

SDS 385: Exercises 6 - The Proximal Gradient Method

October 18, 2016

Professor Scott

Spencer Woody

Problem 1

Proximal Operators

(A) First, we state the definition of the Moreau envelope and the proximal gradient:

$$E_\gamma f(x) = \min_z \left\{ f(z) + \frac{1}{2\gamma} \|z - x\|_2^2 \right\} \leq f(x) \quad (1)$$

$$\text{prox}_\gamma f(x) = \arg \min_z \left\{ f(z) + \frac{1}{2\gamma} \|z - x\|_2^2 \right\} \quad (2)$$

For the local linear approximation of $f(x)$ about the point x_0 ,

$$f(x) \approx \hat{f}(x; x_0) = f(x_0) + (x - x_0)^T \nabla f(x_0) \quad (3)$$

the proximal operator is derived as follows:

$$\text{prox}_\gamma \hat{f}(x; x_0) = \arg \min_x \left\{ f(x_0) + (x - x_0)^T \nabla f(x_0) + \frac{1}{2\gamma} \|x - x_0\|_2^2 \right\} \quad (4)$$

$$= \arg \min_x \left\{ x^T \nabla f(x_0) + \frac{1}{2\gamma} \|x - x_0\|_2^2 \right\} \quad (5)$$

We take the gradient of the objective function with respect to the feature vector x and set it to zero to find the arg min.

$$\nabla f(x_0) - \frac{1}{\gamma}(x - x_0) = 0 \quad (6)$$

$$x = x_0 + \gamma \nabla f(x_0) \quad (7)$$

This is equivalent to gradient-step for $f(x)$ with size γ .

(B) Suppose we have a negative log-likelihood in the form

$$l(x) = \frac{1}{2} x^T P x - q^T x + r. \quad (8)$$

The proximal operator with parameter $1/\gamma$ of $l(x)$ is

$$\text{prox}_{1/\gamma} l(x) = \arg \min_z \left\{ l(z) + \frac{\gamma}{2} \|z - x\|_2^2 \right\} \quad (9)$$

$$= \arg \min_z \left\{ \frac{1}{2} z^T P z - q^T z + r + \frac{\gamma}{2} \|z - x\|_2^2 \right\} \quad (10)$$

Once again, taking the gradient of the objective function with respect to z and setting it to zero yields

$$Pz - q + \gamma(z - x) = 0 \quad (11)$$

$$(P + \gamma I)z = (\gamma x + q) \quad (12)$$

$$z = (P + \gamma I)^{-1}(\gamma x + q) \quad (13)$$

where I is the identity matrix, assuming that $(P + \gamma I)$ is invertible. Now suppose that y are generated

conditionally on x by $(y|x) \sim N(Ax, \Omega^{-1})$. The log-likelihood of x may be expressed by

$$l(x) \propto \frac{1}{2}(y - Ax)^T \Omega (y - Ax) \quad (14)$$

$$= \frac{1}{2}(y^T - x^T A^T) \Omega (y - Ax) \quad (15)$$

$$= \frac{1}{2}(y^T \Omega - x^T A^T \Omega)(y - Ax) \quad (16)$$

$$= \frac{1}{2}(y^T \Omega y - y^T \Omega Ax - x^T A^T \Omega y + x^T A^T \Omega Ax) \quad (17)$$

$$= \frac{1}{2}(y^T \Omega y - 2y^T \Omega Ax + x^T A^T \Omega Ax) \quad (18)$$

$$= \frac{1}{2}x^T A^T \Omega Ax - y^T \Omega Ax + y^T \Omega y \quad (19)$$

Therefore the negative log-likelihood takes the form expressed in (8) with $P = A^T \Omega A$, $q = A^T \Omega^T y$, and $r = y^T \Omega y$.

