

SDS 385: Homework 6

G. Paulon

October 19, 2016

Problem 1. Proximal operators

Let $f(x)$ be a convex function. The **Moreau envelope** $E_\gamma f(x)$ and **proximal operator** $\text{prox}_\gamma f(x)$ for parameter $\gamma > 0$ are defined as

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

$$\text{prox}_\gamma f(x) = \underset{z}{\operatorname{argmin}} \left\{ f(z) + \frac{1}{2\gamma} \|z - x\|_2^2 \right\}.$$

It is easy to see that $E_\gamma f(x) \leq f(x)$.

Moreover, let us suppose that x^* is the minimum point of $f(x)$. This is also the minimum point of the envelope function. In fact,

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

because x^* minimizes both parts of the function. Moreover, let \tilde{x} be another point

$$E_\gamma f(\tilde{x}) = \min_z \left\{ f(z) + \frac{1}{2\gamma} \|z - \tilde{x}\|^2 \right\} \geq \min_z f(z) + \min_{\tilde{z}} \frac{1}{2\gamma} \|\tilde{z} - \tilde{x}\|^2 = f(x^*).$$

The proximal operator of a function evaluated at one point moves the point towards the minimum. More precisely, the transformed point $\text{prox}_\gamma f(x)$ is a compromise between minimizing f and being near to x . The parameter γ controls this trade-off between these two terms. The Moreau envelope, instead, is a regularized version of f . It approximates f from below, and has the same set of minimizing values as f .

(A) *The proximal operator gives a nice interpretation of classical gradient descent. Consider the local linear approximation of $f(\mathbf{x})$ about a point \mathbf{x}_0 :*

$$f(\mathbf{x}) \approx \hat{f}(\mathbf{x}; \mathbf{x}_0) = f(\mathbf{x}_0) + (\mathbf{x} - \mathbf{x}_0)^T \nabla f(\mathbf{x}_0).$$

Derive the proximal operator (with parameter γ) of the linear approximation $\hat{f}(\mathbf{x}; \mathbf{x}_0)$, and show that this proximal operator is identical to a gradient descent step for $f(\mathbf{x})$ of size γ , starting from the point \mathbf{x}_0 .

Computing the proximal operator of $\hat{f}(\mathbf{x}; \mathbf{x}_0)$ is equivalent to moving towards the minimum of $\hat{f}(\mathbf{x}; \mathbf{x}_0)$ staying still close to \mathbf{x}_0 . This makes sense, since the problem consists in minimizing a **local** linear approximation and so it requires to stay close to \mathbf{x}_0 .

$$\text{prox}_\gamma \hat{f}(\mathbf{x}; \mathbf{x}_0) = \underset{\mathbf{x}}{\operatorname{argmin}} \left\{ f(\mathbf{x}_0) + (\mathbf{x} - \mathbf{x}_0)^T \nabla f(\mathbf{x}_0) + \frac{1}{2\gamma} \|\mathbf{x} - \mathbf{x}_0\|_2^2 \right\}$$

This problem can be easily solved by computing the gradient of the objective function

$$\begin{aligned} & \nabla_{\mathbf{x}} \left\{ f(\mathbf{x}_0) + (\mathbf{x} - \mathbf{x}_0)^T \nabla f(\mathbf{x}_0) + \frac{1}{2\gamma} (\mathbf{x} - \mathbf{x}_0)^T (\mathbf{x} - \mathbf{x}_0) \right\} \\ &= \nabla f(\mathbf{x}_0) + \frac{1}{\gamma} (\mathbf{x} - \mathbf{x}_0) \end{aligned}$$

which is then set to 0, yielding

$$\hat{x} = \text{prox}_{\gamma} \hat{f}(x; x_0) = x_0 - \gamma \nabla f(x_0),$$

which is exactly the gradient update. Therefore, the gradient step minimizes a local linear approximation of the function in a neighbourhood of the current iterate (where, presumably, the linear approximation is reasonable).

