

Chapter IV

First-order methods

Skills to acquire

- Compute subdifferentials and proximals operators
- Proximal gradient algorithms
- Applications to well known optimization problems.

We consider a general convex optimization problem

$$\min_{\mathbf{x} \in \mathbb{R}^n} f(\mathbf{x}) \quad (\text{IV.1})$$

This chapter presents first order algorithms, i.e., which make use of the function gradient. Such algorithms have low computational cost per iterations.

1 Gradient descent algorithms

1.1 Description

We suppose that f is differentiable with $\text{dom}(f) = \mathbb{R}^n$. The gradient descent algorithm solves the problem (IV.1) using the core iteration k :

- **Choice of the step:** Choose $\tau^{(k)}$
- **Main iteration:** $\mathbf{x}^{(k+1)} = \mathbf{x}^{(k)} - \tau^{(k)} \nabla f(\mathbf{x}^{(k)})$

The basic algorithm have some variants which depend mainly on the choice of the step $\tau^{(k)}$. A simple *constant step*

$$\tau^{(k)} = \tau, \forall k,$$

can diverge if τ is too big and slow if τ is too small. On the other hand, an optimal *exact line search*

$$\tau^{(k)} = \min_{\tau} f(x^k - \tau \nabla f(x^l))$$

can be hard to compute, computationally demanding, or even does not exist.

Another popular approach is *Backtracking line search*, which consists in the following. At each iteration, start with $\tau^{(k)} = \tau_0$ (e.g., $\tau_0 = 1$). Then, take

$$\tau^{(k)} = \beta \tau^{(k)},$$

til it satisfies the Armijo condition

$$f(x^{(k)} - t^{(k)} \nabla f(x^{(k)})) \leq \left(f(x^{(k)}) - \alpha t^{(k)} \|\nabla f(x)\|_2^2 \right),$$

where $0 < \beta < 1$ and $0 < \alpha \leq 1/2$ are some fixed parameters. One can simply use $\alpha = \beta = 0.5$. These simple τ updates tend to work well in practice.

1.2 Convergence

Gradient descent algorithms are inexpensive as they do not require second derivatives.

Theorem 1.1. *Suppose that f is convex and $\nabla f(x)$ is L -Lipschitz continuous*

$$\|\nabla f(\mathbf{x}) - \nabla f(\mathbf{y})\|_2 \leq L \|\mathbf{x} - \mathbf{y}\|_2, \quad \forall \mathbf{x}, \mathbf{y}.$$

If the optimal value is finite and attained at x^ then the gradient descent converges to x^* , with a rate $O(1/\epsilon)^1$ for constant step size and $O(\log(1/\epsilon))$ for backtracking line search.*

For unconstrained problems, gradient descent is robust and less tuning. Thus, it is still widely used in practice. But, for a non-convex function f , the convergence is often slow and depends on the starting point and the scaling. In addition, it cannot handle nondifferentiable functions. There are many other better algorithms that deal with its shortcomings.

1.3 Accelerated gradient method

Heavy Ball Method In this technique, we add a momentum to gradient descent:

$$x^{(k+1)} = x^{(k)} - \alpha_k \nabla f(x^{(k)}) + \beta_k (x^{(k)} - x^{(k-1)}),$$

where momentum prevents oscillation.

Nesterov Extrapolation (A simple variant) Nesterov's exploration alternates between gradient updates and proper extrapolation

$$\begin{cases} x^{(k+1)} = y^{(k)} - \alpha_k \nabla f(y^{(k)}), \\ y^{(k+1)} = x^{(k+1)} + \frac{k}{k+3} (x^{(k+1)} - x^{(k)}). \end{cases}$$

Each iteration takes nearly same cost as GD but converge faster. Note that it is not a descent method.

¹i.e., to get $f(x^{(k)}) - f(x^*) \leq \epsilon$, we need $O(1/\epsilon)$ iterations.

1.4 Quasi-Newton methods

By successive measurements of the gradient, Quasi-Newton methods build a quadratic model of the objective function that can be used to estimate the inverse of the Hessian. The algorithm is the following.

Algorithm IV.1: Quasi-Newton algorithm

```

/* Initialization                                     */
Choose initial point  $\mathbf{x}^{(0)} \in \mathbb{R}^n$  and a matrix  $H_0 \succ 0$ 
/* Main loop                                           */
for  $k = 0, 1, 2, 3, \dots, N$  do
    /* compute quasi-Newton direction                 */
     $\Delta x_k = -H_k^{-1} \nabla f(x_k)$ 
    /* determine step size (e.g., by backtracking line search) */
    Compute  $t_k$ 
    /* compute  $x_{k+1}$                                    */
     $x_{k+1} = x_k + t_k + \Delta x_k$ 
    /* Update  $H$                                          */
    Compute  $H_{k+1}$ .
end

```

There are different update rules exist for H_{k+1} . One can also estimate H_k^{-1} or a factorization of H^k to simplify the computation of Δx_k . The BFGS method (Broyden-Fletcher-Goldfarb-Shanno) is the most popular quasi-Newton method.

BFGS update

$$H_{k+1} = H_k + \frac{yy^T}{y^T s} - \frac{H_k s s^T H_k}{s^T H_k s},$$

where

$$s = x_{k+1} - x_k, \quad y = \nabla f(x_{k+1}) - \nabla f(x_k).$$

Inverse update

$$H_{k+1}^{-1} = \left(I - \frac{sy^T}{y^T s} \right) H_k^{-1} \left(I - \frac{ys^T}{y^T s} \right) + \frac{ss^T}{y^T s}.$$

1.5 Nondifferentiable or constrained problems

For nondifferentiable or constrained problems, there are many methods that can be used, e.g., subgradient and proximal gradient methods.

2 Sub-gradient methods

Subgradients generalize gradients for optimizing convex functions that are not necessarily differentiable. In many situations, we deal with common operations that preserve convexity but not differentiability (e.g., the max-operation).

