# Chapter 9

# Optimization over a Convex Set

## 9.1 ▪ Stationarity

Throughout this chapter we will consider the constrained optimization problem (P) given by

$$\text{(P)} \quad \begin{array}{ll} \min & f(\mathbf{x}) \\ \text{s.t.} & \mathbf{x} \in C, \end{array}$$

where $f$ is a continuously differentiable function and $C$ is a closed and convex set. In Chapter 2 we discussed the notion of *stationary points* of continuously differentiable functions, which are points in which the gradient vanishes. It was shown that stationarity is a necessary condition for a point to be an unconstrained local optimum point. The situation is more complicated when considering constrained problems of the form (P), where instead of looking at stationary points of a function, we need to consider the notion of stationary points of *a problem*.

**Definition 9.1 (stationary points of constrained problems).** *Let $f$ be a continuously differentiable function over a closed convex set $C$. Then $\mathbf{x}^* \in C$ is called a* **stationary point** *of (P) if $\nabla f(\mathbf{x}^*)^T(\mathbf{x}-\mathbf{x}^*) \geq 0$ for any $\mathbf{x} \in C$.*

 Stationarity actually means that there are no feasible descent directions of $f$ at $\mathbf{x}^*$. This suggests that stationarity is in fact a necessary condition for a local minimum of (P).

**Theorem 9.2 (stationarity as a necessary optimality condition).** *Let $f$ be a continuously differentiable function over a closed convex set $C$, and let $\mathbf{x}^*$ be a local minimum of (P). Then $\mathbf{x}^*$ is a stationary point of (P).*

*Proof.* Let $\mathbf{x}^*$ be a local minimum of (P), and assume in contradiction that $\mathbf{x}^*$ is not a stationary point of (P). Then there exists $\mathbf{x} \in C$ such that $\nabla f(\mathbf{x}^*)^T(\mathbf{x}-\mathbf{x}^*) < 0$. Therefore, $f'(\mathbf{x}^*; \mathbf{d}) < 0$ where $\mathbf{d} = \mathbf{x} - \mathbf{x}^*$, and hence, by Lemma 4.2, it follows that there exists $\varepsilon \in (0,1)$ such that $f(\mathbf{x}^* + t\mathbf{d}) < f(\mathbf{x}^*)$ for all $t \in (0, \varepsilon)$. Since $C$ is convex we have that $\mathbf{x}^* + t\mathbf{d} = (1-t)\mathbf{x}^* + t\mathbf{x} \in C$, leading to the conclusion that $\mathbf{x}^*$ is *not* a local minimum point of (P), in contradiction to the assumption that $\mathbf{x}^*$ is a local minimum point of (P). □

**Example 9.3** ($C = \mathbb{R}^n$)**.** If $C = \mathbb{R}^n$, then the stationary points of (P) are the points $\mathbf{x}^*$ satisfying

$$\nabla f(\mathbf{x}^*)^T (\mathbf{x} - \mathbf{x}^*) \geq 0 \tag{9.1}$$

for all $\mathbf{x} \in \mathbb{R}^n$. Plugging $\mathbf{x} = \mathbf{x}^* - \nabla f(\mathbf{x}^*)$ into (9.1), we obtain that $-\|\nabla f(\mathbf{x}^*)\|^2 \geq 0$, implying that

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

Therefore, it follows that the notion of a stationary point of a function and a stationary point of a *minimization problem* coincide when the problem is unconstrained. ∎

**Example 9.4** ($C = \mathbb{R}_+^n$)**.** Consider the optimization problem

$$\text{(Q)} \quad \begin{array}{ll} \min & f(\mathbf{x}) \\ \text{s.t.} & x_i \geq 0, \quad i = 1, 2, \ldots, n, \end{array}$$

where $f$ is a continuously differentiable function over $\mathbb{R}_+^n$. A vector $\mathbf{x}^* \in \mathbb{R}_+^n$ is a stationary point of (Q) if and only if

$$\nabla f(\mathbf{x}^*)^T (\mathbf{x} - \mathbf{x}^*) \geq 0 \text{ for all } \mathbf{x} \geq \mathbf{0},$$

which is the same as

$$\nabla f(\mathbf{x}^*)^T \mathbf{x} - \nabla f(\mathbf{x}^*)^T \mathbf{x}^* \geq 0 \text{ for all } \mathbf{x} \geq \mathbf{0}. \tag{9.2}$$

We will now use the following technical result:

$$\mathbf{a}^T \mathbf{x} + b \geq 0 \text{ for all } \mathbf{x} \geq \mathbf{0} \text{ if and only if } \mathbf{a} \geq \mathbf{0} \text{ and } b \geq 0.$$

Using this simple result, it follows that (9.2) holds if and only if

$$\nabla f(\mathbf{x}^*) \geq \mathbf{0} \text{ and } \nabla f(\mathbf{x}^*)^T \mathbf{x}^* \leq 0. \tag{9.3}$$

Since $\nabla f(\mathbf{x}^*) \geq \mathbf{0}$ and $\mathbf{x}^* \geq \mathbf{0}$, we can conclude that (9.3) holds if and only if

$$\nabla f(\mathbf{x}^*) \geq \mathbf{0} \text{ and } x_i^* \frac{\partial f}{\partial x_i}(\mathbf{x}^*) = 0, \quad i = 1, 2, \ldots, n.$$

We can compactly write the above condition as follows:

$$\frac{\partial f}{\partial x_i}(\mathbf{x}^*) \begin{cases} = 0, & x_i^* > 0, \\ \geq 0, & x_i^* = 0. \end{cases} \qquad \blacksquare$$

**Example 9.5 (stationarity over the unit-sum set).** Consider the optimization problem

$$\text{(R)} \quad \begin{array}{ll} \min & f(\mathbf{x}) \\ \text{s.t.} & \mathbf{e}^T \mathbf{x} = 1, \end{array}$$

where $f$ is a continuously differentiable function over $\mathbb{R}^n$. The feasible set

$$U = \{\mathbf{x} \in \mathbb{R}^n : \mathbf{e}^T \mathbf{x} = 1\} = \left\{ \mathbf{x} \in \mathbb{R}^n : \sum_{i=1}^n x_i = 1 \right\}$$

is called *the unit-sum set* . A point $\mathbf{x}^* \in U$ is a stationary point of (R) if and only if

(I) $\quad \nabla f(\mathbf{x}^*)^T(\mathbf{x} - \mathbf{x}^*) \geq 0$ for all $\mathbf{x}$ satisfying $\mathbf{e}^T\mathbf{x} = 1$.

We will show that condition (I) is equivalent to the following simple and explicit condition:

$$(\text{II}) \quad \frac{\partial f}{\partial x_1}(\mathbf{x}^*) = \frac{\partial f}{\partial x_2}(\mathbf{x}^*) = \cdots = \frac{\partial f}{\partial x_n}(\mathbf{x}^*).$$

We begin by showing that (II) implies (I). Indeed, assume that $\mathbf{x}^* \in U$ satisfies (II). Then for any $\mathbf{x} \in U$ one has

$$\begin{aligned}
\nabla f(\mathbf{x}^*)^T(\mathbf{x} - \mathbf{x}^*) &= \sum_{i=1}^n \frac{\partial f}{\partial x_i}(\mathbf{x}^*)(x_i - x_i^*) \\
&= \frac{\partial f}{\partial x_1}(\mathbf{x}^*)\left(\sum_{i=1}^n x_i - \sum_{i=1}^n x_i^*\right) = \frac{\partial f}{\partial x_1}(\mathbf{x}^*)(1-1) = 0,
\end{aligned}$$

and in particular $\nabla f(\mathbf{x}^*)^T(\mathbf{x} - \mathbf{x}^*) \geq 0$. We have thus shown that (I) is satisfied. To show the reverse direction, take $\mathbf{x}^* \in U$ that satisfies (I) and assume in contradiction that (II) does not hold. This means that there exist two different indices $i \neq j$ such that

$$\frac{\partial f}{\partial x_i}(\mathbf{x}^*) > \frac{\partial f}{\partial x_j}(\mathbf{x}^*).$$

Define the vector $\mathbf{x} \in U$ as

$$x_k = \begin{cases} x_k^*, & k \notin \{i, j\}, \\ x_i^* - 1, & k = i, \\ x_j^* + 1, & k = j. \end{cases}$$

Then

$$\begin{aligned}
\nabla f(\mathbf{x}^*)^T(\mathbf{x} - \mathbf{x}^*) &= \frac{\partial f}{\partial x_i}(\mathbf{x}^*)(x_i - x_i^*) + \frac{\partial f}{\partial x_j}(\mathbf{x}^*)(x_j - x_j^*) \\
&= -\frac{\partial f}{\partial x_i}(\mathbf{x}^*) + \frac{\partial f}{\partial x_j}(\mathbf{x}^*) \\
&< 0,
\end{aligned}$$

which is a contradiction to the assumption that (I) is satisfied. We thus conclude that (I) implies (II). ∎

**Example 9.6 (stationarity over the unit-ball).** Consider the optimization problem

$$(\text{S}) \quad \begin{array}{ll} \min & f(\mathbf{x}) \\ \text{s.t.} & \|\mathbf{x}\| \leq 1, \end{array}$$

where $f$ is a continuously differentiable function over $B[0,1]$. A point $\mathbf{x}^* \in B[0,1]$ is a stationary point of problem (S) if and only if

$$\nabla f(\mathbf{x}^*)^T(\mathbf{x} - \mathbf{x}^*) \geq 0 \text{ for all } \mathbf{x} \text{ satisfying } \|\mathbf{x}\| \leq 1, \tag{9.4}$$

which is equivalent to

$$\min\{\nabla f(\mathbf{x}^*)^T \mathbf{x} - \nabla f(\mathbf{x}^*)^T \mathbf{x}^* : ||\mathbf{x}|| \leq 1\} \geq 0.$$

We will now use the following fact:

[A] for any $\mathbf{a} \in \mathbb{R}^n$ the optimal value of the problem $\min\{\mathbf{a}^T \mathbf{x} : ||\mathbf{x}|| \leq 1\}$ is $-||\mathbf{a}||$.