- (B) *Many intermediate steps in statistical optimization problems can be written very compactly in terms of proximal operators of log-likelihoods or penalty functions. For example, consider a negative log likelihood of the form*

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

First, show that if we have a Gaussian sampling model of the form $(y|x) \sim N(Ax, \Omega^{-1})$, then our negative log-likelihood can be written in the form given above. Then show that the proximal operator with parameter $1/\gamma$ of $l(x)$ takes the form

$$\text{prox}_{\frac{1}{\gamma}} l(x) = (P + \gamma I)^{-1} (\gamma x - q)$$

assuming the relevant inverse exists.

This is a multivariate normal model, whose likelihood function can be written as

$$L(y; x) \propto \exp \left\{ -\frac{1}{2} (y - Ax)^T \Omega (y - Ax) \right\}.$$

Therefore the negative log-likelihood is

$$\begin{aligned} l(y; x) &\propto \frac{1}{2} (y - Ax)^T \Omega (y - Ax) \\ &= \frac{1}{2} y^T \Omega y + \frac{1}{2} x^T A^T \Omega A x - y^T \Omega A x \\ &= \frac{1}{2} x^T P x + q^T x + r \end{aligned}$$

where

$$\begin{aligned} P &= A^T \Omega A \\ q &= -A^T \Omega y \\ r &= \frac{1}{2} y^T \Omega y. \end{aligned}$$

The proximal operator is then

$$\begin{aligned} \text{prox}_{\frac{1}{\gamma}} l(x) &= \underset{z}{\operatorname{argmin}} \left\{ l(z) + \frac{\gamma}{2} \|z - x\|_2^2 \right\} \\ &= \underset{z}{\operatorname{argmin}} \left\{ \frac{1}{2} z^T P z + q^T z + r + \frac{\gamma}{2} (z - x)^T (z - x) \right\}. \end{aligned}$$

The minimum can be found by computing the gradient

$$\begin{aligned} & \nabla \left\{ \frac{1}{2} \mathbf{z}^T P \mathbf{z} + \mathbf{q}^T \mathbf{z} + r + \frac{\gamma}{2} (\mathbf{z} - \mathbf{x})^T (\mathbf{z} - \mathbf{x}) \right\} \\ &= P \mathbf{z} + \mathbf{q} + \frac{\gamma}{2} 2(\mathbf{z} - \mathbf{x}), \end{aligned}$$

where in the last equality we exploited the fact that P is a symmetric matrix. Therefore the optimum is obtained by setting the gradient equal to 0, that is

$$\begin{aligned} (P + \gamma I) \hat{\mathbf{z}} &= \gamma \mathbf{x} - \mathbf{q} \\ \Rightarrow \hat{\mathbf{z}} &= \text{prox}_{\frac{\gamma}{2}} l(\mathbf{x}) = (P + \gamma I)^{-1} (\gamma \mathbf{x} - \mathbf{q}) \end{aligned}$$

(C) Let $\phi(\mathbf{x}) = \tau \|\mathbf{x}\|_1$. Express the proximal operator of this function in terms of the soft-thresholding function that we learned about in the last set of exercises.

The proximal operator of this function is

$$\text{prox}_{\gamma} \phi(\mathbf{x}) = \underset{\mathbf{z}}{\text{argmin}} \left\{ \tau \|\mathbf{z}\|_1 + \frac{1}{2\gamma} \|\mathbf{z} - \mathbf{x}\|_2^2 \right\}.$$

To solve this minimization problem, as usual, we compute the gradient and we set it to 0. Let us compute the generic partial derivative (element-wise expression since it is a separable objective function):

$$\begin{aligned} & \left. \frac{\partial}{\partial z_i} \left\{ \tau |z_i| + \frac{1}{2\gamma} (z_i - x_i)^2 \right\} \right|_{z_i = \hat{z}_i} \\ &= \tau \text{sign}(\hat{z}_i) + \frac{1}{\gamma} (\hat{z}_i - x_i). \end{aligned}$$

Let us now split the three cases:

- if $\hat{z}_i > 0$, then $\hat{z}_i = x_i - \tau\gamma$ under the constraint $x_i > \tau\gamma$;
- if $\hat{z}_i < 0$, then $\hat{z}_i = x_i + \tau\gamma$ under the constraint $x_i < -\tau\gamma$;
- if $\hat{z}_i = 0$, that means that $-\tau\gamma < x_i < \tau\gamma$.

