# 30. Lasso Regression
November 24, 2025    9:42 AM

Last time (1970): Ridge Regression

Solve: $\hat{\beta}_\lambda^R = \underset{\tilde{\beta} \in \mathbb{R}^P}{\arg\min} \left\{ \sum_{i=1}^{n} (Y_i - X_i^T \tilde{\beta})^2 + \lambda \sum_{j=1}^{P} \tilde{\beta}_j^2 \right\}$

<span style="color:red">Squared Error</span>   <span style="color:red">Quadratic Penalty</span>

- Closed form solution

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

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

- Shrinkage Estimator, send the estimates towards zero
  $\rightarrow$ Variance $\downarrow$   bias $\uparrow$

- Sometimes called $L^2$-regularization
- It <u>Doesn't</u> do variable selection $\rightarrow$ <span style="color:red">Hard to interpret</span>

Lasso (1996) $\Rightarrow$ Least absolute Selection + Shrinkage Operator

$$\hat{\beta}_\lambda^L = \underset{\tilde{\beta} \in \mathbb{R}^P}{\arg\min} \left\{ \sum_{i=1}^{n} (Y_i - X_i^T \tilde{\beta})^2 + \lambda \sum_{j=1}^{P} |\tilde{\beta}_j| \right\}$$

Squared error    Absolute/linear penalty

<span style="color:red">( also called $L^1$-regularization )</span>

- Why is it interesting?
  $\rightarrow$ Does Both Shrinkage + Selection at same time
  $\rightarrow$ Convex optimization Problem <span style="color:red">(Easier to solve than non-convex)</span>
  $\rightarrow$ No closed form solution.

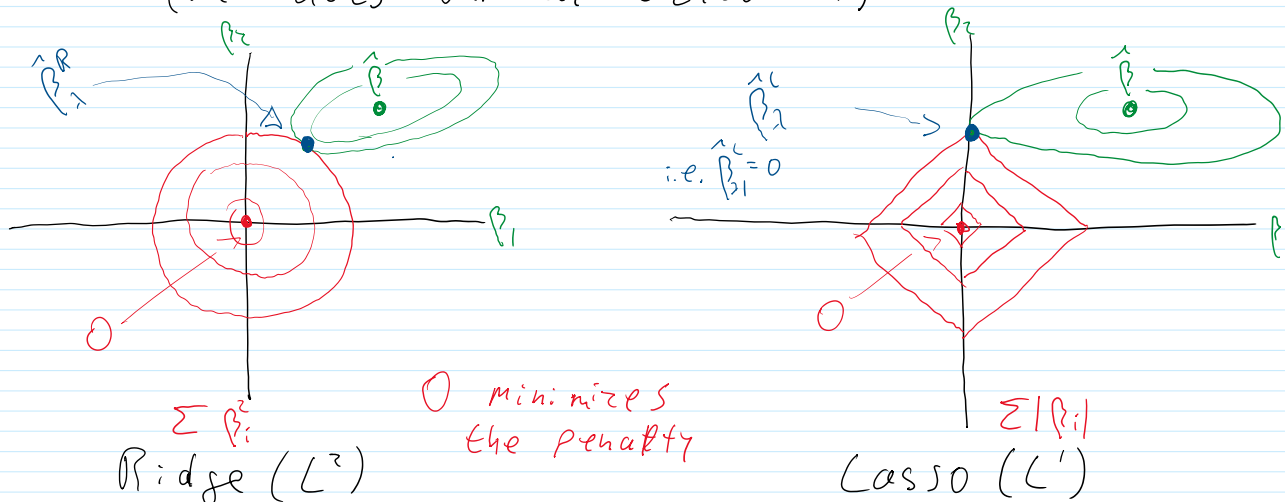  <span style="color:green">* OLS and Ridge are unique in that Most Statistical Machine Learning Methods do not have a nice closed form solution $\Rightarrow$ i.e. Need a Computer!</span>
  $\rightarrow$ This work spawned countless

footersegment type="footer_navigation">STAT 378^J 2025 Page 1

→ This work spawned countless
   Follow up papers on Theory.

Picture version of why Lasso sets some $\hat{\beta}^L_{\lambda j} = 0$
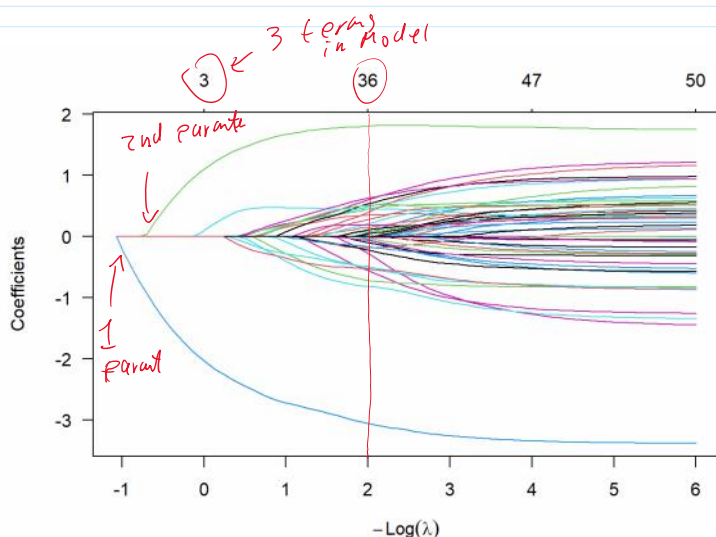   (i.e. does variable selection)



$\sum \beta_i^2$

Ridge ($L^2$)

$O$ minimizes the penalty

$\sum |\beta_i|$

Lasso ($L^1$)

$\hat{\beta}$ Least squares so it is the "best" point
   to Minimize $\sum (Y_i - X_i^T \tilde{\beta})^2$

Elastic Net (2005) : Zou + Hastie
   why not do both Lasso + Ridge!

$$\hat{\beta}^{EN}_{\lambda} = \arg\min_{\tilde{\beta} \in \mathbb{R}^p} \left\{ \sum_{i=1}^n (Y_i - X_i^T \tilde{\beta})^2 + \lambda_1 \sum_{j=1}^p |\tilde{\beta}_j| + \lambda_2 \sum_{j=1}^p \tilde{\beta}_j^2 \right\}$$

Squared Error     $L^1$     $L^2$

penalty

Claim: "Better" than Lasso but still does Shrinkage + Selection



3 terms in Model

2nd param

1 param

"Lasso Paths"

• Each line in a parameter
  that "enters" the Model
  as $\lambda \downarrow$ or $-\log \lambda \uparrow$

• Every choice of $\lambda$
  (cross section of plot)
  is a Lasso regression
  Model.

$\hat{\beta}^R_\lambda = (X^T X + \lambda I)^{-1} X^T Y$     or     $\dfrac{\sum x_i y_i}{\sum x_i^2}$   or   $\dfrac{\sum x_i y_i}{\lambda + \sum x_i^2}$

$$\hat{\beta}_\lambda^R = \left(X^T X + \lambda I\right)^{-1} X^T Y \quad , \quad e.g. \quad \frac{\sum x_i y_i}{\sum x_i^2} \quad \text{or} \quad \frac{\sum x_i y_i}{\sum x_i^2 + \lambda}$$

<span style="color:red">↰ like adding $\lambda$ to each term in the denominator</span>

People have proven "Sparsistency" → sparse + consistency

    if the true $\beta$ has, say, $\beta_i = 0$

    then      $\hat{\beta}_{\lambda\,i}$ will $= 0$ for large enough sample $n$.