Indeed, if $\mathbf{a} = \mathbf{0}$, then obviously the optimal value is $0 = -||\mathbf{a}||$. If $\mathbf{a} \neq \mathbf{0}$, then by the Cauchy–Schwarz inequality we have for any $\mathbf{x} \in B[\mathbf{0}, 1]$

$$\mathbf{a}^T \mathbf{x} \geq -||\mathbf{a}|| \cdot ||\mathbf{x}|| \geq -||\mathbf{a}||,$$

so that

$$\min\{\mathbf{a}^T \mathbf{x} : ||\mathbf{x}|| \leq 1\} \geq -||\mathbf{a}||.$$

On the other hand, the lower bound $-||\mathbf{a}||$ is attained at $\mathbf{x} = -\frac{\mathbf{a}}{||\mathbf{a}||}$, and hence fact [A] is established.

Returing to the characterization of stationary points over the unit ball, by fact [A] it follows that (9.4) is the same as

$$-\nabla f(\mathbf{x}^*)^T \mathbf{x}^* \geq ||\nabla f(\mathbf{x}^*)||. \tag{9.5}$$

However, we have by the Cauchy–Schwarz inequality that the following inequality holds:

$$-\nabla f(\mathbf{x}^*)^T \mathbf{x}^* \leq ||\nabla f(\mathbf{x}^*)|| \cdot ||\mathbf{x}^*|| \leq ||\nabla f(\mathbf{x}^*)||,$$

and we conclude that $\mathbf{x}^*$ is a stationary point of (S) if and only if

$$||\nabla f(\mathbf{x}^*)|| = -\nabla f(\mathbf{x}^*)^T \mathbf{x}^*. \tag{9.6}$$

The above condition is a simple and explicit characterization of the stationarity property. We can, however, find a more informative characterization of stationarity by the following argument: Let $\mathbf{x}^*$ be a point satisfying (9.6). We have two cases:

- If $\nabla f(\mathbf{x}^*) = \mathbf{0}$, then (9.6) holds automatically.

- If $\nabla f(\mathbf{x}^*) \neq \mathbf{0}$, then $||\mathbf{x}^*|| = 1$ since otherwise, if $||\mathbf{x}^*|| < 1$, we have by (once again!) the Cauchy–Schwarz inequality that

$$-\nabla f(\mathbf{x}^*)^T \mathbf{x}^* \leq ||\nabla f(\mathbf{x}^*)|| \cdot ||\mathbf{x}^*|| < ||\nabla f(\mathbf{x}^*)||,$$

which is a contradiction to (9.6). We therefore conclude that when $\nabla f(\mathbf{x}^*) \neq \mathbf{0}$, $\mathbf{x}^*$ is a stationary point if and only if $||\mathbf{x}^*|| = 1$ and

$$-\nabla f(\mathbf{x}^*)^T \mathbf{x}^* = ||\nabla f(\mathbf{x}^*)|| \cdot ||\mathbf{x}^*||. \tag{9.7}$$

The latter equality is equivalent to saying that there exists $\lambda \leq 0$ such that $\nabla f(\mathbf{x}^*) = \lambda \mathbf{x}^*$. To show this, note that if there exists such a $\lambda$, then indeed

$$-\nabla f(\mathbf{x}^*)^T \mathbf{x}^* = -\lambda ||\mathbf{x}^*||^2 \overset{\lambda \leq 0}{=} ||\lambda \mathbf{x}^*|| \cdot ||\mathbf{x}^*|| = ||\nabla f(\mathbf{x}^*)|| \cdot ||\mathbf{x}^*||,$$

so that (9.7) is satisfied. In the other direction, if (9.7) holds, then this means that the Cauchy–Schwarz inequality is satisfied as an equality, and hence, by Lemma 1.5, it follows that there exist $\lambda \in \mathbb{R}$ such that $\nabla f(\mathbf{x}^*) = \lambda \mathbf{x}^*$. The parameter $\lambda$ must be

nonpositive since otherwise the left-hand side of (9.7) would be negative, and the right-hand side would be positive, contradicting the equality in (9.7).

In conclusion, $\mathbf{x}^*$ is a stationary point of (S) if and only if either $\nabla f(\mathbf{x}^*) = \mathbf{0}$ or $\|\mathbf{x}^*\| = 1$ and there exists $\lambda \leq 0$ such that $\nabla f(\mathbf{x}^*) = \lambda \mathbf{x}^*$. ∎

To summarize the four examples, we write explicitly each of the stationarity conditions in the following table.

| feasible set | explicit stationarity condition |
|---|---|
| $\mathbb{R}^n$ | $\nabla f(\mathbf{x}^*) = \mathbf{0}$ |
| $\mathbb{R}^n_+$ | $\frac{\partial f}{\partial x_i}(\mathbf{x}^*) \begin{cases} = 0, & x_i^* > 0 \\ \geq 0, & x_i^* = 0 \end{cases}$ |
| $\{\mathbf{x} \in \mathbb{R}^n : \mathbf{e}^T \mathbf{x} = 1\}$ | $\frac{\partial f}{\partial x_1}(\mathbf{x}^*) = \cdots = \frac{\partial f}{\partial x_n}(\mathbf{x}^*)$ |
| $B[\mathbf{0}, 1]$ | $\nabla f(\mathbf{x}^*) = \mathbf{0}$ or $\|\mathbf{x}^*\| = 1$ and $\exists \lambda \leq 0 : \nabla f(\mathbf{x}^*) = \lambda \mathbf{x}^*$ |

## 9.2 ▪ Stationarity in Convex Problems

Stationarity is a necessary optimality condition for local optimality. However, when the objective function is additionally assumed to be convex, stationarity is a necessary *and sufficient* condition for optimality.

**Theorem 9.7.** *Let $f$ be a continuously differentiable convex function over a closed and convex set $C \subseteq \mathbb{R}^n$. Then $\mathbf{x}^* \in C$ is a stationary point of*

$$\text{(P)} \quad \begin{array}{ll} \min & f(\mathbf{x}) \\ s.t. & \mathbf{x} \in C \end{array}$$

*if and only if $\mathbf{x}^*$ is an optimal solution of (P).*

**Proof.** If $\mathbf{x}^*$ is an optimal solution of (P), then by Theorem 9.2, it follows that $\mathbf{x}^*$ is a stationary point of (P). To prove the sufficiency of the stationarity condition, assume that $\mathbf{x}^*$ is a stationary point of (P), and let $\mathbf{x} \in C$. Then

$$f(\mathbf{x}) \geq f(\mathbf{x}^*) + \nabla f(\mathbf{x}^*)^T(\mathbf{x} - \mathbf{x}^*) \geq f(\mathbf{x}^*),$$

where the first inequality follows from the gradient inequality for convex functions (Theorem 7.6) and the second inequality follows from the definition of a stationary point. We have thus shown that $\mathbf{x}^*$ is the global minimum point of (P), and the reverse direction is established. □

## 9.3 ▪ The Orthogonal Projection Revisited

We can use the stationarity property in order to establish an important property of the orthogonal projection operator. This characterization will be called *the second projection theorem*. Geometrically it states that for a given closed and convex set $C$, $\mathbf{x} \in \mathbb{R}^n$, and $\mathbf{y} \in C$, the angle between $\mathbf{x} - P_C(\mathbf{x})$ and $\mathbf{y} - P_C(\mathbf{x})$ is greater than or equal to 90 degrees. This phenomenon is illustrated in Figure 9.1.
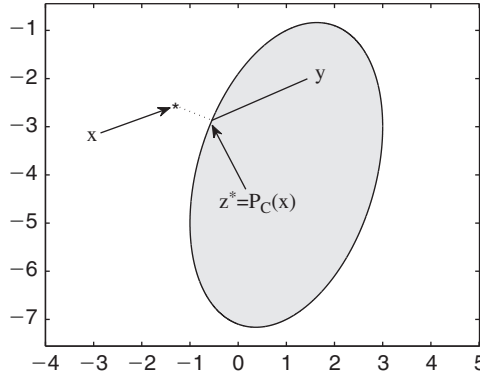
**Figure 9.1.** *The orthogonal projection operator.*

**Theorem 9.8 (second projection theorem).** *Let $C$ be a closed convex set and let $\mathbf{x} \in \mathbb{R}^n$. Then $\mathbf{z} = P_C(\mathbf{x})$ if and only if $\mathbf{z} \in C$ and*

$$(\mathbf{x} - \mathbf{z})^T(\mathbf{y} - \mathbf{z}) \leq 0 \text{ for any } \mathbf{y} \in C. \tag{9.8}$$

*Proof.* $\mathbf{z} = P_C(\mathbf{x})$ if and only if it is the optimal solution of the problem

$$\begin{aligned} \min \quad & g(\mathbf{y}) \equiv \|\mathbf{y} - \mathbf{x}\|^2 \\ \text{s.t.} \quad & \mathbf{y} \in C. \end{aligned}$$

Therefore, by Theorem 9.7 it follows that $\mathbf{z} = P_C(\mathbf{x})$ if and only if

$$\nabla g(\mathbf{z})^T(\mathbf{y} - \mathbf{z}) \geq 0 \text{ for all } \mathbf{y} \in C,$$

which is the same as (9.8). $\quad\square$

Another important property of the orthogonal projection operator is given in the following theorem, which also establishes the so-called *nonexpansiveness* property of $P_C$.