Therefore, a compact form can be written as

$$\hat{\mathbf{z}} = \text{prox}_{\gamma} \phi(\mathbf{x}) = \text{sign}(\mathbf{x}) (|\mathbf{x}| - \tau\gamma)_+ = S_{\tau\gamma}(\mathbf{x}).$$

Problem 2. The proximal gradient method

Suppose that we have some objective function that can be expressed as $f(\mathbf{x}) = l(\mathbf{x}) + \phi(\mathbf{x})$, where $l(\mathbf{x})$ is differentiable but $\phi(\mathbf{x})$ is not. Recall from above the idea of forming a local linear approximation to a function at some point \mathbf{x}_0 and then adding a quadratic regularizer. This gave us an interpretation of gradient descent evaluating the proximal operator of our locally linear approximation. Here, we will apply this idea to the first term in our objective, $l(\mathbf{x})$. Define

$$l(\mathbf{x}) \approx \tilde{l}(\mathbf{x}; \mathbf{x}_0) = l(\mathbf{x}_0) + (\mathbf{x} - \mathbf{x}_0)^T \nabla l(\mathbf{x}_0) + \frac{1}{2\gamma} \|\mathbf{x} - \mathbf{x}_0\|_2^2$$

as our linear approximation to $l(\mathbf{x})$, plus the quadratic regularizer. Now we add in the $\phi(\mathbf{x})$ term to get the approximation for our original objective:

$$f(\mathbf{x}) \approx \tilde{f}(\mathbf{x}; \mathbf{x}_0) = \tilde{l}(\mathbf{x}; \mathbf{x}_0) + \phi(\mathbf{x}).$$

(A) Consider the surrogate optimization problem in which we minimize the approximation $\tilde{f}(\mathbf{x}; \mathbf{x}_0)$ in lieu of our original objective $f(\mathbf{x})$, i.e.

$$\hat{\mathbf{x}} = \operatorname{argmin}_{\mathbf{x}} \left\{ \tilde{l}(\mathbf{x}; \mathbf{x}_0) + \phi(\mathbf{x}) \right\}.$$

Show that the solution to this problem is of the form

$$\hat{\mathbf{x}} = \operatorname{prox}_{\gamma} \phi(\mathbf{u}), \quad \text{where } \mathbf{u} = \mathbf{x}_0 - \gamma \nabla l(\mathbf{x}_0).$$

This is just the proximal operator of the non-smooth part of the objective, $\phi(\mathbf{x})$, evaluated at an intermediate gradient-descent step for the smooth part, $l(\mathbf{x})$.

We have to solve the following minimization problem:

$$\begin{aligned} \hat{\mathbf{x}} &= \operatorname{argmin}_{\mathbf{x}} \left\{ \tilde{l}(\mathbf{x}; \mathbf{x}_0) + \phi(\mathbf{x}) \right\} \\ &= \operatorname{argmin}_{\mathbf{x}} \left\{ l(\mathbf{x}_0) + (\mathbf{x} - \mathbf{x}_0)^T \nabla l(\mathbf{x}_0) + \frac{1}{2\gamma} \|\mathbf{x} - \mathbf{x}_0\|_2^2 + \phi(\mathbf{x}) \right\}. \end{aligned}$$

Since the first term of the sum does not depend on \mathbf{x} , we can get rid of it. For the same reason, we can introduce another term (its purpose will be clear at the end of the proof) $\frac{\gamma}{2} \nabla l(\mathbf{x}_0)^T \nabla l(\mathbf{x}_0)$ which does not depend on \mathbf{x} , and the minimization problem is the same. Therefore