2.1 Sub-gradient

Recall that for differentiable f , that there is a single unique linear tangent $\nabla f(\mathbf{x})$ that under estimates the function (Figure IV.1):

$$f(\mathbf{y}) \geq f(\mathbf{x}) + \nabla f(\mathbf{x})^T(\mathbf{y} - \mathbf{x}) \quad \forall \mathbf{x}, \mathbf{y}.$$

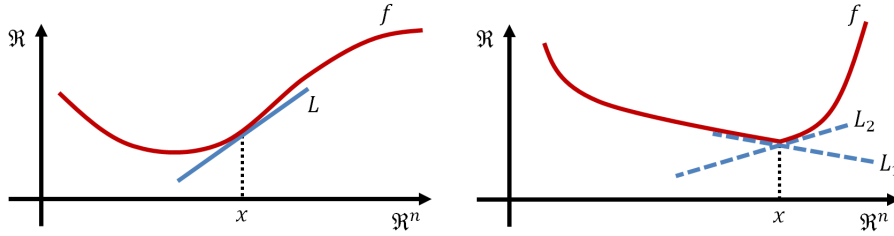


Fig. IV.1: Linear tangents under estimating a function.

Definition 2.1. A **subgradient** of a function f at \mathbf{x} is any $\mathbf{g} \in \mathbb{R}^n$ such that

$$f(\mathbf{y}) \geq f(\mathbf{x}) + \mathbf{g}^T(\mathbf{y} - \mathbf{x}), \quad \forall \mathbf{y}$$

The set of all subgradients of convex f is called the **subdifferential**:

$$\partial f(\mathbf{x}) \stackrel{\text{def}}{=} \left\{ \mathbf{g} \in \mathbb{R}^n ; \forall \mathbf{y}, f(\mathbf{y}) \geq f(\mathbf{x}) + \mathbf{g}^T(\mathbf{y} - \mathbf{x}) \right\}.$$

Remark: Informally, the “size” of $\partial f(\mathbf{x})$ controls how smooth f is at \mathbf{x} . For nonconvex f , subgradients need not exist.

Example 2.1 (absolute value). For $f(x) = |x|$, it is easy to show that for $x \neq 0$, there is a unique subgradient $g = \text{sign}(x)$. For $x = 0$, $\partial f(x) = [-1, 1]$.

Example 2.2 (Subdifferential of norms at 0). Let $f : \mathbb{R}^n \rightarrow \mathbb{R}$ be given by $f(\mathbf{x}) = \|\mathbf{x}\|$. We can show that:

$$\partial f(0) = \mathbb{B}_{\|\cdot\|_*}(0, 1) = \left\{ \mathbf{g} \in \mathbb{R}^n \mid \|\mathbf{g}\|_* \leq 1 \right\},$$

where $\|\cdot\|_*$ denotes the dual norm. In particular, for $g(x) = \|\mathbf{x}\|_1$, $\partial g(0) = [-1, 1]^n$.

Proposition 2.1 (Differentiability). f is differentiable at $\mathbf{x} \iff \partial f(\mathbf{x}) = \{\nabla f(\mathbf{x})\}$.

Example 2.3. Let $f(x) = |\cdot|, x \in \mathbb{R}$.

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

Example 2.4 (Subdifferential of the ℓ_2 -norm). Let $f(\mathbf{x}) = \|\mathbf{x}\|_2$.

$$\partial f(\mathbf{x}) = \begin{cases} \frac{\mathbf{x}}{\|\mathbf{x}\|_2} & \mathbf{x} \neq 0, \\ \mathbb{B}_2(0, 1) & \mathbf{x} = 0. \end{cases}$$

Theorem 2.2. A subdifferential $\partial f(x)$ is closed and convex, even for nonconvex f .

Theorem 2.3 (Existence). Let f be convex with a nonempty $\text{dom}(f)$. Then,

- (a) $x \in \text{relint}(\text{dom}(f)) \implies \partial f(x) \neq \emptyset$.
(b) $\partial f(x) \neq \emptyset$ and bounded $\iff x \in \text{int}(\text{dom}(f))$.

Example 2.5.

$$f(x) = \begin{cases} -\sqrt{x}, & x \geq 0 \\ +\infty, & x < 0 \end{cases}; \quad \partial f(x) = \begin{cases} -\frac{1}{2\sqrt{x}}, & x > 0 \\ \emptyset, & x = 0 \end{cases}$$

Note that 0 is not in the interior of $\text{dom}(f)$.

Theorem 2.4 (Directional Derivatives). Let $\mathbf{x} \in \text{int}(\text{dom}(f))$. Then, the directional derivative $f'(\mathbf{x}; \mathbf{d})$ is finite for all $\mathbf{d} \in \mathbb{R}^n$. In particular, we have

$$f'(\mathbf{x}; \mathbf{d}) = \max_{s \in \partial f(\mathbf{x})} s^T \mathbf{d}$$

Example 2.6 (connection to convex geometry). For a non empty set $C \subseteq \mathbb{R}^n$, $\partial \mathbb{1}_C(\mathbf{x}) = \mathcal{N}_C(\mathbf{x})$. Recall that for $\mathbf{x} \notin C$, $\partial \mathbb{1}_C(\mathbf{x}) = \mathcal{N}_C(\mathbf{x}) = \emptyset$ by convention. As a special case,

$$\partial \mathbb{1}_{\mathbb{B}(0,1)}(\mathbf{x}) = \mathcal{N}_{\mathbb{B}(0,1)}(\mathbf{x}) = \begin{cases} \left\{ \mathbf{y} \in \mathbb{R}^n \mid \|\mathbf{y}\|_* \leq \mathbf{y}^T \mathbf{x} \right\}, & \|\mathbf{x}\| \leq 1 \\ \emptyset, & \|\mathbf{x}\| > 1. \end{cases}$$

2.2 Subdifferential calculus

In this section, we present some practical rules for calculating subgradients. We can distinguish two types of rules giving weak (one subgradient only) or strong (all subgradients) results.

