

$$\text{Ridge: } \|X\theta - Y\|_2^2 + \lambda \|\theta\|_2^2 \rightarrow \hat{\theta}_{\text{ridge}} = (X^T X + \lambda I)^{-1} X^T Y$$

$$\text{Lasso: } \|X\theta - Y\|_2^2 + \lambda \|\theta\|_1 \rightarrow \hat{\theta}_{\text{lasso}} =$$

$\|\theta\|_1 = \sum_{i=1}^n |\theta_i|$  Induces sparsity  $\rightarrow$  Easy to interpret.

$\|\theta\|_2^2 = \sum_{i=1}^n \theta_i^2$  In many cases, has few non-zero entries.

$d > n \rightarrow \text{OLS is not unique}$

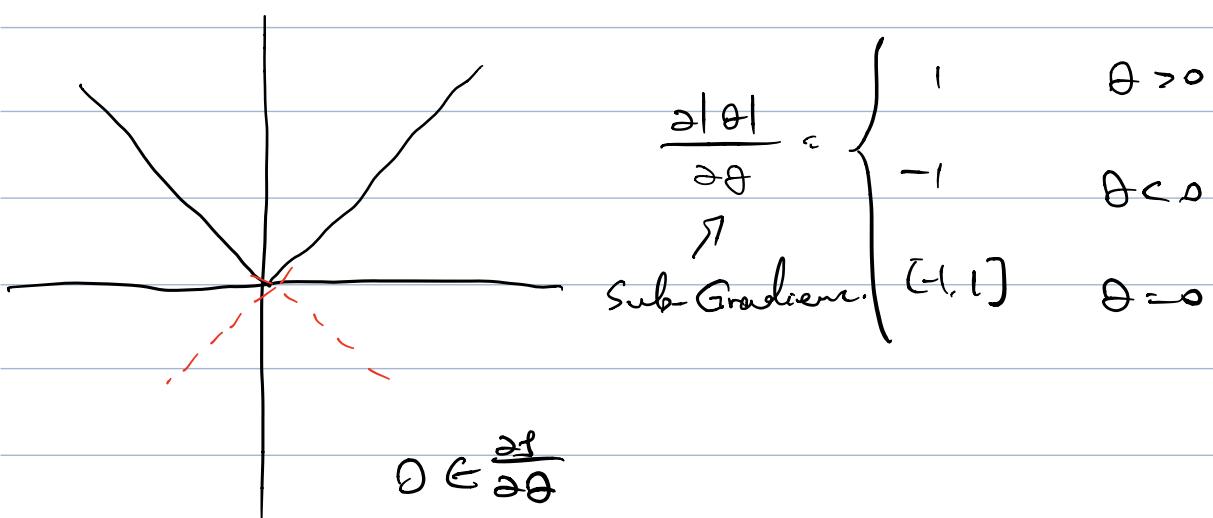
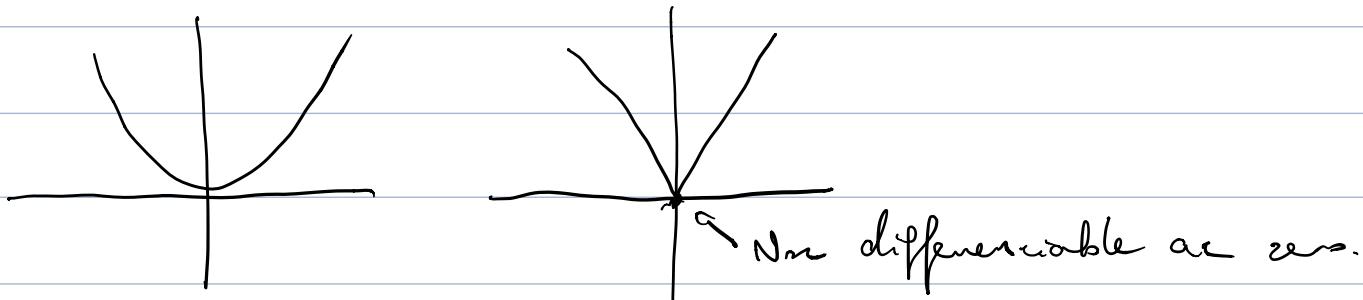
$$\# \text{non-zero entries in } \theta^* = \| \theta^* \|_0 = s$$

scc d, d>n.