$$\begin{aligned} \hat{\mathbf{x}} &= \operatorname{argmin}_{\mathbf{x}} \left\{ l(\mathbf{x}_0) + (\mathbf{x} - \mathbf{x}_0)^T \nabla l(\mathbf{x}_0) + \frac{1}{2\gamma} \|\mathbf{x} - \mathbf{x}_0\|_2^2 + \phi(\mathbf{x}) \right\} \\ &= \operatorname{argmin}_{\mathbf{x}} \left\{ \phi(\mathbf{x}) + (\mathbf{x} - \mathbf{x}_0)^T \nabla l(\mathbf{x}_0) + \frac{1}{2\gamma} (\mathbf{x} - \mathbf{x}_0)^T (\mathbf{x} - \mathbf{x}_0) + \frac{\gamma}{2} \nabla l(\mathbf{x}_0)^T \nabla l(\mathbf{x}_0) \right\} \\ &= \operatorname{argmin}_{\mathbf{x}} \left\{ \phi(\mathbf{x}) + \frac{1}{2\gamma} \left[(\mathbf{x} - \mathbf{x}_0)^T (\mathbf{x} - \mathbf{x}_0) + 2\gamma (\mathbf{x} - \mathbf{x}_0)^T \nabla l(\mathbf{x}_0) + \gamma^2 \nabla l(\mathbf{x}_0)^T \nabla l(\mathbf{x}_0) \right] \right\} \\ &= \operatorname{argmin}_{\mathbf{x}} \left\{ \phi(\mathbf{x}) + \frac{1}{2\gamma} (\mathbf{x} - \mathbf{x}_0 + \gamma \nabla l(\mathbf{x}_0))^T (\mathbf{x} - \mathbf{x}_0 + \gamma \nabla l(\mathbf{x}_0)) \right\} \\ &= \operatorname{argmin}_{\mathbf{x}} \left\{ \phi(\mathbf{x}) + \frac{1}{2\gamma} \|\mathbf{x} - \mathbf{x}_0 + \gamma \nabla l(\mathbf{x}_0)\|_2^2 \right\} \\ &= \operatorname{argmin}_{\mathbf{x}} \left\{ \phi(\mathbf{x}) + \frac{1}{2\gamma} \|\mathbf{x} - (\mathbf{x}_0 - \gamma \nabla l(\mathbf{x}_0))\|_2^2 \right\} \end{aligned}$$

Let us now define $\mathbf{u} = \mathbf{x}_0 - \gamma \nabla l(\mathbf{x}_0)$. Then,

$$\begin{aligned} \hat{\mathbf{x}} &= \operatorname{argmin}_{\mathbf{x}} \left\{ \phi(\mathbf{x}) + \frac{1}{2\gamma} \|\mathbf{x} - (\mathbf{x}_0 - \gamma \nabla l(\mathbf{x}_0))\|_2^2 \right\} \\ &= \operatorname{argmin}_{\mathbf{x}} \left\{ \phi(\mathbf{x}) + \frac{1}{2\gamma} \|\mathbf{x} - \mathbf{u}\|_2^2 \right\} \\ &= \operatorname{prox}_{\gamma} \phi(\mathbf{u}). \end{aligned}$$

(B) The proximal gradient method is an iterative algorithm that consists in the following step

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

Now consider the lasso regression problem:

$$\hat{\beta} = \underset{\beta}{\operatorname{argmin}} \left\{ \frac{1}{2} \|\mathbf{y} - X\beta\|_2^2 + \lambda \|\beta\|_1 \right\}.$$

Using the results on proximal operators you have derived already, write down some concise pseudo-code for using the proximal gradient algorithm to minimize this objective. Identify the primary computational costs of this algorithm.

In the Lasso context, we can set

$$l(\beta) = \frac{1}{2} \|\mathbf{y} - X\beta\|_2^2, \quad \phi(\beta) = \lambda \|\beta\|_1,$$

where the first function is differentiable and the second is not. The proximal gradient algorithm consists then in the minimization of the approximation of the objective $\tilde{l}(\beta; \beta_0) + \phi(\beta)$. The update step takes the following form

$$\beta^{(t+1)} = \text{prox}_{\gamma^{(t)}} \phi(\mathbf{u}^{(t)}) = S_{\gamma^{(t)}\lambda}(\mathbf{u}^{(t)}),$$

where $\mathbf{u}^{(t)} = \beta^{(t)} - \gamma^{(t)} \nabla l(\beta^{(t)})$. Let us recall that

$$\begin{aligned} \nabla l(\beta) &= \nabla \left\{ \frac{1}{2} (\mathbf{y} - X\beta)^T (\mathbf{y} - X\beta) \right\} \\ &= -X^T (\mathbf{y} - X\beta), \end{aligned}$$

and therefore the update step becomes

$$\mathbf{u}^{(t)} = \beta^{(t)} + \gamma^{(t)} X^T (\mathbf{y} - X\beta^{(t)}).$$

Evaluating $\nabla l(\beta)$ requires two matrix-vector multiplications, plus a negligible vector addition. Evaluating the proximal operator of $\phi(\mathbf{u})$ is negligible.

