

SDSC6015 - Semester A, 2025
Stochastic Optimization for Machine Learning
5. Projected Subgradient Method

1 Motivation

Recall that using gradient descent, one can solve an unconstrained, smooth convex optimization problem

$$\underset{\mathbf{x} \in \mathbb{R}^n}{\text{minimize}} f(\mathbf{x}),$$

where $f : \mathbb{R}^n \rightarrow \mathbb{R}$ is convex and differentiable with $\text{dom}(f) = \mathbb{R}^n$. Gradient descent is a line search method in which for each iteration, we do

$$\mathbf{x}_{k+1} = \mathbf{x}_k + \alpha_k \cdot \mathbf{d}_k$$

where

- $k \in \mathbb{Z}_+$: the index of iterations
- \mathbf{x}_k : solution at the k^{th} iteration
- \mathbf{d}_k : direction at the k^{th} iteration
- α_k : stepsize at the k^{th} iteration

Limitation: What if

- We have constraints?
- The function f is not differentiable?

Answer: Projected Subgradient Method!

- We have constraints? Projection.
- The function f is not differentiable? Subgradient.

2 Subgradients: Definitions and Properties

Definition 1 (Proper function). *A function $f : \mathbb{R}^n \rightarrow (-\infty, \infty]$ is called proper if $\text{dom}(f) \neq \emptyset$.*

Definition 2 (Subgradient). *Assume $f : \mathbb{R}^n \rightarrow (-\infty, \infty]$ is a proper convex function and $\mathbf{x}_0 \in \text{dom}(f)$. Then, $\mathbf{y} \in \mathbb{R}^n$ is called a subgradient of f at \mathbf{x}_0 if*

$$f(\mathbf{x}) \geq f(\mathbf{x}_0) + \mathbf{y}^\top (\mathbf{x} - \mathbf{x}_0), \quad \forall \mathbf{x} \in \text{dom}(f).$$

The set $\partial f(\mathbf{x}_0)$ of all subgradients is the subdifferential of f at \mathbf{x}_0 .

Note:

- If f is differentiable at \mathbf{x}_0 , then $\partial f(\mathbf{x}_0) = \{\nabla f(\mathbf{x}_0)\}$.
 - Generally, $\partial f(\mathbf{x}_0)$ can be empty or contain one or infinitely many subgradients.
-
- $\partial f(\mathbf{x}_0)$ is a closed convex set.

2.1 Some Examples

Example 1 (Subdifferential of norms at $\mathbf{0}$). Consider $f : \mathbb{R}^n \rightarrow \mathbb{R}$ and $f(\mathbf{x}) = \|\mathbf{x}\|$, then

$$\partial f(\mathbf{0}) = B_{\|\cdot\|_*}[\mathbf{0}, 1] = \{\mathbf{y} \in \mathbb{R}^n \mid \|\mathbf{y}\|_* \leq 1\},$$

where $\|\cdot\|_*$ is the dual norm of $\|\cdot\|$.

Example 2 (Subdifferential of the L_1 -norm at $\mathbf{0}$). Consider $f : \mathbb{R}^n \rightarrow \mathbb{R}$ and $f(\mathbf{x}) = \|\mathbf{x}\|_1$, then

$$\partial f(\mathbf{0}) = B_{\|\cdot\|_\infty}[\mathbf{0}, 1] = [-1, 1]^n.$$

Example 3 (Subdifferential of the L_2 -norm). Consider $f : \mathbb{R}^n \rightarrow \mathbb{R}$ and $f(\mathbf{x}) = \|\mathbf{x}\|_2$, then

$$\partial f(\mathbf{x}) = \begin{cases} \left\{ \frac{\mathbf{x}}{\|\mathbf{x}\|_2} \right\}, & \mathbf{x} \neq \mathbf{0}, \\ B_{\|\cdot\|_2}[\mathbf{0}, 1], & \mathbf{x} = \mathbf{0}, \end{cases}$$

