

Exercises 5 Sparsity

Natalia Zuniga-Garcia

October 15, 2017

Broadly speaking, there are three ways to think about enforcing sparsity in a statistical model.¹

- The old-school frequentist way, based on classical hypothesis tests (chi-squared tests, F tests, likelihood ratio tests, etc), perhaps coupled with greedy algorithms for model selection.
- The Bayesian way, in which sparsity is baked into the prior distribution on some parameter²
- The new-school frequentist way, based on penalizing the likelihood function in a way that encourages sparsity, and computing a *maximum penalized likelihood* estimate.

These viewpoints all have their advantages and disadvantages. One the major advantages of the penalized-likelihood view of sparsity (especially versus the Bayesian view) is scalability to large data sets. That makes it particularly suited to big-data applications, which is why we'll pursue it here, and for much of the rest of the course.

1 Penalized likelihood and soft thresholding

(A) Define the function

$$S_\lambda(y) = \arg \min_{\theta} \frac{1}{2}(y - \theta)^2 + \lambda|\theta|. \quad (1)$$

The intuition here is that θ is a parameter of a statistical model, and y is data. The first (quadratic) term rewards good fit to the data, while the second term rewards θ for being “simpler” (i.e. closer to zero). $S_\lambda(y)$ returns an estimate for θ that blends these two goals.

First show (in a trivial one- or two-liner) that the quadratic term in the objective above is the negative log likelihood of a Gaussian distribution with mean θ and variance 1.

Solution

The Gaussian PDF is:

$$f(y, \mu, \sigma) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp \frac{(y-\mu)^2}{2\sigma^2}$$

A Gaussian distribution with mean θ and variance 1 is:

$$f(y, \theta, 1) = \frac{1}{\sqrt{2\pi}} \exp \frac{(y-\theta)^2}{2}$$

¹In this context, sparsity means that some of the parameters are zero. This is different from the notion of sparsity considered previously, where some of the *features* in a regression problem were zero.

²E.g. http://www-stat.wharton.upenn.edu/~edgeorge/Research_papers/fastN96.pdf

The likelihood for a single y is:

$$L(\theta|y, 1) = \frac{1}{\sqrt{2\pi}} \exp \frac{(y-\theta)^2}{2}$$

And the negative log-likelihood is:

$$l(\theta|y, 1) = -\left(-\frac{1}{2}\log(2\pi) - \frac{1}{2}(y-\theta)^2 \right) = C + \frac{1}{2}(y-\theta)^2$$

Where, $C = \frac{1}{2}\log(2\pi)$ is a constant. We can observe that the negative log-likelihood contains the quadratic form in Equation (1).

Then prove that

$$S_\lambda(y) = \text{sign}(y) \cdot (|y| - \lambda)_+,$$

where $a_+ = \max(a, 0)$ is the positive part of a . This is a basic exercise in sub-differential calculus³, if you happen to know this subject. But I am not assuming that you do, and you can also prove the statement using ordinary differential calculus if you split it up into cases ($\theta > 0$, $\theta < 0$, and $\theta = 0$).⁴

Solution

First, we try to minimize the objective function by calculation its derivative with respect to θ and equalizing it to zero as follow,

$$\begin{aligned} \frac{\partial}{\partial \theta} S_\lambda(y) = 0 &\Rightarrow \frac{\partial}{\partial \theta} \left(\frac{1}{2}(y-\theta)^2 + \lambda|\theta| \right) = -(y-\theta) + \lambda \text{sign}(\theta) = 0 \\ &\Rightarrow \theta = y - \lambda \text{sign}(\theta) \end{aligned}$$

Now, we split this into three cases:

- Case 1. $\theta > 0$, $\theta = y - \lambda \Rightarrow$ Then we have that $y > \lambda$
- Case 2. $\theta < 0$, $\theta = y + \lambda \Rightarrow$ Then we have that $y < -\lambda$
- Case 3. $\theta = 0$, $y = \lambda \Rightarrow$ Then we have that $|y| < \lambda$

Which can be written as:

$$S_\lambda(y) = \text{sign}(y) \cdot (|y| - \lambda)_+$$

³<http://ocw.mit.edu/courses/electrical-engineering-and-computer-science/6-253-convex-analysis-and-optimization-spring-2012/lecture-notes/MIT6.253S12 lec12.pdf>

⁴If you know about constrained optimization, you could also introduce a slack variable z and write the objective as

$$\frac{1}{2}(y-\theta)^2 + \lambda|z|,$$

and then do the minimization subject to the constraint that $\theta = z$.

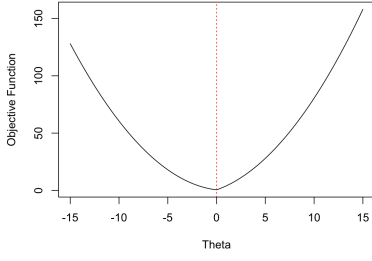
To build intuition, you should try plotting the objective in Equation 1 for various y and λ , and verifying that $S_\lambda(y)$ indeed obtains the minimum.⁵

Solution

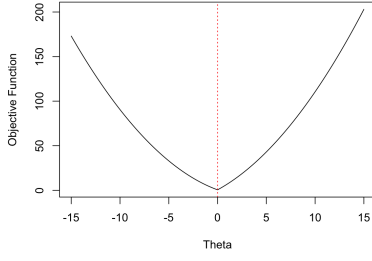
The solution is presented in the file *Part_1A.R*,

```
1 # Plot the Objective function with different y and lambda values
2 y = c(-1, 0, 5)
3 lambda = c(2, 5, 10)
4
5 for (i in 1:length(y)){
6   for (j in 1:length(lambda)){
7     # Objective Function
8     png(filename=paste('P1A', y[i], lambda[j], '.png', sep = ""),
9         width = 15, height = 12, units = "cm", res = 200)
10    values = list(y = y[i], lambda=lambda[j])
11    curve(0.5 * (values$y - x)^2 + values$lambda * abs(x), from=-15,
12        to=15, xname="x", ylab="Objective_Function", xlab="Theta")
13    # S is the Soft Theresholding
14    S <- abline(v = sign(y[i]) * max(abs(y[i]) - lambda[j], 0),
15              col=2, lty=3)
16    dev.off()
17   }
18 }
```

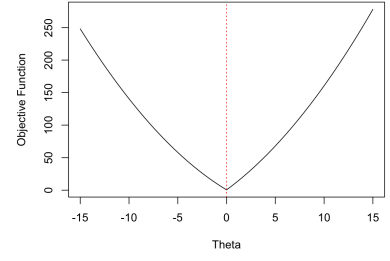
⁵Recall the “curve” function in R, which acts like a graphing calculator.