$\hat{\theta}_{\text{ols}}$   $\hat{\theta}_{\text{ridge}}$   $\hat{\theta}_{\text{lasso}}$

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

$$\underset{\theta}{\operatorname{argmin}} (y - \theta)^2 + \lambda |\theta|$$



Suppose  $\underline{\theta^* > 0}$  (If  $y > \frac{\lambda}{2}$ )

$$2(\theta - y) + \lambda = 0$$

$$\Rightarrow \theta^* = y - \frac{\lambda}{2}$$

$\theta^* < 0$

$$\Rightarrow 2(\theta - y) - \lambda = 0$$

$$\Rightarrow \theta^* = y + \frac{\lambda}{2} \quad (\text{if } y < -\frac{\lambda}{2})$$

$\theta^* = 0$ :

$$\frac{\partial f}{\partial \theta} = 2(\theta - y) + \lambda x = 0 \quad x \in [-1, 1].$$

$$= -2y + \lambda x \quad x \in [-1, 1].$$

$$0 = -2y + \lambda x \Rightarrow x = \frac{2y}{\lambda}$$

$$\frac{2y}{\lambda} \in [-1, 1]$$

$$\Leftrightarrow y \in [-\frac{\lambda}{2}, \frac{\lambda}{2}]$$

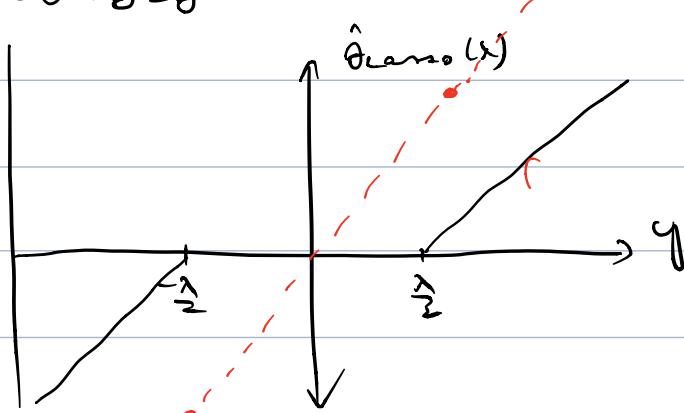
$$\Rightarrow \theta^* = 0$$

$$\underset{\theta}{\operatorname{argmin}} (y - \theta)^2 + \lambda |\theta|$$

$$\theta^* = \begin{cases} y - \frac{\lambda}{2} & \text{if } y > \lambda/2 \\ y + \frac{\lambda}{2} & \text{if } y < -\lambda/2 \\ 0 & \text{if } y \in [-\frac{\lambda}{2}, \frac{\lambda}{2}]. \end{cases}$$

Why?

