

Lecture 12

12.1 Problem Formulation	1
12.2 Subgradient Method with Normalized Stepsizes	2
12.3 Functional Growth	5

12.1 Problem Formulation

We consider the following convex optimization problem:

$$\min_{x \in Q} f(x), \tag{12.1}$$

for a convex set $Q \subseteq \mathbb{R}^n$, and a convex (possibly non-differentiable) function $f : \mathbb{R}^n \rightarrow \mathbb{R}$, which, for simplicity of the presentation, we define over the whole space. In modern applications, the dimension n in problem (12.1) is large: $n \rightarrow \infty$ (*large-scale optimization*), so we cannot directly apply the heavy machinery of the ellipsoid method. Instead, we will analyze a non-smooth analog of the gradient method.

We assume that for any point $x \in Q$ we can compute a subgradient vector $f'(x) \in \partial f(x)$ that satisfies:

$$f(y) \geq f(x) + \langle f'(x), y - x \rangle, \quad x, y \in Q.$$

For the set Q we assume a possibility of computing the projection, in the Euclidean norm:

$$\pi_Q(x) := \arg \min_{y \in Q} \|y - x\|. \tag{12.2}$$

This is a different and more expensive operation than a separation oracle for Q . The possibility of computing projections (12.2) typically means that the set Q is *simple*. Thus, as we did when discussing the fully composite problems in Lecture 9, we assume that the main difficulty of solving problem (12.1) lies in the objective function f , but not in the constraints. However, it is possible to generalize the subgradient method to the case of using only a separation oracle for Q , when the set is specified by a number of black-box functional inequalities.

12.1.1 Optimality Condition

Let us discuss an optimality condition for a point x^* to be a global minimum of (12.1). We have the following generalization of the first-order condition from unconstrained minimization. Let $Df(x^*)[h]$ denote the directional derivative of f along the direction h :

$$Df(x)[h] := \lim_{\alpha \rightarrow 0^+} \frac{1}{\alpha} [f(x + \alpha h) - f(x)].$$

- For differentiable functions: $Df(x)[h] \equiv \langle \nabla f(x), h \rangle$.
- For convex functions, the directional derivative along any direction exists at all interior points of the domain. In this case, we have the following interesting relationship:

$$Df(x)[h] = \max \{ \langle g, h \rangle : g \in \partial f(x) \}.$$

Proposition 12.1.1. *Let x^* be a constrained minimum of f over set Q . Then*

$$Df(x^*)[x - x^*] \geq 0, \quad x \in Q. \quad (12.3)$$

Proof. Indeed, by the definition of the directional derivative, we have for a sufficiently small $\alpha > 0$:

$$Df(x^*)[x - x^*] = \frac{1}{\alpha}[f(x + \alpha h) - f(x)] + o(1) \geq o(1),$$

where $o(1)$ goes to zero when $\alpha \rightarrow 0$. Taking the limit complete the proof. \square

Corollary 12.1.2. *For a constrained minimum of a differentiable function f , we have:*

$$\langle \nabla f(x^*), x - x^* \rangle \geq 0, \quad x \in Q. \quad (12.4)$$

Remark 12.1.3. See also Theorem 10.1.5 in Lecture 10 for a more general optimality condition suitable for additive composite optimization, that generalizes (12.4).

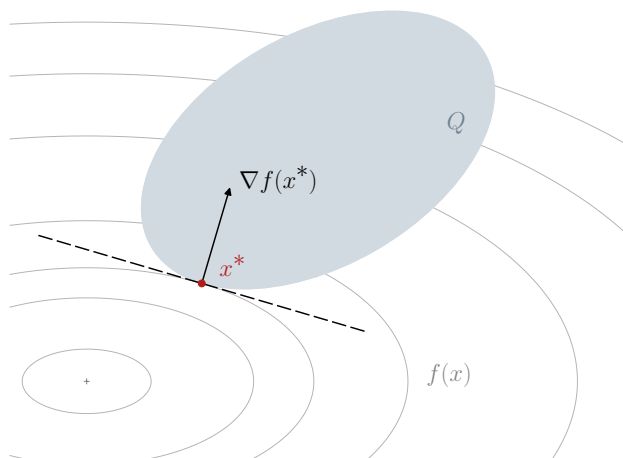


Figure 12.1: The constrained minimum of a function f over a set Q .

