

Exercises 6 The Proximal Gradient Method

Natalia Zuniga-Garcia

October 23, 2017

In this set of exercises, we learn about the proximal gradient algorithm. This is our first all-purpose algorithm capable of handling non-smooth terms that enforce sparsity, like an ℓ^1 regularizer.

1 Proximal operators

First, some definitions. 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

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

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

Intuitively, the Moreau envelope is a regularized version of f . It approximates f from below, and has the same set of minimizing values as f . The proximal mapping returns the value that solves the little minimization problem defined by the Moreau envelope. The objective in this little minimization problem balances two goals: minimizing f , and staying near x . The proximal operator says *where* the minimum occurs, while the Moreau envelope says what the *value* of the minimum is.

Figure 1 shows a simple one-dimensional example of a Moreau envelope and a proximal operator for the absolute-value function.

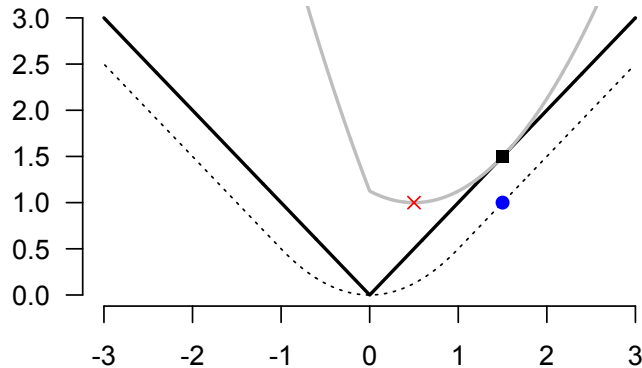


Figure 1: This picture shows a simple example of a proximal operator and Moreau envelope. The solid black line shows the function $f(x) = |x|$, and the dotted line shows the corresponding Moreau envelope $E_\gamma f(x)$ with parameter $\gamma = 1$. The grey line shows the function $|x| + (1/2)(x - x_0)^2$ for $x_0 = 1.5$, whose minimum (shown as a red cross) defines the Moreau envelope and proximal operator. This point has horizontal coordinate $\text{prox}_\gamma f(x_0) = 0.5$ and vertical coordinate $E_\gamma f(x_0) = 1$, and is closer than x_0 to the overall minimum at $x = 0$. The blue circle shows the point $(x_0, E_\gamma f(x_0))$; the vertical coordinate of the blue point is precisely the vertical coordinate of the red point, emphasizing the point-wise construction of the Moreau envelope in terms of a simple optimization problem.

- (A) The proximal operator gives a nice interpretation of classical gradient descent. Consider the local linear approximation of $f(x)$ about a point x_0 :

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

Derive the proximal operator (with parameter γ) of the linear approximation $\hat{f}(x; x_0)$, and show that this proximal operator is identical to a gradient-descent step for $f(x)$ of size γ , starting from the point x_0 . In reflecting on this answer, make sure you feel comfortable with the following statement: “the gradient step minimizes a local linear approximation of the function, subject to a quadratic regularizer that keeps the next iterate close to x_0 (where, presumably, the linear approximation is reasonable).”

Solution

The proximal operator is:

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

Using the local linear approximation of $f(x)$ about a point x_0 :

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

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

$$\begin{aligned} \nabla f(x_0) + \frac{1}{\gamma} (z - x_0) &= 0 \\ z &= x_0 - \gamma \nabla f(x_0) \end{aligned}$$

We can observe that this is the gradient-descent step of size γ .

- (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

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

Show that the proximal operator with parameter $1/\gamma$ of $l(x)$ takes the form

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

assuming the relevant inverse exists.

Solution

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

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

Now we try to minimize by calculation its derivative with respect to z and equalizing it to zero as follow,

$$\begin{aligned} 0 &\stackrel{\text{set}}{=} Pz - q + \gamma(z - x) \\ -Pz - \gamma z &= -q - \gamma x \\ (P + \gamma I)z &= \gamma x + q \end{aligned}$$

Where I is the identity matrix. Assuming that $(P + \gamma I)$ is invertible,

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

Then, we showed that,

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

Then 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. Here A is like our feature matrix, Ω is a covariance matrix, and x is like a regression vector.¹ Specify what P , q , and r are for this negative log likelihood.