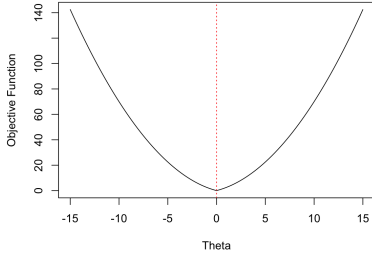
(a) $y = -1$ and $\lambda = 2$



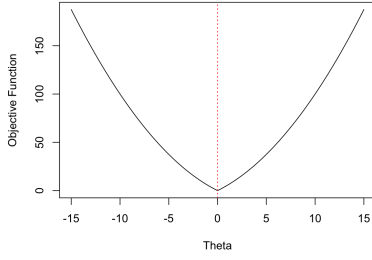
(b) $y = -1$ and $\lambda = 5$



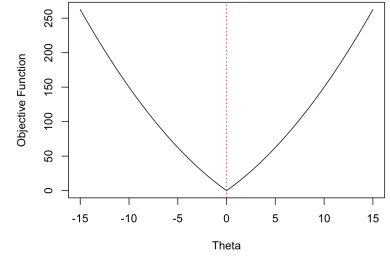
(c) $y = -1$ and $\lambda = 10$



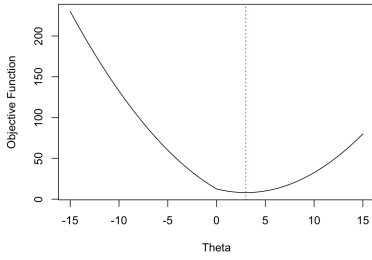
(d) $y = 0$ and $\lambda = 2$



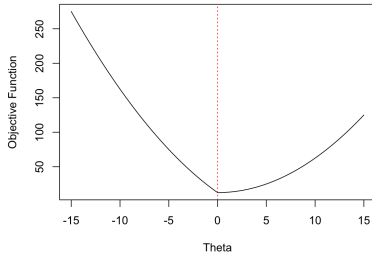
(e) $y = 0$ and $\lambda = 5$



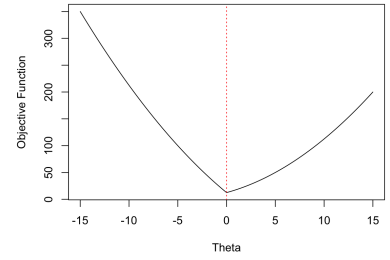
(f) $y = 0$ and $\lambda = 10$



(g) $y = 5$ and $\lambda = 2$



(h) $y = 5$ and $\lambda = 5$



(i) $y = 5$ and $\lambda = 10$

Figure 1: Objective Function and $S_\lambda(y)$ with different values of y and λ

$S_\lambda(y)$ is called the *soft thresholding* function with parameter λ . Plot this as function of y for a few different parameters of λ . You'll see how it encourages sparsity in a "soft" way, especially if you compare it to the hard-thresholding function

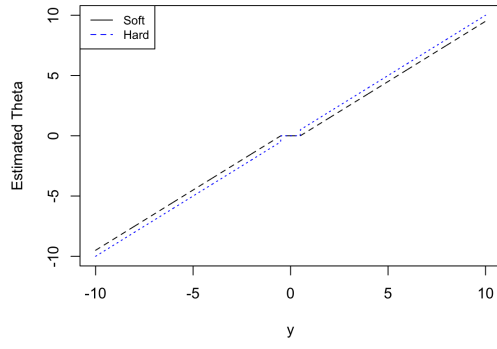
$$H_\lambda(y) = \begin{cases} y & \text{if } y \geq \lambda \\ 0 & \text{otherwise.} \end{cases}$$

Solution

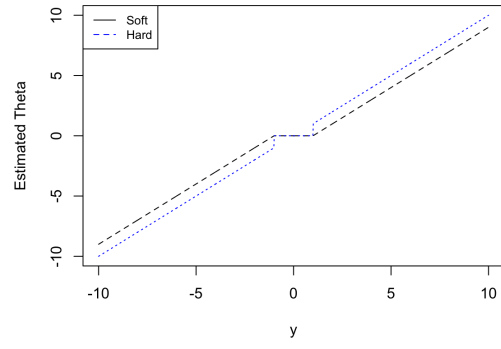
The solution is presented in the file *Part_1A.R*,

```

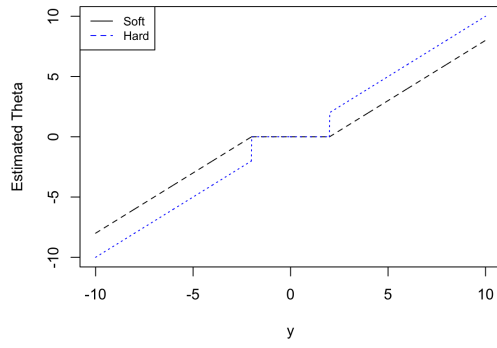
1 # Plot the soft-thresholding and hard-thresholding functions
2 # with different lambda
3 y <- seq(-10,10,length=1000)
4 lambda = c(0.5, 1, 2, 3, 5, 6)
5 for (j in 1:length(lambda)){
6   S <- array(NA, dim = length(y))
7   H <- array(NA, dim = length(y))
8   for (i in 1:length(y)){
9     # Soft-thresholding
10    S[i] <- sign(y[i]) * max(abs(y[i]) - lambda[j], 0)
11    H[i] <- ifelse(abs(y[i])>=lambda[j], y[i], 0)
12  }
13  png(filename=paste('P1AA', lambda[j], '.png', sep = ""),
14       width = 15, height = 12, units = "cm", res = 200)
15  plot(y, drop(S), type = "l", xlab="y", ylab="Estimated_Theta",
16       col="black", ylim=c(-10,10), lty=2)
17  points(y, drop(H),type="l",col="blue",lty=3)
18  legend('topleft',legend = c('Soft','Hard'), col=c("black", "blue"),
19        lty=1:2, cex=0.8)
20  dev.off()
21 }
```



