

Lecture 6. Subgradient Methods

Bao Wang

Department of Mathematics

Scientific Computing and Imaging Institute

University of Utah

Math 5750/6880, Fall 2021

Loss function

So far, we have formulated training machine learning models as

$$\min f(\mathbf{x}) = \frac{1}{N} \sum_{i=1}^N \mathcal{L}_i(\mathbf{x}) + R(\mathbf{x})$$

where \mathbf{x} is the parameter of the machine learning model, $\mathcal{L}_i(\mathbf{x})$ is the loss of the i th training instance, and $R(\mathbf{x})$ is the regularization term.

How to find the optimal \mathbf{x}^* if $R(\mathbf{x})$ is not differentiable everywhere, e.g. ℓ_1 -regularization?

Subgradient methods or proximal gradient methods.

We will focus on the "subgradient-based methods" in this lecture, i.e.,

$$\mathbf{x}^{t+1} = \mathbf{x}^t - \eta_t \mathbf{g}^t, \quad (1)$$

where \mathbf{g}^t is any subgradient of f at \mathbf{x}^t .

Subgradients

We say \mathbf{g} is a subgradient of f at the point \mathbf{x} if

$$f(\mathbf{z}) \geq \underbrace{f(\mathbf{x}) + \mathbf{g}^\top(\mathbf{z} - \mathbf{x})}_{\text{a linear under-estimate of } f}, \quad \forall \mathbf{z}. \quad (2)$$

The set of all subgradients of f at \mathbf{x} is called the subdifferential of f at \mathbf{x} , denoted by $\partial f(\mathbf{x})$.

Example

Let $f(x) = |x|$, then

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

Example (a subgradient of norms at 0)

Let $f(\mathbf{x}) = \|\mathbf{x}\|$ for any norm $\|\cdot\|$, then for any \mathbf{g} obeying $\|\mathbf{g}\|_* \leq 1$, then

$$\mathbf{g} \in \partial f(0),$$

where $\|\cdot\|_*$ is the dual norm of $\|\cdot\|$ (i.e. $\|\mathbf{x}\|_* := \sup_{\mathbf{z}: \|\mathbf{z}\| \leq 1} \langle \mathbf{z}, \mathbf{x} \rangle$).

Proof. To see this, it suffices to prove that

$$f(\mathbf{z}) \geq f(0) + \langle \mathbf{g}, \mathbf{z} - 0 \rangle, \quad \forall \mathbf{z} \Leftrightarrow \langle \mathbf{g}, \mathbf{z} \rangle \leq \|\mathbf{z}\|, \quad \forall \mathbf{z}.$$

This follows from generalized Cauchy-Schwarz, i.e.

$$\langle \mathbf{g}, \mathbf{z} \rangle \leq \|\mathbf{g}\|_* \|\mathbf{z}\| \leq \|\mathbf{z}\|.$$

Basic rules of subgradient methods

Scaling: $\partial(\alpha f) = \alpha \partial f$ for $\alpha > 0$.

Basic rules of subgradient methods

Summation: $\partial(f_1 + f_2) = \partial f_1 + \partial f_2.$

Basic rules of subgradient methods

Affine transformation: if $h(\mathbf{x}) = f(\mathbf{Ax} + \mathbf{b})$, then $\partial h(\mathbf{x}) = \mathbf{A}^\top \partial f(\mathbf{Ax} + \mathbf{b})$.

Basic rules of subgradient methods

Chain rule: suppose f is convex, and g is differentiable, nondecreasing, and convex.
Let $h = g \circ f$, then

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

Basic rules of subgradient methods

Composition: suppose $f(\mathbf{x}) = h(f_1(\mathbf{x}), \dots, f_n(\mathbf{x}))$, where f_i 's are convex, and h is differentiable, nondecreasing, and convex. Let $\mathbf{q} = \nabla h(\mathbf{y})|_{\mathbf{y}=[f_1(\mathbf{x}), \dots, f_n(\mathbf{x})]}$, and $\mathbf{g}_i \in \partial f_i(\mathbf{x})$. Then

$$q_i \mathbf{g}_1 + \dots + q_n \mathbf{g}_n \in \partial f(\mathbf{x})$$

Basic rules of subgradient methods

Pointwise maximum: if $f(\mathbf{x}) = \max_{1 \leq i \leq k} f_i(\mathbf{x})$, then

$$\partial f(\mathbf{x}) = \underbrace{\text{conv} \left\{ \cup \{ \partial f_i(\mathbf{x}) \mid f_i(\mathbf{x}) = f(\mathbf{x}) \} \right\}}_{\text{convex hull of subdifferentials of all active functions}}$$