Solution

The likelihood of the multivariate Gaussian distribution is:

$$L(y|x) = (2\pi)^{n/2} |\Omega|^{1/2} \exp \left\{ -\frac{1}{2} (y - Ax)^T \Omega (y - Ax) \right\}$$

Then, the negative likelihood function is:

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

Where we obtained:

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

¹This is notation from the optimization literature, rather than the stats literature. It's important to be able to translate between both.

- (C) Let $\phi(x) = \tau\|x\|_1$. Express the proximal operator of this function in terms of the soft-thresholding function that we learned about in the last set of exercises. (An element-wise expression is fine.)

Solution

The proximal operator is:

$$\begin{aligned}\text{prox}_{\gamma} \phi(x) &= \arg \min_z \left\{ \phi(z) + \frac{1}{2\gamma} \|z - x\|_2^2 \right\} = \arg \min_z \left\{ \tau\|x\|_1 + \frac{1}{2\gamma} \|z - x\|_2^2 \right\} \\ &= \arg \min_z \left\{ \tau \sum_{i=1}^p |z_i| + \frac{1}{2\gamma} \|z - x\|_2^2 \right\}\end{aligned}$$

Then we have that for each z_i (element-wise):

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

Which is equal to the soft-thresholding function in Exercises 5:

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

Which we proved to be equal to:

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

Thus, the proximal operator can be written as:

$$\text{prox}_{\gamma} \phi(x_i) = \text{sign}(x_i) \cdot (|x_i| - \gamma \tau)_+$$

2 The proximal gradient method

Suppose that we have some objective function that can be expressed as $f(x) = l(x) + \phi(x)$, where $l(x)$ is differentiable but $\phi(x)$ is not. The proximal gradient method is designed for precisely this situation.

Recall from above the idea of forming a local linear approximation to a function at some point 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'll apply this idea to the first term in our objective, $l(x)$. Define

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

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

$$f(x) \approx \tilde{f}(x; x_0) = \tilde{l}(x; x_0) + \phi(x). \quad (3)$$

- (A) Consider the surrogate optimization problem in which we minimize the approximation $\tilde{f}(x; x_0)$ in Equation 3, in lieu of our original objective $f(x)$.

$$\hat{x} = \arg \min_x \left\{ \tilde{l}(x; x_0) + \phi(x) \right\}.$$

Show that the solution to this problem is of the form

$$\hat{x} = \text{prox}_{\gamma} \phi(u), \quad \text{where } u = x_0 - \gamma \nabla l(x_0). \quad (4)$$

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

Solution

We have that,

$$\hat{x} = \arg \min_x \left\{ \tilde{l}(x; x_0) + \phi(x) \right\}$$

Then,

$$\begin{aligned} \hat{x} &= \arg \min_x \left\{ l(x_0) + (x - x_0)^T \nabla l(x_0) + \frac{1}{2\gamma} \|x - x_0\|_2^2 + \phi(x) \right\} \\ &= \arg \min_x \left\{ x^T \nabla l(x_0) + \frac{1}{2\gamma} \|x - x_0\|_2^2 + \phi(x) \right\} = \arg \min_x \left\{ \frac{1}{2\gamma} \|x - x_0\|_2^2 - \frac{1}{2\gamma} 2x^T (\gamma \nabla l(x_0)) + \phi(x) \right\} \\ &= \arg \min_x \left\{ \frac{1}{2\gamma} [\|x - x_0\|_2^2 - 2x^T (\gamma \nabla l(x_0))] + \phi(x) \right\} \\ &= \arg \min_x \left\{ \frac{1}{2\gamma} [(x - x_0)^T (x - x_0) - 2(x - x_0)^T (\gamma \nabla l(x_0)) + (\gamma \nabla l(x_0))^T (\gamma \nabla l(x_0))] + \phi(x) \right\} \\ &= \arg \min_x \left\{ \frac{1}{2\gamma} (x - x_0 + \gamma \nabla l(x_0))^T (x - x_0 + \gamma \nabla l(x_0)) + \phi(x) \right\} \\ &= \arg \min_x \left\{ \frac{1}{2\gamma} \|x - x_0 + \gamma \nabla l(x_0)\|_2^2 + \phi(x) \right\} = \arg \min_x \left\{ \phi(x) + \frac{1}{2\gamma} \|x - (x_0 - \gamma \nabla l(x_0))\|_2^2 \right\} \end{aligned}$$

