

Lecture 13. Proximal Gradient Descent

13.1 Proximal operator and proximal gradient descent

So far we only consider minimizing differentiable functions. If the objective function is not differentiable, clearly neither the gradient descent nor the Newton's method works. In this lecture, we focus on how to solve the optimization problem for much more families of convex functions. We will generalize a method of gradient descent named the *proximal gradient descent* for nondifferentiable functions.

Recall the gradient descent iteration $\mathbf{x}_{k+1} = \mathbf{x}_k - \eta \nabla f(\mathbf{x}_k)$, where we let \mathbf{x}_{k+1} be the minimum point of

$$\hat{f}(\mathbf{x}) = f(\mathbf{x}_k) + \langle \nabla f(\mathbf{x}_k), \mathbf{x} - \mathbf{x}_k \rangle + \frac{1}{2\eta} \|\mathbf{x} - \mathbf{x}_k\|^2.$$

Now we assume $f(\mathbf{x})$ is not differentiable, but $f(\mathbf{x})$ can be divided into two parts:

$$f(\mathbf{x}) = g(\mathbf{x}) + h(\mathbf{x})$$

where $g(\mathbf{x})$ is convex and differentiable, and $h(\mathbf{x})$ is convex but not necessarily differentiable. Then we define

$$\hat{g}(\mathbf{x}) = g(\mathbf{x}_k) + \langle \nabla g(\mathbf{x}_k), \mathbf{x} - \mathbf{x}_k \rangle + \frac{1}{2\eta} \|\mathbf{x} - \mathbf{x}_k\|^2$$

to approximate $g(\mathbf{x})$ and let

$$\mathbf{x}_{k+1} = \arg \min_{\mathbf{x}} \hat{g}(\mathbf{x}) + h(\mathbf{x})$$

to approximate the minimum point of $f(\mathbf{x}) = g(\mathbf{x}) + h(\mathbf{x})$.

Note that

$$\begin{aligned} \langle \nabla g(\mathbf{x}_k), \mathbf{x} - \mathbf{x}_k \rangle + \frac{1}{2\eta} \|\mathbf{x} - \mathbf{x}_k\|^2 &= \frac{1}{2\eta} \langle \mathbf{x} - \mathbf{x}_k + \eta \nabla g(\mathbf{x}_k), \mathbf{x} - \mathbf{x}_k + \eta \nabla g(\mathbf{x}_k) \rangle - \frac{\eta}{2} \|\nabla g(\mathbf{x}_k)\|^2 \\ &= \frac{1}{2\eta} \|\mathbf{x} - (\mathbf{x}_k - \eta \nabla g(\mathbf{x}_k))\|^2 - \frac{\eta}{2} \|\nabla g(\mathbf{x}_k)\|^2 \end{aligned}$$

where $\frac{\eta}{2}\|\nabla g(\mathbf{x}_k)\|^2$ is a constant not depending on \mathbf{x} (assuming that \mathbf{x}_k and η are fixed). So we obtain

$$\mathbf{x}_{k+1} = \arg \min_{\mathbf{x}} \hat{g}(\mathbf{x}) + h(\mathbf{x}) = \arg \min_{\mathbf{x}} \frac{1}{2\eta} \|\mathbf{x} - (\mathbf{x}_k - \eta \nabla g(\mathbf{x}_k))\|^2 + h(\mathbf{x}).$$

Here $\mathbf{x}_k - \eta \nabla g(\mathbf{x}_k)$ is the gradient descent iteration if we would like to optimize only $g(\mathbf{x})$. So roughly speaking, after adding $h(\mathbf{x})$, we hope \mathbf{x}_{k+1} locate near the local minimum of g , and not make h large.

Now we define the *proximal operator* as follows.

