# Optimization and Machine Learning  SI151

Lu Sun

School of Information Science and Technology

ShanghaiTech University

March 4, 2021

Today:
- Linear Methods for Regression II
  - Ridge Regression
  - The Lasso
  - Discussion

Readings:
- The Elements of Statistical Learning (ESL), Chapter 3
- Pattern Recognition and Machine Learning (PRML), Chapter 3

# Introduction

- **Subset selection**
  - retain a subset of the predictors, and discard the rest
  - accuracy and interpretation
  - discrete process
    - variable are either retained or discarded
    - high variance

- **Shrinkage methods**
  - continuous process
    - don't suffer much from high variability
  - ridge regression, lasso, …

# Linear Methods for Regression

--- Ridge Regression

# **Shrinkage Methods** – Ridge Regression

- Shrink the regression coefficients
  - impose a penalty on the size

*loss func.*   *regularization*

P1   $\hat{\beta}^{\text{ridge}} = \underset{\beta}{\arg\min} \left\{ \sum_{i=1}^{N} \left( y_i - \beta_0 - \sum_{j=1}^{p} x_{ij}\beta_j \right)^2 + \lambda \sum_{j=1}^{p} \beta_j^2 \right\}$

  - the larger the value of $\lambda$, the greater the amount of shrinkage
  - the coefficients are shrunk toward zero
- An equivalent expression

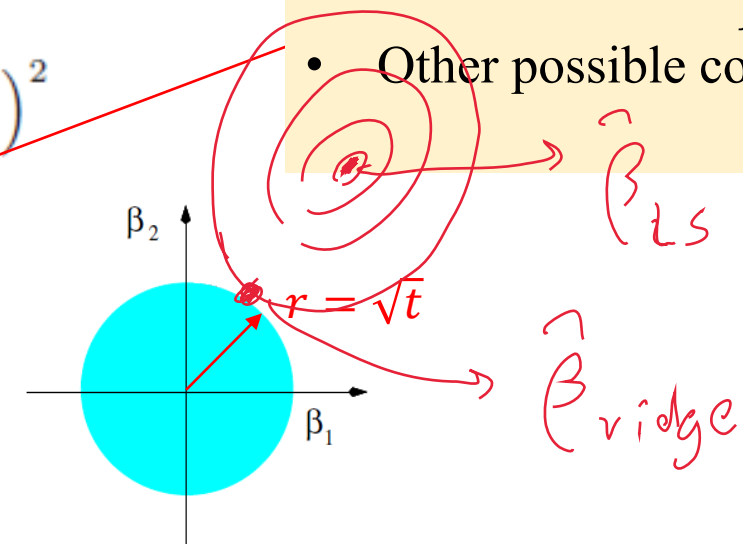P2   $\hat{\beta}^{\text{ridge}} = \underset{\beta}{\arg\min} \sum_{i=1}^{N} \left( y_i - \beta_0 - \sum_{j=1}^{p} x_{ij}\beta_j \right)^2$

subject to $\sum_{j=1}^{p} \beta_j^2 \le t,$

  - One-to-one correspondence between $\lambda$ and $t$

- Squared $\ell_2$-norm on $\beta$

$$\|\beta\|_2^2 = \beta^T \beta = \sum_{j=1}^{p} \beta_j^2$$

- Other possible constraints?

$\hat{\beta}_{LS}$

$r = \sqrt{t}$

$\hat{\beta}_{ridge}$

$\beta_2$

$\beta_1$

# **Shrinkage Methods** – Ridge Regression

- Equivalence between P1 and P2

P1: $\quad \hat{\beta} = \underset{\beta}{\operatorname{argmin}} \|\mathbf{y} - \mathbf{X}\beta\|_2^2 + \lambda\|\beta\|_2^2$

P2: $\quad \tilde{\beta} = \underset{\beta}{\operatorname{argmin}} \|\mathbf{y} - \mathbf{X}\beta\|_2^2, \text{ s.t.} \|\beta\|_2^2 \leq t$

- Goal: $\forall t, \exists \lambda \geq 0: \hat{\beta} = \tilde{\beta}$

**Proof:**

- Step 1: assume that P1 is solved

$$\boxed{\mathbf{X}^T(\mathbf{y} - \mathbf{X}\hat{\beta}) + \lambda\hat{\beta} = 0}$$

- Lagrange form of P2

$$L(\beta, \mu) = \|\mathbf{y} - \mathbf{X}\beta\|_2^2 + \mu(\|\beta\|_2^2 - t)$$

- KKT conditions

  1. $\nabla_\beta L(\tilde{\beta}, \tilde{\mu}) = 0 \implies \boxed{\mathbf{X}^T(\mathbf{y} - \mathbf{X}\tilde{\beta}) + \tilde{\mu}\tilde{\beta} = 0}$
  2. $\tilde{\mu}\left(\|\tilde{\beta}\|_2^2 - t\right) = 0$
  3. $\tilde{\mu} \geq 0$
  4. $\|\tilde{\beta}\|_2^2 \leq t$

- Thus,
  - if
$$t = \|\hat{\beta}\|_2^2$$
  - Then
$$\tilde{\mu} = \lambda, \qquad \tilde{\beta} = \hat{\beta}$$
  - Satisfy the KKT conditions.

- Step 2: conversely, assume that P2 is solved

- The optimal solution $(\tilde{\beta}, \tilde{\mu})$ must satisfies KKT conditions. Therefore, let $\lambda = \tilde{\mu}$, we always have $\hat{\beta} = \tilde{\beta}$.

Strong duality holds for P2:

| $(\tilde{\beta}, \tilde{\mu})$ is the optimal solution of P2 | $\implies$ $\impliedby$ | $(\tilde{\beta}, \tilde{\mu})$ satisfies KKT conditions |
| --- | --- | --- |

# **Shrinkage Methods** – Ridge Regression

## Important notes

- ridge solutions are not equivalent under <span style="color:red">scaling of inputs</span>
  - *standardize* the inputs before solving it
- the intercept $\beta_0$ should be <span style="color:red">left out</span> of the penalty term

Ex. 3.5 →
  - once $x_{ij} - \bar{x}_j$, $\beta_0$ is estimated by $\bar{y} = \frac{1}{N}\sum_{i=1}^{N} y_i$
  - the rest parameters are estimated by the centered data

- Henceforth we assume the data has been <span style="color:blue">standardized</span>
  - **X** has $p$ rather than $p + 1$ columns

**Standardization**

$$x' = \frac{x - \bar{x}}{\sigma}$$

$X = \begin{bmatrix} x_1 & \cdots & x_p \\ \vdots & & \vdots \\ x_1 & \cdots & x_p \end{bmatrix}$   $y$

$\beta_0 = \bar{y}$, $\hat{\beta}$

$x_0 \implies \hat{y}_0$

① standardize $x_0$

② $\hat{y}_0 \leftarrow x_0^T \hat{\beta} + \beta_0$

**Prediction?**

$\hat{y}_0 - \bar{y} = x_0^T \hat{\beta}$

$\hat{y}_0 = x_0^T \hat{\beta} + \beta_0$

# Least Squares

## Training

1. standardize

$(\forall j)$ $\quad x_{ij} \leftarrow \dfrac{x_{ij} - \bar{x}_j}{\sigma_{x_j}}$,

2. $\quad X \leftarrow [\mathbb{1}, X]$

3. $\quad \min_{\beta} \| y - X\beta \|_2^2$

$\Rightarrow \hat{\beta} = (X^TX)^{-1} X^T y$

## Testing $\quad (x_0 \in \mathbb{R}^p)$

1. $x_{0,j} \leftarrow \dfrac{x_{0,j} - \bar{x}_j}{\sigma_{x_j}}$

2. $x_0 \leftarrow (1; x_0)$

3. $\hat{y}_0 \leftarrow x_0^T \hat{\beta}$

---

$D = \{(x_i, y_i)\}_{i=1}^n$

$(x_i \in \mathbb{R}^p, \; y_i \in \mathbb{R})$

$(X \in \mathbb{R}^{n \times p}, \; y \in \mathbb{R}^{n \times 1})$

$\bar{x}_j = \dfrac{1}{n} \sum_{i=1}^n x_{ij}$

---

# Ridge Regression

## Training

1. standardization

$x_{ij} \leftarrow \dfrac{x_{ij} - \bar{x}_j}{\sigma_{x_j}}$

2. centering $y$

$y_i \leftarrow y_i - \bar{y}, \quad (\bar{y} = \dfrac{1}{n} \sum_{i=1}^n y_i)$

3. $\min_{\beta} \| y - X\beta \|_2^2 + \lambda \| \beta \|_2^2$

$\Rightarrow \hat{\beta} = (X^TX + \lambda I)^{-1} X^T y$

$(\hat{\beta}_0 = \bar{y})$

## Testing $\quad (x_0)$

1. $x_{0,j} \leftarrow \dfrac{x_{0,j} - \bar{x}_j}{\sigma_{x_j}}$

2. $\hat{y}_0 \leftarrow x_0^T \hat{\beta} + \underset{\bar{y}}{\hat{\beta}_0}$

# **Shrinkage Methods** – Ridge Regression

*Handwritten annotations:*
① $\hat{\beta}_0 = \bar{y}$
② $y_i \leftarrow y_i - \bar{y}$
$x_i^T \beta$

