

Variable Selection and Regularized Methods ¹

Shaobo Li

University of Kansas

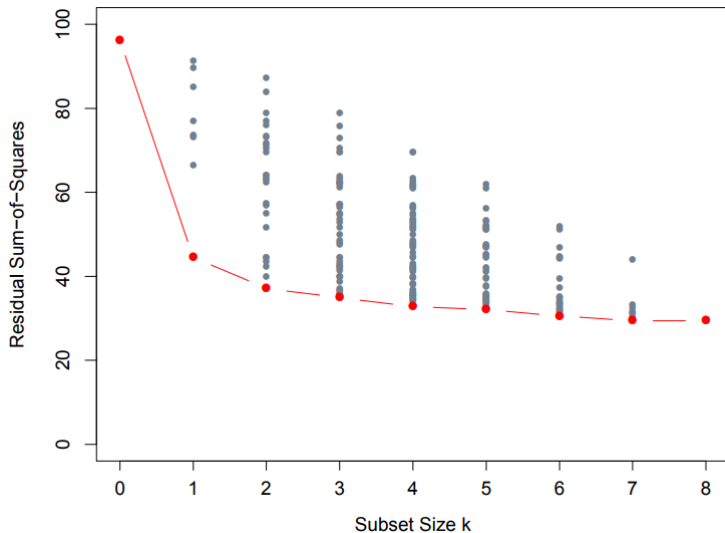
¹Partially based on Hastie, et al. (2009) ESL, and James, et al. (2013) ISLR

- Variable selection – excluding unnecessary variables
 - Interpretation and simplicity
 - Prediction stability and accuracy
 - Bias-variance tradeoff
- Common approaches
 - Subset selection
 - Shrinkage method (also called *regularization*)
 - Dimension reduction (project p predictors to an m -dimensional subspace)
- Some times it is subjective, and needs domain knowledge that certain variables must be in the model.

Best Subset Selection

- Select the best subset of predictors such that the model is optimal in terms of a certain assessment metric
- Computationally expensive even infeasible
 - *leaps and bounds* (an R package “leaps”) algorithm makes it feasible for p as large as 30 or 40.
- Suppose there are 10 predictors. How many models do we need to fit and evaluated?

Illustration of Best Subset Selection



Forward, Backward, and Stepwise Selection

- Computationally less expensive than best subset
- Iteratively adding or dropping one variable at a time
- Forward/backward is **greedy** procedure. That is, they won't adjust any added/dropped variables in previous step
- Stepwise: start with forward, and then iteratively add and drop variables
- Commonly used selection criteria: AIC, BIC
- R package: "step"
- An illustration: [click here](#)

Selection Criteria – AIC, BIC

- Akaike information criterion (AIC), the smaller the better

$$AIC = -2 \log(\hat{L}) + 2p$$

- Bayesian information criteria (BIC), the smaller the better

$$BIC = -2 \log(\hat{L}) + \log(n)p$$

where \hat{L} is estimated likelihood function

- For linear regression, $-2 \log(\hat{L})$ is equivalent to RSS
- BIC weighs more on p comparing to AIC. [What does this mean?](#)

Shrinkage Methods

- Also called penalized estimation, or regularization.
- Shrink the regression coefficients toward 0 by constraints (regularization)
- Estimates are usually *biased*
- A game of **bias-variance tradeoff**
- Shrinkage methods are generally preferred over subset methods. Why?
- We discuss two popular shrinkage methods:
 - Ridge regression
 - LASSO

- Least absolute shrinkage and selection operator (LASSO)
- Introduced by Tibshirani (1996)
- One of the most popular variable selection methods
- It estimates the coefficients and selects variables simultaneously.
- A tuning parameter λ controls the “power” of selection.
- Need to standardize all predictors in shrinkage estimation. [Why?](#)

