

Foundation of Computing and Mathematical Sciences — Optimization —

Tokyo Institute of Technology
Dept. Mathematical and Computing Sciences
MITUHIRO FUKUDA

Fall/Winter Semester of 2012

“In our opinion, the main fact, which should be known to any person dealing with optimization models, is that in general *optimization problems are unsolvable*.”
— Yuri Nesterov

Bibliography

- [DASPREMONT2008] A. d’Aspremont, “Smooth optimization with approximate gradient”, *SIAM Journal on Optimization* **19** (2008), pp. 1171–1183.
- [GK2008] C. C. Gonzaga and E. W. Karas, “Fine tuning Nesterov’s steepest descent algorithm for differentiable convex programming”, *Mathematical Programming*, to appear.
- [LLM2006] G. Lan, Z. Lu, and R. D. C. Monteiro, “Primal-dual first-order methods with $\mathcal{O}(1/\varepsilon)$ iteration-complexity for cone programming”, *Mathematical Programming*, **126** (2011), pp.1–29.
- [NESTEROV2004] Yu. Nesterov, *Introductory Lecture on Convex Optimization: A Basic Course*, (Kluwer Academic Publishers, Boston, 2004).
- [NESTEROV2005] Yu. Nesterov, “Smooth minimization of non-smooth functions”, *Mathematical Programming* **103** (2005), pp. 127–152.
- [NESTEROV2005-2] Yu. Nesterov, “Excessive gap technique in nonsmooth convex minimization”, *SIAM Journal on Optimization* **16** (2005), pp. 669–700.
- [NESTEROV2007] Yu. Nesterov, “Smoothing technique and its applications in semidefinite optimization”, *Mathematical Programming* **110** (2007), pp. 245–259.
- [NESTEROV2009] Yu. Nesterov, “Primal-dual subgradient methods for convex problems”, *Mathematical Programming* **120** (2009), pp. 221–259.
- [NOCEDAL2006] J. Nocedal and S. J. Wright, *Numerical Optimization*, 2nd edition, (Springer, New York, 2006).
- [TSENG2010] P. Tseng, “Approximation accuracy, gradient methods, and error bound for structured convex optimization”, *Mathematical Programming* **12** (2010), pp. 263–295.
- [YUAN2010] Y.-X. Yuan, “A short note on the Q -linear convergence of the steepest descent method”, *Mathematical Programming* **123** (2010), pp. 339–343.

Chapter 1

Nonlinear Optimization

1.1 General minimization problem and terminologies

Definition 1.1.1 We define the *general minimization problem* as follows

$$\begin{cases} \text{minimize} & f(\mathbf{x}) \\ \text{subject to} & f_j(\mathbf{x}) \& 0, \quad j = 1, 2, \dots, m \\ & \mathbf{x} \in S, \end{cases} \quad (1.1)$$

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 1.1.2 The *feasible set* Q of (1.1) 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 $Q \equiv \mathbb{R}^n$, (1.1) is a *unconstrained optimization problem*.
- If $Q \subsetneq \mathbb{R}^n$, (1.1) is a *constrained optimization problem*.
- If all functionals $f(\mathbf{x})$, $f_j(\mathbf{x})$ are differentiable, (1.1) is a *smooth optimization problem*.
- If one of functionals $f(\mathbf{x})$, $f_j(\mathbf{x})$ is non-differentiable, (1.1) is a *non-smooth optimization problem*.
- If all constraints are linear $f_j(\mathbf{x}) = \sum_{i=1}^n [\mathbf{a}]_{ji}[\mathbf{x}]_i + [\mathbf{b}]_j$ ($j = 1, 2, \dots, m$), (1.1) is a *linear constrained optimization problem*.
 - In addition, if $f(\mathbf{x})$ is linear, (1.1) is a *linear programming problem*.
 - In addition, if $f(\mathbf{x})$ is quadratic, (1.1) is a *quadratic programming problem*.
- If $f(\mathbf{x})$, $f_j(\mathbf{x})$ ($j = 1, 2, \dots, m$) are quadratic, (1.1) is a *quadratically constrained quadratic programming problem*.

