

## 3 The Study of Smooth Optimization Problems

---

This chapter plays a central role in this monograph. In its first section, we present smooth optimization problems and deduce existence conditions for minimality. We take this opportunity to prove and discuss the Ekeland Variational Principle and its consequences. We then obtain necessary conditions for optimality as well as sufficient optimality conditions of the first and second-order for smooth objective functions under geometric restrictions (i.e., restrictions of the type  $x \in M$ , where  $M$  is an arbitrary set). The second section is dedicated to the investigation of optimality conditions under functional restrictions (with equalities and inequalities). The main aim is to deduce Karush-Kuhn-Tucker conditions and to introduce and compare several qualification conditions. Special attention is paid to the case of convex and affine data. Subsequently, we derive second-order optimality conditions for the case of functional restrictions. The last section of this chapter includes two examples which show that, for practical problems, the computational challenges posed by the optimality conditions are sometimes not easy to solve.

---

### 3.1 General Optimality Conditions

Let  $U \subset \mathbb{R}^p$  be a nonempty open set,  $f : U \rightarrow \mathbb{R}$  be a function and  $M \subset U$  a nonempty set. We are interested in studying the minimization problem for the function  $f$  when its argument belongs to  $M$ . Formally, we write this problem in the following form:

$$(P) \min f(x), \text{ subject to } x \in M.$$

The function  $f$  is called the objective function, or cost function, and the set  $M$  is called the set of feasible points of the problem  $(P)$ , or the set of constraints, or the set of restrictions.

We should say from the very beginning that we shall study the minimization of  $f$ , but, by virtue of the relation  $\max f = -\min(-f)$ , similar results for its maximization could be obtained. Let us start by defining the notion of a solution associated to the problem  $(P)$ .

**Definition 3.1.1.** *One says that  $\bar{x} \in M$  is a local solution (or, simply, solution) of the problem  $(P)$  or minimum point of the function  $f$  on the set  $M$  if there exists a neighborhood  $V$  of  $\bar{x}$  such that  $f(\bar{x}) \leq f(x)$  for every  $x \in M \cap V$ . If  $V = \mathbb{R}^p$ , one says that  $\bar{x}$  is a global solution of  $(P)$  or minimal global point of  $f$  on  $M$ .*

Of course, for the maximization problem, the corresponding solution is clear. Let us mention that in this chapter we will only deal with the case of smooth functions (up to the order two, i.e.,  $f \in C^2$ ).

**Remark 3.1.2.** We shall distinguish two main situations for the study of problem (P) : (i) the case where  $M = U$  and (ii) the case where  $M$  appears as an intersection between a closed set of  $\mathbb{R}^p$  and  $U$ . In the former case, we say that the optimization problem (P) has no constraints (or restrictions), while in the latter, we call this a problem with constraints. In the case of a problem without constraints, we use as well the term “local minimum point of  $f$ ” instead of local solution. Let us observe that in Definition 3.1.1, if  $\bar{x} \in \text{int } M$ , then  $\bar{x}$  is a local solution of the unconstrained problem (it is enough to take a smaller neighborhood  $V$  such that  $V \subset M$ ). So, in the case of problems with restrictions, the interesting situation is when  $\bar{x} \in \text{bd } M$ . If  $\bar{x} \in \text{int } M$  we say that the restriction is inactive.

In the next sections, the two cases mentioned above will be treated together, but afterwards the discussion will split.

The basis of Optimization Theory consists of two fundamental results: the Weierstrass Theorem which ensures the existence of extrema and the Fermat Theorem (on stationary, or critical points) which gives a necessary condition for a point to be an extremum (without constraints) of a function. The theory follows the main trajectory of these fundamental results: on one hand, the study of existence conditions, and, on the other hand, the study of the (necessary and sufficient) optimality conditions.

We now recall these basic results. The classical Weierstrass Theorem, also given in the first chapter, states that a continuous function on a compact interval has a global minimum point on that interval. We have shown already that some conditions can be relaxed. Here is the theorem again.

**Theorem 3.1.3** (Weierstrass Theorem). *If  $f : K \subset \mathbb{R}^p \rightarrow \mathbb{R}$  is a lower (upper) semicontinuous function and  $K$  is a compact set, then  $f$  is lower (upper) bounded on  $K$  (i.e.,  $f(K)$  is a bounded below (above) set) and  $f$  attains its minimum (maximum) on  $K$ , i.e., there exists  $\bar{x} \in K$  such that  $\inf_{x \in K} f(x) = f(\bar{x})$  ( $\sup_{x \in K} f(x) = f(\bar{x})$ , respectively).*

**Remark 3.1.4.** *The compactness of  $K$  is essential in the Weierstrass Theorem because otherwise the conclusion does not hold. As an example, let us consider the continuous function  $f : (0, 1] \rightarrow \mathbb{R}$ ,  $f(x) = x^{-1} \sin(x^{-1})$ . Clearly,  $\inf_{(0,1]} f(x) = -\infty$ , and  $\sup_{(0,1]} f(x) = +\infty$ .*

We now present Fermat Theorem. The particular case of a function of one real variable was presented in Section 1.3. The proof of the theorem will be given later in this chapter.

**Theorem 3.1.5** (Fermat Theorem). *Let  $S \subset \mathbb{R}^p$  be a set and  $a \in \text{int } S$ . If  $f : S \rightarrow \mathbb{R}$  is of class  $C^1$  in a neighborhood of  $a$ , and  $a$  is a local minimum or maximum point of  $f$ , then  $a$  is also a stationary (or critical) point of  $f$ , i.e.,  $\nabla f(a) = 0$ .*

Some remarks are in order to understand the applications of Fermat Theorem.

**Remark 3.1.6.** 1. *The converse of Fermat Theorem is not true: for instance, the derivative of  $f : \mathbb{R} \rightarrow \mathbb{R}$ ,  $f(x) = x^3$  vanishes at 0, but this point is neither minimum nor maximum of  $f$ .*

2. *The interiority condition for  $a$  is essential since, without this assumption, the conclusion does not hold: take  $f : [0, 1] \rightarrow [0, 1]$ ,  $f(x) = x$  which has at  $\bar{x} = 0$  a minimum point where the derivative is not 0. Therefore, in view of Remark 3.1.2, one can say that Fermat Theorem applies only to unconstrained problems.*

3. *If  $S$  is compact and  $f$  is continuous on  $S$  and differentiable on  $\text{int } S$ , it is possible to have  $\nabla f(x) \neq 0$  for every  $x \in \text{int } S$ , and in such a case the extreme points of  $f$  on  $S$ , which surely exist from Weierstrass Theorem, lie on the boundary of  $S$ .*

We now start our discussion on the existence conditions for the minimum points.

**Theorem 3.1.7.** *Let  $f : \mathbb{R}^p \rightarrow \mathbb{R}$  be a lower semicontinuous function and  $M \subset \mathbb{R}^p$  be a nonempty, closed set. If there exists  $v > \inf_{x \in M} f(x)$  such that the level set of  $f$  relative to  $M$ , i.e.,  $M \cap N_v f = \{x \in M \mid f(x) \leq v\}$ , is bounded, then  $f$  attains its global minimum on  $M$ .*

*Proof* It is obvious that if  $f$  has a global minimum on  $M$ . Similarly, it also has a global minimum on  $M \cap N_v f$ . Since this set is compact and  $f$  is lower semicontinuous, from Weierstrass Theorem 3.1.3, we infer that  $f$  is lower bounded and attains its global minimum on  $M \cap N_v f$ , whence on  $M$ , too.  $\square$

Obviously, in the above result, if  $M$  is bounded, then the hypothesis is automatically fulfilled. The interesting case is where  $M$  is unbounded and in this situation we ensure the boundedness assumption by imposing a certain condition on  $f$ .

**Proposition 3.1.8.** *Let  $f : \mathbb{R}^p \rightarrow \mathbb{R}$  be a function and  $M \subset \mathbb{R}^p$  be a closed, unbounded set. If  $\lim_{x \in M, \|x\| \rightarrow \infty} f(x) = \infty$  (i.e., for every sequence  $(x_k) \subset M$  with  $\|x_k\| \rightarrow \infty$ , one has  $f(x_k) \rightarrow \infty$ ), then the set  $N_v f \cap M$  is bounded for every  $v > \inf_{x \in M} f(x)$ .*

*Proof* Let  $v > \inf_{x \in M} f$ . If the set  $N_v f \cap M$  would be unbounded, then there exists  $(x_k) \subset N_v f \cap M$  with  $\|x_k\| \rightarrow \infty$ . On one hand, by our assumption,  $\lim f(x_k) = \infty$  and, on the other hand,  $f(x_k) \leq v$  for every  $k \in \mathbb{N}$ , which is absurd. So  $N_v f \cap M$  is bounded.  $\square$

If the condition  $\lim_{x \in M, \|x\| \rightarrow \infty} f(x) = \infty$  holds, we say that  $f$  is coercive relative to  $M$ . If  $\lim_{\|x\| \rightarrow \infty} f(x) = \infty$ , then we say that  $f$  is coercive.

Some special conditions to ensure the existence and the attainment of the minimum on a bounded, not closed set can be given as well. We present here such a result which we are going to use in Chapter 6.

**Proposition 3.1.9.** *Let  $D, C \subset \mathbb{R}^p$  be two sets such that  $C$  is compact and  $D \cap C \neq \emptyset$ . Let  $\varphi : D \cap C \rightarrow \mathbb{R}$  be a continuous function. Suppose that the following condition holds: for every sequence  $(x_k) \subset D \cap C$ ,  $x_k \rightarrow \bar{x} \in \text{bd } D \cap C$ , the sequence  $(\varphi(x_k))$  is unbounded above. Then  $\varphi$  is lower bounded and attains its minimum on  $D \cap C$ .*

*Proof* Firstly, let us observe that  $\varphi$  cannot be constant. Let  $x_0 \in D \cap C$  such that  $\varphi(x_0) > \inf_{x \in D \cap C} \varphi(x)$ . It is enough to show that the level set  $A := \{x \in D \cap C \mid \varphi(x) \leq \varphi(x_0)\}$  is compact (see the proof of Theorem 3.1.7). Obviously,  $A$  is bounded (as a subset of  $C$ ). It remains to show that  $A$  is closed. Let  $(x_k) \subset A$ ,  $x_k \rightarrow \bar{x}$ . Suppose, by contradiction, that  $\bar{x} \notin A$ . Then, from the closedness of  $C$ ,

$$\bar{x} \in (\text{cl } D \setminus D) \cap C \subset \text{bd } D \cap C.$$

By assumption,  $(\varphi(x_k))$  is unbounded above, which contradicts the definition of the set  $A$ . Hence,  $A$  is compact and the conclusion follows. □

In general, a global minimum is local minimum, but, of course, the converse is false. We shall derive a condition for the fulfillment of this converse. We now give a necessary and sufficient condition for a point to be a global minimum.

**Theorem 3.1.10.** *Let  $f : \mathbb{R}^p \rightarrow \mathbb{R}$  be a continuous function and  $\bar{x} \in \mathbb{R}^p$ . Then the following assertions are equivalent:*

- (i)  $\bar{x}$  is a global minimum point of  $f$ ;
- (ii) every  $x \in \mathbb{R}^p$  with  $f(x) = f(\bar{x})$  is a local minimum point of  $f$ .

*Proof* The implication (i)  $\Rightarrow$  (ii) is obvious. Suppose that (ii) holds, but  $\bar{x}$  is not a global minimum point. Then, there exists  $u \in \mathbb{R}^p$  with  $f(u) < f(\bar{x})$ . We define  $\varphi : [0, 1] \rightarrow \mathbb{R}$  by  $\varphi(t) := f(t\bar{x} + (1-t)u)$ . The set  $S := \{t \in [0, 1] \mid \varphi(t) = f(\bar{x})\}$  is nonempty ( $1 \in S$ ), closed (from the continuity of  $f$ ) and bounded. Then there is  $t_0 = \min S$ . Clearly,  $t_0 \in (0, 1]$ . From the hypothesis,  $f(t_0\bar{x} + (1-t_0)u) = f(\bar{x})$  tells us that the point  $t_0\bar{x} + (1-t_0)u$  is a local minimum of  $f$ . Consequently, there is  $\varepsilon > 0$  such that for every  $t \in [0, 1] \cap (t_0 - \varepsilon, t_0 + \varepsilon)$ ,  $\varphi(t) \geq \varphi(t_0)$ . Since  $t_0 = \min S$ , if one takes  $t_1 \in [0, 1] \cap (t_0 - \varepsilon, t_0)$ , the strict inequality  $\varphi(t_1) > \varphi(t_0) > \varphi(0)$  holds. The function  $f$  being continuous, it has the Darboux property (or intermediate value property), whence there exists  $t_2 \in (0, t_1)$  with  $\varphi(t_2) = \varphi(t_0) = f(\bar{x})$ , and this contradicts the minimality of  $t_0$ . Therefore, the assumption made was false, hence the conclusion holds. □

Notice that, if  $f$  is not continuous, the result does not hold. For this, it is sufficient to analyze the following function  $f : \mathbb{R} \rightarrow \mathbb{R}$ ,

$$f(x) = \begin{cases} -x - 1, & x \in (-\infty, -1] \\ x + 1, & x \in (-1, 0) \\ -1, & x = 0 \\ -x + 1, & x \in (0, 1) \\ x - 1, & x \in [1, \infty). \end{cases}$$

Let us observe that, by using the function  $f : \mathbb{R} \rightarrow \mathbb{R}, f(x) = x^3$  and  $\bar{x} = 0$ , that one cannot replace in the item (ii) above the local minimality by the stationarity.

The above results ensure sufficient conditions for the existence of minimum points under compactness assumptions for the level sets. Conversely, it is clear that the lower boundedness of the function is a necessary condition for the existence of minimum points, but the boundedness of the level sets is not. For instance, the function  $f : (0, \infty) \rightarrow \mathbb{R}, f(x) = (x - 1)^2 e^{-x}$  attains its minimum at  $\bar{x} = 1$  and the minimal value is 0, but  $N_v f$  is not bounded for every value of  $v > 0 = \inf\{f(x) \mid x \in (0, \infty)\}$ .

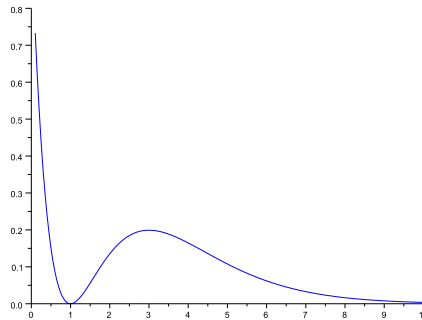


Figure 3.1: The graph of  $(x - 1)^2 e^{-x}$ .

One may like to have a notion of approximate minimum which is advantageous to always exist if  $f$  is lower bounded. We work here with unconstrained problems in order to better illustrate the main ideas.

**Definition 3.1.11.** Let  $f : \mathbb{R}^p \rightarrow \mathbb{R}$  be a lower bounded function and take  $\varepsilon > 0$ . A point  $x_\varepsilon$  is called  $\varepsilon$ -minimum of  $f$  if

$$f(x_\varepsilon) \leq \inf_{x \in \mathbb{R}^p} f(x) + \varepsilon.$$

Since  $\inf_{x \in \mathbb{R}^p} f(x) \in \mathbb{R}$ , the existence of  $\varepsilon$ -minima for every positive  $\varepsilon$  is ensured. We use the generic term of approximate minima for  $\varepsilon$ -minima.

We now present a very important result, the Ekeland Variational Principle, which states that close to an approximate minimum point one can find a genuine minimum point for some perturbation of the initial function. This results was proved by the French mathematician Ivar Ekeland in 1974.

