

Convex Optimisation

Matthew Tam
(University of Melbourne)
matthew.tam@unimelb.edu.au

ACE Network

81 2021

* Assignment 4

- online after lecture today
- due June 3rd @ 5pm

Last few weeks

Frank-Wolfe

$$\min_{x \in C} f(x)$$

differentiable
compact

Proximal gradient descent

$$\min_{x \in \mathbb{R}^n} f(x) + g(x)$$

proper, lsc, convex
differentiable

- FISTA = accelerated version
- activated f via Df
- activated g via prox g

$$x_{k+1} = \text{prox}_{\lambda g}(x_k - \lambda Df(x_k)) \quad \text{THEN.}$$

3.4 Subgradient method

In this section, we consider the problem

$$\min_{x \in \mathbb{R}^n} f(x)$$

where $f: \mathbb{R}^n \rightarrow \mathbb{R}$ is a convex function, but not necessarily differentiable.

If f were differentiable, then we could apply the method from Section 3.2 with $g=0$ to obtain algorithm

$$x_{k+1} = x_k - \lambda \nabla f(x_k).$$

replace this
with $\delta f(x_k)$.

Since our function is not differentiable, we will replace the gradient with subgradients in the above iteration.

$$|x|$$

Example 3.4-1

✓ $x^* = 0$

Consider $\min_{x \in \mathbb{R}} f(x)$ where $f(x) = |x|$ which has solution $x^* = 0$. Then

$$\partial f(x) = \begin{cases} \{+1\} & x > 0 \\ \{-1\} & x < 0 \\ (-1, +1] & x = 0 \end{cases}$$

Let $\lambda > 0$ and $x_0 \in (\frac{1}{3}\lambda, \frac{2}{3}\lambda)$. Consider the iteration

$$x_{n+1} = x_n - \lambda \phi_k \text{ where } \phi_k \in \partial f(x_n)$$

Then

$$x_1 = \begin{cases} x_0 - \lambda(+1) < \frac{2}{3}\lambda - \lambda = -\frac{1}{3}\lambda \\ x_0 - \lambda(+1) > \frac{1}{3}\lambda - \lambda = -\frac{2}{3}\lambda \end{cases}$$

$$\Rightarrow x_1 \in \left(-\frac{2}{3}\lambda, -\frac{1}{3}\lambda\right).$$

By symmetry, it follows that $x_2 \in (\frac{1}{3}\lambda, \frac{2}{3}\lambda)$.

Thus, $|x_k| > \frac{1}{3}\lambda \quad \forall k \in \mathbb{N}$. and

So (x_n) cannot converge to $x^* = 0$. \square

Given an initial $x_0 \in \mathbb{R}^n$, we consider the subgradient method which computes (x_n) according to

$$x_{n+1} = x_n - \lambda_n \cdot \frac{\phi_n}{\|\phi_n\|} \quad (3.18)$$

where $\phi_n \in \partial f(x_n)$.

where $(\lambda_n) \subseteq \mathbb{R}_{++}$ is a specially chosen sequence of stepsizes. In fact, we will require that

$$\lim_{n \rightarrow \infty} \lambda_n = 0 \quad \text{and} \quad \sum_{k=0}^{\infty} \lambda_k = +\infty. \quad (3.19)$$

This is satisfied, for instance, by

$$\lambda_n = \frac{1}{\sqrt{n}} \quad \forall n \in \mathbb{N}.$$

... not

Theorem 3.4.2

accuracy needed.

Let $f: \mathbb{R}^n \rightarrow \mathbb{R}$ be a continuous convex function, and let $\underset{x \in \mathbb{R}^n}{\arg \min} f(x)$. Let (x_k) be the sequence given by (3.18) with stepsizes given by (3.19). Then

$$\min_{0 \leq i \leq k} f(x_i) \rightarrow f(x^*) \text{ as } k \rightarrow \infty.$$

Proof

Suppose by way of a contradiction, that there exists $\tilde{f} \in \mathbb{R}$ such that

$$f(x_k) \geq \tilde{f} > f(x^*) \quad \forall k \in \mathbb{N}.$$

By continuity of f , there exists $\rho > 0$ such that $f(x) \leq \tilde{f}$ for all x such that $\|x - x^*\| \leq \rho$. In particular, for

$$y_k := x^* + \rho \frac{\phi_k}{\|\phi_k\|},$$

we have $f(y_n) \leq \tilde{f}$ since

$$\|y_n - x^*\| = \left\| p \cdot \frac{\phi_k}{\|\phi_k\|} \right\| = p.$$

Since $\phi_k \in \partial f(x_n)$, we have,

$$\begin{aligned} f(y_n) &\geq f(x_n) + \langle \phi_k, y_n - x_n \rangle \\ &\geq \tilde{f} + \langle \phi_k, x^* - x_n + y_n - x^* \rangle \\ &= \tilde{f} + \langle \phi_k, x^* - x_n \rangle + \langle \phi_k, y_n - x^* \rangle \\ &= \tilde{f} + \langle \phi_k, x^* - x_n \rangle + \langle \phi_k, p \cdot \frac{\phi_k}{\|\phi_k\|} \rangle \\ &= \tilde{f} + \langle \phi_k, x^* - x_n \rangle + p \|\phi_k\|. \end{aligned}$$