Rule 1: Conic combination

Theorem 2.5. Let $f, g : \mathbb{R}^n \rightarrow (-\infty, \infty]$ be proper functions and let $\alpha > 0$.

- Multiplication by a Positive Scalar

$$\forall \mathbf{x} \in \text{dom}(f) : \quad \partial(\alpha f(\mathbf{x})) = \alpha \partial f(\mathbf{x}).$$

- *Summation*

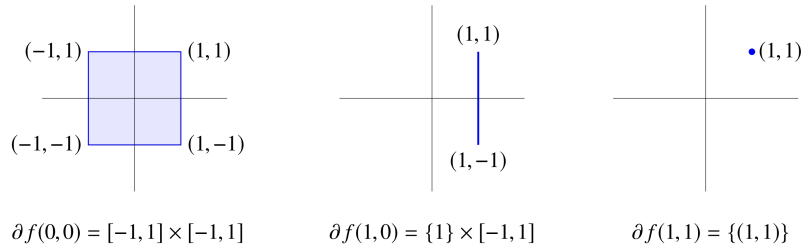
$$\forall \mathbf{x} \in \text{dom}(f) \cap \text{dom}(g) : \quad \partial f(\mathbf{x}) + \partial g(\mathbf{x}) \subseteq \partial(f+g)(\mathbf{x}),$$

$$\forall \mathbf{x} \in \text{int}(\text{dom}(f) \cap \text{dom}(g)) : \quad \partial(f+g)(\mathbf{x}) = \partial f(\mathbf{x}) + \partial g(\mathbf{x}).$$

The first summation rule is weak and the second is strong.

Example 2.7 (Subdifferentials of the ℓ_1 -norm). Consider $f : \mathbb{R}^n \rightarrow \mathbb{R}$ given by $f(\mathbf{x}) = \|\mathbf{x}\|_1 = \sum_{i=1}^n |x_i|$. Then, $f = \sum_{i=1}^n f_i$, where $f_i(\mathbf{x}) = |x_i|$. Hence,

$$\partial f(\mathbf{x}) = \left\{ \mathbf{z} \in \mathbb{R}^n \mid z_i = \text{sign}(x_i) \text{ if } x_i \neq 0, |z_j| \leq 1 \text{ if } x_j = 0 \right\}.$$



Rule 2: Affine composition

Theorem 2.6. Let $f : \mathbb{R}^n \rightarrow (-\infty, \infty]$ be a proper function and $h(\mathbf{x}) = f(A\mathbf{x} + \mathbf{b})$. Assume h is proper, and let $\alpha > 0$.

- *Weak result*

$$\forall \mathbf{x} \in \text{dom}(f) : \quad A^T(\partial f(A\mathbf{x} + \mathbf{b})) \subseteq \partial h(\mathbf{x}).$$

- *Strong result*

$$\forall \mathbf{x} \in \text{int}(\text{dom}(f)) / A(\mathbf{x}) + \mathbf{b} \in \text{int}(\text{dom}(f)) : \quad \partial h(\mathbf{x}) = A^T(\partial f(A\mathbf{x} + \mathbf{b})).$$

Example 2.8. Compute the subdifferential of $f(\mathbf{x}) = \|A\mathbf{x} + \mathbf{b}\|_1$.

$$\partial f(\mathbf{x}) = \sum_{i \in I_1} \text{sign}(\mathbf{a}_i^T \mathbf{x} + b_i) \mathbf{a}_i + \sum_{i \in I_2} [-\mathbf{a}_i, \mathbf{a}_i],$$

where \mathbf{a}_i is the i^{th} column of A , $I_1 = \{i : \mathbf{a}_i^T \mathbf{x} + b_i \neq 0\}$, $I_2 = \{i : \mathbf{a}_i^T \mathbf{x} + b_i = 0\}$. A weak result indicating one possible subgradient is $A^T \text{sign}(A\mathbf{x} + \mathbf{b}) \in \partial f(\mathbf{x})$.

Example 2.9. Compute the subdifferential of $\|A\mathbf{x} + \mathbf{b}\|_2$.

$$\partial f(\mathbf{x}) = \begin{cases} \frac{A^T(A\mathbf{x} + \mathbf{b})}{\|A\mathbf{x} + \mathbf{b}\|_2}, & A\mathbf{x} + \mathbf{b} \neq 0, \\ A^T \mathbb{B}_{\|\cdot\|_2}(0, 1), & A\mathbf{x} + \mathbf{b} = 0. \end{cases}$$

If a weak result is required, then $0 \in \partial f(\mathbf{x})$ for any \mathbf{x} satisfying $A\mathbf{x} + \mathbf{b} = 0$.

Rule 3: Chain rule

Theorem 2.7 (Sub-differential calculus). *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 2.10 (subdifferential of $\|\cdot\|_1^2$). *Consider the function $h : \mathbb{R}^n \rightarrow \mathbb{R}$ given by $h(\mathbf{x}) = \|\mathbf{x}\|_1^2$, which can be written as the composition $h = g \circ f$, where $f(\mathbf{x}) = \|\mathbf{x}\|_1$ and $g(t) = \max\{t, 0\}^2$. Both f and g are real-valued convex functions, and g is nondecreasing and differentiable over \mathbb{R} with derivative $g'(t) = 2\max\{t, 0\}$. Therefore, by the chain rule of subdifferential calculus, for any $x \in \mathbb{R}^n$,*

$$\partial h(\mathbf{x}) = g'(f(\mathbf{x}))\partial f(\mathbf{x}) = 2\max\{0, \|\mathbf{x}\|_1\}\partial f(\mathbf{x}),$$

$$\partial h(\mathbf{x}) = 2\|\mathbf{x}\|_1 \{z \in \mathbb{R}^n : z_i = \text{sign}(x_i), i \in I_1, |z_j| \leq 1, j \in I_2\},$$

where $I_1 = \{i : x_i \neq 0\}$ and $I_2 = \{i : x_i = 0\}$.

