

We know from Corollary 4.5 that the number of iterations of the uniform grid method is at least $(\lfloor L/(2\varepsilon) \rfloor + 2)^n$. Theorem 4.6 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}* .

4.3 Optimality Conditions for Smooth Optimization Problems

Let $f : \mathbb{R}^n \rightarrow \mathbb{R}$ be a differentiable function on \mathbb{R}^n , $\bar{\mathbf{x}} \in \mathbb{R}^n$, and \mathbf{s} be a direction in \mathbb{R}^n such that $\|\mathbf{s}\|_2 = 1$. Consider the local decrease (or increase) of $f(\mathbf{x})$ along \mathbf{s} :

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

Since $f(\bar{\mathbf{x}} + \alpha \mathbf{s}) - f(\bar{\mathbf{x}}) = \alpha \langle \nabla f(\bar{\mathbf{x}}), \mathbf{s} \rangle + o(\|\alpha \mathbf{s}\|_2)$, we have $f'(\bar{\mathbf{x}}; \mathbf{s}) = \langle \nabla f(\bar{\mathbf{x}}), \mathbf{s} \rangle$.

Using the Cauchy-Schwarz inequality $-\|\mathbf{x}\|_2 \|\mathbf{y}\|_2 \leq \langle \mathbf{x}, \mathbf{y} \rangle \leq \|\mathbf{x}\|_2 \|\mathbf{y}\|_2$,

$$f'(\bar{\mathbf{x}}; \mathbf{s}) = \langle \nabla f(\bar{\mathbf{x}}), \mathbf{s} \rangle \geq -\|\nabla f(\bar{\mathbf{x}})\|_2.$$

Choosing in particular the direction $\bar{\mathbf{s}} = -\nabla f(\bar{\mathbf{x}})/\|\nabla f(\bar{\mathbf{x}})\|_2$,

$$f'(\bar{\mathbf{x}}; \bar{\mathbf{s}}) = -\left\langle \nabla f(\bar{\mathbf{x}}), \frac{\nabla f(\bar{\mathbf{x}})}{\|\nabla f(\bar{\mathbf{x}})\|_2} \right\rangle = -\|\nabla f(\bar{\mathbf{x}})\|_2.$$

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

Theorem 4.8 (First-order necessary optimality condition) Let \mathbf{x}^* be a local minimum of the differentiable function $f(\mathbf{x})$. Then

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

Proof:

Let \mathbf{x}^* be the local minimum of $f(\mathbf{x})$. Then, there is $r > 0$ such that for all \mathbf{y} with $\|\mathbf{y} - \mathbf{x}^*\|_2 \leq r$, $f(\mathbf{y}) \geq f(\mathbf{x}^*)$.

Since f is differentiable on \mathbb{R}^n ,

$$f(\mathbf{y}) = f(\mathbf{x}^*) + \langle \nabla f(\mathbf{x}^*), \mathbf{y} - \mathbf{x}^* \rangle + o(\|\mathbf{y} - \mathbf{x}^*\|_2) \geq f(\mathbf{x}^*).$$

Dividing by $\|\mathbf{y} - \mathbf{x}^*\|_2$, and taking the limit $\mathbf{y} \rightarrow \mathbf{x}^*$,

$$\langle \nabla f(\mathbf{x}^*), \mathbf{s} \rangle \geq 0, \quad \forall \mathbf{s} \in \mathbb{R}^n, \quad \|\mathbf{s}\|_2 = 1.$$

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

$$\langle \nabla f(\mathbf{x}^*), \mathbf{s} \rangle = 0, \quad \forall \mathbf{s} \in \mathbb{R}^n, \quad \|\mathbf{s}\|_2 = 1.$$

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

Remark 4.9 For the first-order sufficient optimality condition, we need convexity for the function $f(\mathbf{x})$.

Corollary 4.10 Let \mathbf{x}^* be a local minimum of a differentiable function $f(\mathbf{x})$ subject to linear equality constraints

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

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

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

$$\nabla f(\mathbf{x}^*) = \mathbf{A}^T \boldsymbol{\lambda}^*.$$

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

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

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 4.8, $\phi'(\mathbf{0}) = \mathbf{0}$. That is,

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

Now there is $\mathbf{t}^* \in \mathbb{R}^k$ and $\boldsymbol{\lambda}^* \in \mathbb{R}^m$ such that

$$\nabla f(\mathbf{x}^*) = \sum_{i=1}^k t_i^* \mathbf{u}_i + \mathbf{A}^T \boldsymbol{\lambda}^*.$$

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

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

Therefore, we have the result. ■

The following type of result is called *theorems of the alternative*, and are closely related to duality theory in optimization.

