

## Lecture September 10

Shrinkage Methods: Ridge and Lasso,

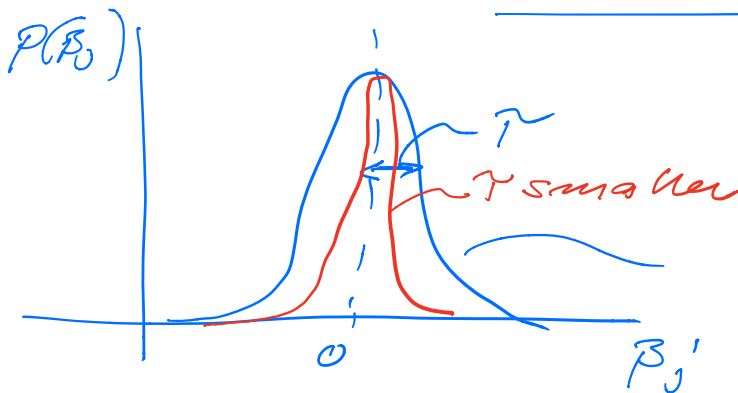
Intuitive understanding:

Bayes' theorem:

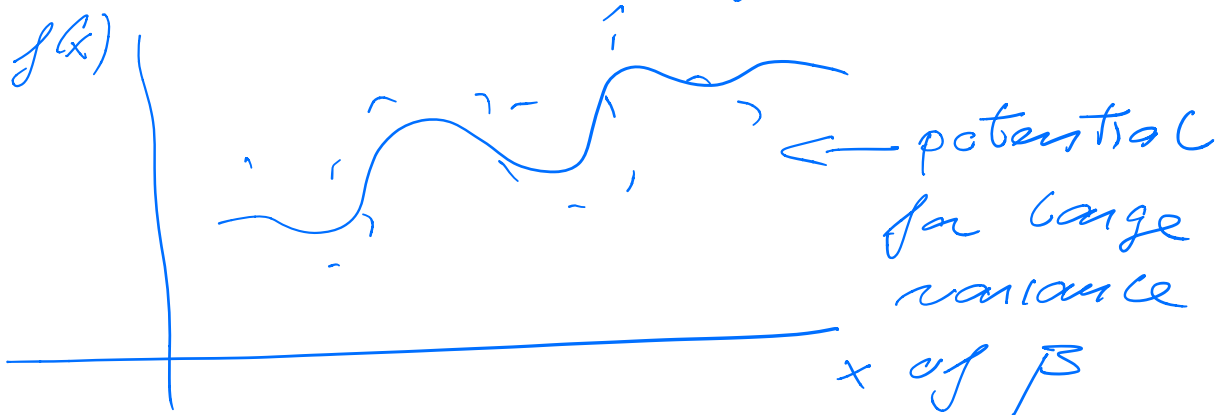
$$P(\beta | X, y) \propto P(y | X, \beta) P(\beta)$$

$$\prod_{i=0}^{n-1} N(y_i | x_i^T \beta)$$

↑  
prior



$$P(\beta) = \prod_{i=0}^{p-1} N(0, \tau^2)$$
$$= \frac{1}{\tau^p} e^{-\beta_j^2 / 2\tau^2}$$



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

Can we shrink this variance!

$$\hat{\beta}^{\text{Ridge}} = \arg \min_{\beta \in \mathbb{R}^p} \frac{1}{n} \sum_{i=0}^{n-1} (y_i - x_i^T \beta)^2 + \lambda \sum_{j=0}^{p-1} \beta_j^2$$

