

This chapter collects some basic results on Lagrangian duality, in particular as it applies to convex programs with a zero duality gap.

## 6.1 The relaxation theorem

Given the problem to find

$$f^* := \inf_{\boldsymbol{x}} f(\boldsymbol{x}), \tag{6.1a}$$

subject to 
$$x \in S$$
, (6.1b)

where  $f : \mathbb{R}^n \to \mathbb{R}$  is a given function and  $S \subseteq \mathbb{R}^n$ , we define a relaxation to (6.1) to be a problem of the following form: find

$$f_R^* := \inf_x \inf_x f_R(x), \tag{6.2a}$$

subject to 
$$x \in S_R$$
, (6.2b)

where  $f_R : \mathbb{R}^n \to \mathbb{R}$  is a function with the property that  $f_R \leq f$  on S, and where  $S_R \supseteq S$ . For this pair of problems, we have the following basic result.

**Theorem 6.1** (Relaxation Theorem) (a) [relaxation]  $f_R^* \leq f^*$ .

(b) [infeasibility] If (6.2) is infeasible, then so is (6.1).

(c) [optimal relaxation] If the problem (6.2) has an optimal solution,  $\boldsymbol{x}_R^*$ , for which it holds that

$$\boldsymbol{x}_{R}^{*} \in S$$
 and  $f_{R}(\boldsymbol{x}_{R}^{*}) = f(\boldsymbol{x}_{R}^{*}),$  (6.3)

then  $\boldsymbol{x}_R^*$  is an optimal solution to (6.1) as well.

**Proof.** The result in (a) is obvious, as every solution feasible in (6.1) is both feasible in (6.2) and has a lower objective value in the latter problem.

The result in (b) follows for similar reasons. For the result in (c), we note that

$$f(\boldsymbol{x}_R^*) = f_R(\boldsymbol{x}_R^*) \le f_R(\boldsymbol{x}) \le f(\boldsymbol{x}), \qquad \boldsymbol{x} \in S,$$

from which the result follows.

This basic result will be utilized both in this chapter and later on to motivate why Lagrangian relaxation, objective function linearization and penalization constitute relaxations, and to derive optimality conditions and algorithms based on them.

## 6.2 Lagrangian duality

In this section we formulate the Lagrangian dual problem and establish its convexity. The Weak Duality Theorem is also established, and we introduce the terms "Lagrangian relaxation," "Lagrange multiplier," and "duality gap."

#### 6.2.1 Lagrangian relaxation and the dual problem

Consider the optimization problem to find

$$f^* := \inf_{x} \inf_{x} f(x),$$
  
subject to  $x \in X,$   
 $g_i(x) \le 0,$   $i = 1, \dots, m,$   
(6.4)

where  $f : \mathbb{R}^n \to \mathbb{R}$  and  $g_i : \mathbb{R}^n \to \mathbb{R}$  (i = 1, 2, ..., m) are given functions, and  $X \subseteq \mathbb{R}^n$ .

For this problem, we assume that

$$-\infty < f^* < \infty, \tag{6.5}$$

that is, that f is bounded from below on the feasible set and the problem has at least one feasible solution.

**Definition 6.2** (Lagrange function, relaxation, multiplier) (a) For an arbitrary vector  $\mu \in \mathbb{R}^m$ , the Lagrange function is

$$L(x, \mu) := f(x) + \sum_{i=1}^{m} \mu_i g_i(x) = f(x) + \mu^{\mathrm{T}} g(x).$$
 (6.6)

(b) Consider the problem to

minimize 
$$L(x, \mu)$$
, (6.7)  
subject to  $x \in X$ .

Whenever  $\mu$  is non-negative, the problem (6.7) is referred to as a Lagrangian relaxation.

(c) We call the vector  $\mu^* \in \mathbb{R}^m$  a Lagrange multiplier vector if it is non-negative and if  $f^* = \inf_{x \in X} L(x, \mu^*)$  holds.

Note that the Lagrangian relaxation (6.7) is a relaxation, in terms of Section 6.1.

**Theorem 6.3 (Lagrange multipliers and global optima)** Let  $\mu^*$  be a Lagrange multiplier vector. Then,  $x^*$  is an optimal solution to (6.4) if and only if  $x^*$  is feasible in (6.4) and

$$x^* \in \arg\min_{x \in X} L(x, \mu^*), \text{ and } \mu_i^* g_i(x^*) = 0, \ i = 1, \dots, m.$$
 (6.8)

**Proof.** If  $x^*$  is an optimal solution to (6.4), then it is in particular feasible, and

$$f^* = f(\boldsymbol{x}^*) \geq L(\boldsymbol{x}^*, \boldsymbol{\mu}^*) \geq \inf_{\boldsymbol{x} \in X} L(\boldsymbol{x}, \boldsymbol{\mu}^*),$$

where the first inequality stems from the feasibility of  $x^*$  and the definition of a Lagrange multiplier vector. The second part of that definition implies that  $f^* = \inf_{x \in X} L(x, \mu^*)$ , so that equality holds throughout in the above line of inequalities. Hence, (6.8) follows.

Conversely, if  $x^*$  is feasible and (6.8) holds, then by the use of the definition of a Lagrange multiplier vector,

$$f(x^*) = L(x^*, \mu^*) = \min_{x \in X} L(x, \mu^*) = f^*,$$

so  $x^*$  is a global optimum.

Let

$$q(\boldsymbol{\mu}) := \inf_{\boldsymbol{x} \in X} L(\boldsymbol{x}, \boldsymbol{\mu})$$
(6.9)

be the Lagrangian dual function, defined by the infimum value of the Lagrange function over X; the Lagrangian dual problem is to

$$\underset{\boldsymbol{\mu}}{\text{maximize } q(\boldsymbol{\mu}), } \tag{6.10}$$

subject to 
$$\mu \geq 0^m$$

143

For some  $\mu$ ,  $q(\mu) = -\infty$  is possible; if it is true for all  $\mu \ge 0^m$ , then

$$q^* := \underset{\boldsymbol{\mu} \ge \mathbf{0}^m}{\operatorname{supremum}} q(\boldsymbol{\mu})$$

equals  $-\infty$ . (We can then say that the dual problem is infeasible.)

The effective domain of q is

$$D_q := \{ \boldsymbol{\mu} \in \mathbb{R}^m \mid q(\boldsymbol{\mu}) > -\infty \}.$$

**Theorem 6.4 (convex dual problem)** The effective domain  $D_q$  of q is convex, and q is concave on  $D_q$ .

**Proof.** Let  $x \in \mathbb{R}^n$ ,  $\mu, \bar{\mu} \in \mathbb{R}^m$ , and  $\alpha \in [0, 1]$ . We have that

$$L(\boldsymbol{x}, \alpha \boldsymbol{\mu} + (1 - \alpha) \bar{\boldsymbol{\mu}}) = \alpha L(\boldsymbol{x}, \boldsymbol{\mu}) + (1 - \alpha) L(\boldsymbol{x}, \bar{\boldsymbol{\mu}}).$$

Take the infimum over  $x \in X$  on both sides; then,

$$\begin{split} \inf_{x \in X} L(x, \alpha \mu + (1 - \alpha)\bar{\mu}) &= \inf_{x \in X} \left\{ \alpha L(x, \mu) + (1 - \alpha)L(x, \bar{\mu}) \right\} \\ &\geq \inf_{x \in X} \alpha L(x, \mu) + \inf_{x \in X} (1 - \alpha)L(x, \bar{\mu}) \\ &= \alpha \inf_{x \in X} L(x, \mu) + (1 - \alpha) \inf_{x \in X} L(x, \bar{\mu}), \end{split}$$

since  $\alpha \in [0, 1]$ , and the sum of infimum values may be smaller than the infimum of the sum, since in the former case we have the possibility to choose different optimal solutions in the two problems. Hence,

$$q(\alpha \boldsymbol{\mu} + (1 - \alpha)\bar{\boldsymbol{\mu}}) \ge \alpha q(\boldsymbol{\mu}) + (1 - \alpha)q(\bar{\boldsymbol{\mu}})$$

holds. This inequality has two implications: if  $\mu$  and  $\bar{\mu}$  lie in  $D_q$ , then so does  $\alpha \mu + (1 - \alpha)\bar{\mu}$ , so  $D_q$  is convex; also, q is concave on  $D_q$ .

That the Lagrangian dual problem always is convex (we indeed maximize a concave function!) is good news, because it means that it can be solved efficiently. What remains is to show how a Lagrangian dual optimal solution can be used to generate a primal optimal solution.

Next, we establish that every feasible point in the Lagrangian *dual* problem always underestimates the objective function value of every feasible point in the *primal* problem; hence, also their optimal values have this relationship.

**Theorem 6.5 (Weak Duality Theorem)** (a) Let x and  $\mu$  be feasible in the problems (6.4) and (6.10), respectively. Then,

$$q(\boldsymbol{\mu}) \leq f(\boldsymbol{x}).$$

In particular,

$$q^* \le f^*.$$

(b) If  $q(\mu) = f(x)$ , then the pair  $(x, \mu)$  is optimal in its respective problem.

**Proof.** For all  $\mu \geq 0^m$  and  $x \in X$  with  $g(x) \leq 0^m$ ,

$$q(\boldsymbol{\mu}) = \inf_{\boldsymbol{z} \in X} L(\boldsymbol{z}, \boldsymbol{\mu}) \leq f(\boldsymbol{x}) + \boldsymbol{\mu}^{\mathrm{T}} \boldsymbol{g}(\boldsymbol{x}) \leq f(\boldsymbol{x}),$$

 $\mathbf{SO}$ 

$$q^* = \sup_{\mu \ge 0^m} q(\mu) \le \inf_{x \in X: g(x) \le 0^m} f(x) = f^*.$$

The result follows.

Weak duality is also a consequence of the Relaxation Theorem: For any  $\mu \geq 0^m$ , let

$$S := X \cap \{ \boldsymbol{x} \in \mathbb{R}^n \mid \boldsymbol{g}(\boldsymbol{x}) \le \boldsymbol{0}^m \},$$
(6.11a)

$$S_R := X, \tag{6.11b}$$

$$f_R := L(\boldsymbol{\mu}, \cdot). \tag{6.11c}$$

Then, the weak duality statement is the result in Theorem 6.1(a).

If our initial feasibility assumption (6.5) is false, then what does weak duality imply? Suppose that  $f^* = -\infty$ . Then, weak duality implies that  $q(\mu) = -\infty$  for all  $\mu \ge 0^m$ , that is, the dual problem is infeasible. Suppose then that  $X \ne \emptyset$  but that  $X \cap \{ x \in \mathbb{R}^n \mid g(x) \le 0^m \}$  is empty. Then,  $f^* = \infty$ , by convention. The dual function satisfies  $q(\mu) < \infty$  for all  $\mu \ge 0^m$ , but it is possible that  $q^* = -\infty, -\infty < q^* < \infty$ , or  $q^* = \infty$  (see [Ber99, Figure 5.1.8]). For linear programs,  $-\infty < q^* < \infty$  implies  $-\infty < f^* < \infty$ ; see below.

If  $q^* = f^*$ , then we say that the *duality gap* (as given by  $\Gamma := f^* - q^*$ ) is zero, or that *there is no duality gap*. If there exists a Lagrange multiplier vector, then by the weak duality theorem, this implies that there is no duality gap. The converse is not true in general: there may be cases where no Lagrange multipliers exist even when there is no duality gap; in that case though, the Lagrangian dual problem cannot have an optimal solution, as implied by the following result.

**Proposition 6.6** (duality gap and the existence of Lagrange multipliers) (a) If there is no duality gap, then the set of Lagrange multiplier vectors equals the set of optimal dual solutions (which however may be empty).

(b) If there is a duality gap, then there are no Lagrange multipliers.

**Proof.** By definition, a vector  $\boldsymbol{\mu}^* \geq \mathbf{0}^m$  is a Lagrange multiplier vector if and only if  $f^* = q(\boldsymbol{\mu}^*) \leq q^*$ , the equality following from the definition of  $q(\boldsymbol{\mu}^*)$  and the inequality from the definition of  $q^*$  as the supremum of  $q(\boldsymbol{\mu})$  over  $\mathbb{R}^m_+$ . By weak duality, this relation holds if and only if there is no duality gap and  $\boldsymbol{\mu}^*$  is an optimal dual solution.

Before moving on, we remark on the statement of the problem (6.4). There are several ways in which the original set of constraints of the problem can be placed either within the definition of the ground set X (which is kept intact), or within the explicit constraints defined by the functions  $g_i$  (which are Lagrangian relaxed). How to distinguish between the two, that is, how to decide whether a constraint should be kept or be Lagrangian relaxed, depends on several factors. For example, keeping more constraints within X may result in a smaller duality gap, and fewer multipliers also result in a simpler Lagrangian dual problem. On the other hand, the Lagrangian subproblem defining the dual function simultaneously becomes more complex and difficult to solve. There are no immediate rules to follow, but experimentation and experience.

#### 6.2.2 Global optimality conditions

The following result characterizes every optimal primal and dual solution. It is however applicable only in the presence of Lagrange multipliers; in other words, the below system (6.12) is consistent if and only if there exists a Lagrange multiplier vector and there is no duality gap.

**Theorem 6.7** (global optimality conditions in the absence of a duality gap) The vector  $(x^*, \mu^*)$  is a pair of primal optimal solution and Lagrange multiplier vector if and only if

$\boldsymbol{\mu}^* \geq \boldsymbol{0}^m,$	(Dual feasibility)	(6.12a)
$oldsymbol{x}^* \in rg\min_{oldsymbol{x} \in X} L(oldsymbol{x},oldsymbol{\mu}^*),$	(Lagrangian optimality)	(6.12b)
$x^* \in X,  g(x^*) \le 0^m,$	(Primal feasibility)	(6.12c)
$\mu_i^* g_i(\boldsymbol{x}^*) = 0,  i = 1, \dots, m.$	$(Complementary \ slackness)$	(6.12d)

**Proof.** Suppose that the pair  $(x^*\mu^*)$  satisfies (6.12). Then, from (6.12a) we have that the Lagrangian problem to minimize  $L(x, \mu^*)$  over  $x \in X$  is a (Lagrangian) relaxation of (6.4). Moreover, according to (6.12b)  $x^*$  solves this problem, (6.12c) shows that  $x^*$  is feasible in (6.4), and (6.12d) implies that  $L(x^*, \mu^*) = f(x^*)$ . The Relaxation Theorem 6.1

then yields that  $x^*$  is optimal in (6.4), which in turn implies that  $\mu^*$  is a Lagrange multiplier vector.

Conversely, if  $(x^*, \mu^*)$  is a pair of optimal primal solution and Lagrange multiplier vector, then they are primal and dual feasible, respectively. The relations (6.12b) and (6.12d) follow from Theorem 6.3.

**Theorem 6.8** (global optimality and saddle points) The vector  $(x^*, \mu^*)$  is a pair of optimal primal solution and Lagrange multiplier vector if and only if  $x^* \in X$ ,  $\mu^* \geq 0^m$ , and  $(x^*, \mu^*)$  is a saddle point of the Lagrangian function on  $X \times \mathbb{R}^m_+$ , that is,

$$L(\boldsymbol{x}^*, \boldsymbol{\mu}) \le L(\boldsymbol{x}^*, \boldsymbol{\mu}^*) \le L(\boldsymbol{x}, \boldsymbol{\mu}^*), \qquad (\boldsymbol{x}, \boldsymbol{\mu}) \in X \times \mathbb{R}^m_+, \qquad (6.13)$$

holds.

**Proof.** We establish that (6.12) and (6.13) are equivalent; Theorem 6.7 then gives the result. The first inequality in (6.13) is equivalent to

$$-\boldsymbol{g}(\boldsymbol{x}^*)^{\mathrm{T}}(\boldsymbol{\mu}-\boldsymbol{\mu}^*) \ge 0, \qquad \boldsymbol{\mu} \in \mathbb{R}^m_+, \qquad (6.14)$$

for the given pair  $(x^*, \mu^*) \in X \times \mathbb{R}^m_+$ . This variational inequality is equivalent to stating that<sup>1</sup>

$$\mathbf{0}^m \ge \boldsymbol{g}(\boldsymbol{x}^*) \perp \boldsymbol{\mu}^* \ge \mathbf{0}^m, \tag{6.15}$$

where  $\perp$  denotes orthogonality: that is, for any vectors  $\boldsymbol{a}, \boldsymbol{b} \in \mathbb{R}^n$ ,  $\boldsymbol{a} \perp \boldsymbol{b}$ means that  $\boldsymbol{a}^T \boldsymbol{b} = 0$ . Because of the sign restrictions posed on  $\boldsymbol{\mu}$  and  $\boldsymbol{g}$ , that is, the vectors  $\boldsymbol{a}$  and  $\boldsymbol{b}$ , the relation  $\boldsymbol{a} \perp \boldsymbol{b}$  actually means that not only does it hold that  $\boldsymbol{a}^T \boldsymbol{b} = 0$  but in fact  $a_i b_i = 0$  must hold for all  $i = 1, \ldots, n$ . This complementarity system is, for the given  $\boldsymbol{\mu}^* \in \mathbb{R}^m_+$ , the same as (6.12a), (6.12c) and (6.12d). The second inequality in (6.13) is equivalent to (6.12b).

The above two theorems also imply that the set of primal-dual optimal solutions  $(x^*, \mu^*)$  is a Cartesian product set,  $X^* \times U^*$ . For example,

 $<sup>^{1}</sup>$ We establish the equivalence between (6.14) and (6.15) as follows. (The proof extends that for line search problems in unconstrained optimization in a footnote in Section 11.3.1.)

