# Continuous Optimization

## Chapter 3: Constrained Optimization

## 1  Definitions

In this chapter we will consider constrained optimization problems with the following shape

$$\begin{aligned} \min \quad & f(x) \\ \text{s.t.} \quad & x \in C \end{aligned} \tag{1}$$

**Definition 1.1** (Convex Set)**.** *A set $C$ is said to be convex if given $x_1, x_2 \in C$ and $\lambda \in [0,1]$, then $\lambda x_1 + (1-\lambda)x_2 \in C$.*

**Definition 1.2** (Convex Function)**.** *A function $f : C \to \mathbb{R}$ defined on a convex set $C$ is said to be convex if given $x_1, x_2 \in C$ and $\lambda \in [0,1]$, then*

$$f(\lambda x_1 + (1-\lambda)x_2) \le \lambda f(x_1) + (1-\lambda)f(x_2).$$

**Definition 1.3** (Strictly Convex Function)**.** *A function $f : C \to \mathbb{R}$ defined on a convex set $C$ is said to be strictly convex if given $x_1, x_2 \in C$ and $\lambda \in [0,1]$, then*

$$f(\lambda x_1 + (1-\lambda)x_2) < \lambda f(x_1) + (1-\lambda)f(x_2).$$

A function is called concave if $-f$ is convex and strictly concave if $-f$ is strictly convex.
Now, given $\Delta_k$ the unit-simplex, that is the subset of $\mathbb{R}^k$ comprising all nonnegative vectors whose sum is 1, i.e.,

$$\{\lambda \in \mathbb{R}^k : \lambda \ge 0, e^t \lambda = 1\},$$

we can provide the following very useful result by Jensen's.

**Theorem 1.1** (Jensen's Inequality)**.** *Let $f : C \to \mathbb{R}$ be a convex function over a convex set $C$. Then for any $x_1, x_2, \ldots, x_k \in C$ and $\lambda \in \Delta_k$ we have*

$$f\left(\sum_{i=1}^{k} \lambda_i x_i\right) \le \sum_{i=1}^{k} \lambda_i f(x_i). \tag{2}$$

*Proof.* We will prove (2) by induction on $k$. For $k = 1$ the result is obvious ($f(x_1) \le f(x_1) \ \forall x_1 \in C$). We now assume that (2) holds for $k$ and we will prove that it also holds for $k+1$. Suppose we have $x_1, x_2, \ldots, x_{k+1} \in C$ and $\lambda \in \Delta_{k+1}$, we will show that $f(z) \le \sum_{i=1}^{k+1} \lambda_i f(x_i)$ with $z = \sum_{i=1}^{k+1} \lambda_i x_i$. If $\lambda_{k+1} = 1$, then $z = x_{k+1}$ and (2) is obvious. If $\lambda_{k+1} < 1$, then

$$\begin{aligned} f(z) &= f\left(\sum_{i=1}^{k} \lambda_i x_i + \lambda_{k+1} x_{k+1}\right) \\ &= f\left((1 - \lambda_{k+1})\sum_{i=1}^{k} \frac{\lambda_i}{1 - \lambda_{k+1}} x_i + \lambda_{k+1} x_{k+1}\right) \\ &\le (1 - \lambda_{k+1})f(v) + \lambda_{k+1} f(x_{k+1}), \end{aligned}$$

with $v = \sum_{i=1}^{k} \frac{\lambda_i}{1 - \lambda_{k+1}} x_i$. Since $\sum_{i=1}^{k} \frac{\lambda_i}{1 - \lambda_{k+1}} = \frac{1 - \lambda_{k+1}}{1 - \lambda_{k+1}} = 1$, it follows that $v$ is a convex combination of $k$ points from $C$, hence by the induction hypotesis we have that $f(v) \le \sum_{i=1}^{k} \frac{\lambda_i}{1 - \lambda_{k+1}} f(x_i)$, which combined with the equality above yields

$$f(z) \le \sum_{i=1}^{k+1} \lambda_i f(x_i).$$

$\square$

# 2 Characterizations of Convex Functions

**Theorem 2.1** (Gradient characterization of convex functions). *Let $f \in \mathrm{C}^1(C)$, where $C$ is convex. Then $f$ is convex over $C$ if and only if*

$$f(x) + \nabla f(x)^T (y - x) \leq f(y) \quad \forall x, y \in C. \tag{3}$$

*Proof.* Exercise. $\qquad\square$

**Proposition 2.1** (Sufficiency of stationarity under convexity). *Let $f \in \mathrm{C}^1(C)$, where $C \subseteq \mathbb{R}^n$ is convex. Suppose that $\nabla f(x^*) = 0$ for some $x^* \in C$. Then $x^*$ is a global minimizer of $f$ over $C$.*