**Theorem 9.9.** *Let $C$ be a nonempty closed and convex set. Then*

1. *for any $\mathbf{v}, \mathbf{w} \in \mathbb{R}^n$*

$$(P_C(\mathbf{v}) - P_C(\mathbf{w}))^T(\mathbf{v} - \mathbf{w}) \geq \|P_C(\mathbf{v}) - P_C(\mathbf{w})\|^2, \tag{9.9}$$

2. **(nonexpansiveness)** *for any $\mathbf{v}, \mathbf{w} \in \mathbb{R}^n$*

$$\|P_C(\mathbf{v}) - P_C(\mathbf{w})\| \leq \|\mathbf{v} - \mathbf{w}\|. \tag{9.10}$$

*Proof.* Recall that by Theorem 9.8 we have that for any $\mathbf{x} \in \mathbb{R}^n$ and $\mathbf{y} \in C$

$$(\mathbf{x} - P_C(\mathbf{x}))^T(\mathbf{y} - P_C(\mathbf{x})) \leq 0. \tag{9.11}$$

Substituting $\mathbf{x} = \mathbf{v}$ and $\mathbf{y} = P_C(\mathbf{w})$, we have

$$(\mathbf{v} - P_C(\mathbf{v}))^T(P_C(\mathbf{w}) - P_C(\mathbf{v})) \leq 0. \tag{9.12}$$

Substituting $\mathbf{x} = \mathbf{w}$ and $\mathbf{y} = P_C(\mathbf{v})$, we obtain

$$(\mathbf{w} - P_C(\mathbf{w}))^T (P_C(\mathbf{v}) - P_C(\mathbf{w})) \leq 0. \tag{9.13}$$

Adding the two inequalities (9.12) and (9.13) yields

$$(P_C(\mathbf{w}) - P_C(\mathbf{v}))^T (\mathbf{v} - \mathbf{w} + P_C(\mathbf{w}) - P_C(\mathbf{v})) \leq 0,$$

and hence,

$$(P_C(\mathbf{v}) - P_C(\mathbf{w}))^T (\mathbf{v} - \mathbf{w}) \geq \|P_C(\mathbf{v}) - P_C(\mathbf{w})\|^2,$$

showing the desired inequality (9.9).

To prove (9.10), note that if $P_C(\mathbf{v}) = P_C(\mathbf{w})$, the inequality is trivial. If $P_C(\mathbf{v}) \neq P_C(\mathbf{w})$, then by the Cauchy–Schwarz inequality we have

$$(P_C(\mathbf{v}) - P_C(\mathbf{w}))^T (\mathbf{v} - \mathbf{w}) \leq \|P_C(\mathbf{v}) - P_C(\mathbf{w})\| \cdot \|\mathbf{v} - \mathbf{w}\|,$$

which combined with (9.9) yields the inequality

$$\|P_C(\mathbf{v}) - P_C(\mathbf{w})\|^2 \leq \|P_C(\mathbf{v}) - P_C(\mathbf{w})\| \cdot \|\mathbf{v} - \mathbf{w}\|.$$

Dividing by $\|P_C(\mathbf{v}) - P_C(\mathbf{w})\|$ implies (9.10).  $\square$

Coming back to stationarity, the next result describes an additional useful representation of stationarity in terms of the orthogonal projection operator.

**Theorem 9.10.** *Let $f$ be a continuously differentiable function defined on the closed and convex set $C$, and let $s > 0$. Then $\mathbf{x}^* \in C$ is a stationary point of the problem*

$$\text{(P)} \quad \begin{array}{ll} \min & f(\mathbf{x}) \\ s.t. & \mathbf{x} \in C \end{array}$$

*if and only if*

$$\mathbf{x}^* = P_C(\mathbf{x}^* - s\nabla f(\mathbf{x}^*)). \tag{9.14}$$

**Proof.** By the second projection theorem (Theorem 9.8), it follows that $\mathbf{x}^* = P_C(\mathbf{x}^* - s\nabla f(\mathbf{x}^*))$ if and only if

$$(\mathbf{x}^* - s\nabla f(\mathbf{x}^*) - \mathbf{x}^*)^T (\mathbf{x} - \mathbf{x}^*) \leq 0 \text{ for any } \mathbf{x} \in C.$$

The above relation is equivalent to

$$\nabla f(\mathbf{x}^*)^T (\mathbf{x} - \mathbf{x}^*) \geq 0 \text{ for any } \mathbf{x} \in C,$$

namely to stationarity.  $\square$

Note that condition (9.14) seems to depend on the parameter $s$, but by its equivalence to stationarity, it is essentially *independent* of $s$.

## 9.4 ▪ The Gradient Projection Method

The stationarity condition

$$\mathbf{x}^* = P_C(\mathbf{x}^* - s\nabla f(\mathbf{x}^*)) \tag{9.15}$$

naturally motivates the following algorithm, called *the gradient projection method*, for solving problem (P). This algorithm can be seen as a fixed point method for solving the equation (9.15).

---

**The Gradient Projection Method**

**Input:** $\varepsilon > 0$ - tolerance parameter.

**Initialization:** Pick $\mathbf{x}_0 \in C$ arbitrarily.
**General step:** For any $k = 0, 1, 2, \ldots$ execute the following steps:

(a) Pick a stepsize $t_k$ by a line search procedure.

(b) Set $\mathbf{x}_{k+1} = P_C(\mathbf{x}_k - t_k \nabla f(\mathbf{x}_k))$.

(c) If $\|\mathbf{x}_k - \mathbf{x}_{k+1}\| \leq \varepsilon$, then STOP, and $\mathbf{x}_{k+1}$ is the output.

---

Obviously, in the unconstrained case, that is, when $C = \mathbb{R}^n$, the gradient projection method is just the gradient method studied in Chapter 4.

There are several strategies for choosing the stepsizes $t_k$. We will consider two choices:

- **constant stepsize** $t_k = \bar{t}$ for all $k$.

- **backtracking**.

We will elaborate on the specific details of the backtracking procedure in what follows.

## 9.4.1 ▪ Sufficient Decrease and the Gradient Mapping

To establish the convergence of the method, we will prove a sufficient decrease lemma for $C^{1,1}$ functions, similarly to the sufficient decrease lemma proved for the gradient method (Lemma 4.23).

**Lemma 9.11 (sufficient decrease lemma for constrained problems).** *Suppose that $f \in C_L^{1,1}(C)$, where $C$ is a closed convex set. Then for any $\mathbf{x} \in C$ and $t \in (0, \frac{2}{L})$ the following inequality holds:*

$$f(\mathbf{x}) - f(P_C(\mathbf{x} - t\nabla f(\mathbf{x}))) \geq t\left(1 - \frac{Lt}{2}\right)\left\|\frac{1}{t}(\mathbf{x} - P_C(\mathbf{x} - t\nabla f(\mathbf{x})))\right\|^2. \tag{9.16}$$

**Proof.** For the sake of simplicity, we will make the notation $\mathbf{x}^+ = P_C(\mathbf{x} - t\nabla f(\mathbf{x}))$. By the descent lemma (Lemma 4.22) we have that

$$f(\mathbf{x}^+) \leq f(\mathbf{x}) + \langle \nabla f(\mathbf{x}), \mathbf{x}^+ - \mathbf{x} \rangle + \frac{L}{2}\|\mathbf{x} - \mathbf{x}^+\|^2. \tag{9.17}$$

From the second projection theorem (Theorem 9.8) we have

$$\langle \mathbf{x} - t\nabla f(\mathbf{x}) - \mathbf{x}^+, \mathbf{x} - \mathbf{x}^+ \rangle \leq 0,$$

from which it follows that

$$\langle \nabla f(\mathbf{x}), \mathbf{x}^+ - \mathbf{x} \rangle \leq -\frac{1}{t}\|\mathbf{x}^+ - \mathbf{x}\|^2,$$

which combined with (9.17) yields

$$f(\mathbf{x}^+) \leq f(\mathbf{x}) + \left(-\frac{1}{t} + \frac{L}{2}\right) \|\mathbf{x}^+ - \mathbf{x}\|^2.$$

Hence, taking into account the definition of $\mathbf{x}^+$, the desired result follows. $\quad\square$

The result of Lemma 9.11 is a generalization of the sufficient decrease property derived in the unconstrained setting in Lemma 4.23. In fact, when $C = \mathbb{R}^n$ the obtained inequality is exactly the same as the one obtained in Lemma 4.23.

It is convenient to define *the gradient mapping* as

$$G_M(\mathbf{x}) = M\left[\mathbf{x} - P_C\left(\mathbf{x} - \frac{1}{M}\nabla f(\mathbf{x})\right)\right],$$

where $M > 0$. We assume that the identities of $f$ and $C$ are clear from the context. Note that in the unconstrained case $G_M(\mathbf{x}) = \nabla f(\mathbf{x})$, so the gradient mapping is an extension of the usual gradient operation. In addition, by Theorem 9.10, $G_M(\mathbf{x}) = \mathbf{0}$ if and only if $\mathbf{x}$ is a stationary point of (P). This means that we can look at $\|G_L(\mathbf{x})\|$ as an *optimality measure*. The result of Lemma 9.11 essentially states that

$$f(\mathbf{x}) - f(P_C(\mathbf{x} - t\nabla f(\mathbf{x}))) \geq t\left(1 - \frac{Lt}{2}\right)\|G_{\frac{1}{t}}(\mathbf{x})\|^2. \tag{9.18}$$

This generalized sufficient decrease property allows us to prove similar results to those proven in the unconstrained case. Before proceeding to the convergence analysis, we will prove an important monotonicity property of the norm of the gradient mapping $G_M(\mathbf{x})$ with respect to the parameter $M$.

