

Introduction to Machine Learning CS182

Lu Sun

School of Information Science and Technology

ShanghaiTech University

October 17, 2023

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 *扔掉不需要的子集*
 - 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 *是连续的过程*
 - **continuous** process *方差会小*
 - don't suffer much from high variability
 - ridge regression, lasso, ...

Linear Methods for Regression

--- Ridge Regression

好截距, 因为惩罚 β_0 无意义.

Shrinkage Methods – Ridge Regression

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

P1

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

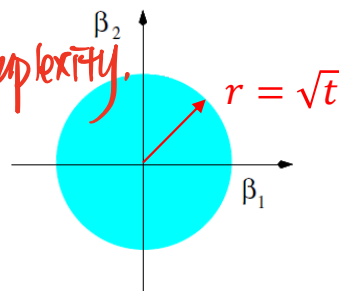
- the larger the value of λ , the greater the amount of shrinkage
- the coefficients are shrunk toward **zero**
- An equivalent expression

P2

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

$$\text{subject to } \sum_{j=1}^p \beta_j^2 \leq t$$

- One-to-one** correspondence between λ and t



正则化不包过拟合的
故要先优化

岭回归

可写为矩阵: $\|Y - \beta - \beta X\|^2$

都是二范数:
半正定矩阵
classifer 简单即可

$\lambda = 100$ 时惩罚过大, 故
 $\lambda = 0$ 时直接RSS拟合

Squared ℓ_2 -norm on β
 $\|\beta\|_2^2 = \beta^T \beta = \sum_{j=1}^p \beta_j^2$
Other possible constraints?

Shrinkage Methods – Ridge Regression *

- Equivalence between P1 and P2

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

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

- Goal: $\forall \lambda, \exists t \geq 0: \hat{\beta} = \tilde{\beta}$ (Step 1)
- $\forall t, \exists \lambda \geq 0: \hat{\beta} = \tilde{\beta}$ (Step 2)

Proof:

- Step 1: assume that P1 is solved

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

- Lagrange form of P2

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

- KKT conditions

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

- Thus,

□ if

$$t = \|\hat{\beta}\|_2^2$$

□ Then

$$\tilde{\mu} = \lambda, \quad \tilde{\beta} = \hat{\beta}$$

□ Satisfy the KKT conditions.

- Step 2: conversely, assume that P2 is solved
- The optimal solution $(\tilde{\beta}, \tilde{\mu})$ must satisfy 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



$(\tilde{\beta}, \tilde{\mu})$ satisfies KKT conditions

Shrinkage Methods – Ridge Regression

Important notes

- ridge solutions are not equivalent under **scaling of inputs**

- *standardize* the inputs before solving it

- the intercept β_0 should be **left out** of the penalty term

Ex. 3.5 → □ once $x_{ij} - \bar{x}_j$, β_0 is estimated by $\bar{y} = \frac{1}{N} \sum_{i=1}^N y_i$ $\beta_0 = \bar{y}$

- the rest parameters are estimated by the centered data \hookrightarrow 把 0 去掉.

- Henceforth we assume the data has been **standardized**

- \mathbf{X} has p rather than $p + 1$ columns

Standardization

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

Training : preprocess : $x' = \frac{x - \bar{x}}{\sigma}$ (对每列, σ 是列方差)
 Centering : $y \rightarrow y - \bar{y}$ $\beta_0 \rightarrow \bar{y}$
 train : $\min(\|y - X\beta\|^2 + \lambda \|\beta\|^2)$

Prediction?

Test : preprocess : $x_0 = \frac{x_0 - \bar{x}}{\sqrt{\text{var}(x_0)}}$ predict : $\hat{y}_i = \hat{\beta}^T X + \beta_0$

Shrinkage Methods – Ridge Regression

- Ridge regression in **matrix** form

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

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

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

$$\begin{aligned} \operatorname{PRSS}(\lambda, \beta) &= (\mathbf{y} - \mathbf{X}\beta)^T (\mathbf{y} - \mathbf{X}\beta) + \lambda \beta^T \beta \\ &= \mathbf{y}^T \mathbf{y} - \beta^T \mathbf{X}^T \mathbf{y} - \mathbf{y}^T \mathbf{X} \beta + \beta^T \mathbf{X}^T \mathbf{X} \beta + \lambda \beta^T \beta \end{aligned}$$

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

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

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

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

一定非奇异
2.2 2.2 2.2

SVD分解之后对非零的 λ_p 一定为非奇异。

(是 $p \times p$ 的矩阵)

Shrinkage Methods – Ridge Regression

Additional insight into ridge regression \rightarrow 对角矩阵, $D^T = D$.

$$\begin{aligned} X^T X &= (UDV^T)^T (UDV^T) = (VDU^T)(UDV^T) = VD^2V^T \\ X^T X + \lambda I_p &= VD^2V^T + \lambda VV^T \\ &= V(D^2 + \lambda I_p)V^T \end{aligned}$$

- Singular value decomposition (SVD) 对角矩阵.

$$U^T U = I_p, V^T V = I_p \quad X = UDV^T$$

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

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

\rightarrow 一定非奇异的.

Least squares

$$\begin{aligned} X\hat{\beta}^{ls} &= X(X^T X)^{-1} X^T y \\ &= U U^T y, \\ &= \sum_{j=1}^p \underbrace{u_j}_{\text{The } j\text{-th column of } U} u_j^T y \end{aligned}$$

The j -th column of U

Ridge regression

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

- 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 8维
 - ridge coefficient estimates
- Effective degree of freedom

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

有效自由度 \Rightarrow 复杂度

$$\begin{aligned} df(\lambda) &= \text{Tr} \left(\mathbf{X}(\mathbf{X}^T \mathbf{X} + \lambda \mathbf{I}_p)^{-1} \mathbf{X}^T \right) \\ &= \text{Tr} \left(\mathbf{U} \mathbf{D} (\mathbf{D}^2 + \lambda \mathbf{I}_p)^{-1} \mathbf{D} \mathbf{U}^T \right) \\ &= \sum_{j=1}^p \frac{d_j^2}{d_j^2 + \lambda} \end{aligned}$$