Example 4 (Subdifferential of the indicator functions). Consider $S \subseteq \mathbb{R}^n$ is nonempty and the indicator function

$$\delta_S(\mathbf{x}) = \begin{cases} 0, & \mathbf{x} \in S, \\ \infty, & \mathbf{x} \notin S. \end{cases}$$

Then

$$\partial \delta_S(\mathbf{x}) = N_S(\mathbf{x}) = \{\mathbf{y} \in \mathbb{R}^n \mid \langle \mathbf{y}, \mathbf{z} - \mathbf{x} \rangle \leq 0, \forall \mathbf{z} \in S\}, \quad \forall \mathbf{x} \in S.$$

For $\mathbf{x} \notin S$, $\partial \delta_S(\mathbf{x}) = N_S(\mathbf{x}) = \emptyset$ by convention. Note that $N_S(\mathbf{x})$ is called the normal cone of S at \mathbf{x} .

Example 5 (Subgradient of the dual function). Consider the minimization problem

$$\min_{\mathbf{x} \in \mathbb{R}^n} \{f(\mathbf{x}) \mid \mathbf{h}(\mathbf{x}) \leq \mathbf{0}, \mathbf{x} \in \mathcal{X}\},$$

where $\emptyset \neq \mathcal{X}$, $f : \mathbb{R}^n \rightarrow \mathbb{R}$, and $\mathbf{h} : \mathbb{R}^n \rightarrow \mathbb{R}^m$. The Lagrange dual function is

$$g(\boldsymbol{\lambda}) = \min_{\mathbf{x} \in \mathcal{X}} \left\{ L(\mathbf{x}, \boldsymbol{\lambda}) \equiv f(\mathbf{x}) + \boldsymbol{\lambda}^\top \mathbf{h}(\mathbf{x}) \right\}.$$

2.2 Properties

Recall the following results from our lecture on Convex Sets!

Definition 3 (Supporting hyperplane). Consider a nonempty set $C \subseteq \mathbb{R}^n$ and a boundary point $\mathbf{x}_0 \in \mathbf{bd}(C)$. If $\mathbf{a} \neq \mathbf{0}$ in \mathbb{R}^n satisfies $\mathbf{a}^\top \mathbf{x} \leq \mathbf{a}^\top \mathbf{x}_0$, $\forall \mathbf{x} \in C$, then $\{\mathbf{x} \mid \mathbf{a}^\top \mathbf{x} = \mathbf{a}^\top \mathbf{x}_0\}$ is called a supporting hyperplane to C at \mathbf{x}_0 .

Theorem 1 (Weak Separating Hyperplane Theorem). Consider any convex set $C \subseteq \mathbb{R}^n$ and a point $\mathbf{x}_0 \in \mathbb{R}^n \setminus C$. Then, there exist $\mathbf{a} \neq \mathbf{0}$ (in \mathbb{R}^n) and $b \in \mathbb{R}$ with

$$\mathbf{a}^\top \mathbf{x} \leq b, \quad \forall \mathbf{x} \in C, \quad \text{and} \quad \mathbf{a}^\top \mathbf{x}_0 \geq b.$$

Theorem 2. For every convex set $C \subseteq \mathbb{R}^n$ and any $\mathbf{x}_0 \in \mathbf{bd}(C)$ there is a supporting hyperplane to C at \mathbf{x}_0 .

Theorem 3. Assume that $f : \mathbb{R}^n \rightarrow \mathbb{R}$ is a proper convex function and $\mathbf{x}_0 \in \text{int}(\text{dom}(f))$. Then, $\partial f(\mathbf{x}_0) \neq \emptyset$.

3 Computing Subgradients

There are rules that are useful for computing subgradients and subdifferentials. When the rules/results are useful for computing some of the subgradients in the subdifferential set, they are referred as “*weak results*”; when the rules/results are useful for obtaining the full characterization of the subdifferential set, they are referred as “*strong results*”.