**Theorem 3.1.12** (Ekeland Variational Principle). *Let  $f : \mathbb{R}^p \rightarrow \mathbb{R}$  be a lower semicontinuous and lower bounded function. Let  $\varepsilon > 0$  and let  $x_\varepsilon$  be an  $\varepsilon$ -minimum of  $f$ . Then, for every  $\delta > 0$  there exists  $\bar{x}_\varepsilon \in \mathbb{R}^p$  that have the following properties:*

$$\begin{aligned} f(\bar{x}_\varepsilon) &\leq f(x_\varepsilon), \\ \|\bar{x}_\varepsilon - x_\varepsilon\| &\leq \delta, \\ f(\bar{x}_\varepsilon) &\leq f(x) + \varepsilon\delta^{-1} \|x - \bar{x}_\varepsilon\|, \quad \forall x \in \mathbb{R}^p. \end{aligned}$$

*Proof* Let us consider the function  $g : \mathbb{R}^p \rightarrow \mathbb{R}$

$$g(x) = f(x) + \varepsilon\delta^{-1} \|x - x_\varepsilon\|.$$

Using the assumptions on  $f$ , we infer that  $g$  is lower semicontinuous and lower bounded. Moreover,  $g$  is coercive, i.e.,

$$\lim_{\|x\| \rightarrow \infty} g(x) = +\infty.$$

From Theorem 3.1.7 and Proposition 3.1.8, the function  $g$  has a global minimum point, which we denote by  $\bar{x}_\varepsilon$ . Consequently,

$$f(\bar{x}_\varepsilon) + \varepsilon\delta^{-1} \|\bar{x}_\varepsilon - x_\varepsilon\| \leq f(x) + \varepsilon\delta^{-1} \|x - x_\varepsilon\|, \quad \forall x \in \mathbb{R}^p. \tag{3.1.1}$$

For  $x = x_\varepsilon$ , we get

$$f(\bar{x}_\varepsilon) \leq f(x_\varepsilon).$$

That is the first relation in the conclusion. On the other hand, using this inequality

$$f(x_\varepsilon) \leq \inf_{x \in \mathbb{R}^p} f(x) + \varepsilon,$$

with the relation (3.1.1), for  $x = x_\varepsilon$ , implies that:

$$\delta^{-1} \|\bar{x}_\varepsilon - x_\varepsilon\| \leq 1.$$

This is the second part of the conclusion. Relation (3.1.1) allows us to write, successively, for any  $x \in \mathbb{R}^p$ ,

$$\begin{aligned} f(\bar{x}_\varepsilon) &\leq f(x) + \varepsilon\delta^{-1} (\|x - x_\varepsilon\| - \|\bar{x}_\varepsilon - x_\varepsilon\|) \\ &\leq f(x) + \varepsilon\delta^{-1} \|x - \bar{x}_\varepsilon\|. \end{aligned}$$

That is the last part of the conclusion. The proof is complete. □

**Remark 3.1.13.** Notice that the Ekeland Variational Principle holds (with minor changes in the proof) if instead of the whole space  $\mathbb{R}^p$  one takes a closed subset of it.

Clearly, the point  $\bar{x}_\varepsilon$  is a global minimum point for the function

$$x \mapsto f(x) + \varepsilon\delta^{-1} \|x - \bar{x}_\varepsilon\|.$$

On the other hand, if we want  $\bar{x}_\varepsilon$  to be close to  $x_\varepsilon$  (i.e.,  $\delta$  to be small), then the perturbation term,  $\varepsilon\delta^{-1} \|\cdot - \bar{x}_\varepsilon\|$  is big. A compromise would be to choose  $\delta := \sqrt{\varepsilon}$ , and in this case one gets the next consequence.

**Corollary 3.1.14.** Let  $f : \mathbb{R}^p \rightarrow \mathbb{R}$  be a lower semicontinuous and lower bounded function. Take  $\varepsilon > 0$  and let  $x_\varepsilon$  be an  $\varepsilon$ -minimum of  $f$ . Then there exists  $\bar{x}_\varepsilon \in \mathbb{R}^p$  having the properties:

$$\begin{aligned} f(\bar{x}_\varepsilon) &\leq f(x_\varepsilon), \\ \|\bar{x}_\varepsilon - x_\varepsilon\| &\leq \sqrt{\varepsilon}, \\ f(\bar{x}_\varepsilon) &\leq f(x) + \sqrt{\varepsilon} \|x - \bar{x}_\varepsilon\|, \quad \forall x \in \mathbb{R}^p. \end{aligned}$$

The Ekeland Variational Principle has many applications. Some of these refer to the same issues of extreme points. The next example of such an application asserts that every differentiable function has approximate critical points (for which the norm of the differential is arbitrarily small).

**Theorem 3.1.15.** Let  $f : \mathbb{R}^p \rightarrow \mathbb{R}$  be a differentiable, lower bounded function. Then for every  $\varepsilon, \delta > 0$ , there exists  $\bar{x}_\varepsilon \in \mathbb{R}^p$  with  $f(\bar{x}_\varepsilon) \leq \inf_{x \in \mathbb{R}^p} f(x) + \varepsilon$  and  $\|\nabla f(\bar{x}_\varepsilon)\| \leq \varepsilon\delta^{-1}$ . In particular, there exists  $(x_n) \subset \mathbb{R}^p$  with

$$f(x_n) \rightarrow \inf_{x \in \mathbb{R}^p} f(x), \quad \nabla f(x_n) \rightarrow 0.$$

*Proof* Let  $x_\varepsilon$  be an  $\varepsilon$ -minimum of  $f$ . According to Theorem 3.1.12, there exists  $\bar{x}_\varepsilon \in \mathbb{R}^p$  with the three mentioned properties. Since  $f(\bar{x}_\varepsilon) \leq f(x_\varepsilon)$ , we infer that  $f(\bar{x}_\varepsilon) \leq \inf_{x \in \mathbb{R}^p} f(x) + \varepsilon$ . Let  $x := \bar{x}_\varepsilon + tu$  with  $u \in \mathbb{R}^p$  and  $t > 0$ . The relation  $f(\bar{x}_\varepsilon) \leq f(x) + \varepsilon\delta^{-1} \|x - \bar{x}_\varepsilon\|$  holds for any  $x \in \mathbb{R}^p$ , so we have

$$\frac{f(\bar{x}_\varepsilon + tu) - f(\bar{x}_\varepsilon)}{t} \geq -\varepsilon\delta^{-1} \|u\|.$$

Passing to the limit with  $t \rightarrow 0$ , we deduce

$$\nabla f(\bar{x}_\varepsilon)(u) \geq -\varepsilon\delta^{-1} \|u\|, \quad \forall u \in \mathbb{R}^p,$$

that is,

$$-\nabla f(\bar{x}_\varepsilon)(u) \leq \varepsilon\delta^{-1} \|u\|, \quad \forall u \in \mathbb{R}^p.$$

Changing  $u$  into  $-u$ , we get

$$\nabla f(\bar{x}_\varepsilon)(u) \leq \varepsilon\delta^{-1} \|u\|, \quad \forall u \in \mathbb{R}^p,$$

whence

$$\|\nabla f(\bar{x}_\varepsilon)(u)\| \leq \varepsilon \delta^{-1} \|u\|, \quad \forall u \in \mathbb{R}^p,$$

and this implies  $\|\nabla f(\bar{x}_\varepsilon)\| \leq \varepsilon \delta^{-1}$ . For  $\varepsilon := n^{-1}$ ,  $\delta = \sqrt{n^{-1}}$ ,  $n \in \mathbb{N}^*$  we obtain the second part of the conclusion.  $\square$

The Ekeland Variational Principle also allows us to prove the equivalence of several existence conditions for minimum points.

**Theorem 3.1.16.** *Let  $f : \mathbb{R}^p \rightarrow \mathbb{R}$  be a differentiable, lower bounded function. The following assertions are equivalent:*

- (i)  $\lim_{\|x\| \rightarrow \infty} f(x) = \infty$ ;
- (ii)  $N_\nu f$  is bounded for any  $\nu > \inf_{x \in \mathbb{R}^p} f(x)$ ;
- (iii) every sequence  $(x_n)$  for which  $(f(x_n))$  is convergent and  $\nabla f(x_n) \rightarrow 0$  has a convergent subsequence.

*Proof* The implication (i)  $\Rightarrow$  (ii) was already proved in Proposition 3.1.8, while (ii)  $\Rightarrow$  (iii) is a consequence of the fact that every bounded sequence has a convergent subsequence. Let us show that (iii)  $\Rightarrow$  (i). Suppose, by way of contradiction, that there exist  $c \in \mathbb{R}$  and a sequence  $(x_n) \subset \mathbb{R}^n$  with  $\|x_n\| \rightarrow \infty$  and  $f(x_n) \leq c$ , for every  $n \in \mathbb{N}$ . Clearly,  $c \geq \inf_{x \in \mathbb{R}^p} f(x)$ . For every  $n \in \mathbb{N}^*$  we choose

$$\varepsilon_n := c + n^{-1} - \inf_{x \in \mathbb{R}^p} f(x) > 0,$$

so,

$$f(x_n) < \inf_{x \in \mathbb{R}^p} f(x) + \varepsilon_n.$$

Let  $\delta_n := 2^{-1} \|x_n\| > 0$ . Like in Theorem 3.1.15 (and its proof) there exists  $\bar{x}_n$  with

$$\begin{aligned} f(\bar{x}_n) &\leq f(x_n) \leq \inf_{x \in \mathbb{R}^p} f(x) + \varepsilon_n, \\ \|\bar{x}_n - x_n\| &\leq \delta_n, \\ \|\nabla f(\bar{x}_n)\| &\leq \varepsilon_n \delta_n^{-1}. \end{aligned}$$

But,

$$\|\bar{x}_n\| \geq \|x_n\| - \|\bar{x}_n - x_n\| \geq \|x_n\| - 2^{-1} \|x_n\| = 2^{-1} \|x_n\|,$$

whence  $\|\bar{x}_n\| \rightarrow \infty$ . On the other hand,

$$\|\nabla f(\bar{x}_n)\| \leq \frac{2}{\|x_n\|} (c + n^{-1} - \inf_{x \in \mathbb{R}^p} f(x)) \rightarrow 0.$$

Since  $(f(\bar{x}_n))$  is bounded, it has a convergent subsequence. From (iii), one deduces that  $(\bar{x}_n)$  should also have such a subsequence, but this is not possible. Consequently, (i) holds.  $\square$

The condition (iii) in the above result is called Palais-Smale condition.

On the basis of Ekeland Variational Principle one obtains a new condition for the existence of the minimum points.



**Theorem 3.1.17.** Let  $\alpha > 0$  and  $f : \mathbb{R}^p \rightarrow \mathbb{R}$  be a lower semicontinuous and lower bounded function. Suppose that for every  $x \in \mathbb{R}^p$  with  $\inf_{x \in \mathbb{R}^p} f(x) < f(x)$  there exists  $z \in \mathbb{R}^p \setminus \{x\}$  such that

$$f(z) < f(x) - \alpha \|z - x\|.$$

Then  $f$  has a minimum global point.

*Proof* Suppose, to obtain a contradiction, that the conclusion does not hold. Then, for every  $x \in \mathbb{R}^p$ ,  $\inf_{x \in \mathbb{R}^p} f(x) < f(x)$ , so, by the assumptions made, there exists  $z_x \in \mathbb{R}^p \setminus \{x\}$  such that

$$f(z_x) < f(x) - \alpha \|z_x - x\|.$$

By the Ekeland Variational Principle for  $\varepsilon > 0$ ,  $\delta > 0$  with  $\varepsilon\delta^{-1} = \alpha$ , there is an element  $u \in \mathbb{R}^p$  with

$$f(u) \leq f(v) + \alpha \|v - u\|, \quad \forall v \in \mathbb{R}^p.$$

Then

$$f(u) \leq f(z_u) + \alpha \|z_u - u\| < f(u),$$

which is absurd. Hence the conclusion hold. □

A straightforward example of a function which satisfies the conditions of the above result is  $f : \mathbb{R} \rightarrow \mathbb{R}$ ,  $f(x) = \beta |x|$  for  $\beta > \alpha$ .

In the second part of this section, we present necessary optimality conditions and sufficient optimality conditions. At first, we deduce necessary optimality conditions that use the ideas developed around the construction and the study of the Bouligand tangent cone.

**Theorem 3.1.18** (First-order necessary optimality condition). *If  $\bar{x}$  is a local solution of  $(P)$  and  $f$  is differentiable at  $\bar{x}$ , then  $\nabla f(\bar{x})(u) \geq 0$  for every  $u \in T_B(M, \bar{x})$ .*

*Proof* Let  $V$  be a neighborhood of  $\bar{x}$  where  $f(\bar{x}) \leq f(x)$  for every  $x \in V \cap M$ . Let  $u \in T_B(M, \bar{x})$ . Then there exists  $(t_n) \subset (0, \infty)$  with  $t_n \rightarrow 0$  and  $(u_n) \rightarrow u$  such that for every  $n$ ,

$$\bar{x} + t_n u_n \in M.$$

Obviously, the sequence  $(t_n u_n)$  converges towards  $0 \in \mathbb{R}^p$  and, for  $n$  large enough,  $\bar{x} + t_n u_n$  belongs to  $V$ . Taking into account the differentiability of  $f$  at  $\bar{x}$ , there exists  $(\alpha_n) \subset \mathbb{R}$ ,  $\alpha_n \rightarrow 0$  such that for every  $n \in \mathbb{N}$ ,

$$f(\bar{x} + t_n u_n) = f(\bar{x}) + t_n \nabla f(\bar{x})(u_n) + t_n \|u_n\| \alpha_n,$$

whence,

$$\nabla f(\bar{x})(u_n) + \|u_n\| \alpha_n \geq 0,$$

for all  $n$  large enough. Passing to the limit for  $n \rightarrow \infty$ , we get the conclusion. □

**Remark 3.1.19.** *The conclusion of Theorem 3.1.18 could equivalently be written as*

$$-\nabla f(\bar{x}) \in N_B(M, \bar{x}).$$

**Remark 3.1.20.** *Taking into account Proposition 2.1.12, if  $\bar{x} \in \text{int } M$  (inactive restriction), Theorem 3.1.18 gives  $\nabla f(\bar{x})(u) \geq 0$  for every  $u \in \mathbb{R}^p$ . The linearity of  $\nabla f(\bar{x})$  implies  $\nabla f(\bar{x}) = 0$ , i.e., the Fermat Theorem on stationary points.*

We present now a second-order necessary optimality condition for the problem without restrictions.

**Theorem 3.1.21** (Second-order necessary optimality condition). *Let  $U \subset \mathbb{R}^p$  be an open set and  $\bar{x} \in U$ . If  $f : U \rightarrow \mathbb{R}$  is of class  $C^2$  on a neighborhood of  $\bar{x}$ , and  $\bar{x}$  is a local minimum point of  $f$ , then  $\nabla f(\bar{x}) = 0$  and  $\nabla^2 f(\bar{x})$  is positive semidefinite (that is,  $\nabla^2 f(\bar{x})(u, u) \geq 0$  for every  $u \in \mathbb{R}^p$ ).*

*Proof* Let  $V \subset U$  be a neighborhood of  $\bar{x}$  such that  $f(\bar{x}) \leq f(x)$  for every  $x \in V$  and  $f$  is of class  $C^2$  on  $V$ . The fact that  $\nabla f(\bar{x}) = 0$  follows from Fermat Theorem. As before, take  $u \in \mathbb{R}^p$  and  $(t_n) \subset (0, \infty)$  with  $t_n \rightarrow 0$ . Taylor Theorem 1.3.4 says that for every  $n \in \mathbb{N}$  there exists  $c_n \in (\bar{x}, \bar{x} + t_n u)$  such that

$$f(\bar{x} + t_n u) - f(\bar{x}) = t_n \nabla f(\bar{x})(u) + \frac{1}{2} t_n^2 \nabla^2 f(c_n)(u, u) = \frac{1}{2} t_n^2 \nabla^2 f(c_n)(u, u).$$

For  $n$  sufficiently large,  $f(\bar{x} + t_n u) - f(\bar{x}) \geq 0$ , whence

$$\nabla^2 f(c_n)(u, u) \geq 0,$$

and passing to the limit as  $n \rightarrow \infty$  we get  $c_n \rightarrow \bar{x}$ . Since  $f$  is of class  $C^2$ , we infer

$$\nabla^2 f(\bar{x})(u, u) \geq 0,$$

whence  $\nabla^2 f(\bar{x})$  is positive semidefinite. □

Obviously, if  $\bar{x} \in \text{int } U$  is a local maximum of  $f$ , then  $\nabla f(\bar{x}) = 0$  and  $\nabla^2 f(\bar{x})$  is negative semidefinite (i.e.,  $\nabla^2 f(\bar{x})(u, u) \leq 0$  for every  $u \in \mathbb{R}^p$ ). Further, if  $\nabla^2 f(\bar{x})$  is neither positive semidefinite, or negative semidefinite, then  $\bar{x}$  is not an extreme point of  $f$ .