First, suppose that (6.15) is fulfilled. Then,  $-g(x^*)^T(\mu - \mu^*) = -g(x^*)^T \mu \ge 0$ , for all  $\mu \ge 0^m$ , that is, (6.14) is fulfilled. Conversely, suppose that (6.14) is fulfilled. Setting  $\mu = 0^m$  yields that  $g(x^*)^T \mu^* \ge 0$ . On the other hand, the choice  $\mu = 2\mu^*$ yields that  $-g(x^*)^T \mu^* \ge 0$ . Hence,  $g(x^*)^T \mu^* = 0$  holds. Last, let  $\mu = \mu^* + e_i$ , where  $e_i$  is the *i*<sup>th</sup> unit vector in  $\mathbb{R}^m$ . Then,  $-g(x^*)^T(\mu - \mu^*) = -g_i(x^*) \ge 0$ . Since this is true for all  $i \in \{1, 2, ..., m\}$  we have obtained that  $-g(x^*) \ge 0^m$ , that is,  $g(x^*) < 0^m$ . We are done.

given any optimal dual solution  $\mu^* \in U^*$ , every optimal primal solution  $x^* \in X^*$  satisfies (6.12). Hence, we can write, for an *arbitrary* dual vector  $\mu^* \in U^*$ ,

$$X^* = \left\{ \boldsymbol{x}^* \in \mathbb{R}^n \mid \boldsymbol{x}^* \text{ satisfies (6.12) for } \boldsymbol{\mu} = \boldsymbol{\mu}^* \right\}$$
$$= \left\{ \left. \boldsymbol{x}^* \in \arg\min_{\boldsymbol{x} \in X} L(\boldsymbol{x}, \boldsymbol{\mu}^*) \right| \boldsymbol{g}(\boldsymbol{x}^*) \leq \boldsymbol{0}^m; \ (\boldsymbol{\mu}^*)^{\mathrm{T}} \boldsymbol{g}(\boldsymbol{x}^*) = \boldsymbol{0} \right\}.$$

We note that structurally similar results to the above two theorems which are valid for the general problem (6.4) with any size of the duality gap can be found in [LaP05].<sup>2</sup>

We finally note a practical connection between the KKT system (5.9) and the above system (6.12). The practical use of the KKT system is normally to investigate whether a primal vector  $\boldsymbol{x}$ —obtained perhaps from a solver for our problem—is a candidate for a locally optimal solution; in other words, we have access to  $\boldsymbol{x}$  and generate a vector  $\boldsymbol{\mu}$  of Lagrange multipliers in the investigation of the KKT system (5.9). In contrast, the system (6.12) is normally investigated in the reverse order; we formulate and solve the Lagrangian dual problem, thereby obtaining an optimal dual vector  $\boldsymbol{\mu}$ . Starting from that vector, we investigate the global optimality conditions stated in (6.12) to obtain, if possible, an optimal primal vector  $\boldsymbol{x}$ . In the section to follow, we show when this is possible, and provide strong connections between the systems (5.9) and (6.12) in the convex and differentiable case.

#### 6.2.3 Strong duality for convex programs

So far the results have been rather non-technical to achieve: the convexity of the Lagrangian dual problem comes with very few assumptions on the original, primal problem, and the characterization of the primaldual set of optimal solutions is simple and also quite easily established. In order to establish *strong duality*, that is, to establish sufficient conditions under which there is no duality gap, however, takes much more. In particular, as is the case with the KKT conditions we need regularity conditions (that is, constraint qualifications), and we also need to utilize separation theorems such as Theorem 4.28. Most importantly, however, is that strong duality is deeply associated with the convexity of the original problem, and it is in particular under convexity that the primal and

<sup>&</sup>lt;sup>2</sup>The system (6.12) is there appended with two relaxation parameters which measure, respectively, the near-optimality of  $\boldsymbol{x}^*$  in the Lagrangian subproblem [that is, the  $\varepsilon$ -optimality in (6.12b)], and the violation of the complementarity conditions (6.12d). The saddle point condition (6.13) is similarly perturbed, and at an optimal solution, the sum of these two parameters equals the duality gap.

dual optimal solutions are linked through the global optimality conditions provided in the previous section. We begin by concentrating on the inequality constrained case, proving this result in detail. We will also specialize the result to quadratic and linear optimization problems.

Consider the inequality constrained convex program (6.4), where  $f : \mathbb{R}^n \to \mathbb{R}$  and  $g_i$  (i = 1, ..., m) are convex functions and  $X \subseteq \mathbb{R}^n$  is a convex set. For this problem, we introduce the following regularity condition, due to Slater (cf. Definition 5.38):

$$\exists x \in X \text{ with } g(x) < 0^m. \tag{6.16}$$

**Theorem 6.9 (Strong Duality, inequality constrained convex programs)** Suppose that the feasibility condition (6.5) and Slater's constraint qualification (6.16) hold for the convex problem (6.4).

(a) There is no duality gap and there exists at least one Lagrange multiplier vector  $\mu^*$ . Moreover, the set of Lagrange multipliers is bounded and convex.

(b) If the infimum in (6.4) is attained at some  $x^*$ , then the pair  $(x^*, \mu^*)$  satisfies the global optimality conditions (6.12).

(c) If further f and g are differentiable at  $x^*$ , then the condition (6.12b) can equivalently be written as the variational inequality

$$\nabla_{\boldsymbol{x}} L(\boldsymbol{x}^*, \boldsymbol{\mu}^*)^{\mathrm{T}}(\boldsymbol{x} - \boldsymbol{x}^*) \ge 0, \qquad \boldsymbol{x} \in X.$$
(6.17)

If, in addition, X is open (such as is the case when  $X = \mathbb{R}^n$ ), then this reduces to the condition that

$$\nabla_{x} L(x^{*}, \mu^{*}) = \nabla f(x^{*}) + \sum_{i=1}^{m} \mu_{i}^{*} \nabla g_{i}(x^{*}) = \mathbf{0}^{n}, \qquad (6.18)$$

and the global optimality conditions (6.12) reduce to the Karush–Kuhn– Tucker conditions stated in Theorem 5.25.

**Proof.** (a) We begin by establishing the existence of a Lagrange multiplier vector (and the presence of a zero duality gap).<sup>3</sup>

First, we consider the following subset of  $\mathbb{R}^{m+1}$ :

$$A := \{ (z_1, \dots, z_m, w)^{\mathrm{T}} | \exists x \in X \text{ with } g_i(x) \le z_i, i = 1, \dots, m; f(x) \le w \}.$$

It is elementary to show that A is convex.

Next, we observe that  $((\mathbf{0}^m)^{\mathrm{T}}, f^*)^{\mathrm{T}}$  is not an interior point of A; otherwise, for some  $\varepsilon > 0$  the point  $((\mathbf{0}^m)^{\mathrm{T}}, f^* - \varepsilon)^{\mathrm{T}} \in A$  holds, which

<sup>&</sup>lt;sup>3</sup>This result is [Ber99, Proposition 5.3.1], whose proof we also utilize.

would contradict the definition of  $f^*$ . Therefore, by the (possibly nonproper) separation result in Theorem 4.28, we can find a hyperplane passing through  $((\mathbf{0}^m)^{\mathrm{T}}, f^*)^{\mathrm{T}}$  such that A lies in one of the two corresponding half-spaces. In particular, there then exists a vector  $(\boldsymbol{\mu}^{\mathrm{T}}, \beta)^{\mathrm{T}} \neq ((\mathbf{0}^m)^{\mathrm{T}}, 0)^{\mathrm{T}}$  such that

$$\beta f^* \le \beta w + \boldsymbol{\mu}^{\mathrm{T}} \boldsymbol{z}, \qquad (\boldsymbol{z}^{\mathrm{T}}, w)^{\mathrm{T}} \in A.$$
(6.19)

This implies that

$$\geq 0; \qquad \boldsymbol{\mu} \geq \boldsymbol{0}^m, \tag{6.20}$$

since for each  $(\boldsymbol{z}^{\mathrm{T}}, w)^{\mathrm{T}} \in A$   $(\boldsymbol{z}^{\mathrm{T}}, w + \gamma)^{\mathrm{T}} \in A$  and  $(z_1, \ldots, z_{i-1}, z_i + \gamma, z_{i+1}, \ldots, z_m, w)^{\mathrm{T}} \in A$  for all  $\gamma > 0$  and  $i = 1, \ldots, m$ .

β

We claim that  $\beta > 0$  in fact holds. Indeed, if it was not the case, then  $\beta = 0$  and (6.19) then implies that  $\boldsymbol{\mu}^{\mathrm{T}} \boldsymbol{z} \geq 0$  for every pair  $(\boldsymbol{z}^{\mathrm{T}}, w)^{\mathrm{T}} \in A$ . But since  $(\boldsymbol{g}(\bar{\boldsymbol{x}})^{\mathrm{T}}, f(\bar{\boldsymbol{x}}))^{\mathrm{T}} \in A$  [where  $\bar{\boldsymbol{x}}$  is such that it satisfies the Slater condition (6.16)], we would obtain that  $0 \leq \sum_{i=1}^{m} \mu_i g_i(\bar{\boldsymbol{x}})$  which in view of  $\boldsymbol{\mu} \geq \mathbf{0}^m$  [cf. (6.20)] and the assumption that  $\bar{\boldsymbol{x}}$  satisfies the Slater condition (6.16) implies that  $\boldsymbol{\mu} = \mathbf{0}^m$ . This means, however, that  $(\boldsymbol{\mu}^{\mathrm{T}}, \beta)^{\mathrm{T}} = ((\mathbf{0}^m)^{\mathrm{T}}, 0)^{\mathrm{T}}$ —a contradiction. We may therefore claim that  $\beta > 0$ . We further, with no loss of generality, assume that  $\beta = 1$ .

Thus, since  $(\boldsymbol{g}(\boldsymbol{x})^{\mathrm{T}}, f(\boldsymbol{x}))^{\mathrm{T}} \in A$  for every  $\boldsymbol{x} \in X$ , (6.19) yields that

$$f^* \leq f(x) + \mu^{\mathrm{T}} g(x), \qquad x \in X.$$

Taking the infimum over  $x \in X$  and using the fact that  $\mu \ge 0^m$  we obtain

$$f^* \leq \inf_{\boldsymbol{x} \in X} \{f(\boldsymbol{x}) + \boldsymbol{\mu}^{\mathrm{T}} \boldsymbol{g}(\boldsymbol{x})\} = q(\boldsymbol{\mu}) \leq \operatorname{supremum}_{\boldsymbol{\mu} \geq \boldsymbol{0}^{m}} q(\boldsymbol{\mu}) = q^*.$$

Using the Weak Duality Theorem 6.5 it follows that  $\mu$  is a Lagrange multiplier vector, and there is no duality gap. This part of the proof is now done.

Take any vector  $\bar{x} \in X$  satisfying (6.16) and a Lagrange multiplier vector  $\mu^*$ . By the definition of a Lagrange multiplier vector,  $f^* \leq L(\bar{x}, \mu^*)$  holds, which implies that

$$\sum_{i=1}^{m} \mu_i^* \le \frac{[f(\bar{x}) - f^*]}{\min_{i=1,\dots,m} \{-g_i(\bar{x})\}}$$

Since  $\mu^* \geq 0^m$ , boundedness follows. As by Proposition 6.6(a) the set of Lagrange multipliers is the set of optimal solutions to the dual problem (6.10), convexity follows from the identification of the dual solution set with the set of vectors  $\mu \in \mathbb{R}^m_+$  for which

 $q(\boldsymbol{\mu}) \geq q^*$ 

holds. This is the upper level set for q at the level  $q^*$ ; this set is convex, by the concavity of q (cf. Theorem 6.4 and Proposition 3.44).

(b) The result follows from Theorem 6.7.

(c) The first part follows from Theorem 4.23, as the Lagrangian function  $L(\cdot, \mu^*)$  is convex. The second part follows by identification.

Consider next the extension of the inequality constrained convex program (6.4) in which we seek to find

$$f^* := \inf_{x} \inf_{x} f(x), \qquad (6.21)$$
  
subject to  $x \in X,$   
 $g_i(x) \le 0, \qquad i = 1, \dots, m,$   
 $\varepsilon_j^{\mathrm{T}} x - d_j = 0, \qquad j = 1, \dots, \ell,$ 

under the same conditions as stated following (6.4), and where  $\varepsilon_j \in \mathbb{R}^n$ ,  $j = 1, \ldots, \ell$ . For this problem, we replace the Slater condition (6.16) with the following (cf. [BSS93, Theorem 6.2.4]):

 $\exists x \in X \text{ with } g(x) < 0^m \text{ and } 0^m \in \operatorname{int} \{ Ex - d \mid x \in X \}, \quad (6.22)$ 

where  $\boldsymbol{E} \in \mathbb{R}^{\ell \times n}$  has rows  $\boldsymbol{\varepsilon}_j^{\mathrm{T}}$ , and  $\boldsymbol{d} = (d_j)_{j \in \{1, \dots, \ell\}} \in \mathbb{R}^{\ell}$ .

Note that in the statement (6.22), the "int" can be stricken whenever X is polyhedral, so that the latter part simply states that Ex = d.

For this problem, the Lagrangian dual problem is to find

$$\begin{aligned} q^* &:= \underset{(\boldsymbol{\mu},\boldsymbol{\lambda})}{\text{subject to } \boldsymbol{\mu} \geq \mathbf{0}^m, \end{aligned} \tag{6.23}$$

where

$$q(\mu, \lambda) := \inf_{x} \operatorname{Infimum} L(x, \mu, \lambda) := f(x) + \mu^{\mathrm{T}} g(x) + \lambda^{\mathrm{T}} (Ex - d),$$
  
subject to  $x \in X$ .

**Theorem 6.10 (Strong Duality, general convex programs)** Suppose that in addition to the feasibility condition (6.5), Slater's constraint qualification (6.22) holds for the problem (6.21).

(a) The duality gap is zero and there exists at least one Lagrange multiplier vector pair  $(\mu^*, \lambda^*)$ .

(b) If the infimum in (6.21) is attained at some  $x^*$ , then the triple

 $(x^*, \mu^*, \lambda^*)$  satisfies the global optimality conditions

$$\mu^* \ge \mathbf{0}^m, \quad (Dual \ feasibility) \tag{6.24a}$$

$$\boldsymbol{x}^* \in \arg\min_{\boldsymbol{x} \in X} L(\boldsymbol{x}, \boldsymbol{\mu}^*, \boldsymbol{\lambda}^*), \quad (Lagrangian optimality)$$
(6.24b)

$$x^* \in X, \ g(x^*) \leq \mathbf{0}^m, \ Ex^* = d,$$
 (Primal feasibility) (6.24c)  
 $u^* a (x^*) = 0, \ i = 1, \dots, m$  (Complementary shedross) (6.24d)

$$\mu_i^* g_i(\boldsymbol{x}^*) = 0, \ i = 1, \dots, m.$$
 (Complementary slackness) (6.24d)

(c) If further f and g are differentiable at  $x^*$ , then the condition (6.24b) can equivalently be written as

$$\nabla_{\boldsymbol{x}} L(\boldsymbol{x}^*, \boldsymbol{\mu}^*, \boldsymbol{\lambda}^*)^{\mathrm{T}}(\boldsymbol{x} - \boldsymbol{x}^*) \ge 0, \qquad \boldsymbol{x} \in X.$$
 (6.25)

If, in addition, X is open (such as is the case when  $X = \mathbb{R}^n$ ), then this reduces to the condition that

$$\nabla_{\boldsymbol{x}} L(\boldsymbol{x}^*, \boldsymbol{\mu}^*, \boldsymbol{\lambda}^*) = \nabla f(\boldsymbol{x}^*) + \sum_{i=1}^m \mu_i^* \nabla g_i(\boldsymbol{x}^*) + \sum_{j=1}^\ell \lambda_j^* \varepsilon_j = \boldsymbol{0}^n, \quad (6.26)$$

and the global optimality conditions (6.24) reduce to the Karush–Kuhn– Tucker conditions stated in Theorem 5.33.

**Proof.** The proof is similar to that of Theorem 6.9.

We finally consider a special case where automatically a regularity condition holds.

Consider the affinely constrained convex program to find

$$f^* := \inf_{x} \inf_{x} f(x), \qquad (6.27)$$
  
subject to  $x \in X,$   
 $a_i^{\mathrm{T}} x - b_i \leq 0, \qquad i = 1, \dots, m,$   
 $\varepsilon_j^{\mathrm{T}} x - d_j = 0, \qquad j = 1, \dots, \ell,$ 

where  $f : \mathbb{R}^n \to \mathbb{R}$  is convex and  $X \subseteq \mathbb{R}^n$  is polyhedral.

**Theorem 6.11 (Strong Duality, affine constraints)** If the feasibility condition (6.5) holds for the problem (6.27), then there is no duality gap and there exists at least one Lagrange multiplier vector.

**Proof.** Again, the proof is similar to that of Theorem 6.9, except that no additional regularity conditions are needed.<sup>4</sup>

<sup>&</sup>lt;sup>4</sup>For a detailed proof, see [Ber99, Proposition 5.2.1]. (The special case where f is moreover differentiable is covered in [Ber99, Proposition 3.4.2].)

The existence of a multiplier vector [which by Proposition 6.6 and the absence of a duality gap implies the existence of an optimal solution to the dual problem (6.10)] does not imply the existence of an optimal solution to the primal problem (6.27) without any additional assumptions (take the minimization of f(x) := 1/x over  $x \ge 1$  for example). However, when f is either weakly coercive, quadratic or linear, the existence results are stronger; see the primal existence results in Theorems 4.6, 4.7, and 6.12 below, for example.

For convex programs where a Slater CQ holds, the Lagrange multipliers defined in this section, and those that appear in the Karush– Kuhn–Tucker conditions, clearly are identical. Next, we specialize the above to linear and quadratic programs.

#### 6.2.4 Strong duality for linear and quadratic programs

The following result will be established and analyzed in detail in Chapter 10 on linear programming duality (cf. Theorem 10.6), but can in fact also be established similarly to above. (See [BSS93, Theorem 2.7.3] or [Ber99, Proposition 5.2.2], for example.) Its proof will however be relegated to that of Theorem 10.6.

**Theorem 6.12 (Strong Duality, linear programs)** Assume, in addition to the conditions of Theorem 6.11, that f is linear, so that (6.27) is a linear program. Then, the primal and dual problems have optimal solutions and there is no duality gap.

The above result states a strong duality result for a general linear program. We next develop an explicit Lagrangian dual problem for a linear program.

Let  $A \in \mathbb{R}^{m \times n}$ ,  $c \in \mathbb{R}^n$ , and  $b \in \mathbb{R}^m$ ; consider the linear program

$$\begin{array}{l} \underset{x}{\text{minimize } c^{1}x,} \quad (6.28)\\ \text{subject to } Ax = b,\\ x \geq 0^{n}. \end{array}$$

If we let  $X := \mathbb{R}^n_+$ , then the Lagrangian dual problem is to

$$\begin{array}{l} \underset{\boldsymbol{\lambda} \in \mathbb{R}^{m}}{\operatorname{maximize}} \quad \boldsymbol{b}^{\mathrm{T}} \boldsymbol{\lambda}, \qquad (6.29) \\ \text{subject to } \quad \boldsymbol{A}^{\mathrm{T}} \boldsymbol{\lambda} < \boldsymbol{c}. \end{array}$$

The reason why we can write it in this form is that

$$q(\boldsymbol{\lambda}) := \inf_{x \ge 0^n} \left\{ \boldsymbol{c}^{\mathrm{T}} \boldsymbol{x} + \boldsymbol{\lambda}^{\mathrm{T}} (\boldsymbol{b} - \boldsymbol{A} \boldsymbol{x}) \right\} = \boldsymbol{b}^{\mathrm{T}} \boldsymbol{\lambda} + \inf_{x \ge 0^n} (\boldsymbol{c} - \boldsymbol{A}^{\mathrm{T}} \boldsymbol{\lambda})^{\mathrm{T}} \boldsymbol{x},$$

so that

$$q(\boldsymbol{\lambda}) = \begin{cases} \boldsymbol{b}^{\mathrm{T}} \boldsymbol{\lambda}, & \text{if } \boldsymbol{A}^{\mathrm{T}} \boldsymbol{\lambda} \leq \boldsymbol{c}, \\ -\infty, & \text{otherwise.} \end{cases}$$

(The infimum is attained at zero if and only if these inequalities are satisfied; otherwise, the inner problem is unbounded below.)