3.1 Multiplication by a positive scalar

Theorem 4. *Let $f : \mathbb{R}^n \rightarrow \mathbb{R}$ be a proper function and let $\alpha > 0$. Then for any $\mathbf{x} \in \text{dom}(f)$*

$$\partial(\alpha f)(\mathbf{x}) = \alpha \partial f(\mathbf{x}).$$

Recall that for a scalar $\beta \in \mathbb{R}$ and a set $S \subseteq \mathbb{R}^n$, the set $\beta S = \{\beta \mathbf{x} \mid \mathbf{x} \in S\}$.

3.2 Summation

Theorem 5. *Let $f_1, f_2 : \mathbb{R}^n \rightarrow \mathbb{R}$ be proper convex functions, and let $\mathbf{x} \in \text{dom}(f_1) \cap \text{dom}(f_2)$.*

(a) *The following inclusion holds:*

$$\partial f_1(\mathbf{x}) + \partial f_2(\mathbf{x}) \subseteq \partial(f_1 + f_2)(\mathbf{x}).$$

(b) *If $\mathbf{x} \in \text{int}(\text{dom}(f_1)) \cap \text{int}(\text{dom}(f_2))$, then*

$$\partial(f_1 + f_2)(\mathbf{x}) = \partial f_1(\mathbf{x}) + \partial f_2(\mathbf{x}).$$

Recall that for two sets $A, B \subseteq \mathbb{R}^n$ that reside in the same space, the sum $A + B$ stands for the Minkowski sum given by $A + B = \{\mathbf{a} + \mathbf{b} \mid \mathbf{a} \in A, \mathbf{b} \in B\}$.

Example 6 (L_1 -norm function). Consider the function $f : \mathbb{R}^n \rightarrow \mathbb{R}$ and $f(\mathbf{x}) = \|\mathbf{x}\|_1$. What is the subdifferential? Could you compute any subgradient?

3.3 Affine Transformation

Theorem 6. Let $f : \mathbb{R}^n \rightarrow \mathbb{R}$ be a proper convex function and $h(\mathbf{x}) = f(\mathbf{A}\mathbf{x} + \mathbf{b})$ for some $\mathbf{A} \in \mathbb{R}^{m \times n}$ and $\mathbf{b} \in \mathbb{R}^m$. Assume that h is proper, meaning that

$$dom(h) = \{\mathbf{x} \in \mathbb{R}^n \mid \mathbf{A}\mathbf{x} + \mathbf{b} \in dom(f)\} \neq \emptyset.$$

(a) For any $\mathbf{x} \in dom(h)$,

$$\mathbf{A}^\top(\partial f(\mathbf{A}\mathbf{x} + \mathbf{b})) \subseteq \partial h(\mathbf{x}).$$

(b) If $\mathbf{x} \in int(dom(h))$ and $\mathbf{A}\mathbf{x} + \mathbf{b} \in int(dom(f))$, then

$$h(\mathbf{x}) = \mathbf{A}^\top(\partial f(\mathbf{A}\mathbf{x} + \mathbf{b})).$$

Note that $\mathbf{A}^\top S = \{\mathbf{A}^\top \mathbf{x} \mid \mathbf{x} \in S\}$.

Example 7. Consider the function $f : \mathbb{R}^n \rightarrow \mathbb{R}$ and $f(\mathbf{x}) = \|\mathbf{A}\mathbf{x} + \mathbf{b}\|_1$. What is the subdifferential?

3.4 Composition

Theorem 7. Let $f : \mathbb{R}^n \rightarrow \mathbb{R}$ be a convex function and $g : \mathbb{R} \rightarrow \mathbb{R}$ be a nondecreasing convex function. Let $\mathbf{x} \in \mathbb{R}^n$ and suppose that g is differentiable at the point $f(\mathbf{x})$. Let $h = g \circ f$. Then

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

Example 8. Consider the function $h : \mathbb{R}^n \rightarrow \mathbb{R}$ and $h(\mathbf{x}) = \|\mathbf{x}\|_1^2$.