Basic rules of subgradient methods

Pointwise supremum: if $f(\mathbf{x}) = \sup_{\alpha \in \mathcal{F}} f_{\alpha}(\mathbf{x})$, then

$$\partial f(\mathbf{x}) = \text{closure} \left(\text{conv} \left\{ \cup \{ \partial f_{\alpha}(\mathbf{x}) \mid f_{\alpha}(\mathbf{x}) = f(\mathbf{x}) \} \right\} \right)$$

Example

Let $f(x) = \max\{f_1(x), f_2(x)\}$ where f_1 and f_2 are differentiable, then

$$\partial f(x) = \begin{cases} \{f_1'(x)\} & \text{if } f_1(x) > f_2(x) \\ [f_1'(x), f_2'(x)] & \text{if } f_1(x) = f_2(x) \\ \{f_2'(x)\} & \text{if } f_1(x) < f_2(x) \end{cases}$$

Example (ℓ_1 norm)

$$f(\mathbf{x}) = \|\mathbf{x}\|_1 = \sum_{i=1}^n \underbrace{|x_i|}_{:=f_i(\mathbf{x})} \quad \text{since}$$

$$\partial f_i(\mathbf{x}) = \begin{cases} \text{sgn}(x_i)\mathbf{e}_i & \text{if } x_i \neq 0 \\ [-1, 1] \cdot \mathbf{e}_i & \text{if } x_i = 0 \end{cases}$$

Note that $f(\mathbf{x}) = \|\mathbf{x}\|_1 = \sum_{i, |x_i| \neq 0} |x_i|$, thus we have

$$\sum_{i, x_i \neq 0} \text{sgn}(x_i)\mathbf{e}_i \in \partial f(\mathbf{x}).$$

How about the subgradient at $\mathbf{x} = \mathbf{0}$?

Example

Let $h(\mathbf{x}) = \|\mathbf{Ax} + \mathbf{b}\|_1$, and denote $f(\mathbf{x}) = \|\mathbf{x}\|_1$ and $\mathbf{A} = [\mathbf{a}_1, \dots, \mathbf{a}_m]^\top$, we have

$$\mathbf{g} = \sum_{i:\mathbf{a}_i^\top \mathbf{x} + b_i \neq 0} \text{sgn}(\mathbf{a}_i^\top \mathbf{x} + b_i) \mathbf{e}_i \in \partial f(\mathbf{Ax} + \mathbf{b})$$

$$\Rightarrow \mathbf{A}^\top \mathbf{g} = \sum_{i:\mathbf{a}_i^\top \mathbf{x} + b_i \neq 0} \text{sgn}(\mathbf{a}_i^\top \mathbf{x} + b_i) \mathbf{a}_i \in \partial h(\mathbf{x}).$$

Example

Consider the piecewise linear function

$$f(\mathbf{x}) = \max_{1 \leq i \leq m} \{\mathbf{a}_i^\top \mathbf{x} + b_i\},$$

pick any \mathbf{a}_j s.t. $\mathbf{a}_j^\top \mathbf{x} + b_j = \max_i \{\mathbf{a}_i^\top \mathbf{x} + b_i\}$, then

$$\mathbf{a}_j \in \partial f(\mathbf{x}).$$

Example (l_∞ norm)

Let $f(\mathbf{x}) = \|\mathbf{x}\|_\infty = \max_{1 \leq i \leq n} |x_i|$, if $\mathbf{x} \neq 0$, then pick any x_j obeying $|x_j| = \max_i |x_i|$ to obtain

$$\text{sgn}(x_j)\mathbf{e}_j \in \partial f(\mathbf{x}).$$

Negative subgradients are not necessarily descent directions

Consider $f(\mathbf{x}) = |x_1| + 3|x_2|$, at $\mathbf{x} = (1, 0)$: $\mathbf{g}_1 = (1, 0) \in \partial f(\mathbf{x})$ and $-\mathbf{g}_1$ is a descent direction; $\mathbf{g}_2 = (1, 3) \in \partial f(\mathbf{x})$ while $-\mathbf{g}_2$ is not a descent direction. This is because f is not continuous at \mathbf{x} , one can change directions significantly without violating the validity of subgradients.