Definition 1.1.3 \mathbf{x}^* is called a *global optimal solution* of (1.1) 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 (1.1) if there exists an open ball $B(\varepsilon) = \{\mathbf{x} \in \mathbb{R}^n \mid \|\mathbf{x} - \mathbf{x}^*\| < \varepsilon\} \subseteq Q$ such that $f(\mathbf{x}^*) \leq f(\mathbf{x})$, $\forall \mathbf{x} \in B(\varepsilon)$. Moreover, $f(\mathbf{x}^*)$ is called a *local optimal value*.

| General Iterative Scheme | |
|--------------------------|--|
| Input: | A starting point \mathbf{x}_0 and an accuracy $\varepsilon > 0$. |
| Initialization: | Set the <i>iteration counter</i> $k := 0$, and the <i>information set</i> $I_{-1} := \emptyset$. |
| MAIN LOOP | |
| 1. | Call oracle \mathcal{O} at \mathbf{x}_k . |
| 2. | Update the information set: $I_k := I_{k-1} \cup (\mathbf{x}_k, \mathcal{O}(\mathbf{x}_k))$. |
| 3. | Apply the rules of the <i>method</i> \mathcal{M} to I_k and compute \mathbf{x}_{k+1} . |
| 4. | Check <i>stopping criterion</i> \mathcal{T}_ε . If Yes , output $\bar{\mathbf{x}}$. Otherwise set $k := k + 1$ and go to Step 1. |

Definition 1.1.4 The *analytical complexity* of a method is the number of calls of an oracle which is required to solve a problem \mathcal{P} up to the given accuracy ε .

Definition 1.1.5 The *arithmetical complexity* of a method is the total number of arithmetic operations (including the work of the oracle and the method) which is required to solve a problem \mathcal{P} up to the given accuracy ε .

Assumption 1.1.6 (Local black box)

1. The only information available for the numerical scheme is the answer of the oracle.
2. The oracle is local, that is, a small variation of the problem far enough from the test point \mathbf{x} does not change the answer at \mathbf{x} .

Definition 1.1.7

1. The *zero-order oracle* returns the value $f(\mathbf{x})$.
2. The *first-order oracle* returns the value $f(\mathbf{x})$, and the gradient $f'(\mathbf{x})$.
3. The *second-order oracle* returns the value $f(\mathbf{x})$, $f'(\mathbf{x})$ and the Hessian $f''(\mathbf{x})$.

1.2 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.

$$\begin{cases} \text{minimize} & f(\mathbf{x}) \\ \text{subject to} & \mathbf{x} \in B_n = \{\mathbf{x} \in \mathbb{R}^n \mid 0 \leq [\mathbf{x}]_i \leq 1, i = 1, 2, \dots, n\}. \end{cases} \quad (1.2)$$

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 (1.2), 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 1.2.1 Let f^* be the global optimal value for (1.2). Then the uniform grid method yields

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

Proof: Let \mathbf{x}^* be a global optimal solution. Then there are coordinates (i_1, i_2, \dots, i_n) such that $\mathbf{x} \equiv \mathbf{x}_{i_1, i_2, \dots, i_n} \leq \mathbf{x}^* \leq \mathbf{x}_{i_1+1, i_2+1, \dots, i_n+1} \equiv \mathbf{y}$. 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 $\hat{\mathbf{x}} = (\mathbf{x} + \mathbf{y})/2$ and form a new point $\tilde{\mathbf{x}}$ as:

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

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

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

Corollary 1.2.2 The analytical complexity of the problem (1.2) 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{C} defined as follows:

| | |
|------------------------------|--|
| Model: | $\min_{\mathbf{x} \in B_n} f(\mathbf{x}),$ |
| Oracle: | $f(\mathbf{x})$ is ℓ_∞ -Lipschitz continuous on B_n . |
| Approximate solution: | zero-order local black box (only function values) Find $\bar{\mathbf{x}} \in B_n$ such that $f(\bar{\mathbf{x}}) - f^* < \varepsilon$ |

Theorem 1.2.3 For $\varepsilon < \frac{L}{2}$, the analytical complexity of class of problems \mathcal{C} using zero-order methods is 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 \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 $B \equiv \{\mathbf{x} \mid \hat{\mathbf{x}} \leq \mathbf{x} \leq \hat{\mathbf{x}} + \mathbf{e}/p\}$ 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 $2p \leq L/\varepsilon$, $B' \subseteq B = \{\mathbf{x} \mid \|\mathbf{x} - \mathbf{x}^*\|_\infty \leq 1/(2p)\}$.

