

## SHRINKAGE METHODS

29/03/21

They reduce prediction error by modifying the LS criterion by constraining / shrinking the coefficients

Two main approaches : ridge regression and lasso

RIDGE REGRESSION (1970 , Tikhonov regularization)

$$\hat{\beta}_{LS} = \underset{\beta}{\operatorname{argmin}} \quad (\underline{Y} - \underline{X}\beta)^t (\underline{Y} - \underline{X}\beta) \quad \rightsquigarrow \boxed{\hat{\beta}_{LS} = (\underline{X}^t \underline{X})^{-1} \underline{X}^t \underline{Y}}$$

$$\hat{\beta}_{RIDGE} = \underset{\beta}{\operatorname{argmin}} \quad (\underline{Y} - \underline{X}\beta)^t (\underline{Y} - \underline{X}\beta) \quad \text{subject to} \quad \sum_{j=1}^p \beta_j^2 \leq c \quad \text{for } c \geq 0$$

CONSTRAINED  
OPTIMIZATION

## RIDGE ESTIMATOR

$$\min_{\beta} (\gamma - x\beta)^t (\gamma - x\beta) \quad \text{st} \quad \sum_{j=1}^p \beta_j^2 \leq c \quad (*)$$

$$\min_{\substack{\beta \\ \sum \beta_j^2 \leq c}} (\gamma - x\beta)^t (\gamma - x\beta) = \min_{\beta} \max_{\lambda \geq 0} \underbrace{(\gamma - x\beta)^t (\gamma - x\beta)}_{\text{Lagrange multiplier}} + \lambda \underbrace{\left( \sum_{j=1}^p \beta_j^2 - c \right)}_{\text{Lagrangian}}$$

$$\hat{\beta}_{\text{Ridge}} = \arg \min_{\beta} (\gamma - x\beta)^t (\gamma - x\beta) + \lambda \underbrace{\sum_{j=1}^p \beta_j^2}_{\text{PENALTY}} \quad \text{for } \lambda \geq 0 \quad (***) \quad \text{PENALIZED INFERENCE}$$

$$\frac{d}{d\beta} \left( (\gamma - x\beta)^t (\gamma - x\beta) + \lambda \sum_{j=1}^p \beta_j^2 \right) = \\ = -2x^t \gamma + 2x^t x \beta + 2\lambda \beta = 0$$

$$\Rightarrow x^t \gamma - x^t x \beta - \lambda \beta = 0$$

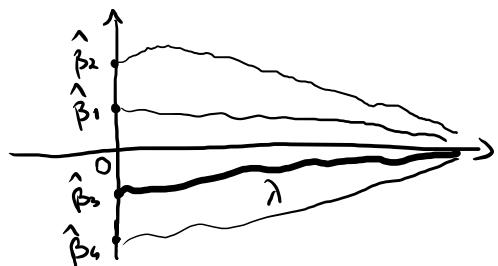
$$\Rightarrow x^t \gamma = (x^t x + \lambda I) \beta \Rightarrow \hat{\beta}_{\text{Ridge}}^{(1)} = (x^t x + \lambda I)^{-1} x^t \gamma$$

## THE ROLE OF $\lambda$ ( $\geq 0$ )

- ①  $\lambda = 0$  ( $c = +\infty$ )  $\rightarrow$  LS (unconstrained) solution
- ② As  $\lambda$  increases ( $c$  decreases), the impact of the shrinkage penalty increases and the ridge coefficients will approach zero
- ③  $\lambda$  stabilizes estimation for big data since it adds  $\lambda$  to the diagonal of  $X^t X$ . So  $X^t X + \lambda I$  is now invertible also for  $n < p$  (big data)

## Remarks

- ① We have an estimate of  $\hat{\beta}$  for each  $\lambda$ , so results are shown for the full path of solutions



