

Proximal Methods

Ran Wang

January 30, 2021

Chapter 1

Proximal Algorithms

1.1 Motivation

In this chapter, we briefly outline the problem of minimizing functions that are not necessarily differentiable. A typical example is the l_1 -regularized problem. For example, the object might look like

$$\min_{\beta} \sum_{i=1}^N (y_i - x_i^T \beta)^2 + \lambda \|\beta\|_1.$$

Here, β is the parameter we want to find. Should we not have $\lambda \|\beta\|_1$, then everything is differentiable and can be solved using quasi-Newton methods, among other things. However, the absolute value function is not differentiable everywhere, which causes problems.

The first solution is to consider *sub-differentials*. Sub-differentials are defined as

$$\partial f(x) = \{y \mid f(z) \geq f(x) + y^T(z - x) \text{ for all } z \in \text{dom } f\},$$

where $\text{dom } f$ is the domain of the function. Note that if a function is differentiable then $\partial f = \{\nabla f\}$. However, in general case, the sub-differential is not a singleton.

For simplicity, we assume all the functions we discuss are subdifferentiable.

1.2 Proximal Algorithms

The proximal operator is defined as

$$\text{prox}_f(v) = \underset{x}{\operatorname{argmin}} (f(x) + (1/2)\|x - v\|_2^2).$$

As simple as the definition might look like, it has quite some nice results. The first one is a fixed-point properties. That is, the point x^* minimizes f if and only if

$$x^* = \text{prox}_f(x^*).$$

Proof. First we show that if x^* is the minimizer, then $x^* = \text{prox}_f(x^*)$. Note that for any x ,

$$f(x) + (1/2) \|x - x^*\|_2^2 \geq f(x^*) = f(x^*) + (1/2) \|x^* - x^*\|_2^2,$$

and thus by definition, $x^* = \text{prox}_f(x^*)$.

Now consider the reverse case, let $\tilde{x} = \text{prox}_f(v)$. Take the subdifferential operator, we see that this is equivalent to

$$0 \in \partial f(\tilde{x}) + (\tilde{x} - v).$$

Taking $\tilde{x} = v = x^*$, it follows that $0 \in \partial f(x^*)$, so x^* minimizes f . \square

The second interesting, and rather surprising fact is that, the proximal operator is actually the resolvent of subdifferential operator. More specifically,

$$\text{prox}_{\lambda f} = (I + \lambda \partial f)^{-1}.$$

Proof. If $z \in (I + \lambda \partial f)^{-1}(x)$, then $0 \in \partial f(z) + (1/\lambda)(z - x)$. This implies $0 \in \partial_z (f(z) + (1/2\lambda)\|z - x\|_2^2)$. Now, since we can prove that $f(z) + (1/2\lambda)\|z - x\|_2^2$ is strongly convex, we can deduce that $z = \underset{u}{\text{argmin}} (f(u) + (1/2\lambda)\|u - x\|_2^2)$. \square

Finally, let us look at the case of minimizing $f + g$ where f is differentiable and g is not. A famous algorithms goes,

$$x^{k+1} := \text{prox}_{\lambda^k g}(x^k - \lambda^k \nabla f(x^k)).$$

To see why, consider the fixed point version that is $x^* = \text{prox}_{\lambda g}(x^* - \lambda \nabla f(x^*))$, we show that if this is true then x is indeed the solution to

$$x^* = \underset{x}{\text{argmin}} f(x) + g(x).$$

Proof. Note that x^* is the minimizer if and only if $0 \in \nabla f(x^*) + \partial g(x^*)$. Now with some straight forward computation,

$$\begin{aligned} 0 &\in \lambda \nabla f(x^*) + \lambda \partial g(x^*) \\ 0 &\in \lambda \nabla f(x^*) - x^* + x^* + \lambda \partial g(x^*) \\ (I + \lambda \partial g)(x^*) &\ni (I - \lambda \nabla f)(x^*) \\ x^* &= (I + \lambda \partial g)^{-1}(I - \lambda \nabla f)(x^*) \\ x^* &= \text{prox}_{\lambda g}(x^* - \lambda \nabla f(x^*)) \end{aligned}$$

□

1.3 Applications

Now let us return to the previous topics. Let us assume that the **negative** log-likelihood function is $f_\beta(X)$, here β is the parameter we want to estimate and X is the data, which include both the dependent and independent variables. The goal now is the minimize the following quantities.

$$f_\beta(X) + \lambda \|\beta\|_1,$$

where $\lambda > 0$ is a positive hyperparameter, $\|\cdot\|$ is the l_1 norm, namely for any $x \in R^N$, $\|x\|_1 = \sum_{i=1}^N |x_i|$. Since the latter term is not differentiable at 0, we are in position to workout the solution as is indicated in the previous section.

Now let $w := \beta^k - \lambda^k \nabla f_{\beta^k}(X)$, where λ^k is the step-size, and β^k is the value of β at step k , and let $g : x \mapsto \lambda \|x\|_1$, we have

$$\begin{aligned} \beta^{k+1} &:= \text{prox}_{\lambda^k g}(w) \\ &= \text{argmin}_x \left(\lambda^k g(x) + \frac{1}{2} \|x - w\|_2^2 \right) \\ &= \text{argmin}_x \left(\lambda \|x\|_1 + \frac{1}{2\lambda_k} \|x - w\|_2^2 \right) \end{aligned}$$

Now it suffices to calculate the for each i , since none of i depends on others. For that, note that (**exercises!**) we have

$$\text{argmin}_z \frac{1}{2} \|\beta - z\|_2^2 + \lambda t \|z\|_1 = S_{\lambda t}(\beta)$$

where we have

$$[S_\lambda(\beta)]_i = \begin{cases} \beta_i - \lambda & \text{if } \beta_i > \lambda \\ 0 & \text{if } -\lambda \leq \beta_i \leq \lambda \\ \beta_i + \lambda & \text{if } \beta_i < -\lambda \end{cases}$$

Therefore we have

$$[\beta^{k+1}]_i = \begin{cases} w_i - \lambda \lambda_k & \text{if } w_i > \lambda \lambda_k \\ 0 & \text{if } -\lambda \lambda_k \leq w_i \leq \lambda \lambda_k \\ w_i + \lambda \lambda_k & \text{if } w_i < -\lambda \lambda_k \end{cases}$$

With that, one can easily implement the algorithm.