

SDS 385: Homework 6

G. Paulon

October 7, 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

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

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

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 $f(x)$ 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 .

Computing the proximal operator of $\hat{f}(x; x_0)$ is equivalent to moving towards the minimum of $\hat{f}(x; x_0)$ staying still close to x_0 . This makes sense, since the problem consists in minimizing a **local** linear approximation and so it requires to stay close to x_0 .

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

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

$$\begin{aligned} & \nabla_x \left\{ f(x_0) + (x - x_0)^T \nabla f(x_0) + \frac{1}{2\gamma} (x - x_0)^T (x - x_0) \right\} \\ &= \nabla f(x_0) + \frac{1}{\gamma} (x - x_0) \end{aligned}$$

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

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

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

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

assuming the relevant inverse exists.

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

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

Therefore the negative log-likelihood is

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

where

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

The proximal operator is then

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

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

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

- (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):

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

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

$$\hat{\mathbf{z}} = \text{prox}_\gamma \phi(\mathbf{x}) = \text{sign}(\mathbf{x}) (|\mathbf{x}| - \tau\gamma \mathbf{1})_+ = S_{\tau\gamma}(\mathbf{x}).$$

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.

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

Show that the solution to this problem is of the form

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

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:

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

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

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

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

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

(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:

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

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.

(C)

Appendix A

R code