Further, why is it that  $\lambda$  here is not restricted in sign? Suppose we were to split the system Ax = b into an inequality system of the form

$$egin{array}{lll} Ax\leq b,\ -Ax\leq -b. \end{array}$$

Let  $((\mu^+)^T, (\mu^-)^T)^T$  be the corresponding vector of multipliers, and take the Lagrangian dual for this formulation. Then, we would have a Lagrange function of the form

$$(\boldsymbol{x}, \boldsymbol{\mu}^+, \boldsymbol{\mu}^-) \mapsto L(\boldsymbol{x}, \boldsymbol{\mu}^+, \boldsymbol{\mu}^-) := \boldsymbol{c}^{\mathrm{T}} \boldsymbol{x} + (\boldsymbol{\mu}^+ - \boldsymbol{\mu}^-)^{\mathrm{T}} (\boldsymbol{b} - \boldsymbol{A} \boldsymbol{x}),$$

and since  $\mu^+ - \mu^-$  can take on any value in  $\mathbb{R}^m$  we can simply replace it with the unrestricted vector  $\lambda \in \mathbb{R}^m$ . This motivates why the multiplier for an equality constraint never is sign restricted; the same was the case, as we saw in Section 5.6, for the multipliers in the KKT conditions.

As applied to this problem, Theorem 6.12 states that if both the primal or dual problems have feasible solutions, then they both have optimal solutions, satisfying strong duality  $(c^{\mathrm{T}}x^* = b^{\mathrm{T}}\lambda^*)$ . On the other hand, if any of the two problems has an unbounded solution, then the other problem is infeasible.

Consider next the quadratic programming problem to

$$\min_{x} \max_{x} \left\{ \frac{1}{2} x^{\mathrm{T}} Q x + c^{\mathrm{T}} x \right\},$$
subject to  $Ax \leq b,$ 

$$(6.30)$$

where  $\boldsymbol{Q} \in \mathbb{R}^{n \times n}$ ,  $\boldsymbol{c} \in \mathbb{R}^n$ ,  $\boldsymbol{A} \in \mathbb{R}^{m \times n}$ , and  $\boldsymbol{b} \in \mathbb{R}^m$ . We develop an explicit dual problem under the assumption that  $\boldsymbol{Q}$  is positive definite.

By Lagrangian relaxing the inequality constraints, we obtain that the inner problem in x is solved by letting

$$\boldsymbol{x} = -\boldsymbol{Q}^{-1}(\boldsymbol{c} + \boldsymbol{A}^{\mathrm{T}}\boldsymbol{\mu}). \tag{6.31}$$

Substituting this expression into the Lagrangian function yields the Lagrangian dual problem to

$$\begin{array}{l} \underset{\mu}{\operatorname{maximize}} & \left\{ -\frac{1}{2} \boldsymbol{\mu}^{\mathrm{T}} \boldsymbol{A} \boldsymbol{Q}^{-1} \boldsymbol{A}^{\mathrm{T}} \boldsymbol{\mu} - (\boldsymbol{b} + \boldsymbol{A} \boldsymbol{Q}^{-1} \boldsymbol{c})^{\mathrm{T}} \boldsymbol{\mu} - \frac{1}{2} \boldsymbol{c}^{\mathrm{T}} \boldsymbol{Q}^{-1} \boldsymbol{c} \right\}, \quad (6.32) \\ \text{subject to } & \boldsymbol{\mu} \geq \boldsymbol{0}^{m}, \end{array}$$

Strong duality follows for this convex primal-dual pair of quadratic programs, in much the same way as for linear programs.

**Theorem 6.13 (Strong Duality, quadratic programs)** For the primal-dual pair of convex quadratic programs (6.30), (6.32), the following holds:

(a) If both problems have feasible solutions, then both problems also have optimal solutions, and the primal problem (6.30) also has a unique optimal solution, given by (6.31) for any optimal Lagrange multiplier vector, and in the two problems the optimal values are equal.

(b) If either of the two problems has an unbounded solution, then the other one is infeasible.

(c) Suppose that Q is positive semi-definite, and that the feasibility condition (6.5) holds. Then, both the problem (6.30) and its Lagrangian dual have nonempty, closed and convex sets of optimal solutions, and their optimal values are equal.

In the result (a) it is important to note that the Lagrangian dual problem (6.32) is not necessarily strictly convex; the matrix  $AQ^{-1}A^{T}$  need not be positive definite, especially so when A does not have full rank. The result (c) extends the strong duality result from linear programming, since Q in (c) can be the zero matrix. In the case of (c) we of course cannot write the Lagrangian dual problem in the form of (6.32) because Q is not necessarily invertible.

## 6.3 Two illustrative examples

**Example 6.14 (an explicit**, differentiable dual problem) Consider the problem to

$$\begin{array}{ll} \underset{x}{\text{minimize}} & f(x) := x_1^2 + x_2^2,\\ \text{subject to} & x_1 + x_2 \geq 4,\\ & x_j \geq 0, \qquad j = 1, 2. \end{array}$$

We consider the first constraint to be the complicated one, and hence define  $g(\mathbf{x}) := -x_1 - x_2 + 4$  and let  $X := \{ (x_1, x_2)^T \mid x_j \ge 0, j = 1, 2 \}$ . Then, the Lagrangian dual function is

$$\begin{aligned} q(\mu) &= \min_{x \in X} L(x, \mu) := f(x) - \mu(x_1 + x_2 - 4) \\ &= 4\mu + \min_{x \in X} \left\{ x_1^2 + x_2^2 - \mu x_1 - \mu x_2 \right\} \\ &= 4\mu + \min_{x_1 \ge 0} \left\{ x_1^2 - \mu x_1 \right\} + \min_{x_2 \ge 0} \left\{ x_2^2 - \mu x_2 \right\}, \qquad \mu \ge 0. \end{aligned}$$

For a fixed  $\mu \ge 0$ , the minimum is attained at  $x_1(\mu) = \frac{\mu}{2}, x_2(\mu) = \frac{\mu}{2}$ .

Substituting this expression into  $q(\mu)$ , we obtain that  $q(\mu) = f(\boldsymbol{x}(\mu)) - \mu(x_1(\mu) + x_2(\mu) - 4) = 4\mu - \frac{\mu^2}{2}$ .

Note that q is strictly concave, and it is differentiable everywhere (due to the fact that f, g are differentiable and  $x(\mu)$  is unique), by Danskin's Theorem 6.17(d).

We have that  $q'(\mu) = 4 - \mu = 0 \iff \mu = 4$ . As  $\mu = 4 \ge 0$ , it is the optimum in the dual problem:  $\mu^* = 4$ ;  $\boldsymbol{x}^* = (x_1(\mu^*), x_2(\mu^*))^{\mathrm{T}} = (2, 2)^{\mathrm{T}}$ . Also,  $f(\boldsymbol{x}^*) = q(\mu^*) = 8$ .

This is an example where the dual function is differentiable, and therefore we can utilize Proposition 6.29(c). In this case, the optimum  $x^*$  is also unique, so it is automatically given as  $x^* = x(\mu)$ .

**Example 6.15 (an implicit**, non-differentiable dual problem) Consider the linear programming problem to

minimize 
$$f(\mathbf{x}) := -x_1 - x_2$$
,  
subject to  $2x_1 + 4x_2 \le 3$ ,  
 $0 \le x_1 \le 2$ ,  
 $0 \le x_2 \le 1$ .

The optimal solution is  $x^* = (3/2, 0)^{\mathrm{T}}, f(x^*) = -3/2$ . Consider Lagrangian relaxing the first constraint, obtaining

$$\begin{split} L(\boldsymbol{x}, \mu) &= -x_1 - x_2 + \mu(2x_1 + 4x_2 - 3); \\ q(\mu) &= -3\mu + \min_{\substack{0 \le x_1 \le 2}} \left\{ (-1 + 2\mu)x_1 \right\} + \min_{\substack{0 \le x_2 \le 1}} \left\{ (-1 + 4\mu)x_2 \right\} \\ &= \begin{cases} -3 + 5\mu, & 0 \le \mu \le 1/4, \\ -2 + \mu, & 1/4 \le \mu \le 1/2, \\ -3\mu, & 1/2 \le \mu. \end{cases} \end{split}$$

Check that  $\mu^* = 1/2$ , and hence that  $q(\mu^*) = -3/2$ . For linear programs, we have strong duality, but how do we obtain the optimal primal solution from  $\mu^*$ ? It is clear that q is non-differentiable at  $\mu^*$ . Let us utilize the characterization given in the system (6.12).

First, at  $\mu^*$ , it is clear that  $X(\mu^*)$  is the set  $\{(2\alpha, 0)^T \mid 0 \le \alpha \le 1\}$ . Among the subproblem solutions, we next have to find one that is primal feasible as well as complementary.

Primal feasibility means that  $2 \cdot 2\alpha + 4 \cdot 0 \leq 3 \iff \alpha \leq 3/4$ .

Further, complementarity means that  $\mu^* \cdot (2x_1^* + 4x_2^* - 3) = 0 \iff \alpha = 3/4$ , since  $\mu^* \neq 0$ . We conclude that the only primal vector that satisfies the system (6.12) together with the dual optimal solution  $\mu^* = 1/2$  is  $x^* = (3/2, 0)^{\mathrm{T}}$ .

In the first example, the Lagrangian dual function is differentiable since  $\boldsymbol{x}(\mu)$  is unique. The second one shows that otherwise, there may be kinks in the function q where there are alternative solutions  $\boldsymbol{x}(\mu)$ ; as a result, to obtain a primal optimal solution becomes more complex. The Dantzig–Wolfe algorithm, for example, represents a means by which to automatize the process that we have just shown; the algorithm generates extreme points of  $X(\mu)$  algorithmically, and constructs the best feasible convex combination thereof, obtaining a primal–dual optimal solution in a finite number of iterations for linear programs.

The above examples motivate a deeper study of the differentiability properties of convex (or, concave) functions in general, and the Lagrangian dual objective function in particular.

# 6.4 \*Differentiability properties of the dual function

We have established that the Lagrangian dual problem (6.10) is a convex one, and further that under some circumstances the primal and dual optimal values are the same. We now turn to study the Lagrangian dual problem in detail, and in particular how it can be solved efficiently. First, we will establish when the dual function q is differentiable. We will see that differentiability holds only in some special cases, in which we can recognize the workings of the Lagrange multiplier method; this classic method was illustrated in Example 6.14. Most often, the function q will however be non-differentiable, and then this method will fail. This means that we must devise a more general numerical method which is not based on gradients but rather subgradients. This type of algorithm is the topic of the next section; we begin by studying the topic of subgradients of convex functions in general.

#### 6.4.1 Subdifferentiability of convex functions

Throughout this section we suppose that  $f : \mathbb{R}^n \to \mathbb{R}$  is a convex function, and study its subdifferentiability properties. We will later on apply our findings to the Lagrangian dual function q, or, rather, its negative -q. We first remark that a finite convex function is automatically continuous (cf. Theorem 4.26).

**Definition 6.16 (subgradient)** Let  $f : \mathbb{R}^n \to \mathbb{R}$  be a convex function. We say that a vector  $g \in \mathbb{R}^n$  is a subgradient of f at  $x \in \mathbb{R}^n$  if

$$f(\boldsymbol{y}) \ge f(\boldsymbol{x}) + \boldsymbol{g}^{\mathrm{T}}(\boldsymbol{y} - \boldsymbol{x}), \qquad \boldsymbol{y} \in \mathbb{R}^{n}.$$
 (6.33)

The set of such vectors g defines the subdifferential of f at x, and is denoted  $\partial f(x)$ .

For concave functions, the reverse inequality of course holds; for simplicity we will refer also to such vectors g as subgradients.

Notice the connection to the characterization of a convex function in  $C^1$  in Theorem 3.40(a). The difference between them is that g is not unique at a non-differentiable point. (Just as the gradient has a role in supporting hyperplanes to the graph of a convex function in  $C^1$ , the role of a subgradient is the same; at a non-differentiable point there are more then one supporting hyperplane to the graph of f.)

We illustrate this in Figure 6.1.



Figure 6.1: Three possible slopes of the convex function f at x.

Notice that a minimum  $x^*$  of f over  $\mathbb{R}^n$  is characterized by the inclusion  $\mathbf{0}^n \in \partial f(x^*)$ ; recognize, again, the similarity to the  $C^1$  case.

We list some additional basic results for convex functions next. Proofs will not be given here; we refer instead to the convex analysis text by Rockafellar [Roc70].

**Proposition 6.17** (properties of a convex function) Let  $f : \mathbb{R}^n \to \mathbb{R}$  be a convex function.

(a) [boundedness of  $\partial f(x)$ ] For every  $x \in \mathbb{R}^n$ ,  $\partial f(x)$  is a nonempty, convex, and compact set. If X is bounded then  $\bigcup_{x \in X} \partial f(x)$  is bounded.

(b) [closedness of  $\partial f$ ] The subdifferential mapping  $x \Rightarrow \partial f(x)$  is closed; in other words, if  $\{x_k\}$  is a sequence of vectors in  $\mathbb{R}^n$  converging to x, and  $g_k \in \partial f(x_k)$  holds for every k, then the sequence  $\{g_k\}$  of subgradients is bounded and every limit point thereof belongs to  $\partial f(x)$ .

(c) [directional derivative and differentiability] For every  $x \in \mathbb{R}^n$ , the directional derivative of f at x in the direction of  $p \in \mathbb{R}^n$  satisfies

$$f'(\boldsymbol{x};\boldsymbol{p}) = \underset{\boldsymbol{g} \in \partial f(\boldsymbol{x})}{\operatorname{maximum}} \boldsymbol{g}^{\mathrm{T}} \boldsymbol{p}.$$
(6.34)

In particular, f is differentiable at x with gradient  $\nabla f(x)$  if and only if it has  $\nabla f(x)$  as its unique subgradient at x; in that case, f'(x; p) = $\nabla f(\boldsymbol{x})^{\mathrm{T}}\boldsymbol{p}.$ 

(d) [Danskin's Theorem—directional derivatives of a convex max function] Let Z be a compact subset of  $\mathbb{R}^m$ , and let  $\phi : \mathbb{R}^n \times Z \to \mathbb{R}$  be continuous and such that  $\phi(\cdot, z) : \mathbb{R}^n \to \mathbb{R}$  is convex for each  $z \in Z$ . Let the function  $f : \mathbb{R}^n \to \mathbb{R}$  be given by

$$f(\boldsymbol{x}) := \underset{\boldsymbol{z} \in Z}{\operatorname{maximum}} \phi(\boldsymbol{x}, \boldsymbol{z}), \qquad \boldsymbol{x} \in \mathbb{R}^{n}.$$
(6.35)

The function f then is convex on  $\mathbb{R}^n$  and has a directional derivative at  $\boldsymbol{x}$  in the direction of  $\boldsymbol{p}$  equal to

$$f'(\boldsymbol{x};\boldsymbol{p}) = \underset{\boldsymbol{z}\in Z(\boldsymbol{x})}{\operatorname{maximum}} \ \phi'(\boldsymbol{x},\boldsymbol{z};\boldsymbol{p}), \tag{6.36}$$

where  $\phi'(x, z; p)$  is the directional derivative of  $\phi(\cdot, z)$  at x in the direction of  $\boldsymbol{p}$ , and  $Z(\boldsymbol{x}) := \{ \boldsymbol{z} \in \mathbb{R}^m \mid \phi(\boldsymbol{x}, \boldsymbol{z}) = f(\boldsymbol{x}) \}.$ 

In particular, if  $Z(\mathbf{x})$  contains a single point  $\bar{\mathbf{z}}$  and  $\phi(\cdot, \bar{\mathbf{z}})$  is differentiable at x, then f is differentiable at x, and  $\nabla f(x) = \nabla_x \phi(x, \bar{z})$ , where  $\nabla_{x}\phi(x,\bar{z})$  is the vector with components  $\frac{\partial\phi(x,\bar{z})}{\partial x_{i}}$ ,  $i = 1, \ldots, n$ . If further  $\phi(\cdot, z)$  is differentiable for all  $z \in Z$  and  $\nabla_{x}\phi(x, \cdot)$  is con-

tinuous on Z for each x, then

$$\partial f(\boldsymbol{x}) = \operatorname{conv} \{ \nabla_{\boldsymbol{x}} \phi(\boldsymbol{x}, \boldsymbol{z}) \mid \boldsymbol{z} \in Z(\boldsymbol{x}) \}, \qquad \boldsymbol{x} \in \mathbb{R}^{n}.$$

**Proof.** (a) This is a special case of [Roc70, Theorem 24.7].

- (b) This is [Roc70, Theorem 24.5].
- (c) This is [Roc70, Theorem 23.4 and 25.1].
- (d) This is [Ber99, Proposition B.25].

Figure 6.2 illustrates the subdifferential of a convex function.

We apply parts of the above results in order to characterize a minimum of a convex function on  $\mathbb{R}^n$ .

**Proposition 6.18** (optimality of a convex function over  $\mathbb{R}^n$ ) Let f:  $\mathbb{R}^n \to \mathbb{R}$  be a convex function. The following three statements are equivalent:

- 1. f is globally minimized at  $x^* \in \mathbb{R}^n$ ;
- 2.  $\mathbf{0}^n \in \partial f(\mathbf{x}^*);$
- 3.  $f'(\boldsymbol{x}^*; \boldsymbol{p}) \geq 0$  for all  $\boldsymbol{p} \in \mathbb{R}^n$ .



Figure 6.2: The subdifferential of a convex function f at x.

**Proof.** We establish the result thus:  $1 \implies 2 \implies 3 \implies 1$ .  $[1 \implies 2]$ : By the statement 1., we have that  $f(y) \ge f(x^*)$  for every  $y \in \mathbb{R}^n$ . This implies that for  $g = 0^n$ , we satisfy the subgradient inequality (6.33). This establishes the statement 2.

