

5 Algorithms for Minimizing Unconstrained Functions

5.1 General Minimization Problem and Terminologies

Definition 5.1 We define the *general minimization problem* as follows

$$\begin{cases} \text{minimize} & f(\mathbf{x}) \\ \text{subject to} & f_j(\mathbf{x}) \ \& \ 0, \quad j = 1, 2, \dots, m \\ & \mathbf{x} \in S, \end{cases} \quad (3)$$

where $f : \mathbb{R}^n \rightarrow \mathbb{R}$, $f_j : \mathbb{R}^n \rightarrow \mathbb{R}$ ($j = 1, 2, \dots, m$), the symbol $\&$ could be $=$, \geq , or \leq , and $S \subseteq \mathbb{R}^n$.

Definition 5.2 The *feasible set* Q of (3) is

$$Q = \{\mathbf{x} \in S \mid f_j(\mathbf{x}) \ \& \ 0, \ (j = 1, 2, \dots, m)\}.$$

In the following items we assume $S \equiv \mathbb{R}^n$.

- If $Q \equiv \mathbb{R}^n$, (3) is a *unconstrained optimization problem*.
- If $Q \subsetneq \mathbb{R}^n$, (3) is a *constrained optimization problem*.
- If all functionals $f(\mathbf{x})$, $f_j(\mathbf{x})$ are differentiable, (3) is a *smooth optimization problem*.
- If one of functionals $f(\mathbf{x})$, $f_j(\mathbf{x})$ is non-differentiable, (3) is a *non-smooth optimization problem*.
- If all constraints are linear $f_j(\mathbf{x}) = \langle \mathbf{a}_j, \mathbf{x} \rangle + b_j$ ($j = 1, 2, \dots, m$), (3) is a *linear constrained optimization problem*.
 - In addition, if $f(\mathbf{x})$ is linear, (3) is a *linear programming problem*.
 - In addition, if $f(\mathbf{x})$ is quadratic, (3) is a *quadratic programming problem*.
- If $f(\mathbf{x})$, $f_j(\mathbf{x})$ ($j = 1, 2, \dots, m$) are quadratic, (3) is a *quadratically constrained quadratic programming problem*.

Definition 5.3 \mathbf{x}^* is called a *global optimal solution* of (3) if $f(\mathbf{x}^*) \leq f(\mathbf{x})$, $\forall \mathbf{x} \in Q$. Moreover, $f(\mathbf{x}^*)$ is called the *global optimal value*. \mathbf{x}^* is called a *local optimal solution* of (3) if there exists an open ball $B(\mathbf{x}^*, \varepsilon) := \{\mathbf{x} \in \mathbb{R}^n \mid \|\mathbf{x} - \mathbf{x}^*\|_2 < \varepsilon\}$ such that $f(\mathbf{x}^*) \leq f(\mathbf{x})$, $\forall \mathbf{x} \in B(\mathbf{x}^*, \varepsilon) \cap Q$. Moreover, $f(\mathbf{x}^*)$ is called a *local optimal value*.

5.2 Complexity Bound for a Global Optimization Problem on the Unit Box

Consider one of the simplest problems in optimization, that is, minimizing a function in the n -dimensional box.

$$\begin{cases} \text{minimize} & f(\mathbf{x}) \\ \text{subject to} & \mathbf{x} \in B_n := \{\mathbf{x} \in \mathbb{R}^n \mid 0 \leq [\mathbf{x}]_i \leq 1, \ i = 1, 2, \dots, n\}. \end{cases} \quad (4)$$

To be coherent, we use the ℓ_∞ -norm:

$$\|\mathbf{x}\|_\infty = \max_{1 \leq i \leq n} |[\mathbf{x}]_i|.$$

Let us also assume that $f(\mathbf{x})$ is *Lipschitz continuous* on B_n :

$$|f(\mathbf{x}) - f(\mathbf{y})| \leq L\|\mathbf{x} - \mathbf{y}\|_\infty, \quad \forall \mathbf{x}, \mathbf{y} \in B_n.$$

Let us define a very simple method to solve (4), the **uniform grid method**.

Given a positive integer $p > 0$,

1. Form $(p+1)^n$ points

$$\mathbf{x}_{i_1, i_2, \dots, i_n} = \left(\frac{i_1}{p}, \frac{i_2}{p}, \dots, \frac{i_n}{p} \right)^T$$

where $(i_1, i_2, \dots, i_n) \in \{0, 1, \dots, p\}^n$.

2. Among all points $\mathbf{x}_{i_1, i_2, \dots, i_n}$, find a point $\bar{\mathbf{x}}$ which has the minimal value for the objective function.
3. Return the pair $(\bar{\mathbf{x}}, f(\bar{\mathbf{x}}))$ as the result.