(C) Let $\phi(x) = \tau \|x\|_1$.

$$\text{prox}_\gamma \phi(x) = \arg \min_z \left\{ \phi(x) + \frac{1}{2\gamma} \|z - x\|_2^2 \right\} \quad (20)$$

$$= \arg \min_z \left\{ \tau \|z\|_1 + \frac{1}{2\gamma} \|z - x\|_2^2 \right\} \quad (21)$$

$$= \arg \min_z \left\{ \tau \sum_{i=1}^p |z_i| + \frac{1}{2\gamma} \|z - x\|_2^2 \right\} \quad (22)$$

We find the element-wise solution to find each z_i .

$$\arg \min_{z_i} \left\{ \tau |z_i| + \frac{1}{2\gamma} (z_i - x_i)^2 \right\} = \arg \min_{z_i} \left\{ \frac{1}{2} (x_i - z_i)^2 + \gamma \tau |z_i| \right\} \quad (23)$$

In the previous set of exercises, we showed that this is equal to the soft-thresholding function

$$\text{sign}(x_i)(|x_i| - \gamma \tau)_+. \quad (24)$$

Problem 2

The proximal gradient method

(A) We want to minimize an objective function $f(x) = l(x) + \phi(x)$ where $l(x)$ is differentiable but $\psi(x)$ is not. We define the linear approximator of $l(x)$ as

$$l(x) \approx \tilde{l}(x; x_0) = l(x_0) + (x - x_0)^T \nabla l(x_0) + \frac{1}{2\gamma} \|x - x_0\|_2^2. \quad (25)$$

Now our original objective function can be approximated with

$$f(x) \approx \tilde{f}(x; x_0) = \tilde{l}(x; x_0) + \phi(x). \quad (26)$$

We look to minimize this approximation of the objective function.

$$\hat{x} = \arg \min_x \{ \tilde{f}(x; x_0) \} \quad (27)$$

$$= \arg \min_x \{ \tilde{l}(x; x_0) + \phi(x) \} \quad (28)$$

$$= \arg \min_x \left\{ l(x_0) + (x - x_0)^T \nabla l(x_0) + \frac{1}{2\gamma} \|x - x_0\|_2^2 + \phi(x) \right\} \quad (29)$$

We can disregard the $l(x_0)$ term in (29) because it does not relate to x . Furthermore we can introduce another term $\frac{\gamma}{2} \nabla l(x_0)^T \nabla l(x_0)$ which also does not relate to x in order to complete the square. Now the objective function becomes

$$\arg \min_x \left\{ \phi(x) + (x - x_0)^T \nabla l(x_0) + \frac{1}{2\gamma} \|x - x_0\|_2^2 + \frac{\gamma}{2} \nabla l(x_0)^T \nabla l(x_0) \right\} \quad (30)$$

$$= \arg \min_x \left\{ \phi(x) + (x - x_0)^T \nabla l(x_0) + \frac{1}{2\gamma} (x - x_0)^T (x - x_0) + \frac{\gamma}{2} \nabla l(x_0)^T \nabla l(x_0) \right\} \quad (31)$$

$$= \arg \min_x \left\{ \phi(x) + \frac{1}{2\gamma} \left(2\gamma (x - x_0)^T \nabla l(x_0) + (x - x_0)^T (x - x_0) + \gamma^2 \nabla l(x_0)^T \nabla l(x_0) \right) \right\} \quad (32)$$

$$= \arg \min_x \left\{ \phi(x) + \frac{1}{2\gamma} \left((x - x_0)^T (x - x_0) + 2\gamma (x - x_0)^T \nabla l(x_0) + \gamma^2 \nabla l(x_0)^T \nabla l(x_0) \right) \right\} \quad (33)$$

$$= \arg \min_x \left\{ \phi(x) + \frac{1}{2\gamma} \left((x - x_0 + \gamma \nabla l(x_0))^T (x - x_0 + \gamma \nabla l(x_0)) \right) \right\} \quad (34)$$

$$= \arg \min_x \left\{ \phi(x) + \frac{1}{2\gamma} \|x - x_0 + \gamma \nabla l(x_0)\|_2^2 \right\} \quad (35)$$

$$= \arg \min_x \left\{ \phi(x) + \frac{1}{2\gamma} \|x - (x_0 - \gamma \nabla l(x_0))\|_2^2 \right\} \quad (36)$$

Therefore we see the solution for \hat{x} is equivalent to

$$\hat{x} = \underset{\gamma}{\text{prox}} \phi(u), \text{ where } u = x_0 - \gamma \nabla l(x_0) \quad (37)$$

(B) The proximal gradient method is an iterative algorithm which expressed as

$$x^{(t+1)} = \underset{\gamma^{(t)}}{\text{prox}} \phi(u^{(t)}), \quad u^{(t)} = x^{(t)} - \gamma^{(t)} \nabla l(x^{(t)}) \quad (38)$$