Corollary 4.11 Given $\mathbf{A} \in \mathbb{R}^{m \times n}$, $\mathbf{b} \in \mathbb{R}^m$, $\mathbf{c} \in \mathbb{R}^n$, $\eta \in \mathbb{R}$, either

$$\left\{ \begin{array}{l} \langle \mathbf{c}, \mathbf{x} \rangle < \eta \\ \mathbf{Ax} = \mathbf{b} \end{array} \right. \text{ has a solution } \mathbf{x} \in \mathbb{R}^n, \quad (3)$$

or

$$\left(\begin{array}{l} \left\{ \begin{array}{l} \langle \mathbf{b}, \boldsymbol{\lambda} \rangle > 0 \\ \mathbf{A}^T \boldsymbol{\lambda} = \mathbf{0} \end{array} \right. \\ \text{or} \\ \left\{ \begin{array}{l} \langle \mathbf{b}, \boldsymbol{\lambda} \rangle \geq \eta \\ \mathbf{A}^T \boldsymbol{\lambda} = \mathbf{c} \end{array} \right. \end{array} \right) \text{ has a solution } \boldsymbol{\lambda} \in \mathbb{R}^m, \quad (4)$$

but never both

Proof:

Let us first show that if $\exists \mathbf{x} \in \mathbb{R}^n$ satisfying (3), $\nexists \boldsymbol{\lambda} \in \mathbb{R}^m$ satisfying (4). Let us assume by contradiction that $\exists \boldsymbol{\lambda}$. Then $\langle \boldsymbol{\lambda}, \mathbf{Ax} \rangle = \langle \boldsymbol{\lambda}, \mathbf{b} \rangle$ and in the homogeneous case it gives $0 = \langle \boldsymbol{\lambda}, \mathbf{b} \rangle > 0$ and in the non-homogeneous case it gives $\eta > \langle \mathbf{c}, \mathbf{x} \rangle = \langle \boldsymbol{\lambda}, \mathbf{b} \rangle \geq \eta$. Both of cases are impossible.

Now, let us assume that $\nexists \mathbf{x} \in \mathbb{R}^n$ satisfying (3). If additionally $\nexists \mathbf{x} \in \mathbb{R}^n$ such that $\mathbf{Ax} = \mathbf{b}$, it means that the columns of the matrix \mathbf{A} do not span the vector \mathbf{b} . Therefore, there is $\mathbf{0} \neq \boldsymbol{\lambda} \in \mathbb{R}^m$ which is orthogonal to all of these columns and $\langle \mathbf{b}, \boldsymbol{\lambda} \rangle \neq 0$. Selecting the correct sign, we constructed a $\boldsymbol{\lambda}$ which satisfies the homogeneous system of (4). Now, if for all \mathbf{x} such that $\mathbf{Ax} = \mathbf{b}$ we have $\langle \mathbf{c}, \mathbf{x} \rangle \geq \eta$, it means that the minimization of the function $f(\mathbf{x}) = \langle \mathbf{c}, \mathbf{x} \rangle$ subject to $\mathbf{Ax} = \mathbf{b}$ has an optimal solution \mathbf{x}^* with $f(\mathbf{x}^*) \geq \eta$ (since $\exists \mathbf{x} \in \mathbb{R}^n$ such that $\mathbf{Ax} = \mathbf{b}$, we can always assume that $m \leq n$ eliminating redundant linear constraints from the system. If $n = m$ and \mathbf{A} is nonsingular, take $\boldsymbol{\lambda} = \mathbf{A}^{-T} \mathbf{c}$. Otherwise, we can eliminate again redundant linear constraint to have $n > m$). From Corollary 4.10, $\exists \boldsymbol{\lambda} \in \mathbb{R}^m$ such that $\mathbf{A}^T \boldsymbol{\lambda} = \mathbf{c}$, and $\langle \mathbf{b}, \boldsymbol{\lambda} \rangle = \langle \mathbf{x}^*, \mathbf{A}^T \boldsymbol{\lambda} \rangle = \langle \mathbf{x}^*, \mathbf{c} \rangle \geq \eta$. ■

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

$$\nabla f(\mathbf{y}) = \nabla f(\bar{\mathbf{x}}) + \nabla^2 f(\bar{\mathbf{x}})(\mathbf{y} - \bar{\mathbf{x}}) + o(\|\mathbf{y} - \bar{\mathbf{x}}\|_2),$$

where $o(r)$ is such that $\lim_{r \rightarrow 0} \|o(r)\|_2/r = 0$ and $o(0) = 0$.

