



Business Analytics

Week 8
Advanced regression

1 May 2018

Outline

Week	Topic	Chapter	Lecturer
1	Introduction	1	Souhaib
2	Statistical learning	2	Souhaib
3	Regression	3	Souhaib
4	Classification	4	Souhaib
5	Clustering	10	Souhaib
Semester break			
6	Model selection and resampling methods	5	Souhaib
7	Dimension reduction	6,10	Souhaib
8	Dimension reduction	6	Souhaib
9	Advanced regression	6	Souhaib
10	Advanced classification	9	Souhaib
11	Tree-based methods	8	Souhaib
12	Project presentation		Souhaib

Regression

$$Y = f(X) + \varepsilon,$$

where $X = (X_1, \dots, X_p)$, $\mathbb{E}[\varepsilon] = 0$ and $\mathbb{E}[\varepsilon^2|X] = \sigma^2$.

$$\underset{h \in \mathcal{H}}{\text{minimize}} \mathbb{E}[(Y - h(X))^2]$$

Linear regression

$$\mathcal{H} = \left\{ h \mid h(\mathbf{X}) = \beta_0 + \sum_{j=1}^p \mathbf{X}_j \beta_j = \boldsymbol{\beta}' \mathbf{X} \right\}$$
$$\underset{\boldsymbol{\beta} \in \mathbb{R}^p}{\text{minimize}} \left\{ \sum_{i=1}^n \left(y_i - \beta_0 - \sum_{j=1}^p \beta_j \mathbf{X}_{ij} \right)^2 \right\}$$

$$\equiv \underset{\boldsymbol{\beta} \in \mathbb{R}^p}{\text{minimize}} (\mathbf{y} - \mathbf{X}\boldsymbol{\beta})'(\mathbf{y} - \mathbf{X}\boldsymbol{\beta})$$

$$\equiv \underset{\boldsymbol{\beta} \in \mathbb{R}^p}{\text{minimize}} \|\mathbf{y} - \mathbf{X}\boldsymbol{\beta}\|_2^2$$

$$\hat{\boldsymbol{\beta}}^{\text{ls}} = (\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{y}$$

Bias and variance in linear regression

Let us assume

$$\mathbf{y} = \mathbf{X}\boldsymbol{\beta}^* + \boldsymbol{\varepsilon},$$

where $\boldsymbol{\varepsilon} \sim \mathcal{N}_n(\mathbf{0}, \sigma^2 \mathbf{I})$.

- What is the bias of $\hat{\boldsymbol{\beta}}^{\text{ls}}$?

$$\mathbb{E}[\hat{\boldsymbol{\beta}}^{\text{ls}}] - \boldsymbol{\beta}^* = \mathbf{0}$$

- What is the variance of $\hat{\boldsymbol{\beta}}^{\text{ls}}$?

$$\begin{aligned}\text{Var}(\hat{\boldsymbol{\beta}}^{\text{ls}}) &= \sigma^2 (\mathbf{X}'\mathbf{X})^{-1} \\ &= \sigma^2 \mathbf{p} \text{ (if } \mathbf{X} \text{ has orthonormal columns)}\end{aligned}$$

Shortcomings in high-dimension

- The shortcomings don't even have to do with the linearity assumption!
- It might happen that the columns of \mathbf{X} are not linearly independent, so that \mathbf{X} is not of full rank. Then $\mathbf{X}'\mathbf{X}$ is singular and the least squares coefficients are not uniquely defined.
- **Predictive ability:** tradeoff between bias and variance.
- **Interpretative ability:** When the number of variables p is large, we may sometimes seek, for the sake of interpretation, a smaller set of *important variables*

Alternatives

- Dimension Reduction: We project the p predictors into a M -dimensional subspace, where $M < p$. Then these M projections are used as predictors to fit a linear regression model by least squares.
- Subset Selection: We identify a subset of the p predictors that we believe to be related to the response. We then fit a model using least squares on the reduced set of variables.
- **Shrinkage**: We fit a model involving all p predictors, but the estimated coefficients are shrunk towards zero relative to the least squares estimates. This shrinkage (also known as regularization) has the effect of reducing variance and can also perform variable selection.

Best subset selection

$$\begin{aligned} & \underset{\beta}{\text{minimize}} \left\{ \sum_{i=1}^n \left(y_i - \beta_0 - \sum_{j=1}^p \beta_j x_{ij} \right)^2 \right\} \\ & \text{subject to } \sum_{j=1}^p \mathbf{I}(\beta_j \neq 0) \leq s, \quad s \geq 0. \end{aligned}$$

where $s \geq 0$ is a tuning parameter.