With lasso linear regression, we find

$$\hat{\beta} = \arg \min_{\beta} \left\{ \frac{1}{2} \|y - X\beta\|_2^2 + \lambda \|\beta\|_1 \right\}. \quad (39)$$

We may use proximal gradient method to accomplish this. Define

$$l(\beta) = \frac{1}{2} \|y - X\beta\|_2^2 = \frac{1}{2} (y - X\beta)^T (y - X\beta) = \frac{1}{2} y^T y - \beta^T X^T y + \frac{1}{2} \beta^T X^T X \beta \quad (40)$$

$$\nabla l(\beta) = X^T X \beta - X^T y = -X^T (y - X\beta) \quad (41)$$

$$\phi(\beta) = \lambda \|\beta\|_1. \quad (42)$$

With this we can frame our proximal gradient method for finding lasso solutions as follows:

$$\beta^{(t+1)} = \underset{\gamma^{(t)}}{\text{prox}} \lambda \|u^{(t)}\|_1 \quad (43)$$

$$u^{(t)} = \beta^{(t)} - \gamma^{(t)} \nabla l(\beta^{(t)}) \quad (44)$$

$$= \beta^{(t)} + \gamma^{(t)} X^T (y - X\beta^{(t)}) \quad (45)$$

Using the results from the previous exercise, we see that the element-wise solution is

$$\beta_i^{(t+1)} = \text{sign}(u_i^{(t)}) (|u_i^{(t)}| - \gamma \lambda)_+ \quad (46)$$

The primary computational cost of this algorithm computing the two matrix products. Finding the proximal gradient is trivial, only involving the sign and max commands within R. Implementation in R gives similar results to the output of the `glmnet` package. See Figure 1. Note that in my code, I scale λ by N , the number of data points.

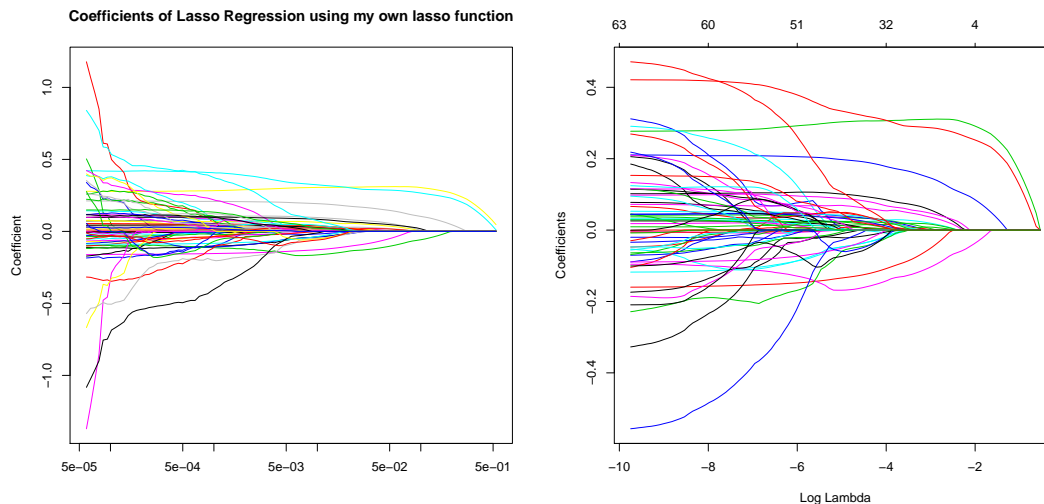


Figure 1: Comparison of output from my proximal gradient method (left) and output of `glmnet` (right)

- (C) Now we implement the accelerated proximal gradient algorithm. The $z^{(t+1)}$ term is an extrapolated version of $\beta^{(t+1)}$ based on the magnitude of the previous step. If we have taken a big step, we are more

confident in going further in the direction we have gone.

$$\beta^{(t+1)} = \underset{\gamma^{(t)}}{\text{prox}} \lambda \|u^{(t)}\|_1 \quad (47)$$

$$u^{(t)} = z^{(t)} - \gamma^{(t)} \nabla l(z^{(t)}) \quad (48)$$

$$s^{(t+1)} = \frac{1 + (1 + 4(s^{(t)}))^{1/2}}{2} \quad (49)$$

$$z^{(t+1)} = \beta^{(t+1)} + \left(\frac{s^{(t)} - 1}{s^{(t+1)}} \right) (\beta^{(t+1)} - \beta^{(t)}) \quad (50)$$