 $[2 \implies 3]$ : We can equivalently write

$$\partial f(\boldsymbol{x}) = \{ \boldsymbol{g} \in \mathbb{R}^n \mid \boldsymbol{g}^{\mathrm{T}} \boldsymbol{p} \leq f'(\boldsymbol{x}; \boldsymbol{p}), \quad \boldsymbol{p} \in \mathbb{R}^n \}.$$

With  $g = 0^n$  this definition immediately yields the statement 3.

 $[3 \implies 1]$ : By the compactness of the subdifferential [cf. Proposition 6.17(a)] and Weierstrass' Theorem 4.6 the maximum in the expression (6.34) is attained at some  $g \in \partial f(x^*)$ . It follows that, in the subgradient inequality (6.33), we get that

$$f(\boldsymbol{x}^* + \boldsymbol{p}) \ge f(\boldsymbol{x}^*) + \boldsymbol{g}^{\mathrm{T}} \boldsymbol{p} \ge f(\boldsymbol{x}^*), \qquad \boldsymbol{p} \in \mathbb{R}^n,$$

which is equivalent to the statement 1.

This result implies that a direction  $p \in \mathbb{R}^n$  is a descent direction with respect to f at x if and only if f'(x; p) < 0 holds. This result cannot be extended to non-convex functions, even when the function f is in  $C^1$ or even  $C^2$ . [Take  $f(x) := x^3$ ; x = 0; p = -1; see also the discussion on saddle points in Example 11.2(b).]

#### 6.4.2 Differentiability of the Lagrangian dual function

We consider the inequality constrained problem (6.4), where we make the following standing assumption:

 $f, g_i (i = 1, ..., m) \in C^0, X$  is nonempty and compact. (6.37)

Under this assumption, the set of solutions to the Lagrangian subproblem,

$$X(\boldsymbol{\mu}) := \arg\min_{\boldsymbol{x} \in X} L(\boldsymbol{x}, \boldsymbol{\mu}), \qquad \boldsymbol{\mu} \in \mathbb{R}^m, \tag{6.38}$$

is nonempty and compact for any choice of dual vector  $\mu$  by Weierstrass' Theorem 4.6. We first develop the subdifferentiability properties of the associated dual function q, stated in (6.9). The first result strengthens Theorem 6.4 under these additional assumptions.

**Proposition 6.19** (subdifferentiability of the dual function) Suppose that, in the problem (6.4), the compactness condition (6.37) holds.

(a) The dual function (6.9) is finite, continuous and concave on  $\mathbb{R}^m$ . If its supremum over  $\mathbb{R}^m_+$  is attained, then the optimal solution set therefore is closed and convex.

(b) The mapping  $\mu \rightleftharpoons X(\mu)$  is closed on  $\mathbb{R}^m$ . If  $X(\bar{\mu})$  is the singleton set  $\{\bar{x}\}$  for some  $\bar{\mu} \in \mathbb{R}^m$ , and for some sequence  $\mathbb{R}^m \supset \{\mu_k\} \to \bar{\mu}$ ,  $x_k \in X(\mu_k)$  for all k, then  $\{x_k\} \to \bar{x}$ .

(c) Let  $\mu \in \mathbb{R}^m$ . If  $x \in X(\mu)$ , then g(x) is a subgradient to q at  $\mu$ , that is,  $g(x) \in \partial q(\mu)$ .

(d) Let  $\boldsymbol{\mu} \in \mathbb{R}^m$ . Then,

$$\partial q(\boldsymbol{\mu}) = \operatorname{conv} \{ \boldsymbol{g}(\boldsymbol{x}) \mid \boldsymbol{x} \in X(\boldsymbol{\mu}) \}.$$

The set  $\partial q(\mu)$  is convex and compact. Moreover, if U is a bounded set, then  $\bigcup_{\mu \in U} \partial q(\mu)$  is also bounded.

(e) The directional derivative of q at  $\mu \in \mathbb{R}^m$  in the direction of  $p \in \mathbb{R}^m$  is

$$q'(\boldsymbol{\mu};\boldsymbol{p}) = \min_{\boldsymbol{g}\in\partial q(\boldsymbol{\mu})} \boldsymbol{g}^{\mathrm{T}}\boldsymbol{p}.$$

**Proof.** (a) Theorem 6.4 stated the concavity of q on its effective domain. Weierstrass' Theorem 4.6 states that q is finite on  $\mathbb{R}^m$ , which is then also its effective domain. The continuity of q follows from that of any finite concave function, as we have already seen in Theorem 4.26. The closedness property of the solution set is a direct consequence of the continuity of q (the upper level set then automatically is closed), and complements the result of Theorem 6.9(a).

(b) Let  $\{\mu_k\}$  be a sequence of vectors in  $\mathbb{R}^m$  converging to  $\bar{\mu}$ , and let  $x_k \in X(\mu_k)$  be arbitrary. Let x be arbitrary in X, and let further  $\bar{x} \in X$  be an arbitrary limit point of  $\{x_k\}$  (at least one exists by the compactness of X). From the property that for all k,

$$L(\boldsymbol{x}_k, \boldsymbol{\mu}_k) \leq L(\boldsymbol{x}, \boldsymbol{\mu}_k),$$

follows, by the continuity of L, that, in the limit of k in the subsequence in which  $\{x_k\}$  converges to  $\bar{x}$ ,

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

so that  $\bar{x} \in X(\bar{\mu})$ , as desired. The special case of a singleton set  $X(\bar{\mu})$  follows.

(c) Let  $\bar{\mu} \in \mathbb{R}^m$  be arbitrary and let  $\bar{x} \in X(\bar{\mu})$ . We have that

$$\begin{split} q(\bar{\boldsymbol{\mu}}) &= \inf_{\boldsymbol{y} \in X} L(\boldsymbol{y}, \bar{\boldsymbol{\mu}}) = f(\boldsymbol{x}) + \bar{\boldsymbol{\mu}}^{\mathrm{T}} \boldsymbol{g}(\boldsymbol{x}) \\ &= f(\boldsymbol{x}) + \boldsymbol{\mu}^{\mathrm{T}} \boldsymbol{g}(\boldsymbol{x}) + (\bar{\boldsymbol{\mu}} - \boldsymbol{\mu})^{\mathrm{T}} \boldsymbol{g}(\boldsymbol{x}) \geq q(\boldsymbol{\mu}) + (\bar{\boldsymbol{\mu}} - \boldsymbol{\mu})^{\mathrm{T}} \boldsymbol{g}(\boldsymbol{x}), \end{split}$$

which implies that  $g(x) \in \partial q(\mu)$ .

(d) The inclusion  $\partial q(\mu) \subseteq \text{conv} \{ g(x) \mid x \in X(\mu) \}$  follows from (c) and the convexity of  $\partial q(\mu)$ . The opposite inclusion follows by applying the Separation Theorem 3.24.<sup>5</sup>

(e) See Proposition 6.17(c).

The result in (c) is an independent proof of the concavity of q on  $\mathbb{R}^m$ . The result (d) is particularly interesting, because by Carathéodory's Theorem 3.8 every subgradient of q at any point  $\mu$  is the convex combination of a finite number (in fact, at most m + 1) of vectors of the form  $g(x^s)$  with  $x^s \in X(\mu)$ . Computationally, this has been utilized to devise efficient (proximal) bundle methods for the Lagrangian dual problem as well as to devise methods to recover primal optimal solutions.

Next, we establish the differentiability of the dual function under additional assumptions.

**Proposition 6.20** (differentiability of the dual function) Suppose that, in the problem (6.4), the compactness condition (6.37) holds.

(a) Let  $\mu \in \mathbb{R}^m$ . The dual function q is differentiable at  $\mu$  if and only if  $\{g(x) \mid x \in X(\mu)\}$  is a singleton set, that is, if the value of the vector of constraint functions is invariant over the set of solutions  $X(\mu)$  to the Lagrangian subproblem. Then, we have that

$$\nabla q(\boldsymbol{\mu}) = \boldsymbol{g}(\boldsymbol{x}),$$

 $<sup>^{5}</sup>$ See [BSS93, Theorem 6.3.7] for a detailed proof.

\*Subgradient optimization methods

for every  $\boldsymbol{x} \in X(\boldsymbol{\mu})$ .

(b) The result in (a) holds in particular if the Lagrangian subproblem has a unique solution, that is,  $X(\mu)$  is a singleton set. In particular, this property is satisfied for  $\mu \ge 0^m$  if further X is a convex set, f is strictly convex on X, and  $g_i$  (i = 1, ..., m) are convex, in which case  $q \in C^1$ .

**Proof.** (a) The concave function q is differentiable at the point  $\mu$  (where it is finite) if and only if its subdifferential  $\partial q(\mu)$  there is a singleton, cf. Proposition 6.17(c).

(b) Under either one of the assumptions stated,  $X(\mu)$  is a singleton, whence the result follows from (a). Uniqueness follows from the convexity of the feasible set and strict convexity of the objective function, according to Proposition 4.10. That  $q \in C^1$  follows from the continuity of g and Proposition 6.19(b).

**Proposition 6.21** (twice differentiability of the dual function) Suppose that, in the problem (6.4),  $X = \mathbb{R}^n$ , and f and  $g_i$  (i = 1, ..., m) are convex functions in  $C^2$ . Suppose that, at  $\mu \in \mathbb{R}^m$ , the solution x to the Lagrangian subproblem not only is unique, but also that the partial Hessian of the Lagrangian is positive definite at the pair  $(x, \mu)$ , that is,

 $\nabla_{xx}^2 L(x,\mu)$  is positive definite.

Then, the dual function q is twice differentiable at  $\mu$ , with

$$abla^2 q(oldsymbol{\mu}) = -
abla oldsymbol{g}(oldsymbol{x})^{\mathrm{T}} [
abla^2_{oldsymbol{xx}} L(oldsymbol{x},oldsymbol{\mu})]^{-1} 
abla oldsymbol{g}(oldsymbol{x})$$

**Proof.** The result follows from the Implicit Function Theorem, which is stated in Chapter 2, applied to the Lagrangian subproblem.<sup>6</sup>

## 6.5 \*Subgradient optimization methods

We begin by establishing the convergence of classic subgradient optimization methods as applied to a general convex optimization problem.

#### 6.5.1 Convex problems

Consider the convex optimization problem to

$$\min_{\boldsymbol{x}} f(\boldsymbol{x}), \tag{6.39a}$$

subject to 
$$x \in X$$
, (6.39b)

<sup>6</sup>See [Ber99, Pages 596-598] for a detailed analysis.

where  $f: \mathbb{R}^n \to \mathbb{R}$  is convex and the set  $X \subseteq \mathbb{R}^n$  is nonempty, closed and convex.

The subgradient projection algorithm is as follows: select  $x_0 \in X$ , and for  $k = 0, 1, \ldots$  generate

$$\boldsymbol{g}_k \in \partial f(\boldsymbol{x}_k), \tag{6.40a}$$

$$\boldsymbol{x}_{k+1} = \operatorname{Proj}_{\boldsymbol{X}} \left( \boldsymbol{x}_k - \alpha_k \boldsymbol{g}_k \right), \qquad (6.40b)$$

where the sequence  $\{\alpha_k\}$  is generated from one of the following three rules:

The first rule is termed the *divergent series* step length rule, and requires that

$$\alpha_k > 0, \quad k = 0, 1, \dots; \qquad \lim_{k \to \infty} \alpha_k = 0; \qquad \sum_{k=0}^{\infty} \alpha_k = +\infty.$$
 (6.41)

The second rule adds to the requirements in (6.41) the square-summable restriction  $\sim$ 

$$\sum_{k=0}^{\infty} \alpha_k^2 < +\infty. \tag{6.42}$$

The conditions in (6.41) allow for convergence to any point from any starting point, since the total step is infinite, but convergence is therefore also quite slow; the additional condition in (6.42) means fast sequences are selected. An instance of the step length formulas which satisfies both (6.41) and (6.42) is the following:

$$\alpha_k = \gamma + \beta / (k+1), \qquad k = 0, 1, \dots,$$

where  $\beta > 0, \gamma \ge 0$ .

The third step length rule is

$$\alpha_k = \theta_k \frac{f(x_k) - f^*}{\|g_k\|^2}, \qquad 0 < \sigma_1 \le \theta_k \le 2 - \sigma_2 < 2, \tag{6.43}$$

where  $f^*$  is the optimal value of (6.39). We refer to this step length formula as the *Polyak step*, after the Russian mathematician Boris Polyak who invented the subgradient method in the 1960s together with Ermol'ev and Shor.

How is convergence established for subgradient optimization methods? As shall be demonstrated in Chapters 11 and 12 convergence of algorithms for problems with a *differentiable* objective function is typically based on generating descent directions, and step length rules that result in the sequence  $\{x_k\}$  of iterates being strictly descending in the

value of f. For the non-differentiable problem at hand, generating descent directions is a difficult task, since it is not true that the negative of an arbitrarily chosen subgradient of f at a non-optimal vector x defines a descent direction.

In bundle methods one gathers information from more than one subgradient (hence the term bundle) around a current iteration point so that a descent direction can be generated, followed by an inexact line search. We concentrate here on the simpler methodology of subgradient optimization methods, in which we apply the formula (6.40) where the step length  $\alpha_k$  is chosen based on very simple rules.

We establish below that if the step length is small enough, an iteration of the subgradient projection method leads to a vector that is closer to the set of optimal solutions. This technical result also motivates the construction of the Polyak step length rule, and hence shows that the convergence of subgradient methods is based on the reduction of the Euclidean distance to the optimal solutions rather than on the reduction of the value of the objective function f.

**Proposition 6.22** (decreasing distance to the optimal set) Suppose that  $x_k \in X$  is not optimal in (6.39), and that  $x_{k+1}$  is given by (6.40) for some step length  $\alpha_k > 0$ .

Then, for every optimal solution  $x^*$  in (6.39),

$$\|x_{k+1} - x^*\| < \|x_k - x^*\|$$

holds for every step length  $\alpha_k$  in the interval

$$\alpha_k \in (0, 2[f(\boldsymbol{x}_k) - f^*] / \|\boldsymbol{g}_k\|^2).$$
(6.44)

**Proof.** We have that

$$\begin{split} \|\boldsymbol{x}_{k+1} - \boldsymbol{x}^*\|^2 &= \|\operatorname{Proj}_X \left( \boldsymbol{x}_k - \alpha_k \boldsymbol{g}_k \right) - \boldsymbol{x}^* \|^2 \\ &= \|\operatorname{Proj}_X \left( \boldsymbol{x}_k - \alpha_k \boldsymbol{g}_k \right) - \operatorname{Proj}_X \left( \boldsymbol{x}^* \right) \|^2 \\ &\leq \|\boldsymbol{x}_k - \alpha_k \boldsymbol{g}_k - \boldsymbol{x}^* \|^2 \\ &= \|\boldsymbol{x}_k - \boldsymbol{x}^* \|^2 - 2\alpha_k (\boldsymbol{x}_k - \boldsymbol{x}^*)^{\mathrm{T}} \boldsymbol{g}_k + \alpha_k^2 \|\boldsymbol{g}_k \|^2 \\ &\leq \|\boldsymbol{x}_k - \boldsymbol{x}^* \|^2 - 2\alpha_k [f(\boldsymbol{x}_k) - f^*] + \alpha_k^2 \|\boldsymbol{g}_k \|^2 \\ &< \|\boldsymbol{x}_k - \boldsymbol{x}^* \|^2, \end{split}$$

where we have utilized the property that the Euclidean projection is nonexpansive (Theorem 4.31), the subgradient inequality (6.33) for convex

functions, and the bounds on  $\alpha_k$  given by (6.44).

Our first convergence result is based on the divergent series step length formula (6.41), and establishes convergence to the optimal solution set  $X^*$  under an assumption on its boundedness. With the other two step length formulas, this condition will be possible to remove.

Recall the definition (3.11) of the minimum distance from a vector to a closed and convex set; our interest is in the distance from an arbitrary vector  $\boldsymbol{x} \in \mathbb{R}^n$  to the solution set  $X^*$ :

$$\operatorname{dist}_{X^*}(x) := \min_{y \in X^*} \|y - x\|.$$

**Theorem 6.23** (convergence of subgradient optimization methods, I) Let  $\{x_k\}$  be generated by the method (6.40), (6.41). If  $X^*$  is bounded and the sequence  $\{g_k\}$  is bounded, then  $\{f(x_k)\} \to f^*$  and  $\{\text{dist}_{X^*}(x_k)\} \to 0$  holds.

**Proof.** We show that the iterates will eventually belong to an arbitrarily small neighbourhood of the set of optimal solutions to (6.39).

Let  $\delta > 0$  and  $B^{\delta} := \{ \boldsymbol{x} \in \mathbb{R}^n \mid ||\boldsymbol{x}|| \leq \delta \}$ . Since f is convex, X is nonempty, closed and convex, and  $X^*$  is bounded, it follows from [Roc70, Theorem 27.2], applied to the lower semi-continuous, proper<sup>7</sup> and convex function  $f + \chi_X^{-8}$  that there exists an  $\varepsilon = \varepsilon(\delta) > 0$  such that the level set  $\{ \boldsymbol{x} \in X \mid f(\boldsymbol{x}) \leq f^* + \varepsilon \} \subseteq X^* + B^{\delta/2}$ ; this level set is denoted by  $X^{\varepsilon}$ . Moreover, since for all k,  $||\boldsymbol{g}_k|| \leq \sup_s \{||\boldsymbol{g}_s||\} < \infty$ , and  $\{\alpha_k\} \to 0$ , there exists an  $N(\delta)$  such that  $\alpha_k ||\boldsymbol{g}_k||^2 \leq \varepsilon$  and  $\alpha_k ||\boldsymbol{g}_k|| \leq \delta/2$  for all  $k \geq N(\delta)$ .

The sequel of the proof is based on induction and is organized as follows. In the first part, we show that there exists a finite  $k(\delta) \geq N(\delta)$ such that  $x_{k(\delta)} \in X^* + B^{\delta}$ . In the second part, we establish that if  $x_k$ belongs to  $X^* + B^{\delta}$  for some  $k \geq N(\delta)$  then so does  $x_{k+1}$ , by showing that either dist<sub>X\*</sub>( $x_{k+1}$ ) < dist<sub>X\*</sub>( $x_k$ ) holds, or  $x_k \in X^{\varepsilon}$  so that  $x_{k+1} \in$  $X^* + B^{\delta}$  since the step taken is not longer than  $\delta/2$ .