$\sum_{j=1}^4 \hat{\beta}_j^2$  decreases with increasing  $\lambda$

- ②  $\lambda$  controls the bias/variance trade off

$\lambda = 0$  OLS (no bias  
high variance)

As  $\lambda \rightarrow \infty$ , the variance of ridge regression decreases and bias increases

## BIAS / VARIANCE OF RIDGE ESTIMATOR

$$\text{BIAS } \hat{\beta}_{\text{RIDGE}} = (X^T X + \lambda I)^{-1} X^T Y$$

$$= (R + \lambda I)^{-1} X^T Y$$

let  $R = X^T X$

$$= (R + \lambda I)^{-1} R R^{-1} X^T Y$$

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

$$= (R(I + \lambda R^{-1}))^{-1} R \hat{\beta}_{LS}$$

$$(AB)^T = B^T A^T$$

$$= (I + \lambda R^{-1})^{-1} \hat{\beta}_{LS}$$

$$E[\hat{\beta}_{\text{RIDGE}}] = (I + \lambda R^{-1})^{-1} E[\hat{\beta}_{LS}] \stackrel{I}{=} \underbrace{(I + \lambda R^{-1})^{-1} I}_{\text{unless } \lambda=0} \beta \neq \beta$$

LS unbiased

BIASED  
with bias  
increasing  
with  $\lambda$

## VARIANCE

$$\hat{\beta}_{\text{RIDGE}} = (I + \lambda R^{-1})^{-1} \hat{\beta}_{\text{LS}}$$

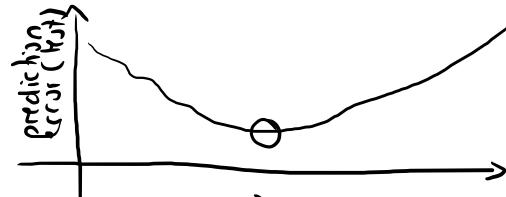
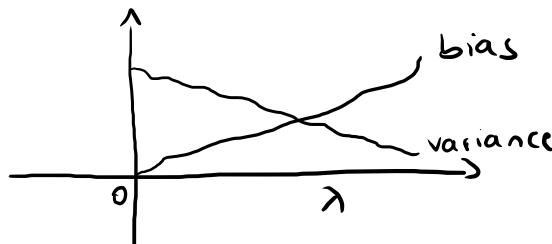
$$\text{Var}(\hat{\beta}_{\text{RIDGE}}) = (I + \lambda R^{-1})^{-1} \underbrace{\text{Var}(\hat{\beta}_{\text{LS}})}_{\sigma^2(X^T X)^{-1}} (I + \lambda R^{-1})^{-1}$$

$$\text{Var}(A Y)$$

$$= A \text{Var}(Y) A^T$$

$$= \sigma^2 (I + \lambda R^{-1})^{-1} (X^T X)^{-1} (I + \lambda R^{-1})^{-1}$$

Decreases to 0 as  $\lambda \rightarrow \infty$



$\lambda$  can be chosen via cross-validation

Fix  $\lambda$   
 Split data in K-folds  
 $\hat{\beta}_{\text{RIDGE}}$  estimated from  $K-1$  folds  
 and make predictions on test set

## SCALING OF DATA

The standard LS estimates are not affected by scaling: if we multiply  $X_j$  by a constant  $c$ ,  $\hat{\beta}_j$  will be scaled by a factor of  $\frac{1}{c}$ , so regardless of how  $X_j$  is scaled,  $X_j \hat{\beta}_j$  remains the same.

This does not happen for ridge regression. This is due to  $\sum_{j=1}^p \beta_j^2$ .

For this reason, ridge regression is applied after standardizing the predictors (to mean 0 and variance 1).

As for the intercept, this is typically estimated beforehand ( $\hat{\beta}_0 = \bar{y}$ , unpenalized). The response variable is then centered to have mean  $x_{r0}$  and  $\beta_1, \dots, \beta_p$  are estimated by ridge regression (to standardized data).