Theorem 5.4 Let $f^* := f(\mathbf{x}^*)$ be the global optimal value for (4). Then the uniform grid method yields

$$f(\bar{\mathbf{x}}) - f(\mathbf{x}^*) \leq \frac{L}{2p}.$$

Proof:

Let \mathbf{x}^* be a global optimal solution. Then there are coordinates (i_1, i_2, \dots, i_n) such that $\mathbf{x} := \mathbf{x}_{i_1, i_2, \dots, i_n} \leq \mathbf{x}^* \leq \mathbf{x}_{i_1+1, i_2+1, \dots, i_n+1} =: \mathbf{y}$. Observe that $[\mathbf{y}]_i - [\mathbf{x}]_i = 1/p$ for $i = 1, 2, \dots, n$ and $[\mathbf{x}^*]_i \in [[\mathbf{x}]_i, [\mathbf{y}]_i]$ ($i = 1, 2, \dots, n$).

Consider $\tilde{\mathbf{x}} = (\mathbf{x} + \mathbf{y})/2$ and form a new point $\tilde{\mathbf{x}}$ as:

$$[\tilde{\mathbf{x}}]_i := \begin{cases} [\mathbf{y}]_i, & \text{if } [\mathbf{x}^*]_i \geq [\tilde{\mathbf{x}}]_i \\ [\mathbf{x}]_i, & \text{otherwise.} \end{cases}$$

It is clear that $|[\tilde{\mathbf{x}}]_i - [\mathbf{x}^*]_i| \leq 1/(2p)$ for $i = 1, 2, \dots, n$. Then $\|\tilde{\mathbf{x}} - \mathbf{x}^*\|_\infty = \max_{1 \leq i \leq n} |[\tilde{\mathbf{x}}]_i - [\mathbf{x}^*]_i| \leq 1/(2p)$. Since $\tilde{\mathbf{x}}$ belongs to the grid,

$$f(\bar{\mathbf{x}}) - f(\mathbf{x}^*) \leq f(\tilde{\mathbf{x}}) - f(\mathbf{x}^*) \leq L\|\tilde{\mathbf{x}} - \mathbf{x}^*\|_\infty \leq L/(2p).$$

■

Let us define our goal

$$\text{Find } \mathbf{x} \in B_n \text{ such that } f(\mathbf{x}) - f(\mathbf{x}^*) < \varepsilon.$$

Corollary 5.5 The number of iterations necessary for the problem (4) for the uniform grid method is at most

$$\left(\left\lfloor \frac{L}{2\varepsilon} \right\rfloor + 2 \right)^n.$$

Proof:

Take $p = \lfloor L/(2\varepsilon) \rfloor + 1$. Then, $p > L/(2\varepsilon)$ and from the previous theorem, $f(\bar{\mathbf{x}}) - f(\mathbf{x}^*) \leq L/(2p) < \varepsilon$. Observe that we constructed $(p+1)^n$ points. ■

Consider the class of problems \mathcal{P} defined as follows:

Model:	$\min_{\mathbf{x} \in B_n} f(\mathbf{x}),$
Oracle:	$f(\mathbf{x})$ is ℓ_∞ -Lipschitz continuous on B_n . Only function values are available
Approximate solution:	Find $\bar{\mathbf{x}} \in B_n$ such that $f(\bar{\mathbf{x}}) - f(\mathbf{x}^*) < \varepsilon$

Theorem 5.6 For $\varepsilon < \frac{L}{2}$, the number of iterations necessary for the class of problems \mathcal{P} using any method which uses only function evaluations is always at least $(\lfloor \frac{L}{2\varepsilon} \rfloor)^n$.

Proof:

Let $p = \lfloor \frac{L}{2\varepsilon} \rfloor$ (which is ≥ 1 from the hypothesis).

Suppose that there is a method which requires $N < p^n$ calls of the oracle to solve the problem in \mathcal{P} .

Then, there is a point $\hat{\mathbf{x}} \in B_n = \{\mathbf{x} \in \mathbb{R}^n \mid 0 \leq [\mathbf{x}]_i \leq 1, i = 1, 2, \dots, n\}$ where there is no test points in the interior of $B := \{\mathbf{x} \mid \hat{\mathbf{x}} \leq \mathbf{x} \leq \hat{\mathbf{x}} + \mathbf{e}/p\}$ where $\mathbf{e} = (1, 1, \dots, 1)^T \in \mathbb{R}^n$.

Let $\mathbf{x}^* := \hat{\mathbf{x}} + \mathbf{e}/(2p)$ and consider the function $\bar{f}(\mathbf{x}) := \min\{0, L\|\mathbf{x} - \mathbf{x}^*\|_\infty - \varepsilon\}$. Clearly, \bar{f} is ℓ_∞ -Lipschitz continuous with constant L and its global minimum is $-\varepsilon$. Moreover, $\bar{f}(\mathbf{x})$ is non-zero valued only inside the box $B' := \{\mathbf{x} \mid \|\mathbf{x} - \mathbf{x}^*\|_\infty \leq \varepsilon/L\}$.

Since $2p \leq L/\varepsilon$, $B' \subseteq \{\mathbf{x} \mid \|\mathbf{x} - \mathbf{x}^*\|_\infty \leq 1/(2p)\} \subseteq B$.

Therefore, $\bar{f}(\mathbf{x})$ is equal to zero to all test points of our method and the accuracy of the method is ε .

If the number of calls of the oracle is less than p^n , the accuracy can not be better than ε . ■

Theorem 5.6 supports the claim that the *general optimization problem are unsolvable*.

