# QTM 347 Machine Learning

## Lecture 11: Lasso and elastic net

Ruoxuan Xiong

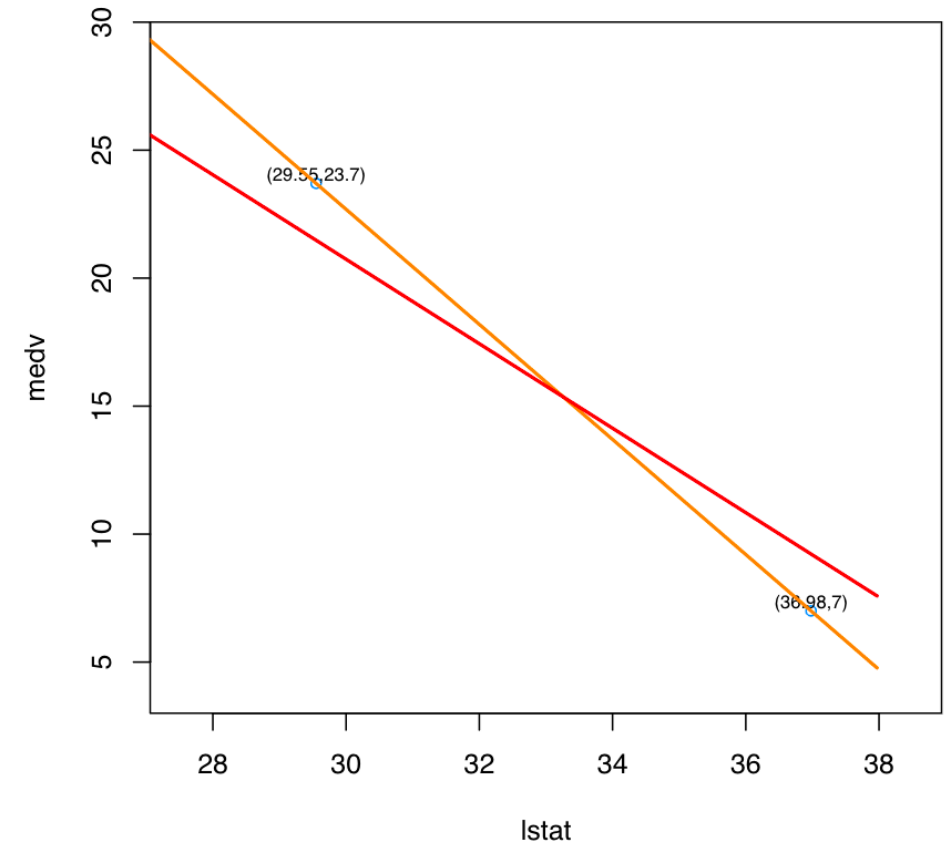Suggested reading: ISL Chapter 6

EMORY

# Lecture plan

- Lasso


- Elastic net

# Ridge regression

- Linear regression minimizes residual sum of squares
  - $RSS = \sum_{i=1}^{n}(medv_i - \beta_0 - lstat \cdot \beta_1)^2$


- Ridge regression minimizes
  - $\sum_{i=1}^{n}(medv_i - \beta_0 - lstat_i \cdot \beta_1)^2 + \lambda \cdot \beta_1^2$

  - $\lambda \geq 0$: tuning hyper-parameter

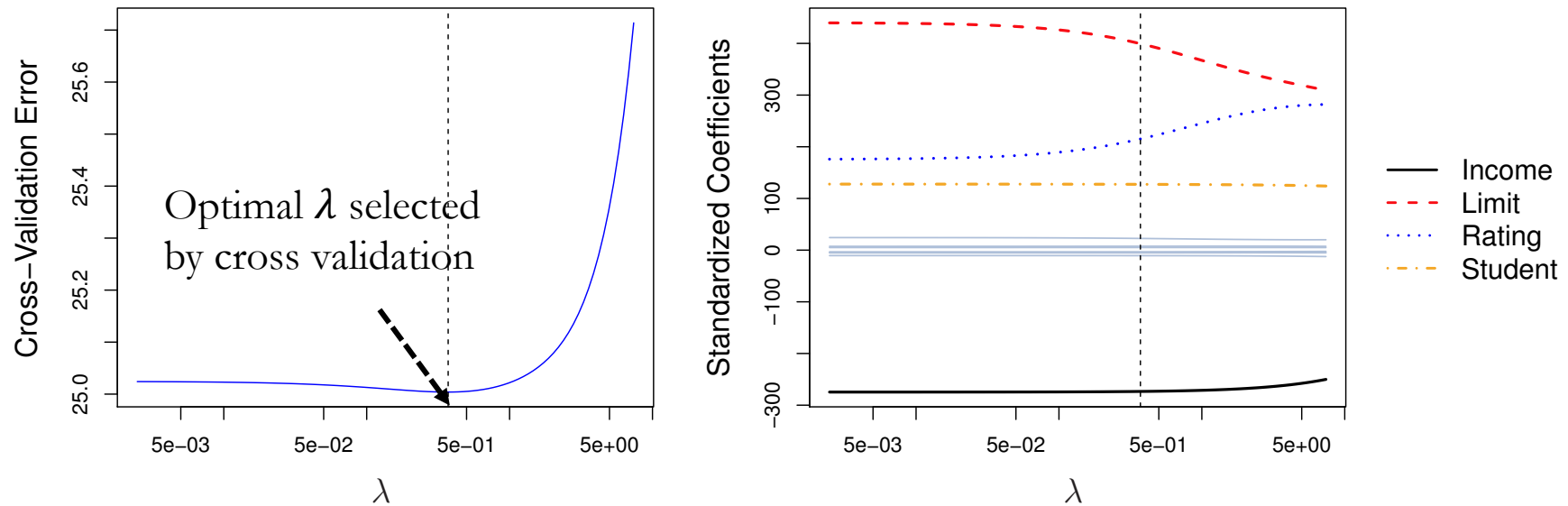# Ridge regression for more than one predictor

- Ridge regression minimizes

$$\sum_{i=1}^{n}\left(Y_i - \beta_0 - \sum_{j=1}^{p}\beta_j X_{i,j}\right)^2 + \lambda\sum_{j=1}^{p}\beta_j^2$$

- $X_{i,j}$: $j$-th predictor of $i$-th observation

- $\|\beta\|_2^2 = \sum_{j=1}^{p}\beta_j^2$: $\|\beta\|_2$ is called the $\ell_2$ norm of $\beta \in \mathbb{R}^p$