*Proof.* Let $z \in C$. Plugging $x = x^*$ and $y = z$ in Theorem 2.1 we obtain that

$$f(z) \geq f(x^*) + \nabla f(x^*)^T (z - x^*),$$

which implies that $f(z) \geq f(x^*)$ because $\nabla f(x^*) = 0$. $\qquad\square$

We note that Proposition 2.1 establishes only the sufficiency of the stationarity condition $\nabla f(x^*) = 0$ for guaranteeing that $x^*$ is a global optimal solution. When $C$ is not the entire space, this condition is not necessary, in fact it might be that the points for which $\nabla f(\cdot) = 0$ are not in $C$. On the other hand, when $C = \mathbb{R}^n$ and $f$ is convex, $\nabla f(x^*) = 0$ is both sufficient and necessary condition for $x^*$ to be a global minimum. We can now establish the conditions under which a twice continuously differentiable function $f$ is convex.

**Theorem 2.2** (Second order characterization of convexity). *Let $f \in \mathrm{C}^2(C)$, where $C \subseteq \mathbb{R}^n$ is convex and open. Thus, we have that $f$ is convex iff $\nabla^2 f(x) \succcurlyeq 0 \quad \forall x \in C$.*

*Proof.* Suppose that $\nabla^2 f(x) \succcurlyeq 0$ for all $x \in C$. We will prove (3) which is enough to establish convexity. Let $x, y \in C$, then by the Mean Value Theorem[2] (Theorem 2.6 from Chapter 1) we get that there exists $z \in [x, y]$ (and hence $z \in C$) for which

$$f(y) = f(x) + \nabla f(x)^T (y - x) + \frac{1}{2}(y - x)^T \nabla^2 f(z)(y - x). \tag{4}$$

Since $\nabla^2 f(z) \succcurlyeq 0$, it follows that $(y - x)^T \nabla^2 f(z)(y - x) \geq 0$, which implies (3). To prove the opposite direction, assume that $f$ is convex over $C$. Let $x \in C$ and $y \in \mathbb{R}^n$. Since $C$ is open, it follows that $x + \lambda y \in C$, for $0 < \lambda < \epsilon$, where $\epsilon$ is a small enough positive constant. Using now the gradient characterization of convex functions (3) we get

$$f(x + \lambda y) \geq f(x) + \lambda \nabla f(x)^T y.$$

In addition, by the quadratic approximation theorem (Theorem 2.4 from Chapter 1), we have that

$$f(x + \lambda y) = f(x) + \lambda \nabla f(x)^T y + \frac{\lambda^2}{2} y^T \nabla^2 f(x) y + o(\lambda^2 \|y\|^2),$$

which combined with the above inequality gives

$$\frac{\lambda^2}{2} y^T \nabla^2 f(x) y + o(\lambda^2 \|y\|^2) \geq 0 \quad \forall \lambda \in (0, \epsilon).$$

Dividing the latter inequality by $\lambda^2$ and taking the limit for $\lambda \to 0^+$, we have

$$y^T \nabla^2 f(x) y \geq 0 \quad \forall y \in \mathbb{R}^n,$$

which concludes the proof. $\qquad\square$

The same theorem works with positive definiteness and strict convexity, meaning also that the minimum in this case is unique.

# 3    Optimization over convex constraints

From now on, we consider (1) where $C$ is convex. On the other hand, we will not always assume also $f$ to be convex. From the convexity of $f$ we have the following two theorems. Notice that the following result is not a direct consequence of Proposition 2.1 as the local (and global) minimum, might be on the boundary of the set and not be stationary (in the sense of unconstrained optimization).

**Theorem 3.1** (global=local in convex optimization). *Let $f : C \to \mathbb{R}$ be a convex function over a convex set $C \subseteq \mathbb{R}^n$. Let $x^* \in C$ be a local minimum of $f$ over $C$. Then $x^*$ is a global minimum of $f$ over $C$.*

*Proof.* Since $x^*$ is a local minimum of $f$ over $C$ there exists $r$ such that $f(x) \geq f(x^*)$ for any $x \in C \cap B[x^*, r]$. Now let $y \in C$ with $y \neq x^*$. We want to show that $f(y) \geq f(x^*)$. Let $\lambda \in (0, 1]$ be such that $x^* + \lambda(y - x^*) \in B[x^*, r]$, for instance $\lambda = \frac{r}{||y-x^*||}$. Now, since $x^* + \lambda(y - x^*) \in B[x^*, r] \cap C$, it follows that $f(x^*) \leq f(x^* + \lambda(y - x^*))$, and hence, by convexity of $f$, also

$$f(x^*) \leq f(x^* + \lambda(y - x^*)) \leq (1 - \lambda)f(x^*) + \lambda f(y)$$