12.2 Subgradient Method with Normalized Stepsizes

12.2.1 Subgradient Step

We will use the optimality condition to analyze one step of the subgradient method. The method is very similar to the basic gradient method. At point $x \in Q$, we compute a subgradient $f'(x) \in \partial f(x)$ and perform the following update, for some $\eta > 0$:

$$x^+ = \arg \min_{y \in Q} \left[m(y) := \eta \langle f'(x), y - x \rangle + \frac{1}{2} \|y - x\|^2 \right]. \quad (12.5)$$

Note that the function $m(\cdot)$ is differentiable and 1-strongly convex. Hence, we have

$$\begin{aligned} m(y) &\geq m(x^+) + \langle \nabla m(x^+), y - x^+ \rangle + \frac{1}{2} \|y - x^+\|^2 \\ &\stackrel{(12.4)}{\geq} m(x^+) + \frac{1}{2} \|y - x^+\|^2. \end{aligned}$$

Expanding the definition of the model gives:

$$\begin{aligned} \frac{1}{2}\|y - x\|^2 + \eta\langle f'(x), y - x \rangle &\geq \frac{1}{2}\|y - x^+\|^2 + \left[\eta\langle f'(x), x^+ - x \rangle + \frac{1}{2}\|x^+ - x\|^2 \right] \\ &\geq \frac{1}{2}\|y - x^+\|^2 - \frac{\eta^2}{2}\|f'(x)\|^2. \end{aligned}$$

Thus, we have proved the following important lemma, which is in the core of analysis of all subgradient methods.

Lemma 12.2.1. *For any $y \in Q$, it holds*

$$\frac{1}{2}\|y - x\|^2 + \frac{\eta^2}{2}\|f'(x)\|^2 \geq \frac{1}{2}\|y - x^+\|^2 + \eta\langle f'(x), x - y \rangle. \quad (12.6)$$

Consequently, substituting $y := x^* \in Q$ (any minimizer), we get:

$$\frac{1}{2}\|x^* - x\|^2 + \frac{\eta^2}{2}\|f'(x)\|^2 \geq \frac{1}{2}\|x^* - x^+\|^2 + \eta\langle f'(x), x - x^* \rangle. \quad (12.7)$$

Note that for convex functions, the inner product in the right hand side of (12.7) is non-negative, and it is bounded by the functional residual:

$$\langle f'(x), x - x^* \rangle \geq f(x) - f^* \geq 0.$$

This reasoning can immediately lead to a convergence rate for the method. But first, let us choose stepsizes in a smart way.

12.2.2 Algorithm

In non-smooth optimization, the main information that is provided by a subgradient $f'(x) \in \partial f(x)$ is of *geometric nature*, as it gives us a certain *separation of the sublevel set* of f at x . At the same time, in contrast to smooth optimization, the magnitude $\|f'(x)\|$ does not reveal much information.

Example 12.2.2. Consider $f(x) = |x|$, $x \in \mathbb{R}$. Then, for any $x \neq 0$, we have $|f'(x)| = 1$, no matter how close we are to the optimum $x^* = 0$.

Therefore, it is natural to *normalize* the subgradient direction: $\frac{f'(x)}{\|f'(x)\|}$, which happens also to equip our method with *universal* (problem-class independent) convergence rates. We consider the following algorithm.

Algorithm 12.1: Subgradient Method.

Initialization: $x_0 \in Q$. Fix $K \geq 1$ and positive parameters $\{\gamma_k\}_{k \geq 0}$.

For $k = 0 \dots K - 1$ iterate:

1. Compute a subgradient: $f'(x_k) \in \partial f(x_k)$

2. Perform the normalized subgradient step:

$$x_{k+1} = \arg \min_{y \in Q} \left[\frac{\gamma_k}{\|f'(x_k)\|} \langle f'(x_k), y - x_k \rangle + \frac{1}{2}\|y - x_k\|^2 \right] = \pi_Q \left(x_k - \frac{\gamma_k}{\|f'(x_k)\|} f'(x_k) \right)$$

Return a point $\bar{x}_K := \arg \min\{f(y) : y \in \{x_0, \dots, x_K\}\}$ or the average $\frac{1}{K} \sum_{i=0}^{K-1} x_i$.