- Ridge regression in matrix form

$$\hat{\beta}^{\text{ridge}} = \text{argmin}_\beta \left\{ \sum_{i=1}^{N} \left( y_i - \beta_0 - \sum_{j=1}^{p} x_{ij}\beta_j \right)^2 + \lambda \sum_{j=1}^{p} \beta_j^2 \right\}$$

$$\hat{\beta}^{ridge} = \text{argmin}_\beta \text{PRSS}(\lambda, \beta) = \text{argmin}_\beta \|\mathbf{y} - \mathbf{X}\beta\|_2^2 + \lambda\|\beta\|_2^2$$

$\beta^T \beta$

- We can rewrite $\text{PRSS}(\lambda, \beta)$ as follows

$$\text{PRSS}(\lambda, \beta) = (\mathbf{y} - \mathbf{X}\beta)^T(\mathbf{y} - \mathbf{X}\beta) + \lambda\beta^T\beta$$

$$= \mathbf{y}^T\mathbf{y} - \beta^T\mathbf{X}^T\mathbf{y} - \mathbf{y}^T\mathbf{X}\beta + \beta^T\mathbf{X}^T\mathbf{X}\beta + \lambda\beta^T\beta$$

- Differentiating $\text{PRSS}(\lambda, \beta)$ w.r.t. $\beta$

$$\frac{\partial \text{PRSS}(\lambda, \beta)}{\partial \beta} = -2\mathbf{X}^T\mathbf{y} + 2(\mathbf{X}^T\mathbf{X} + \lambda\mathbf{I}_p)\beta = 0$$

- The closed form solution $\hat{\beta}^{ridge} = (\mathbf{X}^T\mathbf{X} + \lambda\mathbf{I}_p)^{-1}\mathbf{X}^T\mathbf{y}$

- $\text{rank}(\mathbf{I}_p) = p$
- make the problem nonsingular, even if $\text{rank}(\mathbf{X}) < p$

7

# **Shrinkage Methods** – Ridge Regression

$X = \begin{bmatrix} x_1 & \cdots & x_p \end{bmatrix}$

Column space: $C(X) = \text{span}(\{x_1, \ldots, x_p\})$.

$= \{ \sum_{i=1}^{p} a_i x_i \mid a_i \in \mathbb{R}, \forall i \}$

Additional insight into ridge regression

- Singular value decomposition (SVD)

$$\mathbf{U}^T\mathbf{U} = \mathbf{I}_p, \mathbf{V}^T\mathbf{V} = \mathbf{I}_p \qquad \mathbf{X} = \mathbf{U}\mathbf{D}\mathbf{V}^T$$

- $\mathbf{U} \in \mathbb{R}^{N \times p}$: its columns span the column space ($\mathbb{R}^N$) of $\mathbf{X}$
- $\mathbf{V} \in \mathbb{R}^{p \times p}$: its columns span the row space ($\mathbb{R}^p$) of $\mathbf{X}$
- $\mathbf{D} \in \mathbb{R}^{p \times p}$: diagonal matrix $(d_1 \geq d_2 \geq \cdots \geq d_p \geq 0)$

- Singular values of $\mathbf{X}$
- if $\exists d_j = 0$, $\mathbf{X}$ is singular

Least squares

$$
\begin{aligned}
\mathbf{X}\hat{\beta}^{\text{ls}} &= \mathbf{X}(\mathbf{X}^T\mathbf{X})^{-1}\mathbf{X}^T y \\
&= \mathbf{U}\mathbf{U}^T y, \\
&= \sum_{j=1}^{p} \mathbf{u}_j \mathbf{u}_j^T y
\end{aligned}
$$

$Q$

The $j$-th column of $\mathbf{U}$

Ridge regression

$$
\begin{aligned}
\mathbf{X}\hat{\beta}^{\text{ridge}} &= \mathbf{X}(\mathbf{X}^T\mathbf{X} + \lambda\mathbf{I})^{-1}\mathbf{X}^T y \\
&= \mathbf{U}\,\mathbf{D}(\mathbf{D}^2 + \lambda\mathbf{I})^{-1}\mathbf{D}\,\mathbf{U}^T y \\
&= \sum_{j=1}^{p} \mathbf{u}_j \frac{d_j^2}{d_j^2 + \lambda} \mathbf{u}_j^T y,
\end{aligned}
$$

- shrinkage factor
- smaller $d_j$ leads to a larger shrinkage

# **Shrinkage Methods** – Ridge Regression

- Prostate cancer example
  - #training($N$) = 67, #testing=30
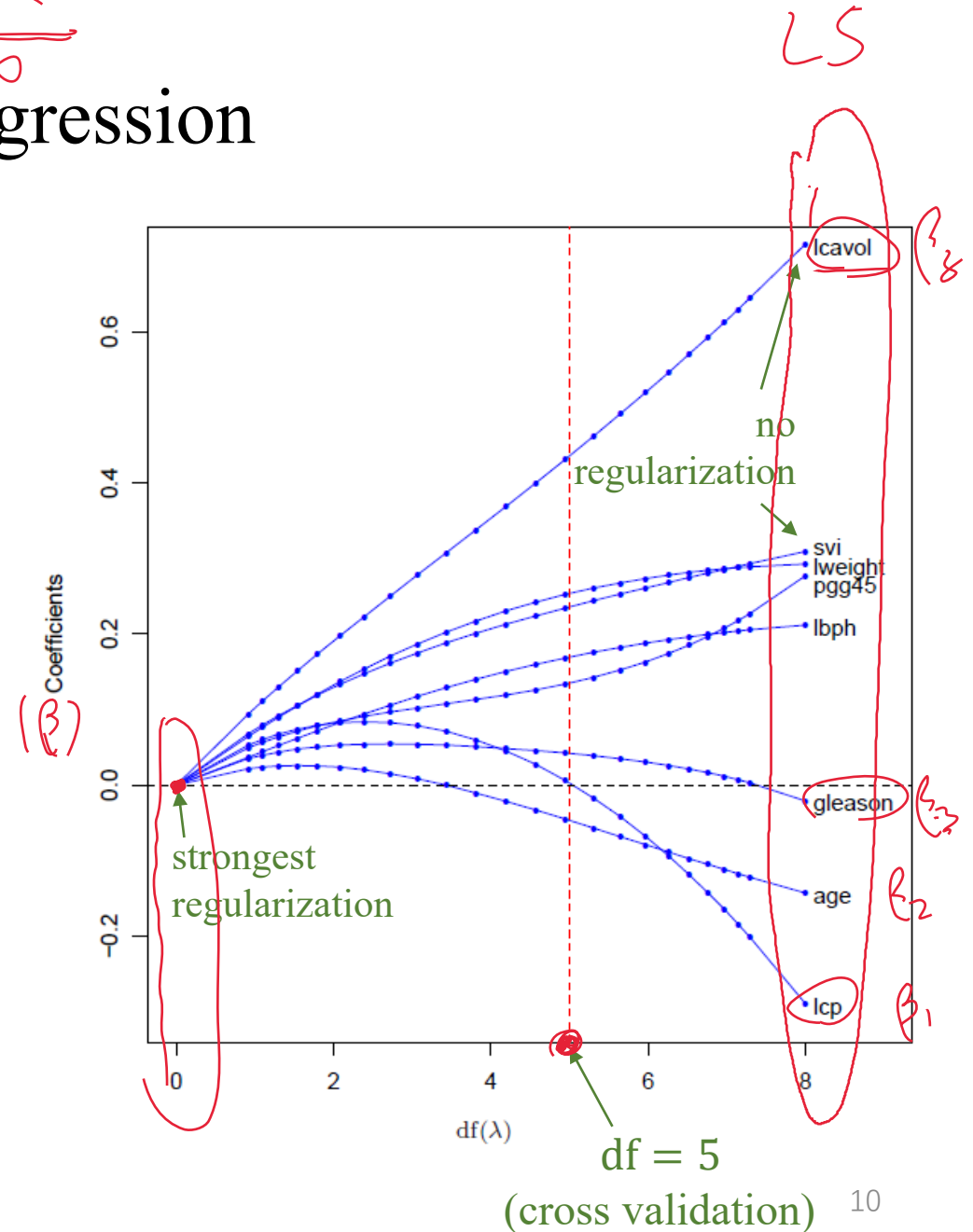  - #variables($p$)=8
  - ridge coefficient estimates
- *Effective degree of freedom*

