

Ch 6.2: Shrinkage - Ridge regression

Lecture 17 - CMSE 381

Michigan State University

::

Dept of Computational Mathematics, Science & Engineering

Wed, Feb 25, 2026

Announcements

Last time:

- Subset selection

This time:

- Ridge regression

Announcements:

- HW #3 due Sunday 3/8

| | | | | | | |
|----|---|------|--|---------|----------------------|----|
| 11 | F | 2/6 | Multiple Logistic Regression / Multinomial Logistic Regression | 4.3.4-5 | HW #2 Due Mon 2/9 | |
| | M | 2/9 | Project Day & Review | | | |
| | W | 2/11 | Midterm #1 | | | |
| 12 | F | 2/13 | Class not held | | | |
| 13 | M | 2/16 | Leave one out and k-fold CV | 5.1.1-3 | | Q5 |
| 14 | W | 2/18 | More k-fold CV | 5.1.4-5 | | |
| 15 | F | 2/20 | k-fold CV for classification | 5.1.5 | | |
| 16 | M | 2/23 | Subset selection | 6.1 | | |
| 17 | W | 2/25 | Shrinkage: Ridge | 6.2.1 | | |
| 18 | F | 2/27 | Shrinkage: Lasso | 6.2.2 | | |
| | M | 3/2 | Spring Break | | | |
| | W | 3/4 | Spring Break | | | |
| | F | 3/6 | Spring Break | | HW #3 Due Sun 3/8 | |
| 19 | M | 3/9 | PCA | 6.3 | | Q6 |
| 20 | W | 3/11 | PCR | 6.3 | | |

Section 1

Last time

Subset selection

Algorithm 6.1 Best subset selection

1. Let \mathcal{M}_0 denote the *null model*, which contains no predictors. This model simply predicts the sample mean for each observation.
 2. For $k = 1, 2, \dots, p$:
 - (a) Fit all $\binom{p}{k}$ models that contain exactly k predictors.
 - (b) Pick the best among these $\binom{p}{k}$ models, and call it \mathcal{M}_k . Here *best* is defined as having the smallest RSS, or equivalently largest R^2 .
 3. Select a single best model from among $\mathcal{M}_0, \dots, \mathcal{M}_p$ using cross-validated prediction error, C_p (AIC), BIC, or adjusted R^2 .
-

Algorithm 6.2 Forward stepwise selection

1. Let \mathcal{M}_0 denote the *null model*, which contains no predictors.
 2. For $k = 0, \dots, p - 1$:
 - (a) Consider all $p - k$ models that augment the predictors in \mathcal{M}_k with one additional predictor.
 - (b) Choose the *best* among these $p - k$ models, and call it \mathcal{M}_{k+1} . Here *best* is defined as having smallest RSS or highest R^2 .
 3. Select a single best model from among $\mathcal{M}_0, \dots, \mathcal{M}_p$ using cross-validated prediction error, C_p (AIC), BIC, or adjusted R^2 .
-

Algorithm 6.3 Backward stepwise selection

1. Let \mathcal{M}_p denote the *full model*, which contains all p predictors.
 2. For $k = p, p - 1, \dots, 1$:
 - (a) Consider all k models that contain all but one of the predictors in \mathcal{M}_k , for a total of $k - 1$ predictors.
 - (b) Choose the *best* among these k models, and call it \mathcal{M}_{k-1} . Here *best* is defined as having smallest RSS or highest R^2 .
 3. Select a single best model from among $\mathcal{M}_0, \dots, \mathcal{M}_p$ using cross-validated prediction error, C_p (AIC), BIC, or adjusted R^2 .
-

What should you learn from this lecture?