L2-norm Regularization

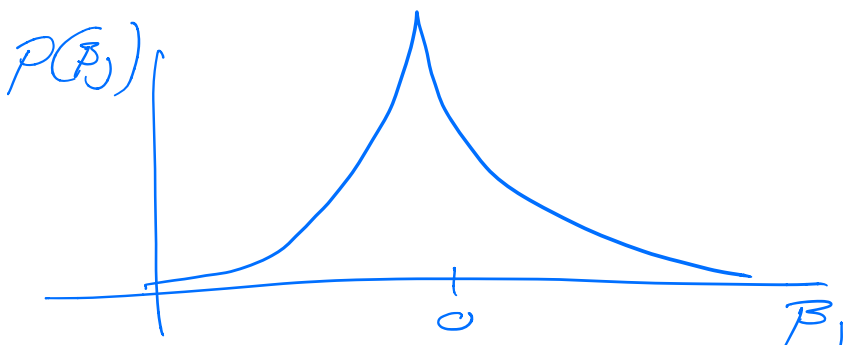
$$\lambda \sim \frac{1}{2\tau^2} \quad \lambda \geq 0$$

$$\Rightarrow \hat{\beta}^{\text{Ridge}} = \underbrace{(X^T X + \lambda I)^{-1}}_{\in \mathbb{R}^{p \times p}} X^T y$$

Lasso: Least absolute shrinkage and selection operator

$p(\beta)$  = double exponential (Laplace) distribution

$$p(\beta_j) = e^{-|\beta_j|/\tau} \quad (e^{-(\beta_j - \mu)/\tau})$$



$$P(\beta/xy) \propto \prod_{i=0}^{n-1} P(y_i/x_i, \beta) \frac{\prod_{j=0}^{p-1} e^{-|\beta_j|/\tau}}{\prod_{j=0}^{p-1} \tau}$$

minimize  $-\log P(\beta/xy)$  to find  
 $\hat{\beta}_{\text{optimal}}$

$$\hat{\beta}^{\text{Lasso}} = \arg \min_{\beta \in \mathbb{R}^p} \frac{1}{n} \sum_{i=0}^{n-1} (y_i - x_i \beta)^2 + \lambda \sum_{j=0}^{p-1} |\beta_j|$$

L1-norm regularization

$$\hat{\beta}^{\text{Lasso}} = \arg \min_{\beta \in \mathbb{R}^p} \frac{1}{n} \|y - X\beta\|_2^2 + \lambda \|\beta\|_1$$

Ridge  $\lambda \geq 0$  given

$$\sum_{j=0}^{p-1} \beta_j^2 \leq s$$

Lasso  $\lambda \geq 0$   $\sum_{j=0}^{p-1} |\beta_j| \leq t$

Simple example to illustrate  
 differences between  $\hat{\beta}^{\text{OLS}}$ ,  
 $\hat{\beta}^{\text{Ridge}}$  and  $\hat{\beta}^{\text{Lasso}}$

$p \times p$  matrix  $X$

$X$  is a diagonal matrix

$$X = \begin{bmatrix} 1 & & & \\ & 1 & & \\ & & \ddots & \\ & & & 1 \end{bmatrix}$$

Skip  $1/n$  in MSE,  $n = p$

OLS :  $\sum_{i=0}^{n-1} (y_i - \beta_i)^2$

$$\hat{\beta}^{OLS} \Rightarrow \hat{\beta}_j^{OLS} = y_j$$

Ridge :  $\sum_{i=0}^{n-1} (y_i - \beta_i)^2 + \lambda \sum_{j=0}^{n-1} \beta_j^2$