- $\beta_0$: mean of $Y_i$

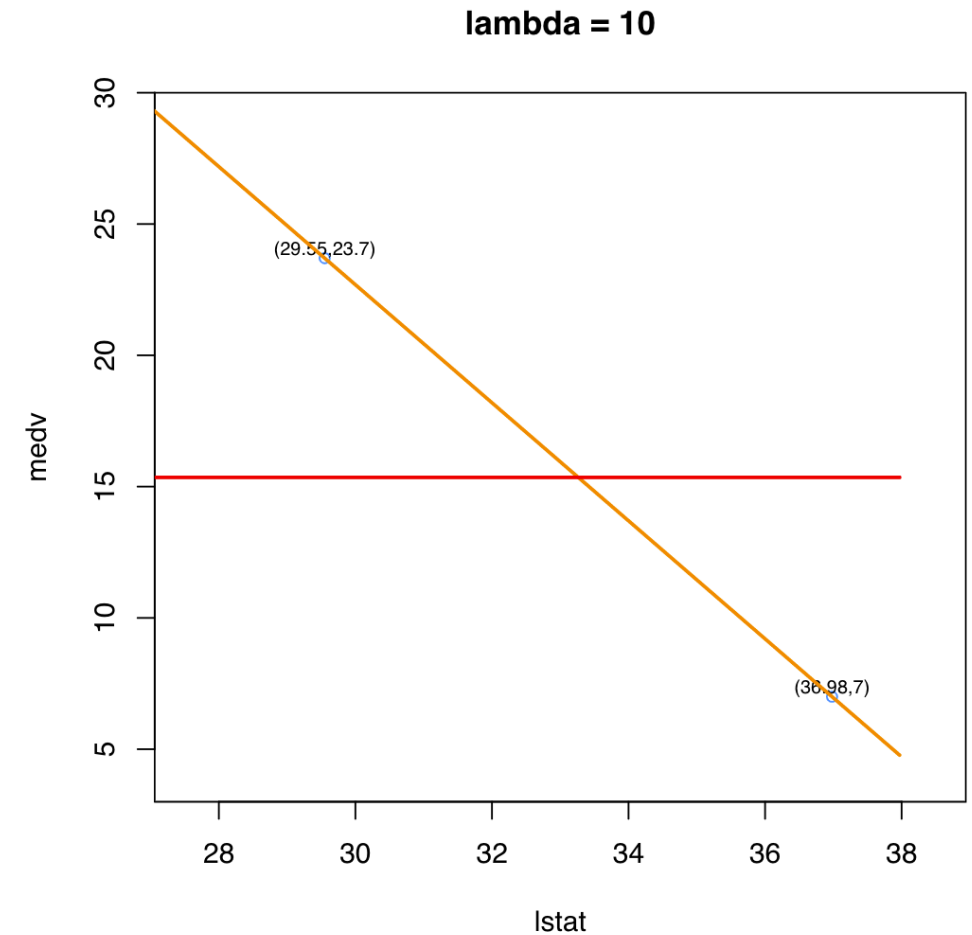- Shrinkage penalty $\lambda$ does not apply to $\beta_0$

EMORY

# Example: Credit card data set (ridge regression)

- Cross validation to choose the optimal $\lambda$

# Lasso

- Lasso: least absolute shrinkage and selection operator

- Lasso minimizes
  - $\sum_{i=1}^{n}(medv_i - \beta_0 - lstat_i \cdot \beta_1)^2 + \lambda \cdot |\beta_1|$

  - $\lambda \geq 0$:  tuning hyper-parameter



**lambda = 10**

(29.55,23.7)

(36.98,7)

lstat

medv

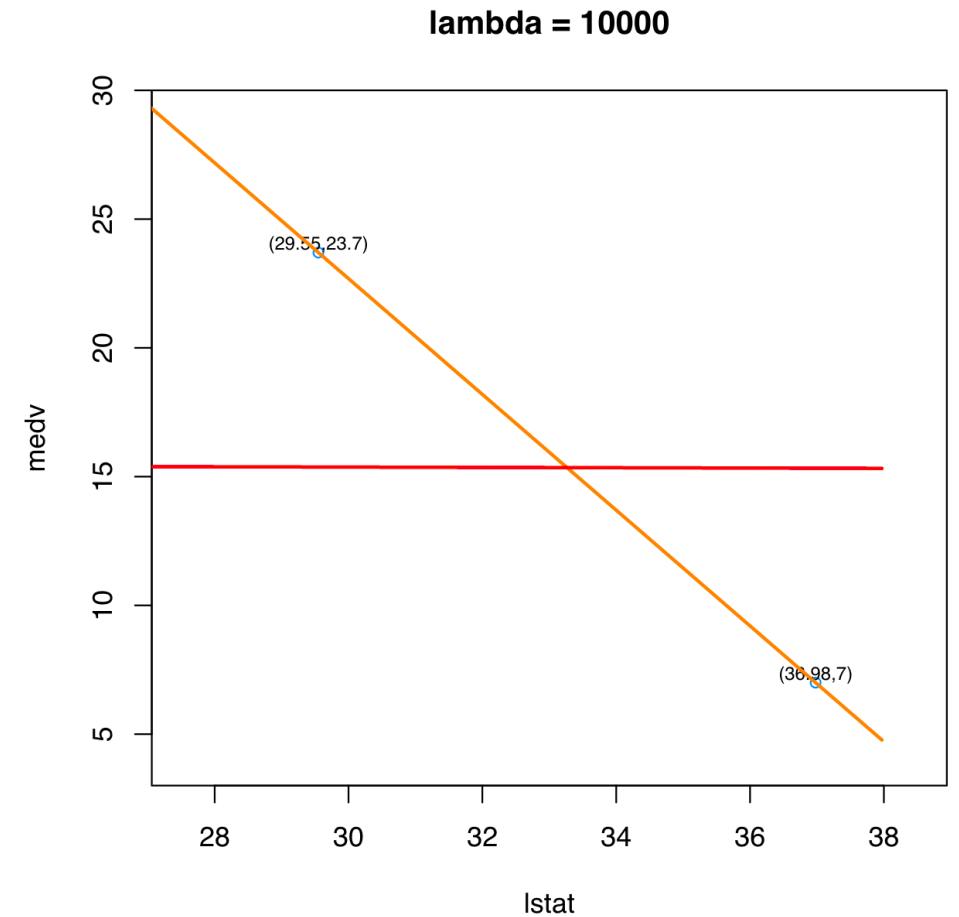# Motivation

- Ridge regression shrinks coefficients to approximately zero, but not exactly zero
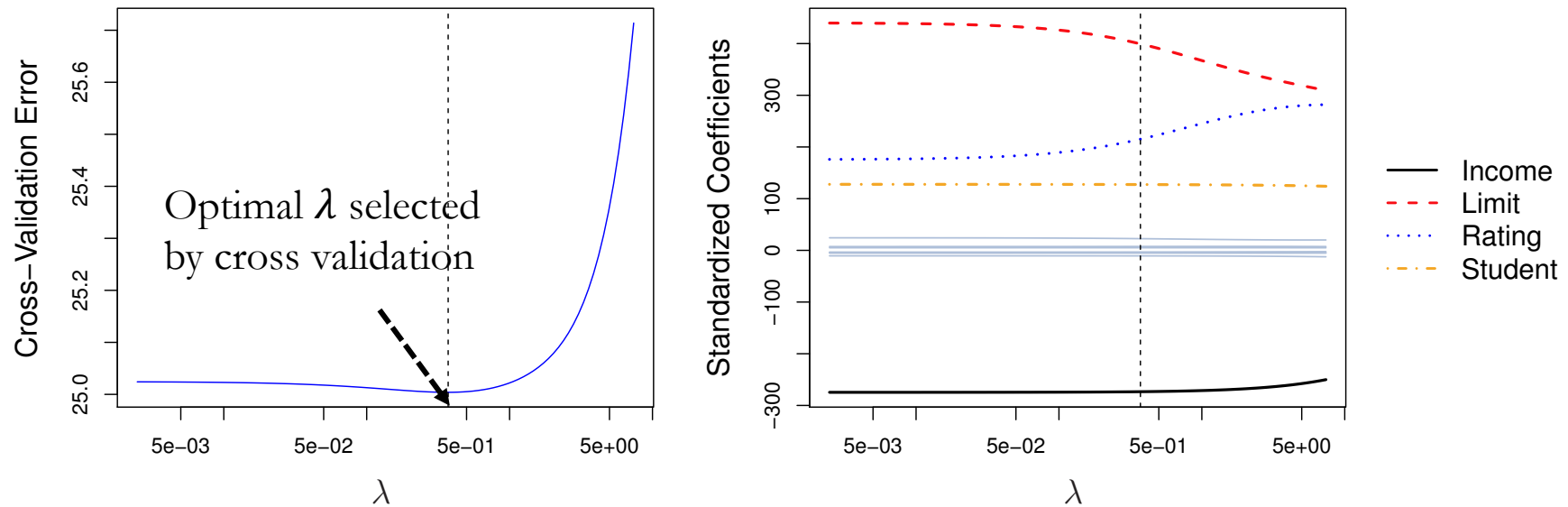
  - $\sum_{i=1}^{n}(medv_i - \beta_0 - lstat_i \cdot \beta_1)^2 + \lambda \cdot \beta_1^2$