$$\text{df}(\lambda) = \sum_{j=1}^{p} \frac{d_j^2}{d_j^2 + \lambda} \in (0, p]$$

A. B. C.

$$Tr(ABC) = Tr(BCA) = Tr(CAB)$$

$$X = UDV^T, V^TV = I_p$$

$$\text{df}(\lambda) = \text{Tr}\left(X(X^TX + \lambda I_p)^{-1}X^T\right)$$
$$= \text{Tr}\left(UD(D^2 + \lambda I_p)^{-1}DU^T\right)$$
$$= \sum_{j=1}^{p} \frac{d_j^2}{d_j^2 + \lambda}$$

Trace equals to sum of eigenvalues

LS

# **Shrinkage Methods** – Ridge Regression

- Prostate cancer example
  - #training($N$) = 67, #testing=30
  - #variables($p$)=8
  - ridge coefficient estimates
- *Effective degree of freedom*

$$\text{df}(\lambda) = \sum_{j=1}^{p} \frac{d_j^2}{d_j^2 + \lambda} \in (0, p]$$

  - $\lambda \to 0, \text{df}(\lambda) = p$ ← no regularization
  - $\lambda \to \infty, \text{df}(\lambda) \to 0$



no regularization

strongest regularization

df = 5
(cross validation)

$(\beta)$

lcavol $\beta_8$

svi
lweight
pgg45

lbph

gleason $\beta_3$

age $\beta_2$

lcp $\beta_1$

$X_1 \quad X_2 \quad X_3 \quad \cdots \quad X_8$
$\beta_1 \quad \beta_2 \quad \beta_3 \quad \cdots \quad \beta_8$
lcp $\quad$ age $\quad$ gleason $\qquad$ lcavol

$\left( y = \beta_0 + \sum_{j=1}^{8} X_j \beta_j \right)$

$$Var(z_j) = \frac{1}{n}(Xv_j)^T(Xv_j) = \frac{1}{n}(u_j d_j)^T(u_j d_j) = \frac{d_j^2}{n}$$

$$\left( \begin{array}{l} E[z_j] = \frac{1}{n}(Xv_j)^T \mathbb{1}_n \\ = \frac{1}{n} v_j^T X^T \mathbb{1}_n = 0 \end{array} \right)$$

# **Shrinkage Methods** – Ridge Regression

Data centering: $X^T \mathbb{1}_n = 0_d$

- **Principal components** in **X**

1. - Sample covariance

   $$X = UDV^T$$

   $$S = \frac{1}{N-1}X^T X = \frac{1}{N-1}VD^2 V^T$$

2. - *Eigen decomposition* of $X^T X$

   - The eigenvector $v_j$ → The $j$-th column of **V**

     - principal components directions of **X**
     - $z_1 = Xv_1$: the first principal component

$$Var(z_j) = Var(Xv_j)$$
$$= Var(u_j d_j)$$
$$= \frac{d_j^2}{N} u_j^T u_j$$
$$= \frac{d_j^2}{N}$$

- $z_1$ has the largest variance
- $z_p$ has the smallest variance



shrinks the coefficients of the low-variance components more than the high-variance components.

11

# Linear Methods for Regression

--- The Lasso

# **Shrinkage Methods** – The Lasso

- The lasso estimate:

$$\hat{\beta}^{\text{lasso}} = \underset{\beta}{\operatorname{argmin}} \left\{ \frac{1}{2} \sum_{i=1}^{N} \left( y_i - \beta_0 - \sum_{j=1}^{p} x_{ij}\beta_j \right)^2 + \lambda \sum_{j=1}^{p} |\beta_j| \right\}.$$

training error — model complexity

$\ell_1$-norm on $\beta$

$$\|\beta\|_1 = \sum_{j=1}^{p} |\beta_j|$$

- the $\ell_2$ ridge penalty is replaced by $\ell_1$ lasso penalty.
- no closed-form solution ($\ell_1$ penalty is nondifferentiable)

- Or equivalently,

$$\hat{\beta}^{\text{lasso}} = \underset{\beta}{\operatorname{argmin}} \sum_{i=1}^{N} \left( y_i - \beta_0 - \sum_{j=1}^{p} x_{ij}\beta_j \right)^2$$

Constraint optimization

$$\text{subject to } \sum_{j=1}^{p} |\beta_j| \leq t.$$

$$\sum_{j=1}^{p} \mathbb{1}_{(\beta_j \neq 0)} \leq t$$

$$\|\beta\|_0$$

- if $t \geq \|\hat{\beta}^{ls}\|_1$, $\hat{\beta}^{lasso} = \hat{\beta}^{ls}$
- if $t = \frac{1}{2}\|\hat{\beta}^{ls}\|_1$, $\hat{\beta}^{ls}$ is shrunk about 50% on average