In fact, many results in this book are generalizations, refinements of these results, or answer to different issues which naturally arise from their analysis.

One may consider whether the converses of these results are true. The answer is negative in both cases: it is sufficient to consider the function  $f : \mathbb{R} \rightarrow \mathbb{R}$ ,  $f(x) = x^3$  and  $\bar{x} = 0$ .

One can however impose supplementary conditions in order to get some equivalences, and the most important of these relates to convex functions.

**Theorem 3.1.22.** *Let  $U \subset \mathbb{R}^p$  be an open convex set and let  $f : U \rightarrow \mathbb{R}$  be a convex differentiable function. The next assertions are equivalent:*

- (i)  $\bar{x}$  is a global minimum point of  $f$  (on  $U$ );
- (ii)  $\bar{x}$  is a local minimum point of  $f$ ;
- (iii)  $\bar{x}$  is a critical point of  $f$  (i.e.,  $\nabla f(\bar{x}) = 0$ ).

*Proof* The implication (i)  $\Rightarrow$  (ii) is obvious for every function, and (ii)  $\Rightarrow$  (iii) follows from Fermat Theorem. Finally, the implication (iii)  $\Rightarrow$  (i) relies on the convexity of  $f$  and follows from Theorem 2.2.10.  $\square$

Therefore, for convex functions, the first-order necessary optimality condition (in the unconstrained case) is also sufficient. In this situation, the second order condition is automatically satisfied (according to Theorem 2.2.10).

Concerning the nature of the extreme points for convex functions, we record here some important aspects.

**Proposition 3.1.23.** *Let  $M \subset \mathbb{R}^p$  be a convex set and let  $f : M \rightarrow \mathbb{R}$  be a convex function. If  $\bar{x} \in M$  is a local minimum point of  $f$  on  $M$ , then  $\bar{x}$  is in fact a global minimum point of  $f$  on  $M$ . If  $u \in \text{int } M$  is a local maximum point of  $f$ , then  $u$  is a global maximum point of  $f$ .*

*Proof* Let  $\bar{x}$  be a local minimum point of  $f$  on  $M$ . Then there exists a convex neighborhood  $V$  of  $\bar{x}$  such that for every  $x \in V \cap M$ ,  $f(\bar{x}) \leq f(x)$ . Let  $x \in M$ . There exists  $\lambda \in (0, 1)$  such that  $y := (1 - \lambda)\bar{x} + \lambda x \in M \cap V$ . Then,

$$f(\bar{x}) \leq f(y) = f((1 - \lambda)\bar{x} + \lambda x) \leq (1 - \lambda)f(\bar{x}) + \lambda f(x),$$

that is,

$$\lambda f(\bar{x}) \leq \lambda f(x),$$

and the conclusion of the first part follows.

For the second part, there is a convex symmetric neighborhood  $V$  of 0 (a ball with the center 0, for instance) such that for every  $v \in V$ ,  $f(u + v) \leq f(u)$  and  $f(u - v) \leq f(u)$ . Then

$$f(u) = f\left(\frac{1}{2}(u + v) + \frac{1}{2}(u - v)\right) \leq \frac{1}{2}f(u + v) + \frac{1}{2}f(u - v) \leq f(u),$$

for all  $v \in V$ . Consequently,  $f(u + v) = f(u)$  for every  $v \in V$ . Therefore,  $u$  is a local (hence global) minimum point of  $f$ .  $\square$

**Proposition 3.1.24.** *Let  $M \subset \mathbb{R}^p$  be a convex set and let  $f : M \rightarrow \mathbb{R}$  be a convex function. If nonempty, the set of minimum points of  $f$  on  $M$  is convex. If, moreover,  $f$  is strictly convex, then this set has at most one element.*

*Proof* From the preceding result, if  $x_1, x_2 \in M$  are (global) minima of  $f$  on  $M$ , then  $f(x_1) = f(x_2)$ . The convexity implies  $f(x) = f(x_1)$  for every  $x \in [x_1, x_2]$ . Therefore, the first part is proved. Suppose now that  $f$  is strictly convex. If this is so, then we would have two different global minima, then  $f(x) < f(x_1)$  for every  $x \in (x_1, x_2)$ , which is not possible.  $\square$

For constrained problems involving convex functions, the first-order necessary optimality condition is, again, a sufficient optimality condition.

**Proposition 3.1.25.** *Let  $U \subset \mathbb{R}^p$  be a convex open set and let  $f : U \rightarrow \mathbb{R}$  be a convex, differentiable function. Let  $M \subset U$  be convex. The element  $\bar{x} \in M$  is a minimum point of  $f$  on  $M$  if and only if*

$$-\nabla f(\bar{x}) \in N(M, \bar{x}).$$

*Proof* Let  $\bar{x} \in M$  be a minimum point of  $f$  on  $M$ . Then, according to Theorem 3.1.18,  $\nabla f(\bar{x})(u) \geq 0$  for every  $u \in T(M, \bar{x})$ , that is

$$-\nabla f(\bar{x}) \in T(M, \bar{x})^- = N(M, \bar{x}).$$

Conversely, we know from the convexity of  $f$  (Theorem 2.2.10), that

$$f(x) \geq f(\bar{x}) + \nabla f(\bar{x})(x - \bar{x}), \quad \forall x \in U.$$

But, using the hypothesis and the convexity of  $M$  (Proposition 2.1.15),

$$-\nabla f(\bar{x}) \in N(M, \bar{x}) = \{u \in \mathbb{R}^p \mid \langle u, x - \bar{x} \rangle \leq 0, \quad \forall x \in M\},$$

whence  $\nabla f(\bar{x})(x - \bar{x}) \geq 0$  for every  $x \in M$ . From these relations we know  $f(x) \geq f(\bar{x})$  for every  $x \in M$ .  $\square$

Coming back to Theorems 3.1.18 and 3.1.21, in order to formulate sufficient optimality conditions, we strengthen the conclusion of these results. The good point is that we get stronger minimality concepts.

**Definition 3.1.26.** *Let  $\alpha > 0$ . One says that  $\bar{x} \in M$  is a strict local solution of order  $\alpha$  for (P), or a strict local minimum point of order  $\alpha$  for  $f$  on  $M$  if there exist two constants  $r, l > 0$  such that for every  $x \in M \cap B(\bar{x}, r)$ ,*

$$f(x) \geq f(\bar{x}) + l \|x - \bar{x}\|^\alpha.$$

The announced results are as follows.

**Theorem 3.1.27.** *Suppose that  $f$  is differentiable at  $\bar{x} \in M$  and*

$$\nabla f(\bar{x})(u) > 0, \quad \forall u \in T_B(M, \bar{x}) \setminus \{0\}.$$

*Then  $\bar{x}$  is a strict local solution of order  $\alpha = 1$  for (P).*

*Proof* Suppose, by way of contradiction, that  $\bar{x}$  is not a strictly local solution of order 1. Then, there exists a sequence  $(x_n) \rightarrow \bar{x}$ ,  $(x_n) \subset M$  such that for every  $n \in \mathbb{N}^*$ ,

$$f(x_n) < f(\bar{x}) + n^{-1} \|x_n - \bar{x}\|.$$

By virtue of this inequality,

$$x_n \neq \bar{x}, \forall n \in \mathbb{N}^*.$$

Since  $f$  is differentiable, there exists a sequence of real numbers  $(y_n) \rightarrow 0$  such that for every  $n \in \mathbb{N}$ ,

$$f(x_n) = f(\bar{x}) + \nabla f(\bar{x})(x_n - \bar{x}) + y_n \|x_n - \bar{x}\|.$$

The combination of these two relation yields

$$n^{-1} \|x_n - \bar{x}\| > \nabla f(\bar{x})(x_n - \bar{x}) + y_n \|x_n - \bar{x}\|,$$

whence, by division with  $\|x_n - \bar{x}\|$ , we deduce

$$n^{-1} > \nabla f(\bar{x}) \left( \frac{x_n - \bar{x}}{\|x_n - \bar{x}\|} \right) + y_n, \forall n \in \mathbb{N}^*. \tag{3.1.2}$$

Since the sequence  $\left( \frac{x_n - \bar{x}}{\|x_n - \bar{x}\|} \right)$  is bounded, there exists a convergent subsequence of it. The limit, denoted by  $u$ , of this subsequence is not zero (being of norm 1) and, furthermore, from  $\|x_n - \bar{x}\| \rightarrow 0$ , we infer that  $u \in T_B(M, \bar{x})$ . Consequently,  $u \in T_B(M, \bar{x}) \setminus \{0\}$ , and passing to the limit in the relation (3.1.2) we have

$$0 \geq \nabla f(\bar{x})(u),$$

which is in contradiction with the hypothesis. □

Notice that for differentiable functions, the concept of a local strict solution of order 1 is specific to the case of active restrictions (that is,  $\bar{x} \in M \setminus \text{int } M$ ): if  $f$  is differentiable at  $\bar{x} \in \text{int } M$ , then  $\bar{x}$  cannot be a local strict solution of order 1. Indeed, if  $\bar{x} \in \text{int } M$  would be local strict solution of order 1, then, on one hand,  $\nabla f(\bar{x}) = 0$  (Fermat Theorem), and, on the other hand,  $\nabla f(\bar{x}) \neq 0$  from the definition of strict solutions.

Concerning second-order optimality conditions, one has the following results.

**Theorem 3.1.28.** *Suppose that  $f$  is of class  $C^2$ ,  $\nabla f(\bar{x}) = 0$  and*

$$\nabla^2 f(\bar{x})(u, u) > 0, \forall u \in T_B(M, \bar{x}) \setminus \{0\}.$$

*Then  $\bar{x}$  is a local strict solution of order  $\alpha = 2$  for problem (P).*

*Proof* As before, one supposes, by contradiction, that the conclusion does not hold. Then there exists a sequence  $(x_n) \rightarrow \bar{x}$ ,  $(x_n) \subset M \setminus \{\bar{x}\}$  such that for every  $n \in \mathbb{N}^*$ ,

$$f(x_n) < f(\bar{x}) + n^{-1} \|x_n - \bar{x}\|^2.$$

From Taylor Theorem 1.3.4, for every  $n \in \mathbb{N}$  there exists  $c_n$  on the segment joining  $\bar{x}$  and  $x_n$  such that

$$\begin{aligned} f(x_n) - f(\bar{x}) &= \nabla f(\bar{x})(x_n - \bar{x}) + \frac{1}{2} \nabla^2 f(c_n)(x_n - \bar{x}, x_n - \bar{x}) \\ &= \frac{1}{2} \nabla^2 f(c_n)(x_n - \bar{x}, x_n - \bar{x}). \end{aligned}$$

We get

$$n^{-1} \|x_n - \bar{x}\|^2 > \frac{1}{2} \nabla^2 f(c_n)(x_n - \bar{x}, x_n - \bar{x}),$$

whence, in order to finish the proof, we divide by  $\|x_n - \bar{x}\|^2$  and we repeat the above arguments.  $\square$

In the unconstrained case, this result gives the following consequence.

**Corollary 3.1.29.** *Let  $U \subset \mathbb{R}^p$  be a nonempty, open set and  $f : U \rightarrow \mathbb{R}$  be a  $C^2$  function. If  $\bar{x} \in U$  is a critical point of  $f$  and  $\nabla^2 f(\bar{x})$  is positive definite (i.e.,  $\nabla^2 f(\bar{x})(u, u) > 0$  for every  $u \in \mathbb{R}^p \setminus \{0\}$ ), then  $\bar{x}$  is a local strict solution of order  $\alpha = 2$  for  $f$ .*

One can identify  $\nabla^2 f(\bar{x})$  with the Hessian matrix of  $f$  at  $\bar{x}$ ,  $\left(\frac{\partial^2 f}{\partial x^i \partial x^j}(\bar{x})\right)_{i,j \in \overline{1,p}}$ , and a sufficient condition for the positive definiteness of it is given by the next criterion, known from linear algebra (Sylvester criterion): all the determinants of the matrices  $\left(\frac{\partial^2 f}{\partial x^i \partial x^j}(\bar{x})\right)_{i,j \in \overline{1,k}}$ ,  $k \in \overline{1,p}$  are strictly positive. Analogously (for  $-f$ ), if the determinants of the matrices  $\left(\frac{\partial^2 f}{\partial x^i \partial x^j}(\bar{x})\right)_{i,j \in \overline{1,k}}$ ,  $k \in \overline{1,p}$  are not zero and change the signs starting with minus, then  $\nabla^2 f(\bar{x})$  is negative definite and  $\bar{x}$  is a maximum point. Furthermore, if all these determinants are not zero, then any other distribution of their signs leads to the conclusion that the reference point is not an extreme point.

## 3.2 Functional Restrictions

The restriction of the problem (P) introduced in the previous section is  $x \in M$ . Many times, in practice this set  $M$  of feasible points is defined by means of functions. Let us consider  $g : \mathbb{R}^p \rightarrow \mathbb{R}^n$  and  $h : \mathbb{R}^p \rightarrow \mathbb{R}^m$  as  $C^1$  functions. As usual,  $g$  and  $h$  can be thought of as  $g = (g_1, g_2, \dots, g_n)$ , and  $h = (h_1, h_2, \dots, h_m)$ , respectively, where  $g_i : \mathbb{R}^p \rightarrow \mathbb{R}$  ( $i \in \overline{1,n}$ ) and  $h_j : \mathbb{R}^p \rightarrow \mathbb{R}$  ( $j \in \overline{1,m}$ ) are  $C^1$  real valued functions.

Let the set of feasible points be defined as:

$$M := \{x \in U \mid g(x) \leq 0, h(x) = 0\} \subset \mathbb{R}^p.$$

Let us observe that we have two types of constraints: equalities and inequalities. Let  $x \in M$ . If for an  $i \in \overline{1,n}$ , one has that  $g_i(x) < 0$ , then the continuity of  $g$  ensures the existence of a neighborhood  $V$  of  $x$  such that  $g_i(y) < 0$  for all  $y \in V$ . Therefore,

when one looks for a certificate that  $x$  is a local solution of  $(P)$ , the restriction  $g_i \leq 0$  does not effectively influence the set of points  $u$  where one should compare  $f(x)$  and  $f(u)$ . For this reason, one says that the restriction  $g_i \leq 0$  is inactive at  $x$  and these kind of restrictions should be eliminated from the discussion. In the opposite case, when  $g_i(x) = 0$ , we call this active (inequality) restriction. For  $\bar{x} \in M$ , we denote the set of indexes corresponding to active inequality type restrictions by

$$A(\bar{x}) = \{i \in \overline{1, n} \mid g_i(\bar{x}) = 0\}.$$

We are now going to present two types of optimality conditions for problem  $(P)$  with functional constraints as described above. These two types of conditions are formally very close, but their differences are important for the detection of extreme points. We start with the Fritz John necessary optimality conditions where the objective function does not play any special role with respect to the functions which define the restrictions. We shall consider the drawbacks of these conditions, and next we shall impose supplementary conditions in order to eliminate them. By this procedure, we get the famous Karush-Kuhn-Tucker necessary optimality conditions which will be extensively used for solving nonlinear optimization problems.

### 3.2.1 Fritz John Optimality Conditions

The result of this subsection refers to necessary optimality conditions for problem  $(P)$  with functional restrictions without any additional assumption to the general framework already described. These conditions were obtained in 1948 by the German mathematician Fritz John.