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

Trace equals to sum of eigenvalues

Shrinkage Methods – Ridge Regression

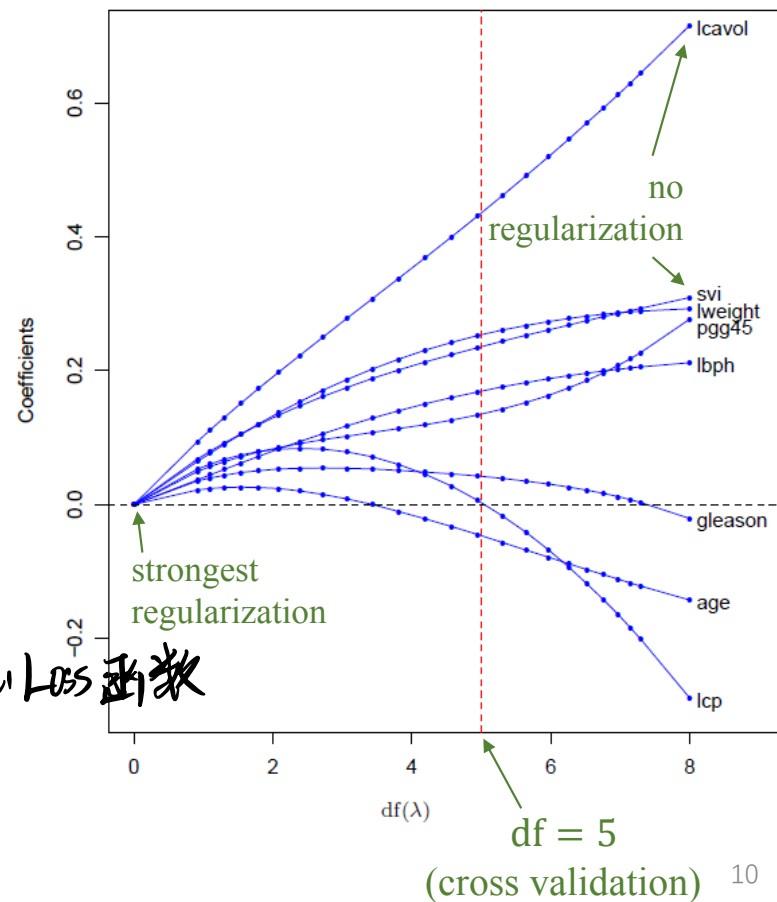
- Prostate cancer example
 - #training(N) = 67, #testing=30
 - #variables(p)=8
 - ridge coefficient estimates

- Effective degree of freedom

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

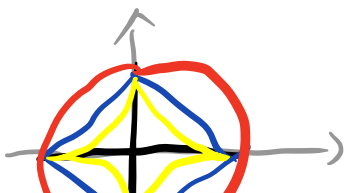
- $\lambda \rightarrow 0, df(\lambda) = p$ — no regularization
- $\lambda \rightarrow \infty, df(\lambda) \rightarrow 0$

把自由度罚没了, 不再关心Loss函数



Linear Methods for Regression

--- The Lasso



$\| \cdot \|_2$ = 范数 (convex)

$\| \cdot \|_1$ = 范数 (convex)

$\| \cdot \|_1 + \| \cdot \|_2$ = 范数 (非凸)

(距零范数最近的 convex 是一范数)

Shrinkage Methods – The Lasso

ℓ_1, ℓ_p 范数, $0 < p < 1$ (非凸)

稀疏化方法

最理想(稀疏化), 但极复杂.

$\|\beta\|_0$: 非零元的个数
从平方和 \rightarrow 绝对值之和

- The **lasso** estimate:

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

$$\ell_1\text{-norm on } \beta$$

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

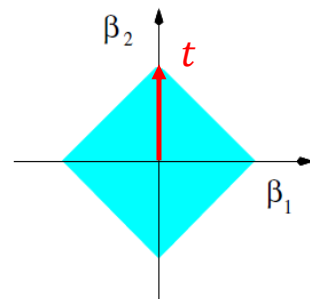
- the ℓ_2 ridge penalty is replaced by ℓ_1 lasso penalty.
- no closed-form solution (ℓ_1 penalty is **nondifferentiable**)
- Or equivalently, \hookrightarrow 在顶点处不可导, 故用梯度下降逼近