- Need to consider $\binom{p}{s}$ models containing s predictors \rightarrow Computationally infeasible when p and s are large + larger the search space, the higher the chance of overfitting
- Stepwise procedures: forward, backward, etc.

Best subset selection

$$\begin{aligned} & \underset{\beta}{\text{minimize}} \left\{ \sum_{i=1}^n \left(y_i - \beta_0 - \sum_{j=1}^p \beta_j x_{ij} \right)^2 \right\} \\ & \text{subject to } \sum_{j=1}^p \mathbf{I}(\beta_j \neq 0) \leq s, \quad s \geq 0. \end{aligned}$$

where $s \geq 0$ is a tuning parameter.

- Need to consider $\binom{p}{s}$ models containing s predictors \rightarrow Computationally infeasible when p and s are large + larger the search space, the higher the chance of overfitting
- Stepwise procedures: forward, backward, etc.

Ridge regression

$$\begin{aligned} & \underset{\beta}{\text{minimize}} \left\{ \sum_{i=1}^n \left(y_i - \beta_0 - \sum_{j=1}^p \beta_j x_{ij} \right)^2 \right\} \\ & \text{subject to } \sum_{j=1}^p \beta_j^2 \leq s \end{aligned}$$

where $s \geq 0$ is a tuning parameter.

■ $s = 0?$

→ $\hat{\beta}^R = (0, \dots, 0)$

■ $s = \infty?$

→ $\hat{\beta}^R = \hat{\beta}^{\text{ls}}$ (least squares)

■ $s \in (0, \infty)$

→ tradeoff

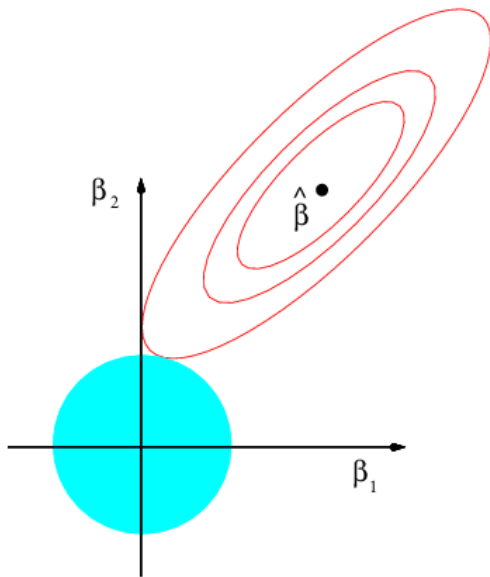
Ridge regression

$$\begin{aligned} & \underset{\beta}{\text{minimize}} \left\{ \sum_{i=1}^n \left(y_i - \beta_0 - \sum_{j=1}^p \beta_j x_{ij} \right)^2 \right\} \\ & \text{subject to } \sum_{j=1}^p \beta_j^2 \leq s \end{aligned}$$

where $s \geq 0$ is a tuning parameter.

- $s = 0?$ $\rightarrow \hat{\beta}^R = (0, \dots, 0)$
- $s = \infty?$ $\rightarrow \hat{\beta}^R = \hat{\beta}^{\text{ls}}$ (least squares)
- $s \in (0, \infty)$ \rightarrow tradeoff

Ridge regression: geometry



Ridge regression: another formulation

$$\underset{\beta}{\text{minimize}} \sum_{i=1}^n \left(y_i - \beta_0 - \sum_{j=1}^p \beta_j x_{ij} \right)^2 + \lambda \sum_{j=1}^p \beta_j^2$$

where $\lambda \geq 0$ is a **tuning parameter**.

- $\lambda = 0?$ $\rightarrow \hat{\beta}^R = \hat{\beta}^{ls}$ (least squares)
- $\lambda = \infty?$ $\rightarrow \hat{\beta}^R = (0, \dots, 0)$
- $\lambda \in (0, \infty)$ \rightarrow tradeoff
- We can solve it by *data augmentation*
- $\hat{\beta}^R = (\mathbf{X}'\mathbf{X} + \lambda \mathbf{I}_p)^{-1} \mathbf{X}'\mathbf{y}$

Ridge regression: another formulation

$$\underset{\beta}{\text{minimize}} \sum_{i=1}^n \left(y_i - \beta_0 - \sum_{j=1}^p \beta_j x_{ij} \right)^2 + \lambda \sum_{j=1}^p \beta_j^2$$

where $\lambda \geq 0$ is a **tuning parameter**.