**Theorem 3.2.1** (Fritz John). *Let  $\bar{x} \in M$  be a solution of  $(P)$ . Then there exist  $\lambda_0 \in \mathbb{R}$ ,  $\lambda_0 \geq 0$ ,  $\lambda = (\lambda_1, \lambda_2, \dots, \lambda_n) \in \mathbb{R}^n$ ,  $\mu = (\mu_1, \mu_2, \dots, \mu_m) \in \mathbb{R}^m$ , with  $\lambda_0 + \|\lambda\| + \|\mu\| \neq 0$  such that*

$$\lambda_0 \nabla f(\bar{x}) + \sum_{i=1}^n \lambda_i \nabla g_i(\bar{x}) + \sum_{j=1}^m \mu_j \nabla h_j(\bar{x}) = 0$$

and

$$\lambda_i \geq 0, \lambda_i g_i(\bar{x}) = 0, \text{ for every } i \in \overline{1, n}.$$

*Proof* Let us take  $\delta > 0$  such that  $D(\bar{x}, \delta) \subset U$  and for every  $x \in M \cap D(\bar{x}, \delta)$ ,  $f(\bar{x}) \leq f(x)$ . For all  $k \in \mathbb{N}^*$  we consider the function  $\varphi_k : D(\bar{x}, \delta) \rightarrow \mathbb{R}$  given by

$$\varphi_k(x) = f(x) + \frac{k}{2} \sum_{i=1}^n (g_i^+(x))^2 + \frac{k}{2} \sum_{j=1}^m (h_j(x))^2 + \frac{1}{2} \|x - \bar{x}\|^2,$$

where  $g_i^+(x) = \max\{g_i(x), 0\}$ . Clearly,  $\varphi_k$  attains its minimum on  $D(\bar{x}, \delta)$  and we denote by  $x_k$  such a minimum point. We also observe that

$$\begin{aligned} 0 \leq \varphi_k(x_k) &= f(x_k) + \frac{k}{2} \sum_{i=1}^n (g_i^+(x_k))^2 + \frac{k}{2} \sum_{j=1}^m (h_j(x_k))^2 + \frac{1}{2} \|x_k - \bar{x}\|^2 \\ &\leq \varphi_k(\bar{x}) = f(\bar{x}). \end{aligned}$$

Since the sequence  $(x_k)$  is bounded and  $f$  is continuous on  $D(\bar{x}, \delta)$ , we infer that  $(f(x_k))$  is also a bounded sequence. Letting  $k \rightarrow \infty$  in the above relation, we get

$$\begin{aligned} \lim_{k \rightarrow \infty} \sum_{i=1}^n (g_i^+(x_k))^2 &= 0 \\ \lim_{k \rightarrow \infty} \sum_{j=1}^m (h_j(x_k))^2 &= 0. \end{aligned}$$

The boundedness of  $(x_k)$  ensures that one can extract a convergent subsequence of it. Without relabeling, we can write  $x_k \rightarrow x^* \in D(\bar{x}, \delta)$ , and the previous relations yield  $x^* \in M$ . Consequently, passing to the limit in the inequality above, we have

$$f(x^*) + \frac{1}{2} \|x^* - \bar{x}\|^2 \leq f(\bar{x}).$$

On the other hand,  $f(\bar{x}) \leq f(x^*)$ , so  $\|x^* - \bar{x}\| = 0$ , that is  $x^* = \bar{x}$ . Therefore  $x_k \rightarrow \bar{x}$ .

An essential remark here is that  $\varphi_k$  is differentiable since the (nondifferentiable) scalar functions  $g_i^+(x)$  are squared, whence  $\nabla (g_i^+(x))^2 = g_i^+(x) \nabla g_i(x)$ . Since  $x_k$  is a minimum for  $\varphi_k$  on  $D(\bar{x}, \delta)$ , we deduce that

$$-\nabla \varphi_k(x_k) \in N(D(\bar{x}, \delta), x_k).$$

For  $k$  sufficiently large,  $x_k$  belongs to the interior of the ball  $D(\bar{x}, \delta)$  and we conclude that for these numbers  $k$ , one has  $N(D(\bar{x}, \delta), x_k) = \{0\}$ . The combination of these facts allow us to write

$$\nabla f(x_k) + k \sum_{i=1}^n g_i^+(x_k) \nabla g_i(x_k) + k \sum_{j=1}^m h_j(x_k) \nabla h_j(x_k) + x_k - \bar{x} = 0, \quad (3.2.1)$$

for every  $k$  large enough. For  $i \in \overline{1, n}$ ,  $j \in \overline{1, m}$ , we denote  $\alpha_i^k := k g_i^+(x_k)$ ,  $\beta_j^k := k h_j(x_k)$  and  $y^k = \sqrt{1 + \sum_{i=1}^n (\alpha_i^k)^2 + \sum_{j=1}^m (\beta_j^k)^2}$ . It is clear that  $y^k > 1$  and we take  $\lambda_0^k := \frac{1}{y^k}$ ,  $\lambda_i^k := \frac{\alpha_i^k}{y^k}$ ,  $\mu_j^k := \frac{\beta_j^k}{y^k}$ . We observe that

$$\left(\lambda_0^k\right)^2 + \sum_{i=1}^n \left(\lambda_i^k\right)^2 + \sum_{j=1}^m \left(\mu_j^k\right)^2 = 1,$$



whence the sequences  $(\lambda_0^k), (\lambda_i^k), (\mu_j^k) (i \in \overline{1, n}, j \in \overline{1, m})$  are bounded. Then there exist subsequences (we keep the indexes) convergent to some real numbers, respectively denoted by

$$\lambda_0, \lambda_1, \lambda_2, \dots, \lambda_n, \mu_1, \mu_2, \dots, \mu_m.$$

These numbers cannot be zero simultaneously. The positivity of the terms of the sequences  $(\lambda_0^k), (\lambda_i^k) (i \in \overline{1, n})$  implies the positivity of their limits  $\lambda_0, \lambda_1, \lambda_2, \dots, \lambda_n$ . Now, we divide relation (3.2.1) by  $y^k$ , and we get

$$\lambda_0^k \nabla f(x_k) + \sum_{i=1}^n \lambda_i^k \nabla g_i(x_k) + \sum_{j=1}^m \mu_j^k \nabla h_j(x_k) + \frac{1}{y^k} (x_k - \bar{x}) = 0.$$

Letting  $k \rightarrow \infty$  we have the first relation in the conclusion. Now we show the second one. Let  $i \in \overline{1, n}$ . If  $\lambda_i = 0$ , there is nothing to prove. Otherwise, if  $\lambda_i > 0$ , from the definition of  $\lambda_i$  we infer that for  $k$  sufficiently large,  $g_i^+(x_k) > 0$ , whence  $g_i^+(x_k) = g_i(x_k)$ . The relation

$$0 < g_i(x_k) \rightarrow g_i(\bar{x}) \leq 0$$

leads us to the conclusion  $g_i(\bar{x}) = 0$ . So, the second part of the conclusion holds and the theorem is completely proved.  $\square$

The relations in the conclusion of Theorem 3.2.1 are called Fritz John necessary optimality conditions. The major drawback of this result is that it does not eliminate the possibility that the real number associated to the objective function (i.e.,  $\lambda_0$ ) can be zero. This means that it would be possible to have too many points where the conditions in the conclusion are satisfied and therefore, in such a case, the result would not give important practical hints on the solutions. For instance, if a feasible point  $x$  satisfies  $\nabla g_i(x) = 0$  for a certain  $i \in A(x)$  or  $\nabla h_j(x) = 0$  for an  $j \in \overline{1, m}$ , then it satisfies Fritz John conditions (with  $\lambda_0 = 0$ ), the objective function being then completely eliminated. In the next subsection we shall impose a condition in order to avoid  $\lambda_0 = 0$ .

Let us first illustrate the possibilities created by Theorem 3.2.1 through two concrete examples.

**Example 3.2.2.** *Let us consider the problem of minimization of  $f : \mathbb{R}^2 \rightarrow \mathbb{R}$ ,*

$$f(x_1, x_2) = (x_1 - 3)^2 + (x_2 - 2)^2$$

*under the restriction  $g(x) \leq 0$ , where  $g : \mathbb{R}^2 \rightarrow \mathbb{R}^4$ ,*

$$g(x_1, x_2) = (x_1^2 + x_2^2 - 5, x_1 + 2x_2 - 4, -x_1, -x_2).$$

*One can observe graphically that  $\bar{x} = (2, 1)$  is solution of the problem and  $A(\bar{x}) = \{1, 2\}$ . We want to verify Fritz John condition at this point. From the second condition, since  $3, 4 \notin A(\bar{x})$ , we get  $\lambda_3 = \lambda_4 = 0$ . Since  $\nabla f(\bar{x}) = (-2, -2)$ ,  $\nabla g_1(\bar{x}) = (4, 2)$ ,  $\nabla g_2(\bar{x}) =$*

$(1, 2)$ , we have to find positive real numbers  $\lambda_0, \lambda_1, \lambda_2 \geq 0$ , not simultaneously zero, such that

$$\lambda_0(-2, -2) + \lambda_1(4, 2) + \lambda_2(1, 2) = (0, 0).$$

We get  $\lambda_1 = \frac{1}{3}\lambda_0$  and  $\lambda_2 = \frac{2}{3}\lambda_0$ , whence, by taking  $\lambda_0 > 0$ , the first Fritz John condition is fulfilled.

Let us now have a look to the point  $x = (0, 0)$ . This time  $A(x) = \{3, 4\}$ , whence  $\lambda_1 = \lambda_2 = 0$ . We have that  $\nabla f(\bar{x}) = (-6, -4)$ ,  $\nabla g_3(\bar{x}) = (-1, 0)$ ,  $\nabla g_4(\bar{x}) = (0, -1)$ . A computation shows that the equation

$$\lambda_0(-6, -4) + \lambda_3(-1, 0) + \lambda_4(0, -1) = (0, 0)$$

has no solution  $(\lambda_0, \lambda_3, \lambda_4)$  different to zero with positive components. Then  $x$  does not fulfill the Fritz John conditions, hence it is not a minimum point for the given problem.

**Example 3.2.3.** Let us consider the problem of minimization of  $f : (0, \infty) \times (0, \infty) \rightarrow \mathbb{R}$ ,  $f(x_1, x_2) = -2x_2$  under the restriction  $g(x) \leq 0$ , where  $g : \mathbb{R}^2 \rightarrow \mathbb{R}^3$ ,  $g(x_1, x_2) = (x_1 - x_2 - 2, -x_1 + x_2 + 2, x_1 + x_2 - 6)$ . The set  $M$  of feasible points is  $[(2, 0), (4, 2)] \setminus \{(2, 0)\}$ , and the minimum point is  $\bar{x} = (4, 2)$ . It is easy to observe that any feasible point satisfies the Fritz John conditions, but the solution is the only point where one can choose  $\lambda_0 \neq 0$ . Indeed, if  $x$  is a feasible point different to  $\bar{x}$ , then  $\lambda_3 = 0$ ,  $\lambda_1 = \lambda_2$  and  $\lambda_0 = 0$ .

### 3.2.2 Karush-Kuhn-Tucker Conditions

As seen before, it is desirable to have a Fritz John type result, but with  $\lambda_0 \neq 0$ . We could directly impose an extra condition in Theorem 3.2.1 in order to ensure this, but we prefer a direct approach because we aim at working with weak assumptions.

Let consider the sets

$$G(\bar{x}) = \left\{ \sum_{i \in A(\bar{x})} \lambda_i \nabla g_i(\bar{x}) + \sum_{j=1}^m \mu_j \nabla h_j(\bar{x}) \mid \lambda_i \geq 0, \forall i \in A(\bar{x}), \mu_j \in \mathbb{R}, \forall j \in \overline{1, m} \right\} \subset \mathbb{R}^p.$$

(where, as usual, we used the identification between  $L(\mathbb{R}^p, \mathbb{R})$  and  $\mathbb{R}^p$ ) and

$$D(\bar{x}) = \{u \in \mathbb{R}^p \mid \nabla g_i(\bar{x})(u) \leq 0, \forall i \in A(\bar{x}) \text{ and } \nabla h_j(\bar{x})(u) = 0, \forall j \in \overline{1, m}\}.$$

Before the main result, we need to shed some light on some important relations for these sets.

**Proposition 3.2.4.** For every  $\bar{x} \in M$  we have:

- (i)  $G(\bar{x}) = D(\bar{x})^-$ ;
- (ii)  $T_B(M, \bar{x}) \subset D(\bar{x})$ .

*Proof* (i) The inclusion  $G(\bar{x}) \subset D(\bar{x})^-$  is obvious, while the reverse one is a direct consequence of Farkas Lemma (Theorem 2.1.8).

(ii) Clearly,  $0 \in D(\bar{x})$ . Let  $u \in T_B(M, \bar{x}) \setminus \{0\}$ . By the definition of tangent vectors, there exist  $(t_n) \subset (0, \infty)$ ,  $t_n \rightarrow 0$  and  $(u_n) \rightarrow u$  such that for every  $n$ ,

$$\bar{x} + t_n u_n \in M.$$

The sequence  $(t_n u_n)$  converges towards 0 in  $\mathbb{R}^p$ . Taking into account the differentiability of  $h$  at  $\bar{x}$ , there exists  $(\alpha_n) \subset \mathbb{R}^p$ ,  $\alpha_n \rightarrow 0$  such that for every  $n \in \mathbb{N}$ ,

$$h(\bar{x} + t_n u_n) = h(\bar{x}) + t_n \nabla h(\bar{x})(u_n) + t_n \|u_n\| \alpha_n.$$

Since  $h(\bar{x} + t_n u_n) = h(\bar{x}) = 0$ , dividing by  $t_n$  and passing to the limit as  $n \rightarrow \infty$ , we get  $\nabla h(\bar{x})(u) = 0$ . Now, for every  $i \in A(\bar{x})$  there exist  $(\alpha_n^i) \subset \mathbb{R}$ ,  $\alpha_n^i \rightarrow 0$  such that for every  $n \in \mathbb{N}$ ,

$$g_i(\bar{x} + t_n u_n) = g_i(\bar{x}) + t_n \nabla g_i(\bar{x})(u_n) + t_n \|u_n\| \alpha_n^i.$$

As before, since  $g_i(\bar{x} + t_n u_n) \leq 0$  and  $g_i(\bar{x}) = 0$ , we have  $\nabla g_i(\bar{x})(u) \leq 0$ , and the proposition is proved.  $\square$

The next example shows that the reverse inclusion in the item (ii) above is false.

**Example 3.2.5.** Let  $g : \mathbb{R}^2 \rightarrow \mathbb{R}$ ,  $g(x_1, x_2) = -x_1 - x_2$ ,  $h : \mathbb{R}^2 \rightarrow \mathbb{R}$ ,  $h(x_1, x_2) = x_1 x_2$  and the feasible point  $\bar{x} = (0, 0)$ . Then:

$$\begin{aligned} D(\bar{x}) &= \{(u_1, u_2) \mid -u_1 - u_2 \leq 0\}, \\ T_B(M, \bar{x}) &= \{(u_1, u_2) \mid u_1 \geq 0, u_2 \geq 0, u_1 u_2 = 0\}. \end{aligned}$$

We establish now a generalized form of a classical result known under the name of Karush-Kuhn-Tucker Theorem, since it was obtained (with stronger assumptions) by the American mathematicians William Karush, Harold William Kuhn and Albert William Tucker. It is interesting to note that William Karush obtained the result in 1939, but the mathematical community become aware of its importance when Harold William Kuhn and Albert William Tucker got the result, in a different way, in 1950.

**Theorem 3.2.6** (Karush-Kuhn-Tucker). *Let  $\bar{x} \in M$  be a solution of the problem (P). Suppose that  $T_B(M, \bar{x})^- = D(\bar{x})^-$ . Then there exist  $\lambda = (\lambda_1, \lambda_2, \dots, \lambda_n) \in \mathbb{R}^n$ ,  $\mu = (\mu_1, \mu_2, \dots, \mu_m) \in \mathbb{R}^m$ , such that*

$$\nabla f(\bar{x}) + \sum_{i=1}^n \lambda_i \nabla g_i(\bar{x}) + \sum_{j=1}^m \mu_j \nabla h_j(\bar{x}) = 0 \tag{3.2.2}$$

and

$$\lambda_i \geq 0, \lambda_i g_i(\bar{x}) = 0, \text{ for every } i \in \overline{1, n}. \tag{3.2.3}$$

*Proof* From Theorem 3.1.18,  $\nabla f(\bar{x})(u) \geq 0$  for every  $u \in T_B(M, \bar{x})$ , whence  $-\nabla f(\bar{x}) \in T_B(M, \bar{x})^-$ . We use now the assumption  $T_B(M, \bar{x})^- = D(\bar{x})^-$  to infer that  $-\nabla f(\bar{x}) \in D(\bar{x})^-$ . From Proposition 3.2.4 (i), we get  $-\nabla f(\bar{x}) \in G(\bar{x})$ . Consequently, there exist  $\lambda_i \geq 0$ ,  $i \in A(\bar{x})$ ,  $\mu_j \in \mathbb{R}$ ,  $j \in \overline{1, m}$  such that  $-\nabla f(\bar{x}) = \sum_{i \in A(\bar{x})} \lambda_i \nabla g_i(\bar{x}) + \sum_{j=1}^m \mu_j \nabla h_j(\bar{x})$ . Now, for indexes  $i \in \overline{1, n} \setminus A(\bar{x})$  we take  $\lambda_i = 0$ , and we obtain the conclusion.  $\square$

If one compares Theorem 3.2.6 and Theorem 3.2.1, one notices the announced difference concerning the real number associated to the objective function.