Let  $x^* \in X^*$  be arbitrary. In every iteration k we then have

$$\|\boldsymbol{x}^* - \boldsymbol{x}_{k+1}\|^2 = \|\boldsymbol{x}^* - \operatorname{Proj}_X (\boldsymbol{x}_k - \alpha_k \boldsymbol{g}_k)\|^2$$
 (6.45a)

$$\leq \|\boldsymbol{x}^* - \boldsymbol{x}_k + \alpha_k \boldsymbol{g}_k\|^2 \tag{6.45b}$$

$$= \|\boldsymbol{x}^* - \boldsymbol{x}_k\|^2 + \alpha_k \left( 2\boldsymbol{g}_k^{\mathrm{T}} \left( \boldsymbol{x}^* - \boldsymbol{x}_k \right) + \alpha_k \|\boldsymbol{g}_k\|^2 \right), \quad (6.45c)$$

<sup>&</sup>lt;sup>7</sup> A proper function is a function which is finite at least at some vector and nowhere attains the value  $-\infty$ . See also Section 1.4.

<sup>&</sup>lt;sup>8</sup>For any set  $S \subset \mathbb{R}^n$  the function  $\chi_S$  is the *indicator function* of the set S, that is,  $\chi_S(\boldsymbol{x}) = 0$  if  $\boldsymbol{x} \in S$ ; and  $\chi_S(\boldsymbol{x}) = +\infty$  if  $\boldsymbol{x} \notin S$ . See also Section 13.1.

\*Subgradient optimization methods

where the inequality follows from the projection property. Now, suppose

$$2 \boldsymbol{g}_{s}^{\mathrm{T}} \left( \boldsymbol{x}^{*} - \boldsymbol{x}_{s} \right) + \alpha_{s} \left\| \boldsymbol{g}_{s} \right\|^{2} < -\varepsilon$$

$$(6.46)$$

for all  $s \ge N(\delta)$ . Then, using (6.45) repeatedly, we obtain that for any  $k \ge N(\delta)$ ,

$$\|\boldsymbol{x}^* - \boldsymbol{x}_{k+1}\|^2 < \|\boldsymbol{x}^* - \boldsymbol{x}_{N(\delta)}\|^2 - \varepsilon \sum_{s=N(\delta)}^k \alpha_s,$$

and from (6.40) it follows that the right-hand side of this inequality tends to minus infinity as  $k \to \infty$ , which clearly is impossible. Therefore,

$$2 \boldsymbol{g}_{k}^{\mathrm{T}} \left( \boldsymbol{x}^{*} - \boldsymbol{x}_{k} \right) + \alpha_{k} \left\| \boldsymbol{g}_{k} \right\|^{2} \geq -\varepsilon$$

$$(6.47)$$

for at least one  $k \geq N(\delta)$ , say  $k = k(\delta)$ . From the definition of  $N(\delta)$ , it follows that  $\boldsymbol{g}_{k(\delta)}^{\mathrm{T}}(\boldsymbol{x}^* - \boldsymbol{x}_{k(\delta)}) \geq -\varepsilon$ . From the definition of a subgradient (cf. Definition 6.16) we have that  $f(\boldsymbol{x}^*) - f(\boldsymbol{x}_{k(\delta)}) \geq \boldsymbol{g}_{k(\delta)}^{\mathrm{T}}(\boldsymbol{x}^* - \boldsymbol{x}_{k(\delta)})$ , since  $\boldsymbol{x}^*, \boldsymbol{x}_{k(\delta)} \in X$ . Hence,  $f(\boldsymbol{x}_{k(\delta)}) \leq f^* + \varepsilon$ , that is,  $\boldsymbol{x}_{k(\delta)} \in X^{\varepsilon} \subseteq X^* + B^{\delta/2} \subset X^* + B^{\delta}$ .

Now, suppose that  $x_k \in X^* + B^{\delta}$  for some  $k \ge N(\delta)$ . If (6.46) holds for s = k, then, by using (6.45), we have that  $||x^* - x_{k+1}|| < ||x^* - x_k||$  for any  $x^* \in X^*$ . Hence,

$$\begin{aligned} \operatorname{dist}_{X^*}(\boldsymbol{x}_{k+1}) &\leq \|\operatorname{Proj}_{X^*}(\boldsymbol{x}_k) - \boldsymbol{x}_{k+1}\| < \|\operatorname{Proj}_{X^*}(\boldsymbol{x}_k) - \boldsymbol{x}_k\| \\ &= \operatorname{dist}_{X^*}(\boldsymbol{x}_k) \leq \delta. \end{aligned}$$

Thus,  $x_{k+1} \in X^* + B^{\delta}$ . Otherwise, (6.47) must hold and, using the same arguments as above, we obtain that  $f(x_k) \leq f^* + \varepsilon$ , i.e.,  $x_k \in X^{\varepsilon} \subseteq x^* + B^{\delta/2}$ . As

$$\begin{aligned} \|\boldsymbol{x}_{k+1} - \boldsymbol{x}_k\| &= \|\operatorname{Proj}_X (\boldsymbol{x}_k - \alpha_k \boldsymbol{g}_k) - \boldsymbol{x}_k\| \le \|\boldsymbol{x}_k - \alpha_k \boldsymbol{g}_k - \boldsymbol{x}_k\| \\ &= \alpha_k \|\boldsymbol{g}_k\| \le \delta/2 \end{aligned}$$

whenever  $k \ge N(\delta)$ , it follows that  $x_{k+1} \in X^* + B^{\delta/2} + B^{\delta/2} = X^* + B^{\delta}$ .

By induction with respect to  $k \ge k(\delta)$ , it follows that  $x_k \in X^* + B^{\delta}$  for all  $k \ge k(\delta)$ . Since this holds for arbitrarily small values of  $\delta > 0$  and f is continuous, the theorem follows.

We next introduce the additional requirement (6.42); the resulting algorithm's convergence behaviour is now much more favourable, and the proof is at the same time less technical.

**Theorem 6.24** (convergence of subgradient optimization methods, II) Let  $\{x_k\}$  be generated by the method (6.40), (6.41), (6.42). If  $X^*$  is nonempty and the sequence  $\{g_k\}$  is bounded, then  $\{f(x_k)\} \to f^*$  and  $\{x_k\} \to x^* \in X^*$  holds.

**Proof.** Let  $x^* \in X^*$  and  $k \ge 1$ . Repeated application of (6.45) yields

$$\|\boldsymbol{x}^{*} - \boldsymbol{x}_{k}\|^{2} \leq \|\boldsymbol{x}^{*} - \boldsymbol{x}_{0}\|^{2} + 2\sum_{s=0}^{k-1} \alpha_{s} \boldsymbol{g}_{s}^{\mathrm{T}} (\boldsymbol{x}^{*} - \boldsymbol{x}_{s}) + \sum_{s=0}^{k-1} \alpha_{s}^{2} \|\boldsymbol{g}_{s}\|^{2} .(6.48)$$

Since  $x^* \in X^*$  and  $g_s \in \partial f(x_s)$  for all  $s \ge 0$  we obtain that

$$f(x_s) \ge f^* \ge f(x_s) + g_s^{\mathrm{T}} (x^* - x_s), \qquad s \ge 0,$$
 (6.49)

and hence that  $\boldsymbol{g}_s^T(\boldsymbol{x}^* - \boldsymbol{x}_s) \leq 0$  for all  $s \geq 0$ . Define  $c := \sup_k \{ \|\boldsymbol{g}_k\| \}$ and  $p = \sum_{k=0}^{\infty} \alpha_k^2$ , so that  $\|\boldsymbol{g}_s\| \leq c$  for any  $s \geq 0$  and  $\sum_{s=0}^{k-1} \alpha_s^2 < p$ . From (6.48) we then conclude that  $\|\boldsymbol{x}^* - \boldsymbol{x}_k\|^2 < \|\boldsymbol{x}^* - \boldsymbol{x}_0\|^2 + pc^2$  for any  $k \geq 1$ , and thus that the sequence  $\{\boldsymbol{x}_k\}$  is bounded.

Assume now that there is no subsequence  $\{\boldsymbol{x}_{k_i}\}$  of  $\{\boldsymbol{x}_k\}$  with  $\{\boldsymbol{g}_{k_i}^{\mathrm{T}}(\boldsymbol{x}^* - \boldsymbol{x}_{k_i})\} \to 0$ . Then there must exist an  $\varepsilon > 0$  with  $\boldsymbol{g}_s^{\mathrm{T}}(\boldsymbol{x}^* - \boldsymbol{x}_s) \leq -\varepsilon$  for all sufficiently large values of s. From (6.48) and the conditions on the step lengths it follows that  $\{\|\boldsymbol{x}^* - \boldsymbol{x}_s\|\} \to -\infty$ , which clearly is impossible. The sequence  $\{\boldsymbol{x}_k\}$  must therefore contain a subsequence  $\{\boldsymbol{x}_{k_i}\}$  such that  $\{\boldsymbol{g}_{k_i}^{\mathrm{T}}(\boldsymbol{x}^* - \boldsymbol{x}_{k_i})\} \to 0$ . From (6.49) it follows that  $\{f(\boldsymbol{x}_{k_i})\} \to f^*$ . The boundedness of  $\{\boldsymbol{x}_k\}$  implies the existence of an accumulation point of the subsequence  $\{\boldsymbol{x}_{k_i}\}$ , say  $\boldsymbol{x}^{\infty}$ . From the continuity of f it follows that  $\boldsymbol{x}^{\infty} \in X^*$ .

To show that  $\boldsymbol{x}^{\infty}$  is the only accumulation point of  $\{\boldsymbol{x}_k\}$ , let  $\delta > 0$ and choose an  $M(\delta)$  such that  $\|\boldsymbol{x}^{\infty} - \boldsymbol{x}_{M(\delta)}\|^2 \leq \delta/2$  and  $\sum_{s=M(\delta)}^{\infty} \alpha_s^2 \leq \delta/(2c^2)$ . Consider any  $k > M(\delta)$ . Analogously to the derivation of (6.48), and using (6.49), we then obtain that

$$\|\boldsymbol{x}^{\infty} - \boldsymbol{x}_{k}\|^{2} \leq \|\boldsymbol{x}^{\infty} - \boldsymbol{x}_{M(\delta)}\|^{2} + \sum_{s=M(\delta)}^{k-1} \alpha_{s}^{2} \|\boldsymbol{g}_{s}\|^{2} < \frac{\delta}{2} + \frac{\delta}{2c^{2}}c^{2} = \delta.$$

Since this holds for arbitrarily small values of  $\delta > 0$ , we are done.

Note that the boundedness condition on  $\{g_k\}$  is fulfilled whenever we know before-hand that the sequence  $\{x_k\}$  is bounded, such as in the case when X itself is bounded; cf. Proposition 6.17(a).

We finally present the even stronger convergence properties of the subgradient projection method using the Polyak step.

**Theorem 6.25** (convergence of subgradient optimization methods, III) Let  $\{x_k\}$  be generated by the method (6.40), (6.43). If  $X^*$  is nonempty then  $\{f(x_k)\} \rightarrow f^*$  and  $\{x_k\} \rightarrow x^* \in X^*$  holds.

**Proof.** From Proposition 6.22 follows that the sequence  $\{||\mathbf{x}_k - \mathbf{x}^*||\}$  is strictly decreasing for every  $\mathbf{x}^* \in X^*$ , and therefore has a limit. By construction of the step length, in which the step lengths are bounded away from zero and  $2[f(\mathbf{x}_k) - f^*]/||\mathbf{g}_k||^2$ , it follows from the proof of Proposition 6.22 that  $\{[f(\mathbf{x}_k) - f^*]^2/||\mathbf{g}_k||^2\} \to 0$  must hold. Since  $\{\mathbf{g}_k\}$  must be bounded due to the boundedness of  $\{\mathbf{x}_k\}$  [Proposition 6.17(a)], we have that  $\{f(\mathbf{x}_k)\} \to f^*$ . Further,  $\{\mathbf{x}_k\}$  is bounded, and due to the continuity property of f every limit point must then belong to  $X^*$ .

It remains to show that there can be only one limit point. This property follows from the monotone decrease of the distance  $||x_k - x^*||$ . In detail, the proof is as follows. Suppose two subsequences of  $\{x_k\}$  exist, such that they converge to two different vectors in  $X^*$ :

$$\{x_{m_i}\} o x_1^*; \qquad \{x_{l_i}\} o x_2^*; \qquad x_1^* 
eq x_2^*.$$

We must then have  $\{\|\boldsymbol{x}_{l_i} - \boldsymbol{x}_1^*\|\} \to \rho > 0$ . Since  $\boldsymbol{x}_1^* \in X^*$  and the distance to  $X^*$  is decreasing,  $\{\|\boldsymbol{x}_k - \boldsymbol{x}_1^*\|\} \to \rho$  holds, and in particular  $\{\|\boldsymbol{x}_{m_i} - \boldsymbol{x}_1^*\|\} \to \rho$ , which is a contradiction.

Contrary to the slow convergence of the subgradient projection algorithms that rely on the divergent series step length rule, under additional conditions on the function f a subgradient algorithm based on the Polyak step length (6.43) is geometrically convergent, in the sense that there exist c > 0 and  $\eta \in (0, 1)$  with

$$\|x_k - x^*\| \le c\eta^k, \qquad k = 0, 1, \dots$$

See Section 6.9 for references to other subgradient algorithms than those presented here.

#### 6.5.2 Application to the Lagrangian dual problem

We remind ourselves that the Lagrangian dual problem is a concave maximization problem, and that the appearance of the dual function is similar to that of the following example:

Let  $h(x) := \min \{h_1(x), h_2(x)\}$ , where  $h_1(x) := 4 - |x|$  and  $h_2(x) := 4 - (x-2)^2$ . Then,

$$h(x) = \begin{cases} 4 - x, & \text{if } 1 \le x \le 4; \\ 4 - (x - 2)^2 & \text{if } x \le 1, \ x \ge 4; \end{cases}$$



Figure 6.3: A convex min-function with three pieces.

cf. Figure 6.3.

The function h is non-differentiable at x = 1 and x = 4, since its graph has non-unique supporting hyperplanes there:

$$\partial h(x) = \begin{cases} \{4-2x\}, & \text{if } x < 1; \\ [-1,2], & \text{if } x = 1; \\ \{-1\}, & \text{if } 1 < x < 4; \\ [-4,-1], & \text{if } x = 4; \\ \{4-2x\}, & \text{if } x > 4. \end{cases}$$

the subdifferential is here either a singleton (at differentiable points) or a closed interval (at non-differentiable points).

Note the monotonically decreasing nature of the relation  $x \mapsto \partial h(x)$ . Note also that  $0 \in \partial h(1)$ , whence  $x^* = 1$  defines a maximum over  $\mathbb{R}$ .

Now, let  $g \in \partial q(\bar{\mu})$ , and let  $U^*$  be the set of optimal solutions to (6.10). Then,

$$U^* \subseteq \{ \boldsymbol{\mu} \in \mathbb{R}^m \mid \boldsymbol{g}^{\mathrm{T}}(\boldsymbol{\mu} - \bar{\boldsymbol{\mu}}) \ge 0 \}.$$

In other words, any subgradient defines a half-space that contains the set of optimal solutions; cf. Figure 6.4. We therefore know that a small enough step in the direction of a subgradient gets us closer to the set of optimal solutions; cf. Proposition 6.22. But again consider Figure 6.4: an arbitrary subgradient, like the on depicted, may not define an ascent direction! As we saw in the previous section, convergence must be based on other arguments, like the decreasing distance to  $U^*$  alluded to above and in the previous section. In the next subsection we discuss in brief the generation of ascent directions.

We consider the Lagrangian dual problem (6.10). We suppose, as in the previous section, that X is compact so that the infimum in (6.9) is attained for every  $\mu \geq 0^m$  (which is the set over which we wish to maximize q) and q is real-valued over  $\mathbb{R}^m_+$ .

\*Subgradient optimization methods



Figure 6.4: The half-space defined by the subgradient g of q at  $\mu$ . Note that the subgradient is not an ascent direction.

In the case of our special concave maximization problem, the iteration has the form

$$\boldsymbol{\mu}_{k+1} = \operatorname{Proj}_{\mathbb{R}^m_+} [\boldsymbol{\mu}_k + \alpha_k \boldsymbol{g}_k] = [\boldsymbol{\mu}_k + \alpha_k \boldsymbol{g}_k]_+$$
$$= (\operatorname{maximum} \{0, (\boldsymbol{\mu}_k)_i + \alpha_k (\boldsymbol{g}_k)_i\})_{i=1}^m, \qquad (6.50)$$

where  $\boldsymbol{g}_k \in \partial q(\boldsymbol{\mu}_k)$  is arbitrarily chosen; we would typically use  $\boldsymbol{g}_k = \boldsymbol{g}(\boldsymbol{x}_k)$ , where  $\boldsymbol{x}_k \in \operatorname{argminimum}_{\boldsymbol{x} \in X} L(\boldsymbol{x}, \boldsymbol{\mu}_k)$ . The projection operation onto the first orthant is, as we can see, very simple.

Replacing the Polyak step (6.43) with the corresponding dual form

$$\alpha_k = \theta_k \frac{q^* - q(\boldsymbol{\mu}_k)}{\|\boldsymbol{g}_k\|^2}, \qquad 0 < \sigma_1 \le \theta_k \le 2 - \sigma_2 < 2, \tag{6.51}$$

convergence will now be a simple consequence of the above theorems.

The compactness condition (6.37) and the fact that the feasible set of (6.4) is nonempty ensure that the problem (6.4) has an optimal solution; in particular, the feasibility condition (6.5) then holds. Further, if we introduce the Slater condition (6.16), we are ensured that there is no duality gap, and that the dual problem (6.10) has a compact set  $U^*$ of optimal solutions. Under these assumptions, we have the following results for subgradient optimization methods.

**Theorem 6.26** (convergence of subgradient optimization methods) Suppose that the problem (6.4) is feasible, and that the compactness condition (6.37) and the Slater condition (6.16) hold.

(a) Let  $\{\boldsymbol{\mu}_k\}$  be generated by the method (6.50), (6.41). Then,  $\{q(\boldsymbol{\mu}_k)\} \rightarrow q^*$ , and  $\{\operatorname{dist}_{U^*}(\boldsymbol{\mu}_k)\} \rightarrow 0$ .

(b) Let  $\{\mu_k\}$  be generated by the method (6.50), (6.41), (6.42). Then,  $\{\mu_k\}$  converges to an optimal solution to (6.10).