  - When $\lambda = 10,000$, $\hat{\beta}_1^R = -0.0062$



**lambda = 10000**

# What if we want to exclude useless variables?

- In the credit data set, the standardized ridge coefficients for variables other than income, limit, rating, and student are nonzero

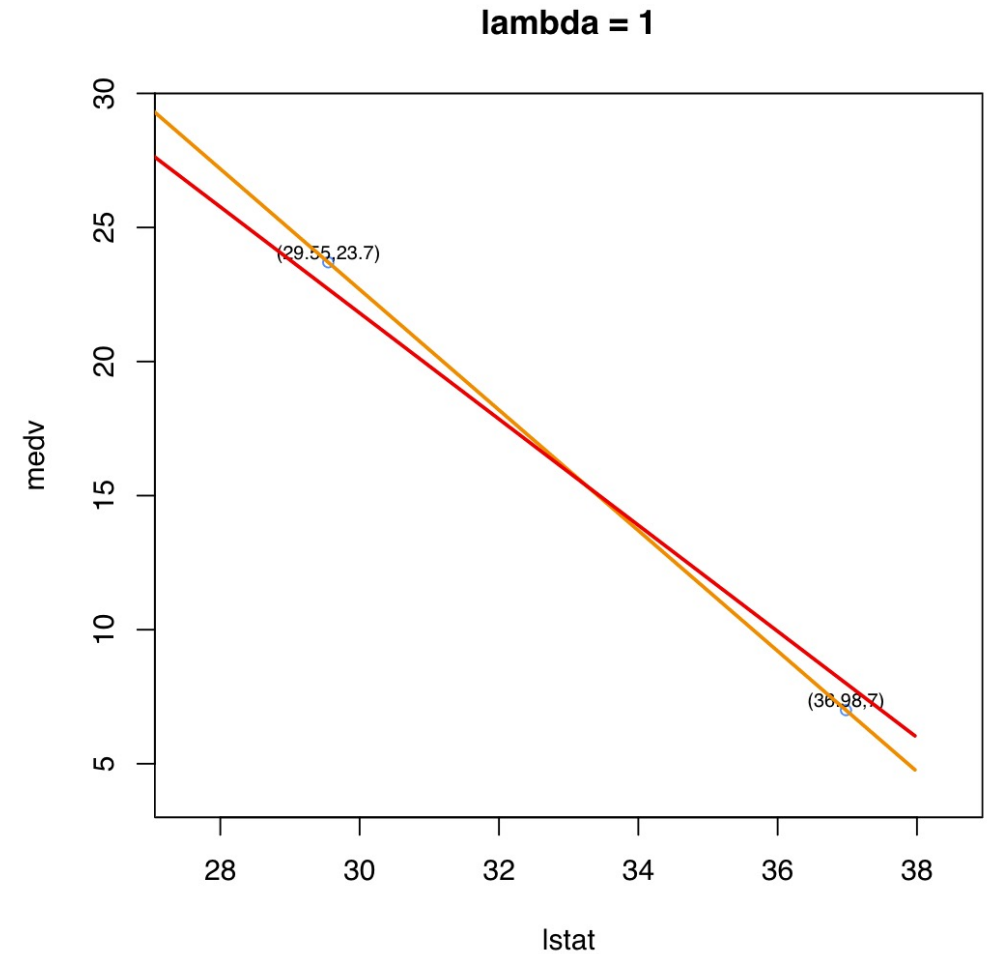- What if we want to perform variable selection?

# Lasso

- Lasso: least absolute shrinkage and selection operator

- Lasso minimizes
  - $\sum_{i=1}^{n}(medv_i - \beta_0 - lstat_i \cdot \beta_1)^2 + \lambda \cdot |\beta_1|$

  - $\lambda \geq 0$:  tuning hyper-parameter

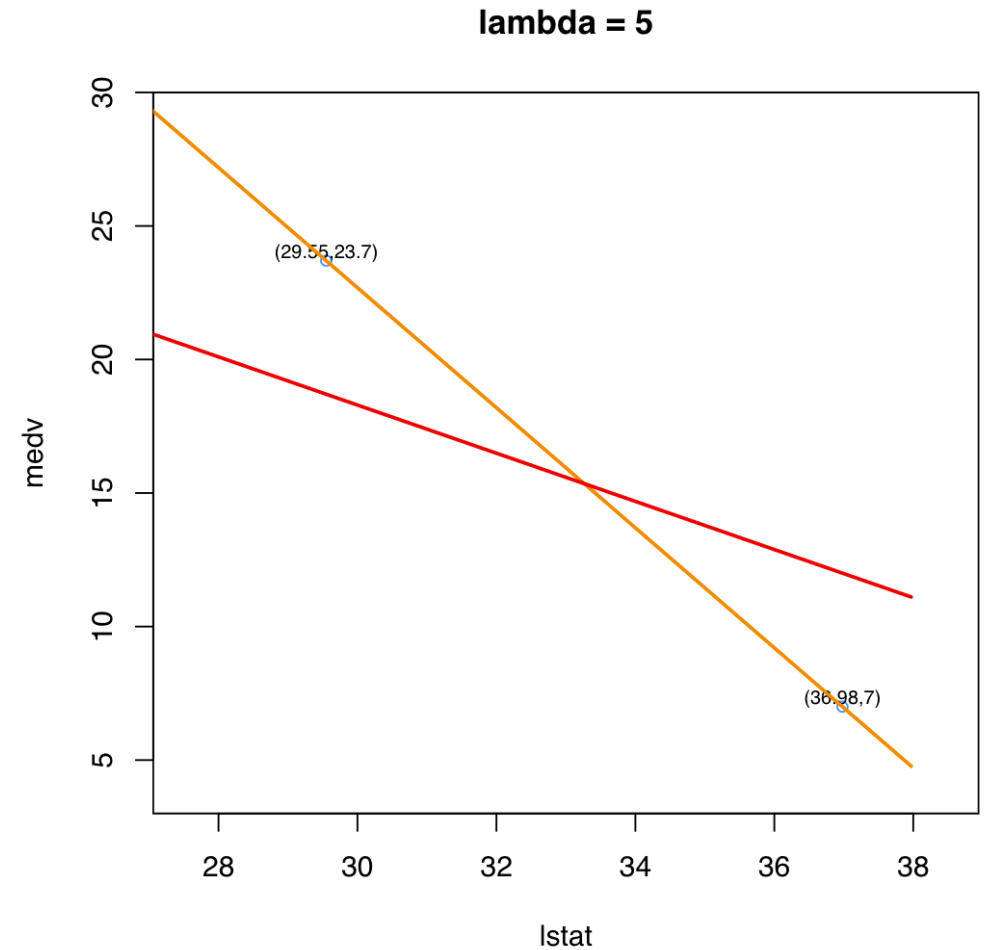# Role of $\lambda$ in Lasso

- Lasso minimizes
  - $\sum_{i=1}^n (medv_i - \beta_0 - lstat_i \cdot \beta_1)^2 + \lambda \cdot |\beta_1|$