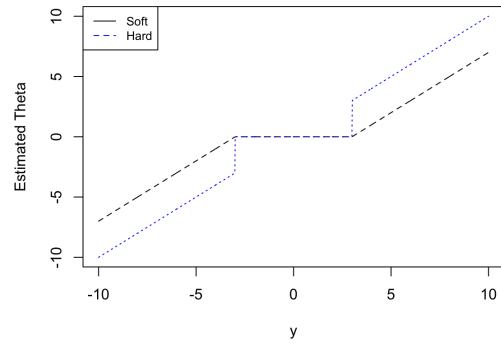
(a) $\lambda = 0.5$



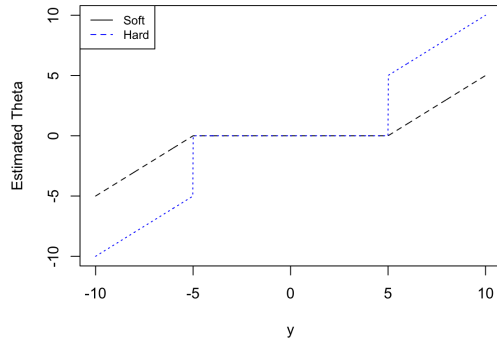
(b) $\lambda = 1$



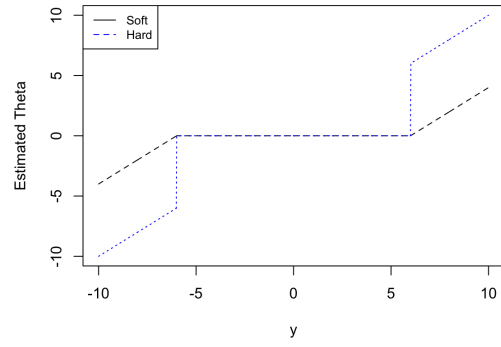
(c) $\lambda = 2$



(d) $\lambda = 3$



(e) $\lambda = 5$



(f) $\lambda = 6$

Figure 2: Soft-thresholding and hard-thresholding functions with different λ values

- (B) Here's a simple toy example to illustrate how soft thresholding can be used to enforce sparsity in statistical models.

Suppose we observe data from the following statistical model:

$$(z_i \mid \theta_i) \sim N(\theta_i, \sigma_i^2).$$

That is, there are n different means θ_i , and we observe 1 normally distributed observation for each one. We allow that each observation has a different variance σ_i^2 ; for now we'll assume these are known. This is called the Gaussian sequence model, or the normal-means problem. It looks like a toy problem, but is surprisingly useful ⁶ in a wide variety of applications, from curve fitting to image denoising to genome-wise association studies.⁷

Now suppose we believe that a lot of the θ_i 's are zero—i.e. that the vector $\theta = (\theta_1, \dots, \theta_n)^T$ is sparse. Consider an estimator for each θ_i of the form

$$\hat{\theta}(y_i) = S_{\lambda\sigma_i^2}(y_i),$$

where S is the soft thresholding operator defined above with parameter $\lambda\sigma_i$. (Side question: why $\lambda\sigma_i^2$ for the soft-thresholding parameter?)

Solution

Because for the objective function we are scaling by the standard deviation:

$$\frac{1}{2} \left(\frac{y - \theta}{\sigma} \right)^2 + \lambda |\theta| \Rightarrow \frac{1}{2} (y - \theta)^2 + \lambda \sigma^2 |\theta|$$

Try the following intuition-building exercise.

1. Choose some sparse vector θ and the corresponding σ_i^2 's (it's OK for them all to be equal). You have freedom in deciding what the nonzero elements of θ should be, and how sparse it should be. (Ideally there would be some tunable parameter that you could use to ratchet up or down the sparsity level.)

Solution

The solution is presented in the file *Part_1B.R*,

```

1 # Create sparse vector theta and sigma vector
2 theta_estimate <- function(N, K){
3   # Create sparse vector theta iid in N(0, 10)
4   # Input:
5   # N = vector size (constant)
6   # K = sparsity level (constant)
7   # Output:
8   # theta = sparse vector (vector N)
9   theta <- rnorm(N, 0, 10)
10  mask <- rbinom(N, 1, K)
11  theta <- theta * mask
12 }
13 sigma <- rep(1, N) # Assumed equal to 1

```

⁶<http://statweb.stanford.edu/~imj/GE06-11-13.pdf>

⁷The simplest example is probably curve-fitting using an orthogonal basis of functions, like Fourier polynomials or wavelets; see the link provided.

2. Simulate one data point $(z_i | \theta_i) \sim N(\theta_i, \sigma_i^2)$ for each θ_i .

```

1 # Simulate a data point for each theta
2 y_estimate <- function(theta, sigma){
3   # Simulate vector of observations
4   # Inputs:
5   # theta = sparse vector (vector N)
6   # sigma = standard deviation (vector N)
7   # Output:
8   # y = vector of observations (vector N)
9   rnorm(length(theta), theta, sigma)
10 }

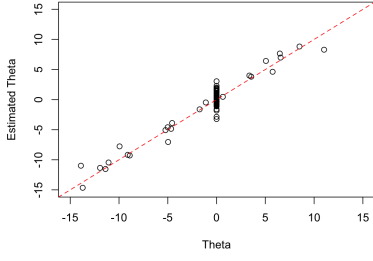
```

3. Compute $\hat{\theta}(y_i) = S_{\lambda\sigma_i^2}(y_i)$ across a discrete grid of different λ values. Plot $\hat{\theta}(y_i)$ versus θ_i , and observe how the soft-thresholding function both *selects* certain θ_i 's by sparsifying the estimate, as well as *shrinks* the nonzero estimates $\hat{\theta}(y_i)$ towards 0 (and towards each other).

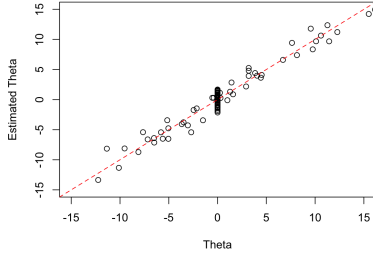
```

1 # Create a function to obtain S lambda (estimated theta)
2 S_lambda <- function(y, lambda){
3   # Estimate theta using the soft-thresholding operator
4   # Input:
5   # y = vector of observations (vector N)
6   # lambda = complexity parameter (constant)
7   # Output:
8   # Theta = estimated theta vector (vector N)
9   sign(y) * pmax(abs(y) - lambda, rep(0, N))
10 }
11 # -----
12 # Estimate theta using S lambda with different lambdas
13 # and sparsity levels
14 N <- 100
15 lambda <- c(0, 1, 2, 5)
16 K <- c(.3, .5, .9)
17 for (j in 1:length(K)){
18   theta <- theta_estimate(N, K[j])
19   y <- y_estimate(theta, sigma)
20   for (i in 1:length(lambda)){
21     theta.hat <- S_lambda(y, lambda[i] * sigma ^ 2)
22     png(filename=paste('P1B_K_', K[j], '_l_', lambda[i], '.png', sep = ""),
23         width = 15, height = 12, units = "cm", res = 200)
24     plot(theta, theta.hat, xlab="Theta", ylab="Estimated Theta",
25         col="black", ylim=c(-15,15), xlim=c(-15,15))
26     abline(a = 0, b = 1, col='red', lty=2)
27     dev.off()
28   }

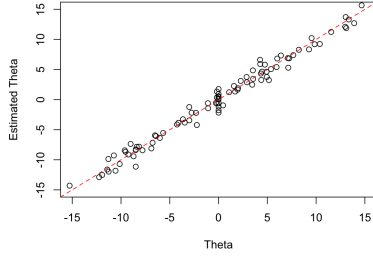
```