$\beta_2$

$\beta_1$

- making $t$ sufficiently small → some coefficients equal to 0

# Shrinkage Methods – The Lasso

- The lasso in matrix form

$$\hat{\beta}^{lasso} = \text{argmin}_\beta \frac{1}{2}\|\mathbf{y} - \mathbf{X}\beta\|_2^2 + \lambda\|\beta\|_1$$
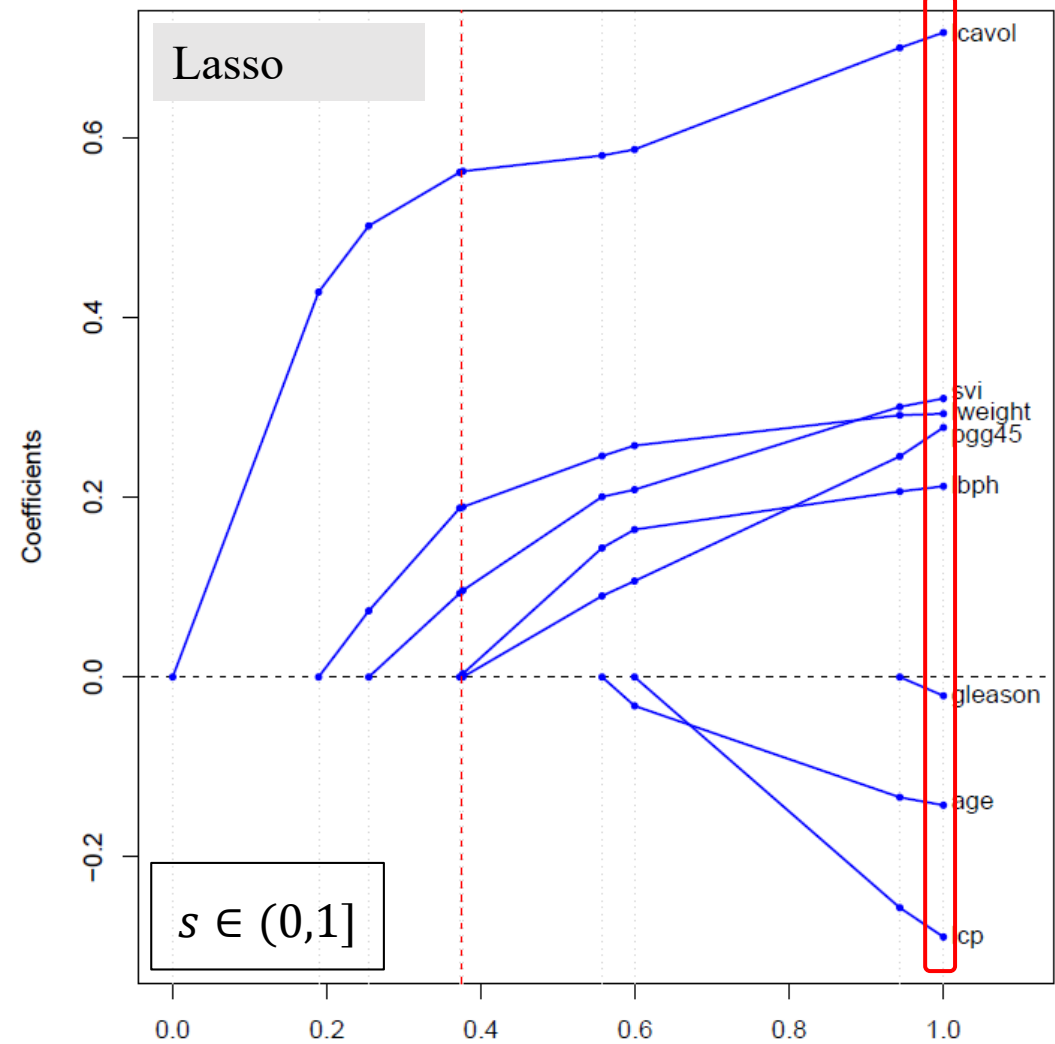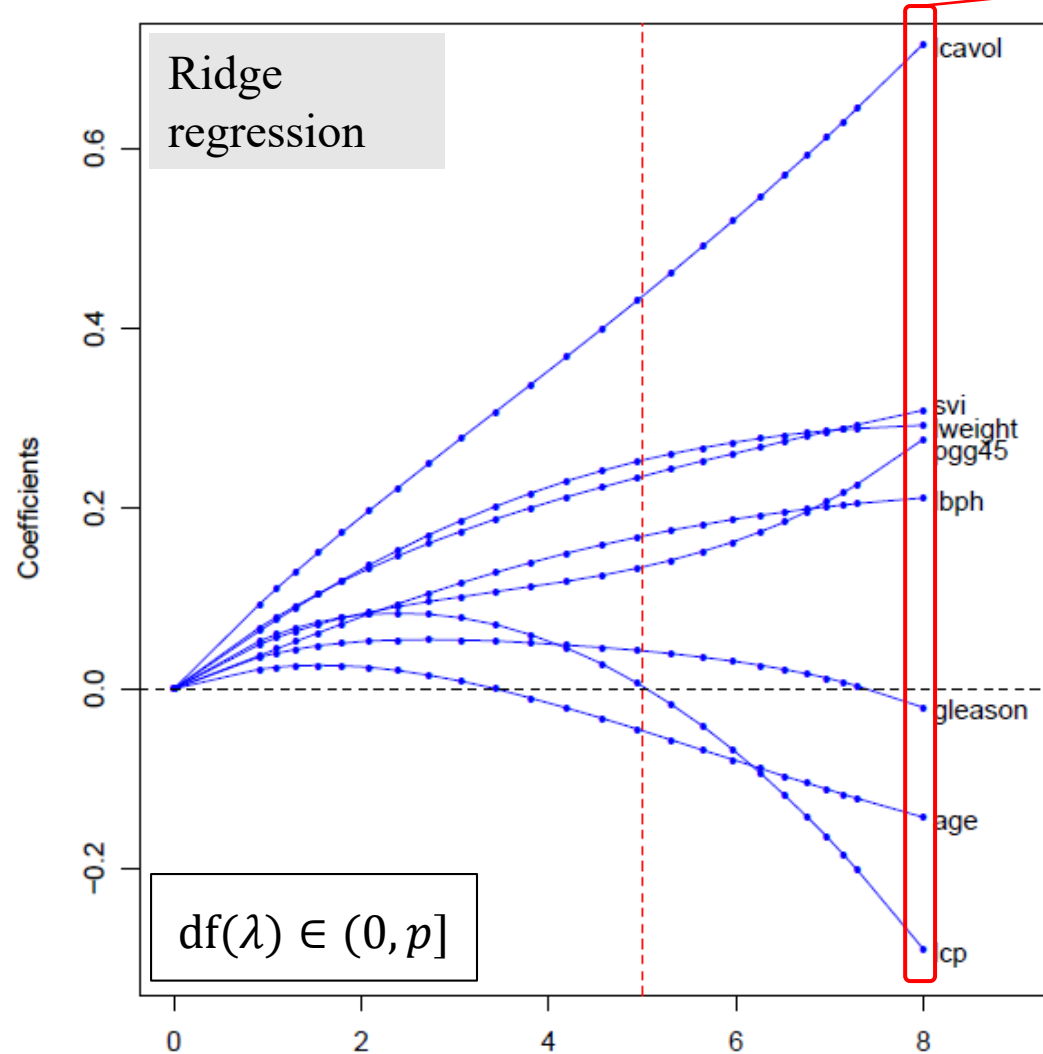
- Prostate cancer example

- The standardized parameter

$$s = t/\|\hat{\beta}^{ls}\|_1 \in (0,1]$$

- $s = 1, \hat{\beta}^{lasso} = \hat{\beta}^{ls}$
- $s \to 0, \hat{\beta}^{lasso} \to 0$
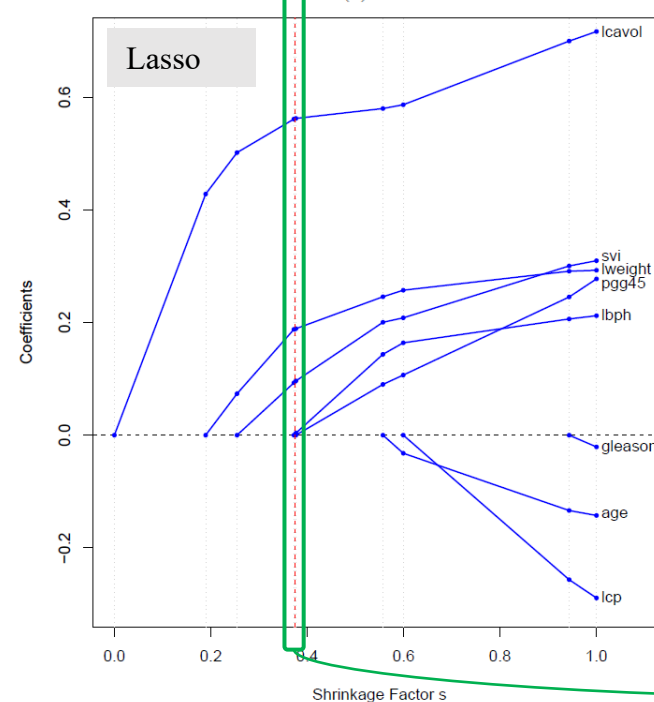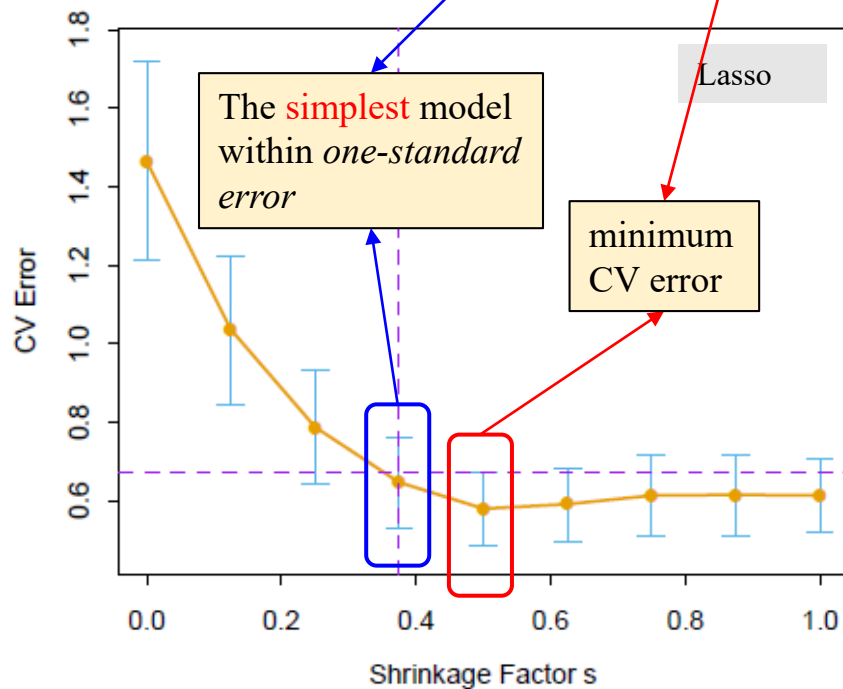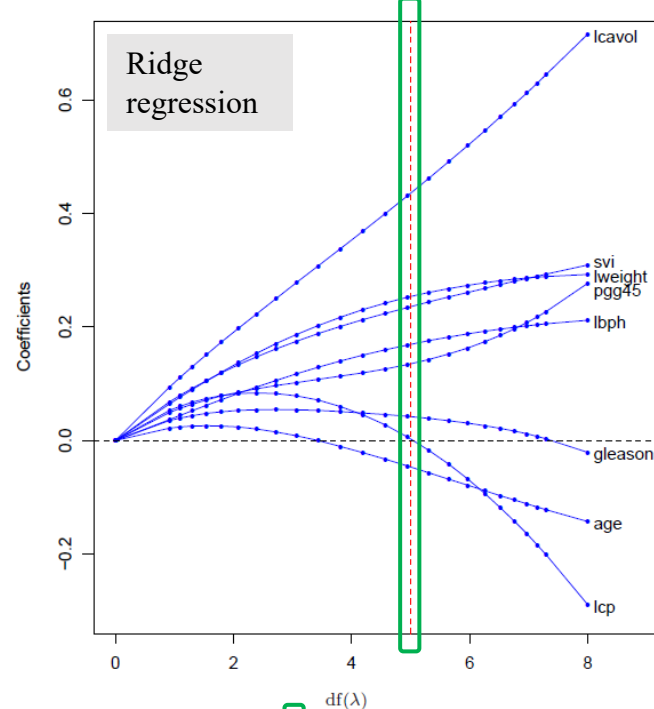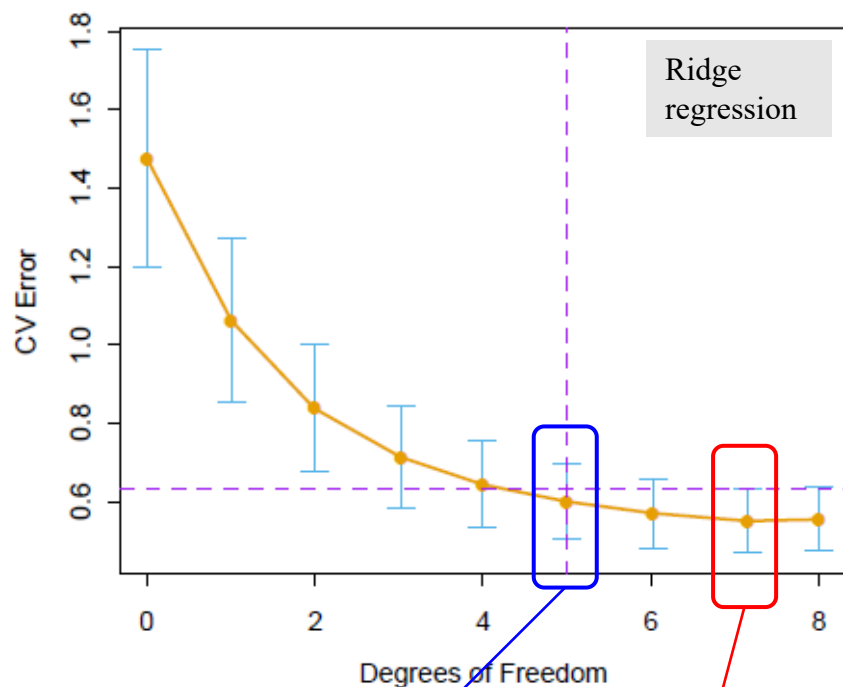- $s \in (0,1), \hat{\beta}_j^{lasso} \in (0, \hat{\beta}_j^{ls}), \forall j$