When implemented in R, the accelerated method converges much faster than the regular method. See Figure 2. My convergence criteria is a relative change in negative log-likelihood of less than 10^{-10} . Holding $\lambda = 0.001$, the accelerated method converges in 616 iterations, while the regular method converges in 19,953 iterations. Comparing the final given values of the coefficients, 63 of 64 covariates have the same sign for their coefficients. The one outstanding covariate, `hd1.tch`, has a coefficient of 0.00267 in the unaccelerated method, and a value of zero in the accelerated method. Figure 3 compares all coefficients for the two methods. Clearly the coefficients are in close agreement when computed under the two different methods.

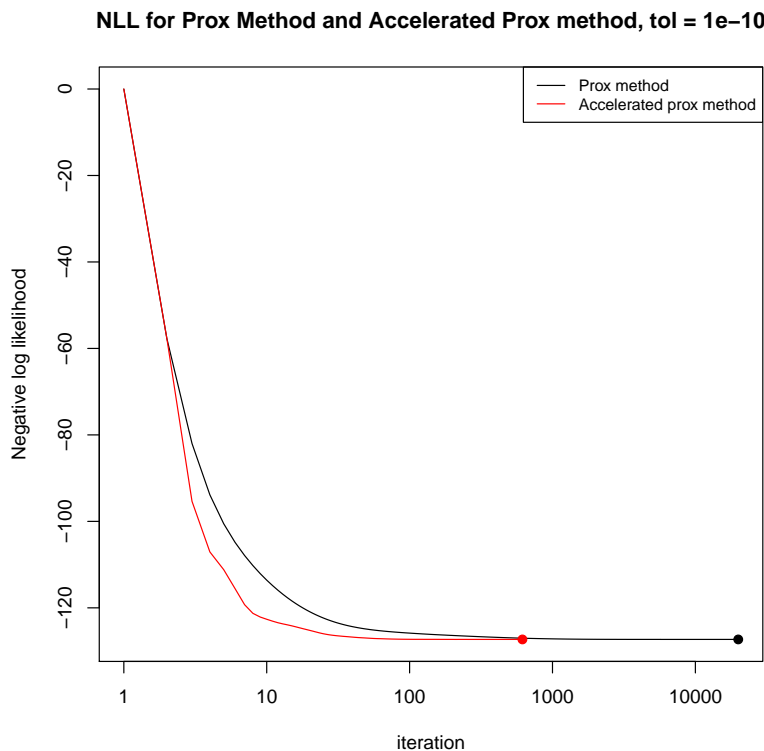


Figure 2: Comparison of convergence for prox method vs. accelerated prox method, $\lambda = 0.001$

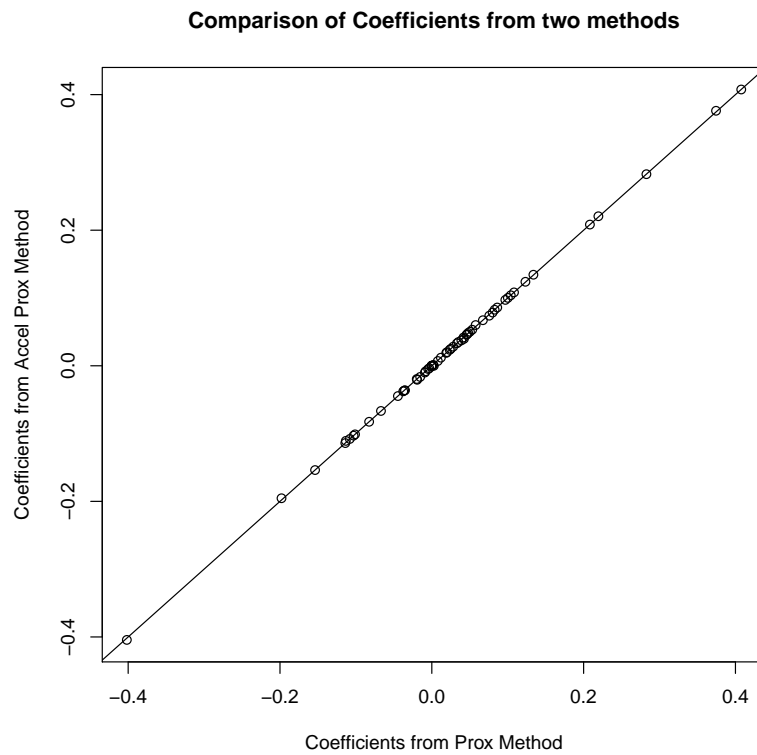


Figure 3: Comparison of coefficients for prox method vs. accelerated prox method, $\lambda = 0.001$