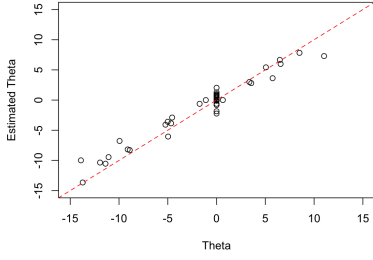
(a) $\lambda = 0$ and $K = 0.3$



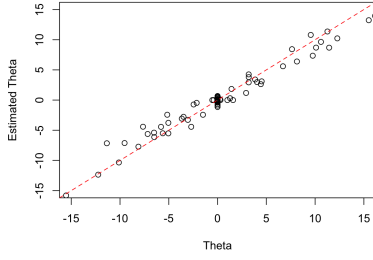
(b) $\lambda = 0$ and $K = 0.5$



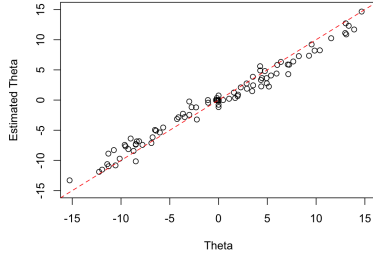
(c) $\lambda = 0$ and $K = 0.9$



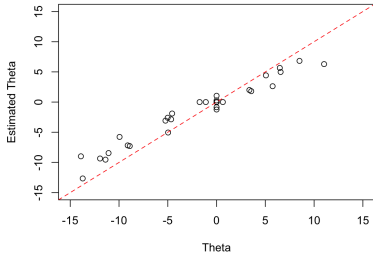
(d) $\lambda = 1$ and $K = 0.3$



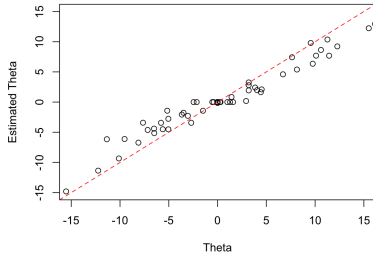
(e) $\lambda = 1$ and $K = 0.5$



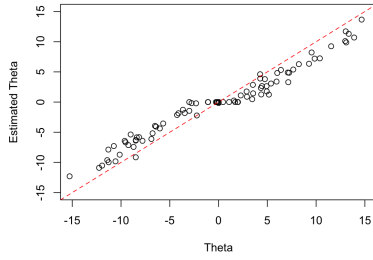
(f) $\lambda = 1$ and $K = 0.9$



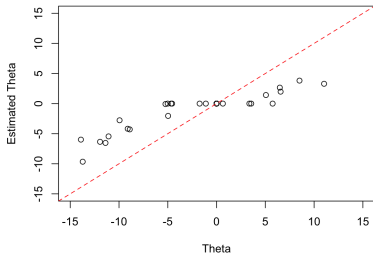
(g) $\lambda = 2$ and $K = 0.3$



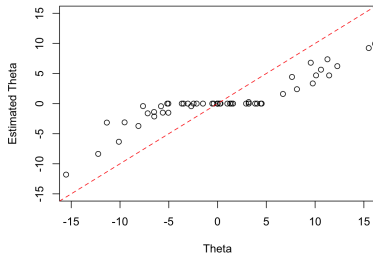
(h) $\lambda = 2$ and $K = 0.5$



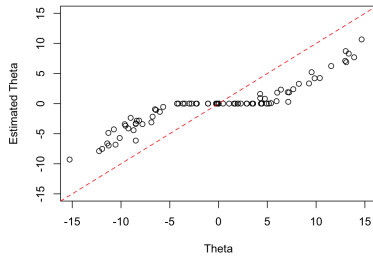
(i) $\lambda = 2$ and $K = 0.9$



(j) $\lambda = 5$ and $K = 0.3$



(k) $\lambda = 5$ and $K = 0.5$



(l) $\lambda = 5$ and $K = 0.9$

Figure 3: $\hat{\theta}$ and θ_i with different λ and sparsity (K) values

4. Plot the mean-squared error of your estimate as a function of λ :

$$\text{MSE}(\lambda) = \frac{1}{n} \sum_{i=1}^n \left\{ \hat{\theta}(y_i) - \theta_i \right\}^2.$$

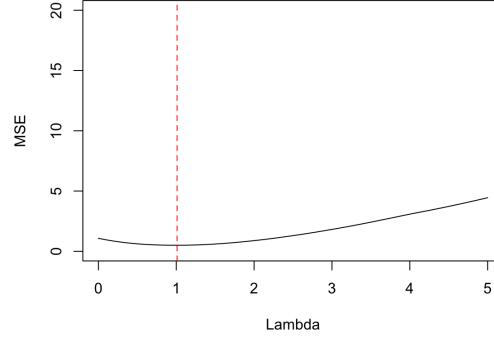
Make sure that the MSE actually obtains a minimum across the grid of λ values you have chosen. Try this for several different configurations of θ .

Note: like I said, this is an intuition-building exercise. You could never choose λ this way for a real problem, because it requires knowledge of the truth.

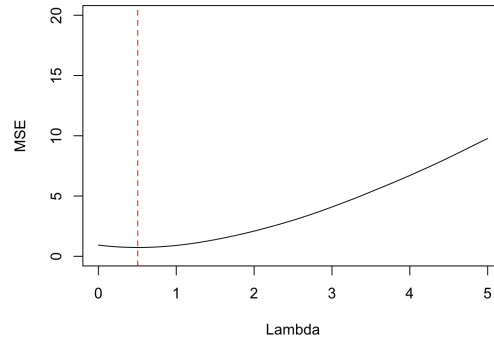
```
1 # Plot the MSE of lambda
2 K <- c(.3, .5, .9)
3 lambda <- seq(0,5,length=100) # More lambdas
4 for (j in 1:length(K)){
5   theta <- theta_estimate(N, K[j])
6   y <- y_estimate(theta, sigma)
7   MSE <- array(NA, dim = length(lambda))
8   for (i in 1:length(lambda)){
9     theta.hat <- S_lambda(y, lambda[i] * sigma ^ 2)
10    MSE[i] <- mean((theta.hat - theta)^2)
11  }
12  png(filename=paste('P1BB_K_',K[j],'.png', sep = ""),
13       width = 15, height = 12, units = "cm", res = 200)
14  plot(lambda, MSE,xlab="Lambda", ylab="MSE", type = "l",
15       col="black", ylim=c(0,20), xlim=c(0,5))
16  lambda.min <- lambda[which.min(MSE)]
17  print(lambda.min)
18  abline(v = lambda.min, col='red', lty=2)
19  dev.off()
20 }
```

How does the optimal λ change as the sparsity in θ changes?⁸

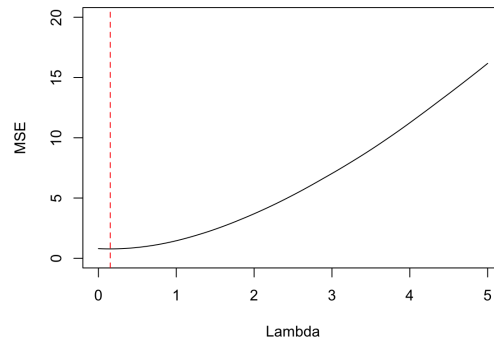
We can observe from Figure 4 that as we decrease the sparsity (greater K), we obtain lower λ values. This means that there is not need to force sparsity.