  - $\lambda = 1 : \hat{\beta}_1^L = -1.978$



lambda = 1
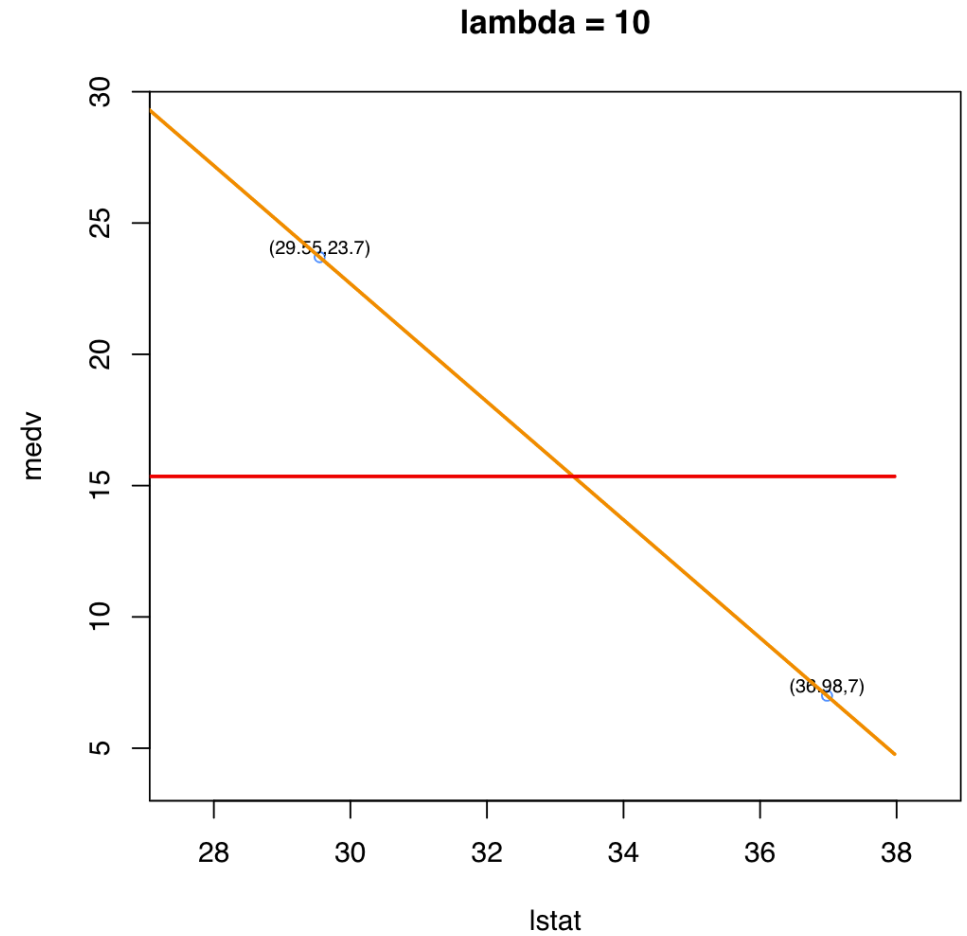
# Role of $\lambda$ in Lasso

- Lasso minimizes
  - $\sum_{i=1}^{n}(medv_i - \beta_0 - lstat_i \cdot \beta_1)^2 + \lambda \cdot |\beta_1|$

  - $\lambda = 5 : \hat{\beta}_1^L = -0.902$

# Role of $\lambda$ in Lasso

- Lasso minimizes
  - $\sum_{i=1}^{n}(medv_i - \beta_0 - lstat_i \cdot \beta_1)^2 + \lambda \cdot |\beta_1|$

  - $\lambda = 10 : \hat{\beta}_1^L = 0$



lambda = 10

# Lasso for more than one predictor
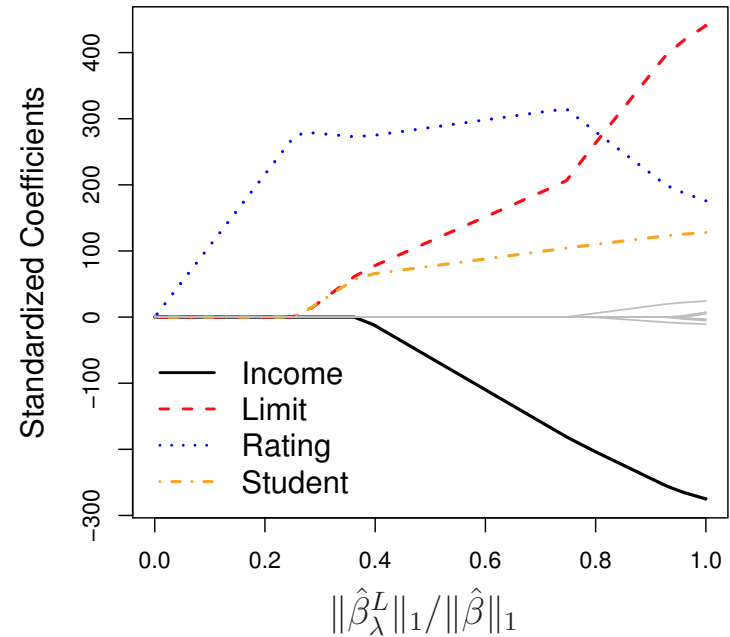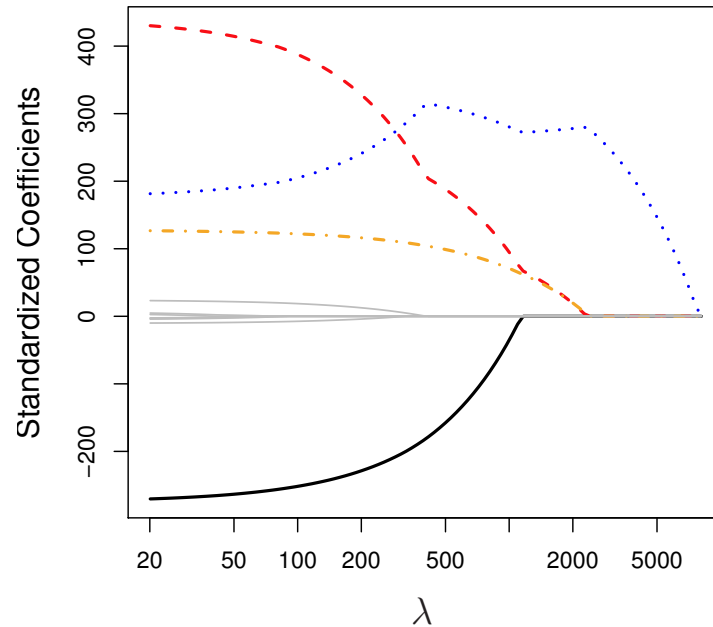
- Lasso minimizes

$$\sum_{i=1}^{n}\left(Y_i - \beta_0 - \sum_{j=1}^{p}\beta_j X_{i,j}\right)^2 + \lambda \sum_{j=1}^{p}|\beta_j|$$