- What is regularization? Why do we need it?
- What are the two basic types of regularization methods? How are they implemented mathematically in linear regression?
- How do you fit a ridge regression model in python?
- How do you control the model flexibility & bias-variance tradeoff when using regularization?
- How do you find the right amount of regularization using cross-validation? How do you do this in python?
- What additional precautions do you need to take when using regularization (compared to least squares)?
- What are the advantages of regularization compared to Least Squares?
- What are the advantages of regularization compared to subset selection?

Section 2

Ridge Regression

- Fit model using all p predictors
- Aim to constrain (regularize) coefficient estimates
- Shrink the coefficient estimates towards 0

$$Y = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \beta_3 X_3 + \beta_4 X_4$$

Ridge regression

Before:

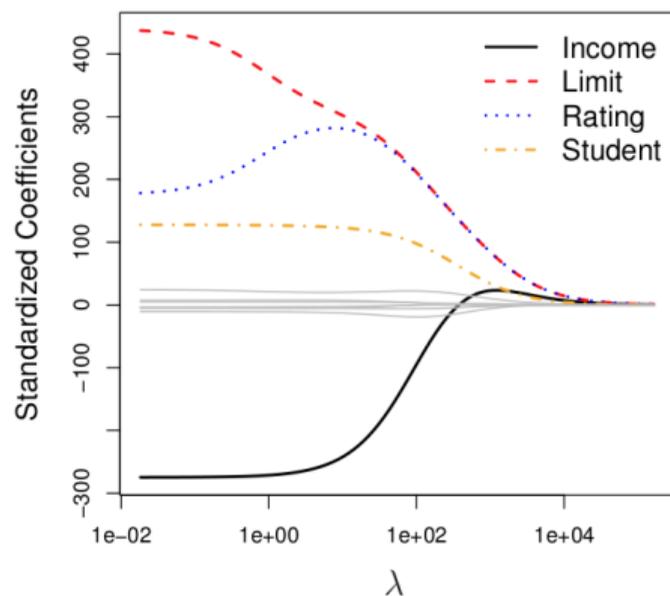
$$RSS = \sum_{i=1}^n \left(y_i - \beta_0 - \sum_{j=1}^p \beta_j x_{ij} \right)^2$$

After:

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

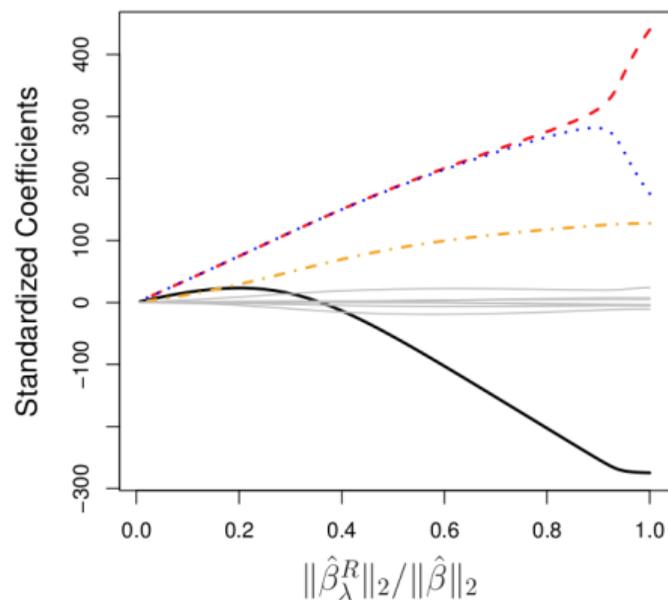
Example from the Credit data

$$RSS + \lambda \sum_{j=1}^p \beta_j^2$$



Same Setting, Different Plot

$$RSS + \lambda \sum_{j=1}^p \beta_j^2 \quad \|\beta\|_2 = \sqrt{\sum_{j=1}^p \beta_j^2}$$



Scale equivariance (or lack thereof)

Scale equivariant: Multiplying a variable by c (cX_i) just returns a coefficient multiplied by $1/c$ ($1/c\beta_i$)

