

Ch 6.2: Shrinkage - The Lasso

Lecture 12 - CMSE 381

Michigan State University

::

Dept of Computational Mathematics, Science & Engineering

March 12, 2024

Last time:

- Ridge Regression

This time:

- The Lasso

Announcements:

- Midterm score release tomorrow
- Can ask questions about grade before the end of recitation this week
- Mean: 72.5 Median: 72.5, curve?

Section 1

Last time

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

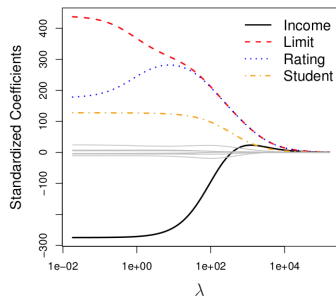
$$Y = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \beta_3 X_3 + \beta_4 X_4$$

- Ridge regression
- Lasso

Ridge regression

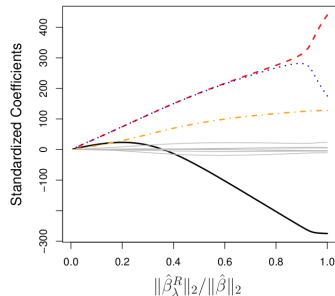
Before:

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



After:

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



Scale equivariance (or lack thereof)

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

- Least squares is scale equivariant
- Ridge regression is not

Solution: standardize predictors

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

- Ex: income variable. If you put in income by \$1,000's of dollars, the answer you would get is quite different. Might emphasize a different coefficient.
- Least squares is scale equivariant
- Ridge regression very much is not
- $X_j \hat{\beta}_{j,\lambda}^R$ depends not only on λ but also on values of other predictors

Section 2

The Lasso

Same goal as before

- Fit model using all p predictors
- Aim to constrain (regularize) coefficient estimates
- Shrink the coefficient estimates towards 0 *Ridge shrunk but DID NOT GET RID OF which is what we actually want*

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

- Ridge regression
- Lasso

The lasso

Least Squares:

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

Ridge:

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

The Lasso:

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

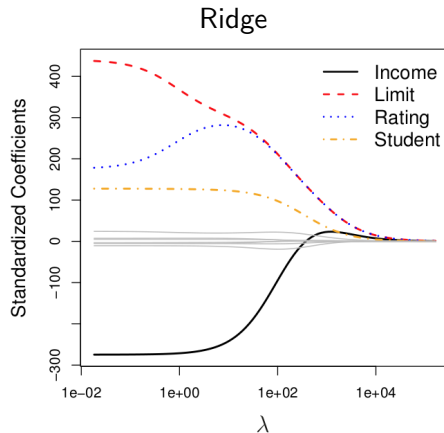
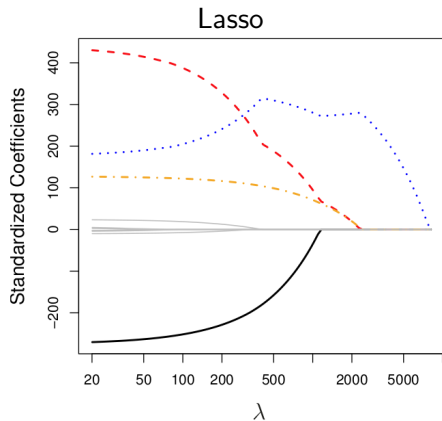
- ℓ_1 penalty instead of ℓ_2 penalty
- ℓ_1 norm: $\|\beta\|_1 = \sum |\beta_j|$
- ℓ_2 norm: $\|\beta\|_2 = \sqrt{\sum \beta_j^2}$

Subsets with lasso

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

- Result of ℓ_1 penalty is that it forces some coefficients to be exactly 0
- Can be interpreted as something more like subset selection, but here called *variable selection*
- Yield are *sparse* models, that is those using only a subset of variables
- Selecting good λ is critical

An example on Credit data set



- $\lambda = 0$: Both are least squares
- $\lambda = \infty$: Everything is 0
-

Why Lasso shrinks coefficients to 0, - a mathematical verification on a simple example

consider a simple special case with

- $n = p$,
- $x_j = e_j$ the j -th standard basis vector
- no intercept.

Then the usual least squares objective becomes

$$\sum_{i=1}^p (y_i - \beta_i)^2$$

The Lasso objective is

$$\sum_{i=1}^p (y_i - \beta_i)^2 + \lambda \sum_{i=1}^p |\beta_i|$$

the Ridge regression objective is

$$\sum_{i=1}^p (y_i - \beta_i)^2 + \lambda \sum_{i=1}^p \beta_i^2$$

- the least squares solution is $\hat{\beta}_j = y_j$
- the ridge solution is $\hat{\beta}_j^R = \frac{y_j}{1+\lambda}$
- the Lasso solution is

$$\hat{\beta}_j^L = \begin{cases} y_j - \frac{\lambda}{2}, & y_j \geq \frac{\lambda}{2} \\ y_j + \frac{\lambda}{2}, & y_j \leq -\frac{\lambda}{2} \\ 0, & |y_j| < \frac{\lambda}{2} \end{cases}$$

Proof of the first case: $y_j \geq \frac{\lambda}{2}$ (the proof only requires the knowledge of calculus).