Definition (Proximal operator)

Given $\mathbf{y} \in \mathbb{R}^n$, let

$$\text{prox}_h(\mathbf{y}) = \arg \min_{\mathbf{x} \in \mathbb{R}^n} \frac{1}{2} \|\mathbf{x} - \mathbf{y}\|^2 + h(\mathbf{x}).$$

Then we can rewrite the proximal gradient descent method as

$$\begin{aligned} \mathbf{x}_{k+1} &= \arg \min_{\mathbf{x}} \frac{1}{2\eta} \|\mathbf{x} - (\mathbf{x}_k - \eta \nabla g(\mathbf{x}_k))\|^2 + h(\mathbf{x}) \\ &= \arg \min_{\mathbf{x}} \frac{1}{2} \|\mathbf{x} - (\mathbf{x}_k - \eta \nabla g(\mathbf{x}_k))\|^2 + \eta h(\mathbf{x}) \\ &= \text{prox}_{\eta h}(\mathbf{x}_k - \eta \nabla g(\mathbf{x}_k)). \end{aligned}$$

Tip

Another viewpoint of the proximal operator is a discretization of the gradient flow. Recall the gradient flow

$$\frac{d}{dt} X(t) = -\nabla f(X(t)).$$

The forward discretization (the *forward Euler method*) is the gradient descent, where

$$\mathbf{x}_{k+1} = \mathbf{x}_k - \eta \nabla f(\mathbf{x}_k).$$

Similarly, we can also try to discretize it in a backward form (the *backward Euler method*):

$$\mathbf{x}_{k+1} = \mathbf{x}_k - \eta \nabla f(\mathbf{x}_{k+1}).$$

However, the iteration becomes difficult since we need to find the point \mathbf{x}_{k+1} satisfying above equations. Actually, this is what the proximal operator is doing. Let

$$\mathbf{x}_{k+1} = \text{prox}_{\eta f}(\mathbf{x}_k) = \arg \min_{\mathbf{x}} \frac{1}{2} \|\mathbf{x} - \mathbf{x}_k\|^2 + \eta f(\mathbf{x}).$$

Then we have $\mathbf{x}_{k+1} - \mathbf{x}_k + \eta \nabla f(\mathbf{x}_{k+1}) = \mathbf{0}$, since

$$\text{LHS} = \nabla \left(\frac{1}{2} \|\mathbf{x} - \mathbf{x}_k\|^2 + \eta f(\mathbf{x}) \right) \Big|_{\mathbf{x}=\mathbf{x}_{k+1}} = \mathbf{0}.$$

13.2 LASSO

The key is that we need to decompose $f(x) = g(x) + h(x)$ properly. Clearly we can set $g(x) \equiv 0$ and $h(x) \equiv f(x)$. But this decomposition is meaningless since we do not know how to compute the proximal operator at all.

Fortunately for some important problems we have a "good" decomposition. For example, we consider the problem of linear regression avoiding overfitting.

Suppose we have a data set $\{(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_m, y_m)\}$ where $\mathbf{x}_i \in \mathbb{R}^n$ and $y_i \in \mathbb{R}$ and we assume that $y_i = \boldsymbol{\beta}^\top \mathbf{x}_i$ for some $\boldsymbol{\beta} \in \mathbb{R}^n$.

Let

$$\mathbf{X} = \begin{pmatrix} \mathbf{x}_1^\top \\ \mathbf{x}_2^\top \\ \vdots \\ \mathbf{x}_m^\top \end{pmatrix} \in \mathbb{R}^{m \times n} \quad \text{and} \quad \mathbf{y} = \begin{pmatrix} y_1 \\ y_2 \\ \vdots \\ y_m \end{pmatrix} \in \mathbb{R}^m.$$

We can compute coefficients $\boldsymbol{\beta}^*$ by solving the following optimization problem (*least square method*)

$$\boldsymbol{\beta}^* = \arg \min_{\boldsymbol{\beta} \in \mathbb{R}^n} \frac{1}{2} \|\mathbf{y} - \mathbf{X}\boldsymbol{\beta}\|^2$$

However, there are usually some redundant coefficients that may cause *overfitting*. So we hope to add some constraints, such as, the number of nonzero entries in $\boldsymbol{\beta}^*$ is at most k . Unfortunately this problem is not a convex optimization after adding this constraint (why?). We still need to approximate it. One idea is to use the

following approximation

$$\boldsymbol{\beta}^* = \arg \min_{\boldsymbol{\beta} \in \mathbb{R}^n} \frac{1}{2} \|\mathbf{y} - \mathbf{X}\boldsymbol{\beta}\|_2^2 + \lambda \|\boldsymbol{\beta}\|_1$$

where $\lambda > 0$ is a parameter. This method is called *LASSO* (*least absolute shrinkage and selection operator*).

