

CS 2601 Linear and Convex Optimization

8. Proximal gradient descent

Bo Jiang

John Hopcroft Center for Computer Science
Shanghai Jiao Tong University

Fall 2022

Outline

- Algorithm and examples
- Convergence analysis

Composite functions

Consider

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

where

- f is convex and smooth
- h is convex but not necessarily smooth

Example. Model with ℓ_1 regularization to promote sparsity,

$$F(\mathbf{x}) = f(\mathbf{x}) + \lambda \|\mathbf{x}\|_1$$

A more concrete example is Lasso in penalized form,

$$F(\mathbf{w}) = \frac{1}{2} \|\mathbf{y} - \mathbf{X}\mathbf{w}\|_2^2 + \lambda \|\mathbf{w}\|_1$$

Proximal gradient descent

Recall for differentiable F , the update rule for gradient descent is

$$\begin{aligned}\mathbf{x}_{k+1} &= \operatorname{argmin}_{\mathbf{x}} \left\{ \underbrace{F(\mathbf{x}_k) + \nabla F(\mathbf{x}_k)^T(\mathbf{x} - \mathbf{x}_k) + \frac{1}{2t_k} \|\mathbf{x} - \mathbf{x}_k\|_2^2}_{\hat{F}(\mathbf{x})} \right\} \\ &= \operatorname{argmin}_{\mathbf{x}} \frac{1}{2t_k} \|\mathbf{x} - (\mathbf{x}_k - t_k \nabla F(\mathbf{x}_k))\|_2^2\end{aligned}$$

加上常数凑平方

Proximal gradient descent uses a similar approximation, **but only for the smooth part f** ,

$$\begin{aligned}\mathbf{x}_{k+1} &= \operatorname{argmin}_{\mathbf{x}} \left\{ \underbrace{f(\mathbf{x}_k) + \nabla f(\mathbf{x}_k)^T(\mathbf{x} - \mathbf{x}_k) + \frac{1}{2t_k} \|\mathbf{x} - \mathbf{x}_k\|_2^2}_{\hat{f}(\mathbf{x})} + h(\mathbf{x}) \right\} \\ &= \operatorname{argmin}_{\mathbf{x}} \left\{ \frac{1}{2} \|\mathbf{x} - (\mathbf{x}_k - t_k \nabla f(\mathbf{x}_k))\|_2^2 + t_k h(\mathbf{x}) \right\}\end{aligned}$$

近邻梯度下降

“Proximal” means we try to make \mathbf{x}_{k+1} stay close to $\mathbf{x}_k - t_k \nabla f(\mathbf{x}_k)$.

Proximal gradient descent (cont'd)

Define the proximal mapping or proximal operator,

$$\text{prox}_h(\mathbf{y}) = \underset{\mathbf{x}}{\operatorname{argmin}} \left\{ \frac{1}{2} \|\mathbf{x} - \mathbf{y}\|_2^2 + h(\mathbf{x}) \right\}$$

so

$$\mathbf{x}_{k+1} = \text{prox}_{t_k h}(\mathbf{x}_k - t_k \nabla f(\mathbf{x}_k))$$

Proximal gradient descent with constant step size

- 1: initialization $\mathbf{x} \leftarrow \mathbf{x}_0 \in \mathbb{R}^n$
- 2: **while** stopping criterion not satisfied **do**
- 3: $\mathbf{x} \leftarrow \text{prox}_{th}(\mathbf{x} - t \nabla f(\mathbf{x}))$
- 4: **end while**
- 5: **return** \mathbf{x}

Note. The proximal operator involves another optimization problem! Fortunately, it depends only on h not on f , and we can compute it in closed form for many important h , e.g. ℓ_1 regularization.