Since $f(\mathbf{x}^t)$ is not necessarily monotone, we will keep track of the best point

$$f^{best,t} := \min_{1 \leq i \leq t} f(\mathbf{x}^i).$$

We also denote by $f^{opt} := \min_{\mathbf{x}} f(\mathbf{x})$ the optimal objective value.

Convex and Lipschitz problems

Clearly, we cannot analyze all nonsmooth functions. A nice class to start with is Lipschitz functions, i.e. the set of all f obeying

$$|f(\mathbf{x}) - f(\mathbf{z})| \leq L_f \|\mathbf{x} - \mathbf{z}\|_2, \quad \forall \mathbf{x} \text{ and } \mathbf{z}.$$

Fundamental inequality for projected subgradient methods

We'd like to optimize $\|\mathbf{x}^{t+1} - \mathbf{x}^*\|_2^2$, but do not have access to \mathbf{x}^* . The key idea is **majorization-minimization**: find another function that majorizes $\|\mathbf{x}^{t+1} - \mathbf{x}^*\|_2^2$, and optimize the majorizing function

Lemma 1. Subgradient update rule (1) obeys

$$\|\mathbf{x}^{t+1} - \mathbf{x}^*\|_2^2 \leq \underbrace{\|\mathbf{x}^t - \mathbf{x}^*\|_2^2}_{\text{fixed}} - 2\eta_t(f(\mathbf{x}^t) - f^{opt}) + \eta_t^2 \|\mathbf{g}^t\|_2^2 \quad (3)$$

majorizing function

Proof of Lemma 1

$$\begin{aligned}\|\mathbf{x}^{t+1} - \mathbf{x}^*\|_2^2 &= \|\mathbf{x}^t - \eta_t \mathbf{g}^t - \mathbf{x}^*\|_2^2 \\ &= \|\mathbf{x}^t - \mathbf{x}^*\|_2^2 - 2\eta_t \langle \mathbf{x}^t - \mathbf{x}^*, \mathbf{g}^t \rangle + \eta_t^2 \|\mathbf{g}^t\|_2^2 \\ &\leq \|\mathbf{x}^t - \mathbf{x}^*\|_2^2 - 2\eta_t (f(\mathbf{x}^t) - f(\mathbf{x}^*)) + \eta_t^2 \|\mathbf{g}^t\|_2^2\end{aligned}$$

where the last line uses the subgradient inequality

$$f(\mathbf{x}^*) - f(\mathbf{x}^t) \geq \langle \mathbf{x}^* - \mathbf{x}^t, \mathbf{g}^t \rangle.$$

Proof of Lemma 1 – Cont'd

The majorizing function in (3) suggests a step size (Polyak's stepsize rule)

$$\eta_t = \frac{f(\mathbf{x}^t) - f^{opt}}{\|\mathbf{g}_t\|_2^2}, \quad (4)$$

which leads to error reduction

$$\|\mathbf{x}^{t+1} - \mathbf{x}^*\|_2^2 \leq \|\mathbf{x}^t - \mathbf{x}^*\|_2^2 - \frac{(f(\mathbf{x}^t) - f(\mathbf{x}^*))^2}{\|\mathbf{g}^t\|_2^2}. \quad (5)$$

The algorithm is useful if f^{opt} is known, and the estimation error is monotonically decreasing with Polyak's stepsize.

Convergence of projected subgradient method with Polyak's stepsize

Theorem 1. Suppose f is convex and L_f -Lipschitz continuous. Then the projected subgradient method (1) with Polyak's stepsize rule obeys

$$f^{best,t} - f^{opt} \leq \frac{L_f \|\mathbf{x}^0 - \mathbf{x}^*\|_2}{\sqrt{t+1}}.$$

The rate $O(1/\sqrt{t})$ is called sublinear convergence rate.