Algorithm 1 Proximal Gradient

```

1: procedure PROXIMAL.GRAIENT( $y, X, \beta^{(1)}, \lambda, \gamma = 10^{-4}$ )
2:   for iter = 2:maxiter do
3:      $\text{gradient} \leftarrow \text{GRAD.LOGLIK}(\beta^{(\text{iter}-1)}, y, X)$ 
4:      $\mathbf{u}^{(t)} \leftarrow \beta^{(\text{iter}-1)} - \gamma \cdot \text{gradient}$ 
5:      $\beta^{(\text{iter})} \leftarrow \text{SOFT.THRESHOLDING}(\mathbf{u}^{(t)}, \gamma\lambda)$ 

```

(C) A cool variation on proximal gradient is called the accelerated proximal gradient algorithm. The following scheme involves a simple extrapolation step based on the previous iteration:

$$\begin{aligned}\beta^{(t+1)} &= \text{prox}_{\gamma^{(t)}} \phi(\mathbf{u}^{(t)}), \quad \mathbf{u}^{(t)} = \mathbf{z}^{(t)} - \gamma^{(t)} \nabla l(\mathbf{z}^{(t)}) \\ s^{(t+1)} &= \frac{1 + (1 + 4(s^{(t)})^2)^{1/2}}{2} \\ \mathbf{z}^{(t+1)} &= \beta^{(t+1)} + \left(\frac{s^{(t)} - 1}{s^{(t+1)}} \right) (\beta^{(t+1)} - \beta^{(t)})\end{aligned}$$

The first step involves the proximal operator of the penalty function, evaluated not at the previous iterate $\beta^{(t)}$, but at an extrapolated version of $\beta^{(t)}$, based on the magnitude of the previous step's update $(\beta^{(t+1)} - \beta^{(t)})$. In this sense, large steps give “momentum” to the next step, where the amount of momentum is modulated by the scalar $s^{(t)}$ terms. Implement this acceleration scheme in your proximal gradient code, and compare its convergence speed to that of the unacceleration version. Does the value of the objective goes down at each and every step? Do you get to the minimum faster?

Algorithm 2 Accelerated Proximal Gradient

```

1: procedure ACC.PROXIMAL.GRADIENT( $y, X, \beta^{(1)}, \lambda, \gamma = 10^{-4}$ )
2:   for iter = 2:maxiter do
3:      $\text{gradient} \leftarrow \text{GRAD.LOGLIK}(\mathbf{z}^{(\text{iter}-1)}, y, X)$ 
4:      $\mathbf{u}^{(t)} \leftarrow \mathbf{z}^{(\text{iter}-1)} - \gamma \cdot \text{gradient}$ 
5:      $\beta^{(\text{iter})} \leftarrow \text{SOFT.THRESHOLDING}(\mathbf{u}^{(t)}, \gamma\lambda)$ 
6:      $s^{(\text{iter})} \leftarrow \frac{1 + \sqrt{1 + 4(s^{(\text{iter}-1)})^2}}{2}$ 
7:      $\mathbf{z}^{(\text{iter})} \leftarrow \beta^{(\text{iter})} + \left( \frac{s^{(\text{iter}-1)} - 1}{s^{(\text{iter})}} \right) (\beta^{(\text{iter})} - \beta^{(\text{iter}-1)})$ 
```

In Figure 1 we see the performance of the two algorithms. We used, as a convergence criterion, the relative decrement of the objective function. The accelerated version converges faster than the normal one. The two different algorithms have then been compared to the results obtained by R's `glmnet` package, showing the exact same estimated β values.