**Lemma 9.12.** *Let $f$ be a continuously differentiable function defined on a nonempty closed and convex set $C$. Suppose that $L_1 \geq L_2 > 0$. Then*

$$\|G_{L_1}(\mathbf{x})\| \geq \|G_{L_2}(\mathbf{x})\| \tag{9.19}$$

*and*

$$\frac{\|G_{L_1}(\mathbf{x})\|}{L_1} \leq \frac{\|G_{L_2}(\mathbf{x})\|}{L_2} \tag{9.20}$$

*for any $\mathbf{x} \in \mathbb{R}^n$.*

**Proof.** Recall that, by the second projection theorem, for any $\mathbf{v} \in \mathbb{R}^n$ and $\mathbf{w} \in C$ the following inequality holds:

$$\langle \mathbf{v} - P_C(\mathbf{v}), P_C(\mathbf{v}) - \mathbf{w} \rangle \geq 0.$$

Plugging $\mathbf{v} = \mathbf{x} - \frac{1}{L_1}\nabla f(\mathbf{x})$ and $\mathbf{w} = P_C(\mathbf{x} - \frac{1}{L_2}\nabla f(\mathbf{x}))$ in the latter inequality it follows that

$$\left\langle \mathbf{x} - \frac{1}{L_1}\nabla f(\mathbf{x}) - P_C\left(\mathbf{x} - \frac{1}{L_1}\nabla f(\mathbf{x})\right), P_C\left(\mathbf{x} - \frac{1}{L_1}\nabla f(\mathbf{x})\right) - P_C\left(\mathbf{x} - \frac{1}{L_2}\nabla f(\mathbf{x})\right) \right\rangle \geq 0,$$

or

$$\left\langle \frac{1}{L_1}G_{L_1}(\mathbf{x}) - \frac{1}{L_1}\nabla f(\mathbf{x}), \frac{1}{L_2}G_{L_2}(\mathbf{x}) - \frac{1}{L_1}G_{L_1}(\mathbf{x}) \right\rangle \geq 0.$$

Exchanging the roles of $L_1$ and $L_2$ yields the following inequality:

$$\left\langle \frac{1}{L_2} G_{L_2}(\mathbf{x}) - \frac{1}{L_2} \nabla f(\mathbf{x}), \frac{1}{L_1} G_{L_1}(\mathbf{x}) - \frac{1}{L_2} G_{L_2}(\mathbf{x}) \right\rangle \geq 0.$$

Multiplying the first inequality by $L_1$ and the second by $L_2$ and adding them we obtain

$$\left\langle G_{L_1}(\mathbf{x}) - G_{L_2}(\mathbf{x}), \frac{1}{L_2} G_{L_2}(\mathbf{x}) - \frac{1}{L_1} G_{L_1}(\mathbf{x}) \right\rangle \geq 0,$$

which, after some expansion of terms, can be seen to be the same as

$$\frac{1}{L_1} \|G_{L_1}(\mathbf{x})\|^2 + \frac{1}{L_2} \|G_{L_2}(\mathbf{x})\|^2 \leq \left( \frac{1}{L_1} + \frac{1}{L_2} \right) G_{L_1}(\mathbf{x})^T G_{L_2}(\mathbf{x}).$$

Using the Cauchy–Schwarz inequality we obtain that

$$\frac{1}{L_1} \|G_{L_1}(\mathbf{x})\|^2 + \frac{1}{L_2} \|G_{L_2}(\mathbf{x})\|^2 \leq \left( \frac{1}{L_1} + \frac{1}{L_2} \right) \|G_{L_1}(\mathbf{x})\| \cdot \|G_{L_2}(\mathbf{x})\|. \qquad (9.21)$$

Note that if $G_{L_2}(\mathbf{x}) = \mathbf{0}$, then by the latter inequality, it follows that $G_{L_1}(\mathbf{x}) = \mathbf{0}$, implying that in this case the inequalities (9.19) and (9.20) hold trivially. Assume then that $G_{L_2}(\mathbf{x}) \neq \mathbf{0}$, and define $t = \frac{\|G_{L_1}(\mathbf{x})\|}{\|G_{L_2}(\mathbf{x})\|}$. Then by (9.21)

$$\frac{1}{L_1} t^2 - \left( \frac{1}{L_1} + \frac{1}{L_2} \right) t + \frac{1}{L_2} \leq 0.$$

Since the roots of the quadratic function are $t = 1, \frac{L_1}{L_2}$, we obtain that

$$1 \leq t \leq \frac{L_1}{L_2},$$

showing that

$$\|G_{L_2}(\mathbf{x})\| \leq \|G_{L_1}(\mathbf{x})\| \leq \frac{L_1}{L_2} \|G_{L_2}(\mathbf{x})\|. \qquad \square$$

### 9.4.2 ▪ Backtracking

Equipped with the definition of the gradient mapping, we are now able to define the backtracking stepsize selection strategy similarly to the definition of the backtracking procedure for descent methods for unconstrained problems (see Chapter 4). The procedure requires three parameters $(s, \alpha, \beta)$, where $s > 0, \alpha \in (0, 1), \beta \in (0, 1)$. The choice of $t_k$ is done as follows: First, $t_k$ is set to be equal to the initial guess $s$. Then while

$$f(\mathbf{x}_k) - f(P_C(\mathbf{x}_k - t_k \nabla f(\mathbf{x}_k))) < \alpha t_k \|G_{\frac{1}{t_k}}(\mathbf{x}_k)\|^2$$

we set $t_k := \beta t_k$. In other words, the stepsize is chosen as $t_k = s\beta^{i_k}$, where $i_k$ is the smallest nonnegative integer for which the condition

$$f(\mathbf{x}_k) - f(P_C(\mathbf{x}_k - s\beta^{i_k} \nabla f(\mathbf{x}_k))) \geq \alpha s\beta^{i_k} \|G_{\frac{1}{s\beta^{i_k}}}(\mathbf{x}_k)\|^2$$

is satisfied.

Note that the backtracking procedure is finite for $C^{1,1}$ functions. Indeed, if $f \in C_L^{1,1}(C)$, then plugging $\mathbf{x} = \mathbf{x}_k$ into (9.18) we obtain

$$f(\mathbf{x}_k) - f(P_C(\mathbf{x}_k - t\nabla f(\mathbf{x}_k))) \geq t\left(1 - \frac{Lt}{2}\right)\left\|G_{\frac{1}{t}}(\mathbf{x}_k)\right\|^2.$$

Hence, since the inequality $t \leq \frac{2(1-\alpha)}{L}$ is the same as $1 - \frac{Lt}{2} \geq \alpha$, we conclude that for any $t \leq \frac{2(1-\alpha)}{L}$ the inequality

$$f(\mathbf{x}_k) - f(P_C(\mathbf{x}_k - t\nabla f(\mathbf{x}_k))) \geq \alpha t\left\|G_{\frac{1}{t}}(\mathbf{x}_k)\right\|^2$$

holds, implying that the backtracking procedure ends when $t_k$ is smaller or equal to $\frac{2(1-\alpha)}{L}$.

We can also compute a lower bound on $t_k$: either $t_k$ is equal to $s$ or the backtracking procedure is invoked, meaning that the stepsize $\frac{t_k}{\beta}$ did not satisfy the backtracking condition given by

$$f(\mathbf{x}_k) - f(P_C(\mathbf{x}_k - t_k\nabla f(\mathbf{x}_k))) \geq \alpha t_k \|G_{\frac{1}{t_k}}(\mathbf{x}_k)\|^2.$$

In particular, by the above discussion we can conclude that $\frac{t_k}{\beta} > \frac{2(1-\alpha)}{L}$, so that $t_k > \frac{2(1-\alpha)\beta}{L}$. To summarize, in the backtracking procedure, the chosen stepsize $t_k$ satisfies

$$t_k \geq \min\left\{s, \frac{2(1-\alpha)\beta}{L}\right\}. \tag{9.22}$$

### 9.4.3 ▪ Convergence of the Gradient Projection Method

We will analyze the convergence of the gradient projection method in the case when $f \in C_L^{1,1}(C)$. We begin with the following lemma showing the sufficient decrease of consecutive function values of the sequence generated by the gradient projection method.

**Lemma 9.13.** *Consider the problem*

$$\text{(P)} \quad \begin{array}{ll} \min & f(\mathbf{x}) \\ s.t. & \mathbf{x} \in C, \end{array}$$

*where $C \subseteq \mathbb{R}^n$ is a nonempty closed and convex set and $f \in C_L^{1,1}(C)$. Let $\{\mathbf{x}_k\}_{k\geq 0}$ be the sequence generated by the gradient projection method for solving problem* (P) *with either a constant stepsize $t_k = \bar{t} \in (0, \frac{2}{L})$ or with a stepsize chosen by the backtracking procedure with parameters $(s, \alpha, \beta)$ satisfying $s > 0, \alpha \in (0,1), \beta \in (0,1)$. Then for any $k \geq 0$*

$$f(\mathbf{x}_k) - f(\mathbf{x}_{k+1}) \geq M\|G_d(\mathbf{x}_k)\|^2, \tag{9.23}$$

*where*

$$M = \begin{cases} \bar{t}\left(1 - \frac{\bar{t}L}{2}\right) & \text{constant stepsize,} \\ \alpha\min\left\{s, \frac{2(1-\alpha)\beta}{L}\right\} & \text{backtracking} \end{cases}$$

*and*

$$d = \begin{cases} 1/\bar{t} & \text{constant stepsize,} \\ 1/s & \text{backtracking.} \end{cases}$$

***Proof.*** The result for the constant stepsize setting follows by plugging $t = \bar{t}$ and $\mathbf{x} = \mathbf{x}_k$ in (9.18). As for the backtracking procedure, by its definition we have

$$f(\mathbf{x}_k) - f(\mathbf{x}_{k+1}) \geq \alpha t_k \|G_{\frac{1}{t_k}}(\mathbf{x}_k)\|^2.$$

The result now follows from the lower bound on $t_k$ given in (9.22) and the fact that $t_k \leq s$, which implies by the monotonocity property of the gradient mapping (Lemma 9.12) that $\|G_{1/t_k}(\mathbf{x}_k)\| \geq \|G_{1/s}(\mathbf{x}_k)\|$.   □

We are now ready to prove the convergence of the norm of the gradient mapping to zero.