Constraint optimization

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

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

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



- making t sufficiently small \rightarrow some coefficients equal to 0

Shrinkage Methods – The Lasso

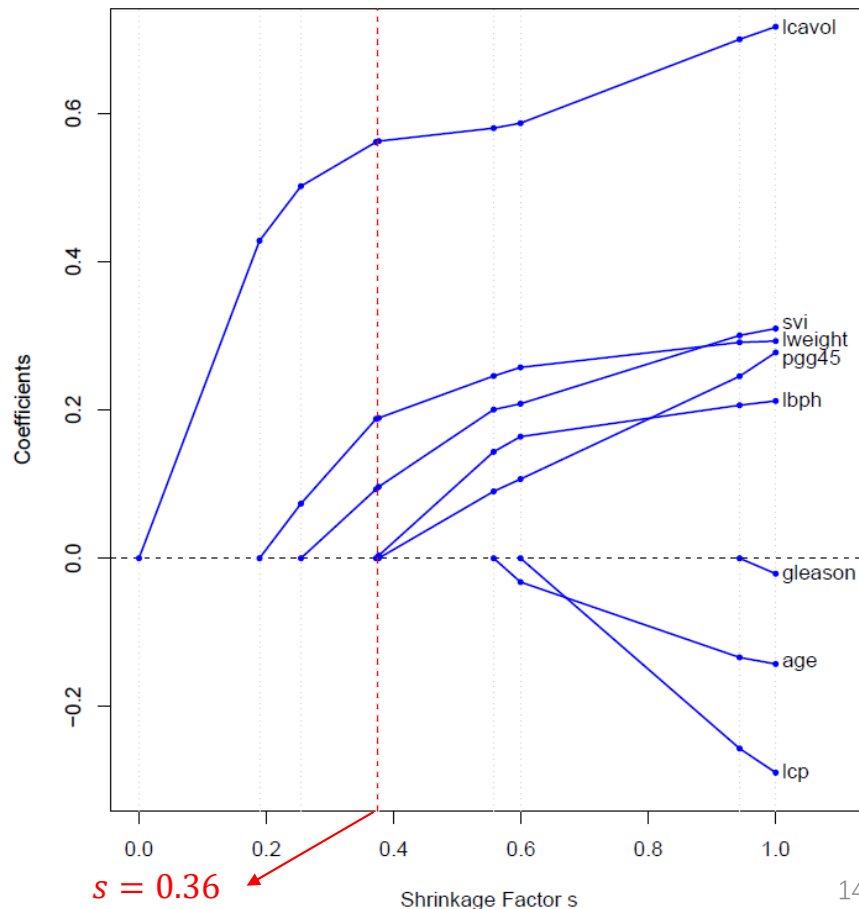
- The lasso in **matrix** form

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

- Prostate cancer example**
- The standardized parameter

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

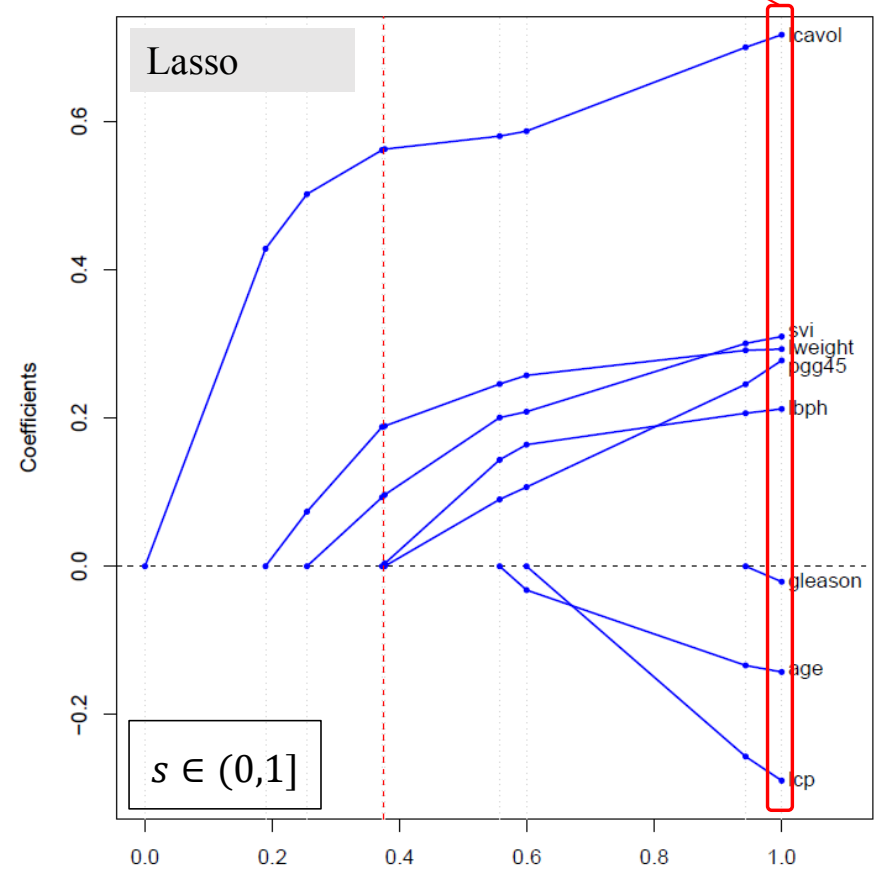
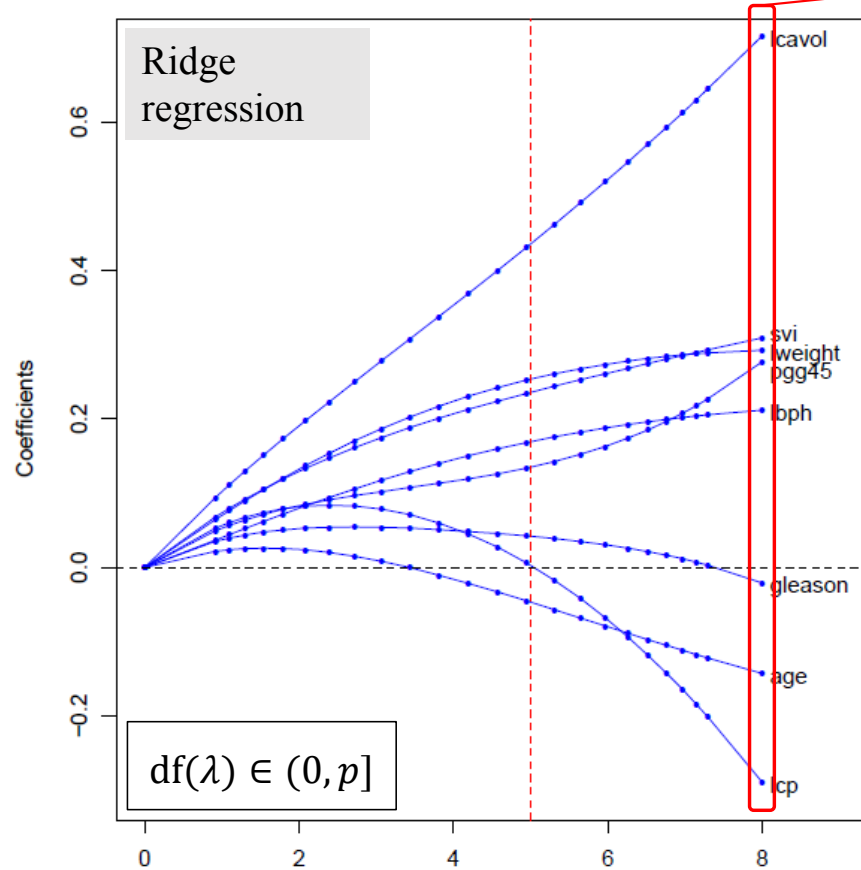
- $s = 1, \hat{\beta}^{lasso} = \hat{\beta}^{ls}$
- $s \rightarrow 0, \hat{\beta}^{lasso} \rightarrow 0$
- $s \in (0,1), \hat{\beta}_j^{lasso} \in (0, \hat{\beta}_j^{ls}), \forall j$