3.5 Maximization

Theorem 8. Let $f_1, f_2, f_3, \dots, f_m : \mathbb{R}^n \rightarrow \mathbb{R}$ be proper convex functions, and define

$$f(\mathbf{x}) = \max\{f_1(\mathbf{x}), f_2(\mathbf{x}), \dots, f_m(\mathbf{x})\}.$$

Let $\mathbf{x} \in \cap_{i=1}^m \text{int}(\text{dom}(f_i))$. Then

$$\partial f(\mathbf{x}) = \text{conv}(\cup_{i \in I(\mathbf{x})} \partial f_i(\mathbf{x})),$$

where $I(\mathbf{x}) = \{i \in [m] \mid f_i(\mathbf{x}) = f(\mathbf{x})\}$

Example 9. Let $f : \mathbb{R}^n \rightarrow \mathbb{R}$ be given by $f(\mathbf{x}) = \max\{x_1, x_2, \dots, x_n\}$.

Example 10. Let $f : \mathbb{R}^n \rightarrow \mathbb{R}$ be given by $f(\mathbf{x}) = \|\mathbf{x}\|_\infty$.

4 Orthogonal Projection

To handle constraints, we make use of the orthogonal projection.

Definition 4. Let $C \subseteq \mathbb{R}^n$ be nonempty closed and convex set. The orthogonal projection mapping

$$P_C(\mathbf{x}) \equiv \arg \min_{\mathbf{y} \in C} \|\mathbf{y} - \mathbf{x}\|.$$

Theorem 9. Let $C \subseteq \mathbb{R}^n$ be nonempty, closed, and convex. Then, for every $\mathbf{x}_0 \in \mathbb{R}^n$, there exists a unique point $\mathbf{x}^* \in C$ that is closest (in the Euclidean norm) to \mathbf{x}_0 .

4.1 Examples

1. If $C = \mathbb{R}_+^n$, $P_C(\mathbf{x}) = \max\{\mathbf{x}, \mathbf{0}\} = [\mathbf{x}]_+$
2. If $C = \text{Box}[\boldsymbol{\ell}, \mathbf{u}] = \{\mathbf{x} \in \mathbb{R}^n \mid \boldsymbol{\ell} \leq \mathbf{x} \leq \mathbf{u}\}$, $P_C(\mathbf{x}) = (\min\{\max\{x_i, \ell_i\}, u_i\})_{i=1}^n$
3. If $C = \{\mathbf{x} \in \mathbb{R}^n \mid \mathbf{A}\mathbf{x} = \mathbf{b}\}$ (\mathbf{A} has full row rank), $P_C(\mathbf{x}) = \mathbf{x} - \mathbf{A}^\top(\mathbf{A}\mathbf{A}^\top)^{-1}(\mathbf{A}\mathbf{x} - \mathbf{b})$
4. If $C = B_{\|\cdot\|_2}[\mathbf{c}, r]$, $P_C(\mathbf{x}) = \mathbf{c} + \frac{r}{\max\{\|\mathbf{x} - \mathbf{c}\|_2, r\}}(\mathbf{x} - \mathbf{c})$
5. If $C = \{\mathbf{x} \mid \mathbf{a}^\top \mathbf{x} \leq \alpha\}$, $P_C(\mathbf{x}) = \mathbf{x} - \frac{[\mathbf{a}^\top \mathbf{x} - \alpha]_+}{\|\mathbf{a}\|^2} \mathbf{a}$

5 Projected Subgradient Method

Consider the following optimization model

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

where (we assume)

- $f : \mathbb{R}^n \rightarrow \mathbb{R}$ is proper closed and convex.
- $C \subseteq \mathbb{R}^n$ is nonempty close and convex.
- $C \subseteq \text{int}(\text{dom}(f))$
- The optimal set is nonempty and denoted by \mathcal{X}^* . The optimal value of the problem is denoted by f^* .
- There exists a constant $L_f > 0$ for which $\|\mathbf{y}\| \leq L_f$ for all $\mathbf{y} \in \partial f(\mathbf{x})$, $\mathbf{x} \in C$.