$s = 0.36$ selected by cross validation

# **Shrinkage Methods** – The Lasso



Least squares

Ridge regression

Lasso

$\mathrm{df}(\lambda) \in (0, p]$

$s \in (0,1]$

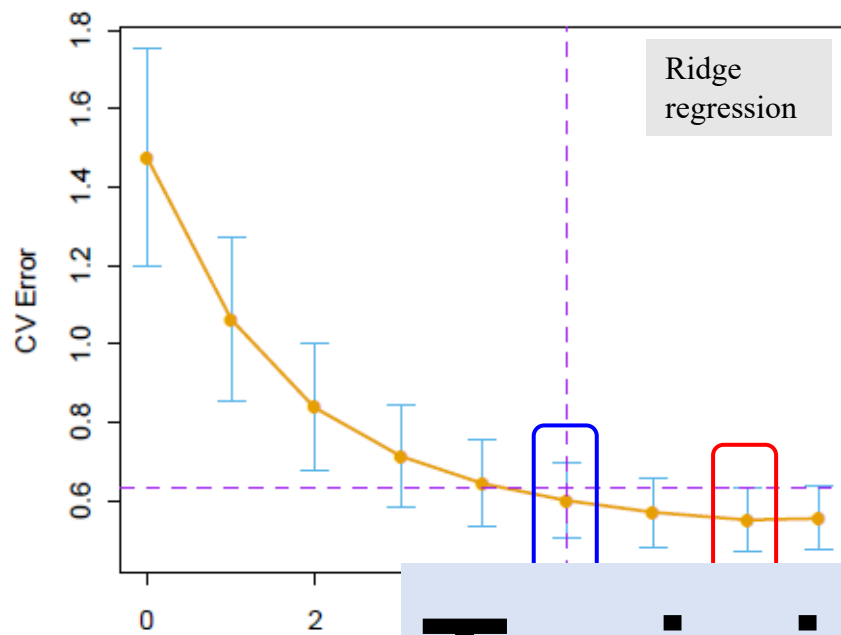Difference: the lasso profiles hit zero, while those for ridge do not.

The simplest model within *one-standard error*

minimum CV error

$df(\lambda) = 5$

$s = 0.36$

| Term | LS | Ridge | Lasso |
|---|---|---|---|
| lcavol | 0.680 | 0.420 | 0.533 |
| lweight | 0.263 | 0.238 | 0.169 |
| age | −0.141 | −0.046 | |
| lbph | 0.210 | 0.162 | 0.002 |
| svi | 0.305 | 0.227 | 0.094 |
| lcp | −0.288 | 0.000 | |
| gleason | −0.021 | 0.040 | |
| pgg45 | 0.267 | 0.133 | |
| Test Error | 0.521 | 0.492 | 0.479 |
| Std Error | 0.179 | 0.165 | 0.164 |

16

**Training set**

**Testing set**

Ridge regression

Ridge regression

Lasso

The simplest model within *one-standard error*

minimum CV error

CV Error

Degree

Shrinkage Factor s

Coefficients

lcavol
svi
lweight
pgg45
lbph
gleason
age
lcp

$df(\lambda) = 5$

| Term | LS | Ridge | Lasso |
|---|---|---|---|
| lcavol | 0.680 | 0.420 | 0.533 |
| lweight | 0.263 | 0.238 | 0.169 |
| lcp | 0.285 | 0.058 | |
| gleason | −0.021 | 0.040 | |
| pgg45 | 0.267 | 0.133 | |
| Test Error | 0.521 | 0.492 | 0.479 |
| Std Error | 0.179 | 0.165 | 0.164 |

- Biased linear methods achieved a better var-bias trade-off
- CV is usually time-consuming
  - e.g. given $s \in [0.1 : 0.1 : 1]$, we need to train the lasso by $10 \times 10 = 100$ times in 10-fold CV.
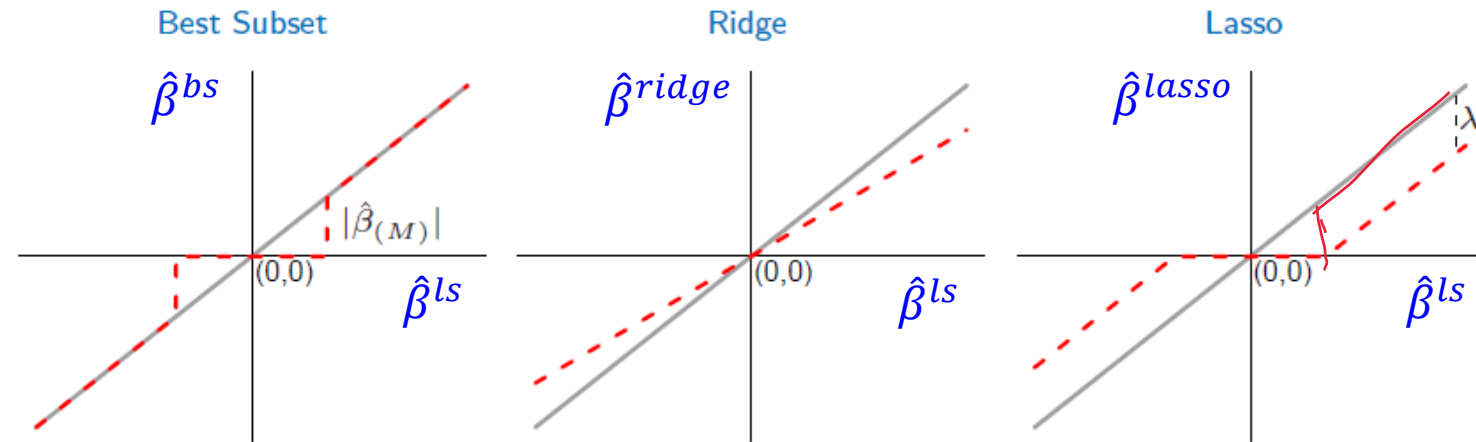
# Linear Methods for Regression