Proximal operator for ℓ_2 regularization

When $h(\mathbf{x}) = \frac{\lambda}{2} \|\mathbf{x}\|_2^2$, the proximal operator is

$$\text{prox}_h(\mathbf{y}) = \underset{\mathbf{x}}{\text{argmin}} \left\{ \overset{\text{fixed part}}{\frac{1}{2} \|\mathbf{x} - \mathbf{y}\|_2^2} + \overset{\text{variant}}{\frac{\lambda}{2} \|\mathbf{x}\|_2^2} \right\} = \frac{\mathbf{y}}{1 + \lambda}$$

Note. ℓ_2 regularization **does not promote sparsity.**

For ℓ_2 regularization, proximal gradient descent turns out to be nothing but gradient descent.

The proximal gradient step with step size t is

$$\mathbf{x}_{k+1} = \frac{\mathbf{x}_k - t \nabla f(\mathbf{x}_k)}{1 + \lambda t} = \mathbf{x}_k - \frac{t}{1 + \lambda t} (\nabla f(\mathbf{x}_k) + \lambda \mathbf{x}_k) = \mathbf{x}_k - \frac{t}{1 + \lambda t} \nabla F(\mathbf{x}_k)$$

exactly the gradient step for $F(\mathbf{x}) = f(\mathbf{x}) + h(\mathbf{x})$ **with step size $\frac{t}{1 + \lambda t}$!**

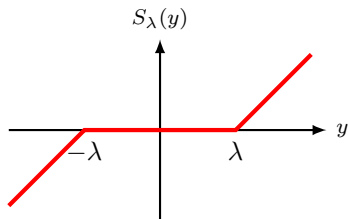
Proximal operator for ℓ_1 regularization

First consider the 1D case, $h(x) = \lambda|x|$ with $\lambda \geq 0$. Need to solve

$$\min_x \quad \frac{1}{2}(x - y)^2 + \lambda|x|$$

The proximal operator is given by the soft-thresholding operator $S_\lambda(y)$,

$$\text{prox}_{\lambda|\cdot|}(y) = S_\lambda(y) = \text{sgn}(y)(|y| - \lambda)^+ = \begin{cases} y - \lambda, & \text{if } y > \lambda \\ 0, & \text{if } -\lambda \leq y \leq \lambda \\ y + \lambda, & \text{if } y < -\lambda \end{cases}$$



Proximal operator for ℓ_1 regularization (cont'd)

Proof. Let $x^* = \text{prox}_h(y)$ to be the minimum of $\frac{1}{2}(x - y)^2 + \lambda|x|$. Then

$$\frac{1}{2}(x^* - y)^2 + \lambda|x^*| \leq \frac{1}{2}(-x^* - y)^2 + \lambda|-x^*| \implies yx^* \geq 0$$

i.e. x^* must have the same sign as y . If $y \geq 0$, then $x^* \geq 0$, so

$$x^* = \operatorname{argmin}_{x \geq 0} \left\{ \frac{1}{2}(x - y)^2 + \lambda x \right\} = \begin{cases} y - \lambda, & \text{if } y > \lambda \\ 0, & \text{if } 0 \leq y \leq \lambda \end{cases}$$

Similarly, if $y \leq 0$, then $x^* \leq 0$, so

$$x^* = \operatorname{argmin}_{x \leq 0} \left\{ \frac{1}{2}(x - y)^2 - \lambda x \right\} = \begin{cases} y + \lambda, & \text{if } y < -\lambda \\ 0, & \text{if } -\lambda \leq y \leq 0 \end{cases}$$

Combining the two cases completes the proof.