Thus, $\lambda f(x^*) \leq \lambda f(y)$, which concludes the proof. $\qquad\square$

**Theorem 3.2** (Convexity of the optimal set in convex optimization). *Let $f : C \to \mathbb{R}$ be a convex function with $C \subseteq \mathbb{R}^n$ convex. Then, the set of optimal solutions of the problem (1), which we denote by $X^*$ is convex. Moreover, if $f$ is strictly convex over $C$, then there exists at most one optimal solution.*

*Proof.* If $X^* = \emptyset$, the result follows trivially. Suppose that $X \neq \emptyset$ and denote the optimal value of $f$ by $f^*$. Let $x, y \in C$ with $\lambda \in [0, 1]$. Then, by convexity $f(\lambda x + (1 - \lambda)y) \leq \lambda f^* + (1 - \lambda)f^* = f^*$, hence $\lambda x + (1 - \lambda)y$ is also optimal, i.e., it belongs to $X^*$, establishing the convexity of $X^*$. Suppose now that $f$ is strictly convex and $X^*$ is nonempty, and suppose by contradiction that there are 2 points $x, y$ in $X^*$. Then $\lambda x + (1 - \lambda)y \in C$, and by the strict convexity of $f$ we have

$$f(\lambda x + (1 - \lambda)y) < \lambda f(x) + (1 - \lambda)f(y) = f^*,$$

which is a contradiction to the fact that $f^*$ is the optimal value. $\qquad\square$

## 3.1    Stationarity

Note that the following definition and the following theorem are given also for the more general case in which $f$ is not convex.

**Definition 3.1** (Stationary points of convex constrained problems). *Let $f \in \mathrm{C}^1(C)$, where $C$ is closed and convex. Then $x^*$ is a stationary point of (1) if*

$$\nabla f(x^*)^T(x - x^*) \geq 0 \ \forall x \in C. \tag{5}$$

In words, this means that there are no feasible descent directions of $f$ at $x^*$. This suggests that stationarity is in fact a necessary condition for a local minimum of (1).

**Theorem 3.3** (Stationarity as necessary optimality condition of a convex constrained problem). *Let $f \in \mathrm{C}^1(C)$, where $C$ is closed and convex and let $x^*$ be a local minimum of (1). Then $x^*$ is a stationary point of (1).*

*Proof.* Let $x^*$ be a local minimum of $f$ and assume by contradiction that is not a stationary point of (1). Then there exists $x \in C$ such that $\nabla f(x^*)(x - x^*) < 0$. Therefore, $f'(x, d) < 0$, where $d = x - x^*$. Hence, by Lemma 1.1 of Chapter 2, there exists $\epsilon \in (0, 1)$ such that $f(x^* + td) < f(x^*) \ \forall t \in (0, \epsilon)$. Since $C$ is convex, we have that $x^* + td = (1 - t)x^* + tx \in C$, leading to the conclusion that $x^*$ is not a local optimum of (1), which is a contradiction. $\qquad\square$

**Theorem 3.4** (Stationarity as necessary and sufficient optimality condition for a convex problem). *Let $f \in \mathrm{C}^1(C)$, where $C$ is closed and convex and $f$ is also convex. Let $x^*$ be a local minimum of (1). Then $x^*$ is a stationary point of (1) iff $x^*$ is a optimal solution of (1).*

*Proof.* The necessity of the stationarity condition follows from Theorem 3.3. To prove the sufficiency, assume that $x^*$ is a stationary point of (1) and let $x \in C$. Then, the gradient characterization of convex functions (3) and stationarity of $x^*$, we get

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

which concludes the proof. $\qquad\square$

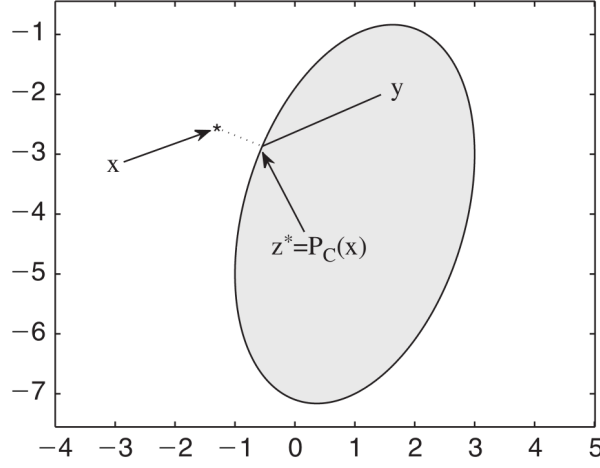Unfortunately, (5) is not an easy condition to check, we need something else.

