Ch 6.2: Shrinkage - Ridge regression Lecture 17 - CMSE 381

Prof. Mengsen Zhang

Michigan State University

:

Dept of Computational Mathematics, Science & Engineering

Wed, Oct 8, 2025

Announcements

Last time:

Subset selection

This time:

Ridge regression

Announcements:

HW #4 due Sunday 10/12

CMSE381 F2025 Schedule : Schedule

Civic	L30 i	_1 2020	_ochedule . ochedule			
	M	9/22	Project Day & Review			
	W	9/24	Midterm #1			
12	F	9/26	Leave one out CV	5.1.1, 5.1.2		
13	М	9/29	k-fold CV	5.1.3		
14	W	10/1	More k-fold CV	5.1.4-5		
15	F	10/3	k-fold CV for classification	5.1.5		
16	М	10/6	Subset selection	6.1		
17	W	10/8	Shrinkage: Ridge	6.2.1		
18	F	10/10	Shrinkage: Lasso	6.2.2	HW #4 Due	
19	М	10/13	PCA	6.3	Sun 10/12	
20	W	10/15	PCR	6.3		
	F	10/17	Review			
	М	10/20	Fall Break			
	W	10/22	Midterm #2			
21	F	10/24	Polynomial & Step Functions	7.1-7.2	HW #5 Due Sun 10/28	
22	М	10/27	Step Functions; Basis functions; Start Splines	7.2-7.4		
23	W	10/29	Regression Splines	7.4		

2/19

. Zhang (MSU-CMSE) Wed, Oct 8, 2025

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 .
- Select a single best model from among M₀,..., M_p using crossvalidated prediction error, C_p (AIC), BIC, or adjusted R².

Algorithm 6.2 Forward stepwise selection

- Let M₀ denote the null model, which contains no predictors.
- 2. For $k = 0, \ldots, p-1$:
 - (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 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 \mathbb{R}^2 .
- Select a single best model from among M₀,..., M_p using crossvalidated prediction error, C_p (AIC), BIC, or adjusted R².

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?

Dr. Zhang (MSU-CMSE) Wed, Oct 8, 2025 5/19

Section 2

Ridge Regression

Goal

- 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:

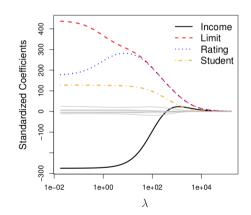
$$RSS = \sum_{i=1}^{n} \left(y_i - \beta_0 - \sum_{j=1}^{p} \beta_j x_{ij} \right)$$

After:

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

Example from the Credit data

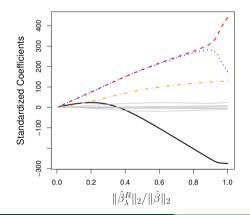
$$RSS + \lambda \sum_{j=1}^{p} \beta_j^2$$



Same Setting, Different Plot

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

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



Test your understanding:PollEv

10 / 19

Wed, Oct 8, 2025

Scale equivavariance (or lack thereof)

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

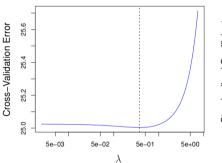
Solution: Standardize predictors

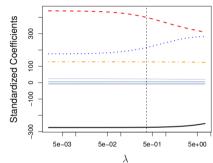
$$\widetilde{x}_{ij} = \frac{x_{ij}}{\sqrt{\frac{1}{n}\sum_{i=1}^{n}(x_{ij} - \overline{x}_{j})^{2}}}$$

Using Cross-Validation to find λ

- ullet 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

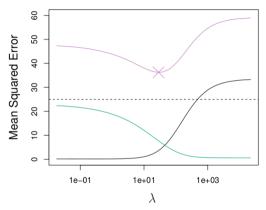




Coding

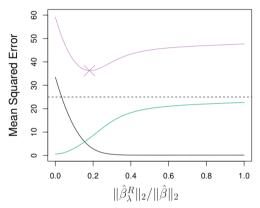
Dr. Zhang (MSU-CMSE)

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

CMSE381	F2025	Schedule	 Schedule

	M	9/22	Project Day & Review			
	W	9/24	Midterm #1			
12	F	9/26	Leave one out CV	5.1.1, 5.1.2		
13	M	9/29	k-fold CV	5.1.3		
14	W	10/1	More k-fold CV	5.1.4-5		
15	F	10/3	k-fold CV for classification	5.1.5		
16	M	10/6	Subset selection	6.1		
17	W	10/8	Shrinkage: Ridge	6.2.1		
18	F	10/10	Shrinkage: Lasso	6.2.2	HW #4 Due	
19	M	10/13	PCA	6.3	Sun 10/12	
20	W	10/15	PCR	6.3		
	F	10/17	Review			
	M	10/20	Fall Break			
	W	10/22	Midterm #2			
21	F	10/24	Polynomial & Step Functions	7.1-7.2	HW #5 Due Sun 10/28	
22	М	10/27	Step Functions; Basis functions; Start Splines	7.2-7.4		
23	W	10/29	Regression Splines	7.4		

- 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?