Proximal operator for ℓ_1 regularization (cont'd)

When $h(\mathbf{x}) = \lambda \|\mathbf{x}\|_1$, the minimization in the proximal operator can be decomposed into n 1D problems, one for each component,

$$\frac{1}{2} \|\mathbf{x} - \mathbf{y}\|_2^2 + \lambda \|\mathbf{x}\|_1 = \sum_{i=1}^n \left\{ \frac{1}{2} (x_i - y_i)^2 + \lambda |x_i| \right\}$$

The proximal operator simply applies the soft-thresholding operator to each component of \mathbf{y} ,

$$\left[\text{prox}_{\lambda \|\cdot\|_1}(\mathbf{y}) \right]_i = [S_\lambda(\mathbf{y})]_i = S_\lambda(y_i) = \begin{cases} y_i - \lambda, & \text{if } y_i > \lambda \\ 0, & \text{if } |y_i| \leq \lambda \\ y_i + \lambda, & \text{if } y_i < -\lambda \end{cases}$$

Note. The soft-thresholding operation gives us some intuition about how ℓ_1 regularization promotes sparsity: small components are set to zero during the optimization process.

Lasso and ISTA

For Lasso in the penalized form,

$$\min_{\mathbf{w}} F(\mathbf{w}) = \frac{1}{2} \|\mathbf{X}\mathbf{w} - \mathbf{y}\|_2^2 + \lambda \|\mathbf{w}\|_1$$

the smooth part $f(\mathbf{w}) = \frac{1}{2} \|\mathbf{X}\mathbf{w} - \mathbf{y}\|_2^2$ has gradient

$$\nabla f(\mathbf{w}) = \mathbf{X}^T(\mathbf{X}\mathbf{w} - \mathbf{y})$$

so the proximal gradient step is

$$\mathbf{w}_{k+1} = S_{\lambda t}(\mathbf{w}_k - t\mathbf{X}^T(\mathbf{X}\mathbf{w}_k - \mathbf{y}))$$

Known as the **iterative soft-thresholding algorithm (ISTA)**.

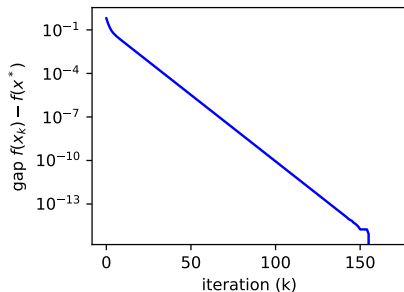
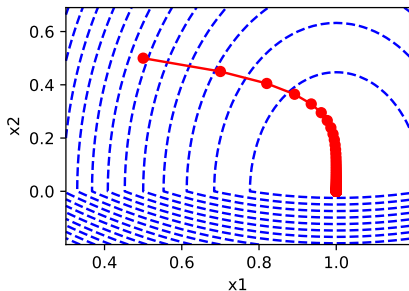
Note. The threshold is λt , not λ .

Lasso and ISTA (cont'd)

Example.

$$\mathbf{X} = \begin{bmatrix} 2 & 0 \\ 0 & 1 \\ 0 & 0 \end{bmatrix}, \quad \mathbf{y} = \begin{bmatrix} 3 \\ 2 \\ 2 \end{bmatrix}, \quad \lambda = 2$$

Step size $t = 0.1$, $\mathbf{w}_0 = (0.5, 0.5)^T$, $\mathbf{w}^* = (1, 0)^T$.



Outline

- Algorithm and examples
- Convergence analysis

Lower bound for strongly convex functions

Lemma. If $\varphi : X \rightarrow \mathbb{R}$ is μ -strongly convex with a minimum \mathbf{x}^* , then

$$\varphi(\mathbf{x}) \geq \varphi(\mathbf{x}^*) + \frac{\mu}{2} \|\mathbf{x} - \mathbf{x}^*\|_2^2, \quad \forall \mathbf{x} \in X$$

Note. For φ is differentiable, this will follow from the quadratic lower bound $\varphi(\mathbf{x}) \geq \varphi(\mathbf{y}) + \nabla \varphi(\mathbf{y})^T (\mathbf{x} - \mathbf{y}) + \frac{\mu}{2} \|\mathbf{x} - \mathbf{y}\|_2^2$ and the optimality condition $\nabla \varphi(\mathbf{x}^*)^T (\mathbf{x} - \mathbf{x}^*) \geq 0$.