This implies

$$\begin{aligned} \langle \phi_k, x_n - x^* \rangle &\geq p \|\phi_k\| + \underbrace{\tilde{f} - f(y_n)}_{\nearrow 0} \\ &\geq p \|\phi_k\| \end{aligned}$$

$$\Rightarrow \frac{\langle \phi_k, x_k - x^* \rangle}{\|\phi_k\|} \rightarrow \rho.$$

Note that $\phi_k \neq 0$ because $x_k \notin \text{argmin}_{x \in \mathbb{R}^n} f(x)$, so division by $\|\phi_k\|$ is valid.

Consequently, we have

$$\begin{aligned}
 & \|x_{k+1} - x^*\|^2 \\
 &= \left\| (x_k - x^*) - \lambda_k \frac{\phi_k}{\|\phi_k\|} \right\|^2 \\
 &= \|x_k - x^*\|^2 - 2\lambda_k \underbrace{\left\langle \frac{\phi_k}{\|\phi_k\|}, x_k - x^* \right\rangle}_{= r_k^2} + \underbrace{\left\| \lambda_k \frac{\phi_k}{\|\phi_k\|} \right\|^2}_{\lambda_k^2} \\
 &= \|x_k - x^*\|^2 - 2\lambda_k \frac{\langle \phi_k, x_k - x^* \rangle}{\|\phi_k\|} + \lambda_k^2 \\
 &\leq \|x_k - x^*\|^2 - 2\lambda_k \rho + \lambda_k^2 \leq \lambda_k \rho
 \end{aligned}$$

Since $\lambda_k \rightarrow 0$, there is a sufficiently large $k_* \in \mathbb{N}$ such that

$$\lambda_n \leq p \quad \forall n \geq k_0.$$

Hence, for $n \geq k_0$, we have

$$\begin{aligned} & \|x_{n+1} - x^*\|^2 \\ & \leq \|x_n - x^*\|^2 - 2\lambda_n p + \lambda_n p \\ & = \|x_n - x^*\|^2 - \lambda_n p. \end{aligned}$$

Equivalently,

$$\|x_{n+1} - x^*\|^2 + \lambda_n p \leq \|x_n - x^*\|^2. \quad \forall n \geq k_0.$$

Telescoping this inequality gives

$$\|x_{n+1} - x^*\|^2 + p \sum_{i=k_0}^n \lambda_i \leq \|x_{k_0} - x^*\|^2 \quad \forall n \geq k_0.$$

This implies

$$p \sum_{i=k_0}^n \lambda_i \leq \|x_{k_0} - x^*\|^2 \quad \forall n \geq k_0$$

Taking the limit as $k \rightarrow \infty$, gives

$$p \sum_{i=k_0}^{\infty} \lambda_i \leq \|x_{k_0} - x^*\|^2 < \infty,$$

which contradicts (3.19). \square

Example 3.4.3

Consider the nonsmooth minimisation problem

$$\min_{x \in \mathbb{R}^n} f(x) \text{ where } f(x) = \|Ax - b\|.$$

where $A \in \mathbb{R}^{m \times n}$ and $b \in \mathbb{R}^m$. This is a bad way to formulate this problem as it is equivalent to the smooth minimisation problem considered last week:

$$\min_{x \in \mathbb{R}^n} \frac{1}{2} \|Ax - b\|^2.$$

Nevertheless, we will use it to illustrate the importance of problem formulation, and the difference between GD and sub-gradient method.

In order to apply the subgradient method, we must compute the subdifferential at a point x such that $Ax \neq b$. To do this, note that

$$f(x) = \left(\|Ax - b\|^2 \right)^{\frac{1}{2}}$$

Then using the chain-rule for differentiation, we have:

$$\begin{aligned} \nabla f(x) &= \frac{1}{2} \left(\|Ax - b\|^2 \right)^{-\frac{1}{2}} \cdot 2A^T(Ax - b) \\ &= \frac{A^T(Ax - b)}{\|Ax - b\|} \end{aligned}$$

By Prop. 2.5.3, we have

$$\partial f(x) = \left\{ \frac{A^T(Ax - b)}{\|Ax - b\|} \right\} \text{ if } Ax \neq b.$$