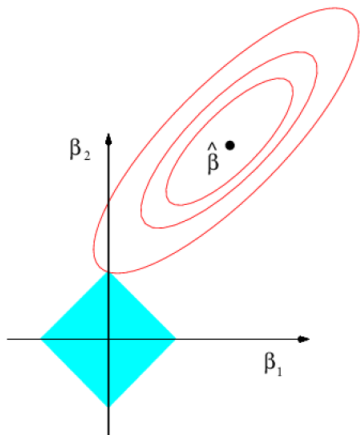
- LASSO solves the (L_1) *penalized least square*

$$\hat{\beta}_{LASSO} = \arg \min_{\beta} \sum_{i=1}^n \left(y_i - \beta_0 - \sum_{j=1}^p \beta_j x_{ij} \right)^2 + \lambda \sum_{j=1}^p |\beta_j|$$

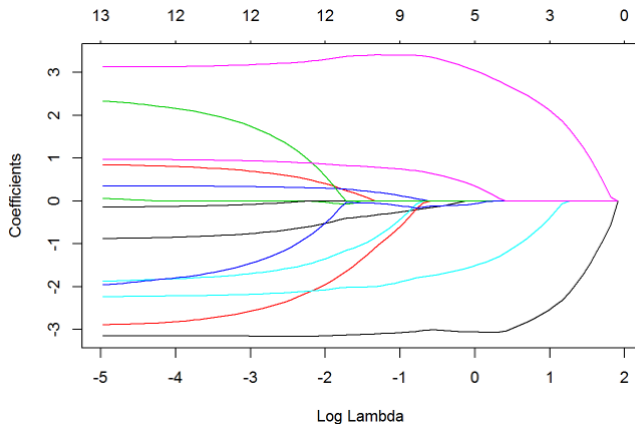
- It is a *convex optimization* problem
- It is equivalent to solve a constrained optimization problem

$$\begin{aligned} & \min \sum_{i=1}^n \left(y_i - \beta_0 - \sum_{j=1}^p \beta_j x_{ij} \right)^2 \\ & \text{s.t. } \sum_{j=1}^p |\beta_j| = a \end{aligned}$$

An Illustration

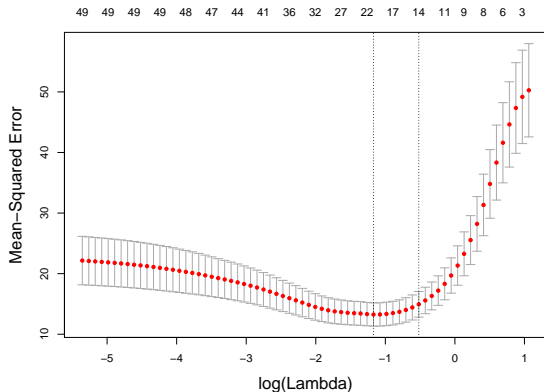


LASSO Regression Solution Path – Boston Housing Data



Tuning Parameter λ Selection

- λ controls the shrinkage level (different λ associates with different estimated model)
- Cross-validation
 - In R, use the function `cv.glm()` in package `glmnet`



High-Dimensional Regression

- Number of predictor is very large (even larger than sample size)
- Ultra-high dimension $p \gg n$
- It is very common for gene expression and image data
- Sparsity assumption: only a few predictors are relevant
- OLS fails when $n < p$. Why?
- LASSO or similar methods provide sparse solution

Elastic Net Regression

- Introduced by Zou and Hastie (2005)
- Combination of Ridge and LASSO

$$\hat{\beta}_{EN} = \arg \min_{\beta} \sum_{i=1}^n (y_i - \mathbf{x}_i^T \beta)^2 + \lambda_1 \sum_{j=1}^p |\beta_j| + \lambda_2 \sum_{j=1}^p \beta_j^2$$

- Convex optimization
- Ridge and LASSO are special cases of Elastic Net
- It incorporates the advantages of both Ridge and LASSO
 - Ridge regression: lower variance; multicollinearity
 - LASSO: variable selection (selects at most n variables if $p > n$)

Some Variants of LASSO

- There are many other type of penalized estimators with different penalty functions that can perform variable selection.
 - Group Lasso (Yuan and Lin, 2006)
 - Adaptive-LASSO (Zou, 2006)
 - SCAD (Fan and Li, 2001)
 - MCP (Zhang, 2010)