Substituting our stepsize choice in the previous lemma, we obtain, for every $k \geq 0$:

$$\frac{\gamma_k^2}{2} + \frac{1}{2}\|x_k - x^*\|^2 \geq \frac{1}{2}\|x_{k+1} - x^*\|^2 + \gamma_k \Delta_k, \quad (12.8)$$

where

$$\Delta_k := \frac{\langle f'(x_k), x_k - x^* \rangle}{\|f'(x_k)\|}$$

is a certain *measure of optimality*. We also denote:

$$\bar{\Delta}_K := \frac{1}{K} \sum_{i=0}^{K-1} \Delta_i \quad \text{and} \quad \Delta_K^* := \min_{0 \leq i \leq K-1} \Delta_i.$$

Clearly, we have $\bar{\Delta}_K \geq \Delta_K^*$.

Telescoping inequality (12.8) for the first $K \geq 1$ iterations of the method, we get the following progress:

$$\frac{1}{2}\|x_0 - x^*\|^2 + \frac{1}{2} \sum_{i=0}^{K-1} \gamma_i^2 \geq \sum_{i=0}^{K-1} \gamma_i \Delta_i \geq \left(\sum_{i=0}^{K-1} \gamma_i \right) \Delta_K^* \quad (12.9)$$

Therefore, we obtain the following bound on our new accuracy measure:

$$\Delta_K^* \leq \frac{\|x_0 - x^*\|^2 + \sum_{i=0}^{K-1} \gamma_i^2}{2 \sum_{i=0}^{K-1} \gamma_i} = \varphi(\gamma_0, \dots, \gamma_{K-1}). \quad (12.10)$$

In principle, we want to choose $\{\gamma_k\}_{k \geq 0}$ such that the right hand side of (12.10) is as small as possible, i.e. to minimize it: $\varphi(\cdot) \rightarrow \min$. Notice that

- $\varphi(\cdot)$ is convex;
- $\varphi(\cdot)$ is symmetric in γ : its value is invariant to any permutation of $\{\gamma_0, \dots, \gamma_{K-1}\}$.

For a convex symmetric function, there is always exists a solution with all the same arguments,

$$\gamma_0^* = \gamma_1^* = \dots = \gamma_{k-1}^* = \gamma.$$

Therefore, a constant choice $\gamma > 0$ is able to give us the best bound in (12.10), when the number of iterations K is fixed. In practice, however, we might want to use a decreasing sequence, e.g. $\gamma_k = O(1/\sqrt{k})$, so not to fix K .

We denote by $R \geq \|x_0 - x^*\|$ any upper bound for the distance from the initial point to any of the solutions. Using this bound in (12.9) along with the constant choice, $\gamma_i \equiv \gamma > 0$, leads to

$$\bar{\Delta}_K \stackrel{(12.9)}{\leq} \frac{\|x_0 - x^*\|^2}{2K\gamma} + \frac{K\gamma}{2} \leq \frac{R^2}{2K} + \frac{K\gamma}{2}.$$

Minimizing the right-hand side in $\gamma > 0$, we obtain the optimal choice

$$\boxed{\gamma := \frac{R}{K^{1/2}}}. \quad (12.11)$$

Thus, we have proved the following result.

Theorem 12.2.3. *Let all γ_k in Algorithm 12.1 be chosen according to (12.11). Then,*

$$\Delta_K^* \leq \bar{\Delta}_K \leq \frac{R}{K^{1/2}}. \quad (12.12)$$

12.2.3 Lipschitz Functions

The question is how to relate our quantities Δ_k with a standard accuracy measure, the functional residual $f(x_k) - f^*$. The simplest reasoning is as follows.

Assume that our subgradients are bounded:

$$\|f'(x)\| \leq M, \quad \forall x \in Q, \quad (12.13)$$

which means that the function f is Lipschitz:

$$|f(y) - f(x)| \leq M\|y - x\|, \quad x, y \in Q.$$

Then, by convexity we immediately obtain the following relationship:

$$\Delta_k := \frac{\langle f'(x_k), x_k - x^* \rangle}{\|f'(x_k)\|} \stackrel{(12.13)}{\geq} \frac{\langle f'(x_k), x_k - x^* \rangle}{M} \geq \frac{f(x_k) - f^*}{M}. \quad (12.14)$$

Corollary 12.2.4. *We have,*

$$f\left(\frac{1}{K} \sum_{i=0}^{K-1} x_i\right) - f^* \leq \frac{1}{K} \sum_{i=0}^{K-1} [f(x_i) - f^*] \stackrel{(12.14)}{\leq} M \cdot \bar{\Delta}_K \leq \frac{MR}{K^{1/2}}. \quad (12.15)$$

Note that we used a conservative choice of stepsizes that requires fixing the number of iterations. Therefore, technically, (12.15) is not a “rate” of convergence, as this bound is achieved only once, for the final output of the algorithm. If we wanted to achieve higher precision, we would need to rerun the method with a smaller stepsize γ . In Lecture 14, we discuss a more advanced approach using *adaptive stepsizes* that solves this issue, as they do not require fixing the number of iterations in advance, and they also work for stochastic problems.