The we have that,

$$\begin{aligned} &= \text{prox}_{\gamma} \phi(x_0 - \gamma \nabla l(x_0)) \\ &= \text{prox}_{\gamma} \phi(u) \end{aligned}$$

Where, $u = x_0 - \gamma \nabla l(x_0)$

- (B) The *proximal gradient* method is an iterative algorithm in which we repeatedly form the approximation in Equation 3 about the current point, and minimize this surrogate function using Equation 4. Written concisely,

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

Now consider the lasso regression problem:

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

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

Now implement the method and apply it to the diabetes data from last week. Make sure you track the convergence of the algorithm, i.e. the objective values over time. Compare your answers to the answers you get from the package software you used last week (i.e. `glmnet` or `scikit-learn`).²

Solution

We use proximal gradient to minimize the objective function

$$\tilde{f}(\beta; \beta_0) = \tilde{l}(\beta; \beta_0) + \phi(\beta)$$

Where,

$$\begin{aligned} l(\beta) &= \frac{1}{2n} \|y - X\beta\|_2^2 = \frac{1}{2n} (y - X\beta)^T (y - X\beta) \\ \nabla l(\beta) &= -\frac{1}{n} X^T (y - X\beta) \\ \phi(\beta) &= \lambda \|\beta\|_1 \end{aligned}$$

We have that,

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

Then, the update is:

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

We also have that,

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

Where element-wise we have,

$$\beta_j^{(t+1)} = \text{sign}(u_j^{(t)}) \cdot (|u_j^{(t)}| - \lambda \gamma^{(t)})_+$$

Which corresponds to the soft-thresholding function $S_\lambda(u_j^{(t)})$.

²Keep in mind that those packages are using a value of λ that is rescaled by a factor of n , since they use the loss function

$$\frac{1}{2n} \|y - X\beta\|_2^2 + \lambda \|\beta\|_1.$$

Pseudo-code for using the proximal gradient algorithm:

For i in 2 to iteration

- Estimate the gradient $\nabla l(\beta) = -\frac{1}{n}X^T(y - X\beta)$
- Estimate $u^{(t)} = \beta^{(t)} - \gamma^{(t)}\nabla l(\beta^{(t)})$
- Update $\beta_j^{(t+1)} = S_\lambda(u_j^{(t)})$

The main computational cost is present in obtaining the gradient due to the matrix multiplication.

The method implementation is presented in the file *Part_2.R*,

```

1  # -----
2  # Estimate the Lasso Objective Function
3
4  Lasso_OF <- function(beta, Y, X, lambda){
5    # Compute the Negative Log-Likelihood function
6    # Input:
7    # beta = regression parameters (vector M)
8    # X = features (matrix M x P)
9    # Y = vector of observations (vector M)
10   # lambda = complexity parameter (scalar)
11   # Output:
12   # Lasso_OF = Lasso Objective Function (scalar)
13   t <- X %*% beta
14   OF <- (0.5/nrow(X))*crossprod(Y - t) + lambda * sum(abs(beta))
15   OF <- as.numeric(OF)
16 }
17
18 # -----
19 # Estimate the Gradient of the Negative Log-Likelihood
20
21 Gradient <- function(beta, Y, X){
22   # Estimate the gradient of the Negative Log-Likelihood of the Logit Model
23   # Input:
24   # beta = regression parameters (vector M)
25   # X = features (matrix M x P)
26   # Y = vector of observations (vector M)
27   # Output:
28   # gradient = gradient (vector M)
29   t <- Y - X %*% beta
30   grad <- -(1/nrow(X)) * crossprod(X, t)
31 }
32
33 # -----
34 # Estimate the Soft-Thresholding Function
35 S_lambda <- function(Y, lambda){
36   # Estimate beta using the soft-thresholding operator
37   # Input:
38   # Y = vector of observations (vector M)
39   # lambda = complexity parameter (constant)
40   # Output:
41   # beta = estimated beta vector (vector M)
42   t <- cbind(rep(0, length(Y)), abs(Y) - lambda)