Figure 1: The orthogonal projection operator.

## 3.2 Orthogonal Projection

We can instead characterize stationary points by using the projection operator. Given a nonempty closed convex set $C$, the orthogonal projection operator $P_C : \mathbb{R}^n \to C$ is defined by

$$P_C(x) = \operatorname{argmin} \left\{ ||x - y||^2 : y \in C \right\} \tag{6}$$

The orthogonal projection operator with input $x$ returns the vector in $C$ that is the closest (in $\ell_2$-norm) to $x$. Note that the orthogonal projection operator is defined as a solution of a convex optimization problem, specifically, a minimization of a convex quadratic function subject to a convex feasibility set. The first orthogonal projection theorem states that the orthogonal projection operator is in fact well-defined, meaning that the optimization problem in (6) has a unique optimal solution.

**Theorem 3.5** (First Projection Theorem). *Let $C$ be a nonempty closed convex set. Then problem* (6) *has a unique optimal solution.*

*Proof.* As $C$ is closed and $||x - y||^2$ is coercive, we have that the problem admits at least one solution (by Theorem 3.8 of Chapter 1). Moreover, $||x - y||^2$ is strictly convex as the objective function is quadratic with positive definite Hessian (the identity). Thus, from Theorem 3.2 we get that (6) has a unique solution. □

The second projection theorem, provides an useful characterization of the projection operator. Geometrically it states that for a given closed and convex set $C$, $x \in \mathbb{R}^n$, and for any $y \in C$, the angle between $x - P_C(x)$ and $y - P_C(x)$ is obtuse. This phenomenon is illustrated in Figure 1.

**Theorem 3.6** (Second Projection Theorem). *Let $C$ be a nonempty closed convex set. Then $z = P_C(x)$ iff*

$$(x - z)^T (y - z) \le 0 \quad \forall \, y \in C. \tag{7}$$

*Proof.* $z = P_C(x)$ iff it is the optimal solution of (6) iff (by Theorem 3.4)

$$\nabla f(z)^T (y - z) \ge 0 \quad \forall \, y \in C,$$

which concludes the proof as $\nabla f(z) = x - z$. □

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 3.7** (Nonexpansiveness of the projection operator). *Let $C$ be a closed and convex set. Then, for any $v, w \in \mathbb{R}^n$*

a)
$$(P_C(v) - P_C(w))^T (v - w) \ge ||P_C(v) - P_C(w)||^2 \tag{8}$$

b)
$$||P_C(v) - P_C(w)|| \le ||v - w||. \tag{9}$$

4

*Proof.* From Theorem 3.6 we have that for any $x \in \mathbb{R}^n$ and $y \in C$

$$(x - P_C(x))^T(y - P_C(x)) \leq 0.$$

Replacing $x = v$ and $y = P_C(w)$ we have

$$(v - P_C(v))^T(P_C(w) - P_C(v)) \leq 0.$$

Replacing, instead, $x = w$ and $y = P_C(v)$

$$(w - P_C(w))^T(P_C(v) - P_C(w)) \leq 0.$$

Now, summing the two inequalities we get

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

and hence,

$$(P_C(v) - P_C(w))^T(v - w) \geq ||P_C(w) - P_C(v)||^2.$$

To prove (9), we note that if $P_C(v) = P_C(w)$, the inequality is trivial. Thus, we assume $P_C(v) \neq P_C(w)$. Then by Cauchy-Schwartz, we have

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

which combined with (8) gives

$$||P_C(v) - P_C(w)||^2 \leq ||P_C(v) - P_C(w)|| \cdot ||v - w||,$$

which concludes the proof as $P_C(v) \neq P_C(w)$. $\qquad\square$

Coming back to stationarity, let us provide the alternative characterization of a stationary point through the projection operator. Notice that this theorem holds also when $f$ is non-convex.

**Theorem 3.8.** *Let $f \in C^1(C)$ with $C$ closed and convex and let $s > 0$. $x^*$ is a stationary point of the problem (1) iff*

$$x^* = P_C(x^* - s\,\nabla f(x^*)). \tag{10}$$

*Proof.* By the second projection theorem (Theorem 3.6), we get that $x^* = P_C(x^* - s\,\nabla f(x^*))$ iff

$$(x^* - s\,\nabla f(x^*) - x^*)^T(x - x^*) \leq 0,$$

which concludes the proof, as $x^*$ is a stationary point when $\nabla f(x^*)^T(x - x^*) \geq 0$ $\qquad\square$

## 3.3 Projected Gradient Method

The characterization of stationary points through equation (10) directly suggest a new algorithm for solving convex constrained optimization methods. As we will see later, this algorithm finds stationary points despite $f$ being convex or not.