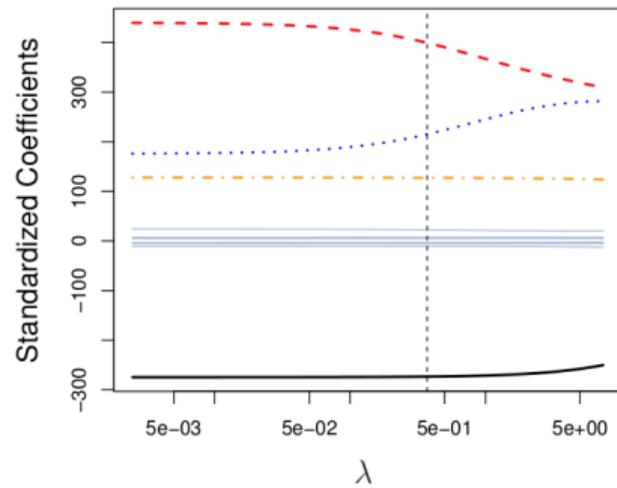
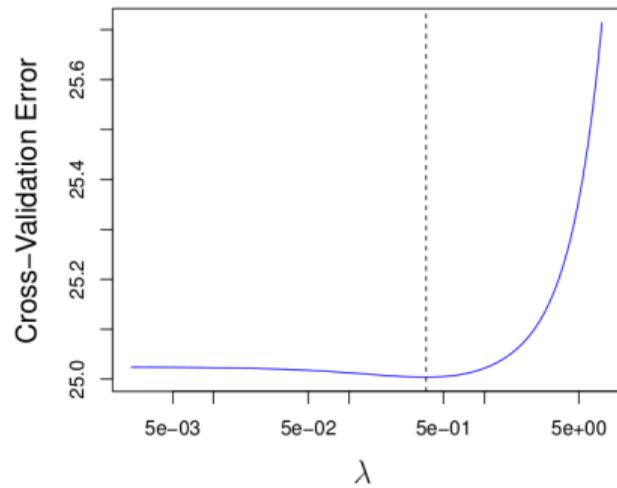
Solution: Standardize predictors

$$\tilde{x}_{ij} = \frac{x_{ij}}{\sqrt{\frac{1}{n} \sum_{i=1}^n (x_{ij} - \bar{x}_j)^2}}$$

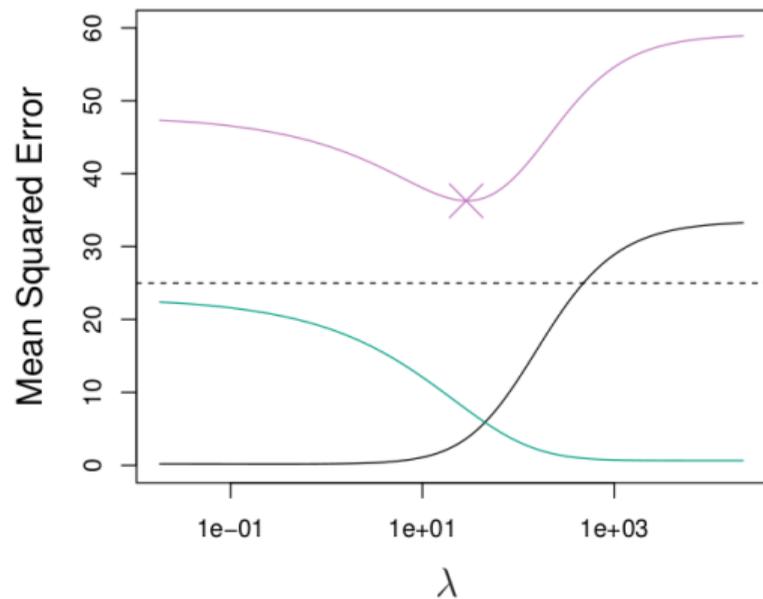
Using Cross-Validation to find λ

- Choose a grid of λ values
- Compute the (k -fold) cross-validation error for each value of λ
- Select the tuning parameter value λ for which the CV error is smallest.
- The model is re-fit using all of the available observations and the selected value of the tuning parameter.