The objective function $f(\boldsymbol{\beta}) = \frac{1}{2} \|\mathbf{y} - \mathbf{X}\boldsymbol{\beta}\|_2^2 + \lambda \|\boldsymbol{\beta}\|_1$ has a clear decomposition: $g(\boldsymbol{\beta}) = \frac{1}{2} \|\mathbf{y} - \mathbf{X}\boldsymbol{\beta}\|_2^2$ and $h(\boldsymbol{\beta}) = \lambda \|\boldsymbol{\beta}\|_1$. The advantage is that the proximal operator is easy to compute with respect to this h .

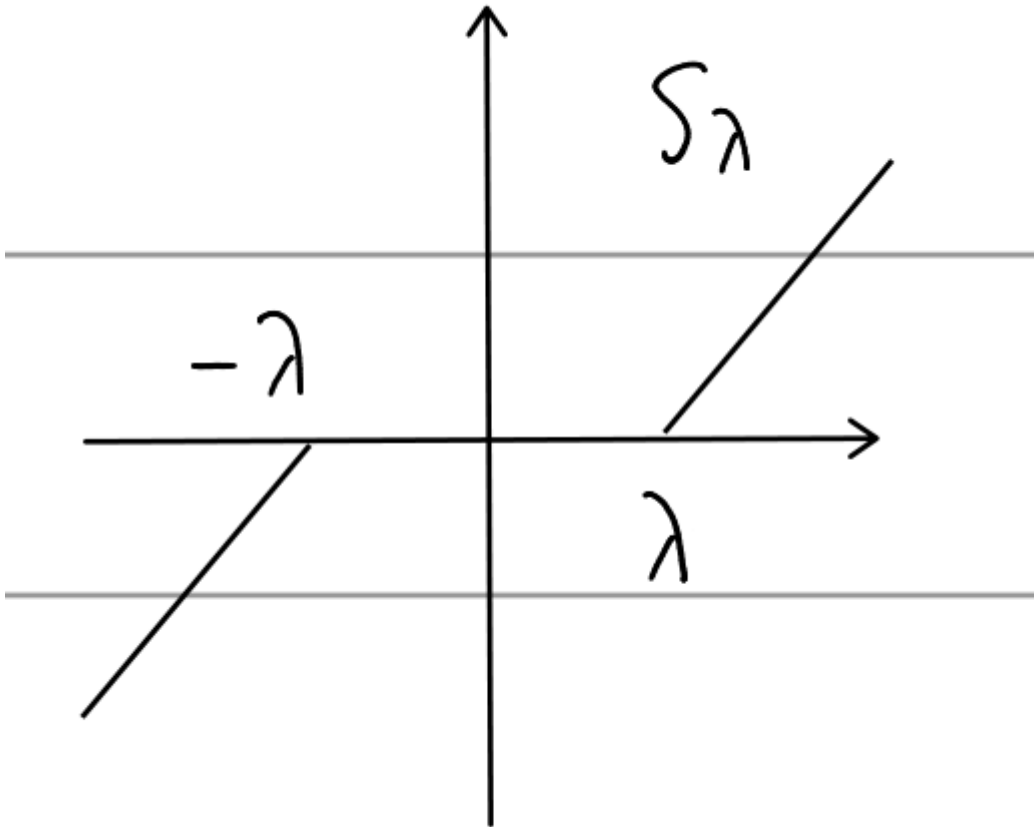
Given $\boldsymbol{\gamma} \in \mathbb{R}^n$, the proximal mapping is

$$\begin{aligned} \text{prox}_h(\boldsymbol{\gamma}) &= \arg \min_{\boldsymbol{\beta} \in \mathbb{R}^n} \frac{1}{2} \|\boldsymbol{\beta} - \boldsymbol{\gamma}\|^2 + \lambda \|\boldsymbol{\beta}\|_1 \\ &= \arg \min_{\boldsymbol{\beta} \in \mathbb{R}^n} \frac{1}{2} ((\beta_1 - \gamma_1)^2 + \dots + (\beta_n - \gamma_n)^2) + \lambda (|\beta_1| + \dots + |\beta_n|). \end{aligned}$$