**Theorem 9.14 (convergence of the gradient projection method).** *Consider the problem*

$$\text{(P)} \quad \begin{array}{ll} \min & f(\mathbf{x}) \\ s.t. & \mathbf{x} \in C, \end{array}$$

*where $C$ is a nonempty closed and convex set and $f \in C_L^{1,1}(C)$ is bounded below. Let $\{\mathbf{x}_k\}_{k \geq 0}$ be the sequence generated by the gradient projection method for solving problem (P) with either a constant stepsize $t_k = \bar{t} \in (0, \frac{2}{L})$ or a stepsize chosen by the backtracking procedure with parameters $(s, \alpha, \beta)$ satisfying $s > 0, \alpha, \beta \in (0, 1)$. Then we have the following:*

(a) *The sequence $\{f(\mathbf{x}_k)\}$ is nonincreasing. In addition, $f(\mathbf{x}_{k+1}) < f(\mathbf{x}_k)$ unless $\mathbf{x}_k$ is a stationary point of (P).*

(b) *$G_d(\mathbf{x}_k) \to \mathbf{0}$ as $k \to \infty$, where $d$ is as defined in Lemma 9.13.*

***Proof.*** (a) By Lemma 9.13 we have that

$$f(\mathbf{x}_k) - f(\mathbf{x}_{k+1}) \geq M\|G_d(\mathbf{x}_k)\|^2. \tag{9.24}$$

Since $M > 0$, it readily follows that $f(\mathbf{x}_k) \geq f(\mathbf{x}_{k+1})$ and equality holds only if $G_d(\mathbf{x}_k) = \mathbf{0}$, which is the same as saying that $\mathbf{x}_k$ is a stationary point of (P).

(b) Since the sequence $\{f(\mathbf{x}_k)\}_{k \geq 0}$ is nonincreasing and bounded below, it converges. Thus, in particular $f(\mathbf{x}_k) - f(\mathbf{x}_{k+1}) \to 0$ as $k \to \infty$, which combined with (9.24) implies that $\|G_d(\mathbf{x}_k)\| \to 0$ as $k \to \infty$.   □

We can also get an estimate for the rate of convergence of the gradient projection method, but in the constrained case it will be in terms of the norm of the gradient mapping.

**Theorem 9.15 (rate of convergence of gradient mapping norms in the gradient projection method).** *Under the setting of Theorem 9.14, let $f^*$ be the limit of the convergent sequence $\{f(\mathbf{x}_k)\}_{k \geq 0}$. Then for any $n = 0, 1, 2, \ldots$*

$$\min_{k=0,1,\ldots,n} \|G_d(\mathbf{x}_k)\| \leq \sqrt{\frac{f(\mathbf{x}_0) - f^*}{M(n+1)}},$$

*where*

$$M = \begin{cases} \bar{t}\left(1 - \frac{\bar{t}L}{2}\right) & \text{constant stepsize,} \\ \alpha \min\left\{s, \frac{2(1-\alpha)\beta}{L}\right\} & \text{backtracking} \end{cases}$$

*and*

$$d = \begin{cases} 1/\bar{t} & \text{constant stepsize,} \\ 1/s & \text{backtracking.} \end{cases}$$

**Proof.** Summing the inequality (9.23) over $k = 0, 1, \ldots, n$ we obtain

$$f(\mathbf{x}_0) - f(\mathbf{x}_{n+1}) \geq M \sum_{k=0}^{n} \|G_d(\mathbf{x}_k)\|^2.$$

Since $f(\mathbf{x}_{n+1}) \geq f^*$, we can thus write

$$f(\mathbf{x}_0) - f^* \geq M \sum_{k=0}^{n} \|G_d(\mathbf{x}_k)\|^2.$$

Finally, using the latter inequality along with the obvious fact that

$$\sum_{k=0}^{n} \|G_d(\mathbf{x}_k)\|^2 \geq (n+1) \min_{k=0,1,\ldots,n} \|G_d(\mathbf{x}_k)\|^2,$$

it follows that

$$f(\mathbf{x}_0) - f^* \geq M(n+1) \min_{k=0,1,\ldots,n} \|G_d(\mathbf{x}_k)\|^2,$$

implying the desired result. $\quad\square$

### 9.4.4 ▪ The Convex Case

When the objective function is assumed to be convex, the constrained problem (P) becomes convex, and stronger convergence results can be established. We will concentrate on the constant stepsize setting and show that the sequence $\{\mathbf{x}_k\}_{k\geq 0}$ converges to an optimal solution of (P). In addition, we will establish a rate of convergence result for the sequence of function values $\{f(\mathbf{x}_k)\}_{k\geq 0}$ to the optimal value.

**Theorem 9.16 (rate of convergence of the sequence of function values).** *Consider the problem*

$$\text{(P)} \quad \begin{array}{ll} \min & f(\mathbf{x}) \\ s.t. & \mathbf{x} \in C, \end{array}$$

*where $C$ is a nonempty closed and convex set and $f \in C_L^{1,1}(C)$ is convex over $C$. Let $\{\mathbf{x}_k\}_{k\geq 0}$ be the sequence generated by the gradient projection method for solving problem* (P) *with a constant stepsize $t_k = \bar{t} \in (0, \frac{1}{L}]$. Assume that the set of optimal solutions, denoted by $X^*$, is nonempty, and let $f^*$ be the optimal value of* (P)*. Then*

(a) *for any $k \geq 0$ and $\mathbf{x}^* \in X^*$*

$$2\bar{t}(f(\mathbf{x}_{k+1}) - f(\mathbf{x}^*)) \leq \|\mathbf{x}_k - \mathbf{x}^*\|^2 - \|\mathbf{x}_{k+1} - \mathbf{x}^*\|^2, \tag{9.25}$$

(b) *for any $n \geq 0$:*

$$f(\mathbf{x}_n) - f^* \leq \frac{\|\mathbf{x}_0 - \mathbf{x}^*\|^2}{2\bar{t}n}. \tag{9.26}$$

***Proof.*** By the descent lemma we have

$$f(\mathbf{x}_{k+1}) \leq f(\mathbf{x}_k) + \left\langle \nabla f(\mathbf{x}_k), \mathbf{x}_{k+1} - \mathbf{x}_k \right\rangle + \frac{L}{2}\|\mathbf{x}_k - \mathbf{x}_{k+1}\|^2. \tag{9.27}$$

Let $\mathbf{x}^*$ be an optimal solution of (P). Then the gradient inequality implies that $f(\mathbf{x}_k) \leq f(\mathbf{x}^*) + \left\langle \nabla f(\mathbf{x}_k), \mathbf{x}_k - \mathbf{x}^* \right\rangle$, which combined with (9.27) yields

$$f(\mathbf{x}_{k+1}) \leq f(\mathbf{x}^*) + \left\langle \nabla f(\mathbf{x}_k), \mathbf{x}_k - \mathbf{x}^* \right\rangle + \left\langle \nabla f(\mathbf{x}_k), \mathbf{x}_{k+1} - \mathbf{x}_k \right\rangle + \frac{L}{2}\|\mathbf{x}_k - \mathbf{x}_{k+1}\|^2. \tag{9.28}$$

By the second projection theorem (Theorem 9.8) we have that

$$\left\langle \mathbf{x}_k - \bar{t}\nabla f(\mathbf{x}_k) - \mathbf{x}_{k+1}, \mathbf{x}^* - \mathbf{x}_{k+1} \right\rangle \leq 0,$$

so that

$$\left\langle \nabla f(\mathbf{x}_k), \mathbf{x}_{k+1} - \mathbf{x}^* \right\rangle \leq \frac{1}{\bar{t}} \left\langle \mathbf{x}_k - \mathbf{x}_{k+1}, \mathbf{x}_{k+1} - \mathbf{x}^* \right\rangle. \tag{9.29}$$

Therefore, combining (9.28), (9.29) and using the inequality $\bar{t} \leq \frac{1}{L}$, we obtain that

$$\begin{aligned}
f(\mathbf{x}_{k+1}) &\leq f(\mathbf{x}^*) + \left\langle \nabla f(\mathbf{x}_k), \mathbf{x}_k - \mathbf{x}^* \right\rangle + \left\langle \nabla f(\mathbf{x}_k), \mathbf{x}_{k+1} - \mathbf{x}_k \right\rangle + \frac{L}{2}\|\mathbf{x}_k - \mathbf{x}_{k+1}\|^2 \\
&= f(\mathbf{x}^*) + \left\langle \nabla f(\mathbf{x}_k), \mathbf{x}_{k+1} - \mathbf{x}^* \right\rangle + \frac{L}{2}\|\mathbf{x}_k - \mathbf{x}_{k+1}\|^2 \\
&\leq f(\mathbf{x}^*) + \frac{1}{\bar{t}}\left\langle \mathbf{x}_k - \mathbf{x}_{k+1}, \mathbf{x}_{k+1} - \mathbf{x}^* \right\rangle + \frac{L}{2}\|\mathbf{x}_k - \mathbf{x}_{k+1}\|^2 \\
&\leq f(\mathbf{x}^*) + \frac{1}{\bar{t}}\left\langle \mathbf{x}_k - \mathbf{x}_{k+1}, \mathbf{x}_{k+1} - \mathbf{x}^* \right\rangle + \frac{1}{2\bar{t}}\|\mathbf{x}_k - \mathbf{x}_{k+1}\|^2 \\
&= f(\mathbf{x}^*) + \frac{1}{2\bar{t}}\left(\|\mathbf{x}_k - \mathbf{x}^*\|^2 - \|\mathbf{x}_{k+1} - \mathbf{x}^*\|^2\right),
\end{aligned}$$

establishing part (a).

To prove part (b), sum the inequalities (9.25) for $k = 0, 1, 2, \ldots, n-1$ and obtain

$$\|\mathbf{x}_n - \mathbf{x}^*\|^2 - \|\mathbf{x}_0 - \mathbf{x}^*\|^2 \leq 2\bar{t}\sum_{k=0}^{n-1}(f(\mathbf{x}^*) - f(\mathbf{x}_{k+1})) \leq 2\bar{t}n(f(\mathbf{x}^*) - f(\mathbf{x}_n)),$$

where we used in the last inequality the monotonicity of the functions values of the generated sequence. Hence,

$$f(\mathbf{x}_n) - f^* \leq \frac{\|\mathbf{x}_0 - \mathbf{x}^*\|^2 - \|\mathbf{x}_n - \mathbf{x}^*\|^2}{2\bar{t}n} \leq \frac{\|\mathbf{x}_0 - \mathbf{x}^*\|^2}{2\bar{t}n},$$

as required. $\quad\square$

Note that when the constant stepsize is chosen as $\bar{t} = \frac{1}{L}$, the result (9.26) then reads as

$$f(\mathbf{x}_n) - f^* \leq \frac{L\|\mathbf{x}_0 - \mathbf{x}^*\|^2}{2n}.$$

The rate of convergence of $O(1/n)$ obtained in Theorem 9.16 is called a *sublinear rate of convergence*. We can also deduce the convergence of the sequence itself, and for that,

we note that by (9.25) the sequence satisfies a property called *Fejér monotonicity*, which we state explicitly in the following lemma.