(a) $K = 0.3$, and Minimum $\lambda = 1.0101$



(b) $K = 0.5$, and Minimum $\lambda = 0.5050$



(c) $K = 0.9$, and Minimum $\lambda = 0.1515$

Figure 4: Mean-squared error (MSE) for different λ and sparsity (K) values

⁸Here it's enough to assess the optimal λ using the good old-fashioned “plot and point” strategy for optimization. That is, you plot the function, point at the minimum, and say proudly, “Here’s the minimum.”

2 The lasso

Although a soft-thresholding approach to the normal-means problem is actually very useful in practice, this utility is not immediately apparent. On the other hand, a generalization of this idea to regression, called the lasso,⁹ both looks and is immediately useful.

Consider the standard linear regression model

$$y = X\beta + e,$$

where y is an n -vector of responses, X is an $n \times p$ features matrix whose i th row x_i is the vector of features for observation i , and e is a vector of errors/residuals.

The lasso involves estimating β as the solution to the penalized least-squares problem¹⁰

$$\hat{\beta} = \arg \min_{\beta} \frac{1}{2} \|y - X\beta\|_2^2 + \lambda \|\beta\|_1,$$

where $\|\beta\|_1$ is the ℓ_1 (pronounced “ell one”) norm of the coefficient vector:

$$\|\beta\|_1 = \sum_{j=1}^p |\beta_j|.$$

As you can see, the penalty function is just like the absolute-value penalty you used on the normal-means problem, generalized to the vector case. And just as in the normal means problem, this penalty function will have the effect of both selecting a set of nonzero β_i ’s (i.e. sparsifying the estimate) as well as shrinking the nonzero β_i ’s toward 0. The bigger λ , the more aggressive the shrinkage effect.

Note: we typically leave the intercept in a lasso fit unpenalized. We can accomplish this by explicitly introducing an intercept, e.g. writing the objective as

$$\frac{1}{2n} \|y - (\alpha \mathbf{1} + X\beta)\|_2^2 + \lambda \|\beta\|_1,$$

where α is a scalar intercept and $\mathbf{1}$ is a vector of all 1’s. Or we can leave the problem in its original form above, and assume that both the response variable y and all columns of the predictor matrix have been standardized have a mean of 0 (in which case there is no need for an explicit intercept). For the rest of these exercises, we’ll assume that the variables have been standardized in this way.

To read more about lasso regression, consult Chapter 3.4.2 of *The Elements of Statistical Learning*, or the original paper¹¹ by Robert Tibshirani.

⁹“Lasso” stands for “least absolute shrinkage and selection operator. It was proposed in a <http://statweb.stanford.edu/~tibs/lasso/lasso.pdf> classic paper by Robert Tibshirani.

¹⁰Note: some will write the “fit” part of the objective function as

$$\frac{1}{2n} \|y - X\beta\|_2^2,$$

where n is the sample size. Translating between these two problems involves multiplying λ by a factor of n .

¹¹<http://statweb.stanford.edu/~tibs/lasso/lasso.pdf>

- (A) For now, we won't worry about *how* the lasso model is fit. Instead, we'll use pre-existing software to fit it:

In R: the package `glmnet`, described here ¹².

In Python: `sklearn.linear_model.Lasso` in scikit-learn ¹³.

That way, we can build intuition for its properties, as well as come to an appreciation for some of the statistical questions that arise in fitting the model. If you want to learn about how these functions fit the lasso model, read about coordinate descent in Chapter 3.8.6 of *The Elements of Statistical Learning*, or in this paper ¹⁴.

Download the data on diabetes progression in 442 adults from the Data folder on the class website. There are two files here.

diabetesY.csv: the response variable for each patient. This is the result of a blood test that provides a quantitative measure of disease progression one year after baseline (e.g. at diagnosis).

diabetesX.csv: 10 baseline patient variables, age, sex, BMI, cholesterol measurements, etc. Also here are all 10 quadratic terms of the form x_{ij}^2 and all 45 possible pairwise interactions of the form $x_{ij} \cdot x_{ik}$. This leads to 65 total variables: 10 linear main effects and 10 quadratic main effects from the baseline variables, and 45 interactions ($45 = 10 \text{ choose } 2$). The 10 baseline variables are standardized to have zero mean and unit Euclidean norm (i.e. the sum of squared entries in the first 10 columns is 1).

Fit the lasso model across a range of λ values (which `glmnet` does automatically) and plot the solution path $\hat{\beta}_\lambda$ as a function of λ , just like Figure 3.10 in *Elements*.¹⁵ (Note: your horizontal axis can just be λ , or $\log \lambda$.)

In addition, you should track the in-sample mean-squared prediction error of the fit across the solution path:

$$\text{MSE}(\hat{\beta}_\lambda) = \frac{1}{n} \sum_{i=1}^n (y_i - x_i^T \hat{\beta}_\lambda)^2 = \frac{1}{n} \|y - X\hat{\beta}_\lambda\|_2^2.$$