# LASSO

Least absolute shrinkage and selection operator

Reduces variance and performs automatic variable selection

$$\hat{\beta}_{\text{LASSO}} = \underset{\beta}{\operatorname{argmin}} \quad (y - X\beta)^t (y - X\beta) + \lambda \sum_{j=1}^p |\beta_j| \quad \text{with } \lambda \geq 0$$

As before, this is equivalent to

$$\begin{aligned} & \min (y - X\beta)^t (y - X\beta) \\ \text{s.t. } & \sum_{j=1}^p |\beta_j| \leq c \end{aligned} \quad ]$$

Compare this with subset selection:

$$\begin{aligned} & \min (y - X\beta)^t (y - X\beta) \\ \text{s.t. } & \sum_{j=1}^p I(\beta_j \neq 0) \leq c \end{aligned} \quad ]$$

## ROLE OF $\lambda$

- ①  $\lambda = 0 \Leftrightarrow LS$
- ② As  $\lambda$  increases ( $c$  decreases), more and more coefficients are shrunk to zero
- ③ One  $\hat{\beta}$  class for each  $\lambda$ , so you can look at the path of solution, but no closed-form analytical solution
- ④ Choice of  $\lambda$  is crucial for bias/variance trade-off.  
So use CV to choose optimal value of  $\lambda$ .

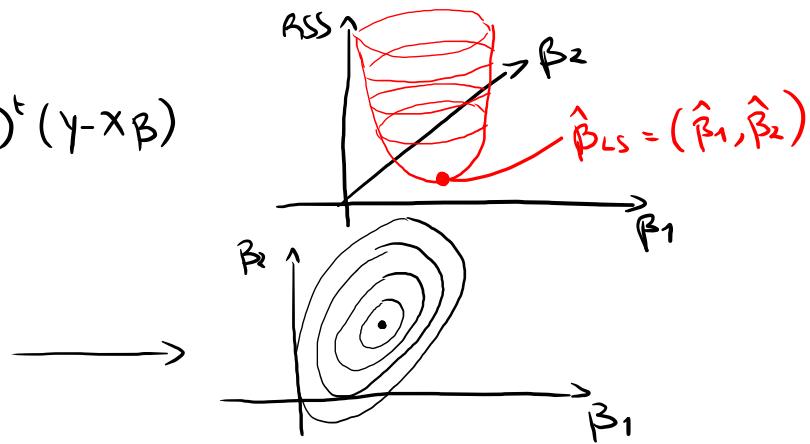
## WHY LASSO ESTIMATES CAN BE EXACTLY ZERO

We say that lasso yields sparse models, that is the models have only a subset of predictors.

In order to see this, let us take two predictors, so we wish to estimate  $\beta_1$  and  $\beta_2$ .

PLOT OF

$$RSS = (\mathbf{y} - \mathbf{x}\boldsymbol{\beta})^t (\mathbf{y} - \mathbf{x}\boldsymbol{\beta})$$