Example 2.11 (Distance to a convex set). *Recall the distance function to a closed and convex set C :*

$$d_C(\mathbf{x}) = \min_{\mathbf{y} \in C} \|\mathbf{y} - \mathbf{x}\|_2 = \|\mathbf{x} - P_C(\mathbf{x})\|_2$$

where $P_C(\mathbf{x})$ is the projection of \mathbf{x} onto C . We will show that

$$\partial d_C(\mathbf{x}) = \begin{cases} \left\{ \frac{\mathbf{x} - P_C(\mathbf{x})}{\|\mathbf{x} - P_C(\mathbf{x})\|_2} \right\}, & \mathbf{x} \notin C, \\ \mathcal{N}_C(\mathbf{x}) \cap \mathbb{B}(0, 1), & \mathbf{x} \in C. \end{cases}$$

We know that the function $\varphi_C(\mathbf{x}) = \frac{1}{2}d_C^2(\mathbf{x})$ is differentiable and

$$\forall \mathbf{x}, \quad \partial \varphi_C(\mathbf{x}) = \{\mathbf{x} - P_C(\mathbf{x})\}.$$

Note that $\varphi_C = g \circ d_C$, where $g(t) = \frac{1}{2}\max\{0, t\}^2$ is a nonincreasing real-valued convex differentiable function. Then by the chain rule:

$$\varphi_C(\mathbf{x}) = d_C(\mathbf{x}) \partial d_C(\mathbf{x}).$$

If $\mathbf{x} \notin C$, then $d_C(\mathbf{x}) = 0$, and thus $\partial d_C(\mathbf{x}) = \left\{ \frac{\mathbf{x} - P_C(\mathbf{x})}{d_C(\mathbf{x})} \right\}$. If $\mathbf{x} \in C$, we show the double inclusion.

Rule 4: Maximization

Theorem 2.8 (max rule of subdifferential calculus). *Let $f_1, f_2, \dots, f_m : \mathbb{R}^n \rightarrow (-\infty, \infty]$ be proper convex functions. Define $f(\mathbf{x}) = \max\{f_1(\mathbf{x}), \dots, f_m(\mathbf{x})\}$. Then,*

$$\forall \mathbf{x} \in \bigcap_{i=1, \dots, m} \text{dom}(f_i) : \quad \partial f(\mathbf{x}) = \text{conv} \left(\bigcup_{i \in I(\mathbf{x})} \partial f_i(\mathbf{x}) \right),$$

where $I(\mathbf{x}) = \{i \mid f_i(\mathbf{x}) = f(\mathbf{x})\}$, the "active" functions at \mathbf{x} . As a weak result; to compute a subgradient at \mathbf{x} , choose any $k \in I(\mathbf{x})$, any subgradient of f_k at \mathbf{x} .

Example 2.12 (max of two functions). Let $f_1, f_2 : \mathbb{R}^N \rightarrow \mathbb{R}$ convex, differentiable, and $f(\mathbf{x}) = \max\{f_1(\mathbf{x}), f_2(\mathbf{x})\}$. Then,

- For $f_1(\mathbf{x}) > f_2(\mathbf{x})$, unique subgradient $g = \nabla f_1(\mathbf{x})$
- For $f_2(\mathbf{x}) > f_1(\mathbf{x})$, unique subgradient $g = \nabla f_2(\mathbf{x})$
- For $f_1(\mathbf{x}) = f_2(\mathbf{x})$, $\partial f(\mathbf{x}) =$ the line segment joining $\nabla f_1(\mathbf{x})$ and $\nabla f_2(\mathbf{x})$.

Example 2.13 (subdifferential of the max function). Let $f : \mathbb{R}^n \rightarrow \mathbb{R}$ be given by $f(\mathbf{x}) = \max\{x_1, x_2, \dots, x_n\}$. Then we get:

$$\partial f(\mathbf{x}) = \text{conv} \left(\bigcup_{i \in I(\mathbf{x})} \partial f_i(\mathbf{x}) \right) = \text{conv} \left(\bigcup_{i \in I(\mathbf{x})} \{\mathbf{e}_i\} \right),$$

$$\partial f(\mathbf{x}) = \left\{ \sum_{i \in I(\mathbf{x})} \lambda_i \mathbf{e}_i : \sum_{i \in I(\mathbf{x})} \lambda_i = 1, \lambda_i \geq 0 \right\}.$$

Example 2.14 (subdifferential of the ℓ_∞ -norm). Let $f : \mathbb{R}^n \rightarrow \mathbb{R}$ be given by $f(\mathbf{x}) = \|\mathbf{x}\|_\infty$.

$$\partial f(\mathbf{x}) = \begin{cases} \mathbb{B}_{\|\cdot\|_1}(0, 1), & \mathbf{x} = 0, \\ \left\{ \sum_{i \in I(\mathbf{x})} \lambda_i \text{sign}(x_i) \mathbf{e}_i : \sum_{i \in I(\mathbf{x})} \lambda_i = 1, \lambda_j \geq 0, j \in I(\mathbf{x}) \right\}, & \mathbf{x} \neq 0. \end{cases}$$

Example 2.15. *piecewise-linear function*

$$f(\mathbf{x}) = \max_{i=1, \dots, m} a_i^T \mathbf{x} + b$$

the subdifferential at \mathbf{x} is a polyhedron

$$\partial f(\mathbf{x}) = \text{conv}\{a_i | i \in I(\mathbf{x})\}$$

Theorem 2.9 (weak maximum rule of subdifferential calculus). Let I be an arbitrary set, and suppose that any $i \in I$ is associated with a proper convex function $f_i : \mathbb{R}^n \rightarrow (-\infty, \infty]$. Let $f(\mathbf{x}) = \max_{i \in I} f_i(\mathbf{x})$. Then for any $\mathbf{x} \in \text{dom}(f)$:

$$\text{conv} \left(\bigcup_{i \in I(\mathbf{x})} \partial f_i(\mathbf{x}) \right) \subseteq \partial f(\mathbf{x}),$$

where $I(\mathbf{x}) = \{i \in I \mid f(\mathbf{x}) = f_i(\mathbf{x})\}$. Usually we get equality, but requires some technical conditions to hold.