- $\lambda = 0?$ $\rightarrow \hat{\beta}^R = \hat{\beta}^{\text{ls}}$ (least squares)
 - $\lambda = \infty?$ $\rightarrow \hat{\beta}^R = (0, \dots, 0)$
 - $\lambda \in (0, \infty)$ \rightarrow tradeoff
-
- We can solve it by *data augmentation*
 - $\hat{\beta}^R = (\mathbf{X}'\mathbf{X} + \lambda \mathbf{I}_p)^{-1} \mathbf{X}'\mathbf{y}$

Ridge regression: another formulation

$$\underset{\beta}{\text{minimize}} \sum_{i=1}^n \left(y_i - \beta_0 - \sum_{j=1}^p \beta_j x_{ij} \right)^2 + \lambda \sum_{j=1}^p \beta_j^2$$

where $\lambda \geq 0$ is a **tuning parameter**.

- $\lambda = 0?$ $\rightarrow \hat{\beta}^R = \hat{\beta}^{\text{ls}}$ (least squares)
- $\lambda = \infty?$ $\rightarrow \hat{\beta}^R = (0, \dots, 0)$
- $\lambda \in (0, \infty)$ \rightarrow tradeoff

- We can solve it by *data augmentation*
- $\hat{\beta}^R = (\mathbf{X}'\mathbf{X} + \lambda \mathbf{I}_p)^{-1} \mathbf{X}'\mathbf{y}$

A Simple special case

$$\underset{\beta}{\text{minimize}} \sum_{j=1}^n (y_j - \sum_{j=1}^p \beta_j \mathbf{x}_{ij})^2 \rightarrow \hat{\beta}^{\text{ls}} = (\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{y}$$

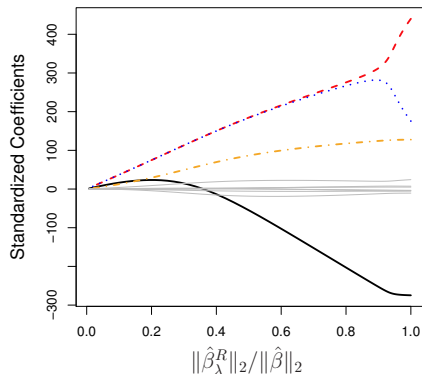
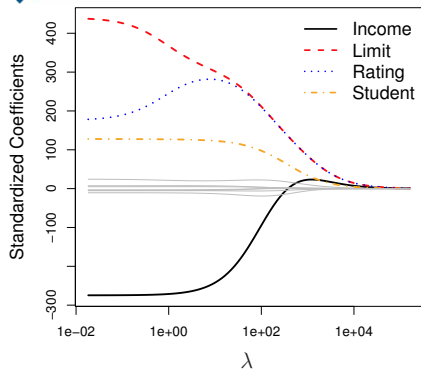
Suppose $n = p$ and $\mathbf{X} = \mathbf{I}_n = \mathbf{I}_p$, then

$$\underset{\beta}{\text{minimize}} \sum_{j=1}^p (y_j - \beta_j)^2 \rightarrow \hat{\beta}_j^{\text{ls}} = y_j$$

$$\underset{\beta}{\text{minimize}} \sum_{j=1}^p (y_j - \beta_j)^2 + \lambda \sum_{j=1}^p \beta_j^2 \rightarrow \hat{\beta}_j^R = \frac{y_j}{(1 + \lambda)} = \frac{\hat{\beta}_j^{\text{ls}}}{(1 + \lambda)}$$

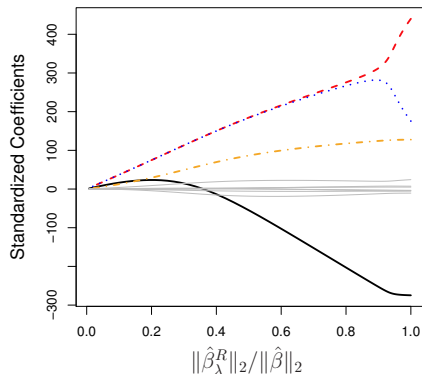
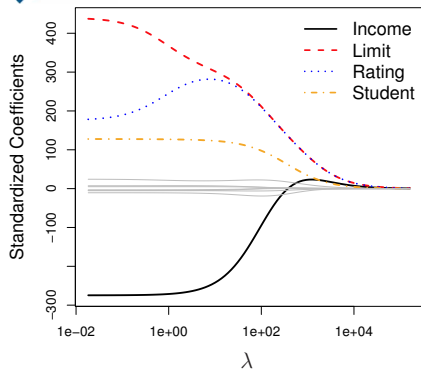
- This illustrates the essential feature of ridge regression: shrinkage. Introducing bias but reducing the variance.
- Each least squares coefficient estimate is shrunk by the **same proportion**.