- $X_{i,j}$: $j$-th predictor of $i$-th observation

- $\|\beta\|_1 = \sum_{j=1}^{p}|\beta_j|$: $\|\beta\|_1$ is called the $\ell_1$ norm of $\beta \in \mathbb{R}^p$

- $\beta_0$: mean of $Y_i$

- Shrinkage penalty $\lambda$ does not apply to $\beta_0$

# Example: Credit card data set (lasso)

- Predict default or not; **11** predictors
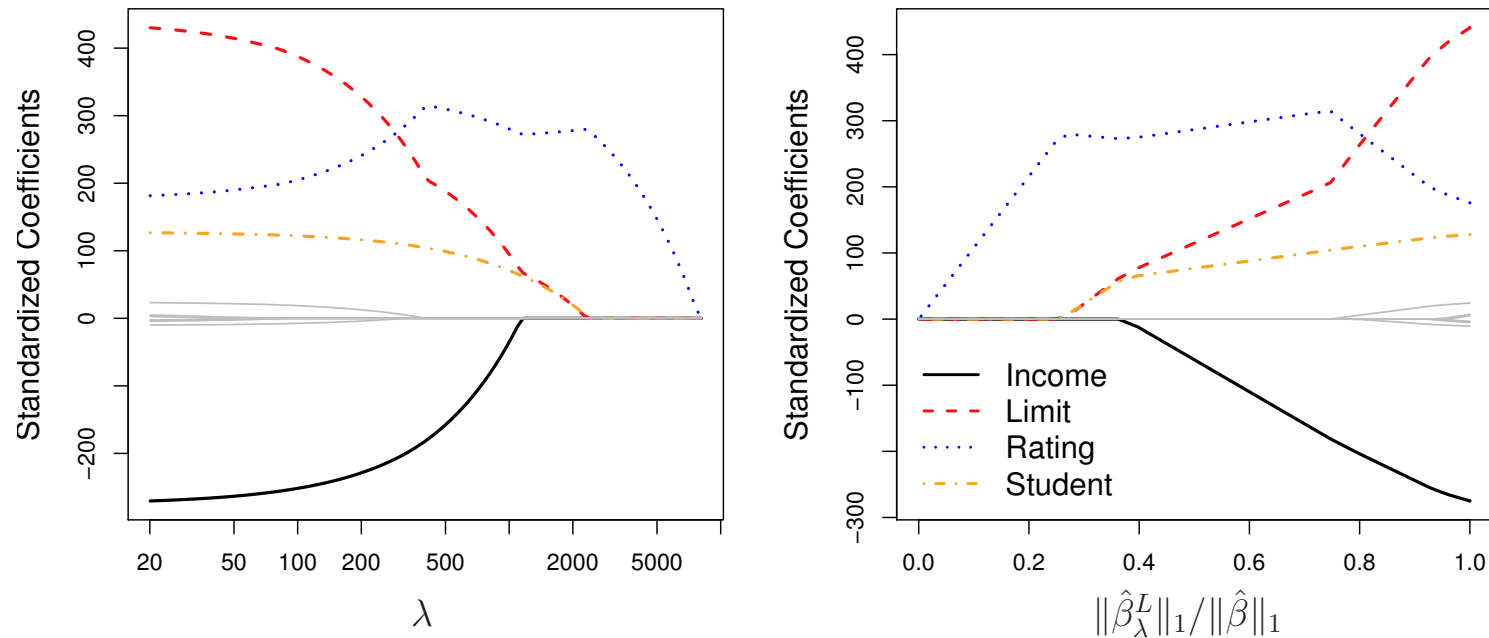  - $\|\beta\|_1 = \sum_{j=1}^{p} |\beta_j|$



**Shrinkage ratio**

- **Shrinkage ratios**: coefficients shrink to zero at varying rates

# Example: Credit card data set (lasso)

- Predict default or not; **11** predictors



- **Variable selection**: As $\lambda$ increases, lasso selects less variables
  - {"empty"} $\to$ {rating} $\to$ {limit, rating, student} $\to$ {income, limit, rating, student}
- **Lasso path**: Different coefficient values by varying $\lambda$