The function  $L : U \times \mathbb{R}^{n+m} \rightarrow \mathbb{R}$ ,

$$L(x, (\lambda, \mu)) := f(x) + \sum_{i=1}^n \lambda_i g_i(x) + \sum_{j=1}^m \mu_j h_j(x)$$

is called the Lagrangian of  $(P)$ . Therefore, the conclusion given by relation (3.2.2) can be written as

$$\nabla_x L(\bar{x}, (\lambda, \mu)) = 0,$$

and the elements  $(\lambda, \mu) \in \mathbb{R}_+^n \times \mathbb{R}^m$  are called Lagrange multipliers. This name is due to the fact that the first time this method was used to investigate constrained optimization problems was given in some of Lagrange's works on calculus of variations problems.

The preceding theorem does not ensure the uniqueness of these multipliers. We denote by  $M(\bar{x})$  the set of Lagrange multipliers at  $\bar{x}$ , i.e.,

$$M(\bar{x}) := \{(\lambda, \mu) \in \mathbb{R}_+^n \times \mathbb{R}^m \mid \nabla_x L(\bar{x}, (\lambda, \mu)) = 0\},$$

where  $\mathbb{R}_+^n := [0, \infty)^n$ .

On the other hand,  $L(x, (\lambda, \mu))$  is an affine function with respect to the variables  $(\lambda, \mu)$ . We can observe the following fact which will appear later in the discussion: if  $\bar{x} \in M$  and  $(\bar{\lambda}, \bar{\mu})$  is a maximum on  $\mathbb{R}_+^n \times \mathbb{R}^m$  for  $(\lambda, \mu) \mapsto L(\bar{x}, (\lambda, \mu))$ , then  $\bar{\lambda}_i g_i(\bar{x}) = 0$  for every  $i \in \overline{1, n}$ .

Theorem 3.2.6 gives necessary optimality conditions for  $(P)$ . If, instead of minimization, we are looking for maximization of the objective function  $f$  under the same constraints, then, from  $\max f = -\min(-f)$ , the necessary condition (3.2.2) can be written as

$$-\nabla f(\bar{x}) + \sum_{i=1}^n \lambda_i \nabla g_i(\bar{x}) + \sum_{j=1}^m \mu_j \nabla h_j(\bar{x}) = 0.$$

Furthermore, let us notice that if one has only equalities as constraints, taking into account that  $h(x) = 0$  is equivalent to  $-h(x) = 0$ , the necessary optimality condition can be written, for both maxima and minima, as

$$\nabla f(\bar{x}) + \sum_{j=1}^m \mu_j \nabla h_j(\bar{x}) = 0.$$

Coming back to the main results, let us observe two more things. Firstly, if the problem has no restrictions (for instance,  $U = M = \mathbb{R}^p$ ), then relation (3.2.2) reduces to the first-order necessary optimality condition (Fermat Theorem):  $\nabla f(\bar{x}) = 0$ . Secondly, the key relation (3.2.2) does not hold without supplementary conditions (here,  $T_B(M, \bar{x})^- = D(\bar{x})^-$ ). To illustrate this consider the following example.

**Example 3.2.7.** Let  $f : \mathbb{R}^2 \rightarrow \mathbb{R}$  and  $g : \mathbb{R}^2 \rightarrow \mathbb{R}^2$  given by  $f(x_1, x_2) = x_1$  and  $g(x_1, x_2) = (-x_2 + (1 - x_1)^3, x_2)$ . It is easy to see that  $\bar{x} = (1, 0)$  is a minimum point of the associated problem, but (3.2.2) does not hold. Clearly, Fritz John conditions are fulfilled for  $\lambda_0 = 0$ .

So, in the next section, every condition which ensures the validity of the Karush-Kuhn-Tucker Theorem is called a qualification condition, and in view of the decisive importance of such requirements, we shall discuss it into detail in the next section.

Before that, let us observe that under certain assumptions, Karush-Kuhn-Tucker conditions (3.2.2) and (3.2.3) are also sufficient for minimality.

**Theorem 3.2.8.** Suppose that  $U$  is convex,  $f$  is convex on  $U$ ,  $h$  is affine and  $g_i, i \in \overline{1, n}$  are convex. Let  $\bar{x} \in M$ . If there exists  $(\lambda, \mu) \in \mathbb{R}^n \times \mathbb{R}^m$  such that (3.2.2) and (3.2.3) hold, then  $\bar{x}$  is a minimum point for (P) (or minimum of  $f$  on  $M$ ).

*Proof* The condition (3.2.2) expresses the fact that

$$\nabla_x L(\bar{x}, (\lambda, \mu)) = 0.$$

Under our assumptions,  $L$  is a convex function in  $x$ , so according to Theorem 3.1.22,  $\bar{x}$  is a minimum (without constraints) of the map  $x \mapsto L(x, (\lambda, \mu))$ . Therefore, for every  $x \in U$ ,

$$L(x, (\lambda, \mu)) = f(x) + \sum_{i=1}^n \lambda_i g_i(x) + \sum_{j=1}^m \mu_j h_j(x) \geq L(\bar{x}, (\lambda, \mu)) = f(\bar{x}).$$

But, for any  $x \in M$ ,

$$\sum_{i=1}^n \lambda_i g_i(x) + \sum_{j=1}^m \mu_j h_j(x) \leq 0,$$

whence  $f(x) \geq f(\bar{x})$ . The proof is complete. □

Concerning the structure of the set of Lagrange multipliers, we have the following result.

**Proposition 3.2.9.** For data with the structure mentioned in the above theorem, the set  $M(\bar{x})$  of the Lagrange multipliers is the same for all minimum points of  $f$  on  $M$ .

*Proof* Clearly,  $M$  is a convex set. Let  $x_1, x_2 \in M$  be two minimum points of  $(P)$ . According to Proposition 3.1.23, one has  $f(x_1) = f(x_2)$ . Let  $(\lambda, \mu) \in M(x_1)$ . Then

$$\nabla f(x_1) + \sum_{i=1}^n \lambda_i \nabla g_i(x_1) + \sum_{j=1}^m \mu_j \nabla h_j(x_1) = 0$$

and

$$\lambda_i \geq 0, \lambda_i g_i(x_1) = 0, \text{ for every } i \in \overline{1, n}.$$

As before,

$$f(x_2) + \sum_{i=1}^n \lambda_i g_i(x_2) \geq f(x_1) = f(x_2).$$

Taking into account the information on the numbers  $\lambda_i$  and  $g_i(x_2)$ , we infer that  $\lambda_i g_i(x_2) = 0$  for every  $i \in \overline{1, n}$ . From

$$L(x_2, (\lambda, \mu)) = f(x_2) = f(x_1) = L(x_1, (\lambda, \mu)),$$

we get that  $x_2$  is a minimum point for the convex function  $L(\cdot, (\lambda, \mu))$  on  $U$ . Hence

$$\nabla f(x_2) + \sum_{i=1}^n \lambda_i \nabla g_i(x_2) + \sum_{j=1}^m \mu_j \nabla h_j(x_2) = 0.$$

We have that  $(\lambda, \mu) \in M(x_2)$ . The other inclusion follows by exchanging  $x_1$  and  $x_2$  in the above proof.  $\square$

We now interpret Theorem 3.2.6 by using the concept of saddle point applied to the Lagrangian function. Firstly, we define the concept.

**Definition 3.2.10.** Let  $X, Y$  be two sets and  $F : X \times Y \rightarrow \mathbb{R}$ . A saddle point of  $F$  is a pair  $(\bar{x}, \bar{y}) \in X \times Y$  with the property that

$$\max_{y \in Y} F(\bar{x}, y) = F(\bar{x}, \bar{y}) = \min_{x \in X} F(x, \bar{y}). \quad (3.2.4)$$

It is clear that the relation (3.2.4) is equivalent to

$$F(\bar{x}, y) \leq F(\bar{x}, \bar{y}) \leq F(x, \bar{y}), \quad \forall (x, y) \in X \times Y$$

and to

$$F(\bar{x}, y) \leq F(x, \bar{y}), \quad \forall (x, y) \in X \times Y.$$

For instance, the point  $(0, 0)$  is a saddle point of  $F : \mathbb{R} \times \mathbb{R} \rightarrow \mathbb{R}$ ,  $F(x, y) = x^2 - y^2$  (the figure below).

The following general result is in order.

**Proposition 3.2.11.** For all saddle points  $(\bar{x}, \bar{y})$  of  $F$ , the value  $F(\bar{x}, \bar{y})$  is constant. If  $(x_1, y_1)$  and  $(x_2, y_2)$  are saddle points, then  $(x_1, y_2)$  and  $(x_2, y_1)$  are saddle points as well.

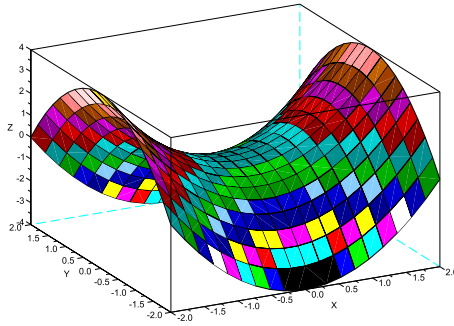


Figure 3.2: A saddle point.

*Proof* The following relations hold:

$$F(x_1, y) \leq F(x_1, y_1) \leq F(x, y_1), \quad \forall (x, y) \in X \times Y$$

$$F(x_2, y) \leq F(x_2, y_2) \leq F(x, y_2), \quad \forall (x, y) \in X \times Y.$$

If, in the first one, we take  $x = x_2$  and  $y = y_2$ , and in the second one we put  $x = x_1$  and  $y = y_1$ , we get  $F(x_1, y_1) = F(x_2, y_2) = F(x_2, y_1) = F(x_1, y_2)$ . Moreover, we can write for every  $(x, y) \in X \times Y$ ,

$$F(x_1, y) \leq F(x_1, y_2) \leq F(x, y_2),$$

whence  $(x_1, y_2)$  is a saddle point. For  $(x_2, y_1)$ , the proof is similar. □

For the general form of problem (P), we consider again the Lagrangian function  $L : U \times (\mathbb{R}_+^n \times \mathbb{R}^m) \rightarrow \mathbb{R}$ ,

$$L(x, (\lambda, \mu)) = f(x) + \sum_{i=1}^n \lambda_i g_i(x) + \sum_{j=1}^m \mu_j h_j(x).$$

**Theorem 3.2.12.** *An element  $(\bar{x}, (\bar{\lambda}, \bar{\mu})) \in U \times (\mathbb{R}_+^n \times \mathbb{R}^m)$  is a saddle point for the Lagrangian function  $L$  if and only if the following relations hold:*

- (i)  $\bar{x}$  is a minimum point for  $L(\cdot, (\bar{\lambda}, \bar{\mu}))$  on the open set  $U$ ;
- (ii)  $\bar{x} \in M$ ;
- (iii)  $\lambda_i g_i(\bar{x}) = 0$ , for every  $i \in \overline{1, n}$ .

*Proof* Let  $(\bar{x}, (\bar{\lambda}, \bar{\mu})) \in U \times (\mathbb{R}_+^n \times \mathbb{R}^m)$  be a saddle point for  $L$ . Then, according to the definition,

$$\max_{(\lambda, \mu) \in \mathbb{R}_+^n \times \mathbb{R}^m} L(\bar{x}, (\lambda, \mu)) = L(\bar{x}, (\bar{\lambda}, \bar{\mu})) = \min_{x \in U} L(x, (\bar{\lambda}, \bar{\mu})).$$

The second part of this relation is equivalent to (i). It remains to be shown that the first equality is equivalent to the combination of (ii) and (iii), and this is based on the fact that  $L$  is affine with respect to  $(\lambda, \mu)$  and, moreover, the particular form of  $\mathbb{R}_+^n \times \mathbb{R}^m$  allows us to easily compute the polar of its Bouligand tangent cone. According to Proposition 3.1.25,  $(\bar{\lambda}, \bar{\mu})$  with the property

$$\max_{(\lambda, \mu) \in \mathbb{R}_+^n \times \mathbb{R}^m} L(\bar{x}, (\lambda, \mu)) = L(\bar{x}, (\bar{\lambda}, \bar{\mu}))$$

is characterized by the relation

$$-\nabla_{(\lambda, \mu)} L(\bar{x}, (\bar{\lambda}, \bar{\mu})) \in N(\mathbb{R}_+^n \times \mathbb{R}^m, (\bar{\lambda}, \bar{\mu})).$$

It is not difficult to see that

$$N(\mathbb{R}_+^n \times \mathbb{R}^m, (\bar{\lambda}, \bar{\mu})) = \{u \in \mathbb{R}^n \mid u_i = 0 \text{ if } \bar{\lambda}_i > 0, u_i \leq 0 \text{ if } \bar{\lambda}_i = 0\} \times \{0\}_{\mathbb{R}^m}.$$

Then

$$\frac{\partial L}{\partial \lambda_i}(\bar{x}, (\bar{\lambda}, \bar{\mu})) = g_i(\bar{x}) : \begin{cases} = 0, & \text{if } \bar{\lambda}_i > 0 \\ \leq 0, & \text{if } \bar{\lambda}_i = 0 \end{cases}, \forall i = \overline{1, n},$$

and

$$\frac{\partial L}{\partial \mu_j}(\bar{x}, (\bar{\lambda}, \bar{\mu})) = h_j(\bar{x}) = 0, \forall j = \overline{1, m}.$$

The proof is complete.  $\square$

**Corollary 3.2.13.** *If  $(\bar{x}, (\bar{\lambda}, \bar{\mu})) \in U \times (\mathbb{R}_+^n \times \mathbb{R}^m)$  is a saddle point for the Lagrangian function  $L$ , then  $\bar{x}$  is a solution of (P).*

*Proof* The preceding result shows that  $\bar{x} \in M$  and

$$f(\bar{x}) = L(\bar{x}, (\bar{\lambda}, \bar{\mu})) \leq L(x, (\bar{\lambda}, \bar{\mu})), \forall x \in U.$$

Since for  $x \in M$ ,

$$L(x, (\bar{\lambda}, \bar{\mu})) \leq f(x),$$

we get  $f(\bar{x}) \leq f(x)$  for every  $x \in M$ .  $\square$

As usual, for convex data the converse holds as well.

**Theorem 3.2.14.** *Suppose that  $U$  is convex,  $f$  is convex on  $U$ ,  $h$  is affine and  $g_i, i \in \overline{1, n}$  are convex. The next relations are equivalent:*

- (i)  $(\bar{x}, (\bar{\lambda}, \bar{\mu})) \in U \times (\mathbb{R}_+^n \times \mathbb{R}^m)$  is a saddle point for the Lagrangian function  $L$ ;
- (ii)  $\bar{x}$  is a minimum point for (P) and  $(\bar{\lambda}, \bar{\mu})$  is a Lagrange multiplier.

*Proof* According to Theorem 3.2.12, relation (i) above is equivalent to all three relations in that result. One applies now Theorem 3.1.22 and the conclusion follows.  $\square$



### 3.2.3 Qualification Conditions

The qualification condition  $T_B(M, \bar{x})^- = D(\bar{x})^-$  imposed in Theorem 3.2.6 is called the Guignard condition at  $\bar{x}$  (after the name of the French mathematician Monique Guignard who proposed it back in 1969) and it is one of the weakest qualification conditions. The difficulty with this condition is that the effective calculations of the involved objects can be tricky in certain situation, and for this reason we want to investigate and to compare it with other qualification conditions as well. Clearly, relation  $T_B(M, \bar{x}) = D(\bar{x})$  is in turn a qualification condition (called the quasiregularity condition), since implies Guignard condition. As expected, the two conditions are not equivalent, as one can see from the next example (see also Example 2.1.7).

