# Business Analytics

**Week 8**
**Advanced regression**

11 September 2017

# Outline

# Regression

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

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

$$m^* = \operatorname*{argmin}_{m \in \mathcal{M}} \mathbb{E}[(Y - m(X))^2]$$

# Linear regression

$$m(X) = \beta_0 + \sum_{j=1}^{p} X_j \beta_j \rightarrow \beta^* = \min_{\boldsymbol{\beta} \in \mathbb{R}^p} \mathbb{E}[(Y - m(X))^2]$$

$$\min_{\boldsymbol{\beta} \in \mathbb{R}^p} \sum_{i=1}^{n} (y_i - \beta_0 - \sum_{j=1}^{p} \beta_j x_{ij})^2 := (\boldsymbol{y} - \boldsymbol{X}\boldsymbol{\beta})'(\boldsymbol{y} - \boldsymbol{X}\boldsymbol{\beta})$$

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

For a truly linear underlying model, the linear regression has expected test error

$$\sigma^2 + 0 + \frac{\sigma^2 p}{n}.$$

# Shortcomings in high-dimension

- The shortcomings don't even have to do with the linearity assumption!

- It might happen that the columns of $X$ are not linearly independent, so that $X$ is not of full rank. Then $X'X$ 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

- <u>Dimension Reduction</u>: 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.

- <u>Subset Selection</u>: 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 shrunken 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

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

$$\text{subject to } \sum_{j=1}^{p} \boldsymbol{I}(\beta_j \neq 0) \leq s, \quad s \geq 0.$$

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

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

$$\text{subject to} \sum_{j=1}^{p} I(\beta_j \neq 0) \leq s, \quad s \geq 0.$$

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
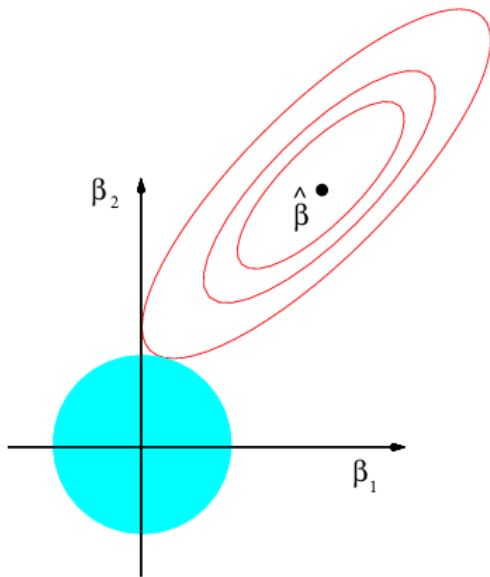
- Stepwise procedures: forward, backward, etc.

# Ridge regression

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

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

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

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

# Ridge regression

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

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

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

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

# Ridge regression: geometry

# Ridge regression: another formulation

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

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

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

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

# Ridge regression: another formulation

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

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

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

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

# Ridge regression: another formulation

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

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

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

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

# A Simple special case

$$\underset{\beta}{\text{minimize}} \sum_{j=1}^{n}(y_i - \sum_{j=1}^{p} \beta_j x_{ij})^2 \rightarrow \hat{\beta}^{\text{ls}} = (\boldsymbol{X}'\boldsymbol{X})^{-1}\boldsymbol{X}'\boldsymbol{y}$$

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

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

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

- This illustrates the essential feature of ridge regression: shrinkage. Introducing bias but reducing the variance.