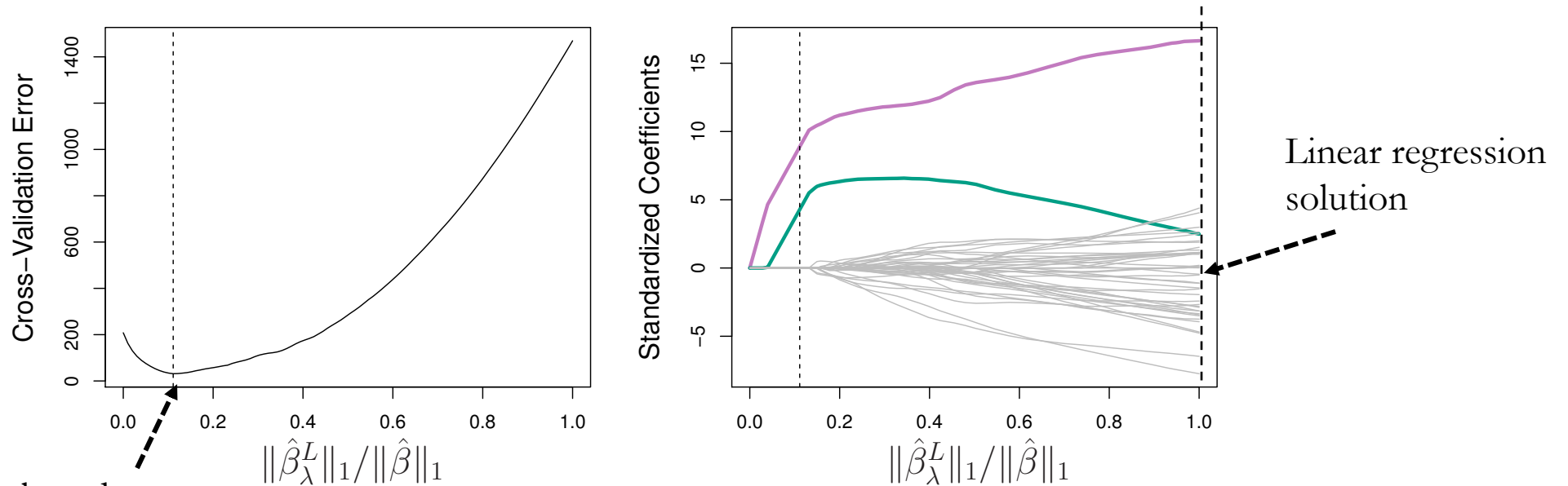# Choose $\lambda$ by cross-validation

- The procedure is the same for ridge and lasso

  1. Choose a grid of $\lambda$ values

  2. Compute the cross-validation error for each $\lambda$ value

  3. Select the $\lambda$ with the smallest cross-validation error

  4. Refit the model using all observations and selected $\lambda$

# Example

- **Simulation I:** Only 2 coefficients are non-zero
  - Simulated data: 45 predictors, 2 out of $\beta_1, \ldots, \beta_{45}$ are nonzero
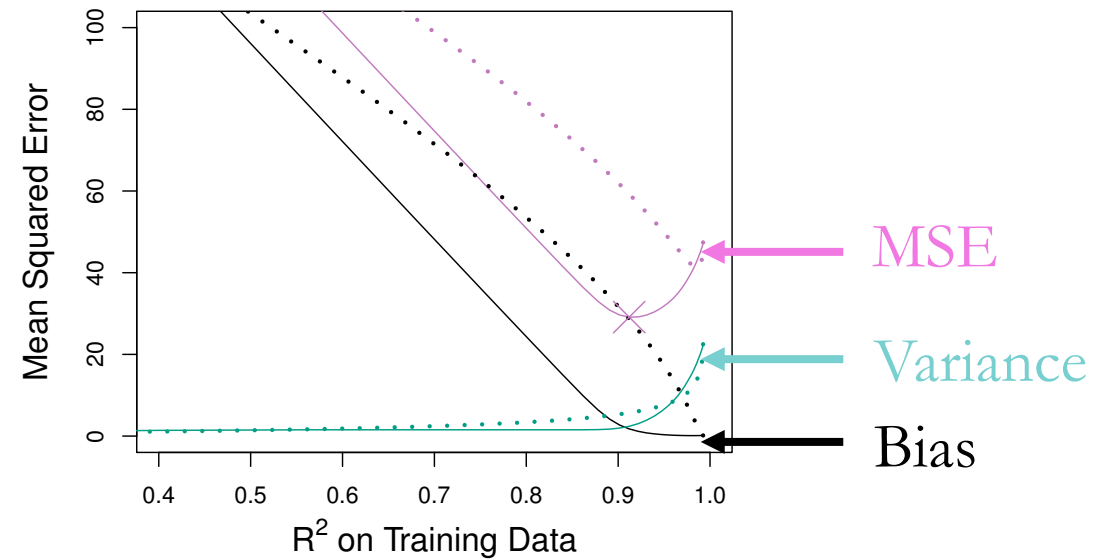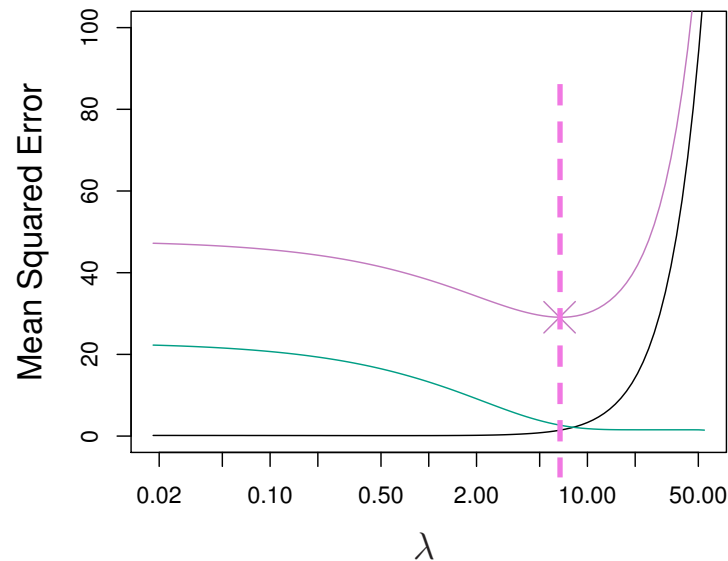  - **10-fold CV to select the lasso regularization parameter**



Optimal $\lambda$ selected by cross-validation

Linear regression solution

# Lasso vs. Ridge regularization

- **Simulation I:** Only 2 coefficients are non-zero
  - Simulated data: 45 predictors, 2 out of $\beta_1, \ldots, \beta_{45}$ are nonzero
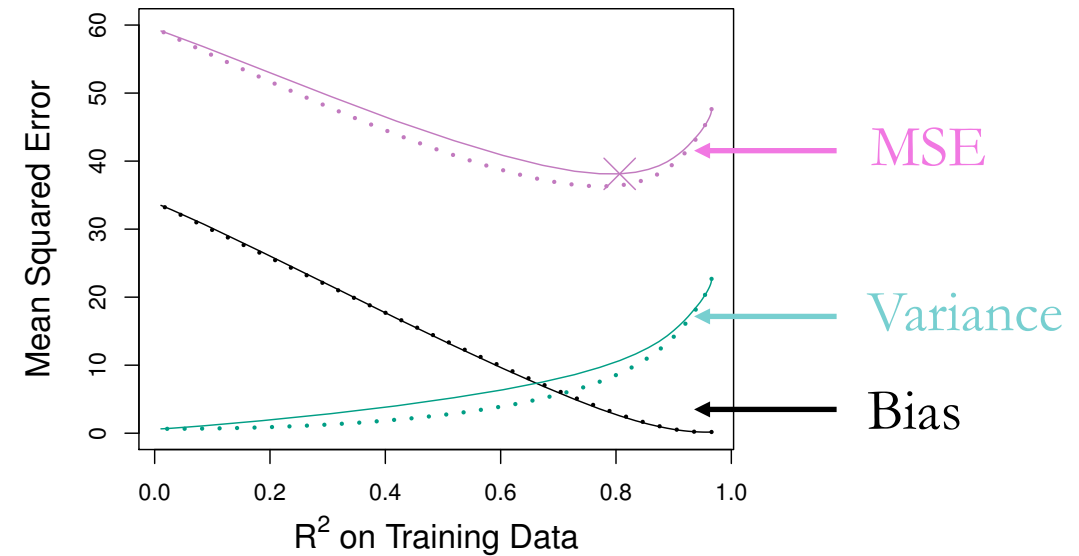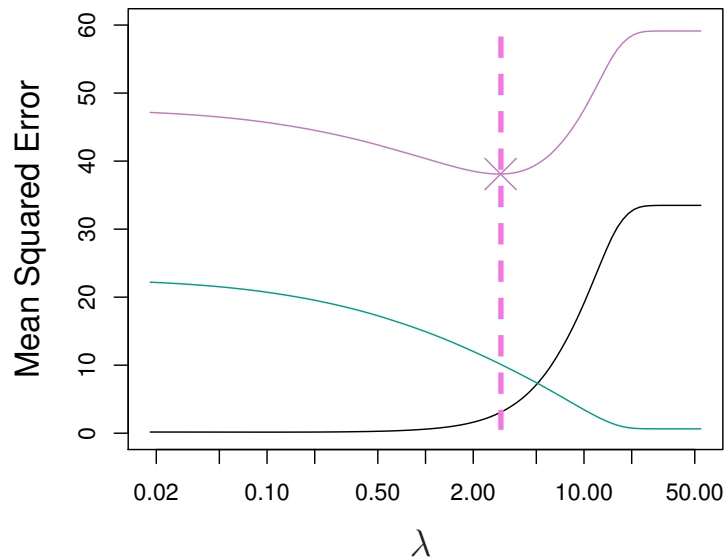
Solid lines (—): Lasso
Dash lines (⋯): Ridge