**Example 3.2.15.** Let  $g : \mathbb{R}^2 \rightarrow \mathbb{R}^2, g(x_1, x_2) = (-x_1, x_2)$  and  $h : \mathbb{R}^2 \rightarrow \mathbb{R}, h(x_1, x_2) = x_1x_2$ . Let us consider the feasible point  $\bar{x} = (0, 0)$ . Then:

$$D(\bar{x}) = \{(u_1, u_2) \mid u_1 \geq 0, u_2 \leq 0\},$$

$$T_B(M, \bar{x}) = \{(u_1, u_2) \mid u_1 \geq 0, u_2 \leq 0, u_1u_2 = 0\}$$

and

$$T_B(M, \bar{x})^- = D(\bar{x})^- = \{(u_1, u_2) \mid u_1 \leq 0, u_2 \geq 0\}.$$

The qualification conditions are linked to the reference point ( $\bar{x}$  in our notation). Every time when no confusion concerning the reference point could appear, we avoid, for simplicity, writing it explicitly.

Two of the most important (from a practical point of view) qualification conditions are listed below. The first one is called the linear independence qualification condition (at  $\bar{x}$ ) and is as follows:

the set  $\{\nabla g_i(\bar{x}) \mid i \in A(\bar{x})\} \cup \{\nabla h_j(\bar{x}) \mid j \in \overline{1, m}\}$  is linearly independent.

The second one is called Mangasarian-Fromovitz qualification condition (at  $\bar{x}$ ):

the set  $\{\nabla h_j(\bar{x}) \mid j \in \overline{1, m}\}$  is linearly independent and  
 $\exists u \in \mathbb{R}^p : \nabla h(\bar{x})(u) = 0$  and  $\nabla g_i(\bar{x})(u) < 0, \forall i \in A(\bar{x})$ .

(The American mathematicians Olvi Leon Mangasarian and Stanley Fromovitz published this condition in 1967.)

We will now establish the relations between these conditions and then show that they are indeed qualification conditions.

**Theorem 3.2.16.** *If the linear independence qualification condition at  $\bar{x} \in M$  holds, then the Mangasarian-Fromovitz qualification condition at  $\bar{x}$  is satisfied.*

*Proof* Without loss of generality, we suppose that  $A(\bar{x}) = \{1, \dots, q\}$ . Let  $T$  be the matrix of dimensions  $(q + m) \times p$  with the lines  $\nabla g_i(\bar{x}), i \in \overline{1, q}, \nabla h_j(\bar{x}), j \in \overline{1, m}$  and let  $b$  be the column vector with  $b_i = -1, i \in \overline{1, q}, b_j = 0, j \in \overline{q+1, q+m}$ . Since the lines of  $T$  are linearly independent, the system  $Td = b$  has a solution. If one denotes by  $u$  such a solution, then

$$\nabla g_i(\bar{x})(u^t) = -1, \forall i \in \overline{1, q} \text{ and } \nabla h_j(\bar{x})(u^t) = 0, \forall j \in \overline{1, m},$$

hence the Mangasarian-Fromovitz condition at  $\bar{x}$  is satisfied.  $\square$

The two conditions are not, however, equivalent.

**Example 3.2.17.** Let  $g_i : \mathbb{R}^2 \rightarrow \mathbb{R}, i \in \overline{1, 3}$  defined by:

$$g_1(x) = (x_1 - 1)^2 + (x_2 - 1)^2 - 2$$

$$g_2(x) = (x_1 - 1)^2 + (x_2 + 1)^2 - 2$$

$$g_3(x) = -x_1$$

and the feasible point  $\bar{x} = (0, 0)$ . Surely, the set

$$\{\nabla g_1(\bar{x}), \nabla g_2(\bar{x}), \nabla g_3(\bar{x}), i \in \overline{1, 3}\}$$

is not linearly independent since it consists of three elements in the two dimensional space  $\mathbb{R}^2$ . On the other hand, for  $u = (1, 0), \nabla g_i(\bar{x})(u) < 0$  for every  $i \in \overline{1, 3}$ .

Theorem 3.2.16 tells us that in order to show that the two conditions above are qualifications conditions, it is enough to show this only for Mangasarian-Fromovitz condition. This becomes obvious if one applies Theorem 3.2.1 and argues by contradiction. Suppose that  $\lambda_0 = 0$ . Then

$$\sum_{i \in A(\bar{x})} \lambda_i \nabla g_i(\bar{x}) + \sum_{j=1}^m \mu_j \nabla h_j(\bar{x}) = 0.$$

We multiply by the vector  $u$  from Mangasarian-Fromovitz condition and deduce that

$$\sum_{i \in A(\bar{x})} \lambda_i \langle \nabla g_i(\bar{x}), u \rangle = 0,$$

whence  $\lambda_i = 0$  for every  $i \in A(\bar{x})$ . Therefore,

$$\sum_{j=1}^m \mu_j \nabla h_j(\bar{x}) = 0,$$

and the linear independence of the gradients  $\{\nabla h_j(\bar{x}) \mid j \in \overline{1, m}\}$  implies that  $\mu_j = 0$  for every  $j \in \overline{1, m}$ . Putting together these remarks, we get the contradiction to  $|\lambda_0| + \|\lambda\| + \|\mu\| \neq 0$ . Consequently,  $\lambda_0 \neq 0$ .

In order to more precisely classify the qualification conditions introduced so far, we show that the Mangasarian-Fromovitz condition implies the quasiregularity condition.

We need an auxiliary result.

**Lemma 3.2.18.** *Let  $\varepsilon > 0$  and  $y : (-\varepsilon, \varepsilon) \rightarrow \mathbb{R}^p$  be a differentiable function such that  $y(0) = \bar{x}$ ,  $y'(0) = u \neq 0$ . Then there exists a sequence  $(x_k) \subset \text{Im } y \setminus \{\bar{x}\}$ ,  $(x_k) \rightarrow \bar{x}$  such that*

$$\frac{x_k - \bar{x}}{\|x_k - \bar{x}\|} \rightarrow \frac{u}{\|u\|}.$$

*Proof* We have

$$\lim_{t \rightarrow 0} \frac{y(t) - \bar{x}}{t} = \lim_{t \rightarrow 0} \frac{y(t) - y(0)}{t} = y'(0) = u \neq 0.$$

In particular, for  $t \neq 0$  sufficiently small one has  $y(t) \neq \bar{x}$ . We consider a sequence  $(t_k) \rightarrow 0$  of positive numbers and we define  $x_k = y(t_k)$ . Then,

$$\frac{x_k - \bar{x}}{\|x_k - \bar{x}\|} = \frac{x_k - \bar{x}}{t_k} \frac{t_k}{\|x_k - \bar{x}\|} \rightarrow \frac{u}{\|u\|}.$$

This ends the proof. □

**Theorem 3.2.19.** *If the Mangasarian-Fromovitz condition is satisfied at  $\bar{x} \in M$ , then  $T_B(M, \bar{x}) = D(\bar{x})$ .*

*Proof* As already observed, one inclusion is always true. We show only the opposite one, so we start with an element  $u \in D(\bar{x})$ . Denote by  $\bar{u} \in \mathbb{R}^p$  the vector given by the Mangasarian-Fromovitz condition. Let  $\lambda \in (0, 1)$  and  $d_\lambda := (1 - \lambda)u + \lambda\bar{u}$ . We show that  $d_\lambda \in T_B(M, \bar{x})$  for every  $\lambda \in (0, 1)$ , and then, taking  $\lambda \rightarrow 0$  and using the closedness of  $T_B(M, \bar{x})$  the conclusion will follow. Suppose that  $d_\lambda \neq 0$ , since otherwise, there is nothing to prove.

Let  $P$  be the operator defined by the matrix (of dimensions  $m \times p$ ) which has on the lines the vectors  $\nabla h_j(\bar{x})$ ,  $j \in \overline{1, m}$  of  $\mathbb{R}^p$ . These vectors are linearly independent and form a basis in the linear space  $\text{Im}(P)$ . Clearly, from the linear independence of  $\nabla h_j(\bar{x})$ ,  $j \in \overline{1, m}$  one deduces that  $m \leq p$ . But  $p = \dim(\text{Im}(P)) + \dim(\text{Ker}(P))$ , and we complete the above linear independent set up to a base of  $\mathbb{R}^p$  with a set of vectors  $\{v_1, v_2, \dots, v_{p-m}\}$ , and we denote by  $Z$  the matrix (of dimensions  $(p - m) \times p$ ) which has on the lines these vectors (which give a base in  $\text{Ker}(P)$ ). Then the square matrix  $\begin{pmatrix} P \\ Z \end{pmatrix}$  is nonsingular. We define  $\varphi : \mathbb{R}^{p+1} \rightarrow \mathbb{R}^p$  by

$$\varphi(x, \tau) = (h(x), (Z(x - \bar{x} - \tau d_\lambda)^t)^t).$$

Then  $\nabla_x \varphi(\bar{x}, 0) = \begin{pmatrix} P \\ Z \end{pmatrix}$  is a nonsingular matrix, whence, from Implicit Functions Theorem (Theorem 1.3.5), there exist  $\varepsilon > 0$  and a differentiable function  $y : (-\varepsilon, \varepsilon) \rightarrow$

$\mathbb{R}^p$  such that

$$\varphi(y(\tau), \tau) = 0,$$

for every  $\tau \in (-\varepsilon, \varepsilon)$ . Then

$$h(y(\tau)) = 0 \text{ and } Z(y(\tau) - \bar{x} - \tau d_\lambda)^t = 0. \quad (3.2.5)$$

At the same time, for every  $\tau \in (-\varepsilon, \varepsilon)$  and every  $x$  close enough to  $\bar{x}$  we have

$$\varphi(x, \tau) = 0 \Rightarrow x = y(\tau).$$

Since  $\varphi(\bar{x}, 0) = 0$ , we infer that  $y(0) = \bar{x}$ . According to the relations (3.2.5), we get, on one hand (by differentiation),

$$Py'(0) = 0,$$

and, on the other hand (by dividing with  $\tau \neq 0$  and passing to the limit),

$$Zy'(0)^t = Zd_\lambda^t.$$

Since  $u, \bar{u} \in D(\bar{x})$ , we get  $P(d_\lambda) = 0$ . We obtain

$$\begin{pmatrix} P \\ Z \end{pmatrix} (y'(0)) = \begin{pmatrix} P \\ Z \end{pmatrix} (d_\lambda),$$

that is  $d_\lambda = y'(0)$ . Using the above lemma, there exists a sequence  $(x_k) \subset \text{Im } y \setminus \{\bar{x}\}$ ,  $(x_k) \rightarrow \bar{x}$  with

$$\frac{x_k - \bar{x}}{\|x_k - \bar{x}\|} \rightarrow \frac{d_\lambda}{\|d_\lambda\|}.$$

Then  $h(x_k) = 0$ . In order to deduce that  $d_\lambda \in T_B(M, \bar{x})$ , it is sufficient to prove that, for  $k$  large enough,  $g(x_k) \leq 0$ . If  $i \notin A(\bar{x})$ , then  $g_i(\bar{x}) < 0$ , and the continuity of  $g_i$  implies that  $g_i(x_k) < 0$  for large  $k$ . If  $i \in A(\bar{x})$ ,  $\langle \nabla g_i(\bar{x}), u \rangle \leq 0$  and  $\langle \nabla g_i(\bar{x}), \bar{u} \rangle < 0$ , hence  $\langle \nabla g_i(\bar{x}), d_\lambda \rangle < 0$ . Since  $g_i$  is smooth (of class  $C^1$ ), there exists a sequence  $(\alpha_k) \rightarrow 0$  such that for every  $k \in \mathbb{N}$ ,

$$g_i(x_k) = g_i(\bar{x}) + \nabla g_i(\bar{x})(x_k - \bar{x}) + \alpha_k \|x_k - \bar{x}\|.$$

Therefore,

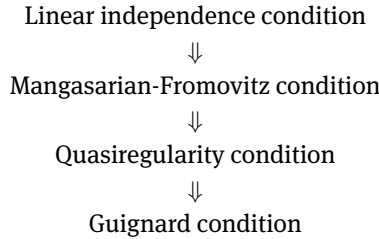
$$\frac{g_i(x_k)}{\|x_k - \bar{x}\|} = \frac{\nabla g_i(\bar{x})(x_k - \bar{x})}{\|x_k - \bar{x}\|} + \alpha_k \xrightarrow{k \rightarrow \infty} \nabla g_i(\bar{x}) \left( \frac{d_\lambda}{\|d_\lambda\|} \right) < 0.$$

Then  $g_i(x_k) < 0$  for sufficiently large  $k$ . Since there are a finite number of indexes  $i$ , we obtain the conclusion.  $\square$

In order to show that all four qualification conditions introduced are different, it remains to prove that the quasiregularity condition does not imply the Mangasarian-Fromovitz condition.

**Example 3.2.20.** Let  $g : \mathbb{R}^2 \rightarrow \mathbb{R}^2$ ,  $g(x_1, x_2) = (-x_1^2 + x_2, -x_1^2 - x_2)$  and the feasible point  $\bar{x} = (0, 0)$ . Then,  $D(\bar{x}) = \{(u_1, 0) \mid u_1 \in \mathbb{R}\}$ . On the other hand, is it easy to check that  $T_B(M, \bar{x}) \supset D(\bar{x})$  (whence the equality holds), but there is no  $\bar{u} \in \mathbb{R}^2$  with  $\nabla g(\bar{x})(\bar{u}) < 0$ .

We have shown the following implications :



and none of the converses hold.

**Remark 3.2.21.** Let us notice that, in particular, Theorem 3.2.19 shows as well that if  $h : \mathbb{R}^p \rightarrow \mathbb{R}$  is a  $C^1$  function, and  $\bar{x} \in \mathbb{R}^p$  has the property that  $\nabla h(\bar{x}) \neq 0$ , then the Bouligand tangent cone to the level curve  $\{x \in \mathbb{R}^p \mid h(x) = h(\bar{x})\}$  at  $\bar{x}$  is the hyperplane  $\{u \in \mathbb{R}^p \mid \nabla h(\bar{x})(u) = 0\}$  (or  $\text{Ker } \nabla h(\bar{x})$ ). Therefore,  $\nabla h(\bar{x})$  is a normal vector to this hyperplane. We recall here that the affine subspace (of  $\mathbb{R}^{p+1}$ ) tangent to the graph of  $h$  at  $(\bar{x}, h(\bar{x}))$  has the equation

$$y = h(\bar{x}) + \nabla h(\bar{x})(x - \bar{x}),$$

and a normal vector to it is  $(\nabla h(\bar{x}), -1)$ .

Therefore, the Mangasarian-Fromovitz condition ensures that the set  $M(\bar{x})$  is nonempty at  $\bar{x}$ , which is local minimum of the problem (P). Moreover, we will now show that this condition implies special properties of the set of Lagrange multipliers.

**Proposition 3.2.22.** If the Mangasarian-Fromovitz condition at  $\bar{x}$  holds, then  $M(\bar{x})$  is convex and compact (in  $\mathbb{R}^{n+m}$ ).