Rule 5: Minimization

Theorem 2.10 (Sub-differential calculus). $f(x) = \inf_y h(x, y)$, h convex. To find a subgradient at \hat{x} (weak result), find \hat{y} that minimizes $h(\hat{x}, y)$ (assuming minimum is attained) and then find subgradient $(g, 0) \in \partial h(\hat{x}, \hat{y})$.

Example 2.16. Find a subgradient of

$$f(x) = \inf_{y \in \mathcal{C}} \|x - y\|_2$$

where \mathcal{C} is a closed convex set.

- if $f(\hat{x}) = 0$ (that is, $\hat{x} \in \mathcal{C}$), take $g = 0$.
- if $f(\hat{x}) > 0$, find projection $\hat{y} = P(\hat{x})$ on \mathcal{C} and take

$$g = \frac{\hat{x} - \hat{y}}{\|\hat{x} - \hat{y}\|_2}$$

2.3 Optimality conditions

Subdifferential sets are extremely useful in characterizing minima points. The following theorem is simple but extremely important.

Theorem 2.11 (Fermat's optimality condition). Let $f : \mathbb{R}^n \rightarrow (-\infty, \infty]$ be a proper convex function. Then,

$$\mathbf{x}^* \in \operatorname{argmin}_{\mathbf{x} \in \mathbb{R}^n} f(\mathbf{x}) \iff 0 \in \partial f(\mathbf{x}^*).$$

Proof. One has

$$\mathbf{x}^* \in \operatorname{argmin} f \iff (\forall \mathbf{y}, f(\mathbf{x}^*) \leq f(\mathbf{y}) + \mathbf{0}^T(\mathbf{x}^* - \mathbf{y})) \iff 0 \in \partial f(\mathbf{x}^*).$$

□

Example 2.17 (medians). Suppose that we are given n different 2D and ordered numbers $a_1 < a_2 < \dots < a_n$. Denote $A = \{a_1, a_2, \dots, a_n\} \subset \mathbb{R}$. The set of possible medians is the optimal solution set of the problem

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

where $f(x) = \sum_{i=1}^n |x - a_i| = f_i(x)$.

$$\partial f(x) = \begin{cases} 2i - n, & x \in (a_i, a_{i+1}), \\ 2i - 1 - n + [-1, 1], & x = a_i, \\ -n, & x < a_1, \\ n, & x > a_n. \end{cases}$$

$0 \in \partial f(a_i)$ if and only if $|2i - 1 - n| \leq 1$, which is equivalent to $\frac{n}{2} \leq i \leq \frac{n}{2} + 1$ and $0 \in \partial f(x)$ for some $x \in (a_i, a_{i+1})$ if and only if $i = \frac{n}{2}$. We can thus conclude that if n is odd, then the only optimal point is $\frac{a_{n+1}}{2}$ and when n is even, the optimal set is the interval $[a_{\frac{n}{2}}, a_{\frac{n}{2}+1}]$, establishing the fact that the optimal set is exactly the set of medians.

Theorem 2.12 (optimality conditions for composite problems). *Let $f, g : \mathbb{R}^n \rightarrow (-\infty, \infty]$ be proper functions such that $\text{dom}(g) \subset \text{int}(\text{dom}(f))$, g is convex. Consider the problem*

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

If $\mathbf{x}^ \in \text{dom}(g)$ is a local optimal solution for which f is differentiable, then*

$$-\nabla f(\mathbf{x}^*) \in \partial g(\mathbf{x}^*).$$

*Such points are called **stationary points**. This condition is sufficient when f is convex. In addition, the stationary points are global optimal solutions.*

Example 2.18 (Convex constraints). *For f convex and differentiable and convex set C :*

$$x^* \in \underset{x \in C}{\text{argmin}} f(x) \iff \exists g \in \partial f(x^*) / -g \in \mathcal{N}_C(x^*).$$

Example 2.19 (Lasso). *Given $y \in \mathbb{R}^n, X \in \mathbb{R}^{n \times p}$:*

$$\underset{\beta}{\text{argmin}} \frac{1}{2} \|y - X\beta\|_2^2 + \lambda \|\beta\|_1,$$

where $\lambda \geq 0$. Subgradient optimality:

$$0 \in \partial \left(\frac{1}{2} \|y - X\beta\|_2^2 + \lambda \|\beta\|_1 \right) \iff 0 \in -X^T(y - X\beta) + \lambda \partial \|\beta\|_1$$

$$0 \in \partial \left(\frac{1}{2} \|y - X\beta\|_2^2 + \lambda \|\beta\|_1 \right) \iff \exists v \in \partial \|\beta\|_1 : X^T(y - X\beta) = \lambda v,$$

$$v_i = \begin{cases} \{1\} & \text{if } \beta_i > 0 \\ \{-1\} & \text{if } \beta_i < 0, \\ = [-1, 1] & \text{if } \beta_i = 0 \end{cases} \quad i = 1, \dots, p$$

Write X_1, \dots, X_p for columns of X . Then our condition reads:

$$\begin{cases} X_i^T(y - X\beta) = \lambda \text{sign}(\beta_i) & \text{if } \beta_i \neq 0 \\ |X_i^T(y - X\beta)| \leq \lambda & \text{if } \beta_i = 0 \end{cases}$$

Note: subgradient optimality conditions don't lead to closed-form expression for a lasso solution. However, they do provide a way to check lasso optimality.