- The **bias**, variance, and MSE are all lower for the lasso

# Lasso vs. Ridge regularization

- **Simulation II:** Most of the coefficients are non-zero
  - Simulated data: 45 predictors $\beta_1, \ldots, \beta_{45}$ are nonzero

Solid lines (—): Lasso
Dash lines (···): Ridge



- The **variance** of ridge regression is smaller
- The **bias** is about the same for both
- Hence the MSE of ridge regression is smaller

# Lasso vs. Ridge regularization

- **Takeaways**: Neither ridge nor the lasso universally dominates

  - Lasso performs better if **a small number of predictors with large coefficients**

  - Ridge performs better if **many predictors with similar coefficients**

  - Select which one by **cross-validation** ☺
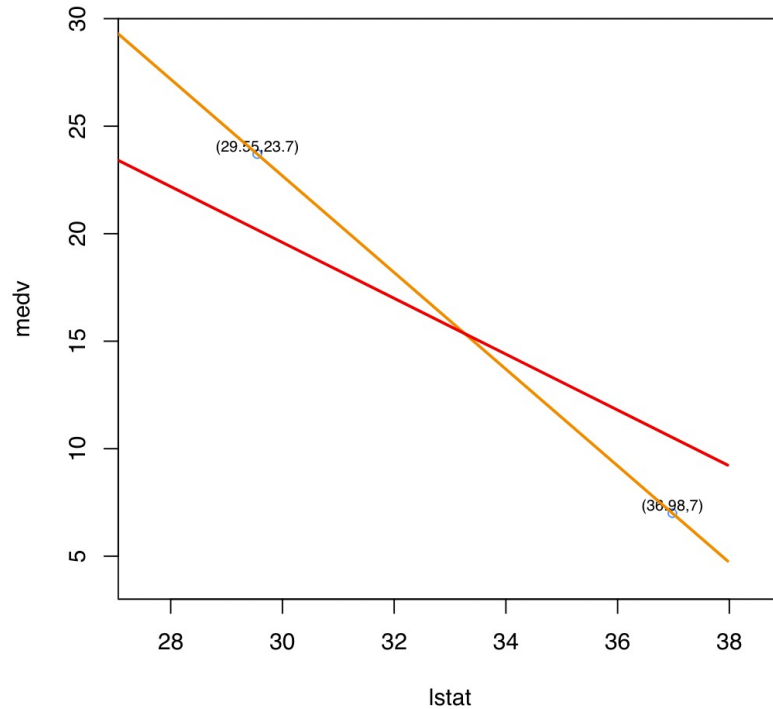
# Lecture plan

- Lasso


- Elastic net

# Elastic net

- Elastic net combines lasso and ridge penalty, and minimizes

  - $\sum_{i=1}^{n}(medv_i - \beta_0 - lstat_i \cdot \beta_1)^2 + \lambda \cdot (1 - \alpha) \cdot \frac{\beta_1^2}{2} + \lambda \cdot \alpha \cdot |\beta_1|$

  - $\lambda \geq 0$:  tuning hyper-parameter

  - $\alpha \in [0,1]$: tuning hyper-parameter

    - $\alpha = 0$: ridge
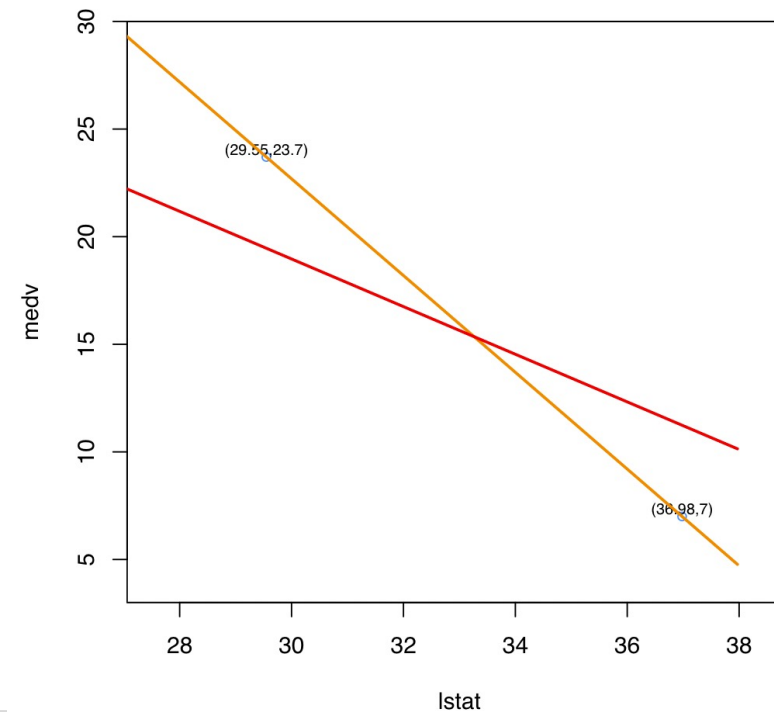
    - $\alpha = 1$: lasso

# Role of $\alpha$ and $\lambda$ in elastic net

- Elastic net combines lasso and ridge penalty, and minimizes
  - $\sum_{i=1}^{n}(medv_i - \beta_0 - lstat_i \cdot \beta_1)^2 + \lambda \cdot (1 - \alpha) \cdot \frac{\beta_1^2}{2} + \lambda \cdot \alpha \cdot |\beta_1|$
  - $\alpha = 0.3, \lambda = 5$: $\hat{\beta}_1^E = -1.299$; $\alpha = 0.7, \lambda = 5$: $\hat{\beta}_1^E = -1.107$