**Lemma 9.17 (Fejér monotonicity of the sequence generated by the gradient projection method).** *Under the setting of Theorem* 9.16 *the sequence* $\{\mathbf{x}_k\}_{k \geq 0}$ *satisfies the following inequality for any* $\mathbf{x}^* \in X^*$ *and* $k \geq 0$*:*

$$\|\mathbf{x}_{k+1} - \mathbf{x}^*\| \leq \|\mathbf{x}_k - \mathbf{x}^*\|.$$

*Proof.* The proof follows by (9.25) and the fact that $f(\mathbf{x}^*) \leq f(\mathbf{x}_{k+1})$. $\quad\square$

We are now ready to prove the convergence of the sequence $\{\mathbf{x}_k\}_{k \geq 0}$ to an optimal solution.

**Theorem 9.18 (convergence of the sequence generated by the gradient projection method).** *Under the setting of Theorem* 9.16 *the sequence* $\{\mathbf{x}_k\}_{k \geq 0}$ *generated by the gradient projection method with a constant stepsize* $t_k = \bar{t} \in (0, \frac{1}{L}]$ *converges to an optimal solution.*

*Proof.* The result follows from the Fejér monotonicity of the sequence. First of all, by part (b) of Theorem 9.16 it follows that $f(\mathbf{x}_k) \to f^*$, where $f^*$ is the optimal value of problem (P), and thus by the continuity of $f$, any limit point of the sequence $\{\mathbf{x}_k\}_{k \geq 0}$ is an optimal solution. In addition, by Fejér monotonicity property, we have for any $\mathbf{x}^* \in X^*$ that $\|\mathbf{x}_k - \mathbf{x}^*\| \leq \|\mathbf{x}_0 - \mathbf{x}^*\|$, implying that the sequence $\{\mathbf{x}_k\}_{k \geq 0}$ is bounded, establishing the fact that the sequence does have at least one limit point. To prove the convergence of the entire sequence, let $\tilde{\mathbf{x}}$ be a limit point of the sequence $\{\mathbf{x}_k\}_{k \geq 0}$, meaning that there exists a subsequence $\{\mathbf{x}_{k_j}\}_{j \geq 0}$ such that $\mathbf{x}_{k_j} \to \tilde{\mathbf{x}}$. Since $\tilde{\mathbf{x}} \in X^*$, it follows by Fejér monotonicity that for any $k \geq 0$

$$\|\mathbf{x}_{k+1} - \tilde{\mathbf{x}}\| \leq \|\mathbf{x}_k - \tilde{\mathbf{x}}\|.$$

Thus, $\{\|\mathbf{x}_k - \tilde{\mathbf{x}}\|\}_{k \geq 0}$ is a nonincreasing sequence which is bounded below (by zero) and hence convergent. Since $\|\mathbf{x}_{k_j} - \tilde{\mathbf{x}}\| \to 0$ as $j \to \infty$, it follows that the whole sequence $\{\|\mathbf{x}_k - \tilde{\mathbf{x}}\|\}_{k \geq 0}$ converges to zero, and consequently $\mathbf{x}_k \to \tilde{\mathbf{x}}$ as $k \to \infty$. $\quad\square$

## 9.5 ▪ Sparsity Constrained Problems

So far we considered problems in which the feasible set $C$ is convex. In this section we will study a class of problems in which the constraint set is nonconvex. Therefore, none of the results derived so far apply for this class of problems, and we will see that beside several similarities between the convex and nonconvex cases, there are substantial differences in the type of results that can be obtained.

### 9.5.1 ▪ Definition and Notation

In this section we consider the problem of minimizing a continuously differentiable objective function subject to a sparsity constraint. More specifically, we consider the problem

$$\text{(S)} \quad \begin{array}{ll} \min & f(\mathbf{x}) \\ \text{s.t.} & \|\mathbf{x}\|_0 \leq s, \end{array}$$

where $f : \mathbb{R}^n \to \mathbb{R}$ is bounded below and a continuously differentiable function whose gradient is Lipschitz with constant $L_f$, $s > 0$ is an integer smaller than $n$, and $\|\mathbf{x}\|_0$ is the so-called $\ell_0$ norm of $\mathbf{x}$, which counts the number of nonzero components in $\mathbf{x}$:

$$\|\mathbf{x}\|_0 = |\{i : x_i \neq 0\}|.$$

It is important to note that the $l_0$ norm is actually *not* a norm. Indeed, it does not even satisfy the homogeneity property since $\|\lambda \mathbf{x}\|_0 = \|\mathbf{x}\|_0$ for $\lambda \neq 0$. We do not assume that $f$ is a convex function. The constraint set is of course nonconvex and most of the analysis of this chapter does not hold for problem (S). Nevertheless, this type of problem is important and has many applications in areas such as compressed sensing, and it is therefore worthwhile to study it and understand how the "classical" results over convex sets can be adjusted in order to be valid for problem (S). First of all, we begin with some notation. For a given vector $\mathbf{x} \in \mathbb{R}^n$, the support set of $\mathbf{x}$ is defined by

$$I_1(\mathbf{x}) \equiv \{i : x_i \neq 0\},$$

and its complement by

$$I_0(\mathbf{x}) \equiv \{i : x_i = 0\}.$$

We denote by $C_s$ the set of vectors $\mathbf{x}$ that are at most $s$-sparse, meaning that they have at most $s$ nonzero elements:

$$C_s = \{\mathbf{x} : \|\mathbf{x}\|_0 \leq s\}.$$

In this notation problem (S) can be written as

$$\min\{f(\mathbf{x}) : \mathbf{x} \in C_s\}.$$

For a vector $\mathbf{x} \in \mathbb{R}^n$ and $i \in \{1, 2, \ldots, n\}$, the $i$th largest absolute value component in $\mathbf{x}$ is denoted by $M_i(\mathbf{x})$, so that in particular

$$M_1(\mathbf{x}) \geq M_2(\mathbf{x}) \geq \cdots \geq M_n(\mathbf{x}).$$

Also, $M_1(\mathbf{x}) = \max_{i=1,\ldots,n} |x_i|$ and $M_n(\mathbf{x}) = \min_{i=1,\ldots,n} |x_i|$.

Our objective now is to study necessary optimality conditions for problem (S), which are similar in nature to the stationarity property. As we will see, the nonconvexity of the constraint in (S) prohibits the usual stationarity conditions to hold, but we will see that other optimality conditions can be constructed.

## 9.5.2 ▪ $L$-Stationarity

We begin by extending the concept of stationarity to the case of problem (S). We will use the characterization of stationarity in terms of the orthogonal projection. Recall that $\mathbf{x}^*$ is a stationary point of the problem of minimizing a continuously differentiable function over a closed and convex set $C$ if and only if

$$\mathbf{x}^* = P_C(\mathbf{x}^* - s\nabla f(\mathbf{x}^*)), \tag{9.30}$$

where $s$ is an arbitrary positive number. It is interesting to note (see Theorem 9.10) that condition (9.30)—although expressed in terms of the parameter $s$—does not actually depend on $s$. Relation (9.30) is no longer a valid condition when $C = C_s$ since the

orthogonal projection is not unique on nonconvex sets and is in fact a multivalued mapping, meaning that its values form a set given by

$$P_{C_s}(\mathbf{x}) = \operatorname{argmin}_{\mathbf{y}}\{\|\mathbf{y} - \mathbf{x}\| : \mathbf{y} \in C_s\}.$$

The existence of vectors in $P_{C_s}(\mathbf{x})$ follows by the closedness of $C_s$ and the coercivity of the function $\|\mathbf{y} - \mathbf{x}\|$ (with respect to $\mathbf{y}$). To extend (9.30) to the sparsity constrained problem (S), we introduce the notion of "$L$-stationarity."

**Definition 9.19.** *A vector $\mathbf{x}^* \in C_s$ is called an $L$-stationary point of* (S) *if it satisfies the relation*

$$[\text{NC}_L] \quad \mathbf{x}^* \in P_{C_s}\left(\mathbf{x}^* - \frac{1}{L}\nabla f(\mathbf{x}^*)\right). \tag{9.31}$$

As already mentioned, the orthogonal projection operator $P_{C_s}(\cdot)$ is not single-valued. In this case, the orthogonal projection $P_{C_s}(\mathbf{x})$ comprises all vectors consisting of the $s$ components of $\mathbf{x}$ with the largest absolute values and with zeros elsewhere. In general, there can be more than one choice to the $s$ largest components in absolute value, and each of these choices gives rise to another vector in the set $P_{C_s}(\mathbf{x})$. For example

$$P_{C_2}((2,1,1)^T) = \left\{(2,1,0)^T, (2,0,1)^T\right\}.$$

Below we will show that under our assumption that $f \in C_{L_f}^{1,1}(\mathbb{R}^n)$, $L$-stationarity is a necessary condition for optimality whenever $L > L_f$.

Before proving this result, we describe a more explicit representation of $[\text{NC}_L]$.