--- Discussion

# **Shrinkage Methods** – Discussion

Orthonormal case ($\mathbf{X}^T\mathbf{X} = \mathbf{I}_p$)

- Best-subset
  - hard-thresholding
  - discontinuity
- Ridge regression
  - proportional shrinkage
- Lasso
  - soft-thresholding



| Estimator | Formula |
|-----------|---------|
| Best subset (size $M$) | $\hat{\beta}_j \cdot I(|\hat{\beta}_j| \geq |\hat{\beta}_{(M)}|)$ |
| Ridge | $\hat{\beta}_j / (1 + \lambda)$ |
| Lasso | $\text{sign}(\hat{\beta}_j)(|\hat{\beta}_j| - \lambda)_+$ |

In this table $\hat{\beta}_j$ represents $\hat{\beta}_j^{ls}$

# **Shrinkage Methods** – Discussion

| Estimator | Formula |
| --- | --- |
| Best subset (size $M$) | $\hat{\beta}_j \cdot I(|\hat{\beta}_j| \geq |\hat{\beta}_{(M)}|)$ |
| Ridge | $\hat{\beta}_j/(1+\lambda)$ |
| Lasso | $\text{sign}(\hat{\beta}_j)(|\hat{\beta}_j| - \lambda)_+$ |

Orthonormal case ($\mathbf{X}^T\mathbf{X} = \mathbf{I}_p$)

- Least squares
$$\hat{\beta}^{ls} = (\mathbf{X}^T\mathbf{X})^{-1}\mathbf{X}^T\mathbf{y} = \mathbf{X}^T\mathbf{y}$$

- Ridge regression
$$\hat{\beta}^{ridge} = (\mathbf{X}^T\mathbf{X} + \lambda\mathbf{I}_p)^{-1}\mathbf{X}^T\mathbf{y}$$
$$= \frac{1}{1+\lambda}\mathbf{X}^T\mathbf{y} = \frac{1}{1+\lambda}\hat{\beta}^{ls}$$

- Best subset
$$\hat{\beta}_j^{bs} = \mathbf{x}_j^T\mathbf{y}, \qquad \forall j$$

- Lasso
$$\text{PRSS}(\beta, \lambda) = \frac{1}{2}\|\mathbf{y} - \mathbf{X}\beta\|_2^2 + \lambda\|\beta\|_1$$
$$= \frac{1}{2}\mathbf{y}^T\mathbf{y} - \beta^T\mathbf{X}^T\mathbf{y} + \frac{1}{2}\beta^T\mathbf{X}^T\mathbf{X}\beta + \lambda\|\beta\|_1$$
$$= \frac{1}{2}\mathbf{y}^T\mathbf{y} - \beta^T\hat{\beta}^{ls} + \frac{1}{2}\beta^T\beta + \lambda\|\beta\|_1$$

- Minimizing $\text{PRSS}(\beta, \lambda)$ is equivalent to
$$\min_{\beta_j} \frac{1}{2}\beta_j^2 - \hat{\beta}_j^{ls}\beta_j + \lambda|\beta_j|, \qquad \forall j$$

- Signs of $\hat{\beta}_j$ and $\hat{\beta}_j^{ls}$ must be the same.
  - $\hat{\beta}_j > 0 \rightarrow \hat{\beta}_j = \hat{\beta}_j^{ls} - \lambda$
  - $\hat{\beta}_j \leq 0 \rightarrow \hat{\beta}_j = \hat{\beta}_j^{ls} + \lambda$

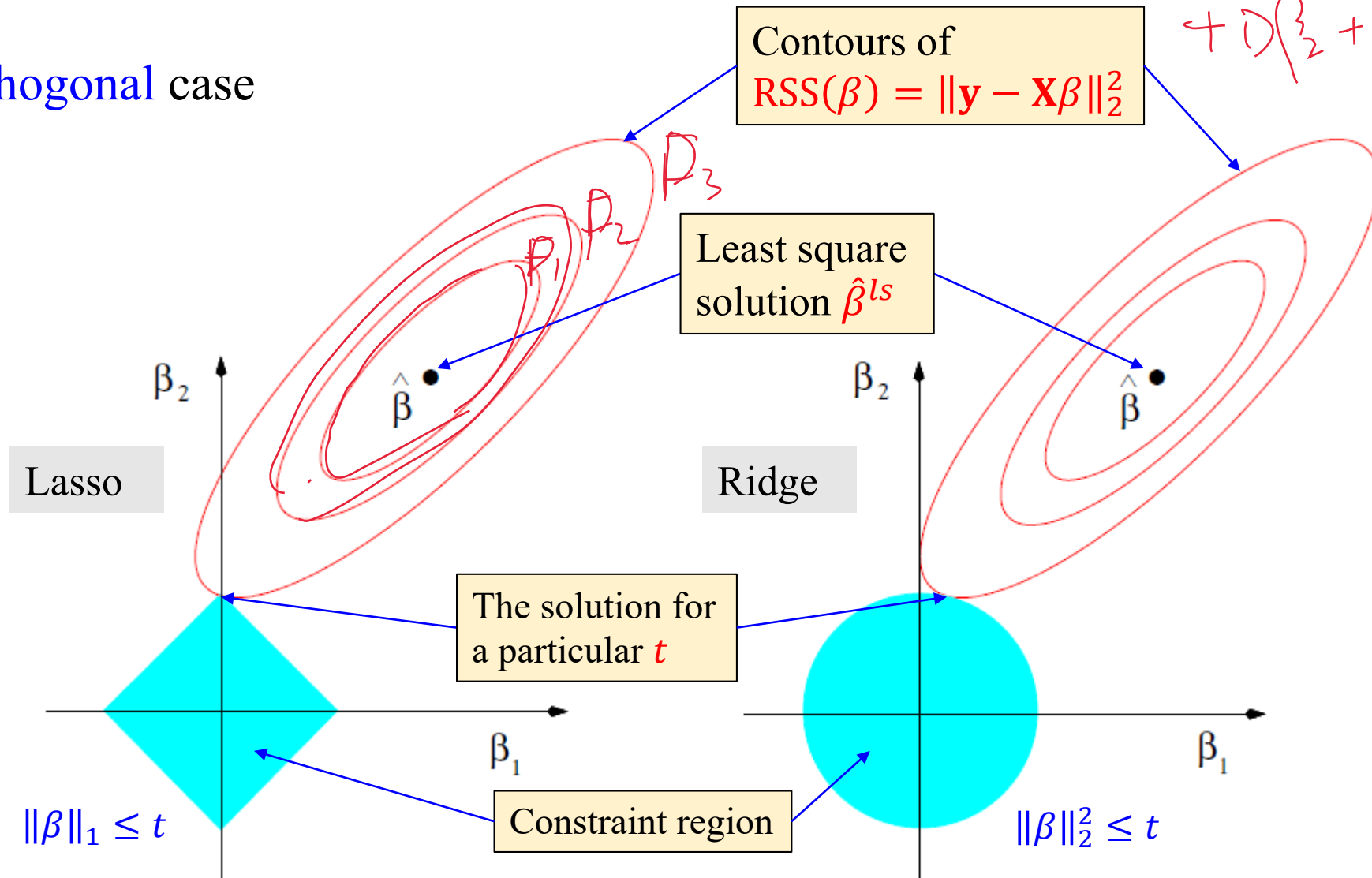- $\hat{\beta}_j^{lasso} = \text{sign}(\hat{\beta}_j^{ls})\left(|\hat{\beta}_j^{ls}| - \lambda\right)_+$