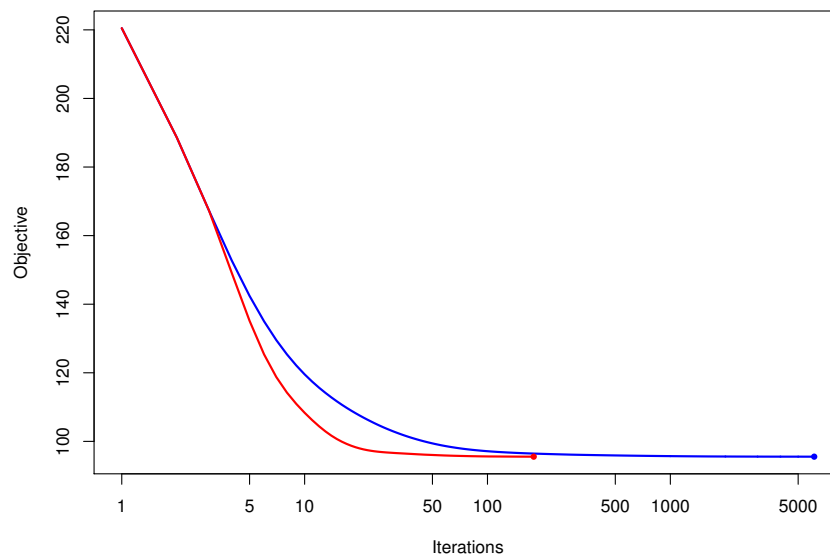


Figure 1: Convergence of the two algorithms. In red, the accelerated version of the proximal gradient.

Appendix A

R code

```
1 log.lik <- function(beta, y, X, lambda){
2   # Computes the objective function of the logit model
3   # -----
4   # Args:
5   #   - beta: regression parameters (length p)
6   #   - y: response vector (length n)
7   #   - X: matrix of the features (n*p)
8   #   - lambda: penalization parameter
9   # Returns:
10  #   - ll: the scalar associated with the objective function at beta
11  # -----
12  Xbeta <- X %*% beta
13  ll <- (1/2)*crossprod(y - Xbeta) + lambda * sum(abs(beta))
14  return(ll)
15 }
16
17 grad.loglik <- function(beta, y, X){
18   # Computes the gradient of negative log-likelihood of the logit model
19   # -----
20   # Args:
21   #   - beta: regression parameters (length p)
22   #   - y: response vector (length n)
23   #   - X: matrix of the features (n*p)
24   # Returns:
25   #   - grad.ll: the gradient vector of the negative log-lik of the data (length p)
26   # -----
27   grad.ll <- array(NA, dim = length(beta))
28   Xbeta <- X %*% beta
29   grad.ll <- - crossprod(X, as.numeric(y - Xbeta))
30   return(grad.ll)
31 }
32
33 soft.thresholding <- function(x, lambda){
34   # Computes the soft thresholding estimator
35   # -----
36   # Args:
37   #   - x: vector of the observations
38   #   - lambda: penalization parameter (threshold)
39   # Returns:
40   #   - theta: the soft thresholding estimator
```

```

41 # -----
42 theta <- sign(x) * pmax(rep(0, length(x)), abs(x) - lambda)
43 return (theta)
44 }