**Lemma 9.20.** *For any $L > 0$, $\mathbf{x}^* \in C_s$ satisfies $[\text{NC}_L]$ if and only if*

$$\left|\frac{\partial f}{\partial x_i}(\mathbf{x}^*)\right| \begin{cases} \leq L M_s(\mathbf{x}^*) & \text{if } i \in I_0(\mathbf{x}^*), \\ = 0 & \text{if } i \in I_1(\mathbf{x}^*). \end{cases} \tag{9.32}$$

**Proof.** $[\text{NC}_L] \Rightarrow$ (9.32). Suppose that $\mathbf{x}^* \in C_s$ satisfies $[\text{NC}_L]$. Note that for any index $j \in \{1, 2, \ldots, n\}$, the $j$th component of any vector in $P_{C_s}(\mathbf{x}^* - \frac{1}{L}\nabla f(\mathbf{x}^*))$ is either zero or equal to $x_j^* - \frac{1}{L}\frac{\partial f}{\partial x_j}(\mathbf{x}^*)$. Now, since $\mathbf{x}^* \in P_{C_s}(\mathbf{x}^* - \frac{1}{L}\nabla f(\mathbf{x}^*))$, it follows that if $i \in I_1(\mathbf{x}^*)$, then $x_i^* = x_i^* - \frac{1}{L}\frac{\partial f}{\partial x_i}(\mathbf{x}^*)$, so that $\frac{\partial f}{\partial x_i}(\mathbf{x}^*) = 0$. If $i \in I_0(\mathbf{x}^*)$, then $|x_i^* - \frac{1}{L}\frac{\partial f}{\partial x_i}(\mathbf{x}^*)| \leq M_s(\mathbf{x}^*)$, which combined with the fact that $x_i^* = 0$ implies that $|\frac{\partial f}{\partial x_i}(\mathbf{x}^*)| \leq L M_s(\mathbf{x}^*)$, and consequently (9.32) holds true.

(9.32) $\Rightarrow [\text{NC}_L]$. Suppose that $\mathbf{x}^* \in C_s$ satisfies (9.32). If $\|\mathbf{x}^*\|_0 < s$, then $M_s(\mathbf{x}^*) = 0$, and by (9.32) it follows that $\nabla f(\mathbf{x}^*) = \mathbf{0}$; therefore, in this case, $P_{C_s}\left(\mathbf{x}^* - \frac{1}{L}\nabla f(\mathbf{x}^*)\right) = P_{C_s}(\mathbf{x}^*)$ is the set $\{\mathbf{x}^*\}$. If $\|\mathbf{x}^*\|_0 = s$, then $M_s(\mathbf{x}^*) \neq 0$ and $|I_1(\mathbf{x}^*)| = s$. By (9.32)

$$\left|x_i^* - \frac{1}{L}\frac{\partial f}{\partial x_i}(\mathbf{x}^*)\right| \begin{cases} = |x_i^*|, & i \in I_1(\mathbf{x}^*), \\ \leq M_s(\mathbf{x}^*), & i \in I_0(\mathbf{x}^*). \end{cases}$$

Therefore, the vector $\mathbf{x}^* - \frac{1}{L}\nabla f(\mathbf{x}^*)$ contains the $s$ components of $\mathbf{x}^*$ with the largest absolute value and all other components are smaller or equal (in absolute value) to them, and consequently $[\text{NC}_L]$ holds. $\square$

Evidently, the condition (9.32) depends on the constant $L$; the condition is more restrictive for smaller values of $L$. This shows that the state of affairs in the nonconvex case is different. Before proceeding to show under which conditions is $L$-stationarity a necessary optimality condition, we will prove the following technical and useful result. The analysis depends heavily on the descent lemma (Lemma 4.22). We recall that the lemma states that if $f \in C_{L_f}^{1,1}(\mathbb{R}^n)$ and $L \geq L_f$, then

$$f(\mathbf{y}) \leq h_L(\mathbf{y}, \mathbf{x}),$$

where

$$h_L(\mathbf{y}, \mathbf{x}) \equiv f(\mathbf{x}) + \langle \nabla f(\mathbf{x}), \mathbf{y} - \mathbf{x} \rangle + \frac{L}{2} \|\mathbf{x} - \mathbf{y}\|^2. \tag{9.33}$$

**Lemma 9.21.** *Suppose that $f \in C_{L_f}^{1,1}(\mathbb{R}^n)$ and that $L > L_f$. Then for any $\mathbf{x} \in C_s$ and $\mathbf{y} \in \mathbb{R}^n$ satisfying*

$$\mathbf{y} \in P_{C_s}\left(\mathbf{x} - \frac{1}{L}\nabla f(\mathbf{x})\right), \tag{9.34}$$

*we have*

$$f(\mathbf{x}) - f(\mathbf{y}) \geq \frac{L - L_f}{2}\|\mathbf{x} - \mathbf{y}\|^2. \tag{9.35}$$

**Proof.** Note that (9.34) can be written as

$$\mathbf{y} \in \operatorname{argmin}_{\mathbf{z} \in C_s} \left\| \mathbf{z} - \left(\mathbf{x} - \frac{1}{L}\nabla f(\mathbf{x})\right)\right\|^2. \tag{9.36}$$

Since

$$\begin{aligned} h_L(\mathbf{z}, \mathbf{x}) &= f(\mathbf{x}) + \langle \nabla f(\mathbf{x}), \mathbf{z} - \mathbf{x} \rangle + \frac{L}{2}\|\mathbf{z} - \mathbf{x}\|^2 \\ &= \frac{L}{2}\left\|\mathbf{z} - \left(\mathbf{x} - \frac{1}{L}\nabla f(\mathbf{x})\right)\right\|^2 + \underbrace{f(\mathbf{x}) - \frac{1}{2L}\|\nabla f(\mathbf{x})\|^2}_{\text{constant w.r.t. } \mathbf{z}}, \end{aligned}$$

it follows that the minimization problem (9.36) is equivalent to

$$\mathbf{y} \in \operatorname{argmin}_{\mathbf{z} \in C_s} h_L(\mathbf{z}, \mathbf{x}).$$

This implies that

$$h_L(\mathbf{y}, \mathbf{x}) \leq h_L(\mathbf{x}, \mathbf{x}) = f(\mathbf{x}). \tag{9.37}$$

Now by the descent lemma we have

$$f(\mathbf{x}) - f(\mathbf{y}) \geq f(\mathbf{x}) - h_{L_f}(\mathbf{y}, \mathbf{x}),$$

which combined with (9.37) and the identity

$$h_{L_f}(\mathbf{x}, \mathbf{y}) = h_L(\mathbf{x}, \mathbf{y}) - \frac{L - L_f}{2}\|\mathbf{x} - \mathbf{y}\|^2$$

yields (9.35).  □

We are now ready to prove the necessity of the $L$-stationarity property under the condition $L > L_f$.

**Theorem 9.22 ($L$-stationarity as a necessary optimality condition).** *Suppose that $f \in C_{L_f}^{1,1}(\mathbb{R}^n)$ and that $L > L_f$. Let $\mathbf{x}^*$ be an optimal solution of* (S). *Then*

(i) $\mathbf{x}^*$ *is an $L$-stationary point,*

(ii) *the set $P_{C_s}(\mathbf{x}^* - \frac{1}{L}\nabla f(\mathbf{x}^*))$ is a singleton.*[3]

**Proof.** We will prove both parts simultaneously. Suppose to the contrary that there exists a vector

$$\mathbf{y} \in P_{C_s}\left(\mathbf{x}^* - \frac{1}{L}\nabla f(\mathbf{x}^*)\right), \tag{9.38}$$

which is different from $\mathbf{x}^*$ ($\mathbf{y} \neq \mathbf{x}^*$). Invoking Lemma 9.21 with $\mathbf{x} = \mathbf{x}^*$, we have

$$f(\mathbf{x}^*) - f(\mathbf{y}) \geq \frac{L - L_f}{2}\|\mathbf{x}^* - \mathbf{y}\|^2,$$

contradicting the optimality of $\mathbf{x}^*$. We conclude that $\mathbf{x}^*$ is the only vector in the set $P_{C_s}(\mathbf{x}^* - \frac{1}{L}\nabla f(\mathbf{x}^*))$. $\quad\square$

To summarize, we have shown that under a Lipschitz condition on $\nabla f$, $L$-stationarity for any $L > L_f$ is a necessary optimality condition.

### 9.5.3 ▪ The Iterative Hard-Thresholding Method

One approach for solving problem (S) is to employ the following natural generalization of the gradient projection algorithm. The method can be seen as a fixed point method aimed at "enforcing" the $L$-stationary condition (9.31):

$$\mathbf{x}^{k+1} \in P_{C_s}\left(\mathbf{x}^k - \frac{1}{L}\nabla f(\mathbf{x}^k)\right), \quad k = 0, 1, 2, \ldots. \tag{9.39}$$

The method is known in the literature as the *iterative hard-thresholding (IHT)* method, and we will adopt this terminology.

> **The IHT Method**
> **Input:** a constant $L > L_f$.
>
> • **Initialization:** Choose $\mathbf{x}_0 \in C_s$.
> • **General step :** $\mathbf{x}^{k+1} \in P_{C_s}\left(\mathbf{x}^k - \frac{1}{L}\nabla f(\mathbf{x}^k)\right), \quad k = 0, 1, 2, \ldots.$

It can be shown that the general step of the IHT method is equivalent to the relation

$$\mathbf{x}^{k+1} \in \operatorname{argmin}_{\mathbf{x} \in C_s} h_L(\mathbf{x}, \mathbf{x}^k), \tag{9.40}$$

where $h_L(\mathbf{x}, \mathbf{y})$ is defined by (9.33). (See also the proof of Lemma 9.21.)

---

[3] A set is called a *singleton* if it contains exactly one element.

Several basic properties of the IHT method are summarized in the following lemma.