We also assume we know a bound $R \geq \|x_0 - x^*\|$. At the same time, it is easy to see that any choice of $R > 0$ in (12.11) will give us a similar bound. For example, in practice, we can set $R := 1$ if no other information is available. However, the closer R is to $\|x_0 - x^*\|$, the faster the method converges.

In the next lecture, we prove that the rate $O\left(\frac{MR}{K^{1/2}}\right)$ is *optimal* for this problem class.

12.3 Functional Growth

Note that to prove our bound (12.15) for the subgradient method on convex Lipschitz functions, it is enough to take much more conservative steps, using the same constant for each iteration in (12.5):

$$\eta := \frac{\gamma}{M} \stackrel{(12.11)}{=} \frac{R}{MK^{1/2}}, \quad (12.16)$$

in contrast to using normalized stepsizes.

Exercise 12.3.1. Check that subgradient steps (12.5) with η given by (12.16) achieve the same bound (12.15) for the functional residual.

In this section, we provide a more advanced analysis that reveals the true power of normalized stepsizes: they enable us to prove convergence rates for problems well beyond the Lipschitz assumption.

12.3.1 Geometric Interpretation

First, let us understand the geometric meaning of the quantity $\Delta_k := \frac{\langle f'(x_k), x_k - x^* \rangle}{\|f'(x_k)\|}$. For simplicity, consider the case of unconstrained optimization, thus

$$Q \equiv \mathbb{R}^n.$$

Staying at point x_k , the subgradient $f'(x_k) \neq 0$ provides us with the supporting hyperplane:

$$L_k = \left\{ y \in \mathbb{R}^n : \langle f'(x_k), x_k - y \rangle = 0 \right\}.$$

Let us look at the optimal solution $x^* \in \mathbb{R}^n$ and find the projection of it onto L_k :

$$\min_{y \in L_k} \|y - x^*\|. \quad (12.17)$$

We take vector $h := \Delta_k \frac{f'(x_k)}{\|f'(x_k)\|}$ and the perturbed solution $y^* := x^* + h$. We have

$$\begin{aligned} \langle f'(x_k), x_k - y^* \rangle &= \langle f'(x_k), x_k - x^* \rangle + \langle f'(x_k), h \rangle \\ &= \Delta_k \|f'(x_k)\| - \Delta_k \frac{\langle f'(x_k), f'(x_k) \rangle}{\|f'(x_k)\|} = 0, \end{aligned}$$

which concludes that $y^* \in L_k$ is the solution to (12.17). Note that

$$\|y^* - x^*\| = \|h\| = \Delta_k.$$

Corollary 12.3.1. Δ_k is the distance from x^* to the hyperplane L_k .

Consider the localizing polyhedron:

$$G_{k+1} = \left\{ y \in \mathbb{R}^n : \langle f'(x_0), x_0 - y \rangle \geq 0, \dots, \langle f'(x_k), x_k - y \rangle \geq 0 \right\}.$$

By convexity $x^* \in G_{k+1}$. Note that Δ_k^* is the minimal distance from x^* to one of the hyperplanes defining the polyhedron. Hence, Δ_k^* is the maximal radius of the Euclidan ball that is contained in the localizer G_{k+1} :

$$\Delta_k^* = \max \left\{ r \geq 0 : B_r(x^*) \subset G_{k+1} \right\}$$

and the subgradient method manages to $\Delta_k^* \rightarrow 0$. Note that such quantity can be used to define an appropriate $\text{size}(\cdot)$ of a set (see previous lecture).

12.3.2 Convergence Rate with Functional Growth

Define the following quantity, called *the growth of f at point x^** :

$$\omega_f(r) := \max \left\{ f(x) - f(x^*) : \|x - x^*\| \leq r \right\}.$$

Clearly, $\omega_f(\cdot)$ is a nondecreasing function of r .