Projected Subgradient Method

Initialization: $\mathbf{x}_0 \in C$

General step: For $k = 0, 1, 2, \dots$

- (a) pick a stepsize $\alpha_k > 0$ and a subgradient $f'(\mathbf{x}_k) \in \partial f(\mathbf{x}_k)$.
- (b) set $\mathbf{x}_{k+1} = P_C(\mathbf{x}_k - \alpha_k f'(\mathbf{x}_k))$.

The sequence generated by the projected subgradient method is $\{\mathbf{x}_k\}_{k \geq 0}$ while the sequence of function values generated by the method is $\{f(\mathbf{x}_k)\}_{k \geq 0}$. This sequence of functions values is not necessarily monotone. Therefore, we would be interested in

$$f_{\text{best}}^k \equiv \min_{n=0, \dots, k} f(\mathbf{x}_n).$$

Lemma 1. Let $\{\mathbf{x}_k\}_{k \geq 0}$ be the sequence generated by the projected subgradient method. Then for any $\mathbf{x}^* \in \mathcal{X}^*$ and $k \geq 0$,

$$\|\mathbf{x}_{k+1} - \mathbf{x}^*\|^2 \leq \|\mathbf{x}_k - \mathbf{x}^*\|^2 - 2\alpha_k(f(\mathbf{x}_k) - f^*) + \alpha_k^2 \|f'(\mathbf{x}_k)\|^2.$$

5.1 Polyak's Stepsize

The Polyak's stepsize rule takes stepsize α_k that minimize the RHS of the result in Lemma 1, and so

$$\alpha_k = \begin{cases} \frac{f(\mathbf{x}_k) - f^*}{\|f'(\mathbf{x}_k)\|^2}, & f'(\mathbf{x}_k) \neq \mathbf{0}, \\ 1, & f'(\mathbf{x}_k) = \mathbf{0}. \end{cases}$$

Theorem 10. Let $\{\mathbf{x}_k\}_{k \geq 0}$ be the sequence generated by the projected subgradient method with Polyak's stepsize rule. Then

- $\|\mathbf{x}_{k+1} - \mathbf{x}^*\|^2 \leq \|\mathbf{x}_k - \mathbf{x}^*\|^2$ for any $k \geq 0$ and $\mathbf{x}^* \in \mathcal{X}^*$;
- $f(\mathbf{x}_k) \rightarrow f^*$ as $k \rightarrow \infty$;
- $f_{best}^k - f^* \leq \frac{L_f d_{\mathcal{X}^*}(\mathbf{x}_0)}{\sqrt{k+1}}$ for any $k \geq 0$. [$d_{\mathcal{X}^*}$ is the distance function.]

5.2 Dynamic Stepsizes

Theorem 11. Let $\{\mathbf{x}_k\}_{k \geq 0}$ be the sequence generated by the projected subgradient method with positive stepsize $\{\alpha_k\}_{k \geq 0}$. If

$$\frac{\sum_{n=0}^k \alpha_n^2}{\sum_{n=0}^k \alpha_n} \rightarrow 0 \quad \text{as } k \rightarrow \infty,$$

then

$$f_{best}^k \rightarrow f^* \quad \text{as } k \rightarrow \infty.$$

Possible choice: $\alpha_k = 1/\sqrt{k+1}$.