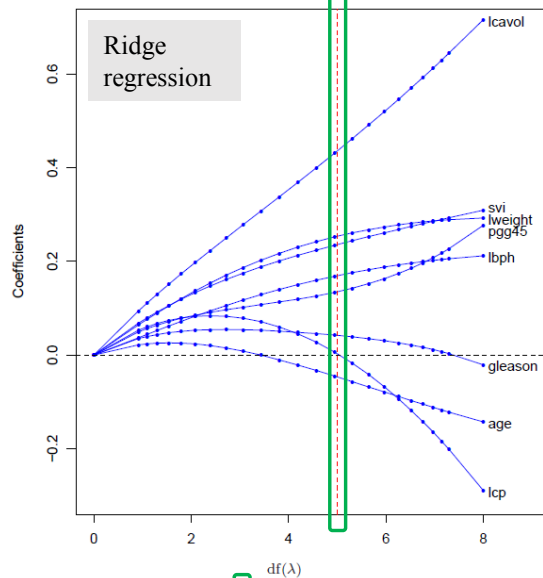
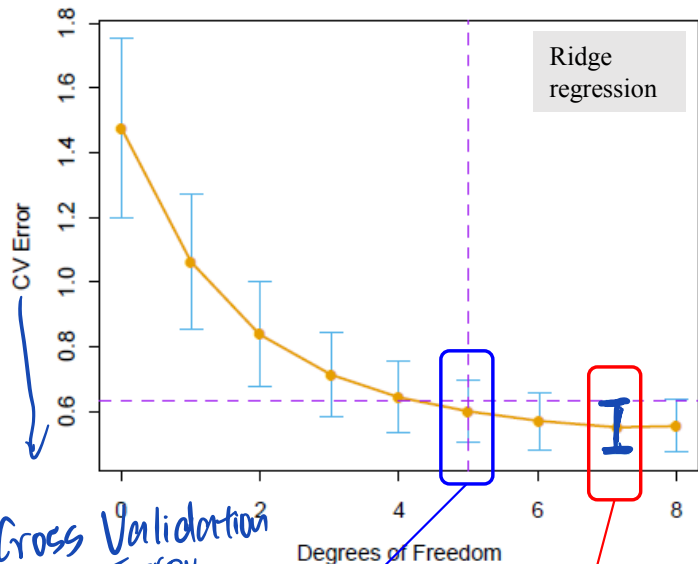
$s = 0.36$
selected by cross validation

Shrinkage Methods – The Lasso

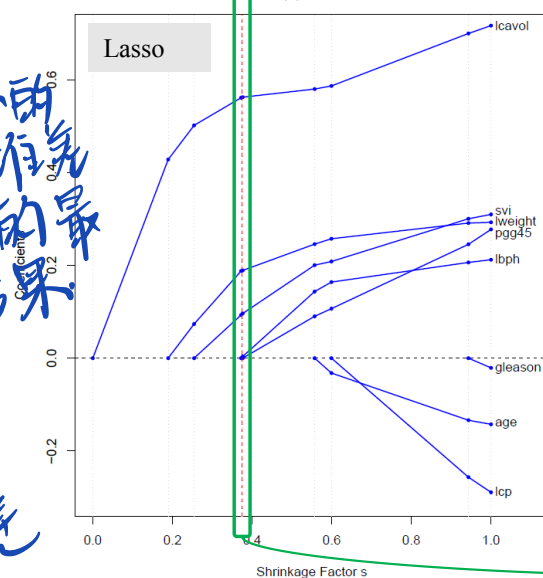
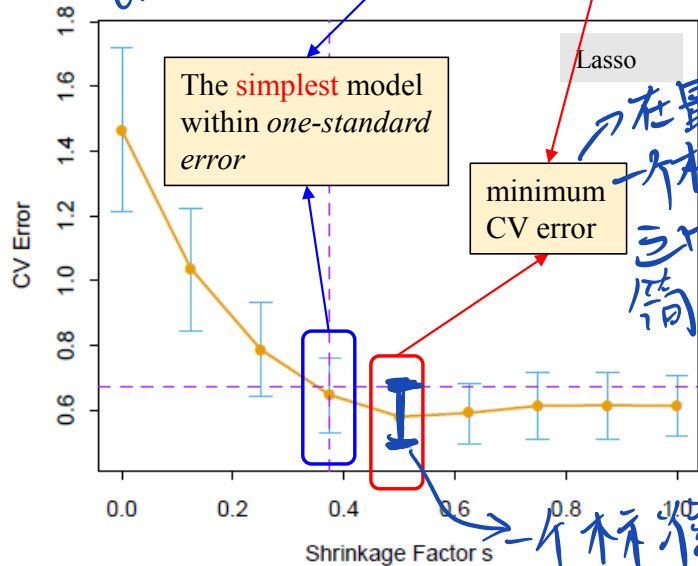
Least squares