Note that in this optimization, all entries in $\boldsymbol{\beta}$ are independent, so we can solve this optimization by solving the following optimization for each entry:

$$\begin{aligned} (\text{prox}_h(\boldsymbol{\gamma}))_i &= \arg \min_{\beta_i \in \mathbb{R}} \frac{1}{2} (\beta_i - \gamma_i)^2 + \lambda |\beta_i| \\ &= \arg \min_{\beta_i} \frac{1}{2} (\beta_i - \gamma_i)^2 + \lambda \cdot \begin{cases} \beta_i & \beta_i \geq 0 \\ -\beta_i & \beta_i < 0 \end{cases} \\ &= \arg \min_{\beta_i} \beta_i^2 - 2\beta_i \cdot \begin{cases} \gamma_i - \lambda & \beta_i \geq 0 \\ \gamma_i + \lambda & \beta_i < 0 \end{cases} \\ &= \mathcal{S}_\lambda(\gamma_i) \triangleq \begin{cases} \gamma_i - \lambda & \gamma_i > \lambda \\ 0 & |\gamma_i| \leq \lambda \\ \gamma_i + \lambda & \gamma_i < -\lambda \end{cases} \end{aligned}$$

where \mathcal{S}_λ is called the soft thresholding operator.



Finally, recall that $g(\boldsymbol{\beta}) = \frac{1}{2} \|\mathbf{y} - \mathbf{X}\boldsymbol{\beta}\|_2^2$, so the iteration of the proximal gradient descent is

$$\boldsymbol{\beta}_{k+1} = \text{prox}_{\eta h}(\boldsymbol{\beta}_k - \eta \nabla g(\boldsymbol{\beta}_k)) = \mathcal{S}_{\lambda\eta}(\boldsymbol{\beta}_k - \eta(-\mathbf{X}^\top(\mathbf{y} - \mathbf{X}\boldsymbol{\beta}_k))).$$

This algorithm is called the *ISTA* (*iterative soft-thresholding algorithm*).

13.3 Correctness and convergence

We now show the correctness and convergence of the proximal gradient descent.

Assume that g is L -smooth and set $\eta \leq 1/L$. Sometimes we would write the proximal gradient descent as the following form

$$\mathbf{x}_{k+1} = \mathbf{x}_k - \eta G_\eta(\mathbf{x}_k),$$

where

$$G_\eta(\mathbf{x}) = \frac{\mathbf{x} - \text{prox}_{\eta h}(\mathbf{x} - \eta \nabla g(\mathbf{x}))}{\eta}.$$

We first show that, $G_\eta(\mathbf{x}_k) = \mathbf{0}$ if and only if \mathbf{x}_k is a minimum point of $f(\mathbf{x})$. A trivial example is $h \equiv 0$. Then we have $\text{prox}_{\eta h}(\mathbf{x}) = \mathbf{x}$ and thus $G_\eta(\mathbf{x}) = \nabla g(\mathbf{x})$. So $G_\eta(\mathbf{x}) = \mathbf{0}$ if and only if \mathbf{x} is a minimizer of $f = g$. The following theorem asserts the general cases.

Theorem

In the proximal gradient descent iterations, $\mathbf{x}_{k+1} = \mathbf{x}_k$ iff $f(\mathbf{x}_k) = f^*$, where f^* is the minimum value of f .