---
**Algorithm 1:** Projected Gradient (PG) Method

    **Input:** $x_0 \in \mathbb{R}^n$, $\epsilon > 0$, $t \in (0, \frac{L}{2})$
1   $k = 0$
2   **while** $||x_{k-1} - x_k|| > \epsilon$ **do**
3     |  $x_{k+1} = P_C(x_k - t\,\nabla f(x_k)$
4     |  $k = k + 1$

---

The proof of convergence of PG is similar to that of GD. In particular, we first prove the Decrease Lemma for constrained optimization problem.

**Lemma 3.1** (Decrease Lemma for Convex Constrained Problems). *Let $f \in C_L^{1,1}(C)$, where $C$ is convex and closed. Then for any $x \in C$ and $t \in (0, \frac{2}{L})$ the following inequality holds*

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

*Proof.* Exercise. $\qquad\square$

It is now convenient to define the gradient mapping as

$$G_M(x) := M\left(x - P_C\left(x - \frac{1}{M}\nabla f(x)\right)\right) \quad \text{with } M > 0. \tag{11}$$

Note that in the unconstrained case $G_M(x) = \nabla f(x)$ so the gradient mapping is an extension of the usual gradient operator. In addition, by Theorem 3.8, $G_M(x) = 0$ iff $x$ is a stationary point of (1). This means that we can look at $\|G_M(x)\|$ as an optimality measure. Moreover, the sufficient decrease stated above can be rewritten as

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

This generalized sufficient decrease property allows us to prove similar results to those proven in the unconstrained case.

**Theorem 3.9** (Convergence of PG method)**.** *Let $f \in \mathrm{C}_\mathrm{L}^{1,1}(C)$, with $C$ closed and convex. Let $\{x_k\}_k$ be a sequence generated by Algorithm 1 for solving (1). Assume that $f$ is bounded below over $C$. Then we have the following*

(a) *The sequence $\{f(x_k)\}_k$ is nonincreasing. In addition, for any $k \geq 0$, $f(x_{k+1}) < f(x_k)$ unless $x_k$ is a stationary point.*

(b) *$G_{\frac{1}{t}}(x_k) \to 0$ as $k \to \infty$.*

Notice that the theorem above only ensures convergence to a stationary point, which in the non-convex case might not be a global minimum. Also, the rate of convergence of PG is the same as that of GD, that is $\mathcal{O}(\frac{1}{\sqrt{T}})$. If we assume $f$ to be convex, we can instead ensure a faster rate of convergence, moreover, thanks to Theorem 3.4 all stationary points of (1) are global minima.

**Theorem 3.10** (Convergence of PG method for convex problems)**.** *Let $f \in \mathrm{C}_\mathrm{L}^{1,1}(C)$ be convex, with $C$ closed and convex. Let $\{x_k\}_k$ be a sequence generated by Algorithm 1 for solving (1). Assume that the set of optimal solutions $X^*$ is nonempty and that $f^*$ is the optimal value. Then we have the following*

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

$$2t(f(x_{k+1}) - f^*) \leq \|x_k - x^*\|^2 - \|x_{k+1} - x^*\|^2,$$

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

$$f(x_n) - f^* \leq \frac{\|x_0 - x^*\|}{2tn}.$$

*Proof.* By the Descent Lemma (for unconstrained optimization, Lemma 2.2 from Chapter 2), we have

$$f(x_{k+1}) \leq f(x_k) + \nabla f(x_k)^T(x_{k+1} - x_k) + \frac{L}{2}\|x_{k+1} - x_k\|^2.$$

Let $x^*$ be a global minimum of (1), then the gradient characterization of convexity (3) implies that $f(x_k) \leq f(x^*) + \nabla f(x_k)^T(x_k - x^*)$, which together with the previous inequality implies that

$$f(x_{k+1}) \leq f(x^*) + \nabla f(x_k)^T(x_k - x^*) + \nabla f(x_k)^T(x_{k+1} - x_k) + \frac{L}{2}\|x_{k+1} - x_k\|^2. \tag{12}$$