Ridge regression: example



While the ridge coefficient estimates tend to **decrease in aggregate** as λ increases, individual coefficients, such as rating and income, may **occasionally increase** as λ increases.

Ridge regression: example

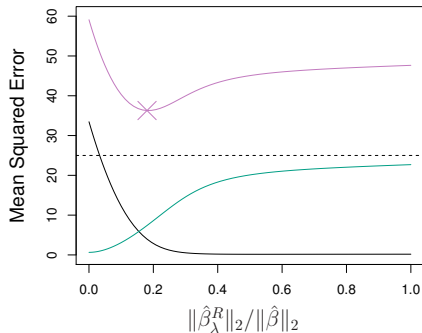
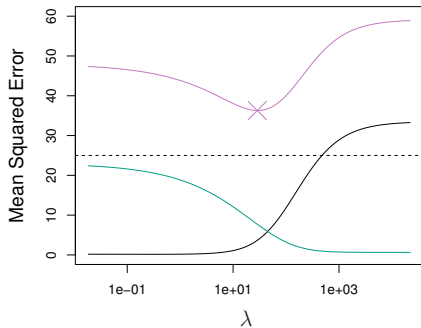


While the ridge coefficient estimates tend to **decrease in aggregate** as λ increases, individual coefficients, such as rating and income, may **occasionally increase** as λ increases.

A note on scaling

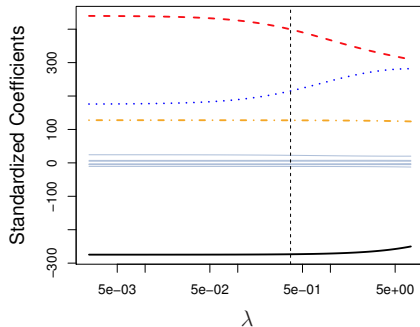
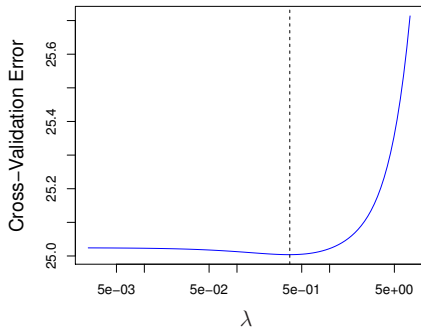
- Standard least squares coefficient estimates are **scale equivariant**
 - multiplying X_j by a constant c simply leads to a scaling of the least squares coefficient estimates by a factor of $1/c$
 - regardless of how the j th predictor is scaled, $X_j\hat{\beta}_j$ will remain the same.
- The ridge regression coefficient estimates **can change substantially** when multiplying a given predictor by a constant
 - This is due to the sum of squared coefficients term in the ridge regression formulation
 - If we use thousands of dollars instead of dollars, it will **not** simply cause the ridge estimate to change by a factor of 1,000

Ridge Regression vs Least Squares



Squared bias (black), variance (green), and test mean squared error (purple)

Selecting the Tuning Parameter



Ridge regression bias and variance

If $\mathbf{R} = \mathbf{X}'\mathbf{X}$:

$$\begin{aligned}\beta_{\lambda}^R &= (\mathbf{X}'\mathbf{X} + \lambda\mathbf{I}_p)^{-1}\mathbf{X}'\mathbf{y} = (\mathbf{R} + \lambda\mathbf{I}_p)^{-1}\mathbf{R}(\mathbf{R}^{-1}\mathbf{X}'\mathbf{y}) \\ &= (\mathbf{R} + \lambda\mathbf{I}_p)^{-1}\mathbf{R}\hat{\beta}^{ls} = [\mathbf{R}(\mathbf{I}_p + \lambda\mathbf{R}^{-1})]^{-1}\mathbf{R}\hat{\beta}^{ls} \\ &= (\mathbf{I}_p + \lambda\mathbf{R}^{-1})\hat{\beta}^{ls}\end{aligned}$$