```

Listing A.1: Auxiliary functions.

```

1 proximal.gradient <- function(y, X, beta0, lambda, gamma = 1E-4, maxiter = 50000, tol = 1E-10){
2   # Proximal gradient algorithm that allows to compute the Lasso regression estimates
3   # -----
4   # Args:
5   #   - y: response vector (length n)
6   #   - X: matrix of the features (n*p)
7   #   - beta0: initial values of beta
8   #   - lambda: penalization parameter for the lasso (threshold)
9   #   - gamma: proximal operator parameter
10  #   - maxiter: maximum number of iterations
11  #   - tol: tolerance defining convergence
12  # Returns:
13  #   - ll: values of the negative log-likelihood for each iteration
14  #   - betas: values of the beta parameters for each iteration
15  # -----
16
17  # Check if the inputs have the correct sizes
18  if ((length(y) != dim(X)[1]) || (length(beta0) != dim(X)[2])){
19    stop("Incorrect input dimensions.")
20  }
21
22  # Initialize the data structures
23  betas <- array(NA, dim=c(maxiter, length(beta0)))
24  betas[1,] <- beta0 # Initial guess
25  ll <- array(NA, dim = maxiter)
26  ll[1] <- log.lik(betas[1,], y, X, lambda)
27
28  for (iter in 2:maxiter){
29
30    # Compute the gradient
31    gradient <- grad.loglik(betas[iter-1, ], y, X)
32    u <- betas[iter-1, ] - gamma * gradient
33
34    # Update beta
35    betas[iter, ] <- soft.thresholding(u, gamma * lambda)
36
37    # Compute log-likelihood
38    ll[iter] <- log.lik(betas[iter,], y, X, lambda)
39
40    # Convergence check
41    if ( abs(ll[iter-1] - ll[iter]) / (ll[iter-1] + 1E-3) < tol ){
42      cat('Algorithm has converged after', iter, 'iterations')
43      ll <- ll[1:iter]
44      betas <- betas[1:iter, ]
45      break;
46    }
47    else if ( (iter == maxiter) & (abs(ll[iter-1] - ll[iter]) / (ll[iter-1] + 1E-3) >= tol) ){
48      print('WARNING: algorithm has not converged')
49      break;
50    }

```

```

51 }
52 return(list("ll" = ll, "betas" = betas[iter,]))
53 }

```

Listing A.2: Proximal gradient algorithm.

```

1  acc.proximal.gradient <- function(y, X, beta0, lambda, gamma = 1E-4, maxiter = 50000, tol = 1E
   -10){
2    # Accelerated proximal gradient algorithm that allows to compute the Lasso regression estimates
3    # -----
4    # Args:
5    #   - y: response vector (length n)
6    #   - X: matrix of the features (n*p)
7    #   - beta0: initial values of beta
8    #   - lambda: penalization parameter for the lasso (threshold)
9    #   - gamma: proximal operator parameter
10   #   - maxiter: maximum number of iterations
11   #   - tol: tolerance defining convergence
12   # Returns:
13   #   - ll: values of the objective function for each iteration
14   #   - betas: values of the beta parameters for each iteration
15   # -----
16
17   # Check if the inputs have the correct sizes
18   if ((length(y) != dim(X)[1]) || (length(beta0) != dim(X)[2])){
19     stop("Incorrect input dimensions.")
20   }
21
22   # Initialize the data structures
23   betas <- array(NA, dim=c(maxiter, length(beta0)))
24   betas[1,] <- beta0 # Initial guess
25   z <- array(NA, dim=c(maxiter, length(beta0)))
26   z[1,] <- beta0
27   s <- array(NA, dim = maxiter)
28   s[1] <- 1
29   ll <- array(NA, dim = maxiter)
30   ll[1] <- log.lik(betas[1,], y, X, lambda)
31
32   for (iter in 2:maxiter){
33
34     # Compute the gradient
35     gradient <- grad.loglik(z[iter-1, ], y, X)
36     u <- z[iter-1, ] - gamma * gradient
37
38     # Update beta
39     betas[iter, ] <- soft.thresholding(u, gamma * lambda)
40     s[iter] <- (1 + sqrt(1 + 4 * s[iter-1]^2))/2
41     z[iter, ] <- betas[iter, ] + ((s[iter-1] - 1)/s[iter]) * (betas[iter, ] - betas[iter-1, ])
42
43     # Compute log-likelihood
44     ll[iter] <- log.lik(betas[iter,], y, X, lambda)
45
46     # Convergence check
47     if ( abs(ll[iter-1] - ll[iter]) / (ll[iter-1] + 1E-3) < tol ){
48       cat('Algorithm has converged after', iter, 'iterations')
49       ll <- ll[1:iter]
50       betas <- betas[1:iter, ]

```

```
51     break;
52 }
53 else if ( (iter == maxiter) & (abs(l1[iter-1] - l1[iter]) / (l1[iter-1] + 1E-3) >= tol) ){
54     print('WARNING: algorithm has not converged')
55     break;
56 }
57 }
58 return(list("l1" = l1, "betas" = betas[iter,]))
59 }
```

Listing A.3: Accelerated proximal gradient algorithm.