Example 2.20 (Soft-thresholding). *Simplified lasso problem with $X = I$:*

$$\underset{\beta}{\text{argmin}} \frac{1}{2} \|y - \beta\|_2^2 + \lambda \|\beta\|_1,$$

$$\beta = S_\lambda(y); \quad [S_\lambda(y)]_i = \begin{cases} y_i - \lambda & \text{if } y_i > \lambda \\ 0 & \text{if } -\lambda \leq y_i \leq \lambda \\ y_i + \lambda & \text{if } y_i < -\lambda \end{cases} \quad i = 1, \dots, n$$

Example 2.21. Consider the problem

$$\underset{x \in \mathbb{R}}{\text{minimize}} \quad f(x) + \lambda \|x\|_1 \quad (\text{IV.2})$$

where $f : \mathbb{R}^n \rightarrow (-\infty, \infty]$ is an extended real-valued function. A point $x^* \in \text{int}(\text{dom}(f))$ in which f is differentiable is a stationary point if $-\nabla f(x^*) \in \lambda \partial g(x^*)$, where $g(\cdot) = \|\cdot\|_1$. So, we obtain that x^* is a stationary point if

$$\frac{\partial f(x^*)}{\partial x_i} \begin{cases} = -\lambda & \text{if } x_i^* > 0 \\ = \lambda & \text{if } x_i^* < 0 \\ \in [-\lambda, \lambda] & \text{if } x_i^* = 0 \end{cases}$$

This is a necessary condition for x^* to be a local minimum. If f is also convex, then it is a necessary and sufficient condition for x^* to be a global optimal solution.

2.4 Sub-gradient descent

When f is non differentiable, one can use a sub-gradient instead of the gradient to define a minimization.

Algorithm IV.2: Sub-gradient descent

```

/* Main loop                                     */
for  $k = 0, 1, 2, 3, \dots, N$  do
    /* Main iteration                             */
     $\mathbf{x}^{(k+1)} = \mathbf{x}^{(k)} - \tau^{(k)} g^{(\ell)}$  where  $g^{(\ell)} \in \partial f(\mathbf{x}^{(\ell)})$ .
    /* Stopping criterion                         */
    if  $\|\mathbf{x}^{(k+1)} - \mathbf{x}^{(k)}\|_2 < \epsilon$  then
        Stop

```

But, subgradients are not necessarily descent directions as shown on the example below. For instance, for $f(\mathbf{x}) = |x_1| + 3|x_2|$ at $\mathbf{x} = (1, 0)$, we have:

- $\mathbf{g}_1 = (1, 0) \in \partial f(\mathbf{x})$, and $-\mathbf{g}_1$ is a descent direction
- $\mathbf{g}_2 = (1, 3) \in \partial f(\mathbf{x})$, but $-\mathbf{g}_2$ is not a descent direction

So we keep track of best iterate $x_{best}^{(k)}$ among $x^{(0)}, \dots, x^{(k)}$ so far, i.e., $f(x_{best}^{(k)}) = \min_{i=0, \dots, k} f(x^{(i)})$.

Step size rules

- fixed step: τ_k constant
- fixed length: $\tau^{(k)} g^{(\ell)} = \|\mathbf{x}^{(k+1)} - \mathbf{x}^{(k)}\|_2$ is constant
- diminishing: $\tau_k \rightarrow 0$ and $\sum_{k=0}^{\infty} \tau_k = \infty$.

Theorem 2.13. If $\sum_{\ell} \tau_{\ell} = +\infty$ and $\sum_{\ell} \tau_{\ell}^2 < +\infty$, then $x^{(\ell)}$ converges to a minimizer of f .

Analysis The subgradient method is not a descent method but the distance to the optimal goes to zero. As for the gradient method, the convergence is ensured when f is L-Lipschitz near the optimum. But, in general, this method performs poorly and convergence can be very slow. However, it handles general nondifferentiable convex problem, especially when ∂f is easy to compute, and often leads to very simple algorithms. Another option could be trying a splitting strategy as we will show later.

3 Proximal Algorithms

Proximal mapping are fundamental tools in convex optimization. Let's start with a proximal view of gradient descent.

3.1 Proximal Mapping

Consider the minimization problem:

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

We perform quadratic approximation, replacing usual Hessian $\nabla^2 f(x)$ by $\frac{1}{t}I$,

$$f(y) \approx f(x) + \nabla f(x)^T(y - x) + \frac{1}{2t}\|y - x\|.$$

This can be seen as a linear approximation of f for which we add a proximity term to x , with weight $1/(2t)$. So, the update rule can be written as

$$\mathbf{x}^{(k+1)} = \operatorname{argmin}_{\mathbf{x}} \left\{ f(\mathbf{x}^{(k)}) + \nabla f(\mathbf{x}^{(k)})^T(\mathbf{x} - \mathbf{x}^{(k)}) + \frac{1}{2t^{(k)}}\|\mathbf{x} - \mathbf{x}^{(k)}\|_2^2 \right\}.$$

Definition 3.2 (proximal mapping). *Given a function $f : \mathbb{R}^n \rightarrow (-\infty, \infty]$, the proximal mapping of f is the operator given by*

$$\operatorname{prox}_f(\mathbf{x}) \stackrel{\text{def.}}{=} \operatorname{argmin}_{\mathbf{z}} \frac{1}{2}\|\mathbf{x} - \mathbf{z}\|_2^2 + f(\mathbf{z}). \quad (\text{IV.3})$$

Example 3.1. $f(x) = \lambda|x|$. $\operatorname{prox}_f(x) = \begin{cases} x - \lambda, & x > \lambda \\ x + \lambda, & x < -\lambda \\ 0, & |x| \leq \lambda. \end{cases}$ *i.e., $\operatorname{prox}_f = S_\lambda$ the soft-thresholding operator.*