```

```

43 prox <- sign(Y) * apply(t, 1, max)
44 }
45
46 # -----
47 # The Proximal Gradient Method Implementation
48
49 proximal_gradient <- function(Y, X, lambda, gamma, iter, tol){
50   # Estimate Lasso regression using The Proximal Gradient Method
51   # Input:
52   # Y = vector of observations (vector M)
53   # X = features (matrix M x P)
54   # lambda = complexity parameter (constant)
55   # gamma = proximal operator parameter (scalar)
56   # iter = number of iterations (scalar = N)
57   # tol = tolerance
58   # Output:
59   # beta = estimated beta vector (vector P)
60   # Lasso_OF = negative log-likelihood per iteration (vector N)
61
62   P <- dim(X)[2]
63   betas <- array(NA, dim=c(iter, P))
64   betas[1,] <- rep(0, P) # Initial guess is zero
65   OF <- array(NA, dim = iter)
66   OF[1] <- Lasso_OF(betas[1,], Y, X, lambda)
67
68   for (i in 2:iter){
69     gradient <- Gradient(betas[i-1,], Y, X) # Step 1 of my pseudocode
70     u <- betas[i-1,] - gamma * gradient # Step 2 of my pseudocode
71     betas[i,] <- S_lambda(u, gamma * lambda) # Step 3 of my pseudocode
72     OF[i] <- Lasso_OF(betas[i,], Y, X, lambda)
73
74     # Convergence Check
75     e <- abs(OF[i-1] - OF[i]) / (OF[i-1] + tol)
76     if (e < tol){
77       cat('Converged at_', i)
78       OF <- OF[1:i]
79       betas <- betas[1:i, ]
80       break
81     }
82     else if ((i == iter) & (e) >= tol){
83       print('Did not Converge')
84       break
85     }
86   }
87   return(list("Lasso_OF" = OF, "betas" = betas[i,]))
88 }
89
90 # -----

```


- (C) A cool variation on proximal gradient is called the *accelerated* proximal gradient algorithm. The following scheme (due to Nesterov, who's super famous in this area) involves a simple extrapolation step based on the previous iteration:

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

The first step involves the proximal operator of the penalty function, evaluated not at the previous iterate $x^{(t)}$, but at an extrapolated version of $x^{(t)}$, based on the magnitude of the previous step's update. 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?

Note: you can use this acceleration trick in ordinary gradient descent, too.

Solution

Pseudo-code for using the accelerated proximal gradient algorithm:

For i in 2 to iteration

- Estimate the gradient $\nabla l(\beta) = -\frac{1}{n} X^T (y - X\beta)$
- Estimate $s^{(t+1)} = \frac{1 + (1 + 4s_t^2)^{1/2}}{2}$
- Estimate $u^{(t)} = z^{(t)} - \gamma^{(t)} \nabla l(\beta^{(t)})$
- Update $\beta_j^{(t+1)} = S_\lambda(u_j^{(t)})$
- Estimate $z^{(t+1)} = x^{(t+1)} + \left(\frac{s^{(t)} - 1}{s^{(t+1)}} \right) (x^{(t+1)} - x^{(t)})$