$$\theta \in \frac{\partial f}{\partial \theta} \Big|_{\theta=\theta^*}$$



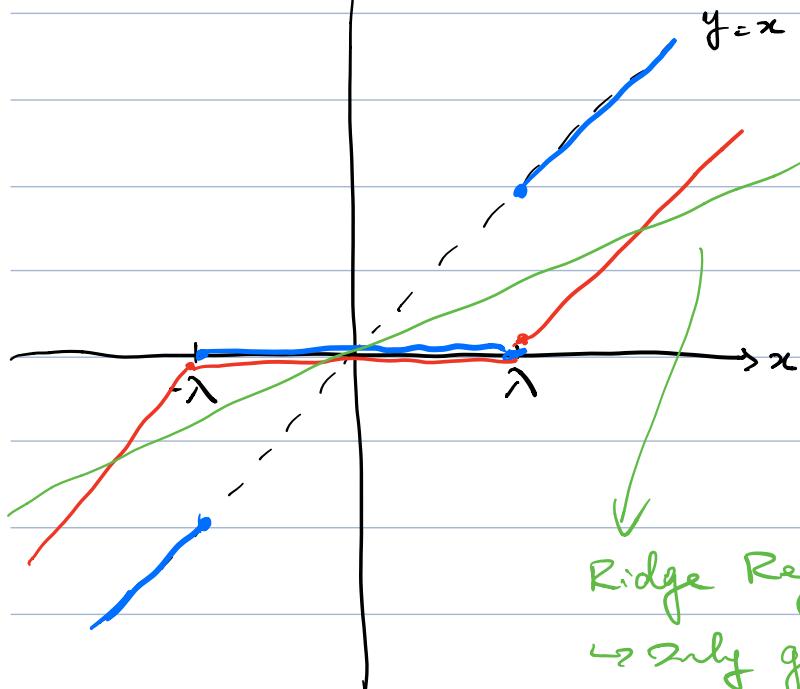
## \* Soft Thresholding

$$\hat{\theta} = \text{sign}(y) \cdot (|y| - \frac{\lambda}{2})_+$$

$$(x)_\pm = \begin{cases} x & x > 0 \\ 0 & x=0 \\ -x & x < 0 \end{cases}$$

$$\text{soft}(y, \frac{\lambda}{2})$$

## Soft Thresholding



Hard( $x, \lambda$ )

Hard Thresholding

Soft( $x, \lambda$ )

Soft Thresholding

- Continous

- No sudden jump

Ridge Regression

→ Only gives sparse solution as 0.

→ Does not make it exactly zero.

→ Only reduces magnitude.

$$f(\theta) = (y - \theta)^2 + \lambda |\theta|$$

$$0 \in \frac{\partial f}{\partial \theta}$$

$$f(\theta) = \sum_{i=1}^n (x_i^\top \theta - y_i)^2 + \lambda \sum_{j=1}^d (\theta_j)$$

$$= \|X\theta - Y\|_2^2 + \lambda \|\theta\|_1$$

$\theta_i \rightarrow$  Variable } → Reduces to one-dimensional

Fix other  $\theta_j$ 's } Loss → closed form solution as shown  
before.

Coordinate-descent [like iterative algorithm]

Each reserve step: new  $\hat{\beta}$ ,

$$1. R(\beta) = \|X\beta - Y\|_2^2$$

$$R(\beta) = R(\beta^{ols}) + (\beta^{ols} - \beta)^T \frac{X^T X}{n} (\beta^{ols} - \beta)$$

$$\text{LHS} = \underbrace{\|X\beta - X\hat{\beta} + X\hat{\beta} - Y\|_2^2}_n = \underbrace{\|X\beta - X\hat{\beta}\|_2^2}_n + \underbrace{\|X\hat{\beta} - Y\|_2^2}_n$$

$$+ 2 \underbrace{(X\beta - X\hat{\beta})^T (X\hat{\beta} - Y)}_n$$

$$= \|X(\beta - \hat{\beta})\|$$

$$= (-\beta + \hat{\beta})^T X^T X (-\beta + \hat{\beta})$$

$$\hookrightarrow \underbrace{(\beta - \hat{\beta})^T X^T (X\hat{\beta} - Y)}_0 \quad X^T Y = X^T Y$$

$$X^T (X\hat{\beta} - Y) = 0$$

$$R(\beta) \geq R(\hat{\beta}_{ols}) \quad \forall \beta$$

$\rightarrow \hat{\beta}_{ols}$  is the minimizer of training error

$$R(\beta) = R(\hat{\beta}) + C$$

$$\Leftrightarrow \left\{ (\beta - \hat{\beta})^T \frac{X^T X}{n} (\beta - \hat{\beta}) = C \right\}$$

Quadratic Form:

$$(v - \hat{v})^T A (v - \hat{v}) = C$$

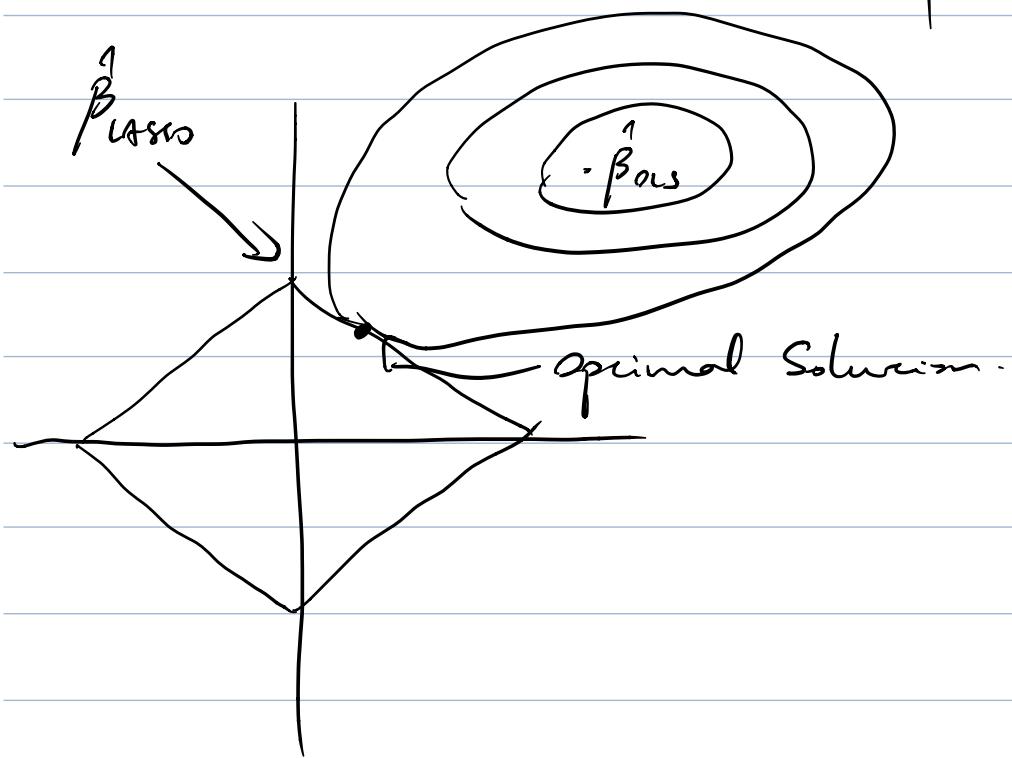
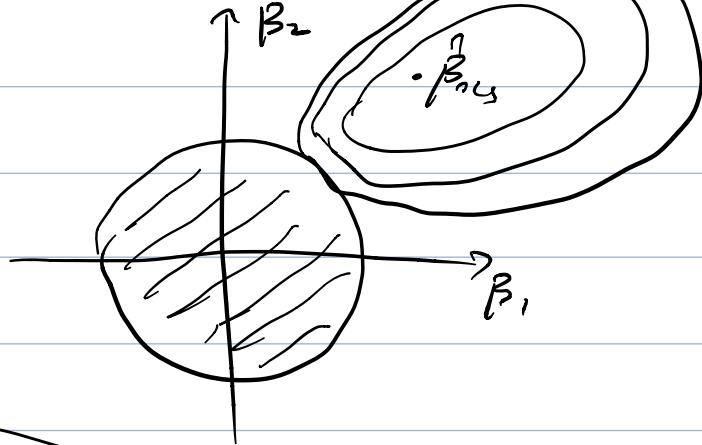
$$\frac{X^T X}{n} = I \Rightarrow \|\beta - \hat{\beta}\|_2^2 = C$$

$$v^T A v = \frac{v^T (A + A^T) v}{2}$$

$$|\beta_1| + |\beta_2| \leq \mu$$

$$\|X\beta - Y\|_2^2 + \lambda \|\beta\|_2$$

$$\min \|X\beta - Y\|_2^2 \text{ s.t. } \|\beta\|_1 \leq \mu$$



$$4. \|X\beta - Y\|_2^2 + \lambda \|\beta\|_1$$

Show  $\beta = 0$  is the optimal Lasso estimator if  $\lambda \geq \|X\theta - Y\|_\infty$

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

$$\lambda \rightarrow \infty \Rightarrow \hat{\beta}_\lambda = 0$$

$$\lambda > \|X\theta - Y\|_\infty \Rightarrow \hat{\beta}_{(\text{LASSO})} = 0$$

$f(0) < f(\beta)$  for any  $\beta \neq 0$ .