*Proof* According to the definition of  $M(\bar{x})$ , an element  $(\lambda, \mu) \in M(\bar{x})$  satisfies

$$\nabla f(\bar{x}) + \sum_{i=1}^n \lambda_i \nabla g_i(\bar{x}) + \sum_{j=1}^m \mu_j \nabla h_j(\bar{x}) = 0$$

and

$$\lambda_i \geq 0, \lambda_i g_i(\bar{x}) = 0, \text{ for every } i \in \overline{1, n}.$$

Therefore, checking the convexity and the closedness of  $M(\bar{x})$  is straightforward. We will now show that  $M(\bar{x})$  is bounded. Let, from the Mangasarian-Fromovitz condition,

$u \in \mathbb{R}^p$  such that

$$\nabla h(\bar{x})(u) = 0 \text{ and } \nabla g_i(\bar{x})(u) < 0, \forall i \in A(\bar{x}).$$

Then, for every  $(\lambda, \mu) \in M(\bar{x})$ ,

$$\nabla f(\bar{x})(u) + \sum_{i \in A(\bar{x})} \lambda_i \nabla g_i(\bar{x})(u) + \sum_{j=1}^m \mu_j \nabla h_j(\bar{x})(u) = 0,$$

whence

$$\sum_{i \in A(\bar{x})} \lambda_i (-\nabla g_i(\bar{x})(u)) = \nabla f(\bar{x})(u),$$

from where we deduce

$$\sum_{i \in A(\bar{x})} \lambda_i \min_{i \in A(\bar{x})} (-\nabla g_i(\bar{x})(u)) \leq \nabla f(\bar{x})(u),$$

so

$$\sum_{i \in A(\bar{x})} \lambda_i \leq \frac{\nabla f(\bar{x})(u)}{\min_{i \in A(\bar{x})} (-\nabla g_i(\bar{x})(u))}.$$

Since the right-hand side is constant, we deduce that the set of multipliers associated to inequalities constraints is bounded. Suppose, by contradiction, that there exists a sequence  $(\mu_k)_{k \in \mathbb{N}} \subset \mathbb{R}^m$  unbounded (without loss of generality, we can suppose that  $\|\mu_k\| \rightarrow \infty$ ) and a sequence  $(\lambda_k)_{k \in \mathbb{N}} \subset \mathbb{R}_+^n$  such that  $(\lambda_k, \mu_k) \in M(\bar{x})$ . Then, for every  $k \in \mathbb{N}$ ,

$$\nabla f(\bar{x}) + \sum_{i \in A(\bar{x})} (\lambda_i)_k \nabla g_i(\bar{x}) + \sum_{j=1}^m (\mu_j)_k \nabla h_j(\bar{x}) = 0.$$

We divide by  $\|\mu_k\|$  and we infer that

$$\|\mu_k\|^{-1} \nabla f(\bar{x}) + \sum_{i \in A(\bar{x})} \|\mu_k\|^{-1} (\lambda_i)_k \nabla g_i(\bar{x}) + \sum_{j=1}^m \|\mu_k\|^{-1} (\mu_j)_k \nabla h_j(\bar{x}) = 0. \quad (3.2.6)$$

From the previous step of the proof we have

$$\sum_{i \in A(\bar{x})} \|\mu_k\|^{-1} (\lambda_i)_k \nabla g_i(\bar{x}) \xrightarrow{k \rightarrow \infty} 0$$

and,

$$\|\mu_k\|^{-1} \nabla f(\bar{x}) \xrightarrow{k \rightarrow \infty} 0.$$

On the other hand, the sequence  $(\|\mu_k\|^{-1} \mu_k)$  is bounded (in  $\mathbb{R}^m$ ), whence, without relabeling, we can suppose that  $(\|\mu_k\|^{-1} \mu_k)$  is convergent towards a limit denoted by  $\bar{\mu} \in \mathbb{R}^m \setminus \{0\}$ . Passing to the limit in (3.2.6), we get

$$\sum_{j=1}^m \bar{\mu}_j \nabla h_j(\bar{x}) = 0.$$

Since  $\bar{\mu} \neq 0$ , this is in contradiction to the linear independence assumed in the Mangasarian-Fromovitz condition. Then  $M(\bar{x})$  is bounded.  $\square$

Let us discuss now two special cases of the problem data.

Firstly, we consider the situation where the inequality restrictions are convex functions, while the equality restriction is affine. The Slater condition takes place if  $h$  is affine,  $g_i$ ,  $i \in \overline{1, n}$  are convex and there exists  $u \in \mathbb{R}^p$  such that  $h(u) = 0$  and  $g(u) < 0$ . This condition was introduced in 1950 by the American mathematician Morton Slater.

**Theorem 3.2.23.** *The Slater condition implies  $T(M, x) = D(x)$  for every  $x \in M$  whence, in particular, is a qualification condition.*

*Proof* Let  $\bar{x} \in M$ . The inclusion  $T(M, \bar{x}) \subset D(\bar{x})$  is always true. Let  $v \in D(\bar{x})$ . By the Slater condition (using the convexity of  $g_i$ ) we deduce (by virtue of Theorem 2.2.10) that

$$0 > g_i(u) \geq g_i(\bar{x}) + \nabla g_i(\bar{x})(u - \bar{x}),$$

whence, for  $i \in A(\bar{x})$ ,  $\nabla g_i(\bar{x})(u - \bar{x}) < 0$ . We denote  $w := u - \bar{x}$ , and for  $\lambda \in (0, 1)$ , we define

$$w_\lambda := (1 - \lambda)v + \lambda w.$$

We show that  $w_\lambda \in T(M, \bar{x})$  for every  $\lambda \in (0, 1)$ . For  $i \in A(\bar{x})$ ,

$$\nabla g_i(\bar{x})(v) \leq 0, \quad \nabla g_i(\bar{x})(w) < 0,$$

hence  $\nabla g_i(\bar{x})(w_\lambda) < 0$ . By Taylor's Formula, there exists  $t > 0$  such that  $g_i(\bar{x} + tw_\lambda) < g_i(\bar{x}) = 0$  for every  $i \in A(\bar{x})$ . Let  $(t_k) \subset (0, \infty)$ ,  $t_k \rightarrow 0$ . Then

$$x_k := (1 - t_k)\bar{x} + t_k(\bar{x} + tw_\lambda) = \bar{x} + t_k tw_\lambda \xrightarrow{k \rightarrow \infty} \bar{x}.$$

In order for the conclusion to follow, we need to show that for  $k$  sufficiently large all  $(x_k)$  are in  $M$ . As usual, for  $i \notin A(\bar{x})$ , the continuity of  $g$  ensures this, while for  $i \in A(\bar{x})$ , we have

$$g_i(x_k) \leq (1 - t_k)g_i(\bar{x}) + t_k g_i(\bar{x} + tw_\lambda) < 0.$$

Since  $h$  is affine and  $h(\bar{x}) = 0$ , we get

$$h(x_k) = h(\bar{x} + t_k tw_\lambda) = t_k t \nabla h(\bar{x})(w_\lambda).$$

But  $v \in D(\bar{x})$ ,  $\nabla h(\bar{x})(v) = 0$ , so

$$\nabla h(\bar{x})(w_\lambda) = \lambda \nabla h(\bar{x})(w) = \lambda \nabla h(\bar{x})(u - \bar{x}) = \lambda h(u) = 0.$$

Therefore,  $h(x_k) = 0$  for any  $k$ , so, finally,  $(x_k)_{k \geq k_0} \subset M$ , and this means that  $w_\lambda \in T(M, \bar{x})$ . We let now  $\lambda \rightarrow 0$ ; since  $T(M, \bar{x})$  is closed, we get  $v \in T(M, \bar{x})$ , and the proof is complete.  $\square$

We consider now the case of affine restrictions. Take a matrix  $A$  of dimensions  $n \times p$ , a matrix  $B$  of dimensions  $m \times p$  and  $b \in \mathbb{R}^n, c \in \mathbb{R}^m$ . Therefore the set  $M$  become  $M = \{x \in \mathbb{R}^p \mid Ax^t \leq b^t, Bx^t = c^t\}$ , where the relationship " $\leq$ " is understood in the componentwise sense. Hence  $g(x) = (Ax^t - b^t)^t, h(x) = (Bx^t - c^t)^t$ .

**Theorem 3.2.24.** *In the above conditions and notation, the quasiregularity condition is automatically fulfilled.*

*Proof* As before, it is enough to prove that  $D(\bar{x}) \subset T(M, \bar{x})$ . Without loss of generality, one can suppose that  $A(\bar{x}) = \bar{1}, \bar{n}$ . Let  $v \in D(\bar{x})$ . Then  $Av^t \leq 0, Bv^t = 0$ . If  $v = 0$ , there is nothing to prove. Otherwise, we define

$$x_k := \bar{x} + \frac{1}{k}v, \forall k \in \mathbb{N}^*.$$

The relations

$$Ax_k^t \leq b^t, Bx_k^t = c^t, x_k \rightarrow \bar{x}$$

show that  $v \in T(M, \bar{x})$ . □

For affine restrictions, every minimum point of  $(P)$  satisfies the conclusions of Theorem 3.2.6.

### 3.3 Second-order Conditions

In this section we obtain second-order optimality conditions for the optimization problem with functional constraints and to this end, we assume that the data are  $C^2$  functions. Let  $\bar{x} \in M$  and  $(\bar{\lambda}, \bar{\mu}) \in \mathbb{R}^{n+m}$  be a vector which satisfies Karush-Kuhn-Tucker conditions (i.e., the conclusion of Theorem 3.2.6). We define the set of critical directions  $C(\bar{x}, (\bar{\lambda}, \bar{\mu}))$  as the set of vectors  $u \in \mathbb{R}^p$  for which

$$\begin{cases} \nabla g_i(\bar{x})(u) = 0, & \text{if } i \in A(\bar{x}) \text{ and } \bar{\lambda}_i > 0, \\ \nabla g_i(\bar{x})(u) \leq 0, & \text{if } i \in A(\bar{x}) \text{ and } \bar{\lambda}_i = 0, \\ \nabla h_j(\bar{x})(u) = 0, & \text{for every } j \in \bar{1}, \bar{m}. \end{cases}$$

Clearly,  $C(\bar{x}, (\bar{\lambda}, \bar{\mu}))$  is a cone.

**Remark 3.3.1.** *Obviously,  $C(\bar{x}, (\bar{\lambda}, \bar{\mu})) \subset D(\bar{x})$ . In particular, under quasiregularity qualification condition, i.e.,  $T_B(M, \bar{x}) = D(\bar{x})$ , one has the inclusion  $C(\bar{x}, (\bar{\lambda}, \bar{\mu})) \subset T_B(M, \bar{x})$ . Moreover, if one has only equalities constraints, then one has the equality, since  $\bar{\lambda}$  does not intervene in such a case.*

**Theorem 3.3.2.** *Let  $\bar{x} \in M$  be a solution of the problem  $(P)$  and  $(\bar{\lambda}, \bar{\mu}) \in \mathbb{R}^{n+m}$  a vector which satisfies Karush-Kuhn-Tucker conditions. If the linear independence condition*



holds at  $\bar{x}$ , then

$$\nabla_{xx}^2 L(\bar{x}, (\bar{\lambda}, \bar{\mu}))(u, u) \geq 0$$

for every  $u \in C(\bar{x}, (\bar{\lambda}, \bar{\mu}))$ .

*Proof* Without loss of generality, we suppose that all the inequality constraints are active. We split the proof into several steps.

At the first step, we repeat, with some modifications, several arguments from the proof of Theorem 3.2.19 in order to get a sequence of feasible points with special properties. Let  $d \in D(\bar{x})$ , and let  $P$  be the operator defined by the matrix (of dimensions  $(n+m) \times p$ ) with the lines consisting of vectors  $\nabla g_i(\bar{x})$ ,  $i \in \overline{1, n}$ ,  $\nabla h_j(\bar{x})$ ,  $j \in \overline{1, m}$  in  $\mathbb{R}^p$ . These vectors are linearly independent and form a basis in the linear subspace  $\text{Im}(P)$ . Let us denote by  $Z$  the matrix (of dimensions  $(p - (n+m)) \times p$ ) whose lines are some vectors that form a basis in  $\text{Ker}(P)$ . The the square matrix  $\begin{pmatrix} P \\ Z \end{pmatrix}$  is nonsingular. We define  $\varphi : \mathbb{R}^{p+1} \rightarrow \mathbb{R}^p$  by

$$\varphi(x, \tau) = ((g(x), h(x)) - \tau(Pd)^t, (Z(x - \bar{x} - \tau d)^t)^t).$$

Then  $\nabla_x \varphi(\bar{x}, 0) = \begin{pmatrix} P \\ Z \end{pmatrix}$  is nonsingular and, from Implicit Function Theorem (i.e., Theorem 1.3.5), there exists  $\varepsilon > 0$  and a differentiable function  $y : (-\varepsilon, \varepsilon) \rightarrow \mathbb{R}^p$  such that

$$\varphi(y(\tau), \tau) = 0,$$

for every  $\tau \in (-\varepsilon, \varepsilon)$ . Moreover, for every  $\tau \in (-\varepsilon, \varepsilon)$  and every  $x$  close enough to  $\bar{x}$ ,

$$\varphi(x, \tau) = 0 \Rightarrow x = y(\tau).$$

Let  $(t_k) \subset (0, \infty)$ ,  $(t_k) \rightarrow 0$ . Then, using the fact that  $\varphi(y(t_k), t_k) = 0$ , there exists, for every  $k$  large enough,  $z_k = y(t_k)$  such that

$$\begin{aligned} g_i(z_k) &= t_k \nabla g_i(\bar{x})(d) \leq 0, \quad \forall i \in \overline{1, n} \\ h_j(z_k) &= t_k \nabla h_j(\bar{x})(d) = 0, \quad \forall j \in \overline{1, m}. \end{aligned} \tag{3.3.1}$$

Therefore,  $(z_k) \subset M$  and the sequence  $(t_k^{-1}(z_k - \bar{x}))$  is convergent. We show that  $t_k^{-1}(z_k - \bar{x}) \rightarrow d$ , which would imply  $d \in T_B(M, \bar{x})$ . According to the Taylor Theorem, from  $\varphi(z_k, t_k) = 0$ , there exists  $(\mu_k) \subset \mathbb{R}^{n+m}$ ,  $\mu_k \rightarrow 0$  such that

$$0 = \left( (P(z_k - \bar{x} - t_k d)^t)^t, (Z(z_k - \bar{x} - t_k d)^t)^t \right) + \|z_k - \bar{x}\| \mu_k,$$

whence

$$\left( \frac{z_k - \bar{x}}{t_k} - d \right)^t = \begin{pmatrix} P \\ Z \end{pmatrix}^{-1} (-t_k^{-1} \|z_k - \bar{x}\| \mu_k),$$

from where, after passing to the limit, one gets the announced relation.

Let now  $u \in C(\bar{x}, (\bar{\lambda}, \bar{\mu})) \subset D(\bar{x})$ . We use now the above construction of the sequence  $(z_k) \rightarrow \bar{x}$  corresponding to  $u$ . We have

$$L(z_k, (\bar{\lambda}, \bar{\mu})) = f(z_k) + \sum_{i=1}^n \bar{\lambda}_i g_i(z_k) + \sum_{j=1}^m \bar{\mu}_j h_j(z_k) = f(z_k) - t_k \sum_{i \in A(\bar{x})} \bar{\lambda}_i \nabla g_i(\bar{x})(u) = f(z_k).$$

From Taylor second-order condition, there exists  $(y_k) \rightarrow 0$  such that for every  $k$ ,

$$\begin{aligned} L(z_k, (\bar{\lambda}, \bar{\mu})) &= L(\bar{x}, (\bar{\lambda}, \bar{\mu})) + \nabla_x L(\bar{x}, (\bar{\lambda}, \bar{\mu}))(z_k - \bar{x}) \\ &\quad + \frac{1}{2} \nabla_{xx} L(\bar{x}, (\bar{\lambda}, \bar{\mu}))(z_k - \bar{x}, z_k - \bar{x}) + y_k \|z_k - \bar{x}\|^2. \end{aligned}$$

But, from the Karush-Kuhn-Tucker conditions,  $L(\bar{x}, (\bar{\lambda}, \bar{\mu})) = f(\bar{x})$  and  $\nabla_x L(\bar{x}, (\bar{\lambda}, \bar{\mu})) = 0$ , whence

$$f(z_k) = f(\bar{x}) + \frac{1}{2} \nabla_{xx} L(\bar{x}, (\bar{\lambda}, \bar{\mu}))(z_k - \bar{x}, z_k - \bar{x}) + y_k \|z_k - \bar{x}\|^2.$$

Since  $z_k \rightarrow \bar{x}$  and  $\bar{x}$  is a solution for the problem (P), we obtain  $f(z_k) - f(\bar{x}) \geq 0$  for every sufficiently large  $k$ . Then

$$\frac{1}{2} \nabla_{xx} L(\bar{x}, (\bar{\lambda}, \bar{\mu}))(z_k - \bar{x}, z_k - \bar{x}) + y_k \|z_k - \bar{x}\|^2 \geq 0.$$

We divide by  $t_k^2$  and we pass to the limit in order to get

$$\nabla_{xx} L(\bar{x}, (\bar{\lambda}, \bar{\mu}))(u, u) \geq 0.$$

The proof is complete. □

We formulate now a converse of the previous result. As shown before, the sufficient optimality condition returns a stronger type of solution (i.e., strict solution).