We can relate our measure of optimality Δ_k with the functional residual. By convexity, we know that

$$\Delta_k = \frac{\langle f'(x_k), x_k - x^* \rangle}{\|f'(x_k)\|} \geq \frac{f(x_k) - f(x^*)}{\|f'(x_k)\|}.$$

Now we present a more advanced reasoning. Consider the perturbation, as we fixed before:

$$h := \frac{\Delta_k}{\|f'(x_k)\|} f'(x_k).$$

Then,

$$\begin{aligned} \langle f'(x_k), x_k - x^* \rangle &= \langle f'(x_k), x_k - x^* - h \rangle + \langle f'(x_k), h \rangle \\ &\geq f(x_k) - f(x^* + h) + \|f'(x_k)\| \Delta_k \\ &= f(x_k) - f^* + (f^* - f(x^* + h)) + \|f'(x_k)\| \Delta_k. \end{aligned}$$

Note that $\langle f'(x_k), x_k - x^* \rangle = \|f'(x_k)\| \Delta_k$. Rearranging the terms, we obtain

$$f(x_k) - f^* \leq f(x^* + h) - f^*.$$

This inequality has a clear geometric meaning: due to convexity, the function value at $f(x_k)$ is better than that at $f(x^* + h)$, as shown on Fig. 12.2. We established the following result.

Theorem 12.3.2. *For the result of the subgradient method, it holds:*

$$f(\bar{x}_K) - f^* \leq \omega_f(\Delta_K^*), \quad \text{and} \quad \Delta_K^* \leq \frac{R}{K^{1/2}}. \quad (12.18)$$

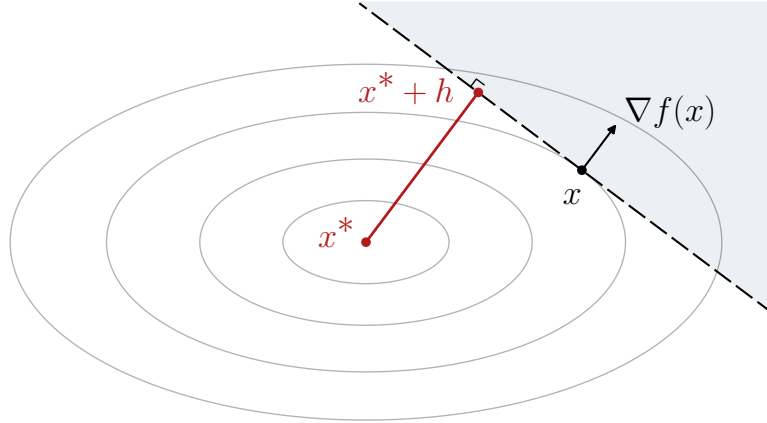


Figure 12.2: Point $x^* + h$ is the projection of the unconstrained minimum x^* to the supporting hyperplane defined by a subgradient at x .

12.3.3 Example: Lipschitz Functions

Proposition 12.3.3. *Consider Lipschitz functions: $f(x + h) - f(x) \leq M\|h\|$. Then, we have*

$$\omega_f(r) \leq Mr. \quad (12.19)$$

Hence, we immediately recover the bound from the previous reasoning:

$$f(\bar{x}_K) - f^* \stackrel{(12.18)}{\leq} \omega_f\left(\frac{R}{K^{1/2}}\right) \stackrel{(12.19)}{\leq} \frac{MR}{K^{1/2}}. \quad (12.20)$$

To obtain an ε -solution we need $\left\lceil \frac{MR}{\varepsilon} \right\rceil^2$ oracle calls.

Let us consider the following objective:

$$f(x) = \max_{1 \leq i \leq m} [\langle a_i, x \rangle - b_i].$$

Then

$$\begin{aligned} f(y) - f(x) &= \max_{1 \leq i \leq m} [\langle a_i, x \rangle - b_i] - \max_{1 \leq j \leq m} [\langle a_j, y \rangle - b_j] \\ &= \max_{1 \leq i \leq m} [\langle a_i, x \rangle - b_i - \max_{1 \leq j \leq m} [\langle a_j, y \rangle - b_j]] \\ &\leq \max_{1 \leq i \leq m} [\langle a_i, x - y \rangle] \leq \max_{1 \leq i \leq m} \|a_i\| \cdot \|x - y\|. \end{aligned}$$

Hence, $M = \max_{1 \leq i \leq m} \|a_i\|$.