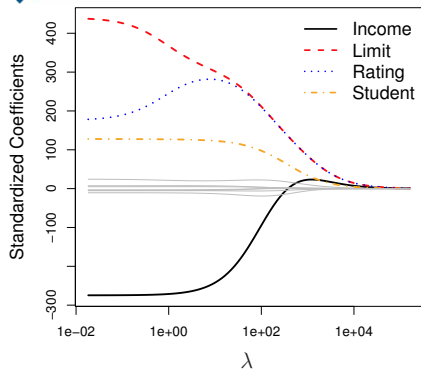- Each least squares coefficient estimate is shrunken by the **same proportion**.
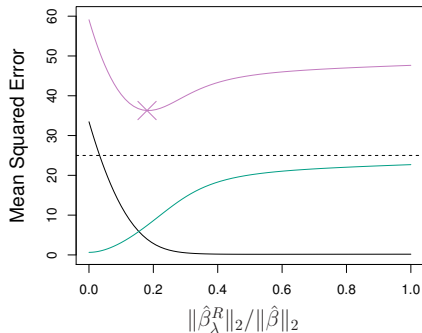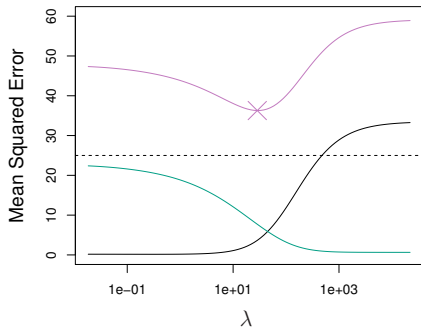
# Ridge regression: example



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

# Ridge regression: example
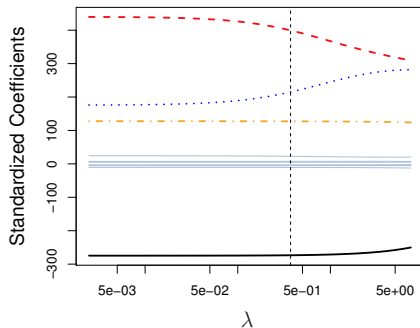


While the ridge coefficient estimates tend to **decrease in aggregate** as $\lambda$ increases, individual coefficients, such as rating and income, may **occasionally increase** as $\lambda$ increases.
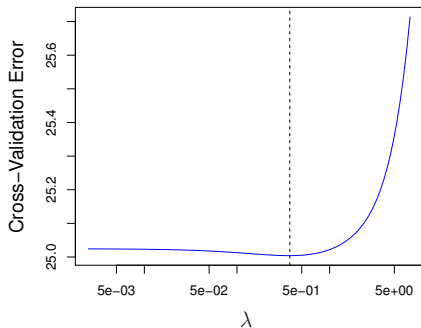
# A note on scaling

- 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)

# Ridge regression bias and

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

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

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

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

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

# Singular Value Decomposition

## Singular Value Decomposition (SVD)

$$X = UDV'$$

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

Note: $XV = UD$

# Least squares regression and SVD

If $\boldsymbol{X} = \boldsymbol{U}\boldsymbol{D}\boldsymbol{V}'$, then

$$\hat{\beta}^{\mathsf{ls}} = \boldsymbol{X}(\boldsymbol{X}'\boldsymbol{X})^{-1}\boldsymbol{X}' = \boldsymbol{V}\boldsymbol{D}^{-1}\boldsymbol{U}'\boldsymbol{y}$$
$$= \boldsymbol{V}\boldsymbol{D}^{-2}\boldsymbol{D}\boldsymbol{U}'\boldsymbol{y}$$

$$\hat{\beta}^{R}_{\lambda} = (\boldsymbol{X}'\boldsymbol{X} + \lambda\boldsymbol{I}_p)^{-1}\boldsymbol{X}'\boldsymbol{y}$$
$$= (\boldsymbol{V}\boldsymbol{D}^2\boldsymbol{V}' + \lambda\boldsymbol{V}\boldsymbol{V}')^{-1}\boldsymbol{V}\boldsymbol{D}\boldsymbol{U}'\boldsymbol{y}$$
$$= \left(\boldsymbol{V}(\boldsymbol{D}^2 + \lambda)\boldsymbol{V}'\right)^{-1}\boldsymbol{V}\boldsymbol{D}\boldsymbol{U}'\boldsymbol{y}$$
$$= \boldsymbol{V}(\boldsymbol{D}^2 + \lambda)^{-1}\boldsymbol{V}'\boldsymbol{V}\boldsymbol{D}\boldsymbol{U}'\boldsymbol{y}$$
$$= \boldsymbol{V}(\boldsymbol{D}^2 + \lambda)^{-1}\boldsymbol{D}\boldsymbol{U}'\boldsymbol{y}$$