¹²https://web.stanford.edu/~hastie/glmnet/glmnet_alpha.html

¹³<http://scikit-learn.org/stable/index.html>

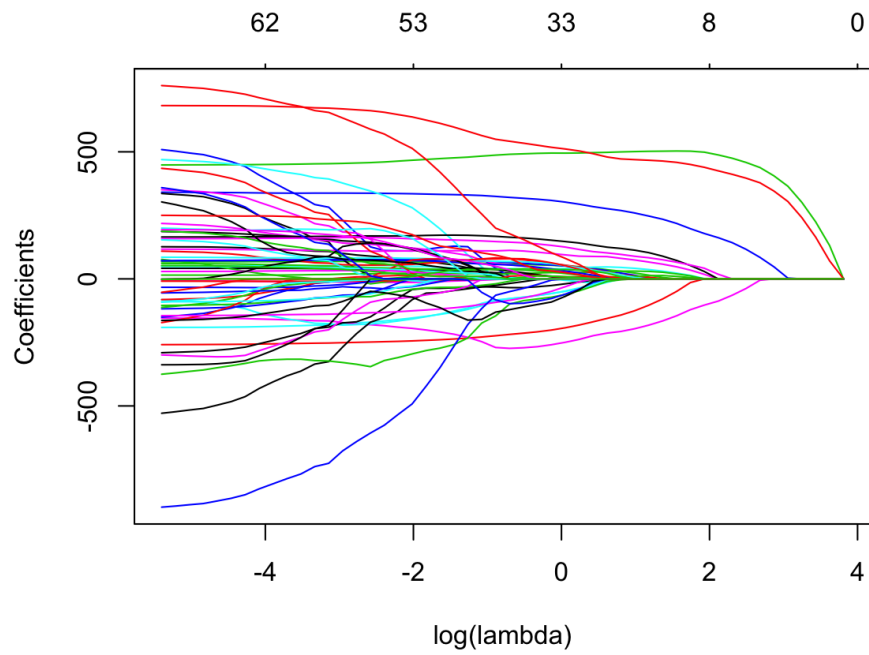
¹⁴<http://arxiv.org/pdf/0708.1485.pdf>

¹⁵For Python users, there is a “warm start” option to the lasso fitter, which can use the last solution as an initial guess for the next call. This is great for fitting a solution path.

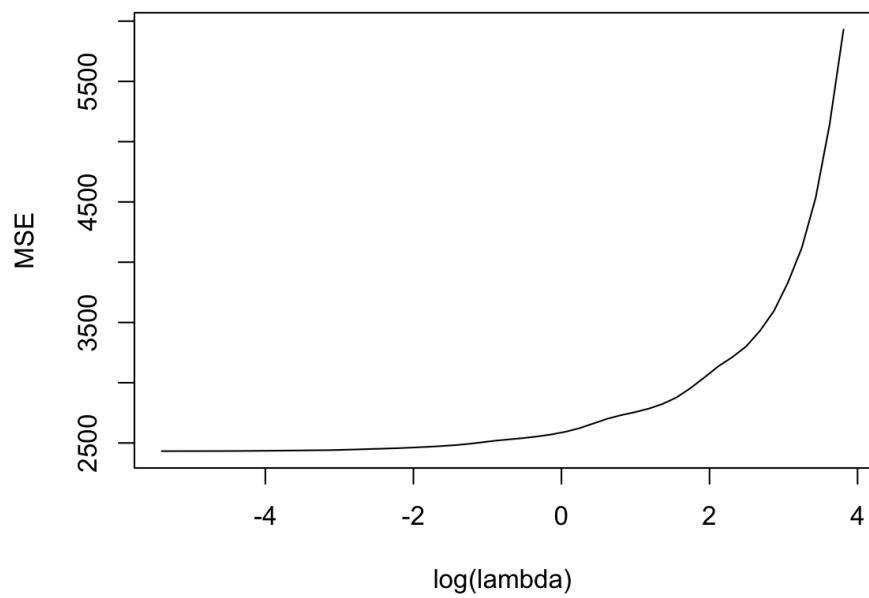
Solution

The solution is presented in the file *Part_2.R*,

```
1 # -----
2 # Load Diabetes Data
3 X = read.csv('diabetesX.csv')
4 X = as.matrix(X)
5 Y = read.csv('diabetesY.csv',header=F)
6 Y = Y[,1]
7 # -----
8 # Fit Lasso Regression Using glmnet and plot the solutions path
9 N <- 50 # Number of lambdas
10 fit <- glmnet(X, Y, alpha = 1, nlambda = N)
11 png(filename=paste('P2A_Lasso','.png', sep = ""), width = 15,
12     height = 12, units = "cm", res = 200)
13 plot(fit, xvar = 'lambda', xlab = 'log(lambda)')
14 dev.off()
15 # -----
16 # In-sample mean square prediction error (MSE)
17 Y.prediction <- predict.glmnet(fit, newx = X)
18 MSE <- array(NA, dim = N)
19 for (i in 1:N){
20   MSE[i] <- mean((Y.prediction[,i] - Y)^2)
21 }
22 png(filename=paste('P2A_MSE1','.png', sep = ""), width = 15,
23     height = 12, units = "cm", res = 200)
24 plot(log(fit$lambda), MSE, type = "l", xlab = 'log(lambda)')
25 dev.off()
```



(a) Profiles of lasso coefficients



(b) In-sample mean-squared prediction error of the fit

Figure 5: Lasso Regression

- (B) A natural way to choose λ is to minimize the expected out-of-sample prediction error. Suppose that (x_*, y_*) is a future data point from the same population/data-generating process as the original data. The goal would be to make the expected error

$$\text{MOOSE}(\hat{\beta}_\lambda) = E \{ (y_* - \hat{y}_*)^2 \} = E \{ (y_* - x_*^T \hat{\beta}_\lambda)^2 \}$$

as small as possible. Here the expected value is taken under what probability distribution generates (x, y) pairs, and “MOOSE” stands for mean out-of-sample squared error.¹⁶ Of course, we don’t have any “future data” lying around, so we have to estimate this quantity using the data we have. The in-sample mean-squared error, $\text{MSE}(\hat{\beta}_\lambda)$, is generally an optimistic estimate of this quantity: out-of-sample error tends to be worse, on average, than in-sample error, and we need some way of quantifying how much worse.

One way is using cross validation. If you’re unfamiliar with this concept, read about it in Chapter 7.10 of *Elements* (you may need to track back to the beginning of chapter 7 to pick up their notation). The essential idea is to split your data set into “training” and “testing” sets. Then you fit the model on the training set, and compute the average out-of-sample prediction error on the test set. This provides an estimate of the MOOSE. (You actually average this estimate over multiple such train/test splits, to reduce the influence of randomness in the split itself.)

I leave it to you to figure out the details of cross validation, including how big your training and testing sets should be. Again, Chapter 7.10 *Elements* has good guidelines here. Note: I expect you to write your own code for doing the train/test splits and computing out-of-sample error, wrapped around the base function for fitting the lasso regression model. If you’re in R, you can use the `cv.glmnet` function as a sanity check on, but not a replacement for, your results.

¹⁶MOOSE is not a standard acronym, but shouldn’t it be?!