Example 3.2. *Consider the following functions from \mathbb{R} to \mathbb{R} , $\lambda > 0$ and $\mu \in \mathbb{R}$,*

$$g_1(x) = 0; \quad g_2(x) = \begin{cases} 0, & x \neq 0 \\ -\lambda, & x = 0 \end{cases}; \quad g_3(x) = \begin{cases} 0, & x \neq 0 \\ \lambda, & x = 0 \end{cases}; \quad g_4(x) = \begin{cases} \mu x, & x \geq 0 \\ \infty, & x < 0, \end{cases};$$

$$\operatorname{prox}_{g_1}(x) = \{x\}; \quad \operatorname{prox}_{g_2}(x) = \begin{cases} \{0\}, & |x| < \sqrt{2\lambda} \\ \{x\}, & |x| > \sqrt{2\lambda}, \end{cases}; \quad \operatorname{prox}_{g_3}(x) = \begin{cases} \{x\}, & x \neq 0 \\ \emptyset, & x = 0 \end{cases}$$

$$\operatorname{prox}_{g_4}(x) = 0.$$

Example 3.3 (Convex Quadratic). Let $f : \mathbb{R}^n \rightarrow \mathbb{R}$ be given by $f(\mathbf{x}) = \frac{1}{2} \mathbf{x}^T A \mathbf{x} + \mathbf{b}^T \mathbf{x} + c$, where $A \in \mathbb{S}_+^n$, $\mathbf{b} \in \mathbb{R}^n$, and $c \in \mathbb{R}$. Then,

$$\text{prox}_f(\mathbf{x}) = (A + I)^{-1}(\mathbf{x} - \mathbf{b}).$$

Theorem 3.14. If f is proper closed and convex, then $\text{prox}_f(x)$ is a singleton. It can be evaluated efficiently for many widely used functions (in particular, regularizers).

3.2 Prox Calculus Rules

Rule 1: Separable functions

Theorem 3.15. Suppose that $f : \mathbb{R}^{n_1} \times \dots \times \mathbb{R}^{n_m} \rightarrow (-\infty, \infty]$ is given by

$$f(\mathbf{x}_1, \dots, \mathbf{x}_m) = \sum_{i=1}^m f_i(\mathbf{x}_i), \quad \mathbf{x}_i \in \mathbb{R}^{n_i}$$

Then for any $(\mathbf{x}_1, \dots, \mathbf{x}_m) \in \mathbb{R}^{n_1} \times \dots \times \mathbb{R}^{n_m}$,

$$\text{prox}_f(\mathbf{x}_1, \dots, \mathbf{x}_m) = \text{prox}_{f_1}(\mathbf{x}_1) \times \dots \times \text{prox}_{f_m}(\mathbf{x}_m).$$

Example 3.4 (ℓ_1 -norm). Suppose that $g : \mathbb{R}^n \rightarrow \mathbb{R}$ is given by $g(x) = \lambda \|\mathbf{x}\|_1$, where $\lambda > 0$. Then, $g(\mathbf{x}) = \sum_{i=1}^n \varphi(x_i)$, and prox_g is the soft thresholding.

Rule 2: Scaling and translation

Theorem 3.16. Let $g : \mathbb{R}^n \rightarrow (-\infty, \infty]$ be a proper function, $\lambda \neq 0$, and $\mathbf{a} \in \mathbb{R}^n$. Define $f(\mathbf{x}) = g(\lambda \mathbf{x} + \mathbf{a})$ and $h(\mathbf{x}) = \lambda g(\mathbf{x}/\lambda)$. Then

$$\text{prox}_f(\mathbf{x}) = \frac{1}{\lambda}(\text{prox}_{\lambda^2 g}(\lambda \mathbf{x} + \mathbf{a}) - \mathbf{a}),$$

$$\text{prox}_h(\mathbf{x}) = \lambda(\text{prox}_{g/\lambda}(\mathbf{x}/\lambda)).$$

Example 3.5. Consider the function $f : \mathbb{R}^n \rightarrow (-\infty, \infty]$ given for any

$$f(x) = \begin{cases} \mu x, & 0 \leq x \leq \alpha, \\ \infty & \text{else,} \end{cases}$$

where $\mu \in \mathbb{R}$ and $\alpha \geq 0$.

$$\text{prox}_f(x) = \min\{\max\{x - \mu, 0\}, \alpha\}.$$

Rule 3: Norm composition

Theorem 3.17. Let $f : \mathbb{R}^m \rightarrow (-\infty, \infty]$ be given by $f(\mathbf{x}) = g(\|\mathbf{x}\|)$ where $g : \mathbb{R}^m \rightarrow (-\infty, \infty]$ is a proper closed convex function satisfying $\text{dom}(g) \subseteq [0, \infty)$. Then

$$\text{prox}_f(\mathbf{x}) = \begin{cases} \text{prox}_g(\|\mathbf{x}\|) \frac{\mathbf{x}}{\|\mathbf{x}\|}, & \mathbf{x} \neq 0, \\ \{u \in \mathbb{R}^n : \|u\| = \text{prox}_g(0)\}, & \mathbf{x} = 0. \end{cases}$$

Example 3.6 (prox of Euclidean norm). Let $f(\mathbf{x}) = \lambda \|\mathbf{x}\|$.

$$\text{prox}_f(\mathbf{x}) = \left(1 - \frac{\lambda}{\max\{\|\mathbf{x}\|, \lambda\}} \mathbf{x}\right)$$

Rule 4: Quadratic addition

Theorem 3.18 (Quadratic addition). *If $f(x) = g(x) + \frac{\rho}{2}\|x - a\|_2^2$, then*

$$\text{prox}_f(x) = \text{prox}_{\frac{1}{1+\rho}g}\left(\frac{1}{1+\rho}\mathbf{x} + \frac{\rho}{1+\rho}\mathbf{a}\right)$$

Theorem 3.19. *Properties*

- **Contraction:**

$$\|\text{prox}_f(x) - \text{prox}_f(y)\| \leq \|x - y\|.$$

- **Subgradient characterization:**

$$u = \text{prox}_f(x) \iff x - u \in \partial f(u)$$

- **Interpretation of prox via resolvent of subdifferential operator:**

$$z = \text{prox}_f(x) \iff z = \underbrace{(\text{Id} + \partial f^{-1})}_{\text{Resolvent of operator of } \partial f}(x)$$

3.3 Algorithms

3.3.1 Proximal Point Algorithm

One has the following equivalence

$$\begin{aligned} x^* \in \text{argmin } f &\iff 0 \in \partial f(x^*) \iff x^* \in (\text{Id} + \tau \partial f)(x^*) \\ &\iff x^* = (\text{Id} + \tau \partial f)^{-1}(x^*) = \text{Prox}_{\tau f}(x^*). \end{aligned}$$

This suggest the following *proximal point algorithm*.

$$\boxed{x^{(k+1)} = \text{Prox}_{\tau_\ell f}(x^{(k)})}.$$

On contrast to the gradient descent fixed point scheme, the proximal point method is converging for any sequence of steps.