Theorem 12. Let $\{\mathbf{x}_k\}_{k \geq 0}$ be the sequence generated by the projected subgradient method with positive stepsizes $\alpha_k = 1/(\|f'(\mathbf{x}_k)\|\sqrt{k+1})$ if $f'(\mathbf{x}_k) \neq 0$ and $\alpha_k = 1/L_f$ otherwise. Then for any $k \geq 1$

$$f_{best}^k - f^* \leq \frac{L_f}{2} \frac{\|\mathbf{x}_0 - \mathbf{x}^*\|^2 + 1 + \log(k+1)}{\sqrt{k+1}}$$

and

$$f(\mathbf{x}_{(k)}) - f^* \leq \frac{L_f}{2} \frac{\|\mathbf{x}_0 - \mathbf{x}^*\|^2 + 1 + \log(k+1)}{\sqrt{k+1}}$$

where $\mathbf{x}_{(k)} = \frac{1}{\sum_{n=0}^k \alpha_n} \sum_{n=0}^k \alpha_n \mathbf{x}_n$.

6 Special Case: Convex Feasibility Problem

Let $S_1, S_2, \dots, S_m \in \mathbb{R}^n$ be close and convex sets. Assume that

$$S \equiv \cap_{i=1}^m S_i \neq \emptyset.$$

The *convex feasibility problem* is the problem of finding a point \mathbf{x} in the intersection $\cap_{i=1}^m S_i$. One can formulate this problem as

$$\min_{\mathbf{x} \in \mathbb{R}^n} \left\{ f(\mathbf{x}) \equiv \max_{i \in [m]} d_{S_i}(\mathbf{x}) \right\},$$

where $d_{S_i}(\mathbf{x}) = \min_{\mathbf{y} \in \mathbb{R}^n} \{ \|\mathbf{x} - \mathbf{y}\| \mid \mathbf{y} \in S_i \}$, for all $i \in [m]$. We assume that the intersection is nonempty, and we have $f^* = 0$. One can show that $L_f = 1$.

Greedy Projection Algorithm

Input: m nonempty closed and convex sets S_1, S_2, \dots, S_m

Initialization: pick $\mathbf{x}_0 \in \mathbb{R}^n$.

General step: For $k = 0, 1, 2, \dots$, compute

$$\mathbf{x}_{k+1} = P_{S_{i_k}}(\mathbf{x}_k),$$

where $i_k \in \arg \max_{i \in [m]} d_{S_i}(\mathbf{x}_k)$.

When $m = 2$, the algorithm amounts the alternating projection method.

Alternating Projection Method

Input: two nonempty closed and convex sets S_1, S_2 .

Initialization: pick $\mathbf{x}_0 \in S_2$ arbitrarily.

General step: For $k = 0, 1, 2, \dots$, compute

$$\mathbf{x}_{k+1} = P_{S_2}(P_{S_1}(\mathbf{x}_k)).$$

Special Case: Consider the system of linear equalities and inequalities (standard form of LP):

$$\mathbf{A}\mathbf{x} = \mathbf{b}, \quad \mathbf{x} \geq \mathbf{0},$$

where $\mathbf{A} \in \mathbb{R}^{m \times n}$ has full row rank and $\mathbf{b} \in \mathbb{R}^m$. Define

$$S_1 = \{\mathbf{x} \in \mathbb{R}^m \mid \mathbf{A}\mathbf{x} = \mathbf{b}\}, \quad S_2 = \mathbb{R}_+^n.$$

As we know, we have

$$P_{S_1}(\mathbf{x}) = \mathbf{x} - \mathbf{A}^\top (\mathbf{A}\mathbf{A}^\top)^{-1}(\mathbf{A}\mathbf{x} - \mathbf{b}), \quad P_{S_2}(\mathbf{x}) = [\mathbf{x}]_+.$$

Alternating Projection Method (System of Linear Equalities and Inequalities)

Initialization: pick $\mathbf{x}_0 \in \mathbb{R}_+^n$ arbitrarily.

General step: For $k = 0, 1, 2, \dots$, compute

$$\mathbf{x}_{k+1} = \left[\mathbf{x}_k - \mathbf{A}^\top (\mathbf{A}\mathbf{A}^\top)^{-1}(\mathbf{A}\mathbf{x}_k - \mathbf{b}) \right]_+$$