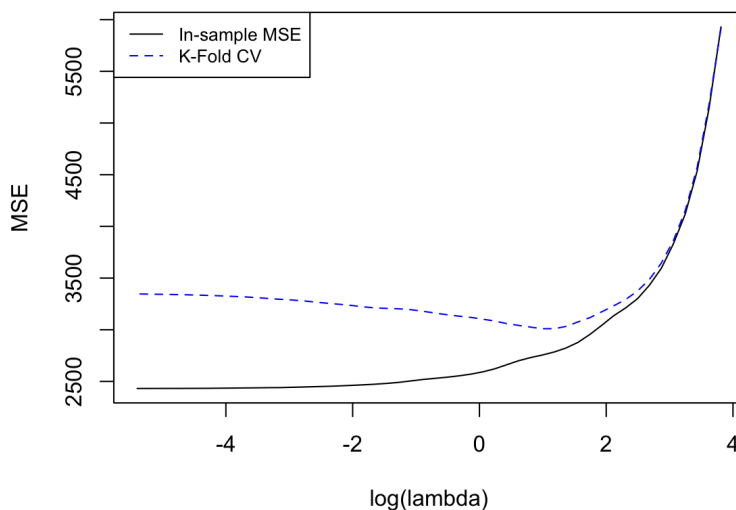
Solution

The solution is presented in the file *Part_2.R*,

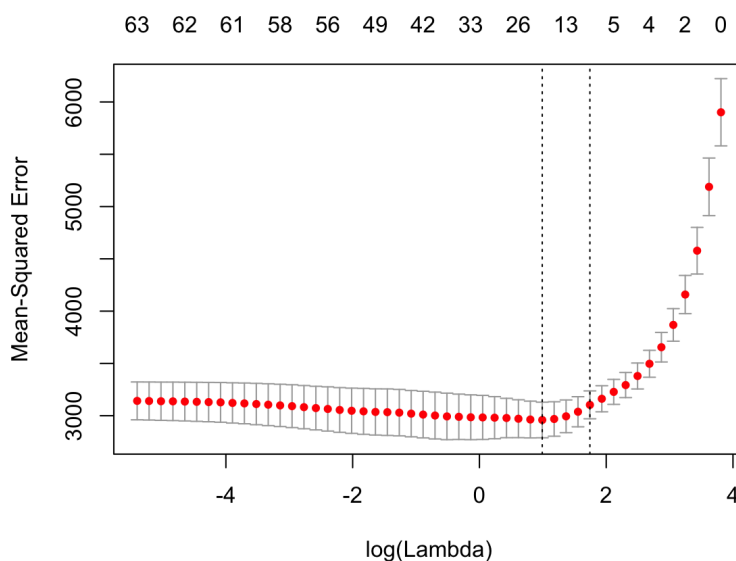
```
1 # K-fold Cross-Validation Section 7.10 of the textbook "Elements..."
2 K_fold <- function(X, Y, K, lambda){
3   # Compute the K-fold Cross-Validation
4   # Input:
5   # X = features (matrix M x P)
6   # y = vector of observations (vector M)
7   # K = number of folds (constant, typically 5 or 10)
8   # lambda = set of complexity parameters (vector N)
9   # Output:
10  # MSE = Mean Square Error (vector N)
11
12  # We split the data into K roughly equal-sized parts
13  # using an indexing function
14  partition <- cut(1:dim(X)[1], breaks = K, labels=FALSE)
15  # Randomization
16  partition <- sample(partition)
17
18  N <- length(lambda)
19  MSE <- array(NA, dim = c(K, N))
20  for (i in 1:K){
21    validation <- partition == i # Validation Kth part
22    X.validation <- X[validation,]
23    Y.validation <- Y[validation]
24    train <- !validation # Train part
25    X.train <- X[train,]
26    Y.train <- Y[train]
27    fit <- glmnet(X.train, Y.train, alpha = 1, lambda = lambda)
28    Y.prediction <- predict.glmnet(fit, newx = X.validation)
29    for (j in 1:N){
30      MSE[i, j] <- mean(((Y.prediction[,j] - Y.validation)^2))
31    }
32  }
33  colMeans(MSE)
34 }
35 # -----
36 # Applying the cross validation using K-fold function
37 K <- 5
38 MSE_Kfold <- K_fold(X, Y, K, fit$lambda)
39 png(filename=paste('P2A_MSE2', '.png', sep = ""), width = 15,
40      height = 12, units = "cm", res = 200)
41 plot(log(fit$lambda), MSE, type = "l", xlab = 'log(lambda)')
42 points(log(fit$lambda), MSE_Kfold, type="l", col="blue", lty=2)
43 legend('topleft', legend = c('In-sample MSE', 'K-Fold CV'),
44        col=c("black", "blue"), lty=1:2, cex=0.8)
45 dev.off()
46 # -----
47 # Sanity Check for K-fold CV
48 MSE_check <- cv.glmnet(X, Y, lambda = fit$lambda, nfolds = K)
49 png(filename=paste('P2A_MSE3', '.png', sep = ""), width = 15,
50      height = 12, units = "cm", res = 200)
51 plot(MSE_check)
52 dev.off()
```

Plot your cross-validated estimate of $\text{MOOSE}(\hat{\beta}_\lambda)$ across the solution path, as a function of λ . How does it compare with the *in-sample* mean-squared error from (A)?

We can observe the behavior similar to Figure 7.1 of the textbook (page 220). The in-sample MSE decreases as λ decreases (over-fitting?), while the MOOSE presents a MSE minimum near $\log(\lambda) = 1.367$ or $\lambda = 3.9222$.



(a) K-Fold Cross-Validation MOOSE compared to in-sample MSE



(b) K-fold cross-validation MOOSE obtained using CV.glmnet (sanity check)

Figure 6: Cross validation

- (C) Cross validation is one particularly simple way, based on the idea of resampling your data, to estimate the generalization error of a model. (Its reliance on resampling means that it shares a lot in common with the bootstrap¹⁷ as a way to estimate the standard error of a model parameter.)

However, cross-validation isn't the only way to estimate generalization error. Another such way is called the C_p statistic, proposed by Collin Mallows (and therefore often called *Mallows' C_p*). For a linear regression model, the C_p statistic is defined as

$$C_p(\hat{\beta}_\lambda) = \text{MSE}(\hat{\beta}_\lambda) + 2 \cdot \frac{s_\lambda}{n} \hat{\sigma}^2,$$

where s_λ is the degrees of freedom of the fit (i.e. the number of nonzero parameters selected at that particular value of λ), and $\hat{\sigma}^2 = \text{var}(\epsilon)$ is an estimate of the residual variance. You can interpret the C_p statistic as the in-sample mean-squared error, plus a penalty for in-sample optimism. As you add parameters to a regression model, MSE goes down and the penalty goes up.

Note 1: In general, the main tradeoffs between cross-validation and the C_p statistic are these:

- Cross validation is more robust, in the sense that it does not require that you believe the model is right in order to provide a decent estimate of generalization error. But it is more computationally expensive, and less statistically efficient if the underlying model is approximately right.
- C_p (or related statistics like AIC, BIC, etc) has the advantage that it does not require the repeated “split and refit” of cross validation, and is thus less computationally expensive. It is also more statistically efficient if the model is right. But it is not as robust as cross validation to violations of the underlying modeling assumptions (like linearity).