(c) Let  $\{\mu_k\}$  be generated by the method (6.50), (6.51). Then,  $\{\mu_k\}$  converges to an optimal solution to (6.10).

**Proof.** The results follow from Theorems 6.23, 6.24, and 6.25, respectively. Note that in the first two cases, boundedness conditions were assumed for  $X^*$  and the sequence of subgradients. The corresponding conditions for the Lagrangian dual problem are fulfilled under the CQs imposed, since they imply that the search for an optimal solution is done over a compact set; cf. Theorem 6.9(a) and its proof.

#### 6.5.3 The generation of ascent directions

Proposition 6.18 shows that the existence of a descent direction with respect to the convex function  $f : \mathbb{R}^n \to \mathbb{R}$  at some  $\bar{x} \in \mathbb{R}^n$  hinges on the existence of some vector  $\bar{p} \in \mathbb{R}^n$  such that  $f'(\bar{x}; \bar{p}) < 0$ . According to the definition of the directional derivative and the compactness of  $\partial f(\bar{x})$ , this is equivalent to the statement that  $g^T \bar{p} \leq \varepsilon < 0$  for every  $g \in \partial f(\bar{x})$ . In the context of Lagrangian duality we show below how we can generate an ascent directions for q at some  $\mu \in \mathbb{R}^m$ .

**Definition 6.27** (steepest ascent direction) Suppose that the problem (6.4) is feasible, and that the compactness condition (6.37) holds. Consider the Lagrangian dual problem (6.10), and let  $\mu \in \mathbb{R}^m$ . A vector  $\bar{p} \in \mathbb{R}^m$  is a steepest ascent direction if

$$q'(\boldsymbol{\mu}; \bar{\boldsymbol{p}}) = \underset{\|\boldsymbol{p}\| \leq 1}{\operatorname{maximum}} q'(\boldsymbol{\mu}; \boldsymbol{p})$$

holds.

**Proposition 6.28** (the shortest subgradient yields the steepest ascent direction) Suppose that the problem (6.4) is feasible, and that the compactness condition (6.37) holds. Consider the Lagrangian dual problem (6.10). The direction  $\bar{p}$  of steepest ascent with respect to q at  $\mu$  is given below, where  $\bar{g} \in \partial q(\mu)$  is the shortest subgradient in  $\partial q(\mu)$  with respect to the Euclidean norm:

$$ar{p} = egin{cases} \mathbf{0}^m, & ext{if } ar{m{g}} = \mathbf{0}^m, \ ar{m{g}} = ar{m{g}}, \ ar{m{g}} = ar{m{g}}, & ext{if } ar{m{g}} 
eq \mathbf{0}^m, \ ar{m{g}} 
eq$$

**Proof.** By Definition 6.27 and Proposition 6.19(e), the following string of equalities and inequalities can easily be verified:

$$\underset{\|\boldsymbol{p}\|\leq 1}{\operatorname{maximum}} q'(\boldsymbol{\mu}; \boldsymbol{p}) = \underset{\|\boldsymbol{p}\|\leq 1}{\operatorname{maximum}} \underset{\boldsymbol{g}\in\partial q(\boldsymbol{\mu})}{\operatorname{minfmum}} \boldsymbol{g}^{\mathrm{T}}\boldsymbol{p}$$

$$\leq \underset{\boldsymbol{g}\in\partial q(\boldsymbol{\mu})}{\operatorname{minfmum}} \underset{\|\boldsymbol{p}\|\leq 1}{\operatorname{maximum}} \boldsymbol{g}^{\mathrm{T}}\boldsymbol{p}$$

$$= \underset{\boldsymbol{g}\in\partial q(\boldsymbol{\mu})}{\operatorname{minfmum}} \|\boldsymbol{g}\|$$

$$= \|\bar{\boldsymbol{g}}\|.$$

$$(6.52)$$

If we can construct a direction  $\bar{p}$  such that  $q'(\mu; \bar{p}) = ||\bar{g}||$  then by (6.52)  $\bar{p}$  is the steepest ascent direction. If  $\bar{g} = \mathbf{0}^m$  then for  $\bar{p} = \mathbf{0}^m$  we obviously have that  $q'(\mu; \bar{p}) = ||\bar{g}||$ . Suppose then that  $\bar{p} \neq \mathbf{0}^m$ , and let  $\bar{p} := \bar{g}/||\bar{g}||$ . Note that

$$q'(\boldsymbol{\mu}; \boldsymbol{p}) = \inf_{\boldsymbol{g} \in \partial_q(\boldsymbol{\mu})} \boldsymbol{g}^{\mathrm{T}} \boldsymbol{\bar{p}} = \inf_{\boldsymbol{g} \in \partial_q(\boldsymbol{\mu})} \frac{\boldsymbol{\bar{g}}^{\mathrm{T}} \boldsymbol{g}}{\|\boldsymbol{g}\|}$$
$$= \frac{1}{\|\boldsymbol{g}\|} \inf_{\boldsymbol{g} \in \partial_q(\boldsymbol{\mu})} \{\|\boldsymbol{\bar{g}}\|^2 + \boldsymbol{\bar{g}}^{\mathrm{T}}(\boldsymbol{g} - \boldsymbol{\bar{g}})\}$$
$$= \|\boldsymbol{\bar{g}}\| + \frac{1}{\|\boldsymbol{g}\|} \inf_{\boldsymbol{g} \in \partial_q(\boldsymbol{\mu})} \boldsymbol{\bar{g}}^{\mathrm{T}}(\boldsymbol{g} - \boldsymbol{\bar{g}}).$$
(6.53)

Since  $\bar{\boldsymbol{g}}$  is the shortest vector in  $\partial q(\boldsymbol{\mu})$ , then, by the variational inequality characterization of the projection of  $\boldsymbol{0}^m$  onto  $\partial q(\boldsymbol{\mu})$  established in Theorem 4.23, we obtain that  $\bar{\boldsymbol{g}}^{\mathrm{T}}(\boldsymbol{g}-\bar{\boldsymbol{g}}) \geq 0$  for every  $\boldsymbol{g} \in \partial q(\boldsymbol{\mu})$ . Hence, infimum $_{\boldsymbol{g} \in \partial q(\boldsymbol{\mu})} \bar{\boldsymbol{g}}^{\mathrm{T}}(\boldsymbol{g}-\bar{\boldsymbol{g}}) = 0$  is achieved at  $\bar{\boldsymbol{g}}$ . From (6.53) it then follows that  $q'(\boldsymbol{\mu}, \bar{\boldsymbol{p}}) = \|\bar{\boldsymbol{g}}\|$ . We are done.

## 6.6 \*Obtaining a primal solution

It remains for us to show how an optimal dual solution  $\mu^*$  can be translated into an optimal primal solution  $x^*$ . Obviously, convexity and strong duality will be needed in general, if we are to be able to utilize the primal-dual optimality characterization in Theorem 6.7. It turns out that the generation of a primal optimum is automatic if q is differentiable at  $\mu^*$ , which is also the condition under which the famous Lagrange multiplier method works. Unfortunately, in many cases, such as for most non-strictly convex optimization problems (like linear programming), this will not be the case, and then the translation work becomes more complex.

We start with the ideal case.

#### 6.6.1 Differentiability at the optimal solution

The following results summarize the optimality conditions for the Lagrangian dual problem (6.10), and their consequences for the availability of a primal optimal solution in the absence of a duality gap.

**Proposition 6.29** (optimality conditions for the dual problem) Suppose that, in the problem (6.4), the compactness condition (6.37) holds. Suppose further that the Lagrangian dual problem has an optimal solution,  $\mu^*$ .

(a) The dual optimal solution is characterized by the inclusion

$$\mathbf{0}^m \in -\partial q(\boldsymbol{\mu}^*) + N_{\mathbb{R}^m_+}(\boldsymbol{\mu}^*). \tag{6.54}$$

In other words, there then exists  $\gamma^* \in \partial q(\mu^*)$ —an optimality-characterizing subgradient of q at  $\mu^*$ —such that

$$\mathbf{0}^m \le \boldsymbol{\mu}^* \perp \boldsymbol{\gamma}^* \le \mathbf{0}^m. \tag{6.55}$$

There exists a finite set of solutions  $x^i \in X(\mu^*)$  (i = 1, ..., k) where  $k \leq m + 1$  such that

$$\boldsymbol{\gamma}^* = \sum_{i=1}^k \alpha_i \boldsymbol{g}(\boldsymbol{x}^i); \qquad \sum_{i=1}^k \alpha_i = 1; \qquad \alpha_i \ge 0, \ i = 1, \dots, k.$$
 (6.56)

Hence, we have that

$$\sum_{i=1}^{k} \alpha_i \mu_i^* g_i(\boldsymbol{x}^i) = 0, \qquad j = 1, \dots, m.$$
(6.57)

(b) If there is a duality gap, then q is non-differentiable at  $\mu^*$ .

(c) If q is differentiable at  $\mu^*$ , then there is no duality gap. Further, any vector in  $X(\mu^*)$  then solves the primal problem (6.4).

**Proof.** (a) The first result is a direct statement of the optimality conditions of the convex and subdifferentiable program (6.10); the complementarity conditions in (6.55) are an equivalent statement of the inclusion in (6.54).

The second result is an application of Carathéodory's Theorem 3.8 to the compact and convex set  $\partial q(\mu^*)$ .

(b) The result is established once (c) is.

(c) Let  $\bar{x}$  be any vector in  $X(\mu^*)$  for which  $\nabla q(\mu^*) = g(\bar{x})$  holds, cf. Proposition 6.20(a). We obtain from (6.55) that

$$\mathbf{0}^m \leq \boldsymbol{\mu}^* \perp \boldsymbol{g}(\bar{\boldsymbol{x}}) \leq \mathbf{0}^m.$$

Hence, the pair  $(\boldsymbol{\mu}, \bar{\boldsymbol{x}})$  fulfills all the conditions stated in (6.12), so that, by Theorem 6.7,  $\bar{\boldsymbol{x}}$  is an optimal solution to (6.4).

\*Obtaining a primal solution

**Remark 6.30** (the non-coordinability phenomenon and decomposition algorithms) Many interesting problems do not comply with the conditions in (c); for example, linear programming is one where the Lagrangian dual problem often is non-differentiable at every dual optimal solution.<sup>9</sup> This is sometimes called the *non-coordinability phenomenon* (cf. [Las70, DiJ79]). It was in order to cope with this phenomenon that Dantzig–Wolfe decomposition ([DaW60, Las70]) and other column generation algorithms, Benders decomposition ([Ben62, Las70]) and generalized linear programming were developed; noticing that the convex combination of a finite number of candidate primal solutions are sufficient to verify an optimal primal–dual solution [cf. (6.57)], methodologies were developed to generate those vectors algorithmically. See also [LPS99] for overviews on the subject of generating primal optimal solutions from dual optimal ones, and [BSS93, Theorem 6.5.2] for an LP procedure that provides primal feasible solutions for convex programs.

Note that the equation (6.57) in (a) reduces to the complementarity condition that  $\mu_i^* g_i(\bar{x}) = 0$  holds, for the averaged solution,  $\bar{x} := \sum_{i=1}^k \alpha_i x^i$ , whenever all the functions  $g_i$  are affine.

#### 6.6.2 Everett's Theorem

The next result shows that the solution to the Lagrangian subproblem solves a perturbed version of the original problem. We state the result for the general problem to find

$$f^* := \inf_{x} \inf_{x} f(x), \qquad (6.58)$$
  
subject to  $x \in X,$   
 $g_i(x) \le 0, \qquad i = 1, \dots, m,$   
 $h_j(x) = 0, \qquad j = 1, \dots, \ell,$ 

where  $f : \mathbb{R}^n \to \mathbb{R}$ ,  $g_i : \mathbb{R}^n \to \mathbb{R}$  (i = 1, 2, ..., m), and  $h_j : \mathbb{R}^n \to \mathbb{R}$  $(j = 1, 2, ..., \ell)$  are given functions, and  $X \subseteq \mathbb{R}^n$ .

**Theorem 6.31 (Everett's Theorem)** Let  $(\mu, \lambda) \in \mathbb{R}^m_+ \times \mathbb{R}^\ell$ . Consider the Lagrangian subproblem to

$$\underset{x \in X}{\text{minimize}} \left\{ f(x) + \boldsymbol{\mu}^{\mathrm{T}} \boldsymbol{g}(x) + \boldsymbol{\lambda}^{\mathrm{T}} \boldsymbol{h}(x) \right\}.$$
(6.59)

Suppose that  $\bar{x}$  is an optimal solution to this problem, and let  $\mathcal{I}(\mu) \subseteq \{1, \ldots, m\}$  denote the set of indices *i* for which  $\mu_i > 0$ .

<sup>&</sup>lt;sup>9</sup>In other words, even if a Lagrange multiplier vector is known, the Lagrangian subproblem may not identify a primal optimal solution.

(a)  $\bar{x}$  is an optimal solution to the perturbed primal problem to

$$\begin{array}{ll} \underset{x}{\operatorname{minimize}} & f(x), & (6.60) \\ \text{subject to} & x \in X, \\ & g_i(x) \leq g_i(\bar{x}), & i \in \mathcal{I}(\bar{x}), \\ & h_j(x) = h_j(\bar{x}), & j = 1, \dots, \ell. \end{array}$$

(b) If  $\bar{x}$  is feasible in (6.58) and  $\mu^{\mathrm{T}} g(\bar{x}) = 0$  holds, then  $\bar{x}$  solves (6.58), and the pair  $(\mu, \lambda)$  then solves the Lagrangian dual problem.

**Proof.** (a) The proof proceeds by showing that the triple  $(\bar{x}, \mu, \lambda)$  is a saddle point of the function  $(x, \mu, \lambda) \mapsto f(x) + \mu^{\mathrm{T}}[g(x) - g(\bar{x})] + \lambda^{\mathrm{T}}[h(x) - h(\bar{x})]$  over  $X \times \mathbb{R}^m_+ \times \mathbb{R}^\ell$ .

Let x satisfy the constraints of (6.60). Since we have that  $h(x) = h(\bar{x})$  and  $\mu^{T}g(x) \leq \mu^{T}g(\bar{x})$ , the optimality of  $\bar{x}$  in (6.59) yields

$$\begin{split} f(x) + \mu^{\mathrm{T}} g(\bar{x}) + \lambda^{\mathrm{T}} h(\bar{x}) &\geq f(x) + \mu^{\mathrm{T}} g(x) + \lambda^{\mathrm{T}} h(x) \\ &\geq f(\bar{x}) + \mu^{\mathrm{T}} g(\bar{x}) + \lambda^{\mathrm{T}} h(\bar{x}), \end{split}$$

which shows that  $f(x) \ge f(\bar{x})$ . We are done.

(b)  $\boldsymbol{\mu}^{\mathrm{T}} \boldsymbol{g}(\bar{\boldsymbol{x}}) = 0$  implies that  $g_i(\bar{\boldsymbol{x}}) = 0$  for  $i \in \mathcal{I}(\boldsymbol{\mu})$ ; from (a)  $\bar{\boldsymbol{x}}$  solves the problem to

$\min_{\boldsymbol{x}} \operatorname{minimize} \ f(\boldsymbol{x}),$		(6.61)
subject to $x \in X$ ,		
$g_i(\boldsymbol{x}) \leq 0,$	$i \in \mathcal{I}(\bar{x}),$	
$h_j(x) = 0,$	$j=1,\ldots,\ell.$	

In particular, then, since the feasible set of (6.58) is contained in that of (6.61) and  $\bar{x}$  is feasible in the former,  $\bar{x}$  must also solve (6.58). That the pair  $(\mu, \lambda)$  solves the dual problem follows by the equality between the primal and dual objective functions at  $(\bar{x}, \mu, \lambda)$ , and weak duality.

One important consequence of the result is that if the right-hand side perturbations  $g_i(\bar{x})$  and  $h_i(\bar{x})$  all are close to zero, the vector  $\bar{x}$  being near-feasible might mean that it is in fact acceptable as an approximate solution to the original problem. (This interpretation hinges on the dualized constraints being *soft* constraints, in the sense that a small violation is acceptable. See Section 1.8 for an introduction to the topic of soft constraints.)

## 6.7 \*Sensitivity analysis

#### 6.7.1 Analysis for convex problems

Consider the inequality constrained convex program (6.4), where  $f : \mathbb{R}^n \to \mathbb{R}$  and  $g_i$  (i = 1, ..., m) are convex functions and  $X \subseteq \mathbb{R}^n$  is a convex set. Suppose that the problem (6.4) is feasible, and that the compactness condition (6.37) and Slater condition (6.16) hold. This is the classic case where there exist multiplier vectors  $\mu^*$ , according to Theorem 6.9, and strong duality holds.

For certain types of problems where the duality gap is zero and where there exist primal-dual optimal solutions, we have access to a beautiful theory of sensitivity analysis. The classic meaning of the term is the answer to the following question: what is the rate of change in  $f^*$  when a constraint right-hand side changes? This question answers important practical questions, like the following in manufacturing: If we buy one unit of additional resource at a given price, or if the demand of a product that we sell increases by a certain amount, then how much additional profit do we make?

We will here provide a basic result which states when this sensitivity analysis of the optimal objective value can be performed for the problem (6.4), and establish that the answer is determined precisely by the value of the Lagrange multiplier vector  $\mu^*$ , provided that it is unique.

**Definition 6.32 (perturbation function)** Consider the function  $p : \mathbb{R}^m \to \mathbb{R} \cup \{\pm \infty\}$  defined by

$$p(\boldsymbol{u}) := \inf_{\boldsymbol{x}} \inf_{\boldsymbol{x}} f(\boldsymbol{x}), \tag{6.62}$$

subject to 
$$x \in X$$
, (6.63)

 $g_i(\boldsymbol{x}) \leq u_i, \quad i=1,\ldots,m, \qquad \boldsymbol{u} \in \mathbb{R}^m;$ 

it is called the perturbation function, or primal function, associated with the problem (6.4). Its effective domain is the set  $P := \{ u \in \mathbb{R}^m \mid p(u) < +\infty \}$ .