$$\hat{\beta}_i^{Ridge} = \frac{y_i}{1 + \lambda}$$

$$\lambda = 0 \Rightarrow \hat{\beta}_i^{Ridge} = \hat{\beta}_i^{OLS}$$

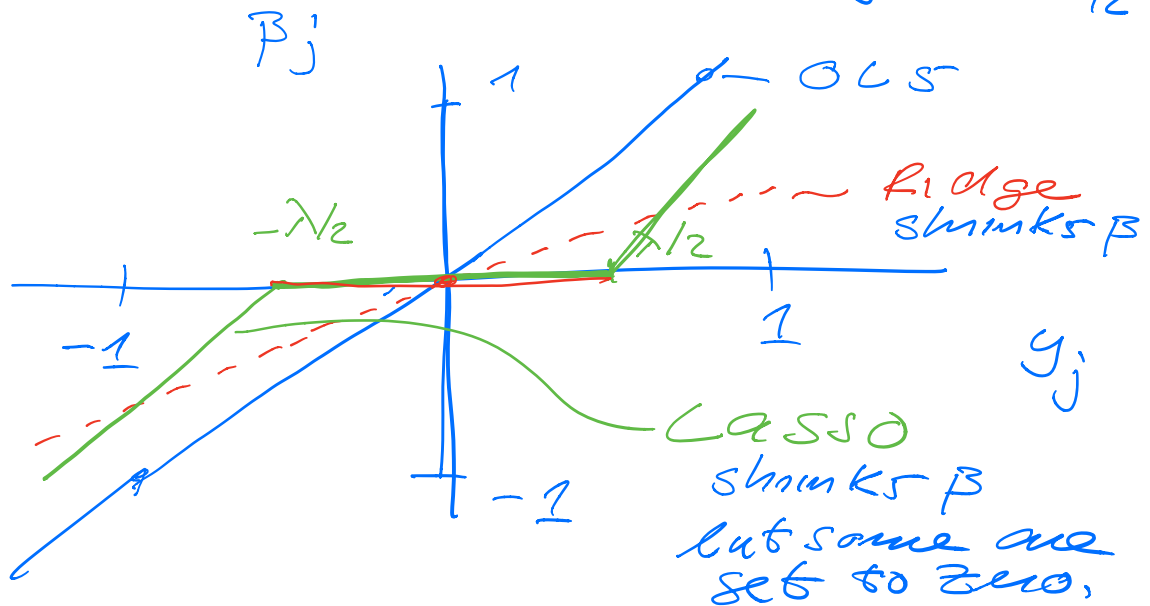
lim  $\lambda \rightarrow \infty$   $\hat{\beta}_i^{Ridge} \rightarrow 0$

Lasso

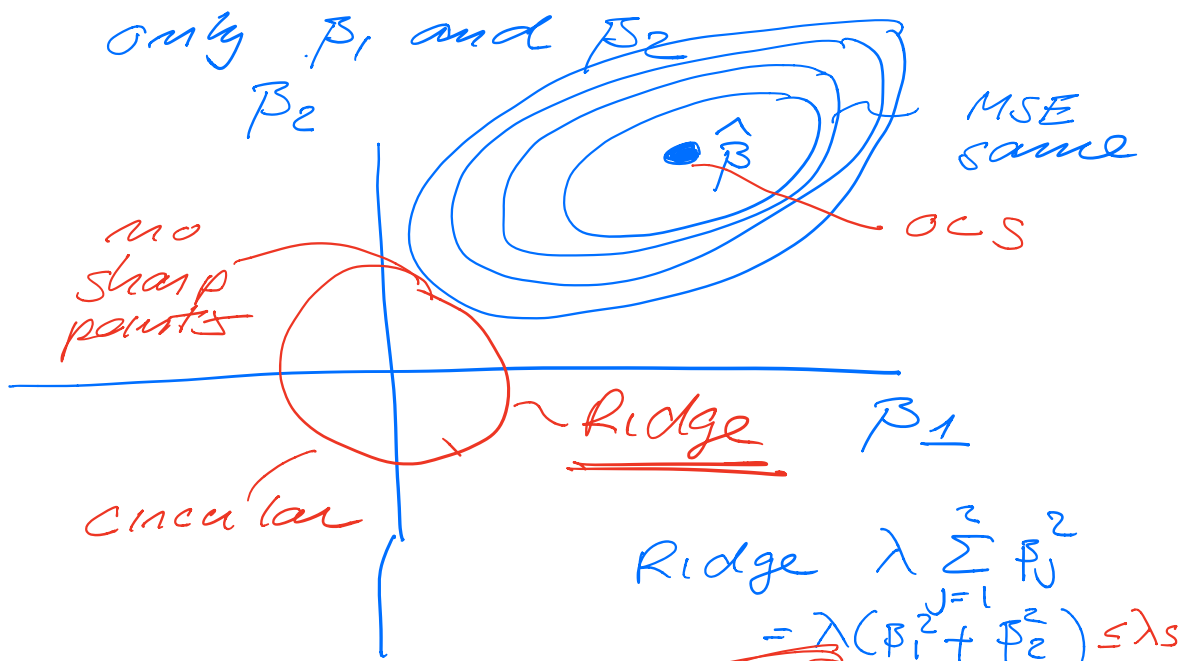
$$\frac{\partial}{\partial \beta_i} \left( \sum_{i=0}^{n-1} (y_i - \beta_i)^2 + \lambda \sum_{i=0}^{n-1} |\beta_i| \right) = 0$$