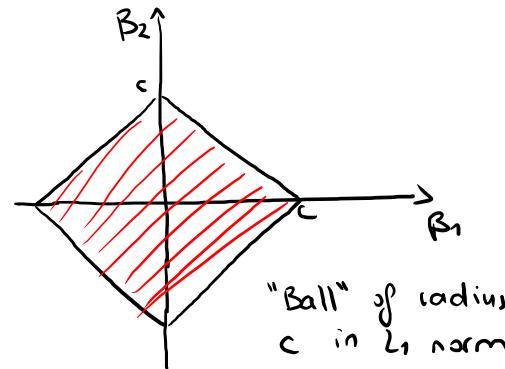
## CONSTRAINED OPTIMIZATION

### LASSO

$$\sum_{i=1}^p |\beta_i| \leq c$$

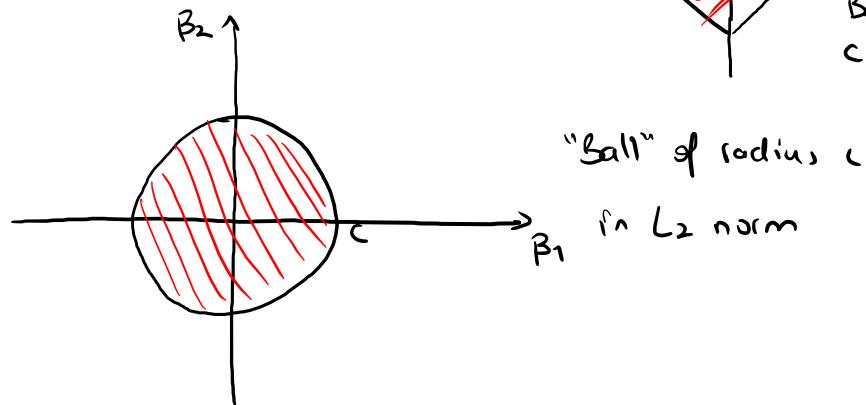
$$|\beta_1| + |\beta_2| \leq c$$

$$\Leftrightarrow \begin{cases} \beta_1 + \beta_2 \leq c & \beta_1, \beta_2 \geq 0 \\ \beta_1 - \beta_2 \leq c & \beta_1 \geq 0, \beta_2 \leq 0 \\ -\beta_1 + \beta_2 \leq c & \beta_1 \leq 0, \beta_2 \geq 0 \\ -\beta_1 - \beta_2 \leq c & \beta_1 \leq 0, \beta_2 \leq 0 \end{cases}$$



### RIDGE

$$\beta_1^2 + \beta_2^2 \leq c$$

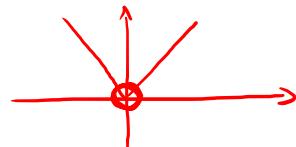


## EXAMPLE

Let's assume that we have one predictor ( $p=1$ )

$$\hat{\beta}_{\text{Lasso}} = \underset{\beta}{\operatorname{argmin}} \quad (\underbrace{y - x^t \beta}_{\sum_{i=1}^n (y_i - x_i \beta)^2})^t (y - x^t \beta) + \lambda |\beta|$$

$$|\beta| = \begin{cases} \beta & \beta \geq 0 \\ -\beta & \beta \leq 0 \end{cases}$$



DERIVATIVE  
wrt  
 $\beta$

$$-2 x^t (y - x^t \beta) + \lambda \operatorname{sign}(\beta) = 0$$

where  $\operatorname{sign}(\beta) = \begin{cases} 1 & \beta > 0 \\ -1 & \beta < 0 \\ [-1, 1] & \beta = 0 \text{ (sub-differential)} \end{cases}$

$$x^t(y - x\hat{\beta}_L) - \frac{\lambda}{2} \operatorname{sign}(\beta) = 0$$

$$\hat{\beta}_L > 0 \Rightarrow x^t(y - x\hat{\beta}_L) = \frac{\lambda}{2}$$

$$\Rightarrow x^t y - x^t x \hat{\beta}_L = \frac{\lambda}{2} \Rightarrow (x^t y - \frac{\lambda}{2}) = x^t x \hat{\beta}_L$$