See the “Estimating prediction error” paper below for more detail.

Note 2: You need to plug an estimate of σ^2 into the C_p statistic. A typical way to proceed is to use the unbiased estimate of σ^2 arising from the ordinary-least-squares solution $\hat{\beta}_{\text{OLS}}$:

$$\hat{\sigma}^2 = \frac{1}{n-p} \sum_{i=1}^n (y_i - x_i^T \hat{\beta}_{\text{OLS}})^2,$$

i.e. the model fit to all p variables. Note: this assumes no intercept due to variable centering; if you have an intercept, subtract an extra degree of freedom in the denominator (i.e. $n - p - 1$).

Another natural estimator for the residual variance σ^2 is to use, for each λ , the formula

$$\hat{\sigma}_\lambda^2 = \frac{1}{n - s_\lambda} \sum_{i=1}^n (y_i - x_i^T \hat{\beta}_\lambda)^2,$$

where again s_λ is the degrees of freedom of the fit. This formula parallels the usual formula given for an unbiased estimate of the error variance for the OLS model.

To be honest, I am not sure if there is any theory supporting the use of one or the other of these variance estimators in the context of the C_p statistic. Mostly I have seen people use the former estimator, based on the OLS fit. The relevant papers here are these, also linked from the class website:

¹⁷[https://en.wikipedia.org/wiki/Bootstrapping_\(statistics\)](https://en.wikipedia.org/wiki/Bootstrapping_(statistics))

- Degrees of freedom of the lasso fit ¹⁸.
- Estimating the residual variance from the lasso fit ¹⁹.
- Estimating prediction error ²⁰. This paper has a much more extensive discussion and list of references about the idea of the C_p statistic and its generalizations.

Solution

The solution is presented in the file *Part.2.R*,

```

1 # Collin Mallows's Cp Statistic
2 Cp_Mallows <- function(X, Y, lambda){
3   # Compute the Collin Mallows's Cp Statistic
4   # Input:
5   # X = features (matrix M x P)
6   # y = vector of observations (vector M)
7   # lambda = set of complexity parameters (vector N)
8   # Output:
9   # Cp = Collin Mallows's Cp Statistic (vector N)
10
11   N <- length(lambda)
12   M <- dim(X)[1]
13   P <- dim(X)[2]
14   # Obtain the unbiased sigma-square from the OLS
15   fit <- glmnet(X, Y, alpha = 1, lambda = 0)
16   Y.OLS <- predict.glmnet(fit, newx = X)
17   sigma.square <- sum((Y - Y.OLS) ^ 2) / (M - P)
18   # In-sample mean square prediction error
19   fit <- glmnet(X, Y, alpha = 1, lambda = lambda)
20   Y.prediction <- predict.glmnet(fit, newx = X)
21   MSE <- array(NA, dim = N)
22   for (i in 1:N){
23     MSE[i] <- mean((Y.prediction[,i] - Y)^2)
24   }
25   # Estimate the Cp Statistics
26   MSE + 2 * (fit$df / M) * sigma.square
27 }
28 #-----
29 # Applying the Collin Mallows's Cp Statistic function
30 MSE_Cp <- Cp_Mallows(X, Y, fit$lambda)
31 png(filename=paste('P2A_MSE4', '.png', sep = ""), width = 15,
32     height = 12, units = "cm", res = 200)
33 plot(log(fit$lambda), MSE, type = "l", xlab = 'log(lambda)')
34 points(log(fit$lambda), MSE_Kfold, type="l", col="blue", lty=2)
35 points(log(fit$lambda), MSE_Cp, type="l", col="red")
36 legend('topleft', legend = c('In-sample MSE', 'K-Fold CV', 'Mallow Cp'),
37     col=c("black", "blue", "red"), lty=1:2, cex=0.8)
38 dev.off()

```

¹⁸<https://projecteuclid.org/euclid.aos/1194461726>

¹⁹<https://arxiv.org/abs/1311.5274>

²⁰https://people.eecs.berkeley.edu/~jordan/sail/readings/archive/efron_Cp.pdf

Compute (and plot) the C_p statistic across the solution path, as a function of λ . How does it compare to the in-sample MSE, and to the cross-validated estimate of generalization error from (B)? Show these all on the same plot. Do they lead to similar choices of λ ?

Figure 7 shows the results for Mallows C_p compared to K-fold cross-validation MOOSE and in-sample MSE. We can observe that Mallows C_p follows the K-fold method and presents similar minimum value of λ but it's less computationally expensive.

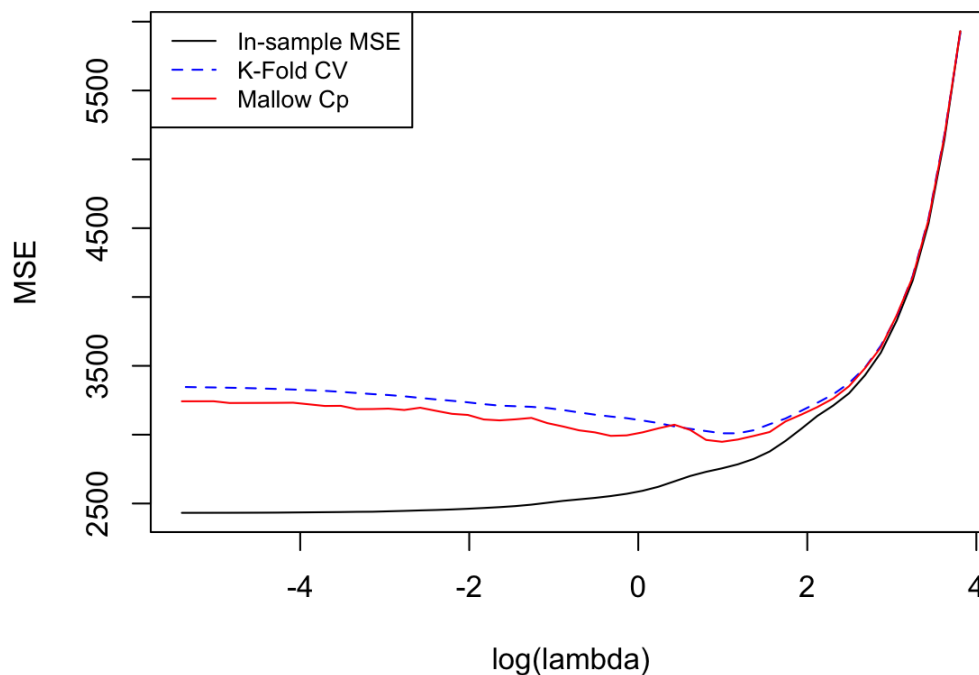


Figure 7: Mallows C_p statistic compared to K-fold cross-validation MOOSE and in-sample MSE