Under the above convexity conditions, we can establish that p is a convex function. Indeed, it holds that for any value of the Lagrange

multiplier vector  $\mu^*$  for the problem (6.4) that

$$\begin{split} q(\boldsymbol{\mu}^{*}) &= \inf_{\substack{x \in X}} \{f(x) + (\boldsymbol{\mu}^{*})^{\mathrm{T}} g(x)\} \\ &= \inf_{\{(u,x) \in P \times X | g(x) \leq u\}} \{f(x) + (\boldsymbol{\mu}^{*})^{\mathrm{T}} g(x)\} \\ &= \inf_{\{(u,x) \in P \times X | g(x) \leq u\}} \{f(x) + (\boldsymbol{\mu}^{*})^{\mathrm{T}} u\} \\ &= \inf_{\substack{u \in P}} \inf_{\substack{x \in X | g(x) \leq u\}}} \{f(x) + (\boldsymbol{\mu}^{*})^{\mathrm{T}} u\}. \end{split}$$

Since  $\mu^*$  is assumed to be a Lagrange multiplier vector, we have that  $q(\mu^*) = f^* = p(\mathbf{0}^m)$ . By the definition of infimum, then, we have that

$$p(\mathbf{0}^m) \leq p(\mathbf{u}) + (\boldsymbol{\mu}^*)^{\mathrm{T}} \boldsymbol{u}, \qquad \boldsymbol{u} \in \mathbb{R}^m$$

that is,  $-\mu^*$  (notice the sign!) is a subgradient of p at  $u = 0^m$  (see Definition 6.16). Moreover, by the result in Proposition 6.17(c), p is differentiable at  $0^m$  if and only if p is finite in a neighbourhood of  $0^m$  and  $\mu^*$  is a unique Lagrange multiplier vector, that is, the Lagrangian dual problem (6.10) has a unique optimal solution. We have therefore proved the following result:

**Proposition 6.33** (a sensitivity analysis result) Suppose that in the inequality constrained problem (6.4),  $f : \mathbb{R}^n \to \mathbb{R}$  and  $g_i : \mathbb{R}^n \to \mathbb{R}$ (i = 1, ..., m) are convex functions and  $X \subseteq \mathbb{R}^n$  is a convex set. Suppose that the problem (6.4) is feasible, and that the compactness assumption (6.37) and Slater condition (6.16) hold. Suppose further that the perturbed problem defined in (6.62) has an optimal solution in a neighbourhood of  $u = 0^m$ , and that on the set of primal-dual optimal solutions to (6.4)-(6.10), the dual optimal solution  $\mu^*$  is unique. Then, the perturbation function p is differentiable at  $u = 0^m$ , and

$$\nabla p(\mathbf{0}^m) = -\boldsymbol{\mu}^*$$

holds.

It is intuitive that the sign of  $\nabla p(\mathbf{0}^m)$  should be non-positive; if a right-hand side of the (less-than) inequality constraints in (6.4) increases, then the feasible set becomes larger. [This means that we might be able to find feasible vectors  $\boldsymbol{x}$  in the new problem with  $f(\boldsymbol{x}) < f^*$ , where  $f^* = p(0)$  is the optimal value of the minimization problem (6.4).]

The result specializes immediately to linear programming problems, which is the problem type where this type of analysis is most often utilized. The proof of differentiability of the perturbation function at zero for that special case can however be done much more simply. (See Section 10.3.1.)

\*Sensitivity analysis

#### 6.7.2 Analysis for differentiable problems

There exist local versions of the analysis valid also for non-convex problems, where we are interested in the effect of a problem perturbation on a KKT point. A special such analysis was recently performed by Bertsekas [Ber04], in which he shows that even when the problem is non-convex and the set of Lagrange multipliers are not unique, a sensitivity analysis is available as long as data is differentiable. Suppose then that in the problem (6.4) the functions f and  $g_i$ ,  $i = 1, \ldots, m$  are in  $C^1$  and that X is nonempty. We generalize the concept of a Lagrange multiplier vector to here mean that it is a vector  $\mu^*$  associated with a local minimum  $x^*$  such that

$$\left(\nabla f(\boldsymbol{x}^*) + \sum_{i=1}^m \mu_i^* \nabla g_i(\boldsymbol{x}^*)\right)^{\mathrm{T}} \boldsymbol{p} \ge 0, \qquad \boldsymbol{p} \in T_X(\boldsymbol{x}^*), \qquad (6.64a)$$

$$\mu_i^* \ge 0, \qquad i = 1, \dots, m, \qquad (6.64b)$$

$$\mu_i^* = 0, \qquad i \notin \mathcal{I}(\boldsymbol{x}^*), \qquad (6.64c)$$

where  $T_X(\boldsymbol{x}^*)$  is the tangent cone to X at  $\boldsymbol{x}^*$  (cf. Definition 5.2). Note that under an appropriate CQ this is equivalent to the KKT conditions, in which case we are simply requiring here that  $\boldsymbol{x}^*$  is a local minimum.

In the below result we utilize the notation

$$g_i^+(x) := ext{maximum} \{0, g_i(x)\}, \qquad i = 1, \dots, m,$$

and let  $g^+(x)$  be the *m*-vector of elements  $g_i^+(x)$ , i = 1, ..., m.

**Theorem 6.34** (sensitivity from the minimum norm multiplier vector) Suppose that  $x^*$  is a local minimum in the problem (6.4), and that the set of Lagrange multipliers is nonempty. Let  $\mu^*$  denote the Lagrange multiplier vector of minimum Euclidean norm. Then, for every sequence  $\{x_k\} \subset X$ of infeasible vectors such that  $\{x_k\} \to x^*$  we have that

$$f(\boldsymbol{x}^*) - f(\boldsymbol{x}_k) \le \|\boldsymbol{\mu}^*\| \cdot \|\boldsymbol{g}^+(\boldsymbol{x}_k)\| + o(\|\boldsymbol{x}_k - \boldsymbol{x}^*\|).$$
(6.65)

Furthermore, if  $\mu^* \neq 0^m$  and  $T_X(x^*)$  is convex, the above inequality is sharp in the sense that there exists a sequence of infeasible vectors  $\{x_k\} \subset X$  such that

$$\lim_{k \to \infty} \frac{f(\boldsymbol{x}^*) - f(\boldsymbol{x}_k)}{\|\boldsymbol{g}^+(\boldsymbol{x}_k)\|} = \|\boldsymbol{\mu}^*\|,$$

and for this sequence

$$\lim_{k \to \infty} \frac{g_i^+(\boldsymbol{x}_k)}{\|\boldsymbol{g}^+(\boldsymbol{x}_k)\|} = \frac{\mu_i^*}{\|\boldsymbol{\mu}^*\|}, \qquad i = 1, \dots, m,$$

holds.

#### 

Theorem 6.34 establishes the optimal rate of cost improvement with respect to infeasible constraint perturbations (in effect, those that imply an enlargement of the feasible set).

We finally remark that under stronger conditions still, even the optimal solution  $x^*$  is differentiable. Such a result is reminiscent to the Implicit Function Theorem, which however only covers equality systems. If we are to study the sensitivity of  $x^*$  to changes in the right-hand sides of inequality constraints as well, then the analysis becomes complicated due to the fact that we must be able to predict if some active constraints may become inactive in the process. In some circumstances, different directions of change in the right-hand sides may cause different subsets of the active constraints  $\mathcal{I}(x^*)$  at  $x^*$  to become inactive, and this would most probably then be a non-differentiable point. A sufficient condition (but not necessary, at least in the case of linear constraints) for this to not happen is when  $x^*$  is strictly complementary, that is, when there exists a multiplier vector  $\mu^*$  with  $\mu_i^* > 0$  for every  $i \in \mathcal{I}(x^*)$ .

## 6.8 Applications

We provide two example applications of Lagrangian duality. The first describes the primal-dual relationship between currents and voltages in an electrical network of devices (voltage sources, diodes, and resistors); this application illustrates that Lagrange multipliers often have direct interpretations. The second application concerns a classic combinatorial optimization problem: the traveling salesman problem. We show how to approximately solve this problem through Lagrangian relaxation and subgradient optimization.

#### 6.8.1 Electrical networks

An electrical network (or, circuit) is an interconnection of analog electrical elements such as resistors, inductors, capacitors, diodes, and transistors. Its size varies from the smallest integrated circuit to an entire electricity distribution network. A circuit is a network that has at least one closed loop. A network is a connection of 2 or more simple circuit elements, and may not be a circuit. The goal when designing electrical networks for signal processing is to apply a predefined operation on *potential differences* (measured in *volts*) or *currents* (measured in *amperes*). Typical functions for these electrical networks are amplification, oscillation and analog linear algorithmic operations such as addition, subtraction, multiplication, and division. In the case of power distribution

#### Applications

networks, engineers design the circuit to transport energy as efficiently as possible while at the same time taking into account economic factors, network safety and redundancy. These networks use components such as power lines, cables, circuit breakers, switches and transformers.

To design any electrical circuits, electrical engineers need to be able to predict the voltages and currents in the circuit. Linear circuits (that is, an electrical network where all elements have a linear current–voltage relation) can be quite easily analyzed through the use of complex numbers and systems of linear equations,<sup>10</sup> while nonlinear elements require a more sophisticated analysis. The classic electrical laws describing the equilibrium state of an electrical network are due to G. Kirchhoff [Kir1847]; referred to as Kirchhoff's circuit laws they express in a mathematical form the conservation of charge and energy.<sup>11</sup>

Formally, we let an electrical circuit be described by branches (or, links) connecting nodes. We present a simple example where the only devices are voltage sources, resistors, and diodes. The resulting equilibrium conditions will be shown to be represented as the solution to a strictly convex quadratic program. In general, devices such as resistors can be non-linear, but linearity is assumed throughout this section.

- A voltage source maintains a constant branch voltage  $v_s$  irrespective of the branch current  $c_s$ . The power absorbed by the device is  $-v_s c_s$ .
- A diode permits the branch current  $c_d$  to flow in one direction only, but consumes no power regardless of the current or voltage on the branch. Denoting the branch voltage by  $v_d$ , the direction condition can be stated as a complementarity condition:

$$c_d \ge 0; \qquad v_d \ge 0; \qquad v_d c_d = 0.$$
 (6.66)

• A resistor consumes power in relation with its resistance, denoted by  $R_r$ . We recognize the following law describing the relationship between the branch current and voltage in a linear resistor:

$$v_r = -R_r c_r. \tag{6.67}$$

The power consumed is given by

$$-v_r c_r = \frac{v_r^2}{R_r} = R_r c_r^2, ag{6.68}$$

where we have utilized (6.67) to derive two alternative relations.

<sup>&</sup>lt;sup>10</sup>For such networks already Maxwell [Max1865] had stated equilibrium conditions. <sup>11</sup>These laws can be derived from Maxwell's equations, but Kirchhoff preceded Maxwell and derived his equations from work done by G. Ohm.

We must be careful about the direction of flow of currents and voltages, and thus define, for each type of device, a node–branch incidence matrix of the form

$$n_{ij} := \begin{cases} -1, & \text{if branch } j \text{ has node } i \text{ as its origin,} \\ 1, & \text{if branch } j \text{ ends in node } i, \\ 0, & \text{otherwise.} \end{cases}$$

The interpretation of a current flow variable is that the direction is from the negative to the positive terminal of the device, that is, from the origin to the ending node of the branch; a negative variable value will therefore correspond to a flow in the opposite direction. Note that for the diodes, the latter is not allowed, as seen in (6.66).

For the three types of devices we hence yield incidence matrices denoted by  $N_S$ ,  $N_R$ , and  $N_D$ , creating a partitioned matrix  $N = [N_S N_D N_R]$ . Similarly, we let  $\boldsymbol{c} = (\boldsymbol{c}_S^{\mathrm{T}}, \boldsymbol{c}_D^{\mathrm{T}}, \boldsymbol{c}_R^{\mathrm{T}})^{\mathrm{T}}$  and  $\boldsymbol{v} = (\boldsymbol{v}_S^{\mathrm{T}}, \boldsymbol{v}_D^{\mathrm{T}}, \boldsymbol{v}_R^{\mathrm{T}})^{\mathrm{T}}$ represent the vectors of branch currents and voltages. We also let  $\boldsymbol{p} = (\boldsymbol{p}_S^{\mathrm{T}}, \boldsymbol{p}_D^{\mathrm{T}}, \boldsymbol{p}_R^{\mathrm{T}})^{\mathrm{T}}$  denote the vector of node potentials. Before stating the optimization problem whose minimum describes the equilibrium of the system, we recall the two fundamental equilibrium laws:

Kirchhoff's current law: The sum of all currents entering a node is equal to the sum of all currents leaving the node. In other words, Nc = 0, or,<sup>12</sup>

$$N_S \boldsymbol{c}_S + \boldsymbol{N}_D \boldsymbol{c}_D + \boldsymbol{N}_R \boldsymbol{c}_R = \boldsymbol{0}. \tag{6.69}$$

Kirchhoff's voltage law: The difference between the node potentials at the ends of each branch is equal to the branch voltage. In other words,  $N^{T}p = v$ , or,<sup>13</sup>

$$\boldsymbol{N}_{S}^{\mathrm{T}}\boldsymbol{p} = \boldsymbol{v}_{S}, \qquad (6.70\mathrm{a})$$

$$\boldsymbol{N}_D^{\mathrm{T}} \boldsymbol{p} = \boldsymbol{v}_D, \qquad (6.70\mathrm{b})$$

$$\boldsymbol{N}_{\boldsymbol{B}}^{\mathrm{T}}\boldsymbol{p} = \boldsymbol{v}_{\boldsymbol{B}}.$$
 (6.70c)

We summarize the equations representing the characteristics of the electrical devices as follows: For the diodes, (6.66) yields

$$\boldsymbol{v}_D \ge \boldsymbol{0}; \qquad \boldsymbol{c}_D \ge \boldsymbol{0}; \qquad \boldsymbol{v}_D^{\mathrm{T}} \boldsymbol{c}_D = 0.$$
 (6.71)

For the resistors, (6.67) yields

$$\boldsymbol{v}_R = -\boldsymbol{R}\boldsymbol{c}_R,\tag{6.72}$$

 $<sup>^{12}</sup>$  This law is also referred to as the first law, the point rule, the junction rule, and the node law.

 $<sup>^{13}\,\</sup>mathrm{This}$  law is a corollary to Ohm's law, and is also referred to as the loop law.

 $\boldsymbol{R}$  being the diagonal matrix with elements equal to the values  $R_r$ .

Hence, (6.69)-(6.72) represent the equilibrium conditions of the circuit. We will now describe the optimization problem whose optimality conditions are, precisely, (6.69)-(6.72) [note that  $v_S$  is fixed]:

minimize 
$$\frac{1}{2} \boldsymbol{c}_{R}^{\mathrm{T}} \boldsymbol{R} \boldsymbol{c}_{R} - \boldsymbol{v}_{S}^{\mathrm{T}} \boldsymbol{c}_{S},$$
 (6.73)  
subject to  $\boldsymbol{N}_{S} \boldsymbol{c}_{S} + \boldsymbol{N}_{D} \boldsymbol{c}_{D} + \boldsymbol{N}_{R} \boldsymbol{c}_{R} = \boldsymbol{0},$   
 $-\boldsymbol{c}_{D} \leq \boldsymbol{0}.$ 

In the problem (6.73) we wish to determine branch currents  $c_S$ ,  $c_D$ , and  $c_R$  so as to minimize the sum of half the energy absorbed in the resistors and the energy loss of the voltage source. Note the sign condition on the diode currents.

Note that this is a convex program with linear constraints, and thus the KKT conditions are both necessary and sufficient for the global optimality of the currents. It is instrumental to check that the KKT conditions for (6.73) are given by (6.69)–(6.72), where the Lagrange multipliers are given by  $(\boldsymbol{p}^T, \boldsymbol{v}_D^T)^T$ .

In the discussion terminating in the Strong Duality Theorem 6.13, we showed that the Lagrangian dual of a strictly convex quadratic optimization problem is yet another convex quadratic optimization problem. In our case, following that development, we can derive the following dual optimization problem in terms of the node potentials p (notice, again, that  $v_S$  is fixed):

maximize 
$$-\frac{1}{2}\boldsymbol{v}_{R}^{\mathrm{T}}\boldsymbol{R}^{-1}\boldsymbol{v}_{R},$$
 (6.74)  
subject to  $\boldsymbol{N}_{D}^{\mathrm{T}}\boldsymbol{p} = \boldsymbol{v}_{S},$   
 $\boldsymbol{N}_{D}^{\mathrm{T}}\boldsymbol{p} - \boldsymbol{v}_{D} = \boldsymbol{0},$   
 $\boldsymbol{N}_{R}^{\mathrm{T}}\boldsymbol{p} - \boldsymbol{v}_{R} = \boldsymbol{0},$   
 $\boldsymbol{v}_{D} \geq \boldsymbol{0}.$ 

In the dual problem (6.74) the matrix  $\mathbf{R}^{-1}$  is the diagonal matrix of conductances. The objective function is equivalent to the minimization of the power absorbed by the resistors, and we wish to determine the branch voltages  $\mathbf{v}_D$  and  $\mathbf{v}_R$ , and the potential vector  $\mathbf{p}$ .

Verify that the KKT conditions for this problem, again, reduce to the equilibrium conditions (6.69)-(6.72). In other words, the Lagrange multipliers for the dual problem (6.74) are the (primal) branch currents.

Finally, let us note that by Theorem 6.13(a) the two problems (6.73) and (6.74) have the same objective value at optimality. That is,

$$\frac{1}{2}\boldsymbol{c}_{R}^{\mathrm{T}}\boldsymbol{R}\boldsymbol{c}_{R}+\frac{1}{2}\boldsymbol{v}_{R}^{\mathrm{T}}\boldsymbol{R}^{-1}\boldsymbol{c}_{R}-\boldsymbol{v}_{S}^{\mathrm{T}}\boldsymbol{c}_{S}=0.$$

By (6.71)–(6.72), the above equation reduces to

$$\boldsymbol{v}_{S}^{\mathrm{T}}\boldsymbol{c}_{S} + \boldsymbol{v}_{D}^{\mathrm{T}}\boldsymbol{c}_{D} + \boldsymbol{v}_{R}^{\mathrm{T}}\boldsymbol{c}_{R} = 0$$

which is precisely the principle of energy conservation.

### 6.8.2 A Lagrangian relaxation of the traveling salesman problem

Lagrangian relaxation has shown to be remarkably efficient for some combinatorial optimization problems. This is surprising when taking into account that such problems are integer or mixed-integer problems, which suffer from non-zero duality gaps in general. What then lies behind their popularity?

- One can show that Lagrangian relaxation of an integer program is always at least as good as that of a *continuous relaxation*<sup>14</sup> (in the sense that the value of  $f_R$  is higher for Lagrangian relaxation than for a continuous relaxation);
- Together with heuristics for finding primal feasible solution, good feasible solutions are often found;
- The Lagrangian relaxed problems can be made computationally much simpler than the original problem, while still keeping a lot of the structure of the original problem.