Ridge Regression

- Recall least square. We solve the optimization

$$\hat{\beta}_{LS} = \arg \min_{\beta} \sum_{i=1}^n \left(y_i - \beta_0 - \sum_{j=1}^p \beta_j x_{ij} \right)^2$$

- Ridge regression solves a (L_2) *penalized least square*

$$\hat{\beta}_{Ridge} = \arg \min_{\beta} \sum_{i=1}^n \left(y_i - \beta_0 - \sum_{j=1}^p \beta_j x_{ij} \right)^2 + \lambda \sum_{j=1}^p \beta_j^2$$

- λ is a tuning parameter, called shrinkage parameter
- Writing in matrix form, we can get the analytical solution

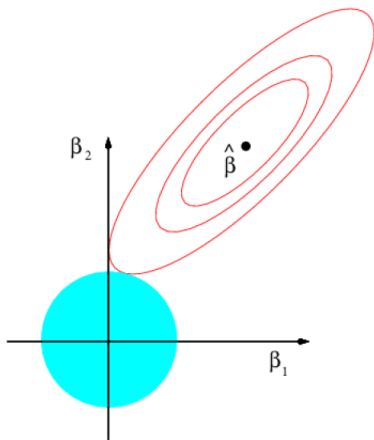
$$\hat{\beta}_{Ridge} = (\mathbf{X}^T \mathbf{X} + \lambda \mathbf{I})^{-1} \mathbf{X}^T \mathbf{y} \quad (\text{Exercise: Show it!})$$

- It is equivalent to solve a constrained optimization problem

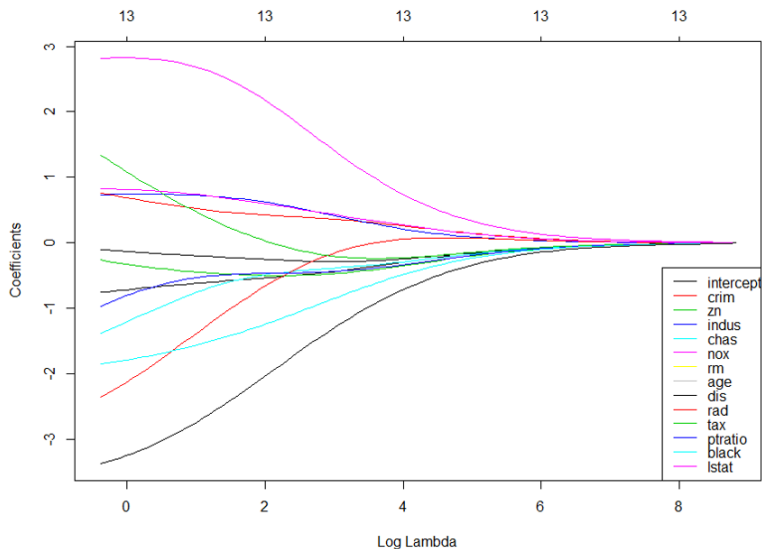
$$\begin{aligned} \min \sum_{i=1}^n \left(y_i - \beta_0 - \sum_{j=1}^p \beta_j x_{ij} \right)^2 \\ \text{s.t. } \sum_{j=1}^p \beta_j^2 = a \end{aligned}$$

- a corresponds to the tuning parameter λ

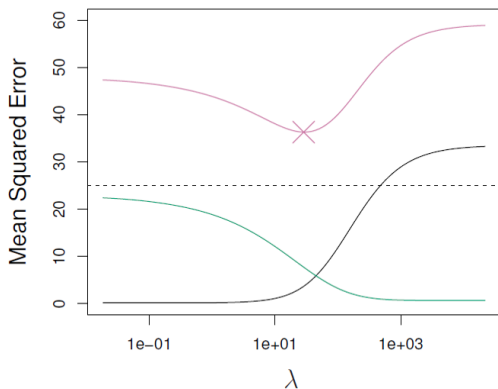
An Illustration



Ridge Regression Solution Path – Boston Housing Data



Bias-Variance Tradeoff



Simulated data with $n=50$ observations, $p=45$ predictors, all having nonzero coefficients. Squared bias (black), variance (green), and test mean squared error (purple) for the ridge regression predictions on a simulated data set.