By the second projection theorem (7) applied on the projected point $x_{k+1}$, we have that

$$(x_k - t\nabla f(x_k) - x_{k+1})^T(x^* - x_{k+1}) \leq 0$$

if and only if

$$\nabla f(x_k)^T(x_{k+1} - x^*) + \frac{1}{t}(x_k - x_{k+1})^T(x^* - x_{k+1}) \leq 0$$

if and only if

$$\nabla f(x_k)^T(x_{k+1} - x^*) \leq \frac{1}{t}(x_k - x_{k+1})^T(x_{k+1} - x^*).$$

Therefore, from the above inequality, (12) and $t \leq \frac{1}{L}$, we get

$$
\begin{aligned}
f(x_{k+1}) &\leq f(x^*) + \nabla f(x_k)^T(x_k - x^*) + \nabla f(x_k)^T(x_{k+1} - x_k) + \frac{L}{2}\|x_{k+1} - x_k\|^2 \\
&= f(x^*) + \nabla f(x_k)^T(x_{k+1} - x^*) + \frac{L}{2}\|x_{k+1} - x_k\|^2 \\
&\leq f(x^*) + \frac{1}{t}(x_k - x_{k+1})^T(x_{k+1} - x^*) + \frac{L}{2}\|x_{k+1} - x_k\|^2 \\
&\leq f(x^*) + \frac{1}{t}(x_k - x_{k+1})^T(x_{k+1} - x^*) + \frac{1}{2t}\|x_{k+1} - x_k\|^2 \\
&= f(x^*) + \frac{1}{2t}(x_k - x_{k+1})^T(x_{k+1} - x^* + x_k - x^*) \\
&= f(x^*) + \frac{1}{2t}(x_k - x_{k+1} + x^* - x^*)^T(x_{k+1} - x^* + x_k - x^*) \\
&= f(x^*) + \frac{1}{2t}(x_k - x^*)^T(x_{k+1} - x^* + x_k - x^*) + \frac{1}{2t}(x^* - x_{k+1})^T(x_{k+1} - x^* + x_k - x^*) \\
&= f(x^*) + \frac{1}{2t}\left(\|x_k - x^*\|^2 + (x_k - x^*)^T(x_{k+1} - x^*) - (x_k - x^*)^T(x_{k+1} - x^*) - \|x_{k+1} - x^*\|^2\right) \\
&= f(x^*) + \frac{1}{2t}\left(\|x_k - x^*\|^2 - \|x_{k+1} - x^*\|^2\right)
\end{aligned}
$$

establishing part (a). To achieve (b), we sum the inequalities (a) for $k = 0, 1, \ldots, n-1$ and obtain

$$
\|x_n - x^*\|^2 - \|x_0 - x^*\|^2 \leq 2t \sum_{k=0}^{n-1}(f(x^*) - f(x_{k+1})) \leq 2tn(f(x^*) - f(x_n)),
$$

where in the last inequality we used the fact that $f(x_{k+1}) \leq f(x_k)$, which, in turn, is a consequence of the Descent Lemma and the fact that $t \in (0, \frac{1}{L})$. Thus,

$$
f(x_n) - f(x^*) \leq \frac{\|x_0 - x^*\|^2 - \|x_n - x^*\|^2}{2tn} \leq \frac{\|x_0 - x^*\|^2}{2tn}.
$$

$\square$

# 4  KKT Conditions

In this chapter we will derive the necessary optimality conditions, i.e., Karush-Kunh-Tucker conditions, for the most general case where $C$ is possibly nonconvex. In particular, we consider problems of the following shape

$$
\begin{aligned}
\min \quad & f(x) \\
\text{s.t.} \quad & g_i(x) \leq 0, \quad i = 0, \ldots, m,
\end{aligned}
\tag{13}
$$

where $f, g_i \in \mathrm{C}^1(\mathbb{R})$ but possibly not convex. Notice that this class of problems is very general, as equality constraints can be included observing that $h(x) = 0$ can be replaced by 2 inequalities $h(x) \leq 0$ and $-h(x) \leq 0$. From now on $C := \{x \in \mathbb{R}^n : g_i(x) \leq 0, \ i \in [m]\}$.

**Definition 4.1** (Feasible Descent Direction). *A vector $d$ is called feasible descent direction at $x \in C$ if $\nabla f(x)^T d < 0$ and there exists $\epsilon > 0$ such that $x + td \in C$ for all $t \in [0, \epsilon]$.*

Obviously, a necessary local optimality condition of a point x is that it does not have any feasible descent directions.