The objective function of Lasso for this simple case is

$$obj \equiv \sum_{i=1}^p (y_i - \beta_i)^2 + \lambda \sum_{i=1}^p |\beta_i| = (y_1 - \beta_1)^2 + \lambda |\beta_1| + C$$

where C contains terms in obj irrelevant to β_1 . Based on the fact that $|\beta_1| = \beta_1$ if $\beta_1 \geq 0$, and $|\beta_1| = -\beta_1$ if $\beta_1 \leq 0$. We can remove the absolute value on β_1 and arrive at the expression

$$obj = \begin{cases} (y_1 - \beta_1)^2 + \lambda \beta_1 + C & \text{if } \beta_1 \geq 0 \\ (y_1 - \beta_1)^2 - \lambda \beta_1 + C & \text{if } \beta_1 \leq 0 \end{cases}$$

Taking derivative with respect to β_1 :

$$\frac{\partial obj}{\partial \beta_1} = \begin{cases} 2(\beta_1 - y_1) + \lambda & \text{if } \beta_1 \geq 0 \\ 2(\beta_1 - y_1) - \lambda & \text{if } \beta_1 \leq 0 \end{cases}$$

The obj has two pieces. For the piece on $\beta \geq 0$, local min occurs at $\frac{\partial obj}{\partial \beta_1} = 0 \Leftrightarrow \beta_1 = y_1 - \frac{\lambda}{2}$. This is also the global min of obj on the interval $\beta \geq 0$. For the piece of obj on $\beta \leq 0$, there is no local min, so the global min must occur at the boundary $-\infty$ or 0 . Comparing the objective values at the three candidate points $\beta_1 = y_1 - \frac{\lambda}{2}, 0, -\infty$, we found the $\beta_1 = y_1 - \frac{\lambda}{2}$ is the global min.

Scale equivariance (or lack thereof)

Scale equivariant: Multiplying a variable by c just returns a coefficient multiplied by $1/c$

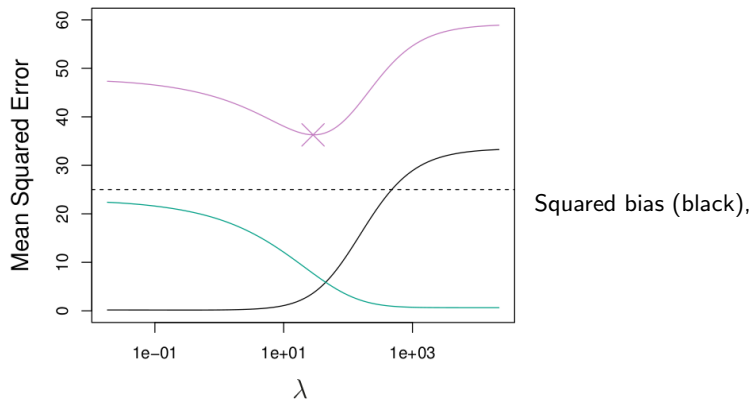
Least squares **is** scale equivariant.
Ridge/Lasso **are very much not**.

Solution: standardize predictors

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

- Ex: income variable.
- Least squares is scale equivariant
- Ridge regression very much is not
- $X_j \hat{\beta}_{j,\lambda}^R$ depends not only on λ but also on values of other predictors

Bias-Variance tradeoff



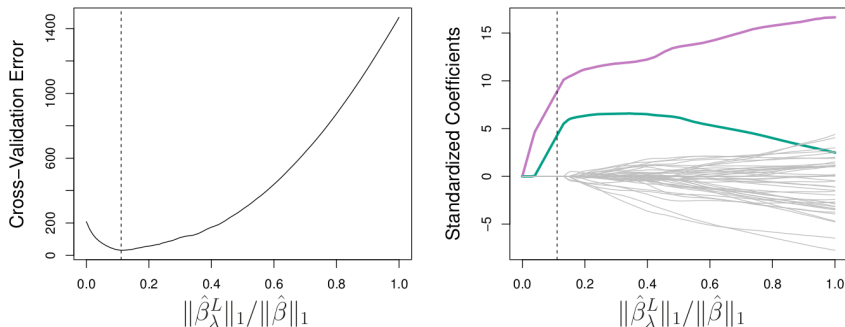
variance (green), and test mean squared error (purple) for simulated data.

Horizontal dashed line is minimum possible test MSE

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.