**Theorem 3.3.3.** *Let  $\bar{x} \in M$  and  $(\bar{\lambda}, \bar{\mu}) \in \mathbb{R}^{n+m}$  a vector which satisfies Karush-Kuhn-Tucker conditions. Suppose that*

$$\nabla_{xx}^2 L(\bar{x}, (\bar{\lambda}, \bar{\mu}))(u, u) > 0$$

*for every  $u \in C(\bar{x}, (\bar{\lambda}, \bar{\mu})) \setminus \{0\}$ . Then  $\bar{x}$  is a local strict solution of second order of (P).*

*Proof* Since the set  $C(\bar{x}, (\bar{\lambda}, \bar{\mu})) \cap \{u \in \mathbb{R}^p \mid \|u\| = 1\}$  is compact, and  $C(\bar{x}, (\bar{\lambda}, \bar{\mu}))$  is a cone, the relation  $\nabla_{xx}^2 L(\bar{x}, (\bar{\lambda}, \bar{\mu}))(u, u) > 0$  for every  $u \in C(\bar{x}, (\bar{\lambda}, \bar{\mu})) \setminus \{0\}$  is equivalent to the existence of a strictly positive number  $\rho$  with the property

$$\nabla_{xx}^2 L(\bar{x}, (\bar{\lambda}, \bar{\mu}))(u, u) \geq \rho \|u\|^2, \forall u \in C(\bar{x}, (\bar{\lambda}, \bar{\mu})).$$

Suppose that there exists  $(z_k) \rightarrow \bar{x}$ ,  $(z_k) \subset M$  such that

$$f(z_k) < f(\bar{x}) + k^{-1} \|z_k - \bar{x}\|^2,$$

for every  $k$  sufficiently large. Then, without loss of generality, we suppose that  $\|z_k - \bar{x}\|^{-1}(z_k - \bar{x}) \rightarrow d \in T_B(M, \bar{x}) \setminus \{0\} \subset D(\bar{x})$ . On the other hand,

$$L(z_k, (\bar{\lambda}, \bar{\mu})) = f(z_k) + \sum_{i \in A(\bar{x})} \bar{\lambda}_i g_i(z_k) \leq f(z_k),$$

and, as before, there exists  $(y_k) \rightarrow 0$  such that for every  $k$ ,

$$L(z_k, (\bar{\lambda}, \bar{\mu})) = f(\bar{x}) + \frac{1}{2} \nabla_{xx} L(\bar{x}, (\bar{\lambda}, \bar{\mu}))(z_k - \bar{x}, z_k - \bar{x}) + y_k \|z_k - \bar{x}\|^2. \quad (3.3.2)$$

Suppose that  $d \notin C(\bar{x}, (\bar{\lambda}, \bar{\mu}))$ . Then there exists  $i_0 \in A(\bar{x})$  with  $\bar{\lambda}_{i_0} \nabla g_{i_0}(\bar{x})(d) < 0$ . For the other indices  $i \in A(\bar{x})$  we have  $\bar{\lambda}_i \nabla g_i(\bar{x})(d) \leq 0$ . Then, it exists  $(\tau_k) \rightarrow 0$  such that for every  $k$ ,

$$\begin{aligned} \bar{\lambda}_{i_0} g_{i_0}(z_k) &= \bar{\lambda}_{i_0} g_{i_0}(\bar{x}) + \bar{\lambda}_{i_0} \nabla g_{i_0}(\bar{x})(z_k - \bar{x}) + \tau_k \bar{\lambda}_{i_0} \|z_k - \bar{x}\| \\ &= \|z_k - \bar{x}\| \bar{\lambda}_{i_0} \nabla g_{i_0}(\bar{x}) \left( \frac{z_k - \bar{x}}{\|z_k - \bar{x}\|} \right) + \tau_k \bar{\lambda}_{i_0} \|z_k - \bar{x}\|. \end{aligned}$$

Hence

$$\begin{aligned} L(z_k, (\bar{\lambda}, \bar{\mu})) &= f(z_k) + \sum_{i \in A(\bar{x})} \bar{\lambda}_i g_i(z_k) \leq f(z_k) + \bar{\lambda}_{i_0} g_{i_0}(z_k) \\ &= f(z_k) + \|z_k - \bar{x}\| \bar{\lambda}_{i_0} \nabla g_{i_0}(\bar{x}) \left( \frac{z_k - \bar{x}}{\|z_k - \bar{x}\|} \right) + \tau_k \bar{\lambda}_{i_0} \|z_k - \bar{x}\|. \end{aligned}$$

From (3.3.2), we get

$$\begin{aligned} f(\bar{x}) + \frac{1}{2} \nabla_{xx} L(\bar{x}, (\bar{\lambda}, \bar{\mu}))(z_k - \bar{x}, z_k - \bar{x}) + y_k \|z_k - \bar{x}\|^2 \\ \leq f(z_k) + \|z_k - \bar{x}\| \bar{\lambda}_{i_0} \nabla g_{i_0}(\bar{x}) \left( \frac{z_k - \bar{x}}{\|z_k - \bar{x}\|} \right) + \tau_k \bar{\lambda}_{i_0} \|z_k - \bar{x}\|. \end{aligned}$$

Furthermore,

$$\begin{aligned} \lim_k \|z_k - \bar{x}\|^{-1} \nabla_{xx} L(\bar{x}, (\bar{\lambda}, \bar{\mu}))(z_k - \bar{x}, z_k - \bar{x}) \\ = \lim_k \|z_k - \bar{x}\| \nabla_{xx} L(\bar{x}, (\bar{\lambda}, \bar{\mu})) \left( \frac{z_k - \bar{x}}{\|z_k - \bar{x}\|}, \frac{z_k - \bar{x}}{\|z_k - \bar{x}\|} \right) = 0. \end{aligned}$$

After relabeling, one can see that there exists  $(v_k) \rightarrow 0$  such that

$$f(z_k) \geq f(\bar{x}) - \|z_k - \bar{x}\| \bar{\lambda}_{i_0} \nabla g_{i_0}(\bar{x}) \left( \frac{z_k - \bar{x}}{\|z_k - \bar{x}\|} \right) + v_k \|z_k - \bar{x}\|.$$

From the assumption made,  $f(z_k) < f(\bar{x}) + k^{-1} \|z_k - \bar{x}\|^2$ , whence

$$f(\bar{x}) + k^{-1} \|z_k - \bar{x}\|^2 \geq f(\bar{x}) - \|z_k - \bar{x}\| \bar{\lambda}_{i_0} \nabla g_{i_0}(\bar{x}) \left( \frac{z_k - \bar{x}}{\|z_k - \bar{x}\|} \right) + v_k \|z_k - \bar{x}\|,$$

that is

$$k^{-1} \|z_k - \bar{x}\| \geq -\bar{\lambda}_{i_0} \nabla g_{i_0}(\bar{x}) \left( \frac{z_k - \bar{x}}{\|z_k - \bar{x}\|} \right) + \nu_k.$$

Passing to the limit, we arrive at a contradiction to the relation  $\bar{\lambda}_{i_0} \nabla g_{i_0}(\bar{x})(d) < 0$ . Consequently,  $d \in C(\bar{x}, (\bar{\lambda}, \bar{\mu})) \setminus \{0\}$ , whence  $\nabla_{xx}^2 L(\bar{x}, (\bar{\lambda}, \bar{\mu}))(d, d) \geq \rho$ . Since  $L(z_k, (\bar{\lambda}, \bar{\mu})) \leq f(z_k)$ , coming back to (3.3.2), we can write

$$f(z_k) \geq f(\bar{x}) + \frac{1}{2} \nabla_{xx} L(\bar{x}, (\bar{\lambda}, \bar{\mu}))(z_k - \bar{x}, z_k - \bar{x}) + y_k \|z_k - \bar{x}\|^2.$$

But  $\nabla_{xx}^2 L(\bar{x}, (\bar{\lambda}, \bar{\mu}))$  is continuous, whence for  $k$  large enough,

$$\nabla_{xx} L(\bar{x}, (\bar{\lambda}, \bar{\mu}))(z_k - \bar{x}, z_k - \bar{x}) > 2^{-1} \rho \|z_k - \bar{x}\|^2.$$

Finally,

$$f(\bar{x}) + k^{-1} \|z_k - \bar{x}\|^2 \geq f(\bar{x}) + \frac{\rho}{4} \|z_k - \bar{x}\|^2 + y_k \|z_k - \bar{x}\|^2,$$

and a new contradiction occurs. □

### 3.4 Motivations for Scientific Computations

In this section, we examine the computational limits of the theoretical results from the previous sections. In some cases, it is not possible to get the exact solution of the optimization problems. The theory leads us to solve some nonlinear (systems of) equations, which do not admit analytical expressions for the solutions. This motivates us to subsequently study numerical algorithms for solving such equations, and this will be done in Chapter 6.

**1. (Least squares method)** We discuss now a special case of an optimization problem without restrictions. This problem belongs to the general approach called the method of least squares which is designed for the interpretation of numerical data issued from experiments in physics, biology, astronomy, chemistry. From historical perspective, this kind of problem arose from the study of the movements of planets and in questions linked to navigation techniques. The mathematician who founded this method is considered to be Gauss, but the method was published for the first time by Legendre.

In few words, the method of least squares refers to the following situation: we dispose of a data set  $v_1, v_2, \dots, v_N$  obtained after some measurements made at the (different) moments  $t_1, t_2, \dots, t_N$ . The objective is to determine the best model function of the form  $t \mapsto \varphi(t, x)$  (where  $x = (x_1, x_2, \dots, x_k)$  are parameters to optimize) which fits with the measurements. Therefore, for every  $i \in \overline{1, N}$ , one defines the residual at the moment  $t_i$  as the absolute difference between the measurement  $v_i$  and the value of the model at same time:

$$r_i := |v_i - \varphi(t_i, x)|,$$

and now the problem is to minimize the function

$$f(x) = \sum_{i=1}^N r_i^2 = \sum_{i=1}^N [v_i - \varphi(t_i, x)]^2.$$

It should be said that another possible objective function (even more natural to be considered) would be

$$\sum_{i=1}^N |v_i - \varphi(t_i, x)|$$

but this construction does not preserve the differentiability. For this reason, one prefers the sum of the squares of the residuals, whence the name of the method.

Let us now consider the simplest case of a linear dependence. Let us suppose that one has made  $N$  measurements at the different moments of time  $t_1, t_2, \dots, t_N > 0$  and, correspondingly, one has the values  $v_1, v_2, \dots, v_N$ . We know that the dependence between these two sets of data is linear, and we are interested in obtaining a line which better fits the collection of observations. Let be a line  $v = at + b$ . As above, the residual at the moment  $t_i$  is  $|v_i - (at_i + b)|$  and in order to “measure” the sum of these residuals, we consider the function  $f : \mathbb{R}^2 \rightarrow \mathbb{R}$ ,

$$f(a, b) = \sum_{i=1}^N [v_i - (at_i + b)]^2.$$

The line with respect to which this sum of residuals will be the smallest, will be that which we seek. Then, we arrive at the problem of minimization (without restrictions) of the function  $f$ . We compute the partial derivatives of  $f$ :

$$\begin{aligned} \frac{\partial f}{\partial a}(a, b) &= \sum_{i=1}^N 2(-t_i) [v_i - (at_i + b)] \\ \frac{\partial f}{\partial b}(a, b) &= \sum_{i=1}^N -2 [v_i - (at_i + b)], \end{aligned}$$

and the calculus of critical points is reduced to computation of the solutions of the system:

$$\begin{cases} \left( \sum_{i=1}^N t_i^2 \right) a + \left( \sum_{i=1}^N t_i \right) b = \sum_{i=1}^N t_i v_i \\ \left( \sum_{i=1}^N t_i \right) a + Nb = \sum_{i=1}^N v_i. \end{cases}$$

The determinant of this system is

$$\Delta := N \left( \sum_{i=1}^N t_i^2 \right) - \left( \sum_{i=1}^N t_i \right)^2 = \left( \sum_{i=1}^N 1^2 \right) \left( \sum_{i=1}^N t_i^2 \right) - \left( \sum_{i=1}^N t_i \right)^2.$$

From the Hölder inequality, this number is positive (equality would be possible only if all the values  $t_i$  are equal, but this is not possible). Then the system admits a unique solution:

$$\begin{pmatrix} a \\ b \end{pmatrix} = \begin{pmatrix} \sum_{i=1}^N t_i^2 & \sum_{i=1}^N t_i \\ \sum_{i=1}^N t_i & N \end{pmatrix}^{-1} \begin{pmatrix} \sum_{i=1}^N t_i v_i \\ \sum_{i=1}^N v_i \end{pmatrix}.$$

Let us observe that, furthermore, the function is coercive since  $\lim_{\|(a,b)\| \rightarrow \infty} f(a,b) = \infty$ , hence, according to Theorem 3.1.7, it admits a minimum point which is necessarily the critical point determined above. Therefore this pair  $(a, b)$  is the solution of the problem.

Another remark is that an important part of the above calculations can be repeated with obvious changes if one supposes a dependence of the type  $v = a \cdot p(t) + b \cdot q(t)$ , where  $p, q : \mathbb{R} \rightarrow \mathbb{R}$ . One obtains the system

$$\begin{pmatrix} \sum_{i=1}^N p^2(t_i) & \sum_{i=1}^N p(t_i)q(t_i) \\ \sum_{i=1}^N p(t_i)q(t_i) & \sum_{i=1}^N q^2(t_i) \end{pmatrix} \begin{pmatrix} a \\ b \end{pmatrix} = \begin{pmatrix} \sum_{i=1}^N p(t_i)v_i \\ \sum_{i=1}^N q(t_i)v_i \end{pmatrix}.$$

Again, the Hölder inequality ensures that the associated matrix is invertible if and only if  $(p^2(t_i))_{i=1, \dots, N}$  and  $(q^2(t_i))_{i=1, \dots, N}$  are not proportional.

In general, for more complicated models (nonlinear in  $x$ ) the method of least squares does not have an easily computable solution and this will be an impetus for us to study several algorithm in order to get good approximation of solution in fast computational time.

**2. (The projection on a closed convex set)** Let  $a_1, a_2, \dots, a_p \in (0, \infty)$  and the set (generalized ellipsoid)

$$M := \left\{ x \in \mathbb{R}^p \mid \sum_{i=1}^p \left( \frac{x_i}{a_i} \right)^2 \leq 1 \right\}.$$

Obviously, this is a convex and compact set. Let  $v \notin M$ . From Theorem 2.1.5, there exists  $\bar{v} \in M$ , the projection of  $v$  on  $M$ . Again, we want to find an expression of this element.

As before,  $\bar{v}$  is the solution of the minimization problem of  $f(x) = \|x - v\|^2$  under the restriction  $x \in M$ . If it would exist a solution  $\bar{x} \in \text{int } M$ , then  $\nabla f(\bar{x}) = 0$ , whence  $\bar{x} - v = 0$ , that is  $v \in M$ , which is false. Therefore, the restriction is active in  $\bar{v}$ , that is  $\sum_{i=1}^p \left( \frac{\bar{v}_i}{a_i} \right)^2 = 1$ . Moreover, the function which defines (with inequality) the constraint  $x \in M$ , i.e.,  $g : \mathbb{R}^p \rightarrow \mathbb{R}$ ,

$$g(x) := \sum_{i=1}^p \left( \frac{x_i}{a_i} \right)^2 - 1,$$

is convex, and the Slater condition holds. Moreover,  $f$  is convex as well, so we can conclude that  $\bar{v}$  is a solution of the problem if and only of there exists  $\lambda \geq 0$  such that

for all indexes  $i$ ,

$$\bar{v}_i - v_i + \lambda \frac{\bar{v}_i}{a_i^2} = 0.$$

Since  $v \neq \bar{v}$ ,  $\lambda > 0$  and

$$\bar{v}_i = \frac{a_i^2 v_i}{a_i^2 + \lambda}.$$

On the other hand, by  $\sum_{i=1}^p \left(\frac{\bar{v}_i}{a_i}\right)^2 = 1$ ,

$$\sum_{i=1}^p \frac{a_i^2 v_i^2}{(a_i^2 + \lambda)^2} = 1,$$

so finding  $\lambda$  (and then  $\bar{v}$ ) requires solving the above equation. Let us remark that the equation has a unique solution, since the mapping

$$0 \leq \lambda \mapsto \sum_{i=1}^p \frac{a_i^2 v_i^2}{(a_i^2 + \lambda)^2}$$

is strictly decreasing, its value at 0 is strictly greater than 1 (notice that  $v \notin M$ ), while its limit at  $+\infty$  is 0. So, to get  $\bar{v}$  one must solve an algebraic equation of degree  $2p$ , and, in general, this is impossible. We will be interested in approximation methods of the solutions of nonlinear equations, and this will be one of the subjects of Chapter 6.