If $\mathbf{W}_{\lambda} = (\mathbf{I}_p + \lambda\mathbf{R}^{-1})$:

$$E[\beta_{\lambda}^R] = E[\mathbf{W}_{\lambda}\hat{\beta}^{ls}] = \mathbf{W}_{\lambda}\beta \stackrel{\lambda \neq 0}{\neq} \beta$$

$$\text{Var}(\beta_{\lambda}^R) = \text{Var}(\mathbf{W}_{\lambda}\hat{\beta}^{ls}) = \mathbf{W}_{\lambda}\text{Var}(\hat{\beta}^{ls})\mathbf{W}_{\lambda}' \preceq \text{Var}(\hat{\beta}^{ls})$$

Singular Value Decomposition

Singular Value Decomposition (SVD)

$$\mathbf{X} = \mathbf{U}\mathbf{D}\mathbf{V}'$$

- \mathbf{X} is $n \times p$ matrix
- \mathbf{U} is $n \times r$ matrix with orthonormal columns ($\mathbf{U}'\mathbf{U} = \mathbf{I}$)
- \mathbf{D} is $r \times r$ diagonal matrix with diagonal entries $d_1, \geq d_2 \geq \dots \geq d_p \geq 0$ called the singular values of \mathbf{X} .
- \mathbf{V} is $p \times r$ matrix with orthonormal columns ($\mathbf{V}'\mathbf{V} = \mathbf{I}$).

Note: $\mathbf{XV} = \mathbf{UD}$

Least squares regression and SVD

If $\mathbf{X} = \mathbf{UDV}'$, then

$$\begin{aligned}\hat{\beta}^{\text{ls}} &= (\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{y} = \mathbf{VD}^{-1}\mathbf{U}'\mathbf{y} \\ &= \mathbf{V}\mathbf{D}^{-2}\mathbf{DU}'\mathbf{y}\end{aligned}$$

$$\begin{aligned}\hat{\beta}_{\lambda}^R &= (\mathbf{X}'\mathbf{X} + \lambda\mathbf{I}_p)^{-1}\mathbf{X}'\mathbf{y} \\ &= (\mathbf{VD}^2\mathbf{V}' + \lambda\mathbf{VV}')^{-1}\mathbf{VDU}'\mathbf{y} \\ &= \left(\mathbf{V}(\mathbf{D}^2 + \lambda)\mathbf{V}'\right)^{-1}\mathbf{VDU}'\mathbf{y} \\ &= \mathbf{V}(\mathbf{D}^2 + \lambda)^{-1}\mathbf{V}'\mathbf{VDU}'\mathbf{y} \\ &= \mathbf{V}(\mathbf{D}^2 + \lambda)^{-1}\mathbf{DU}'\mathbf{y}\end{aligned}$$

Ridge regression and SVD

$$\hat{\mathbf{y}}^{\text{ls}} = \mathbf{X}\hat{\boldsymbol{\beta}}^{\text{ls}} = \mathbf{U}\mathbf{U}'\mathbf{y}$$

$\mathbf{U}'\mathbf{y}$ are the coordinates of \mathbf{y} with respect to the orthonormal basis \mathbf{U}

$$\begin{aligned}\hat{\mathbf{y}}_{\lambda}^{\text{R}} &= \mathbf{X}\hat{\boldsymbol{\beta}}_{\lambda}^{\text{R}} = \mathbf{X}\mathbf{V}(\mathbf{D}^2 + \lambda\mathbf{I})^{-1}\mathbf{D}\mathbf{U}'\mathbf{y} = \mathbf{U}\mathbf{D}(\mathbf{D}^2 + \lambda\mathbf{I})^{-1}\mathbf{D}\mathbf{U}'\mathbf{y} \\ &= \mathbf{U} \text{diag}\left(\frac{d_j^2}{d_j^2 + \lambda}\right)\mathbf{U}'\mathbf{y}\end{aligned}$$

- Since $\lambda \geq 0$, we have $d_j^2/(d_j^2 + \lambda) \leq 1$. A shrinkage is applied to the coordinates by the factors $d_j^2/(d_j^2 + \lambda)$.
- A greater amount of shrinkage is applied to the coordinates of basis vectors with smaller d_j^2 .
- The derived variable $\mathbf{z}_j = \mathbf{X}\mathbf{v}_j = \mathbf{u}_j d_j$ is the j th PC of \mathbf{X} . We project \mathbf{y} onto these components with large d_j , and shrinks the coefficients of low-variance components.