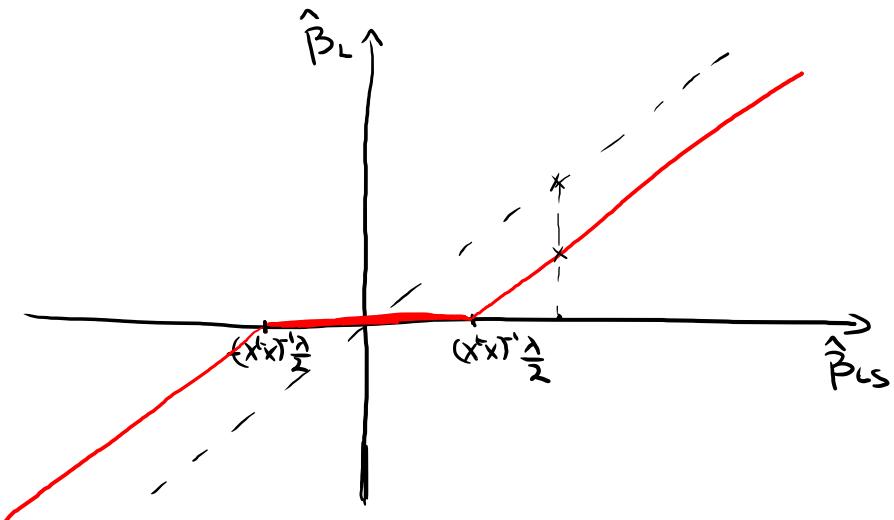
$$\Rightarrow \hat{\beta}_L = (x^t x)^{-1}(x^t y - \frac{\lambda}{2}) = \underbrace{(x^t x)^{-1} x^t y}_{\hat{\beta}_{LS}} - (x^t x)^{-1} \frac{\lambda}{2} = \\ = \underline{\hat{\beta}_{LS} - (x^t x)^{-1} \frac{\lambda}{2}}$$

$$\hat{\beta}_L < 0 \Rightarrow x^t(y - x\hat{\beta}_L) = -\frac{\lambda}{2} \Rightarrow \hat{\beta}_L = \hat{\beta}_{LS} + (x^t x)^{-1} \frac{\lambda}{2}$$

$$\underbrace{\hat{\beta}_L}_{\uparrow} = 0 \Rightarrow x^t(y - x\hat{\beta}_L) \in [-\frac{\lambda}{2}, \frac{\lambda}{2}]$$

$$\Rightarrow x^t y - x^t x \hat{\beta}_L \in [-\frac{\lambda}{2}, \frac{\lambda}{2}]$$

$$\Rightarrow \underbrace{(x^t x)^{-1} x^t y}_{\hat{\beta}_{LS}} - \cancel{(x^t x)^{-1} x^t} \hat{\beta}_L \in \left[ -\cancel{(x^t x)^{-1} \frac{\lambda}{2}}, \cancel{(x^t x)^{-1} \frac{\lambda}{2}} \right] \\ \Rightarrow \hat{\beta}_L \in \left[ -\cancel{(x^t x)^{-1} \frac{\lambda}{2}}, \cancel{(x^t x)^{-1} \frac{\lambda}{2}} \right]$$



$\textcircled{P=1}$

$$(X^T X)^{-1} = \frac{1}{\sum_{i=1}^n x_i^2}$$

SOFT-THRESHOLDING Lasso penalty translates the value towards zero, making them exactly zero when small enough

## LASSO ALGORITHM (COORDINATE DESCENT)

obj function  $f(\beta) = (y - x\beta)^T (y - x\beta) + \lambda \sum_{j=1}^p |\beta_j|$   
"  $f(\beta_1, \beta_2, \dots, \beta_p)$

Optimize w.r.t to  $\beta_j$  while keeping all the other values fixed to the value in the previous iteration.

The algorithms are extremely efficient !

## FINAL COMMENTS

- ① Both ridge and lasso admit solutions also when  $X^T X$  is not invertible so they can be used on high-dimensional data ( $p \gg n$ )
- ② Regularization reduces prediction error, in terms of reducing variance but at the expense of bias.  
 $\lambda$  choice is crucial and needs CV
- ③ Lasso produces sparse solutions compared to ridge, though with higher variance.
- ④ Correlated predictors  $\rightarrow$  ridge likely to do better  
Few associated predictors  $\rightarrow$  lasso might do better and be more useful than ridge.
- ⑤ Thus have been extended in different directions  $\begin{cases} \text{penalties} \\ \text{glm} \end{cases}$