Convergence of subgradient method with Polyak's stepsize

Proof. We have seen from (5) that

$$\begin{aligned}(f(\mathbf{x}^t) - f(\mathbf{x}^*))^2 &\leq \{\|\mathbf{x}^t - \mathbf{x}^*\|_2^2 - \|\mathbf{x}^{t+1} - \mathbf{x}^*\|_2^2\} \|\mathbf{g}^t\|_2^2 \\ &\leq \{\|\mathbf{x}^t - \mathbf{x}^*\|_2^2 - \|\mathbf{x}^{t+1} - \mathbf{x}^*\|_2^2\} L_f^2.\end{aligned}$$

Applying it recursively for all iterations (from 0th to t th) and summing them up yield

$$\sum_{k=0}^t (f(\mathbf{x}^k) - f(\mathbf{x}^*))^2 \leq \left\{ \|\mathbf{x}^0 - \mathbf{x}^*\|_2^2 - \|\mathbf{x}^{t+1} - \mathbf{x}^*\|_2^2 \right\} L_f^2.$$

therefore

$$(t+1)(f^{best,t} - f^{opt})^2 \leq \|\mathbf{x}^0 - \mathbf{x}^*\|_2^2 L_f^2$$

which concludes the proof.

Polyak's stepsize rule requires knowledge of f^{opt} , which is often unknown *a priori*. We might often need simpler rules for setting stepsizes.

How about the other stepsize rules?

Convergence of subgradient methods for convex and Lipschitz functions

Theorem 2. [Subgradient methods for convex and Lipschitz functions] Suppose f is convex and L_f -Lipschitz continuous. Then the projected subgradient update rule (1) obeys

$$f^{best,t} - f^{opt} \leq \frac{\|\mathbf{x}^0 - \mathbf{x}^*\|_2^2 + L_f^2 \sum_{i=0}^t \eta_i^2}{2 \sum_{i=0}^t \eta_i}.$$

Convergence of subgradient methods for convex and Lipschitz functions – General step size

Proof. Applying Lemma 1 recursively gives

$$\|\mathbf{x}^{t+1} - \mathbf{x}^*\|_2^2 \leq \|\mathbf{x}^0 - \mathbf{x}^*\|_2^2 - 2 \sum_{i=0}^t \eta_i (f(\mathbf{x}^i) - f^{opt}) + \sum_{i=0}^t \eta_i^2 \|\mathbf{g}^i\|_2^2.$$

Rearranging terms, we are left with

$$\begin{aligned} 2 \sum_{i=0}^t \eta_i (f(\mathbf{x}^i) - f^{opt}) &\leq \|\mathbf{x}^0 - \mathbf{x}^*\|_2^2 - \|\mathbf{x}^{t+1} - \mathbf{x}^*\|_2^2 + \sum_{i=0}^t \eta_i^2 \|\mathbf{g}^i\|_2^2 \\ &\leq \|\mathbf{x}^0 - \mathbf{x}^*\|_2^2 + L_f^2 \sum_{i=0}^t \eta_i^2. \end{aligned}$$

Thus

$$f^{best,t} - f^{opt} \leq \frac{\sum_{i=0}^t \eta_i (f(\mathbf{x}^i) - f^{opt})}{\sum_{i=0}^t \eta_i} \leq \frac{\|\mathbf{x}^0 - \mathbf{x}^*\|_2^2 + L_f^2 \sum_{i=0}^t \eta_i^2}{2 \sum_{i=0}^t \eta_i}.$$

Other stepsize rules

Constant step size $\eta_t \equiv \eta$:

$$\lim_{t \rightarrow \infty} f^{best,t} - f^{opt} \leq \frac{L_f^2 \eta}{2},$$

i.e. may converge to non-optimal points. (Note that $2 \sum_{i=0}^t \eta_i = \infty$)

Diminishing step size obeying

$$\sum_t \eta_t^2 \leq \infty \text{ and } \sum_t \eta_t \rightarrow \infty : \lim_{t \rightarrow \infty} f^{best,t} - f^{opt} = 0,$$

i.e. converges to optimal points.