如果连续可以直接用

upperbound证明, 但是不需要

Proof. Since φ is μ -strongly convex, there exists a convex $\tilde{\varphi}(\mathbf{x})$ s.t.

$$\varphi(\mathbf{x}) = \tilde{\varphi}(\mathbf{x}) + \frac{\mu}{2} \|\mathbf{x} - \mathbf{x}^*\|^2$$

后项展开, 不会影响 convexity

Fix \mathbf{x} and let $\mathbf{x}_t = t\mathbf{x} + \bar{t}\mathbf{x}^* = \mathbf{x}^* + t(\mathbf{x} - \mathbf{x}^*)$, $t \in [0, 1]$. By convexity of $\tilde{\varphi}$,

$$\varphi(\mathbf{x}^*) \leq \varphi(\mathbf{x}_t) = \tilde{\varphi}(\mathbf{x}_t) + \frac{\mu}{2} \|\mathbf{x}_t - \mathbf{x}^*\|_2^2 \leq t\tilde{\varphi}(\mathbf{x}) + \bar{t}\tilde{\varphi}(\mathbf{x}^*) + \frac{\mu t^2}{2} \|\mathbf{x} - \mathbf{x}^*\|_2^2$$

Since $\tilde{\varphi}(\mathbf{x}^*) = \varphi(\mathbf{x}^*)$,

$$\tilde{\varphi}(\mathbf{x}) \geq \varphi(\mathbf{x}^*) - \frac{\mu t}{2} \|\mathbf{x} - \mathbf{x}^*\|_2^2, \quad \forall t \in [0, 1] \implies \tilde{\varphi}(\mathbf{x}) \geq \varphi(\mathbf{x}^*)$$

Convergence analysis

Theorem. Let $F(\mathbf{x}) = f(\mathbf{x}) + h(\mathbf{x})$, where h is convex, f is L -smooth and m -strongly convex with $m \geq 0$. Let \mathbf{x}^* be the minimum of F . The sequence $\{\mathbf{x}_k\}$ produced by proximal gradient descent with constant step size $t = \frac{1}{L}$ has the following properties.

1. $F(\mathbf{x}_{k+1}) \leq F(\mathbf{x}_k)$ and

$$F(\mathbf{x}_k) - F(\mathbf{x}^*) \leq \frac{L}{2k} \|\mathbf{x}^* - \mathbf{x}_0\|_2^2$$

2. $\|\mathbf{x}_{k+1} - \mathbf{x}^*\|_2^2 \leq (1 - \frac{m}{L}) \|\mathbf{x}_k - \mathbf{x}^*\|_2^2$, and hence

$$\|\mathbf{x}_k - \mathbf{x}^*\|_2^2 \leq (1 - \frac{m}{L})^k \|\mathbf{x}_0 - \mathbf{x}^*\|_2^2$$

Note. If $t \in (0, \frac{1}{L}]$, then the conclusions hold with L replaced by $\frac{1}{t}$, since an L -smooth function is also $\frac{1}{t}$ -smooth.

Note. 2 implies $F(\mathbf{x}_k) - F(\mathbf{x}^*) \leq C(1 - \frac{m}{L})^{\frac{k}{2}} \|\mathbf{x}^* - \mathbf{x}_0\|_2$ for C -Lipschitz continuous F .

Proof

Let

skipped

$$\widehat{F}(\mathbf{x}) = f(\mathbf{x}_k) + \nabla f(\mathbf{x}_k)^T(\mathbf{x} - \mathbf{x}_k) + \frac{L}{2}\|\mathbf{x} - \mathbf{x}_k\|_2^2 + h(\mathbf{x})$$

Since f is L -smooth and m -strongly convex,

$$\frac{m}{2}\|\mathbf{x} - \mathbf{x}_k\|_2^2 \leq f(\mathbf{x}) - f(\mathbf{x}_k) - \nabla f(\mathbf{x}_k)^T(\mathbf{x} - \mathbf{x}_k) \leq \frac{L}{2}\|\mathbf{x} - \mathbf{x}_k\|_2^2$$

