# Business Analytics

**Week 10**
**Advanced regression**

10 October 2016

# Outline

| Week | Topic | Chapter | Lecturer |
|------|-------|---------|----------|
| 1 | Introduction to business analytics & R | 1 | Souhaib |
| 2 | Statistical learning | 2 | Souhaib |
| 3 | Regression for prediction | 3 | Souhaib |
| 4 | Resampling | 5 | Souhaib |
| 5 | Dimension reduction | 6,10 | Souhaib |
| 6 | Visualization | | Di |
| 7 | Visualization | | Di |
| 8 | Classification | 4,8 | Di |
| 9 | Classification | 4,9 | Di |
| | - | | |
| 10 | Classification | 8 | Souhaib |
| 11 | Advanced regression | 6 | Souhaib |
| 12 | Clustering | 10 | Souhaib |

# Regression

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

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

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

# Linear regression

$$m(X) = \beta_0 + \sum_{j=1}^{p} X_j \beta_j$$

$$\min_{\beta \in \mathbb{R}^p} \mathbb{E}[(Y - m(X))^2]$$

$$\text{RSS} = \sum_{i=1}^{n}(y_i - \beta_0 - \sum_{j=1}^{p}\beta_j x_{ij})^2 = (\boldsymbol{y} - \boldsymbol{X}\boldsymbol{\beta})'(\boldsymbol{y} - \boldsymbol{X}\boldsymbol{\beta})$$

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

# 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

- **Subset Selection**
- **Dimension Reduction**
- **Shrinkage**

# Best subset selection

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

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

Need to consider $\binom{p}{s}$ models containing $s$ predictors
$\rightarrow$ Computationally infeasible when $p$ is large

# Best subset selection

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

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

Need to consider $\binom{p}{s}$ models containing $s$ predictors
$\rightarrow$ Computationally infeasible when $p$ is large

# Ridge regression

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

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

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

# Ridge regression

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

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

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

# 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**.

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

# 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)
- $\lambda = \infty$?  $\rightarrow \hat{\boldsymbol{\beta}}^{\mathsf{R}} = (0, \dots, 0)$
- $\lambda \in (0, \infty)$  $\rightarrow$ tradeoff

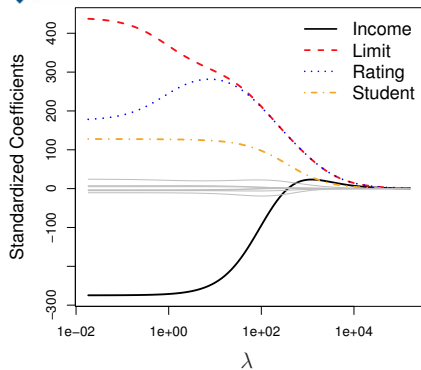# p-norm

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

$$\|\boldsymbol{x}\|_p = \left( \sum_{j=1}^{p} |x_j|^p \right)^{1/p}$$
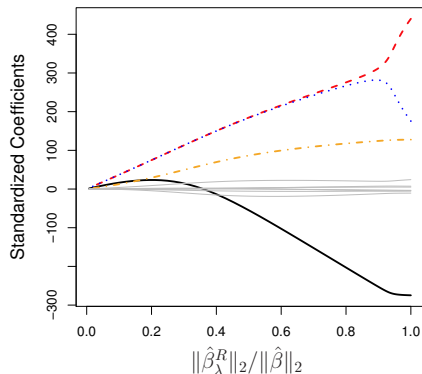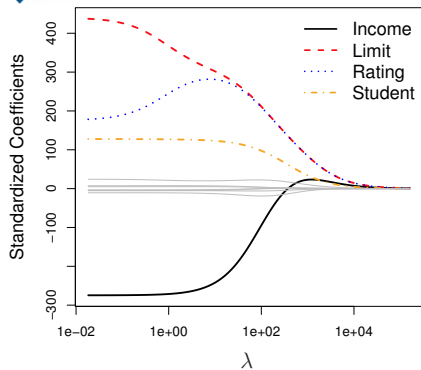
- $p = 1$: $L_1$ norm
- $p = 2$: $L_2$ norm, Euclidean norm
- $p = \infty$: $L_\infty$ norm, uniform norm:
  $\|x\|_\infty = \max\{|x_1|, \ldots, |x_p|\}$.

# 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.
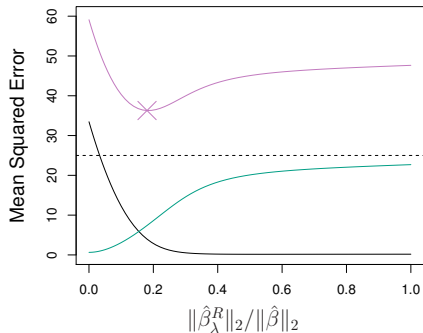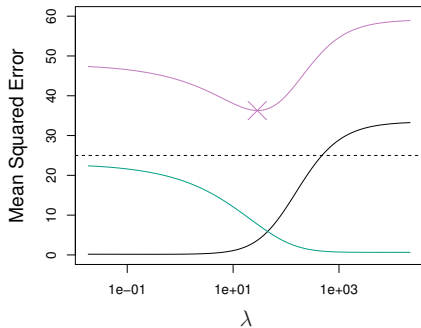
# Ridge regression: example



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

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

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

$$\begin{aligned}
E[\beta_\lambda^R] &= E[(\boldsymbol{I}_p + \lambda\boldsymbol{R}^{-1})\hat{\beta}^{ls}] \\
&= (\boldsymbol{I}_p + \lambda\boldsymbol{R}^{-1})\beta \\
&\overset{\lambda \neq 0}{\neq} \beta
\end{aligned}$$

# 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: $XV = UD$

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

Note that $\boldsymbol{U}'\boldsymbol{y}$ are the coordiantes of $\boldsymbol{y}$ with respect to the orthonormal basis $\boldsymbol{U}$.

# Ridge regression and SVD

$$\hat{\mathbf{y}}^{\mathsf{R}} = \mathbf{X}\hat{\beta}^{\mathsf{R}} = \mathbf{X}(\mathbf{X}'\mathbf{X} + \lambda\mathbf{I})^{-1}\mathbf{X}'\mathbf{y}$$

$$= \mathbf{U}\mathbf{D}(\mathbf{D}^2 + \lambda\mathbf{I})^{-1}\mathbf{D}\mathbf{U}'\mathbf{y}$$

$$= \sum_{j=1}^{p} \mathbf{u}_j \frac{d_j^2}{d_j^2 + \lambda} \mathbf{u}_j'\mathbf{y}$$

where the $\mathbf{u_j}$ are the columns of $\mathbf{U}$. Note that since $\lambda \geq 0$, we have $d_j^2/(d_j^2 + \lambda) \leq 1$.

Ridge regression shrinks the coordinates by the factors $d_j^2/(d_j^2 + \lambda)$. This means that a greater amount of shrinkage is applied to the coordinates of basis vectors with smaller $d_j^2$.