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 1.2.3 supports our initial claim that the *general optimization problem are unsolvable*.

Example 1.2.4 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 |
| 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 analytical complexity estimate is multiplied by 100.
- If we multiply ε by 2, the arithmetic complexity is reduced by 1000.

We know from Corollary 1.2.2 that the analytical complexity for the uniform grid method is $(\lfloor L/(2\varepsilon) \rfloor + 2)^n$. Theorem 1.2.3 showed that any method with zero-order oracle 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 \mathcal{C}* .

1.3 Optimality conditions for unconstrained minimization problems

Let $f(\mathbf{x})$ be differentiable at $\bar{\mathbf{x}}$. Then for $\mathbf{y} \in \mathbb{R}^n$, we have

$$f(\mathbf{y}) = f(\bar{\mathbf{x}}) + \langle f'(\bar{\mathbf{x}}), \mathbf{y} - \bar{\mathbf{x}} \rangle + o(\|\mathbf{y} - \bar{\mathbf{x}}\|),$$

where $o(r)$ is some function of $r > 0$ such that

$$\lim_{r \rightarrow 0} \frac{1}{r} o(r) = 0, \quad o(0) = 0.$$

Let \mathbf{s} be a direction in \mathbb{R}^n such that $\|\mathbf{s}\| = 1$. Consider the local decrease of $f(\mathbf{x})$ along \mathbf{s} :

$$\Delta(\mathbf{s}) = \lim_{\alpha \rightarrow 0} \frac{1}{\alpha} [f(\bar{\mathbf{x}} + \alpha \mathbf{s}) - f(\bar{\mathbf{x}})].$$

Since $f(\bar{\mathbf{x}} + \alpha \mathbf{s}) - f(\bar{\mathbf{x}}) = \alpha \langle f'(\bar{\mathbf{x}}), \mathbf{s} \rangle + o(\|\alpha \mathbf{s}\|)$, we have $\Delta(\mathbf{s}) = \langle f'(\bar{\mathbf{x}}), \mathbf{s} \rangle$.

Using the Cauchy-Schwartz inequality $-\|\mathbf{x}\|\|\mathbf{y}\| \leq \langle \mathbf{x}, \mathbf{y} \rangle \leq \|\mathbf{x}\|\|\mathbf{y}\|$,

$$\Delta(\mathbf{s}) = \langle f'(\bar{\mathbf{x}}), \mathbf{s} \rangle \geq -\|f'(\bar{\mathbf{x}})\|.$$

Choosing the direction $\bar{\mathbf{s}} = -f'(\bar{\mathbf{x}})/\|f'(\bar{\mathbf{x}})\|$,

$$\Delta(\bar{\mathbf{s}}) = -\left\langle f'(\bar{\mathbf{x}}), \frac{f'(\bar{\mathbf{x}})}{\|f'(\bar{\mathbf{x}})\|} \right\rangle = -\|f'(\bar{\mathbf{x}})\|.$$

Thus, the direction $-f'(\bar{\mathbf{x}})$ is the direction of the *fastest local decrease of $f(\mathbf{x})$ at point $\bar{\mathbf{x}}$* .

Theorem 1.3.1 (First-order necessary optimality condition) Let \mathbf{x}^* be a local minimum of the differentiable function $f(\mathbf{x})$. Then

$$f'(\mathbf{x}^*) = \mathbf{0}.$$

Proof: Let \mathbf{x}^* be the local minimum of $f(\mathbf{x})$. Then, there is $r > 0$ such that for all \mathbf{y} with $\|\mathbf{y} - \mathbf{x}^*\| \leq r$, $f(\mathbf{y}) \geq f(\mathbf{x}^*)$.