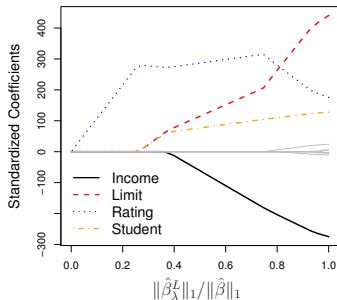
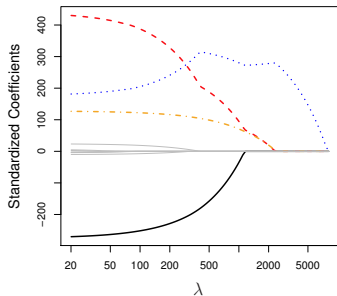
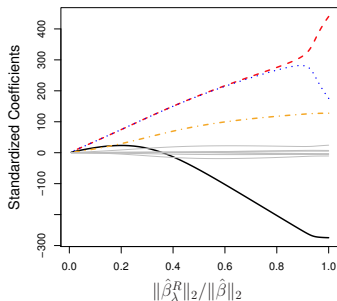
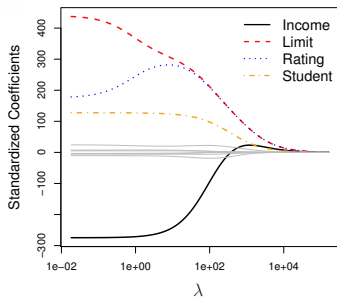
Summary

$$\hat{\mathbf{y}}^{\text{ls}} = \mathbf{U}\mathbf{U}'\mathbf{y}$$

$$\hat{\mathbf{y}}_{\lambda}^{\text{R}} = \mathbf{U} \text{diag} \left(\frac{d_j^2}{d_j^2 + \lambda} \right) \mathbf{U}'\mathbf{y}$$

$$\hat{\mathbf{y}}_k^{\text{PCR}} = \mathbf{U} \text{diag} \left(\underbrace{1, \dots, 1}_k, \underbrace{0, \dots, 0}_{p-k} \right) \mathbf{U}'\mathbf{y}$$

Another shrinkage method



LASSO regression

LASSO: Least Absolute Shrinkage and Selection Operator

$$\begin{aligned} & \underset{\beta}{\text{minimize}} \left\{ \sum_{i=1}^n \left(y_i - \beta_0 - \sum_{j=1}^p \beta_j x_{ij} \right)^2 \right\} \\ & \text{subject to } \sum_{j=1}^p |\beta_j| \leq s \end{aligned}$$

■ $s = 0?$

→ $\hat{\beta}^R = (0, \dots, 0)$

■ $s = \infty?$

→ $\hat{\beta}^R = \hat{\beta}^{\text{ls}}$ (least squares)

■ $s \in (0, \infty)$

→ tradeoff

LASSO regression

LASSO: Least Absolute Shrinkage and Selection Operator

$$\begin{aligned} & \underset{\beta}{\text{minimize}} \left\{ \sum_{i=1}^n \left(y_i - \beta_0 - \sum_{j=1}^p \beta_j x_{ij} \right)^2 \right\} \\ & \text{subject to } \sum_{j=1}^p |\beta_j| \leq s \end{aligned}$$

- $s = 0?$ $\rightarrow \hat{\beta}^R = (0, \dots, 0)$
- $s = \infty?$ $\rightarrow \hat{\beta}^R = \hat{\beta}^{\text{ls}}$ (least squares)
- $s \in (0, \infty)$ \rightarrow tradeoff

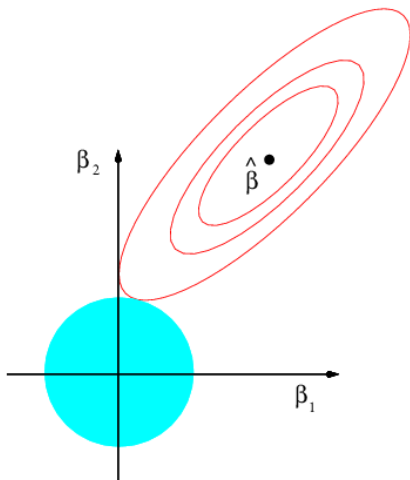
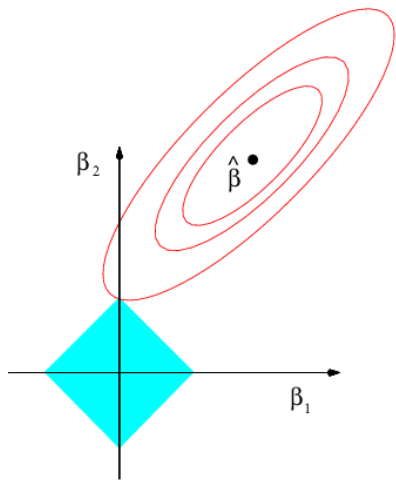
LASSO: another formulation

$$\underset{\beta}{\text{minimize}} \sum_{i=1}^n \left(y_i - \beta_0 - \sum_{j=1}^p \beta_j x_{ij} \right)^2 + \lambda \sum_{j=1}^p |\beta_j|$$

where $\lambda \geq 0$ is a **tuning parameter**.