# Ridge regression and SVD

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

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

$\hat{\mathbf{y}}^{\text{R}}_{\lambda} = \mathbf{X}\hat{\beta}^{\text{R}}_{\lambda} = \mathbf{X}\mathbf{V}(\mathbf{D}^2 + \lambda)^{-1}\mathbf{D}\mathbf{U}'\mathbf{y} = \mathbf{U}\mathbf{D}(\mathbf{D}^2 + \lambda\mathbf{I})^{-1}\mathbf{D}\mathbf{U}'\mathbf{y}$

$\quad = \mathbf{U} \, \text{diag}\left(\dfrac{d_j^2}{d_j^2 + \lambda}\right)\mathbf{U}'\mathbf{y}$

- 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.

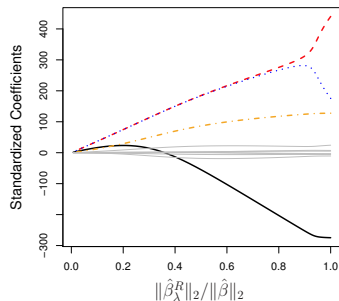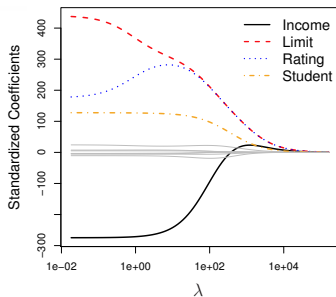# Summary

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

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

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

# Another shrinkage method

# LASSO regression

LASSO: Least Absolute Shrinkage and Selection Operator

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

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

- $s = 0$?      $\rightarrow \hat{\beta}^{\text{R}} = (0, \ldots, 0)$

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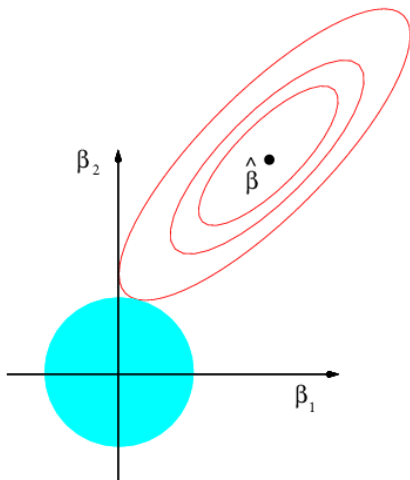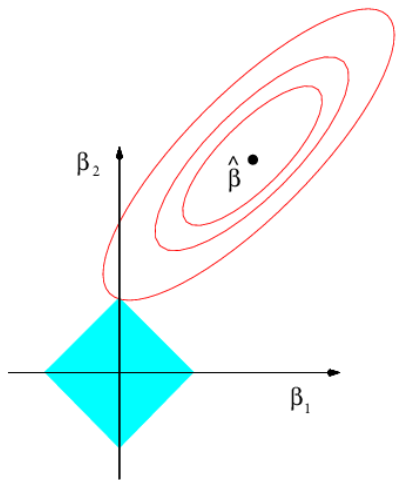
- $s \in (0, \infty)$      $\rightarrow$ tradeoff

# LASSO regression

LASSO: Least Absolute Shrinkage and Selection Operator

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

$$\text{subject to} \sum_{j=1}^{p} |\beta_j| \le s$$

- $s = 0$? $\rightarrow \hat{\boldsymbol{\beta}}^{\mathsf{R}} = (0, \ldots, 0)$

- $s = \infty$? $\rightarrow \hat{\boldsymbol{\beta}}^{\mathsf{R}} = \hat{\boldsymbol{\beta}}^{\mathsf{ls}}$ (least squares)

- $s \in (0, \infty)$ $\rightarrow$ tradeoff

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

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

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

# LASSO vs Ridge: geometry

# Sparsity

We shall say that a signal $x \in R^n$ is **sparse**, when most of the entries of $x$ **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
  $\rightarrow$ 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 $p >> N$ 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 $p >> N$, 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**.