Theorem 4.12 (Second-order necessary optimality condition) Let \mathbf{x}^* be a local minimum of a twice continuously differentiable function $f(\mathbf{x})$. Then

$$\nabla f(\mathbf{x}^*) = 0, \quad \nabla^2 f(\mathbf{x}^*) \succeq \mathbf{O}.$$

Proof:

Since \mathbf{x}^* is a local minimum of $f(\mathbf{x})$, $\exists r > 0$ such that for all $\mathbf{y} \in \mathbb{R}^n$ which satisfy $\|\mathbf{y} - \mathbf{x}^*\|_2 \leq r$, $f(\mathbf{y}) \geq f(\mathbf{x}^*)$.

From Theorem 4.8, $\nabla f(\mathbf{x}^*) = 0$. Then

$$f(\mathbf{y}) = f(\mathbf{x}^*) + \frac{1}{2} \langle \nabla^2 f(\mathbf{x}^*)(\mathbf{y} - \mathbf{x}^*), \mathbf{y} - \mathbf{x}^* \rangle + o(\|\mathbf{y} - \mathbf{x}^*\|_2^2) \geq f(\mathbf{x}^*).$$

And $\langle \nabla^2 f(\mathbf{x}^*)\mathbf{s}, \mathbf{s} \rangle \geq 0$, $\forall \mathbf{s} \in \mathbb{R}^n$ with $\|\mathbf{s}\|_2 = 1$. ■

Theorem 4.13 (Second-order sufficient optimality condition) Let the function $f(\mathbf{x})$ be twice continuously differentiable on \mathbb{R}^n , and let \mathbf{x}^* satisfy the following conditions:

$$\nabla f(\mathbf{x}^*) = 0, \quad \nabla^2 f(\mathbf{x}^*) \succ \mathbf{O}.$$

Then, \mathbf{x}^* is a strict local minimum of $f(\mathbf{x})$.

Proof:

In a small neighborhood of \mathbf{x}^* , function $f(\mathbf{x}^*)$ can be represented as:

$$f(\mathbf{y}) = f(\mathbf{x}^*) + \frac{1}{2} \langle \nabla^2 f(\mathbf{x}^*)(\mathbf{y} - \mathbf{x}^*), \mathbf{y} - \mathbf{x}^* \rangle + o(\|\mathbf{y} - \mathbf{x}^*\|_2^2).$$

Since $o(r)/r \rightarrow 0$, there is a $\bar{r} > 0$ such that for all $r \in [0, \bar{r}]$,

$$|o(r)| \leq \frac{r}{4} \lambda_1(\nabla^2 f(\mathbf{x}^*)),$$

where $\lambda_1(\nabla^2 f(\mathbf{x}^*))$ is the smallest eigenvalue of the symmetric matrix $\nabla^2 f(\mathbf{x}^*)$ which is positive. Then

$$f(\mathbf{y}) \geq f(\mathbf{x}^*) + \frac{1}{2} \lambda_1(\nabla^2 f(\mathbf{x}^*)) \|\mathbf{y} - \mathbf{x}^*\|_2^2 + o(\|\mathbf{y} - \mathbf{x}^*\|_2^2).$$

W.L.O.G, considering that $\bar{r} < 1$, $|o(r^2)| \leq r^2 \lambda_1(\nabla^2 f(\mathbf{x}^*))/4$ for $r \in [0, \bar{r}]$, finally we arrived at

$$f(\mathbf{y}) \geq f(\mathbf{x}^*) + \frac{1}{4} \lambda_1(\nabla^2 f(\mathbf{x}^*)) \|\mathbf{y} - \mathbf{x}^*\|_2^2 > f(\mathbf{x}^*).$$
 ■

4.4 Algorithms for Minimizing Smooth Functions

4.4.1 Steepest Descent Method

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

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

We consider four strategies for the step-size h_k :