# Role of $\alpha$ and $\lambda$ in elastic net
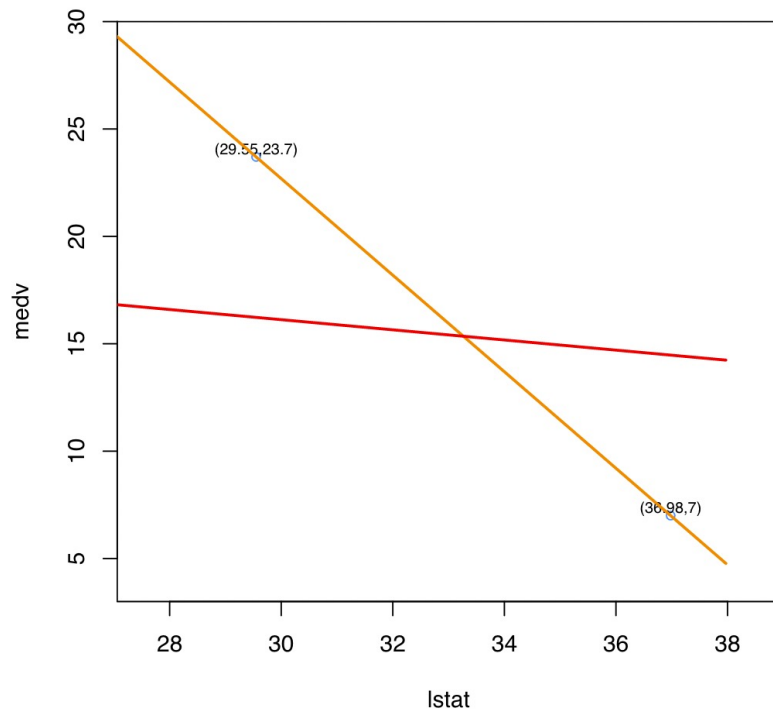
- Elastic net combines lasso and ridge penalty, and minimizes
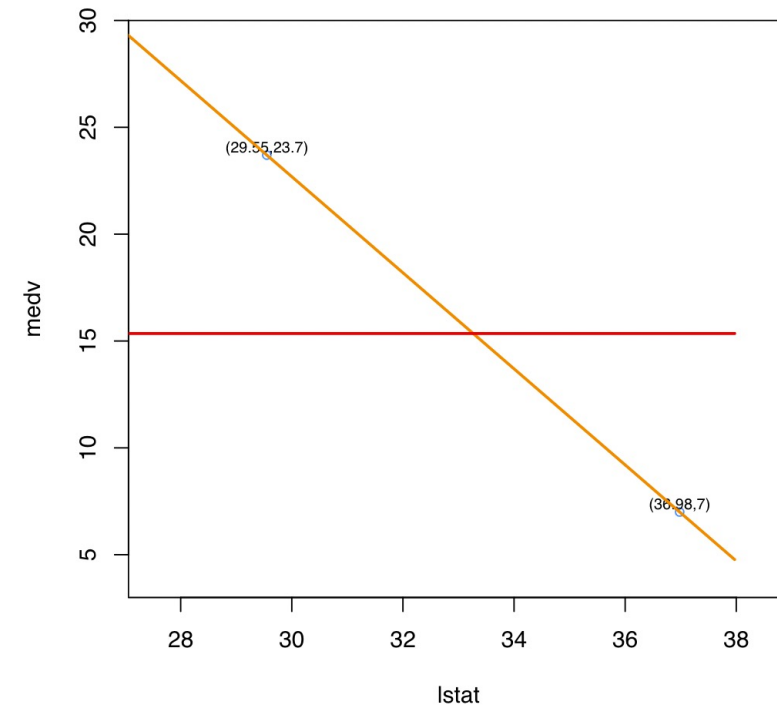
  - $\sum_{i=1}^{n}(medv_i - \beta_0 - lstat_i \cdot \beta_1)^2 + \lambda \cdot (1-\alpha) \cdot \frac{\beta_1^2}{2} + \lambda \cdot \alpha \cdot |\beta_1|$

  - $\alpha = 0.3, \lambda = 20$: $\hat{\beta}_1^E = -0.236$; $\alpha = 0.7, \lambda = 20$: $\hat{\beta}_1^E = 0$

# Role of $\alpha$ and $\lambda$ in elastic net

- Elastic net combines lasso and ridge penalty, and minimizes

  - $\sum_{i=1}^{n}(medv_i - \beta_0 - lstat_i \cdot \beta_1)^2 + \lambda \cdot (1-\alpha) \cdot \frac{\beta_1^2}{2} + \lambda \cdot \alpha \cdot |\beta_1|$

  - $\alpha = 0.3, \lambda = 50$: $\hat{\beta}_1^E = 0$; $\alpha = 0.7, \lambda = 50$: $\hat{\beta}_1^E = 0$


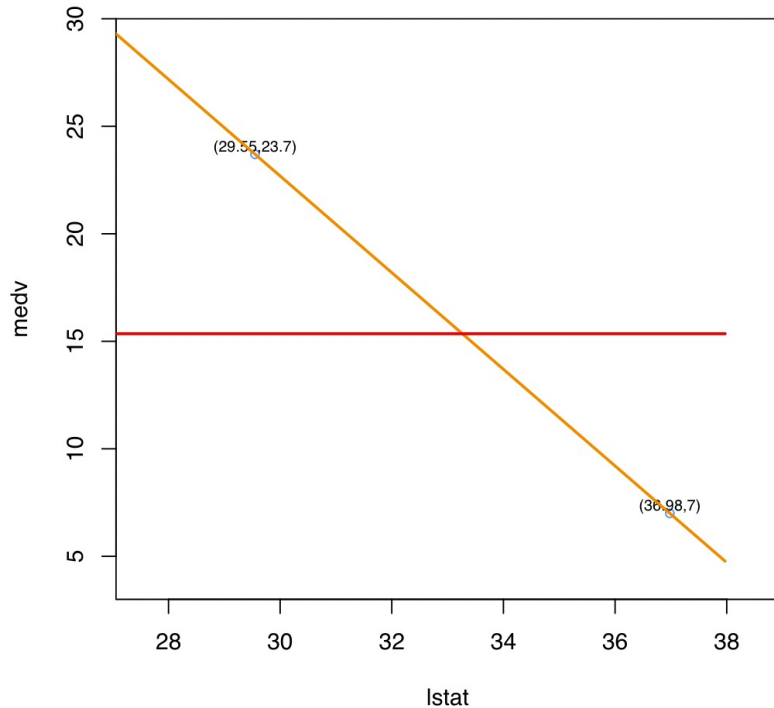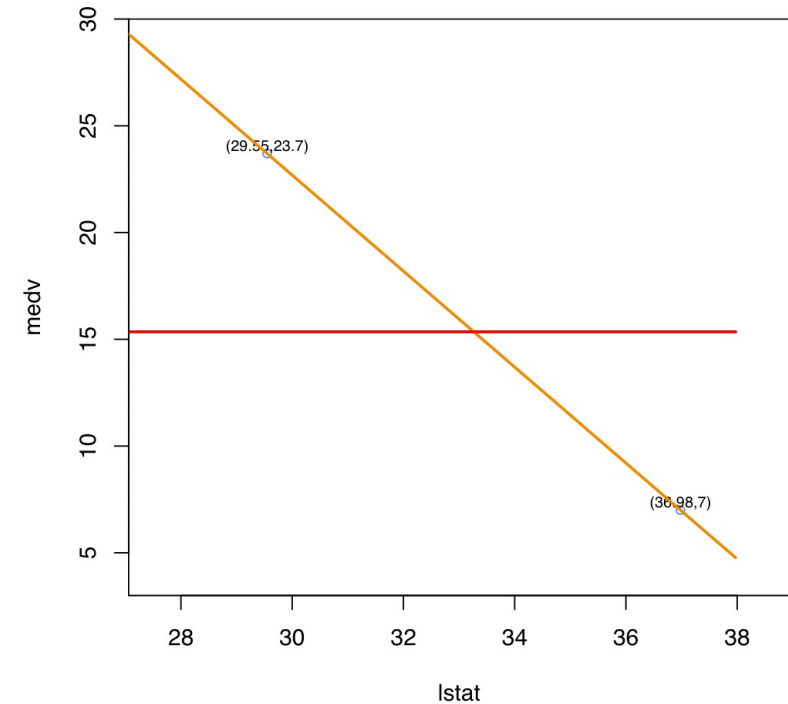
alpha = 0.3, lambda = 50

alpha = 0.7, lambda = 50

# Choose $\alpha$ and $\lambda$ by cross-validation

- The procedure is the same for ridge and lasso

1. Choose a grid of $\alpha$ values and a grid of $\lambda$ values

2. Compute the cross-validation error for each $(\alpha, \lambda)$ value

3. Select the $(\alpha, \lambda)$ with the smallest cross-validation error

4. Refit the model using all observations and selected $(\alpha, \lambda)$