Then  $\beta = 0$  is the optimal estimator.

$$\|Y\|_2^2 \leq \|X\beta - Y\|_2^2 + \lambda \|\beta\|_1$$

$$\underbrace{\|X\beta\|_2^2 + \|Y\|_2^2}_{\geq 0} - 2\beta^T X^T Y + \lambda \|\beta\|_1 = \text{RHS}$$

$$RHS \geq \|Y\|_2^2 + \lambda \|\beta\|_1 - 2\beta^T X^T Y$$

$$|\beta^T v| \leq \|\beta\|_1 \|v\|$$

$$LHS = \left| \sum_{i=1}^d \beta_i v_i \right|$$

Hölder's Inequality

$$\leq \sum_{i=1}^d |\beta_i| |v_i|$$

Young's Inequality

$$\leq \left( \sum_{i=1}^d |\beta_i| \right) \max_j (v_i)$$

Cauchy-Schwarz.

$$= \|\beta\|_1 \|v\|_\infty$$

$$RHS \geq \underbrace{\|Y\|_2^2}_{> \|Y\|_2^2} + \|\beta\|_1 (\lambda - 2\|X^T Y\|_\infty)$$

$$\text{If } \lambda > 2\|X^T Y\|_\infty$$

$$\Rightarrow \hat{\beta}_\lambda(\text{Lasso}) = 0$$

4 & 5

$$(X_1\theta_1 + X_2\theta_2 + X_2^T\theta_2 - y)^2 + \lambda(\theta_1^2 + \theta_2^2) + \lambda\|\theta_2\|_2^2$$

$$(X_1\theta_1 + X_2\theta_2 + X_2^T\theta_2 - y)^2 + \lambda(|\theta_1| + |\theta_2|) + \lambda\|\theta_2\|_2^2$$

$$= (X(\theta_1 + \theta_2) + \tilde{X}^T\tilde{\theta} - y)^2 + \lambda(\theta_1^2 + \theta_2^2) + \lambda\|\tilde{\theta}\|_2^2$$

$$= (X(\theta_1 + \theta_2) + \tilde{X}^T\tilde{\theta} - y)^2 + \lambda(|\theta_1| + |\theta_2|) + \lambda\|\tilde{\theta}\|_2^2$$

$$\Rightarrow \theta_1 = \frac{\theta^*}{2}, \quad \theta_2 = \frac{\theta^*}{2} \quad \text{is the unique solution}$$

$$\Rightarrow \theta_1 = \theta^*, \quad \theta_2 = 0$$

$$\theta_1 = \theta^*, \quad \theta_2 = \theta^*$$

& Many more

&

\* Lasso can give sparse solutions even with collinear/correlated features but Ridge can't.

} Non-Unique but sparse  
solutions possible.

$$\|X\beta - Y\|_2^2 + \lambda \|\beta\| \rightarrow \hat{\beta}(\lambda) \text{ be the solution for problem}$$

$$\|X\hat{\beta} - Y\|_2^2 \text{ s.t. } \|\beta\|_1 \leq \mu$$

then we want to find a  $\mu$  s.t.  $\hat{\beta}(\lambda)$  is optimal

$$\|X\hat{\beta} - Y\|_2^2 \leq \|X\beta - Y\|_2^2 \text{ for any } \|\beta\|_1 \leq \mu.$$

$$\text{Claim: } \mu = \|\hat{\beta}(\lambda)\|_1$$

Proof: We know  $\|\hat{\beta}(\lambda)\|_1$  is optimal for  $P_1$

$$\|X\hat{\beta} - Y\|_2^2 + \lambda \|\hat{\beta}\|_1 \leq \|X\beta - Y\|_2^2 + \lambda \|\beta\|_1, \forall \beta \in$$

$$\text{for } \|\beta\|_1 \leq \mu = \|\hat{\beta}\|_1 \leq \|X\beta - Y\|_2^2 + \lambda \|\beta\|_1 \leq \mu$$

because  $T_0 < 0$ .

Optimality

Condition for  $P_2$

## Kernel Ridge Regression.

$$\begin{aligned} \hat{\beta}_{\text{ridge}} &= \underbrace{(X^T X + \lambda I_n)}_{d \times d}^{-1} X^T Y && \text{preferable if } n \gg d \\ &= X^T \underbrace{(X X^T + \lambda I_n)}_{n \times n}^{-1} Y && \text{preferable if } n \ll d \end{aligned}$$

$$X X^T = \begin{bmatrix} x_1^T x_1 & x_1^T x_2 & \dots & \dots \end{bmatrix}$$

$$X = \begin{bmatrix} x_1^T \\ \vdots \\ x_n^T \end{bmatrix}$$

$$X^T X = \sum_{i=1}^n x_i x_i^T$$