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



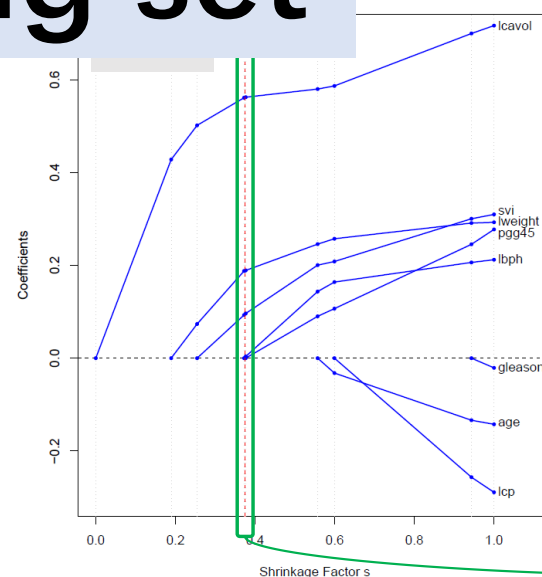
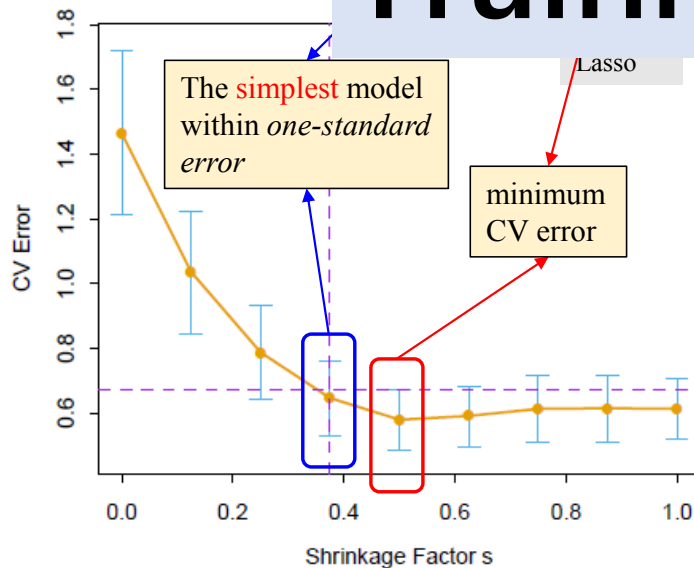
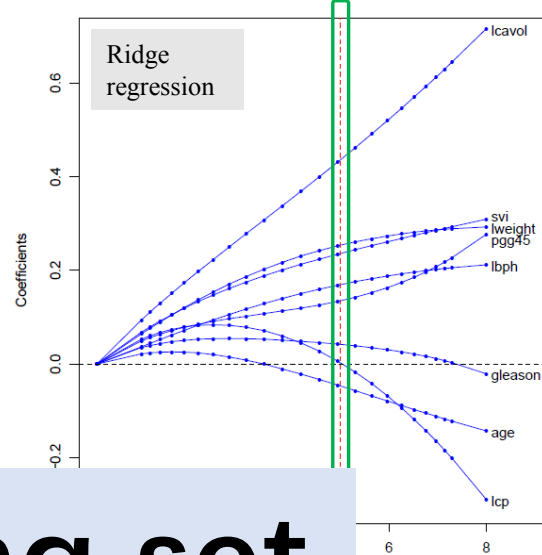
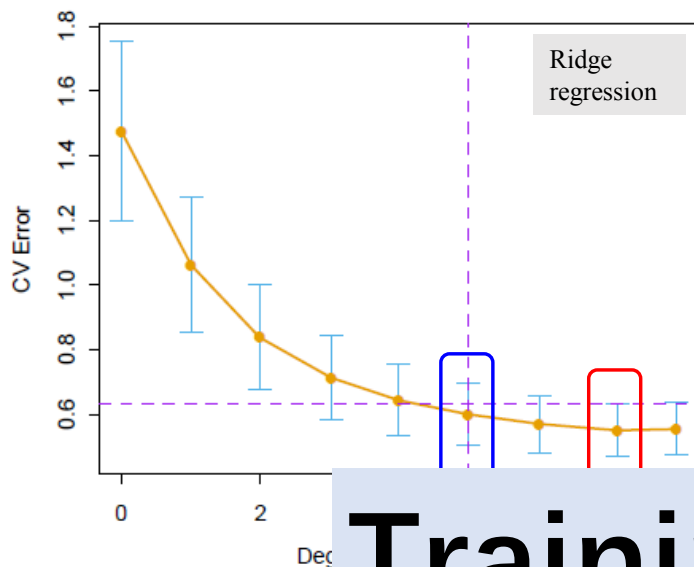
Cross Validation Error



$$df(\lambda) = 5$$

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

$$s = 0.36$$



$df(\lambda) = 5$

Term	LS	Ridge	Lasso
lccavol	0.680	0.420	0.533
lweight	0.263	0.238	0.160
lcp	0.128	0.088	0.088
gleason	-0.021	0.040	0.040
pgg45	0.267	0.133	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.

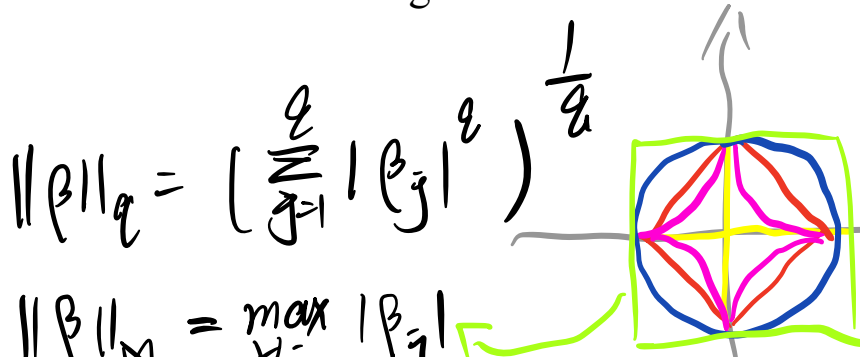
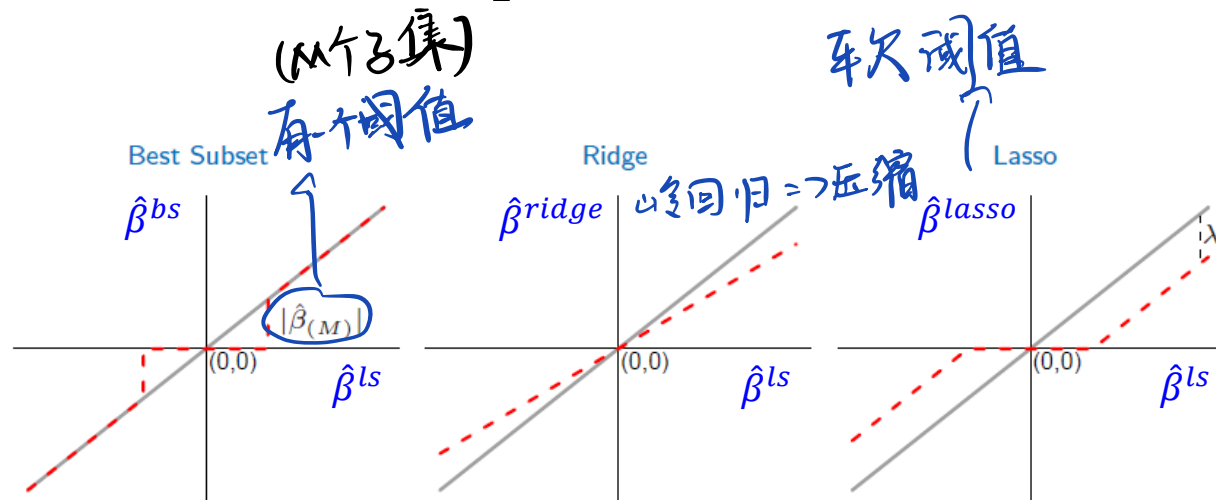
Linear Methods for Regression

--- Discussion

Shrinkage Methods – Geometric Interpretation

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 $\hat{\beta}_j^{ls}$

表示取正的部分

若小于0则取0.