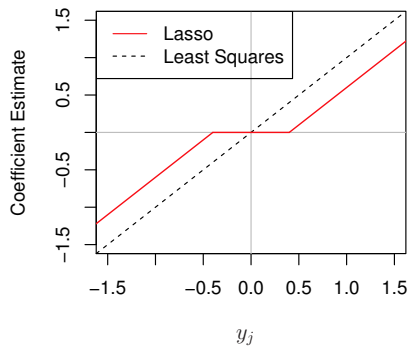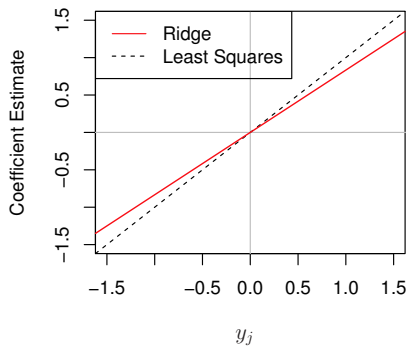**Left**: Lasso. **Right**: Lasso (solid) and ridge (dashed).

# A Simple special case: lasso

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

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

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 **shrunken entirely to zero**.
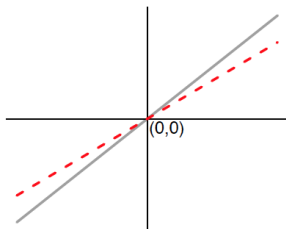
# A Simple special case

# A Simple special case

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

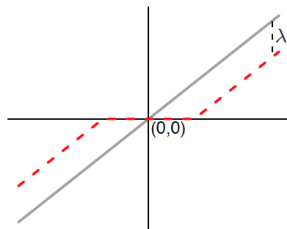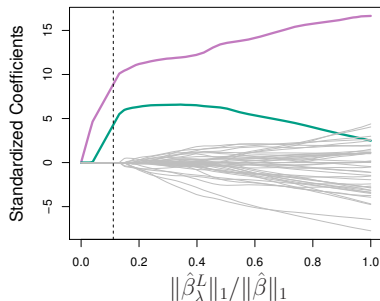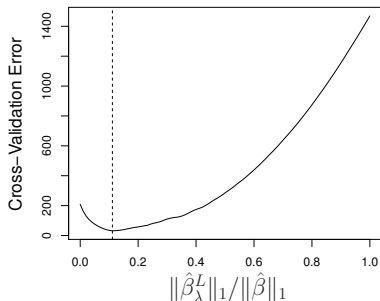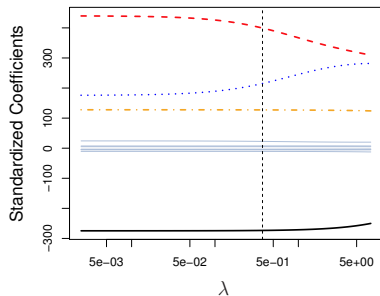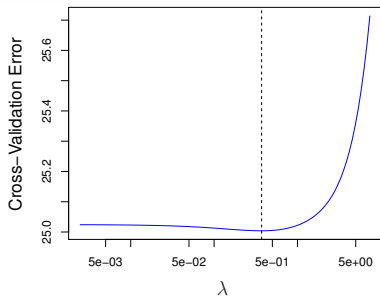| Estimator | Formula |
|---|---|
| Best subset (size $M$) | $\hat{\beta}_j \cdot I(|\hat{\beta}_j| \geq |\hat{\beta}_{(M)}|)$ |
| Ridge | $\hat{\beta}_j/(1 + \lambda)$ |
| Lasso | $\text{sign}(\hat{\beta}_j)(|\hat{\beta}_j| - \lambda)_+$ |



Best Subset

Ridge

Lasso

# Selecting the Tuning Parameter

# q-norm

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

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

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

# q-norm



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