Since f is differentiable,

$$f(\mathbf{y}) = f(\mathbf{x}^*) + \langle f'(\mathbf{x}^*), \mathbf{y} - \mathbf{x}^* \rangle + o(\|\mathbf{y} - \mathbf{x}^*\|) \geq f(\mathbf{x}^*).$$

Dividing by $\|\mathbf{y} - \mathbf{x}^*\|$, and taking the limit $\mathbf{y} \rightarrow \mathbf{x}^*$,

$$\langle f'(\mathbf{x}^*), \mathbf{s} \rangle \geq 0, \quad \forall \mathbf{s}, \quad \|\mathbf{s}\| = 1.$$

Consider the opposite direction $-\mathbf{s}$, and then we conclude that

$$\langle f'(\mathbf{x}^*), \mathbf{s} \rangle = 0, \quad \forall \mathbf{s}, \quad \|\mathbf{s}\| = 1.$$

Choosing $\mathbf{s} = \mathbf{e}_i$ ($i = 1, 2, \dots, n$), we conclude that $f'(\mathbf{x}^*) = 0$. ■

Corollary 1.3.2 Let \mathbf{x}^* be a local minimum of a differentiable function $f(\mathbf{x})$ subject to linear equality constraints

$$\mathbf{x} \in \mathcal{L} \equiv \{\mathbf{x} \in \mathbb{R}^n \mid \mathbf{A}\mathbf{x} = \mathbf{b}\} \neq \emptyset,$$

where $\mathbf{A} \in \mathbb{R}^{m \times n}$, $\mathbf{b} \in \mathbb{R}^m$, $m < n$.

Then, there exists a vector of multipliers $\boldsymbol{\lambda}^*$ such that

$$f'(\mathbf{x}^*) = \mathbf{A}^T \boldsymbol{\lambda}^*.$$

Proof: Consider the vectors \mathbf{u}_i ($i = 1, 2, \dots, k$) with $k \geq n - m$ which form an orthonormal basis of the null space of \mathbf{A} . Then, $\mathbf{x} \in \mathcal{L}$ can be represented as

$$\mathbf{x} = \mathbf{x}(\mathbf{t}) \equiv \mathbf{x}^* + \sum_{i=1}^k t_i \mathbf{u}_i, \quad \mathbf{t} \in \mathbb{R}^k.$$

Moreover, the point $\mathbf{t} = \mathbf{0}$ is the local minimal solution of the function $\phi(\mathbf{t}) = f(\mathbf{x}(\mathbf{t}))$.

From Theorem 1.3.1, $\phi'(\mathbf{0}) = \mathbf{0}$. That is,

$$\frac{d\phi}{dt_i}(\mathbf{0}) = \langle f'(\mathbf{x}^*), \mathbf{u}_i \rangle = 0, \quad i = 1, 2, \dots, k.$$

Now there is \mathbf{t}^* and $\boldsymbol{\lambda}^*$ such that

$$f'(\mathbf{x}^*) = \sum_{i=1}^k t_i^* \mathbf{u}_i + \mathbf{A}^T \boldsymbol{\lambda}^*.$$

For each $i = 1, 2, \dots, k$,

$$\langle f'(\mathbf{x}^*), \mathbf{u}_i \rangle = t_i^* = 0.$$

Therefore, we have the result. ■

If $f(\mathbf{x})$ is twice differentiable at $\bar{\mathbf{x}}$, then for $\mathbf{y} \in \mathbb{R}^n$, we have

$$f'(\mathbf{y}) = f'(\bar{\mathbf{x}}) + f''(\bar{\mathbf{x}})(\mathbf{y} - \bar{\mathbf{x}}) + \mathbf{o}(\|\mathbf{y} - \bar{\mathbf{x}}\|),$$

where $\mathbf{o}(r)$ is such that $\lim_{r \rightarrow 0} \|\mathbf{o}(r)\|/r = 0$ and $\mathbf{o}(0) = \mathbf{0}$.

Theorem 1.3.3 (Second-order necessary optimality condition) Let \mathbf{x}^* be a local minimum of a twice continuously differentiable function $f(\mathbf{x})$. Then

$$f'(\mathbf{x}^*) = 0, \quad f''(\mathbf{x}^*) \succeq \mathbf{O}.$$