- $\lambda = 0?$ $\rightarrow \hat{\beta}^R = \hat{\beta}^{\text{ls}}$ (least squares)
- $\lambda = \infty?$ $\rightarrow \hat{\beta}^R = (0, \dots, 0)$
- $\lambda \in (0, \infty)$ \rightarrow tradeoff

LASSO vs Ridge: geometry



Sparsity

We shall say that a signal $\mathbf{x} \in R^n$ is **sparse**, when most of the entries of \mathbf{x} **vanish**. Formally, we shall say that a signal is s -sparse if it has **at most s nonzero entries**. One can think of an s -sparse signal as having only s degrees of freedom.

- L_q regularization with $q > 1$ does not provide sparse estimate
→ e.g. ridge regression
- For $q < 1$, the solutions are sparse but the problem is **not convex** and this makes the optimisation very challenging computationally.
- The value $q = 1$ is the smallest value that yields a **convex problem**.

The bet on sparsity principle

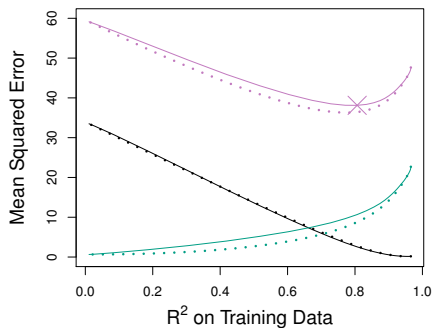
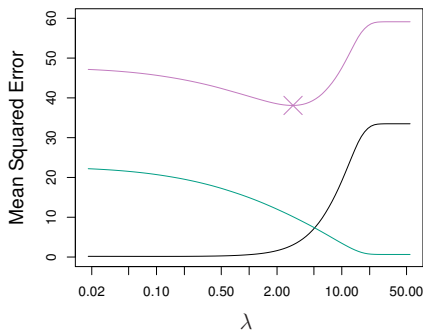
If $p \gg N$ and the true model **is sparse**, so that only $k < N$ parameters are actually nonzero in the true underlying model, then it turns out that we can estimate the parameters effectively, using the lasso and related methods.

if $p \gg N$, and the true model **is not sparse**, then the number of samples N is too small to allow for accurate estimation of the parameters (The amount of information per parameter is N/p)

Use a procedure that does well in sparse problems, since no procedure does well in dense problems

Lasso vs ridge regression

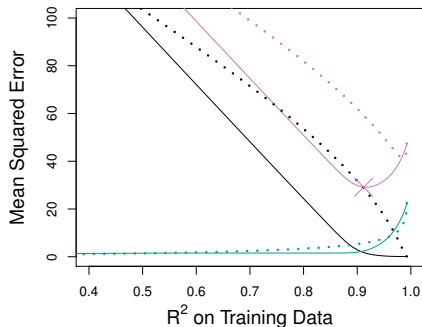
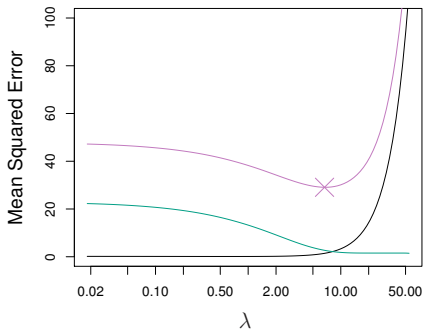
A simulated data set containing $p = 45$ predictors and $n = 50$ observations where **all 45 predictors are related to the response**.



Left: Lasso. **Right:** Lasso (solid) and ridge (dashed).

Lasso vs ridge regression

Now the response is a function of **only 2 out of 45 predictors**.



Left: Lasso. **Right:** Lasso (solid) and ridge (dashed).