LOOCV choice of λ for ridge regression and Credit data

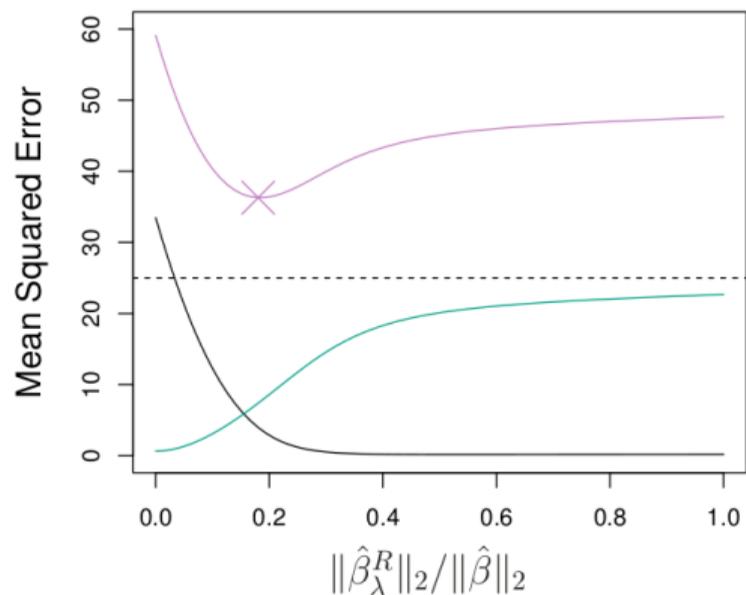


Bias-Variance tradeoff



Squared bias (black), variance (green), and test mean squared error (purple) for simulated data.

More Bias-Variance Tradeoff



Squared bias (black), variance (green), and test mean squared error (purple) for simulated data.

Advantages of Ridge

Ridge vs. Least Squares:

Ridge vs. Subset Selection:

Look back and look ahead

| | | | | | | |
|----|---|------|--|---------|----------------------|----|
| 11 | F | 2/6 | Multiple Logistic Regression / Multinomial Logistic Regression | 4.3.4-5 | HW #2 Due Mon 2/9 | |
| | M | 2/9 | Project Day & Review | | | |
| | W | 2/11 | Midterm #1 | | | |
| 12 | F | 2/13 | Class not held | | | |
| 13 | M | 2/16 | Leave one out and k-fold CV | 5.1.1-3 | | Q5 |
| 14 | W | 2/18 | More k-fold CV | 5.1.4-5 | | |
| 15 | F | 2/20 | k-fold CV for classification | 5.1.5 | | |
| 16 | M | 2/23 | Subset selection | 6.1 | | |
| 17 | W | 2/25 | Shrinkage: Ridge | 6.2.1 | | |
| 18 | F | 2/27 | Shrinkage: Lasso | 6.2.2 | | |
| | M | 3/2 | Spring Break | | | |
| | W | 3/4 | Spring Break | | | |
| | F | 3/6 | Spring Break | | HW #3 Due Sun 3/8 | |
| 19 | M | 3/9 | PCA | 6.3 | | |
| 20 | W | 3/11 | PCR | 6.3 | | Q6 |

- What is regularization? Why do we need it?
- What are the two basic types of regularization methods? How are they implemented mathematically in linear regression?
- How do you fit a ridge regression model in python?
- How do you control the model flexibility & bias-variance tradeoff when using regularization?
- How do you find the right amount of regularization using cross-validation? How do you do this in python?
- What additional precautions do you need to take when using regularization (compared to least squares)?
- What are the advantages of regularization compared to Least Squares?
- What are the advantages of regularization compared to subset selection?