10-fold CV choice of λ for lasso and simulated data



- Data generated with $p = 45$ variables, $n = 50$ data points, only 2 variables actually involved
- The vertical dashed line: choice of λ
- Right: Two colored lines are predictors related to the response (signal), grey are unrelated (noise)
- Correctly gives much larger coefficient estimates to the two signal predictors, but also the min CV error corresponds to a set of coefficient estimates for which only the signal variables are non-zero.
- Note that least squares solution doesn't even emphasize one of the variables (green one)

Section 3

Optimization Formulation

Another formulation for Ridge Regression

Find β to minimize

RSS

Find β to minimize:

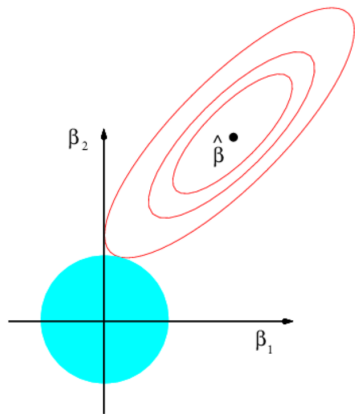
$$RSS + \lambda \sum_{j=1}^p \beta_j^2$$

subject to

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

- For every λ , there is an s to make the right side solution the same as the one on the left.
- Think of s as a budget for how large $\|\beta_j\|_2$ can be.
- Large enough s and least squares answer is available, so would just return least squares.
- Smaller s means we might not get quite the optimal RSS , but do well enough. See next slide.

Visualization using disks



- Red lines are the RSS levelsets
- Blue ball has radius s^2 .

Find β to minimize

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

subject to

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

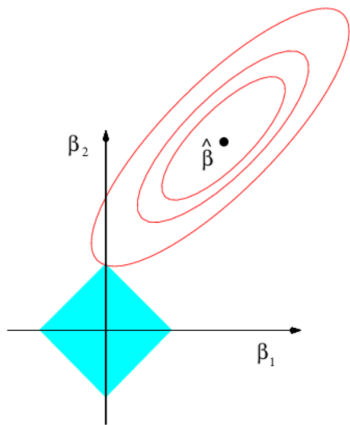
- If s is big enough, encloses the minimum of RSS function
- Smaller s means have to be in the blue ball

What about ℓ_1 ?

$$\|\beta\|_1 = \sum |\beta_i|$$

What does the set of points (β_1, β_2) for which $\|(\beta_1, \beta_2)\|_1 \leq s$ look like?

Same game for Lasso



- Diamond sides at $(0, \pm s)$, $(\pm s, 0)$

Find β to minimize

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

subject to

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

- If s is big enough, encloses the minimum of RSS function
- Smaller s means have to be in the blue ball
- Min likely to happen at corner, on axis, where something is 0

Same game for subset selection

Find β to minimize

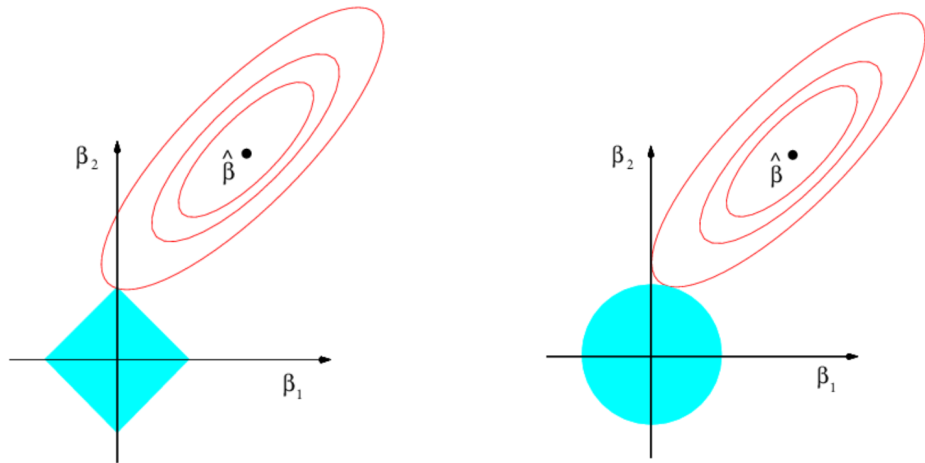
$$RSS = \sum_{i=1}^n \left(y_i - \beta_0 - \sum_{j=1}^p \beta_j x_{ij} \right)^2$$

subject to

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

- Similar format but computationally infeasible since still requires checking all subsets containing s predictors
- Interpretation: ridge regression and lasso are computationally feasible versions of subset selection

Using this visual to understand why lasso gets us zero values



Because ridge regression has circular constraint with no sharp points, intersection doesn't happen on axis line. But because of the diamond shape of ℓ_1 , much better chance of that.

Least Squares:

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

Ridge:

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

The Lasso:

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

Summary

Find β to minimize

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

subject to:

Least Squares:

No constraints

Ridge:

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

The Lasso:

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

Also, find best choice of λ or s using CV