# **Shrinkage Methods** – Discussion

$$= \|y - X\beta\|_2^2$$

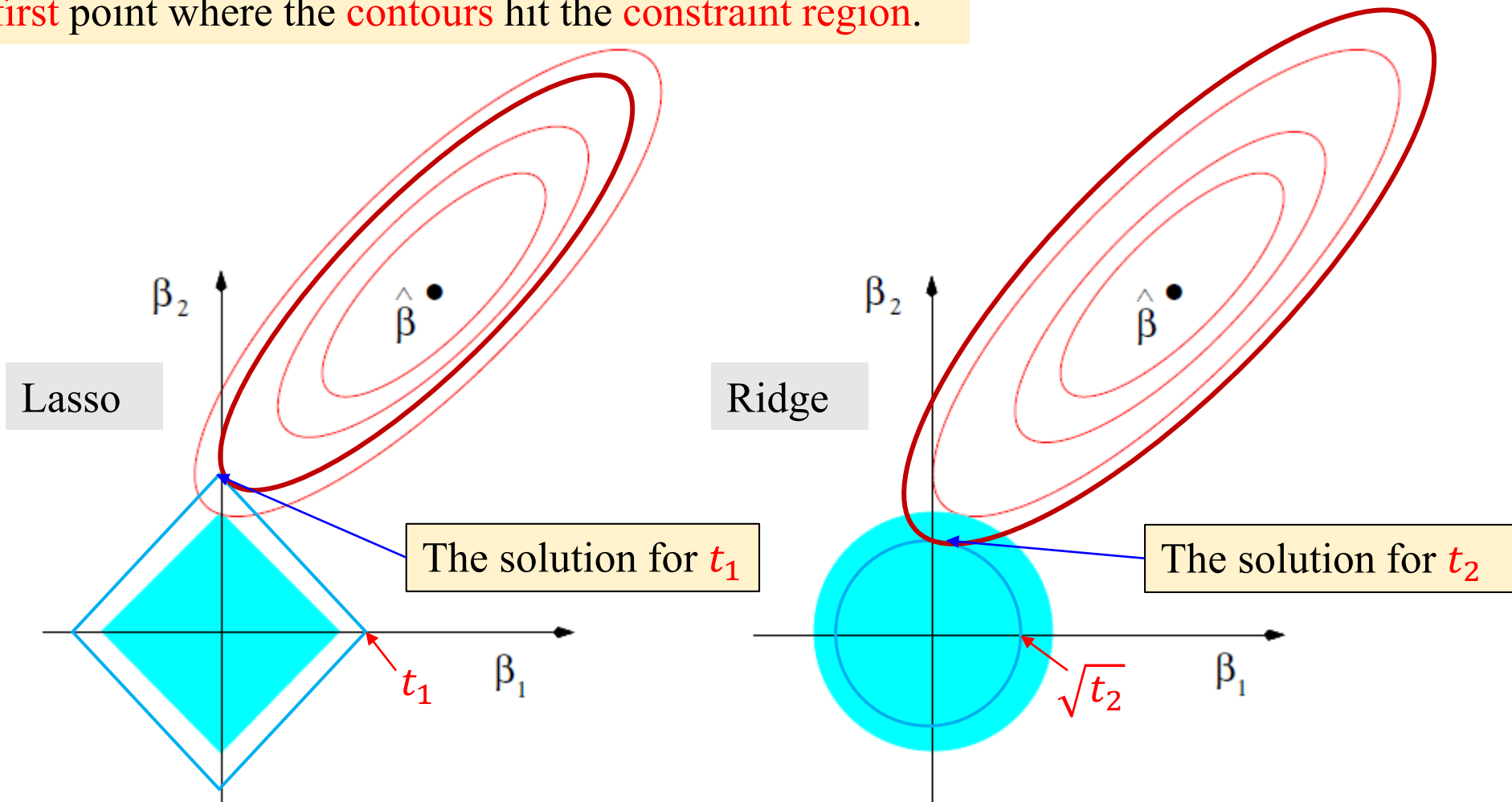$$RSS(\beta) = A\beta_1^2 + B\beta_2^2 + C\beta_1$$

$$+ D\beta_2 + E = F$$

Nonorthogonal case

Contours of
$\text{RSS}(\beta) = \|\mathbf{y} - \mathbf{X}\beta\|_2^2$

Least square
solution $\hat{\beta}^{ls}$

$P_1$ $P_2$ $P_3$

Lasso

Ridge

$\beta_2$

$\hat{\beta}$

$\beta_2$

$\hat{\beta}$

The solution for
a particular $t$

$\beta_1$

$\beta_1$

$\|\beta\|_1 \leq t$

Constraint region

$\|\beta\|_2^2 \leq t$

21

# **Shrinkage Methods** – Discussion

Lasso & Ridge regression:
Find the first point where the contours hit the constraint region.

# **Shrinkage Methods** – Discussion

$$Y = X^T\beta + \varepsilon, \quad \varepsilon \sim N(0, \delta^2)$$

Ridge and Lasso in the Bayes framework
- Suppose a Gaussian conditional distribution

$$\Pr(Y|X, \beta) = \mathcal{N}(X^T\beta, \sigma^2)$$

$$\Pr(Y|X, \beta) = \frac{1}{\sigma\sqrt{2\pi}}\exp\left(-\frac{1}{2}\left(\frac{Y - X^T\beta}{\sigma}\right)^2\right)$$

- Log-likelihood

MLE

$$\ell(\beta) = \ln \Pr(\mathbf{y}|\mathbf{X}, \beta)$$

$$= \sum_{i=1}^{N} \ln \Pr(y_i|x_i, \beta)$$

Constant $\leftarrow$ $$= \boxed{-\frac{N}{2}\log(2\pi) - N\log\sigma} - \boxed{\frac{1}{2\sigma^2}\sum_{i=1}^{N}(y_i - x_i^T\beta)^2}$$

MLE:
$$\hat{\beta}^{ls} = \text{argmax}_\beta\, \ell(\beta)$$
$$= \text{argmin}_\beta\, \|\mathbf{y} - \mathbf{X}\beta\|_2^2$$

- Maximum a posterior (MAP)

Posterior

$$\hat{\beta} = \text{argmax}_\beta\, \boxed{\Pr(\beta|\mathbf{X}, \mathbf{y})} = \text{argmax}_\beta\, \frac{\boxed{\Pr(\mathbf{y}|\mathbf{X}, \beta)}\boxed{\Pr(\beta)}}{\boxed{\Pr(\mathbf{X}, \mathbf{y})}}$$

Prior

Likelihood

Irrelevant with $\beta$

Posterior $\propto$ Likelihood $\times$ Prior

# **Shrinkage Methods** – Discussion

Ridge and Lasso in the Bayes framework

$$\text{MLE}: \hat{\beta}^{MLE} = \text{argmax}_\beta \Pr(\mathbf{y}|\mathbf{X}, \beta) \longleftarrow \text{Least squares}$$

$$\text{MAP}: \hat{\beta}^{MAP} = \text{argmax}_\beta \Pr(\mathbf{y}|\mathbf{X}, \beta)\Pr(\beta) \longleftarrow \text{Ridge \& Lasso}$$

- Ridge regression
  - MAP with a prior $\Pr(\beta) = \mathcal{N}(\beta|0, \frac{1}{\lambda}\mathbf{I}_p)$   Gaussian distribution

$$\hat{\beta}^{ridge} = \text{argmax}_\beta \ln\big(\Pr(\mathbf{y}|\mathbf{X}, \beta)\Pr(\beta)\big)$$

$$= \text{argmax}_\beta \ln\left(\prod_{i=1}^N \mathcal{N}(y_i|x_i^T\beta, \sigma^2) \times \mathcal{N}(\beta|0, \frac{1}{\lambda}\mathbf{I}_p)\right)$$

- Lasso
  - MAP with a prior $\Pr(\beta) = \frac{\lambda}{2}e^{-\lambda\|\beta\|_1}$   Laplacian distribution

$$\hat{\beta}^{lasso} = \text{argmax}_\beta \ln\left(\prod_{i=1}^N \mathcal{N}(y_i|x_i^T\beta, \sigma^2) \times \frac{\lambda}{2}e^{-\lambda\|\beta\|_1}\right)$$

# **Shrinkage Methods** – Discussion

$$\text{Ridge:} \quad ||\beta||_2^2 = \sum_{j=1}^{p} |\beta_j|^2$$

$$\text{Lasso:} \quad ||\beta||_1 = \sum_{j=1}^{p} |\beta_j|^1$$

Generalization of Ridge and Lasso

- Consider the criterion ($q \geq 0$)

$$||\beta||_q$$

$$\tilde{\beta} = \underset{\beta}{\text{argmin}} \left\{ \sum_{i=1}^{N} \left( y_i - \beta_0 - \sum_{j=1}^{p} x_{ij}\beta_j \right)^2 + \lambda \sum_{j=1}^{p} |\beta_j|^q \right\}$$
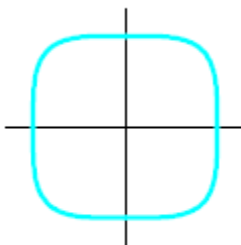
- $q = 0$, best subset
- $q = 1$, lasso
- $q = 2$, ridge regression