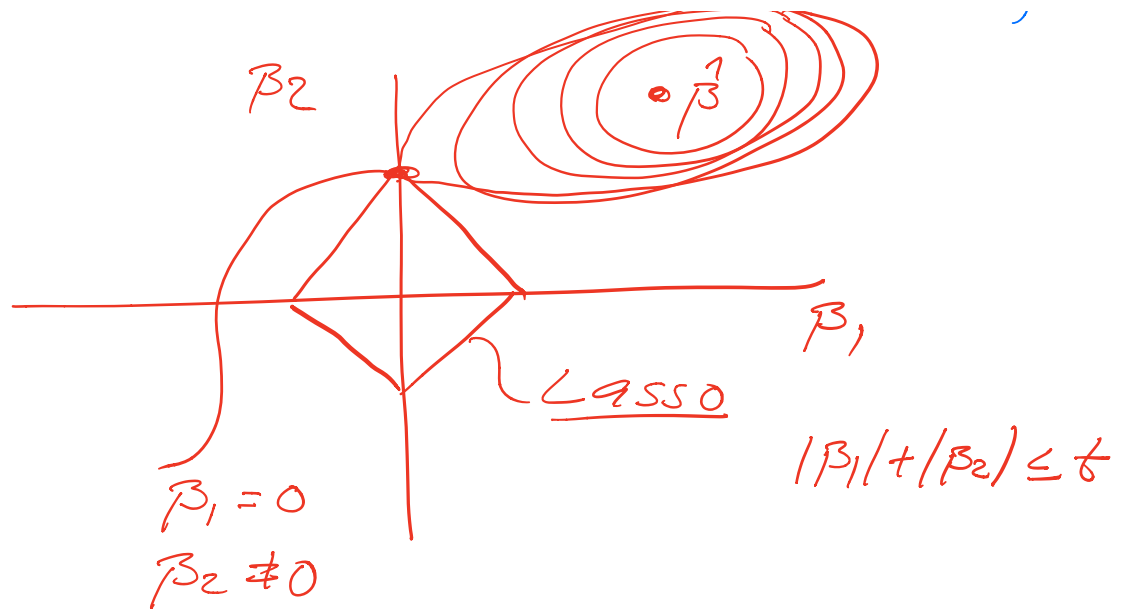
$$-2 \sum_i (y_i - \beta_i) + \lambda \sum_i \frac{\beta_i}{|\beta_i|} = 0$$

$$\beta_{Lasso} = \begin{cases} y_i - \lambda/2 & \text{if } y_i > \lambda/2 \\ y_i + \lambda/2 & \text{if } y_i < -\lambda/2 \\ 0 & \text{if } |y_i| \leq \lambda/2 \end{cases}$$



Example 2





## Ridge analysis-

- Singular value decomposition,
- covariance and correlation matrix,