression
md.ridge = glmnet(x, y, alpha=0)
plot(md.ridge, las=1)
```

Output From
"glm net"



Each line corresponds to a $\hat{\beta}_{\lambda i}^R$

\leftarrow If $\lambda \approx 0$ then $\tilde{\beta}_{\lambda}^R \approx \tilde{\beta}$

Interpolating between $\tilde{\beta}$, 0, 1, and 0

Next time: go From Ridge to Lasso estimator

Next time: go From Ridge to Lasso estimator

$$\lambda \sum \tilde{\beta}_i^2$$

$$\lambda \sum |\beta_i|$$

Lasso = Least Absolute Shrinkage + Selection Operator

$$|\beta|$$

$$\downarrow 0$$

$$\text{set} = 0$$

↳ Do variable selection, but still maintains a convex optimization problem.