**Lemma 4.1.** *Let $x^*$ be a local optimum of (13), then there are no feasible descent directions at $x^*$.*

*Proof.* The proof goes by contradiction and follows directly from the definition of feasible descent direction and directional derivative. $\square$

**Definition 4.2** (Active Constraints). *Let $g_i(x) \leq 0$, $i \in [m]$ be a set of inequalities. The active constraints at $\bar{x}$ are the constraints satisfied as equalities at $\bar{x}$. The set of active constraints is denoted by $I(x) := \{i \in [m] : g_i(x) = 0\}$.*
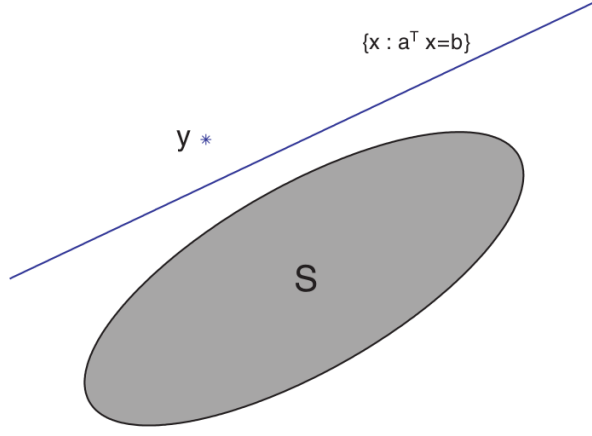
Figure 2: Strict separation of point from a closed and convex set.

**Lemma 4.2.** *Let $x^*$ be a local minimum of the problem* (13) *and let $I(x^*)$ be the set of active constraints at $x^*$. Then, there does not exists a vector $d \in \mathbb{R}^n$ such that*

$$\nabla f(x^*)^T d < 0,$$
$$\nabla g_i(x^*)^T d < 0, \quad i \in I(x^*).$$

*Proof.* Suppose by contradiction that $d$ satisfies the system of inequalities above. Then it follows that there exists $\epsilon_1 > 0$ such that $f(x^* + td) < f(x^*)$ and $g_i(x + td) < g(x^*) = 0$ for any $t \in (0, \epsilon_1)$ and $i \in I(x^*)$. For any $i \notin I(x^*)$, we have $g_i(x^*) < 0$ and hence, by continuity of $g_i$ it follows that there exists $\epsilon_2 > 0$ such that $g_i(x^* + td) < 0$ for any $t \in (0, \epsilon_2)$ and $i \notin I$. We can thus conclude that

$$\nabla f(x^* + td)^T d < f(x^*),$$
$$\nabla g_i(x^* + td)^T d < 0, \quad i \in [m],$$

for all $t \in (0, \min\{\epsilon_1, \epsilon_2\})$, which is a contradiction to the local optimality of $x^*$. $\square$

We have thus shown that a necessary optimality condition for local optimality is the infeasibility of a certain system of strict inequalities. On the other hand, similarly to the stationarity condition, this system is difficult to use in practice. We will state now the Fritz-John conditions.

**Theorem 4.1** (Fritz-John Conditions)**.** *Let $x^*$ be a local minimum of the problem* (13)*. Then there exists multipliers $\lambda_0, \ldots, \lambda_1, \ldots, \lambda_m$ such that they are not all zeros and such that*

$$\lambda_0 \nabla f(x^*) + \sum_{i=1}^m \lambda_i \nabla g_i(x^*) = 0, \tag{14}$$
$$\lambda_i g_i(x^*) = 0 \quad i = 1, \ldots, m.$$

In order to prove this theorem we need a rather large digression into the Alternative Theorems. We begin with a very simple yet powerful result on convex sets, namely the separation theorem between a point and a closed convex set. This result will be the basis for all the optimality conditions that will be discussed later on.

**Theorem 4.2** (Strict Separation Theorem)**.** *Let $C$ be a closed and convex set and let $y \notin C$. Then there exists $p \in \mathbb{R}^n \setminus \{0\}$ and $\alpha \in \mathbb{R}^n$ such that*

$$p^T y > \alpha \qquad and \qquad p^T x \leq \alpha \ \forall x \in C.$$

*Proof.* By the second projection theorem, the vector $\bar{x} = P_C(y) \in C$ satisfies

$$(y - \bar{x})^T (x - \bar{x}) \leq 0 \ \forall x \in C$$

which is the same as

$$(y - \bar{x})^T x \leq (y - \bar{x})^T \bar{x} \ \forall x \in C.$$

Denote $p = y - x \neq 0$ (since $y \notin C$) and $\alpha = (y - \bar{x})^T \bar{x}$. Then we have that $p^T x \leq \alpha \ \forall x \in C$. On the other hand,

$$p^T y = (y - \bar{x})^T y = (y - \bar{x})^T (y - \bar{x}) + (y - \bar{x})^T \bar{x} = ||y - \bar{x}||^2 + \alpha > \alpha,$$

and the result is established. $\square$

Now, before going on with two more alternative theorems, we need to show that the conic hull of a fine set is closed and convex.