Plugging into $\widehat{F}(\mathbf{x})$ yields

$$F(\mathbf{x}) \leq \widehat{F}(\mathbf{x}) \leq F(\mathbf{x}) + \frac{L - m}{2}\|\mathbf{x} - \mathbf{x}_k\|_2^2$$

Note $\widehat{F}(\mathbf{x})$ is L -strongly convex, and $\mathbf{x}_{k+1} = \operatorname{argmin} \widehat{F}(\mathbf{x})$ when $t = \frac{1}{L}$.

$$\begin{aligned} F(\mathbf{x}_{k+1}) &\leq \widehat{F}(\mathbf{x}_{k+1}) \stackrel{(\star)}{\leq} \widehat{F}(\mathbf{x}) - \frac{L}{2}\|\mathbf{x} - \mathbf{x}_{k+1}\|_2^2 \\ &\leq F(\mathbf{x}) + \frac{L - m}{2}\|\mathbf{x} - \mathbf{x}_k\|_2^2 - \frac{L}{2}\|\mathbf{x} - \mathbf{x}_{k+1}\|_2^2 \end{aligned}$$

where (\star) uses the previous lemma.

Proof (cont'd)

The previous slide shows

$$F(\mathbf{x}_{k+1}) \leq F(\mathbf{x}) + \frac{L-m}{2} \|\mathbf{x} - \mathbf{x}_k\|_2^2 - \frac{L}{2} \|\mathbf{x} - \mathbf{x}_{k+1}\|_2^2$$

1. Setting $\mathbf{x} = \mathbf{x}_k$ shows

$$F(\mathbf{x}_{k+1}) \leq F(\mathbf{x}_k) - \frac{L}{2} \|\mathbf{x}_k - \mathbf{x}_{k+1}\|_2^2 \leq F(\mathbf{x}_k)$$

Since $m \geq 0$, setting $\mathbf{x} = \mathbf{x}^*$ shows

$$F(\mathbf{x}_{i+1}) - F(\mathbf{x}^*) \leq \frac{L}{2} \|\mathbf{x}^* - \mathbf{x}_i\|_2^2 - \frac{L}{2} \|\mathbf{x}^* - \mathbf{x}_{i+1}\|_2^2$$

Since $F(\mathbf{x}_{i+1}) \leq F(\mathbf{x}_i)$ by part 1, summing over i from 0 to $k-1$,

$$\begin{aligned} F(\mathbf{x}_k) - F(\mathbf{x}^*) &\leq \frac{1}{k} \sum_{i=0}^{k-1} [F(\mathbf{x}_i) - F(\mathbf{x}^*)] \\ &\leq \frac{L}{2k} \|\mathbf{x}^* - \mathbf{x}_0\|_2^2 - \frac{L}{2k} \|\mathbf{x}^* - \mathbf{x}_k\|_2^2 \leq \frac{L}{2k} \|\mathbf{x}^* - \mathbf{x}_0\|_2^2 \end{aligned}$$

Proof (cont'd)

2. Setting $\mathbf{x} = \mathbf{x}^*$ shows

$$\begin{aligned}\|\mathbf{x}^* - \mathbf{x}_{k+1}\|_2^2 &\leq \left(1 - \frac{m}{L}\right)\|\mathbf{x}^* - \mathbf{x}_k\|_2^2 - \frac{2}{L}\underbrace{[F(\mathbf{x}_{k+1}) - F(\mathbf{x}^*)]}_{\geq 0} \\ &\leq \left(1 - \frac{m}{L}\right)\|\mathbf{x}^* - \mathbf{x}_k\|_2^2\end{aligned}$$

so

$$\|\mathbf{x}^* - \mathbf{x}_k\|_2^2 \leq \left(1 - \frac{m}{L}\right)^k \|\mathbf{x}^* - \mathbf{x}_0\|_2^2$$