Example 5.7 Consider a problem defined by the following parameters. $L = 2$, $n = 10$, and $\varepsilon = 0.01$ (1%).

lower bound $(L/(2\varepsilon))^n$: 10^{20} calls of the oracle
computational complexity of the oracle	: at least n arithmetic operations
total complexity	: 10^{21} arithmetic operations
CPU	: 1GHz or 10^9 arithmetic operations per second
total time	: 10^{12} seconds
one year	: $\leq 3.2 \times 10^7$ seconds
we need	: ≥ 10000 years

- If we change n by $n + 1$, the # of calls of the oracle is multiplied by 100.
- If we multiply ε by 2, the arithmetic complexity is reduced by 1000.

We know from Corollary 5.5 that the number of iterations of the uniform grid method is at least $(\lfloor L/(2\varepsilon) \rfloor + 2)^n$. Theorem 5.6 showed that any method which uses only function evaluations requires at least $(\lfloor L/(2\varepsilon) \rfloor)^n$ calls to have a better performance than ε . If for instance we take $\varepsilon = \mathcal{O}(L/n)$, these two bounds coincide up to a constant factor. In this sense, the uniform grid method is an *optimal method for the class of problems \mathcal{P}* .

5.3 Steepest Descent Method

Consider $f : \mathbb{R}^n \rightarrow \mathbb{R}$ a differentiable function in its domain.

Steepest Descent Method	
Choose:	$\mathbf{x}_0 \in \mathbb{R}^n$
Iterate:	$\mathbf{x}_{k+1} = \mathbf{x}_k - h_k f'(\mathbf{x}_k), k = 0, 1, \dots$

We consider four strategies for the step-size h_k :

1. Constant Step

The sequence $\{h_k\}_{k=0}^\infty$ is chosen in *advance*. For example

$$h_k := h > 0,$$

$$h_k := \frac{h}{\sqrt{k+1}}.$$

This is the simplest strategy.

2. Exact Line Search (Cauchy Step-Size)

The sequence $\{h_k\}_{k=0}^{\infty}$ is chosen such that

$$h_k := \arg \min_{h \geq 0} f(\mathbf{x}_k - hf'(\mathbf{x}_k)).$$

This choice is only theoretical since even for the one dimensional case, it is very difficult and expensive.

3. Goldstein-Armijo Rule

Find a sequence $\{h_k\}_{k=0}^{\infty}$ such that

$$\begin{aligned} \alpha \langle f'(\mathbf{x}_k), \mathbf{x}_k - \mathbf{x}_{k+1} \rangle &\leq f(\mathbf{x}_k) - f(\mathbf{x}_{k+1}), \\ \beta \langle f'(\mathbf{x}_k), \mathbf{x}_k - \mathbf{x}_{k+1} \rangle &\geq f(\mathbf{x}_k) - f(\mathbf{x}_{k+1}), \end{aligned}$$

where $0 < \alpha < \beta < 1$ are fixed parameters.

Since $f(\mathbf{x}_{k+1}) = f(\mathbf{x}_k - h_k f'(\mathbf{x}_k))$,

$$f(\mathbf{x}_k) - \beta h_k \|f'(\mathbf{x}_k)\|_2^2 \leq f(\mathbf{x}_{k+1}) \leq f(\mathbf{x}_k) - \alpha h_k \|f'(\mathbf{x}_k)\|_2^2.$$

The acceptable steps exist unless $f(\mathbf{x}_{k+1}) = f(\mathbf{x}_k - hf'(\mathbf{x}_k))$ is not bounded from below.

4. Barzilai-Borwein Step-Size¹

Let us define $\mathbf{s}_{k-1} := \mathbf{x}_k - \mathbf{x}_{k-1}$ and $\mathbf{y}_{k-1} := f'(\mathbf{x}_k) - f'(\mathbf{x}_{k-1})$. Then, we can define the Barzilai-Borwein (BB) step sizes $\{h_k^1\}_{k=1}^{\infty}$ and $\{h_k^2\}_{k=1}^{\infty}$:

$$\begin{aligned} h_k^1 &:= \frac{\|\mathbf{s}_{k-1}\|_2^2}{\langle \mathbf{s}_{k-1}, \mathbf{y}_{k-1} \rangle}, \\ h_k^2 &:= \frac{\langle \mathbf{s}_{k-1}, \mathbf{y}_{k-1} \rangle}{\|\mathbf{y}_{k-1}\|_2^2}. \end{aligned}$$

The first step-size is the one which minimizes the following secant condition $\|\frac{1}{h}\mathbf{s}_{k-1} - \mathbf{y}_{k-1}\|_2^2$ while the second one minimizes $\|\mathbf{s}_{k-1} - h\mathbf{y}_{k-1}\|_2^2$.

Now, consider the problem

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

where $f \in \mathcal{C}_L^{1,1}(\mathbb{R}^n)$, and $f(\mathbf{x})$ is bounded from below.

¹J. Barzilai and J. M. Borwein, "Two-point step size gradient methods," *IMA Journal of Numerical Analysis*, **8** (1988), pp. 141–148.