Shrinkage Methods – Geometric Interpretation

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 正则化

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

$$= \frac{1}{1+\lambda} \mathbf{X}^T \mathbf{y} = \frac{1}{1+\lambda} \hat{\beta}^{ls}$$

- Best subset

$$\hat{\beta}_j^{bs} = \mathbf{x}_j^T \mathbf{y},$$

$\forall j$ 只需要把 $\hat{\beta}_j^{bs}$ 排序
找前 m 个

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	$\text{sign}(\hat{\beta}_j)(\hat{\beta}_j - \lambda)_+$

- Lasso

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

常数, 可移

$$= \frac{1}{2} \mathbf{y}^T \mathbf{y} - \beta^T \mathbf{X}^T \mathbf{y} + \frac{1}{2} \beta^T \mathbf{X}^T \mathbf{X} \beta + \lambda \|\beta\|_1$$

$$= \frac{1}{2} \mathbf{y}^T \mathbf{y} - \beta^T \hat{\beta}^{ls} + \frac{1}{2} \beta^T \beta + \lambda \|\beta\|_1$$

- Minimizing $\text{PRSS}(\beta, \lambda)$ is equivalent to

$$\min_{\beta_j} \frac{1}{2} \beta_j^2 - \hat{\beta}_j^{ls} \beta_j + \lambda |\beta_j|, \quad \forall j$$

Signs of $\hat{\beta}_j$ and $\hat{\beta}_j^{ls}$ must be the same. \rightarrow 需保证同号

- $\hat{\beta}_j > 0 \rightarrow \hat{\beta}_j = \hat{\beta}_j^{ls} - \lambda$
- $\hat{\beta}_j \leq 0 \rightarrow \hat{\beta}_j = \hat{\beta}_j^{ls} + \lambda$

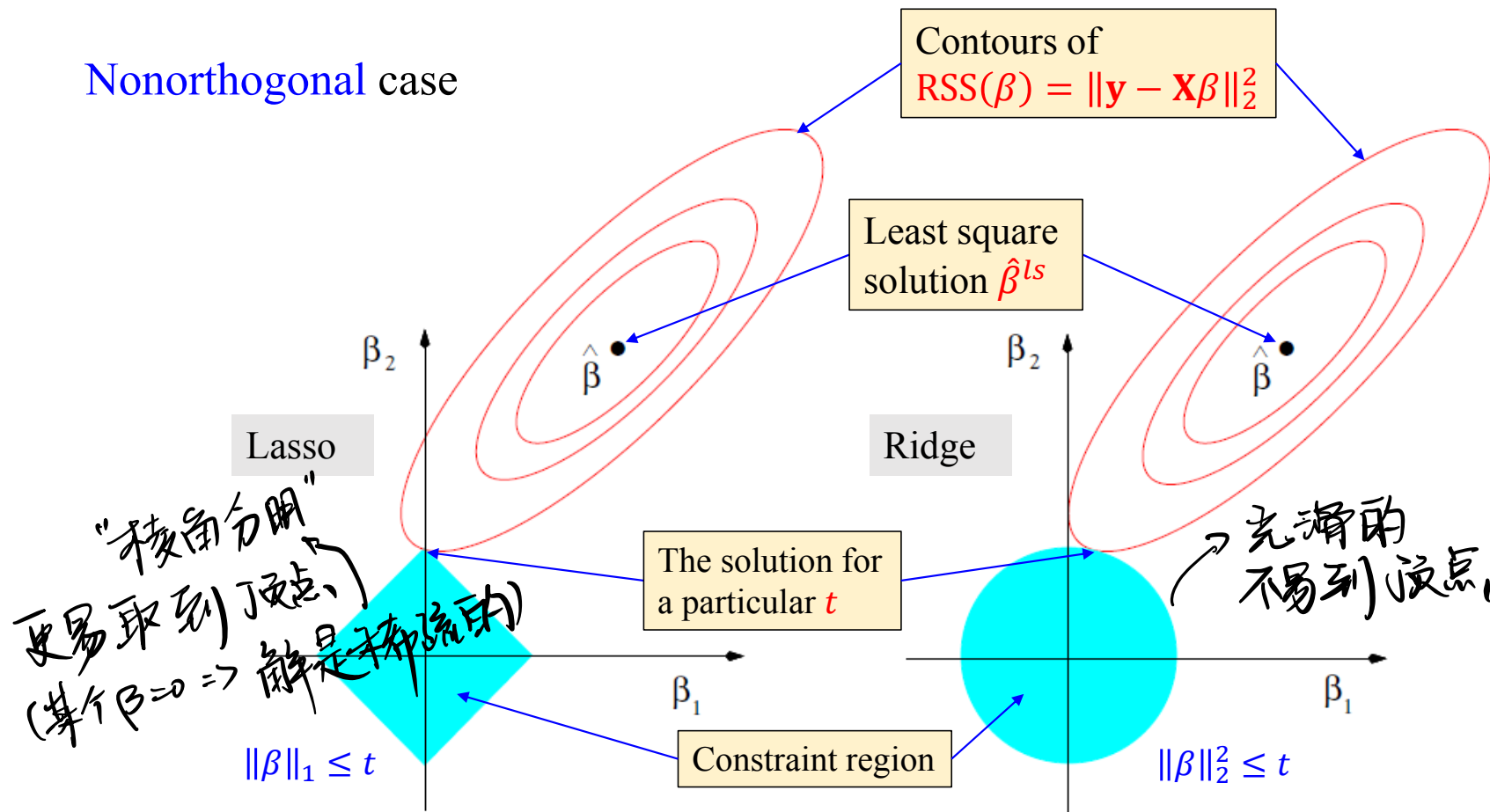
$$\hat{\beta}_j^{lasso} = \text{sign}(\hat{\beta}_j^{ls})(|\hat{\beta}_j^{ls}| - \lambda)_+$$

Partial Gradient.