The "only if" direction is easy, since we have

$$\hat{f}(\mathbf{x}) = \hat{g}(\mathbf{x}) + h(\mathbf{x}) = g(\mathbf{x}_k) + \langle \nabla g(\mathbf{x}_k), \mathbf{x} - \mathbf{x}_k \rangle + \frac{1}{2\eta} \|\mathbf{x} - \mathbf{x}_k\|^2 + h(\mathbf{x}) \geq g(\mathbf{x}) + h(\mathbf{x}) =$$

by L -smoothness of g . It yields that

$$f(\mathbf{x}_{k+1}) \leq \hat{f}(\mathbf{x}_{k+1}) = \min_{\mathbf{x}} \hat{f}(\mathbf{x}) \leq \hat{f}(\mathbf{x}_k) = f(\mathbf{x}_k).$$

If $f(\mathbf{x}_k) = f^*$, all inequalities are tight, so $\hat{f}(\mathbf{x}_k) = \min \hat{f}(\mathbf{x})$. However, because h is convex and $\|\mathbf{x} - \mathbf{x}_k\|^2$ is strictly convex, \hat{f} is also strictly convex, and thus has a unique minimizer \mathbf{x}_{k+1} , which gives that $\mathbf{x}_k = \mathbf{x}_{k+1}$.

Now we would like to show that if $\mathbf{x}_k = \mathbf{x}_{k+1}$ then $f(\mathbf{x}_k) = f^*$. In other words, for all $\mathbf{x} \in \mathbb{R}^n$, $f(\mathbf{x}) \geq f(\mathbf{x}_k)$. How can we show this? A naive idea is to show that $g(\mathbf{x}) \geq g(\mathbf{x}_k)$ and $h(\mathbf{x}) \geq h(\mathbf{x}_k)$. However this idea looks too good to be true.

A reasonable method is to use convexity. Note that

$g(\mathbf{x}) \geq g(\mathbf{x}_k) + \langle \nabla g(\mathbf{x}_k), \mathbf{x} - \mathbf{x}_k \rangle$. So we hope

$$h(\mathbf{x}) \geq h(\mathbf{x}_k) - \langle \nabla g(\mathbf{x}_k), \mathbf{x} - \mathbf{x}_k \rangle.$$

This inequality looks like the first order condition for h . Here $-\nabla g$ performs like the gradient of h . Since h is not differentiable, we would introduce the notion of *subgradients*.

Definition

Let $f : \mathbb{R}^n \rightarrow \mathbb{R}$ be a convex function and $\mathbf{x} \in \mathbb{R}^n$. We say $\mathbf{v} \in \mathbb{R}^n$ is a *subgradient* of f at \mathbf{x} , denoted by $\mathbf{v} \in \partial f(\mathbf{x})$, if $\forall \mathbf{y} \in \mathbb{R}^n$,
 $f(\mathbf{y}) \geq f(\mathbf{x}) + \langle \mathbf{v}, \mathbf{y} - \mathbf{x} \rangle$.

If f is convex, then subgradients always exist (but may not be unique). Just consider the supporting hyperplane of the epigraph of f .

Lemma

If $\mathbf{z} = \text{prox}_{\eta h}(\mathbf{y})$, then $\frac{\mathbf{y} - \mathbf{z}}{\eta} \in \partial h(\mathbf{z})$.

This lemma immediately implies that if $\mathbf{x}_k = \mathbf{x}_{k+1} = \text{prox}_{\eta h}(\mathbf{x}_k - \eta \nabla g(\mathbf{x}_k))$, then $-\nabla g(\mathbf{x}_k) \in \partial h(\mathbf{x}_k)$. Thus

$$\forall \mathbf{x} \in \mathbb{R}^n, \quad f(\mathbf{x}) = g(\mathbf{x}) + h(\mathbf{x}) \geq g(\mathbf{x}_k) + \langle \nabla g(\mathbf{x}_k), \mathbf{x} - \mathbf{x}_k \rangle + h(\mathbf{x}_k) - \langle \nabla g(\mathbf{x}_k), \mathbf{x} - \mathbf{x}_k \rangle$$

Proof of the Lemma

By the definition of \mathbf{z} , we have that for all $\mathbf{x} \in \mathbb{R}^n$,

$$\frac{1}{2\eta} \|\mathbf{x} - \mathbf{y}\|^2 + h(\mathbf{x}) \geq \frac{1}{2\eta} \|\mathbf{z} - \mathbf{y}\|^2 + h(\mathbf{z}).$$

Our goal is to show that for all $\mathbf{x} \in \mathbb{R}^n$, $h(\mathbf{x}) \geq h(\mathbf{z}) + \frac{1}{\eta} \langle \mathbf{y} - \mathbf{z}, \mathbf{x} - \mathbf{z} \rangle$.