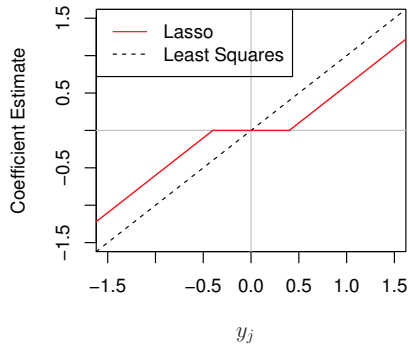
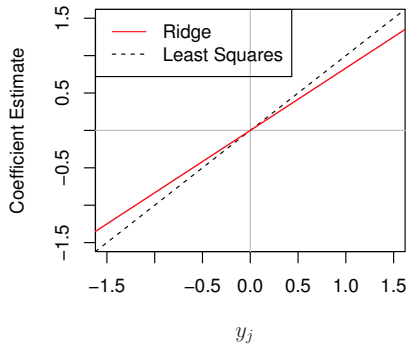
A Simple special case: lasso

$$\underset{\beta}{\text{minimize}} \sum_{j=1}^p (y_j - \beta_j)^2 + \lambda \sum_{j=1}^p |\beta_j|$$

$$\hat{\beta}_j^L = \begin{cases} y_j - \frac{\lambda}{2} & \text{if } y_j > \frac{\lambda}{2}; \\ y_j + \frac{\lambda}{2} & \text{if } y_j < -\frac{\lambda}{2}; \\ 0 & \text{if } |y_j| \leq \frac{\lambda}{2}. \end{cases}$$

The lasso shrinks each least squares coefficient towards zero by a **constant amount**, $\lambda/2$. The least squares coefficients that are less than $\lambda/2$ in absolute value are **shrunk entirely to zero**.

A Simple special case

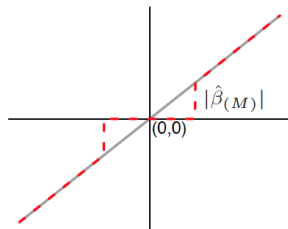


A Simple special case

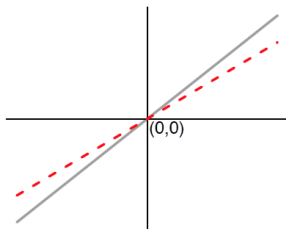
$\hat{\beta}_j$ (OLS estimate) and $\hat{\beta}_{(M)}$ (M th largest coefficient)

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

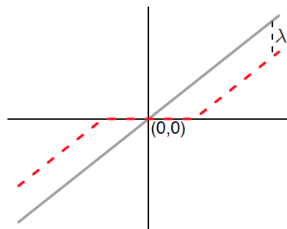
Best Subset



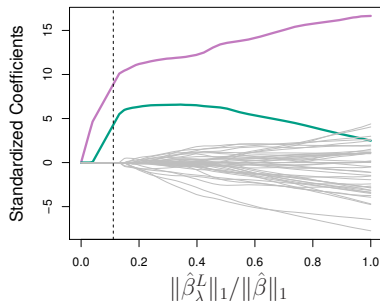
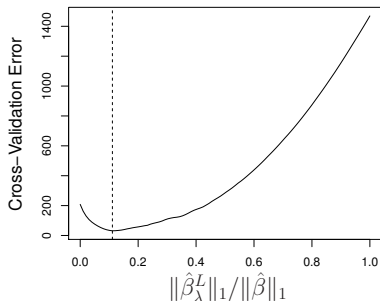
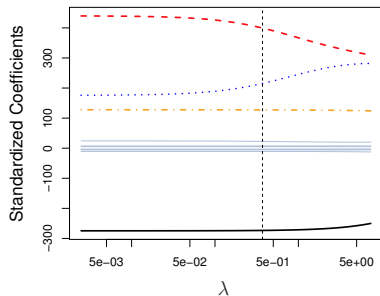
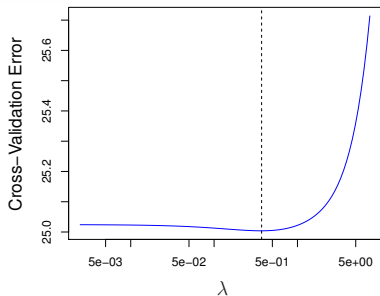
Ridge



Lasso



Selecting the Tuning Parameter



q-norm

Let $q \geq 1$ be a real number. The q -norm of $\mathbf{x} = (x_1, \dots, x_p)$ is given by

$$\|\mathbf{x}\|_q = \left(\sum_{j=1}^p |x_j|^q \right)^{1/q}$$

- $q = 1$: L_1 norm
- $q = 2$: L_2 norm, Euclidean norm
- $q = \infty$: L_∞ norm, uniform norm:
 $\|\mathbf{x}\|_\infty = \max\{|x_1|, \dots, |x_p|\}.$

q-norm