Unit ball:

$$||\beta||_q \leq 1$$

Convex ($q \geq 1$)

Non-convex ($q < 1$)

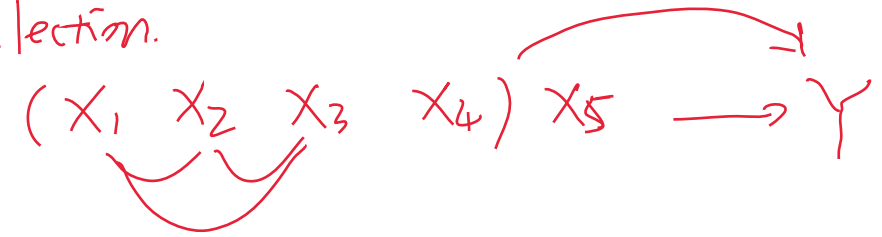| $q = 4$ | $q = 2$ | $q = 1$ | $q = 0.5$ | $q = 0.1$ |

Ridge     Lasso

- Nondifferentiable
- Penalize some coefficients to 0

*Contours of constant value of $\sum_j |\beta_j|^q$ for given values of q.*

# **Shrinkage Methods** – Discussion

*[handwritten: Group selection.]*

*[handwritten: $(X_1 \ X_2 \ X_3 \ X_4) \ X_5 \longrightarrow Y$]*

Generalization of Ridge and Lasso

*[handwritten: Ridge:]*

- Consider the criterion ($q \geq 0$)

*[handwritten: Lasso:]*

*[handwritten: E-net:]*

$$\tilde{\beta} = \underset{\beta}{\operatorname{argmin}}\left\{\sum_{i=1}^{N}(y_i - \beta_0 - \sum_{j=1}^{p}x_{ij}\beta_j)^2 + \lambda\sum_{j=1}^{p}|\beta_j|^q\right\}$$

- $q = 0$, best subset
- $q = 1$, lasso
- $q = 2$, ridge regression

*[handwritten: Trace lasso]*

- $q \in (1,2)$: a compromise between lasso and ridge regression
  - $|\beta_j|^q$ is differentiable at $0 \rightarrow$ hard to set $\beta_j = 0, \forall j$

*[handwritten: Group lasso (exclusive)]*
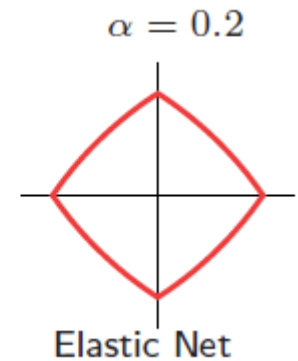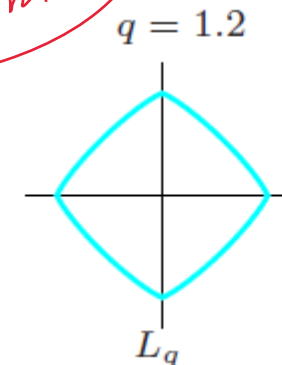
*[handwritten: K-support norm]*

- Elastic-net

$$\min_{\beta}\sum_{i=1}^{N}(y_i - \sum_{j=1}^{p}x_{ij}\beta_j)^2 + \lambda\sum_{j=1}^{p}(\alpha\beta_j^2 + (1-\alpha)|\beta_j|)$$

- $\ell_2$ shrinks the coefficients of correlated predictors
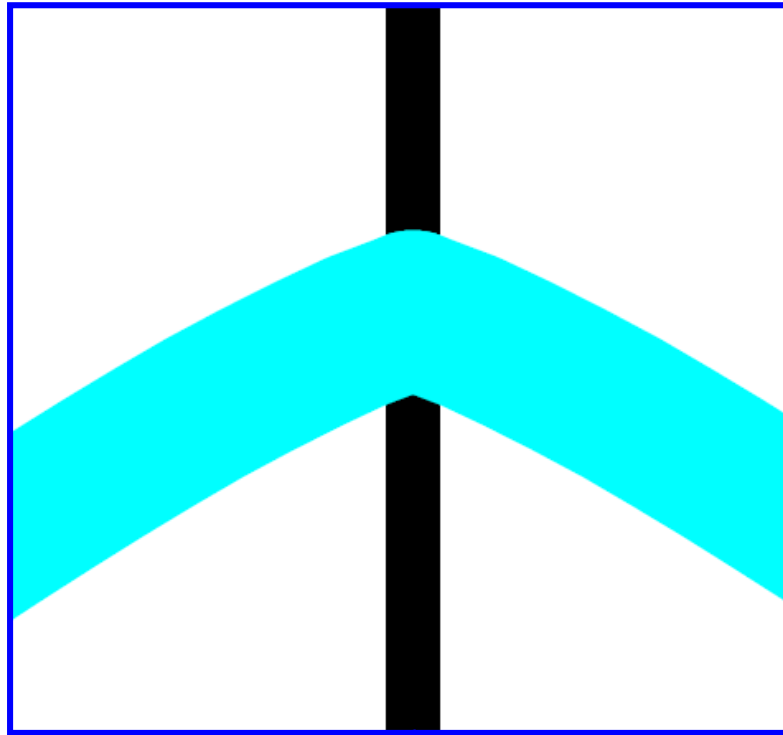- $\ell_1$ selects groups of correlated predictors

$q = 1.2$  $\alpha = 0.2$
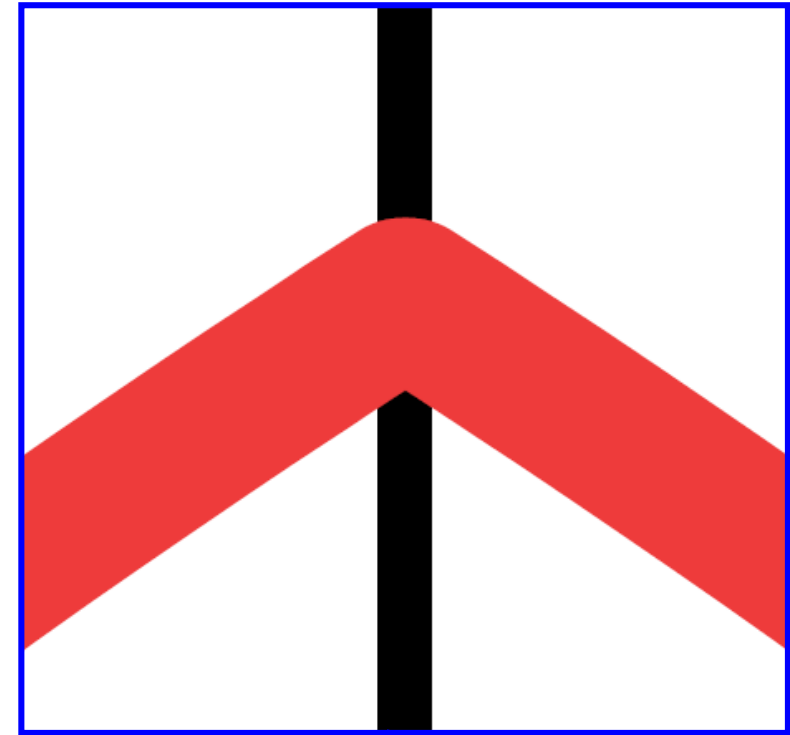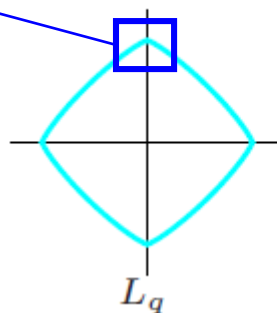
$L_q$  Elastic Net

# **Shrinkage Methods** – Discussion
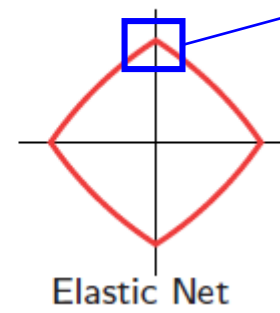


$q = 1.2$  $\alpha = 0.2$

$L_q$  Elastic Net

The elastic-net has sharp (non-differentiable) corners