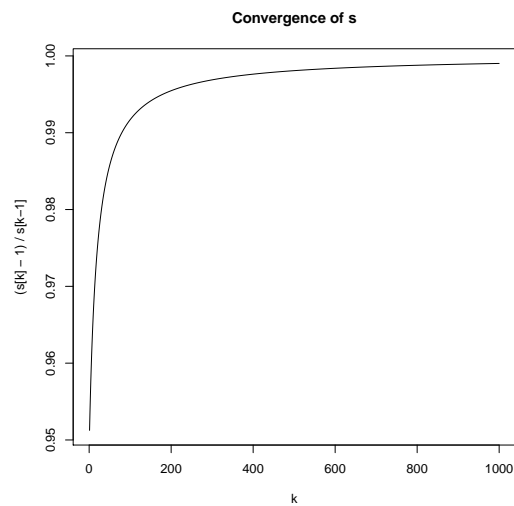


Figure 4: Convergence of $\frac{s^{(t)} - 1}{s^{(t+1)}}$

R script for myfuns.R

```
##### Created by Spencer Woody on 17 Oct 2016 #####

nll <- function(X, Y, beta) {
  mynll <- t(beta) %*% (0.5 * crossprod(X) %*% beta - crossprod(X, Y))
  return(mynll)
}

grad <- function(X, Y, beta) {
  myprod <- Y - X%*%beta
  mygrad <- - crossprod(X, myprod)
  return(mygrad)
}

prox.ell1 <- function(u, lambda) {
  uprime <- abs(u) - lambda
  uprime.zeros <- cbind(rep(0, length(u)), uprime)
  prox <- sign(u) * apply(uprime.zeros, 1, max)
  return(prox)
}

proxlasso <- function(X, Y, lambda = 1e-4, gamma = 1e-4, beta0 = NA, numsteps = 50000, tol = 1e-10) {
  # If there is no beta0, make it a series of zeros
  if (is.na(beta0)) {
    beta0 <- rep(0, ncol(X))
  }

  # Standardize lambda
  lambda = lambda * nrow(X)

  # Initialize matrix for
  beta <- matrix(rep(NA, ncol(X) * (numsteps+1)), nrow = ncol(X))
  beta[, 1] <- beta0

  # Begin trace of nll
  nlltrace <- rep(NA, (numsteps + 1))
  nlltrace[1] <- nll(X, Y, beta0) + lambda * sum(abs(beta0))

  # Initialize convergence message in case convergence not reached
  message <- "Convergence not reached..."

  for (t in 1:numsteps) {
    u <- beta[, t] - gamma * grad(X, Y, beta[, t])

    beta[, t + 1] <- prox.ell1(u, gamma * lambda)

    nlltrace[t + 1] <- nll(X, Y, beta[, t + 1]) + lambda * sum(abs(beta[, t + 1]))

    # Convergence check
  }
}
```



```

    nllchange <- nlltrace[t + 1] / nlltrace[t] - 1
    if (abs(nllchange) < tol) {
55      # Remove excess betas and nll
      beta <- beta[, -((t+2):ncol(beta))]
      nlltrace <- nlltrace[-((t+2):length(nlltrace))]
      # Update convergence message
      message <- sprintf("Convergence reached after %i iterations", (t+1))
60      break
    }
  }

  mylist <- list("finalbeta" = beta[, ncol(beta)], "beta" = beta, "nll" = nlltrace, "conv" = message)
65  return(mylist)
}

70 proxlassofast <- function(X, Y, lambda = 1e-4, gamma = 1e-4, beta0 = NA, numsteps = 5000, tol = 1e-4) {
  # If there is no beta0, make it a series of zeros
  if (is.na(beta0)) {
    beta0 <- rep(0, ncol(X))
  }

75  # Standardize lambda
  lambda = lambda * nrow(X)

  # Create s vector
80  s <- rep(NA, numsteps + 1)
  s[1] <- 10

  for (j in 2:length(s)) {
85    s[j] <- (1 + (1 + 4 * s[j - 1] ^ 2) ^ 0.5) / 2
  }

  # Initialize z matrix
  z <- matrix(0, nrow = ncol(X), ncol = (numsteps + 1))

90  # Initialize matrix for
  beta <- matrix(rep(NA, ncol(X) * (numsteps+1)), nrow = ncol(X))
  beta[, 1] <- beta0

  # Begin trace of nll
95  nlltrace <- rep(NA, (numsteps + 1))
  nlltrace[1] <- nll(X, Y, beta0) + lambda * sum(abs(beta0))

  # Initialize convergence message in case convergence not reached
  message <- "Convergence not reached..."
100

  for (t in 1:numsteps) {
    u <- z[, t] - gamma * grad(X, Y, z[, t])

    beta[, t + 1] <- prox.ell1(u, gamma * lambda)
105
  }
}