The method implementation is presented in the file *Part_2.R*,

```

1  # -----
2  # The Accelerated Proximal Gradient Method Implementation
3
4  ac_proximal_gradient <- function(Y, X, lambda, gamma, iter, tol){
5    # Estimate Lasso regression using The Accelerated Proximal Gradient Method
6    # Input:
7    # Y = vector of observations (vector M)
8    # X = features (matrix M x P)
9    # lambda = complexity parameter (constant)
10   # gamma = proximal operator parameter (scalar)
11   # iter = number of iterations (scalar = N)
12   # tol = tolerance
13   # Output:
14   # beta = estimated beta vector (vector P)
15   # Lasso_OF = negative log-likelihood per iteration (vector N)
16
17   P <- dim(X)[2]
18   betas <- array(NA, dim=c(iter, P))
19   betas[1,] <- rep(0, P) # Initial guess is zero
20   OF <- array(NA, dim = iter)
21   OF[1] <- Lasso_OF(betas[1,], Y, X, lambda)
22   s <- rep(NA, iter)
23   s[1] <- 1
24   z <- rep(0, P)
25
26   for (i in 2:iter){
27     gradient <- Gradient(betas[i-1,], Y, X) # Step 1 of my pseudocode
28     u <- z - gamma * gradient # Step 2 of my pseudocode
29     betas[i,] <- S_lambda(u, gamma * lambda) # Step 3 of my pseudocode
30     s[i] <- (1 + sqrt(1 + 4 * s[i-1]^2))/2 # Step 4 of my pseudocode
31     # Step 5 of my pseudocode
32     z <- betas[i, ] + ((s[i-1] - 1)/s[i]) * (betas[i, ] - betas[i-1, ])
33     OF[i] <- Lasso_OF(betas[i,], Y, X, lambda)
34
35     # Convergence Check
36     e <- abs(OF[i-1] - OF[i]) / (OF[i-1] + tol)
37     if (e < tol){
38       cat('Converged at_', i)
39       OF <- OF[1:i]
40       betas <- betas[1:i, ]
41       break
42     }
43     else if ((i == iter) & (e) >= tol){
44       print('Did not Converge')
45       break
46     }
47   }
48   return(list("Lasso_OF" = OF, "betas" = betas[i,]))
49 }

```

Convergence Comparison

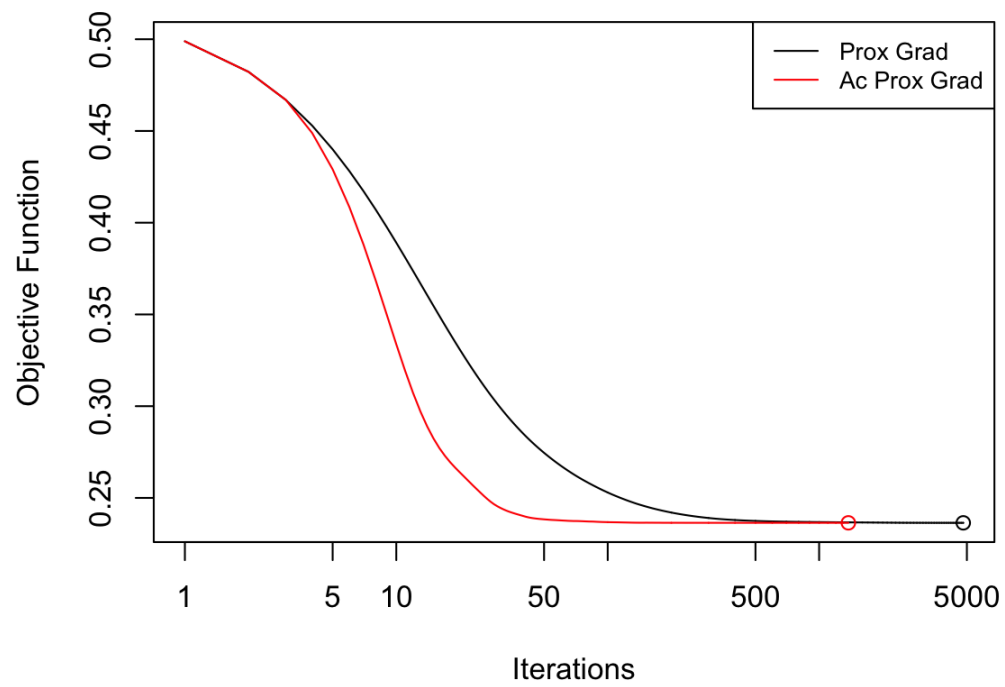


Figure 2: Proximal Gradient Convergency Comparison