**Definition 4.3** (Conic Hull). *Let $S \subseteq \mathbb{R}^n$. Then the conic hull of $S$, denoted by $\mathrm{cone}(S)$, is the set comprising all the conic combinations of vectors from $S$:*

$$\mathrm{cone}(S) := \left\{ \sum_{i=1}^{k} \lambda_i x_i : x_1, \ldots, x_k \in S, \lambda \in \mathbb{R}_+^k, k \in \mathbb{N} \right\}$$

**Lemma 4.3.** *Let $a_1, a_2, \ldots a_k \in \mathbb{R}^n$. Then $\mathrm{cone}(\{a_1, \ldots, a_k\})$ is closed and convex.*

*Proof.* Exercise. $\qquad\square$

We can now go on with the next alternative theorem.

**Lemma 4.4** (Farkas' lemma, second formulation). *Let $c \in \mathbb{R}^n$ and $A \in \mathbb{R}^{m \times n}$. Then the following two claims are equivalent:*

   *M. The implication $Ax \leq 0 \Rightarrow c^T x \leq 0$ holds true.*

   *N. There exists $y \in \mathbb{R}_+^n$ such that $A^T y = c$.*

*Proof.* Suppose that system $N$ is feasible. To see that the implication $M$ holds, suppose that $Ax \leq 0$ for some $x \in \mathbb{R}^n$. Then, multiplying this inequality from the left by $y^T$ (a valid operation since $y \geq 0$) yields

$$y^T A x \leq 0,$$

which concludes the thesis by noticing that $c^T = y^T A$.

The reverse direction is not so obvious. Suppose that the implication $M$ is satisfied, and let us show that system $N$ is feasible. Suppose in contradiction that system $N$ is infeasible, and consider the following set

$$S = \{x \in \mathbb{R}^n : x = A^T y \; y \in \mathbb{R}_+^n\},$$

which is closed and convex thanks to Lemma 4.3. The infeasibility of B means that $c \notin S$. By Theorem 4.2, it follows that there exists a vector $p \in \mathbb{R}^n \setminus \{0\}$ and $\alpha \in \mathbb{R}$ such that $p^T c > \alpha$ and

$$p^T x \leq \alpha \; \forall x \in S \tag{15}$$

Since $0 \in S$, from (15) we have that $\alpha \geq 0$ and hence also $p^T c > 0$. In addition, (15) is equivalent to

$$p^T A^T y \leq \alpha \;\; \forall y \geq 0$$

or to

$$(Ap)^T y \leq \alpha \;\; \forall y \geq 0.$$

Let us now prove that $Ap \leq 0$ (notice that this means component-wise). By contradiction, if there was an index $i \in \{1, 2, \ldots, m\}$ such that $(Ap)_i > 0$, then for $y = \beta e_i$ we would have that $(Ap)^T y = \beta (Ap)_i$ which is an expression that goes to $\infty$ as $\beta \to \infty$, and, thus, cannot be bounded by a constant $\alpha$. At this point we have found a system for which $Ap \leq 0$ and $p^T c = c^T p > 0$, which contradicts the implication $M$. $\qquad\square$

In order to prove Gordon's alternative theorem, we are going to use Farkas' lemma in the following formulation

**Lemma 4.5** (Farkas' lemma, first formulation). *Let $c \in \mathbb{R}^n$ and $A \in \mathbb{R}^{m \times n}$. Then exactly one of the following system has a solution:*

   *I. $Ax \leq 0, c^T x > 0$.*

   *II. $A^T y = c, y \geq 0$.*

To show that the two formulations are equivalent, let us notice that $II$ is equivalent to $N$ and let us write down the truth table of $M$ and $I$. In particular, let us call $M_1$ the statement $Ax \leq 0$ and $M_2$ the statement $c^T x \leq 0$ and notice that $I = M_1 \wedge \bar{M_2}$.

| $M_1$ | $M_2$ | M= $M_1 \Rightarrow M_2$ | $I = M_1 \wedge \bar{M_2}$ |
|---|---|---|---|
| F | F | T | F |
| F | T | T | F |
| T | F | F | T |
| T | T | T | F |

In particular, this means that the two formulations are equivalent as the first formulation (Lemma 4.5) states that exactly one between $I$ and $II$ has solutions while the second formulation (Lemma 4.4) states that $M$ and $N$ are equivalent.

# References