**Lemma 9.23.** *Suppose that $f \in C_{L_f}^{1,1}(\mathbb{R}^n)$ and that $f$ is lower bounded. Let $\{\mathbf{x}^k\}_{k \geq 0}$ be the sequence generated by the IHT method with a constant stepsize $\frac{1}{L}$, where $L > L_f$. Then*

(a) $f(\mathbf{x}^k) - f(\mathbf{x}^{k+1}) \geq \frac{L - L_f}{2} \|\mathbf{x}^k - \mathbf{x}^{k+1}\|^2$,

(b) $\{f(\mathbf{x}^k)\}_{k \geq 0}$ *is a nonincreasing sequence,*

(c) $\|\mathbf{x}^k - \mathbf{x}^{k+1}\| \to 0$,

(d) *for every $k = 0, 1, 2, \ldots$, if $\mathbf{x}^k \neq \mathbf{x}^{k+1}$, then $f(\mathbf{x}^{k+1}) < f(\mathbf{x}^k)$.*

**Proof.** Part (a) follows by substituting $\mathbf{x} = \mathbf{x}^k, \mathbf{y} = \mathbf{x}^{k+1}$ into (9.35). Part (b) follows immediately from part (a). To prove (c), note that since $\{f(\mathbf{x}_k)\}_{k \geq 0}$ is a nonincreasing sequence, which is also lower bounded (by the fact that $f$ is bounded below), it follows that it converges to some limit $\ell$ and hence $f(\mathbf{x}_k) - f(\mathbf{x}_{k+1}) \to \ell - \ell = 0$ as $k \to \infty$. Therefore, by part (a) and the fact that $L > L_f$, the limit $\|\mathbf{x}^k - \mathbf{x}^{k+1}\| \to 0$ holds. Finally, (d) is a direct consequence of (a).  $\square$

As already mentioned, the IHT algorithm can be viewed as a fixed point method for solving the condition for $L$-stationarity. The following theorem states that all accumulation points of the sequence generated by the IHT method with constant stepsize $\frac{1}{L}$ are indeed $L$-stationary points.

**Theorem 9.24.** *Suppose that $f \in C_{L_f}^{1,1}(\mathbb{R}^n)$ and that $f$ is lower bounded. Let $\{\mathbf{x}^k\}_{k \geq 0}$ be the sequence generated by the IHT method with stepsize $\frac{1}{L}$, where $L > L_f$. Then any accumulation point of $\{\mathbf{x}^k\}_{k \geq 0}$ is an $L$-stationary point.*

**Proof.** Suppose that $\mathbf{x}^*$ is an accumulation point of the sequence. Then there exists a subsequence $\{\mathbf{x}^{k_n}\}_{n \geq 0}$ that converges to $\mathbf{x}^*$. By Lemma 9.23

$$f(\mathbf{x}^{k_n}) - f(\mathbf{x}^{k_n+1}) \geq \frac{L - L_f}{2} \|\mathbf{x}^{k_n} - \mathbf{x}^{k_n+1}\|^2. \tag{9.41}$$

Since $\{f(\mathbf{x}^{k_n})\}_{n \geq 0}$ and $\{f(\mathbf{x}^{k_n+1})\}_{n \geq 0}$, as nonincreasing and lower bounded sequences, both converge to the same limit, it follows that $f(\mathbf{x}^{k_n}) - f(\mathbf{x}^{k_n+1}) \to 0$ as $n \to \infty$, which combined with (9.41) yields that $\mathbf{x}^{k_n+1} \to \mathbf{x}^*$ as $n \to \infty$. Recall that for all $n \geq 0$

$$\mathbf{x}^{k_n+1} \in P_{C_s}\left(\mathbf{x}^{k_n} - \frac{1}{L}\nabla f(\mathbf{x}^{k_n})\right).$$

Let $i \in I_1(\mathbf{x}^*)$. By the convergence of $\mathbf{x}^{k_n}$ and $\mathbf{x}^{k_n+1}$ to $\mathbf{x}^*$, it follows that there exists an $N$ such that

$$x_i^{k_n}, x_i^{k_n+1} \neq 0 \text{ for all } n > N,$$

and therefore, for $n > N$,

$$x_i^{k_n+1} = x_i^{k_n} - \frac{\partial f}{\partial x_i}(\mathbf{x}^{k_n}).$$

Taking $n$ to $\infty$ and using the continuity of $f$, we obtain that

$$\frac{\partial f}{\partial x_i}(\mathbf{x}^*) = 0.$$

Now let $i \in I_0(\mathbf{x}^*)$. If there exist an infinite number of indices $k_n$ for which $x_i^{k_n+1} \neq 0$, then as in the previous case, we obtain that $x_i^{k_n+1} = x_i^{k_n} - \frac{\partial f}{\partial x_i}(\mathbf{x}^{k_n})$ for these indices, implying (by taking the limit) that $\frac{\partial f}{\partial x_i}(\mathbf{x}^*) = 0$. In particular, $|\frac{\partial f}{\partial x_i}(\mathbf{x}^*)| \leq L M_s(\mathbf{x}^*)$. On the other hand, if there exists an $M > 0$ such that for all $n > M$ $x_i^{k_n+1} = 0$, then

$$\left| x_i^{k_n} - \frac{1}{L}\frac{\partial f}{\partial x_i}(\mathbf{x}^{k_n}) \right| \leq M_s\left( \mathbf{x}^{k_n} - \frac{1}{L}\nabla f(\mathbf{x}^{k_n}) \right) = M_s(\mathbf{x}^{k_n+1}).$$

Thus, taking $n$ to infinity, while exploiting the continuity of the function $M_s$, we obtain that

$$\left| \frac{\partial f}{\partial x_i}(\mathbf{x}^*) \right| \leq L M_s(\mathbf{x}^*),$$

and hence, by Lemma 9.20, the desired result is established.   $\square$

## Exercises

9.1. Let $f$ be a continuously differentiable convex function over a closed and convex set $C \subseteq \mathbb{R}^n$. Show that $\mathbf{x}^* \in C$ is an optimal solution of the problem

$$\text{(P)} \quad \min\{f(\mathbf{x}) : \mathbf{x} \in C\}$$

if and only if

$$\langle \nabla f(\mathbf{x}), \mathbf{x}^* - \mathbf{x} \rangle \leq 0 \text{ for all } \mathbf{x} \in C.$$

9.2. Consider the Huber function

$$H_\mu(\mathbf{x}) = \begin{cases} \frac{\|\mathbf{x}\|^2}{2\mu}, & \|\mathbf{x}\| \leq \mu, \\ \|\mathbf{x}\| - \frac{\mu}{2} & \text{else}, \end{cases}$$

where $\mu > 0$ is a given parameter. Show that $H_\mu \in C_{\frac{1}{\mu}}^{1,1}$.

9.3. Consider the minimization problem

$$\text{(Q)} \quad \begin{array}{ll} \min & 2x_1^2 + 3x_2^2 + 4x_3^2 + 2x_1x_2 - 2x_1x_3 - 8x_1 - 4x_2 - 2x_3 \\ \text{s.t.} & x_1, x_2, x_3 \geq 0. \end{array}$$

(i) Show that the vector $(\frac{17}{7}, 0, \frac{6}{7})^T$ is an optimal solution of (Q).

(ii) Employ the gradient projection method via MATLAB with constant stepsize $\frac{1}{L}$ ($L$ being the Lipschitz constant of the gradient of the objective function). Show the function values of the first 100 iterations and the produced solution.

9.4. Consider the minimization problem

$$(P) \quad \min\{f(\mathbf{x}) : \mathbf{a}^T \mathbf{x} = 1, \mathbf{x} \in \mathbb{R}^n\},$$

where $f$ is a continuously differentiable function over $\mathbb{R}^n$ and $\mathbf{a} \in \mathbb{R}^n_{++}$. Show that $\mathbf{x}^*$ satisfying $\mathbf{a}^T \mathbf{x}^* = 1$ is a stationary point of (P) if and only if

$$\frac{\frac{\partial f}{\partial x_1}(\mathbf{x}^*)}{a_1} = \frac{\frac{\partial f}{\partial x_2}(\mathbf{x}^*)}{a_2} = \cdots = \frac{\frac{\partial f}{\partial x_n}(\mathbf{x}^*)}{a_n}.$$

9.5. Consider the minimization problem

$$(P) \quad \min\{f(\mathbf{x}) : \mathbf{x} \in \Delta_n\},$$

where $f$ is a continuously differentiable function over $\Delta_n$. Prove that $\mathbf{x}^* \in \Delta_n$ is a stationary point of (P) if and only if there exists $\mu \in \mathbb{R}$ such that

$$\frac{\partial f}{\partial x_i}(\mathbf{x}^*) \begin{cases} = \mu, & x_i^* > 0, \\ \geq \mu, & x_i^* = 0. \end{cases}$$

9.6. Let $S \subseteq \mathbb{R}^n$ be a nonempty closed and convex set, and let $f \in C_L^{1,1}(S)$ be a convex function over $S$. Assume that the optimal value of the problem

$$(P) \quad \min\{f(\mathbf{x}) : \mathbf{x} \in S\},$$

denoted by $f^*$ is real. Prove that for any $\mathbf{x} \in S$ the following inequality holds:

$$f(\mathbf{x}) - f^* \geq \frac{1}{2L} \|G_L(\mathbf{x})\|^2,$$

where $G_L(\mathbf{x}) = L[\mathbf{x} - P_S(\mathbf{x} - \frac{1}{L} \nabla f(\mathbf{x}))]$.

9.7. Let $f \in C_L^{1,1}(S)$ be a strongly convex function over a nonempty closed convex set $S \subseteq \mathbb{R}^n$ with strong convexity parameter $\sigma > 0$, that is,

$$f(\mathbf{y}) \geq f(\mathbf{x}) + \langle \nabla f(\mathbf{x}), \mathbf{y} - \mathbf{x} \rangle + \frac{\sigma}{2} \|\mathbf{x} - \mathbf{y}\|^2 \text{ for all } \mathbf{x}, \mathbf{y} \in S.$$

Consider the problem

$$(P) \quad \begin{array}{ll} \min & f(\mathbf{x}) \\ \text{s.t.} & \mathbf{x} \in S, \end{array}$$

where $S$ is a closed and convex set. Let $\{\mathbf{x}_k\}_{k \geq 0}$ be the sequence generated by the gradient projection method for solving problem (P) with a constant stepsize $t_k = \frac{1}{L}$. Let $\mathbf{x}^*$ be the optimal solution of (P), and let $f^*$ be the optimal value of (P). Prove that there exists a constant $c \in (0, 1)$ such that for any $k \geq 0$

$$\|\mathbf{x}_{k+1} - \mathbf{x}^*\| \leq c \|\mathbf{x}_k - \mathbf{x}^*\|.$$

Find an explicit expression for $c$.