Proximal operator: 优化 $\min \|y - \beta\|_2 + \lambda \|\beta\|_1$
 用来优化 Lasso

Shrinkage Methods – Geometric Interpretation

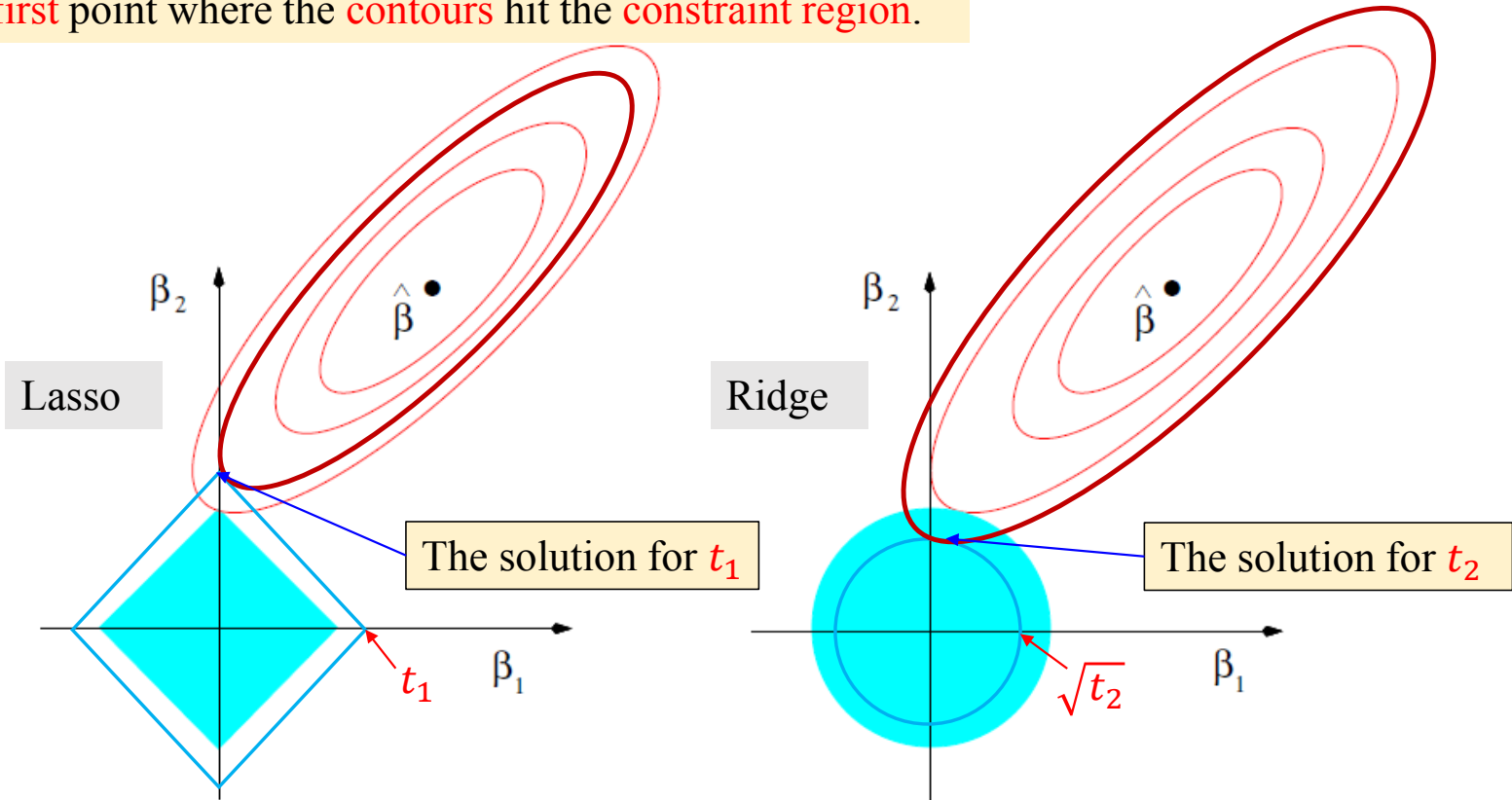
Nonorthogonal case



Shrinkage Methods – Geometric Interpretation

Lasso & Ridge regression:

Find the **first** point where the **contours** hit the **constraint region**.



Shrinkage Methods – Probabilistic Interpretation

Ridge and Lasso in the **Bayes** framework

- Suppose a Gaussian conditional distribution

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

均值 方差

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

- Log-likelihood

$$\begin{aligned} \ell(\beta) &= \ln \Pr(\mathbf{y}|\mathbf{X}, \beta) \\ &= \sum_{i=1}^N \ln \Pr(y_i|x_i, \beta) \end{aligned}$$

MLE:

$$\begin{aligned} \hat{\beta}^{ls} &= \operatorname{argmax}_{\beta} \ell(\beta) \\ &= \operatorname{argmin}_{\beta} \|\mathbf{y} - \mathbf{X}\beta\|_2^2 \end{aligned}$$

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

- Maximum a posterior (**MAP**)

$$\hat{\beta} = \operatorname{argmax}_{\beta} \underbrace{\Pr(\beta|\mathbf{X}, \mathbf{y})}_{\text{Posterior}} = \operatorname{argmax}_{\beta} \frac{\underbrace{\Pr(\mathbf{y}|\mathbf{X}, \beta)}_{\text{Likelihood}} \underbrace{\Pr(\beta)}_{\text{Prior}}}{\underbrace{\Pr(\mathbf{X}, \mathbf{y})}_{\text{Irrelevant with } \beta}}$$

$$\text{Posterior} \propto \text{Likelihood} \times \text{Prior}$$

Shrinkage Methods – Probabilistic Interpretation

Ridge and Lasso in the **Bayes** framework

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

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

- Ridge regression

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

$$\begin{aligned} \hat{\beta}^{ridge} &= \operatorname{argmax}_{\beta} \ln(\Pr(\mathbf{y}|\mathbf{X}, \beta) \Pr(\beta)) \\ &= \operatorname{argmax}_{\beta} \ln\left(\prod_{i=1}^N \mathcal{N}(y_i|x_i^T \beta, \sigma^2) \times \mathcal{N}(\beta|0, \frac{1}{\lambda} \mathbf{I}_p)\right) \end{aligned}$$