```

```
z[, t + 1] <- beta[, t + 1] + (s[t] - 1) / s[t + 1] * (beta[, t + 1] - beta[, t])

nlltrace[t + 1] <- nll(X, Y, beta[, t + 1]) + lambda * sum(abs(beta[, t + 1]))

110 # Convergence check
nllchange <- nlltrace[t + 1] / nlltrace[t] - 1
if (abs(nllchange) < tol) {
  # Remove excess betas and nll
  beta <- beta[, -((t+2):ncol(beta))]
115 nlltrace <- nlltrace[-((t+2):length(nlltrace))]
  # Update convergence message
  message <- sprintf("Convergence reached after %i iterations", (t+1))
  break
}
120 }

mylist <- list("finalbeta" = beta[, ncol(beta)], "beta" = beta, "nll" = nlltrace, "conv" = message)
return(mylist)
}
```

R script for e6.R

```
##### Created by Spencer Woody on 17 Oct 2016 #####

library(MASS)
library(glmnet)
5 library(ggplot2)

source("myfuns.R")

# Read in data, scale everything
10 X <- as.matrix(read.csv("diabetesX.csv", header = TRUE))
Y <- as.numeric(unlist(read.csv("diabetesY.csv", header = FALSE)))

X <- scale(X)
Y <- scale(Y)
15

lambdas <- seq(5e-5, 0.5, length.out = 50)

# Plot coefficients from lasso regression
20 betamat <- matrix(rep(NA, length(lambdas)*ncol(X)), ncol = length(lambdas))

for (i in 1:length(lambdas)) {
  mylasso <- proxlassofast(X, Y, lambda = lambdas[i])
  25 betamat[, i] <- mylasso$finalbeta
}

ybound <- c(min(betamat), max(betamat))

30 pdf("myplot.pdf")
plot(lambdas, betamat[1, ],
     ylim = ybound,
     log = "x",
     main = "Coefficients of Lasso Regression using my own lasso function",
     35 xlab = "",
     ylab = "Coefficient",
     type = "l",
     col = 2)
for (j in 2:nrow(betamat)) {
  40 lines(lambdas, betamat[j, ], col = (j + 4))
}
dev.off()

weight <- 0
45
for (k in 2:length(s)) {
  weight[k - 1] <- (s[k] - 1) / (s[k-1])
}

50 pdf("splot.pdf")
plot(weight, type = "l", main = "Convergence of s", xlab = "k", ylab = "(s[k] - 1) / s[k-1]")
dev.off()
```

```
#####

55 # Compare ordinary and accelerated versions of proximal gradient

m <- -20

60 mylasso <- proxlasso(X, Y, gamma = 0.0002, lambda = 0.001, tol = 10^m)
mylasso2 <- proxlassofast(X, Y, gamma = 0.0002, lambda = 0.001, tol = 10^m)

pdf("slowfast.pdf")
plot(mylasso$nll, type = "l", col = "black",
65 log = "x",
xlab = "iteration",
ylab = "Negative log likelihood",
main = sprintf("NLL for Prox Method and Accelerated Prox method, tol = 1e%i", m))
lines(mylasso2$nll, col = "red")
70 points(length(mylasso$nll), mylasso$nll[length(mylasso$nll)], pch = 19, col = "black")
points(length(mylasso2$nll), mylasso2$nll[length(mylasso2$nll)], pch = 19, col = "red")
legend("topright", legend=c("Prox method", "Accelerated prox method"),
      col=c("black", "red"), lty=c(1,1), cex=0.8)
dev.off()

75 print(length(mylasso$nll))
print(length(mylasso2$nll))

sum(sign(mylasso$finalbeta) == sign(mylasso2$finalbeta))
80 ncol(X)