Theorem 3.20. *If $0 < \tau_{\min} \leq \tau_\ell \leq \gamma_{\max} < +\infty$, then $x^{(\ell)} \rightarrow x^*$ a minimizer of f .*

3.3.2 Proximal gradient descent (Forward-Backward)

It is not always possible to compute the proximal operator of the objective function. We will consider here a class of problems by imposing some structure on the function to be minimized. We consider functions F of the form

$$\min_{\mathbf{x} \in \mathbb{R}^n} F(\mathbf{x}) := f(\mathbf{x}) + g(\mathbf{x}),$$

where f is convex differentiable and g is convex, not necessarily differentiable.

$$\begin{aligned}
x^* \in \operatorname{argmin} f + g &\Leftrightarrow 0 \in \nabla f(x^*) + \partial g(x^*) \\
&\Leftrightarrow x^* - \tau \nabla f(x^*) \in (\operatorname{Id} + \tau \partial g)(x^*) \\
&\Leftrightarrow x^* = (\operatorname{Id} + \tau \partial g)^{-1} \circ (\operatorname{Id} - \tau \nabla f)(x^*).
\end{aligned}$$

This fixed point suggests the following *Forward-Backward algorithm*.

$$\boxed{\mathbf{x}^{(k+1)} = \operatorname{prox}_{t_k g}(\mathbf{x}^{(k)} - t_k \nabla f(\mathbf{x}^{(k)}))}.$$

As for the gradient descent, we can apply a Nesterov acceleration to define the *Accelerated proximal gradient method*

$$\begin{aligned}
\boxed{v = x^{(k-1)} + \frac{k-2}{k+1} (x^{(k-1)} - x^{(k-2)})} \\
\boxed{x^{(k)} = \operatorname{Prox}_{t_k} (v - t_k \nabla g(x))}
\end{aligned}$$

Example 3.7 (Lasso). Given $y \in \mathbb{R}^n, X \in \mathbb{R}^{n \times p}$:

$$\min_{\beta} \underbrace{\frac{1}{2} \|y - X\beta\|_2^2}_{f(\beta)} + \underbrace{\lambda \|\beta\|_1}_{g(\beta)},$$

$$\operatorname{prox}_{t g}(x) = \operatorname{argmin}_z \frac{1}{2t} \|\beta - z\|_2^2 + \lambda \|z\|_1 = S_{\lambda t}(\beta)$$

Hence proximal gradient update is:

$$\beta^+ = S_{\lambda t}(\beta + t X^T (y - X\beta)) \quad (\text{IV.4})$$

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

Applying acceleration gives us ***FISTA*** (*F* is for *Fast*): (Beck and Teboulle '09)

$$v = \beta^{(k-1)} + \frac{k-2}{k+1} (\beta^{(k-1)} - \beta^{(k-2)})$$

$$\beta^{(k)} = \operatorname{prox}_{\lambda t_k} (v + t_k X^T (y - Xv)), \quad k = 1, 2, 3, \dots$$

or

$$\begin{aligned}
x^{(k+1)} &= \operatorname{prox}_{\eta_k h}(y^k - \eta_k \nabla f(y^k)) \\
y^{(k+1)} &= x^{(k+1)} + \frac{\theta_k - 1}{\theta_{k+1}} (x^{(k+1)} - x^{(k)}) \\
y^0 &= x^0, \theta_0 = 1, \theta_{k+1} = \frac{1 + \sqrt{1 + 4\theta_k^2}}{2}
\end{aligned}$$

Example 3.8 (Projected gradient descent). Given closed, convex set $C \in \mathbb{R}$,

$$\min_{x \in C} g(x) \iff \min_x g(x) + \mathbb{1}_C(x) \quad (\text{IV.5})$$

$$x^+ = P_C(x - t \nabla g(x)) \quad (\text{IV.6})$$

3.3.3 Douglas-Rachford

We reconsider here the structured minimization problem

$$\underset{\mathbf{x} \in \mathbb{R}^n}{\text{minimize}} \quad F(\mathbf{x}) := f(\mathbf{x}) + g(\mathbf{x}), \quad (\text{IV.7})$$

where f, g are closed convex, but on contrary to the Forward-Backward, no smoothness is imposed on f . We here suppose that we can compute easily the proximal map of f and g .

Algorithm IV.3: Douglas-Rachford method

```

Data:  $\min_x f(\mathbf{x}) + g(\mathbf{x})$ 
/* Main loop */
for  $k = 0, 1, 2, 3, \dots, N$  do
    /* Douglas-Rachford iterations */
     $u^{(k+1)} = \text{prox}_g(x^{(k)} + w^{(k)})$ 
     $x^{(k+1)} = \text{prox}_f(u^{(k+1)} - w^{(k)})$ 
     $w^{(k+1)} = w^{(k)} + (x^{(k+1)} - u^{(k+1)})$ 
    /* Stopping criterion */
    if  $\|x^{(k+1)} - x^{(k)}\|_2 < \varepsilon$  then
         $\perp$  Stop

```

Note that it is of course possible to inter-change the roles of f and g , which defines another set of iterations.