For the sake of contradiction, assume that there exists $\mathbf{w} \in \mathbb{R}^n$ such that $h(\mathbf{w}) < h(\mathbf{z}) + \frac{1}{\eta} \langle \mathbf{y} - \mathbf{z}, \mathbf{w} - \mathbf{z} \rangle$. Define

$$\delta = \frac{h(\mathbf{z}) - h(\mathbf{w}) + \frac{1}{\eta} \langle \mathbf{y} - \mathbf{z}, \mathbf{w} - \mathbf{z} \rangle}{\|\mathbf{w} - \mathbf{z}\|} > 0.$$

Then $h(\mathbf{w}) = h(\mathbf{z}) + \frac{1}{\eta} \langle \mathbf{y} - \mathbf{z}, \mathbf{w} - \mathbf{z} \rangle - \delta \|\mathbf{w} - \mathbf{z}\|$.

Let $\theta \in (0, 1)$ and $\mathbf{x} = \theta \mathbf{w} + (1 - \theta) \mathbf{z}$. By convexity we have

$$\begin{aligned} h(\mathbf{x}) &\leq \theta h(\mathbf{w}) + (1 - \theta) h(\mathbf{z}) \\ &= h(\mathbf{z}) + \frac{\theta}{\eta} \langle \mathbf{y} - \mathbf{z}, \mathbf{w} - \mathbf{z} \rangle - \theta \delta \|\mathbf{w} - \mathbf{z}\| \\ &= h(\mathbf{z}) + \frac{1}{\eta} \langle \mathbf{y} - \mathbf{z}, \mathbf{x} - \mathbf{z} \rangle - \delta \|\mathbf{x} - \mathbf{z}\|. \end{aligned}$$

Since

$$\|\mathbf{x} - \mathbf{y}\|^2 = \|\mathbf{x} - \mathbf{z}\|^2 + \|\mathbf{z} - \mathbf{y}\|^2 + 2 \langle \mathbf{x} - \mathbf{z}, \mathbf{z} - \mathbf{y} \rangle,$$

it follows that

$$\begin{aligned}\frac{1}{2\eta} \|\mathbf{x} - \mathbf{y}\|^2 + h(\mathbf{x}) &= \frac{1}{2\eta} \|\mathbf{x} - \mathbf{z}\|^2 + \frac{1}{2\eta} \|\mathbf{y} - \mathbf{z}\|^2 + \frac{1}{\eta} \langle \mathbf{x} - \mathbf{z}, \mathbf{z} - \mathbf{y} \rangle + h(\mathbf{x}) \\ &\leq \frac{1}{2\eta} \|\mathbf{x} - \mathbf{z}\|^2 + \frac{1}{2\eta} \|\mathbf{y} - \mathbf{z}\|^2 + h(\mathbf{z}) - \delta \|\mathbf{x} - \mathbf{z}\| \\ &< \frac{1}{2\eta} \|\mathbf{z} - \mathbf{y}\|^2 + h(\mathbf{z})\end{aligned}$$

if $\|\mathbf{x} - \mathbf{z}\|$ is sufficiently small ($< 2\delta\eta$). This contradicts to the definition of \mathbf{z} .

Finally, we present the results of convergence rate.

Theorem

Suppose g is L -smooth and we set $\eta \leq 1/L$. Then we have

$$f(\mathbf{x}_T) - f^* \leq \frac{\|\mathbf{x}_0 - \mathbf{x}^*\|^2}{2T\eta}.$$

If g is further μ -strongly convex, then

$$\|\mathbf{x}_{k+1} - \mathbf{x}^*\|^2 \leq (1 - \mu\eta) \|\mathbf{x}_k - \mathbf{x}^*\|^2$$

If L -smoothness is unknown, we can use the exact/backtracking line search, and the results to the convergence rate are also the same as the rate of gradient descent.