pdf("compcoefs.pdf")
plot(mylasso$finalbeta, mylasso2$finalbeta,
      main = "Comparison of Coefficients from two methods",
85 xlab = "Coefficients from Prox Method",
ylab = "Coefficients from Accel Prox Method")
abline(0, 1)
dev.off()

90 mylasso$finalbeta
mylasso2$finalbeta

# > print(length(mylasso$nll))
# [1] 19953
95 # > print(length(mylasso2$nll))
# [1] 616

# > mylasso$finalbeta
# [1] 0.023560578 -0.153789540 0.282694838 0.208274783 -0.103033550
100 # [6] 0.000000000 -0.113948268 0.053230187 0.407621844 0.041754067
# [11] 0.046650203 0.033109963 -0.008095386 0.057726330 0.000000000
# [16] -0.019305929 0.080070994 0.219445533 0.067286417 0.103984948
# [21] -0.004789084 0.011970890 0.000000000 -0.108030864 0.082766137
# [26] 0.086187888 0.034960689 0.028160431 0.045641503 0.050263648
105 # [31] 0.108286794 -0.101133491 0.000000000 -0.044642852 -0.037578683
```

```

# [36] 0.024963870 0.100191687 -0.113246424 0.075705929 0.038941220
# [41] -0.009111873 0.007940555 0.020034844 0.047921927 0.001962857
# [46] 0.000000000 -0.019820825 -0.003312034 -0.082569062 0.042338316
# [51] 0.000000000 -0.197988211 -0.401744208 -0.036250001 -0.066970573
110 # [56] -0.015327546 0.374411995 0.000000000 0.002665353 0.123333658
# [61] 0.096873190 -0.035289003 0.133837012 0.019053591
# > mylasso2$finalbeta
# [1] 0.023538900 -0.153795619 0.282633738 0.208328306 -0.103014142
# [6] 0.000000000 -0.113982017 0.053249091 0.407576274 0.041765184
115 # [11] 0.046624516 0.033198478 -0.008081280 0.060115879 0.000000000
# [16] -0.020599252 0.078493833 0.220584671 0.067223390 0.103999771
# [21] -0.004744999 0.011962682 0.000000000 -0.107915299 0.082692339
# [26] 0.085975450 0.035031585 0.028175283 0.045726185 0.050280271
# [31] 0.108269248 -0.101125076 0.000000000 -0.044626911 -0.037644834
120 # [36] 0.024955140 0.100180730 -0.110735176 0.073923786 0.037456934
# [41] -0.009942927 0.007065391 0.020099087 0.048478022 0.001293112
# [46] 0.000000000 -0.019492307 -0.003648320 -0.082558916 0.039746085
# [51] 0.000000000 -0.195383539 -0.404309136 -0.036472332 -0.066619737
# [56] -0.016430649 0.376181451 0.000000000 0.000000000 0.123918165
125 # [61] 0.097184059 -0.036062764 0.134341872 0.019007764

#####
# > which(beta[, ncol(beta)] == 0)
# [1] 13 14 15 16 17 21 23 26 31 35 36 38 39 40 41 42 43 45 47 48 50 54 55 56 60
130 # [26] 61 64

plot(nll1, type = "l")
lines(nll2, col = "red")
135

numlambdas <- 100

140 # glmnet requires X to be a matrix, Y to be a vector
fit1 <- glmnet(X, Y, nlambdas = 100)

# Make plot of
pdf("glmnetplot.pdf")
145 plot(fit1, xvar = "lambda")
dev.off()

# In-sample MSE
MSE.lasso <- rep(NA, numlambdas)
150 for (j in 1:numlambdas) {
  Y.hats <- X %*% fit1$beta[, j]
  MSE.lasso[j] <- sum((Y - Y.hats)^2) / N2
}

155 ### Cross validation

```