- Lasso

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

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

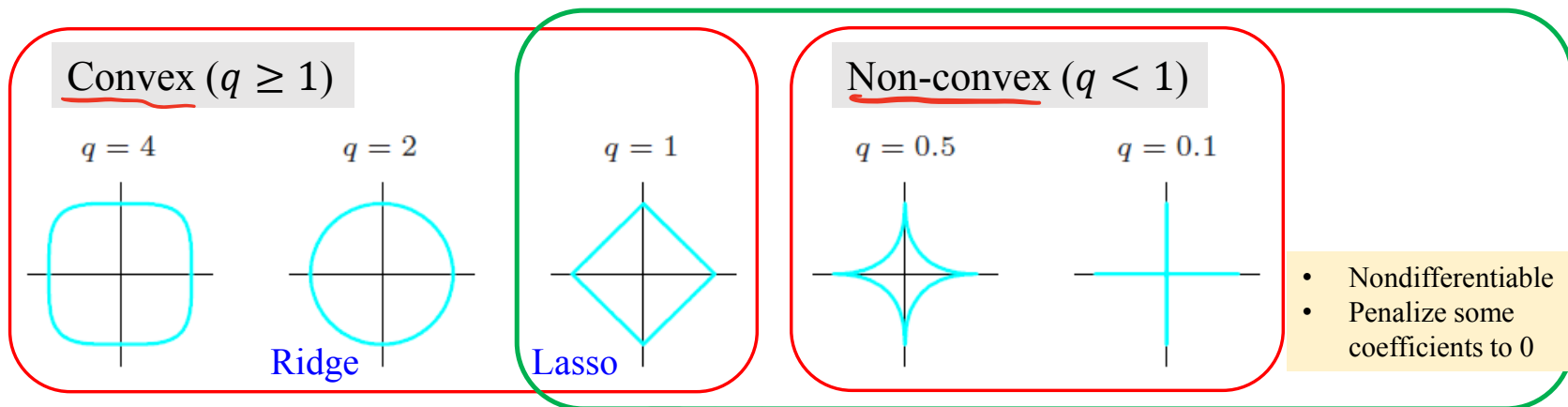
Shrinkage Methods – Generalization

Generalization of Ridge and Lasso

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

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

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



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

若 $q < 1$ 则无距离上的意义 (如: 欧氏距离 $q=2$)
 但求导 \rightarrow 某些系数为 0

因为不满足三角形不等式 $\|x+y\| \leq \|x\| + \|y\|$

Shrinkage Methods – Generalization

岭回归：若相关则成一个组。 L_{lasso} ：无关的置零。

Generalization of Ridge and Lasso

相关的随机抽一个放
剩下的置0。

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

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

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

- $q \in (1,2)$: a compromise between lasso and ridge regression

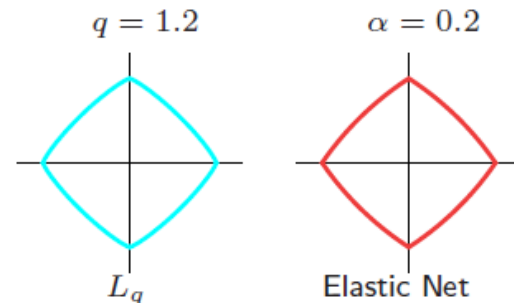
$|\beta_j|^q$ is differentiable at 0 \rightarrow hard to set $\beta_j = 0, \forall j$

- Elastic-net

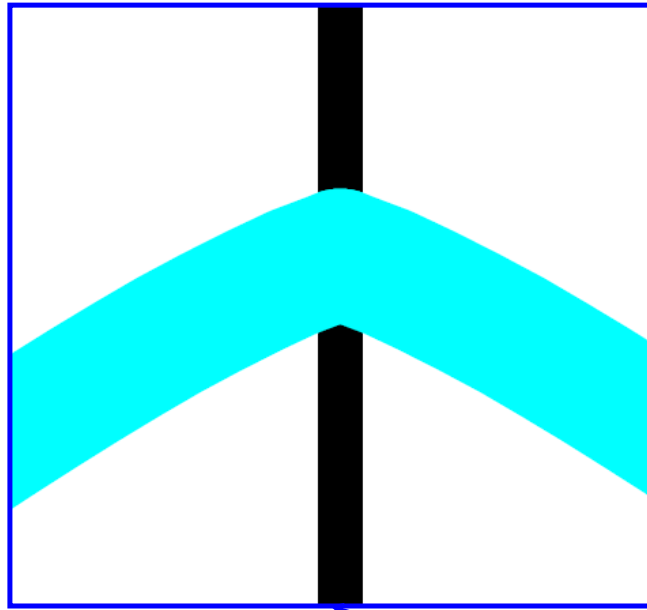
$$\|\cdot\|_2^2 + \lambda \|\cdot\|_1$$

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

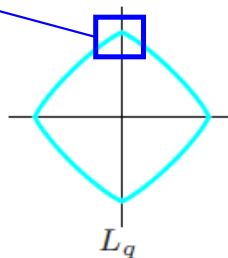
- ℓ_2 shrinks the coefficients of correlated predictors
- ℓ_1 selects groups of correlated predictors



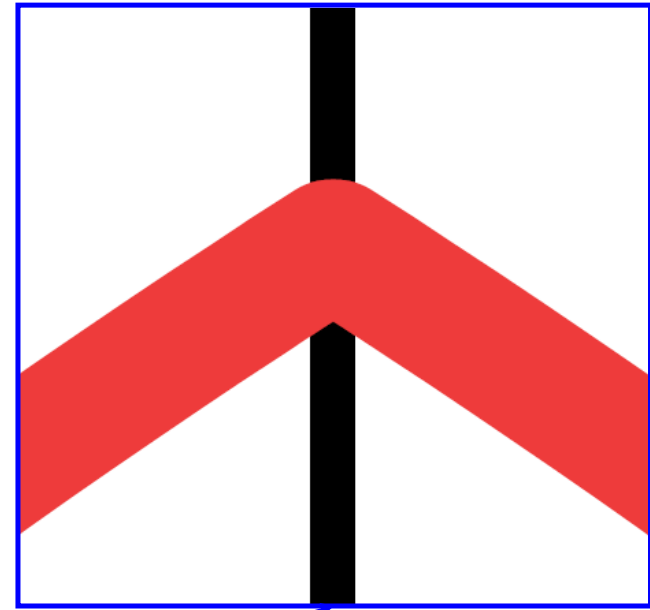
Shrinkage Methods – Generalization



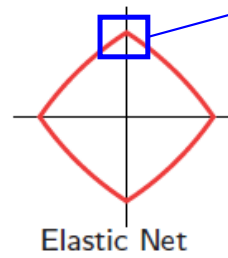
$q = 1.2$



L_q



$\alpha = 0.2$



Elastic Net

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