At every point x , a subgradient can be computed as $f'(x) := a_i$, where i is the active index.

12.3.4 Example: Smooth Functions

Proposition 12.3.4. *Let f be differentiable and smooth. Then,*

$$\begin{aligned} \omega_f(r) &= \max\{f(x^* + h) - f(x^*) : \|h\| \leq r\} \\ &\leq \max\{\langle \nabla f(x^*), h \rangle + \frac{L}{2} \|h\|^2 : \|h\| \leq r\} \\ &\leq \|\nabla f(x^*)\| r + \frac{L}{2} r^2. \end{aligned} \tag{12.21}$$

Substituting this value into our bound, we get:

$$f(\bar{x}_K) - f^* \stackrel{(12.18)}{\leq} \omega_f\left(\frac{R}{K^{1/2}}\right) \stackrel{(12.21)}{\leq} \frac{\|\nabla f(x^*)\| R}{K^{1/2}} + \frac{L R^2}{2K}.$$

Hence, in case of smooth functions and small $\|\nabla f(x^*)\|$, our method automatically receives a faster rate of convergence than that one from (12.20). For unconstrained minimization, $\nabla f(x^*) = 0$.

12.3.5 Example: Maximum of Smooth Functions

Let

$$f(x) = \max_{1 \leq i \leq m} [f_i(x)], \tag{12.22}$$

where each $f_i(x) = \frac{1}{2} \langle A_i x, x \rangle - \langle b_i, x \rangle$ is a convex quadratic function, or, more generally, each f_i is convex and has the Lipschitz gradient with constant L_i .

Then,

$$\begin{aligned} f_i(x) &\leq f_i(x^*) + \langle \nabla f_i(x^*), x - x^* \rangle + \frac{L_i}{2} \|x - x^*\|^2 \\ &\leq f_i(x^*) + \|\nabla f_i(x^*)\| \cdot \|x - x^*\| + \frac{L_i}{2} \|x - x^*\|^2. \end{aligned}$$

And we obtain that

$$\omega_f(r) \leq \max_{1 \leq i \leq m} \|\nabla f_i(x^*)\| \cdot r + \max_{1 \leq i \leq m} L_i \cdot \frac{r^2}{2}. \tag{12.23}$$

Corollary 12.3.5. *It holds:*

$$f(\bar{x}_K) - f^* \stackrel{(12.18)}{\leq} \omega_f\left(\frac{R}{K^{1/2}}\right) \stackrel{(12.23)}{\leq} \frac{\max_{1 \leq i \leq m} \|\nabla f_i(x^*)\| R}{K^{1/2}} + \frac{\max_{1 \leq i \leq m} L_i R^2}{2K}.$$

We see that exactly the same subgradient method with normalized steps will work on different subclasses of non-smooth convex problems, even though the objective in (12.22) is not Lipschitz.

Literature

The subgradient method was discovered by Naum Shor in 1962 [Sho12].

Another powerful selection of stepsizes for the subgradient method includes Polyak’s stepsizes [Pol87, HK19] and their recent modification suitable for constrained and composite optimization [Nes24].

See also Section 7 in [Nem95] and Section 3.2 in [Nes18] for the design of subgradient methods for constrained problems with functional inequalities.

An advanced analysis of the convergence rate for the last iterate x_k of the subgradient method, instead of \bar{x}_k , was developed recently in [ZG25].

- [HK19] Elad Hazan and Sham Kakade. Revisiting the polyak step size. *arXiv preprint arXiv:1905.00313*, 2019.
- [Nem95] Arkadi Nemirovski. *Information-based complexity of convex programming*. Lecture notes, 1995.
- [Nes18] Yurii Nesterov. *Lectures on convex optimization*. Springer, 2018.
- [Nes24] Yurii Nesterov. Primal subgradient methods with predefined step sizes. *Journal of Optimization Theory and Applications*, 203(3):2083–2115, 2024.
- [Pol87] Boris T Polyak. *Introduction to optimization*. 1987.
- [Sho12] Naum Zuselevich Shor. *Minimization methods for non-differentiable functions*. Springer Science & Business Media, 2012.
- [ZG25] Moslem Zamani and François Glineur. Exact convergence rate of the last iterate in subgradient methods. *SIAM Journal on Optimization*, 35(3):2182–2201, 2025.