Other stepsize rules

Optimal choice?

$$\eta_t = \frac{1}{\sqrt{t}}, \quad f^{best,t} - f^{opt} \lesssim \frac{\|\mathbf{x}^0 - \mathbf{x}^*\|_2^2 + L_f^2 \log t}{\sqrt{t}},$$

i.e. attains ϵ -accuracy within about $O(1/\epsilon^2)$ iterations (ignoring the log factor).

Strongly convex and Lipschitz problems

Strongly convex and Lipschitz problems

Theorem 3. [Subgradient methods for strongly convex and Lipschitz functions] Let f be μ -strongly convex and L_f -Lipschitz continuous over \mathcal{C} . If $\eta_t \equiv \eta = \frac{2}{\mu(t+1)}$, then

$$f^{best,t} - f^{opt} \leq \frac{2L_f^2}{\mu} \cdot \frac{1}{t+1}.$$

Strongly convex and Lipschitz problems (Proof)

Proof of Theorem 3. When f is μ -strongly convex, we can improve Lemma 1 to (exercise)

$$\|\mathbf{x}^{t+1} - \mathbf{x}^*\|_2^2 \leq (1 - \mu\eta_t)\|\mathbf{x}^t - \mathbf{x}^*\|_2^2 - 2\eta_t(f(\mathbf{x}^t) - f^{opt}) + \eta_t^2\|\mathbf{g}^t\|_2^2$$

$$\Rightarrow f(\mathbf{x}^t) - f^{opt} \leq \frac{1 - \mu\eta_t}{2\eta_t}\|\mathbf{x}^t - \mathbf{x}^*\|_2^2 - \frac{1}{2\eta_t}\|\mathbf{x}^{t+1} - \mathbf{x}^*\|_2^2 + \frac{\eta_t}{2}\|\mathbf{g}^t\|_2^2$$

Strongly convex and Lipschitz problems (Proof Cont'd)

Since $\eta_t = 2/(\mu(t+1))$, we have

$$f(\mathbf{x}^t) - f^{opt} \leq \frac{\mu(t-1)}{4} \|\mathbf{x}^t - \mathbf{x}^*\|_2^2 - \frac{\mu(t+1)}{4} \|\mathbf{x}^{t+1} - \mathbf{x}^*\|_2^2 + \frac{1}{\mu(t+1)} \|\mathbf{g}^t\|_2^2$$

and hence

$$t(f(\mathbf{x}^t) - f^{opt}) \leq \frac{\mu t(t-1)}{4} \|\mathbf{x}^t - \mathbf{x}^*\|_2^2 - \frac{\mu t(t+1)}{4} \|\mathbf{x}^{t+1} - \mathbf{x}^*\|_2^2 + \frac{1}{\mu} \|\mathbf{g}^t\|_2^2.$$

Summing over all iterations before t , we get

$$\sum_{k=0}^t k(f(\mathbf{x}^k) - f^{opt}) \leq 0 - \frac{\mu t(t+1)}{4} \|\mathbf{x}^{t+1} - \mathbf{x}^*\|_2^2 + \frac{1}{\mu} \sum_{k=0}^t \|\mathbf{g}^k\|_2^2 \leq \frac{t}{\mu} L_f^2.$$

$$\Rightarrow f^{best,k} - f^{opt} \leq \frac{L_f^2}{\mu} \frac{t}{\sum_{k=0}^t k} \leq \frac{2L_f^2}{\mu} \frac{1}{t+1}.$$

Subgradient method summary

	stepsize rule	convergence rate	iteration complexity
convex & Lipschitz problems	$\eta_t \asymp \frac{1}{\sqrt{t}}$	$O\left(\frac{1}{\sqrt{t}}\right)$	$O\left(\frac{1}{\varepsilon^2}\right)$
strongly convex & Lipschitz problems	$\eta_t \asymp \frac{1}{t}$	$O\left(\frac{1}{t}\right)$	$O\left(\frac{1}{\varepsilon}\right)$

In contrast, gradient descent is much faster!