#### 6.8.2.1 The traveling salesman problem

Let the graph  $\mathcal{G} = (\mathcal{N}, \mathcal{L})$  be defined by a number of cities (or, nodes)  $i \in \mathcal{N}$  and undirected links in between subsets of pairs of them:  $(i, j) \in \mathcal{L} \subseteq \mathcal{N} \times \mathcal{N}$ . Notice that the links (i, j) and (j, i) are identical, and are in  $\mathcal{L}$  represented by one non-directed link only.

Let  $c_{ij}$  denote the distance between the cities i and j,  $\{i, j\} \subset \mathcal{N}$ . We introduce the following binary variables:

$$x_{ij} := \begin{cases} 1, & \text{if link } (i,j) \text{ is part of the TSP tour,} \\ 0, & \text{otherwise,} \end{cases} \quad (i,j) \in \mathcal{L}.$$

<sup>&</sup>lt;sup>14</sup> The continuous relaxation amounts to removing the integrality conditions, replacing, for example,  $x_j \in \{0, 1\}$  by  $x_j \in [0, 1]$ .

With these definitions, the undirected traveling salesman problem (TSP) is to

$$\min_{x} \sum_{(i,j)\in\mathcal{L}} c_{ij} x_{ij}, \qquad (6.75a)$$

subject to 
$$\sum_{(i,j)\in\mathcal{L}:\{i,j\}\subset S} x_{ij} \le |\mathcal{S}| - 1, \qquad \mathcal{S}\subset\mathcal{N}, \quad (6.75b)$$

$$\sum_{(i,j)\in\mathcal{L}}^{j,j\in\mathcal{L}} x_{ij} = n, \qquad (6.75c)$$

$$\sum_{i \in \mathcal{N}: (i,j) \in \mathcal{L}} x_{ij} = 2, \qquad j \in \mathcal{N}, \tag{6.75d}$$

$$x_{ij} \in \{0, 1\}, \quad (i, j) \in \mathcal{L}.$$
 (6.75e)

The constraints have the following interpretation: (6.75b) implies that there can be no *sub-tours*, that is, a tour where fewer than n cities are visited (if  $S \subset \mathcal{N}$  then there can be at most |S| - 1 links between nodes in the set S, where |S| is the cardinality-number of members-of the set S); (6.75c) implies that in total n cities must be visited; and (6.75d) implies that each city is connected to two others, such that we make sure to arrive from one city and leave for the next.

This problem is NP-hard, which implies that there is no known polynomial algorithm for solving it. We resort therefore to the use of relaxation techniques, in particular Lagrangian relaxation. We have more than one alternative relaxation to perform: If we Lagrangian relax the tree constraints (6.75b) and (6.75c) the remaining problem is a 2-matching problem; it can be solved in polynomial time. If we instead Lagrangian relax the degree constraints (6.75d) for every node except for one node the remaining problem is a 1-MST problem, that is, a special type of minimum spanning tree problem.

The following definition is classic: a Hamiltonian path (respectively, cycle) is a path (respectively, cycle) which passes every node in the graph exactly once. Every Hamiltonian cycle is a Hamiltonian path from a node s to another node, t, followed by a link (t, s); a subgraph which consists of a spanning tree plus an extra link such that all nodes have degree two. This is then a feasible solution to the TSP.

A 1-MST problem is the problem to find an MST in the graph that excludes node s, followed by the addition of the two least expensive links from node s to that tree. If all nodes happen to get degree two, then the 1-MST solution is a traveling salesman tour (that is, a Hamiltonian cycle). The idea behind solving the Lagrangian dual problem is then to find proper multiplier values such that the Lagrangian relaxation will produce feasible solutions.

#### 6.8.2.2 Lagrangian relaxation of the traveling salesman problem

Suppose that we Lagrangian relax the degree constraints (3), except for node 1. We assume that the starting node for the trip, node  $s \in \mathcal{N}$ , and all the links in  $\mathcal{L}$  connected to it, have been removed temporarily (in the 1-MST, this data is re-introduced later), but without changing the notation to reflect this.

The subproblem is the following: a 1-MST defined by

$$q(\boldsymbol{\lambda}) = \min_{\boldsymbol{x}} \min_{\boldsymbol{x}} \sum_{(i,j)\in\mathcal{L}} c_{ij} x_{ij} + \sum_{j\in\mathcal{N}} \lambda_j \left( 2 - \sum_{i\in\mathcal{N}:(i,j)\in\mathcal{L}} x_{ij} \right)$$
$$= 2 \sum_{j\in\mathcal{N}} \lambda_j + \min_{\boldsymbol{x}} \max_{(i,j)\in\mathcal{L}} (c_{ij} - \lambda_i - \lambda_j) x_{ij}.$$

We see immediately the role of the Lagrange multipliers: a high (low) value of the multiplier  $\lambda_j$  makes node j attractive (unattractive) in the above 1-MST problem, and will therefore lead to more (less) links being attached to it.

When solving the Lagrangian dual problem, we will use the class of subgradient optimization methods, an overview of which is found in Section 6.5.

What is the updating step in the subgradient method, and what is its interpretation? It is as usual an update in the direction of a subgradient, that is, the direction of

$$h_i(\boldsymbol{x}(\boldsymbol{\lambda})) := 2 - \sum_{i \in \mathcal{N}: (i,j) \in \mathcal{L}} x_{ij}(\boldsymbol{\lambda}), \quad i \in \mathcal{N},$$

where the value of  $x_{ij} \in \{0, 1\}$  is the solution to the 1-MST solution with link costs  $c_{ij} - \lambda_i - \lambda_j$ . We see from the direction formula that

$$\lambda_j^{\text{new}} := \lambda_j + \alpha \left( 2 - \sum_{i \in \mathcal{N}: (i,j) \in \mathcal{L}} x_{ij}(\boldsymbol{\lambda}) \right), \qquad j \in \mathcal{N},$$

where  $\alpha > 0$  is a step length. It is interesting to investigate what the update means:

current degree at node 
$$j: \begin{cases} >2 \Longrightarrow \lambda_j \downarrow (\text{link cost }\uparrow) \\ =2 \Longrightarrow \lambda_j - (\text{link cost constant}) \\ <2 \Longrightarrow \lambda_j \uparrow (\text{link cost }\downarrow) \end{cases}$$

In other words, the updating formula in a subgradient method is such that the link cost in the 1-MST subproblem is shifted upwards

(downwards) if there are too many (too few) links connected to node j in the 1-MST. We are hence adjusting the *node prices* of the nodes in such a way as to try to influence the 1-MST problem to always choose 2 links per node to connect to.

#### 6.8.2.3 A feasibility heuristic

A feasibility heuristic takes the optimal solution from the Lagrangian minimization problem over x and adjusts it such that a feasible solution to the original problem is constructed. As one cannot predict if, or when, a primal feasible solution will be found directly from the subproblem, the heuristic will provide a solution that can be used in place of an optimal one, should one not be found. Moreover, as we know from Lagrangian duality theory, we then have access to both lower and upper bounds on the optimal value  $f^*$  of the original problem, and so we have a quality measure of the feasible solutions found.

A feasibility heuristic which can be used together with our Lagrangian heuristic is as follows.

Identify a path in the 1-MST with many links. Then form a subgraph with the remaining nodes and find a path that passes all of them. Put the two paths together in the best way. The resulting path is a Hamiltonian cycle, that is, a feasible solution.

#### 6.8.2.4 The Philips example

In 1987–1988 an M.Sc. project was performed at the department of mathematics at Linköping University, in cooperation with the company Philips, Norrköping. The project was initiated with the goal to improve the current practice of solving a production planning problem.

The problem was as follows: Philips produce circuit boards, perhaps several hundreds or thousands of the same type. There is a new batch of patterns (holes) to be drilled every day, and perhaps even several such batches per day.

In order to speed up the production process the drilling machine is connected to a microcomputer that selects the ordering of the holes to be drilled automatically, given their coordinates. The algorithm for performing the sorting used to be a simple sorting operation that found, for every fixed x-coordinate, the corresponding y-coordinates and sorted them in increasing order. The movement of the drill was therefore from left to right, and for each fixed x-coordinate the movement was vertical. The time it took to drill the holes on one circuit board was, however, far too long, simply because the drill traveled around a lot without performing any tasks, following a path that was too long. (On the other

hand, the actual ordering was very fast to produce!) All in all, the complete batch production took too long because of the poorly planned drill movement.

It was observed that the production planning problem is a traveling salesman problem, where the cities are the holes to be drilled, and the distances between them correspond to the Euclidean distances between them. Therefore, an efficient TSP heuristic was devised and implemented, for use in conjunction with the microcomputer. In fact, it was based on precisely the above Lagrangian relaxation, a subgradient optimization method, and a graph-search type heuristic of the form discussed above.

A typical run with the algorithm took a few minutes, and was always stopped after a fixed number of subgradient iterations; the generation of feasible solutions with the above-mentioned graph search technique was performed at every Kth iteration, where K was set to a value strictly larger than one. (Moreover, feasible solutions were not generated during the first iterations of the dual procedure, because of the poor quality of  $\lambda_k$  for low values of k; it is often the case that the traveling salesman tour resulting from the heuristic is better when the multipliers are nearoptimal in the Lagrangian dual problem.)

In one of the examples implemented it was found that the optimal path length was in the order to 2 meters, and that the upper and lower bounds on  $f^*$  produced lead to the conclusion that the relative error of the path length of the best feasible solution found was less than 7 %, a quite good result, also showing that the duality gap for the problem at hand (together with the Lagrangian relaxation chosen) is quite small.

After implementing the new procedure, Philips could report an increase in production by some 70 %. Hence, the slightly longer time it took to provide a better production plan, that is, the traveling salesman tour for the drill to follow, was more than well compensated by the fact that the drilling could be done much faster.

Here is hence an interesting case where Lagrangian relaxation helped to solve a large-scale, complex and difficult problem by utilizing problem structure.

## 6.9 Notes and further reading

Lagrangian duality has been developed in many sources, including early developments by Arrow, Hurwicz, and Uzawa [AHU58], Everett [Eve63], and Falk [Fal67], and later on by Rockafellar [Roc70]. Our development follows to a large extent that of portions of the text books by Bertsekas [Ber99], Bazaraa et al. [BSS93], and Rockafellar [Roc70].

The Relaxation Theorem 6.1 can almost be considered to be folklore, and can be found in a slightly different form in [Wol98, Proposition 2.3].

The differentiability properties of convex functions were developed largely by Rockafellar [Roc70], whose text we mostly follow.

Subgradient methods were developed in the Soviet Union in the 1960s, predominantly by Ermol'ev, Polyak, and Shor. Text book treatments of subgradient methods are found, for example, in [Sho85, HiL93, Ber99]. Theorem 6.23 is essentially due to Ermol'ev [Erm66]; the proof stems from [LPS96]. Theorem 6.24 is due to Shepilov [She76]; finally, Theorem 6.25 is due to Polyak [Pol69].

Everett's Theorem is due to Everett [Eve63].

Theorem 6.34 stems from [Ber04, Proposition 1.1].

That the equilibrium conditions of an electrical or hydraulic network are attained as the minimum of the total energy loss were known more than a century ago. Mathematical programming models for the electrical network equilibrium problems described in Section 6.8.1 date at least as far back as to Duffin [Duf46, Duf47] and d'Auriac [dAu47]. Duffin constructs his objective function as a sum of integrals of resistance functions. The possibility of viewing the equilibrium problem in at least two related, dual, ways as that of either finding the optimal flows of currents or the optimal potentials was also known early in the analysis of electrical networks; these two principles are written out in [Cro36] in work on pipe networks, and explicitly stated as a pair of primal-dual quadratic programming problems in [Den59]; we followed his development, as represented in [BSS93, Section 1.2.D].

The traveling salesman problem is an essential model problem in combinatorial optimization. Excellent introductions to the field can be found in [Law76, PaS82, NeW88, Wol98, Sch03]. It was the work in [HWC74, Geo74, Fis81, Fis85], among others, in the 1970s and 1980s on the traveling salesman problem and its relatives that made Lagrangian relaxation and subgradient optimization popular, and it remains most popular within the combinatorial optimization field.

#### 6.10 Exercises

**Exercise 6.1** (numerical example of Lagrangian relaxation) Consider the convex problem to

minimize 
$$\frac{1}{x_1} + \frac{4}{x_2}$$
,  
subject to  $x_1 + x_2 \le 4$ ,  
 $x_1, x_2 \ge 0$ .

(a) Lagrangian relax the first constraint, and write down the resulting implicit dual objective function and the dual problem. Motivate why the relaxed problem always has a unique optimum, whence the dual objective function is everywhere differentiable.

(b) Solve the implicit Lagrangian dual problem by utilizing that the gradient to a differentiable dual objective function can be expressed by using the functions that are involved in the relaxed constraints and the unique solution to the relaxed problem.

(c) Write down an explicit Lagrangian dual problem, that is, a dual problem only in terms of the Lagrange multipliers. Solve it, and confirm the results in (b).

(d) Find the original problem's optimal solution.

(e) Show that strong duality holds.

Exercise 6.2 (global optimality conditions) Consider the problem to

minimize 
$$f(x) := x_1 + 2x_2^2 + 3x_3^3$$
,  
subject to  $x_1 + 2x_2 + x_3 \le 3$ ,  
 $2x_1^2 + x_2 \ge 2$ ,  
 $2x_1 + x_3 = 2$ ,  
 $x_j \ge 0$ ,  $j = 1, 2, 3$ .

(a) Formulate the Lagrangian dual problem that results from Lagrangian relaxing all but the sign constraints.

(b) State the global primal-dual optimality conditions.

Exercise 6.3 (Lagrangian relaxation) Consider the problem to

minimize 
$$f(x) := x_1^2 + 2x_2^2$$
,  
subject to  $x_1 + x_2 \ge 2$ ,  
 $x_1^2 + x_2^2 \le 5$ .

Find an optimal solution through Lagrangian duality.

**Exercise 6.4 (Lagrangian relaxation)** In many circumstances it is of interest to calculate the Euclidean projection of a vector onto a subspace. Especially, consider the problem to find the Euclidean projection of the vector  $\boldsymbol{y} \in \mathbb{R}^n$  onto the null space of the matrix  $\boldsymbol{A} \in \mathbb{R}^{m \times n}$ , that is, to find an  $\boldsymbol{x} \in \mathbb{R}^n$  that solves the problem to

minimize 
$$f(x) := \frac{1}{2} ||y - x||^2$$
,  
subject to  $Ax = 0^m$ ,

where  $\boldsymbol{A}$  is such that rank  $\boldsymbol{A} = m$ .

The solution to this problem is classic: the projection is given by

$$\boldsymbol{x}^* = \boldsymbol{y} - \boldsymbol{A}^{\mathrm{T}} (\boldsymbol{A} \boldsymbol{A}^{\mathrm{T}})^{-1} \boldsymbol{A} \boldsymbol{y}.$$

If we let  $\boldsymbol{P} := \boldsymbol{I}^n - \boldsymbol{A}^{\mathrm{T}} (\boldsymbol{A} \boldsymbol{A}^{\mathrm{T}})^{-1} \boldsymbol{A}$ , where  $\boldsymbol{I}^n \in \mathbb{R}^{n \times n}$  is the unit matrix, be the projection matrix, the formula is simply  $\boldsymbol{x}^* = \boldsymbol{P} \boldsymbol{y}$ .

Your task is to derive this formula by utilizing Lagrangian duality. Motivate every step made by showing that the necessary properties are fulfilled.

[Note: This exercise is similar to that in Example 5.51, but utilizes Lagrangian duality rather than the KKT conditions to derive the projection formula.]

**Exercise 6.5** (Lagrangian relaxation, exam 040823) Consider the following linear optimization problem:

minimize 
$$f(x, y) := x - 0.5y$$
,  
subject to  $-x + y \le -1$ ,  
 $-2x + y \le -2$ ,  
 $(x, y) \in \mathbb{R}^2_+$ .

(a) Show that the problem satisfies Slater's constraint qualification. Derive the Lagrangian dual problem corresponding to the Lagrangian relaxation of the two linear inequality constraints, and show that its set of optimal solutions is convex and bounded.

(b) Calculate the set of subgradients of the Lagrangian dual function at the dual points  $(1/4, 1/3)^{T}$  and  $(1, 0)^{T}$ .

**Exercise 6.6 (Lagrangian relaxation)** Provide an explicit form of the Lagrangian dual problem for the problem to

$$\begin{array}{ll} \text{minimize} & \sum_{i=1}^{m} \sum_{j=1}^{n} x_{ij} \ln x_{ij} \\ \text{subject to} & \sum_{i=1}^{m} x_{ij} = b_j, \qquad j = 1, \dots, n, \\ & \sum_{j=1}^{n} x_{ij} = a_i, \qquad i = 1, \dots, m, \\ & x_{ij} \ge 0, \qquad i = 1, \dots, m, \quad j = 1, \dots, n, \end{array}$$

where  $a_i > 0$ ,  $b_j > 0$  for all i, j, and where the linear equalities are Lagrangian relaxed.

Exercise 6.7 (Lagrangian relaxation) Given is the problem to

minimize 
$$f(x) = 2x_1^2 + x_2^2 + x_1 - 3x_2,$$
 (6.76a)  
subject to  $x_1^2 + x_2 \ge 8,$  (6.76b)

$$x_1 \in [1, 3],$$
 (6.76c)

 $x_2 \in [2,5].$  (6.76d)

Lagrangian relax the constraint (6.76b) with a multiplier  $\mu$ . Formulate the Lagrangian dual problem and calculate the dual function's value at  $\mu = 1$ ,  $\mu = 2$ , and  $\mu = 3$ . Within which interval lies the optimal value  $f^*$ ? Also, draw the dual function.

**Exercise 6.8 (Lagrangian duality for integer problems)** Consider the primal problem to

minimize 
$$f(x)$$
,  
subject to  $g(x) \le 0^m$ ,  
 $x \in X$ ,

where  $X \subseteq \mathbb{R}^n$ ,  $f : \mathbb{R}^n \to \mathbb{R}$ , and  $g : \mathbb{R}^n \to \mathbb{R}^m$ . If the restrictions  $g(x) \leq \mathbf{0}^m$  are complicating side constraints which are Lagrangian relaxed, we obtain the Lagrangian dual problem to

(a) Suppose that the set X is finite (for example, consisting of a finite number of integer vectors). Denote the elements of X by  $x^p$ ,  $p = 1, \ldots, P$ . Show that the dual objective function is piece-wise linear. How many linear segments can it have, at most? Why is it not always built up by that many segments?

[Note: This property holds regardless of any properties of f and g.] (b) Illustrate the result in (a) on the linear 0/1 problem to find

where the first constraint is considered complicating.

(c) Suppose that the function f and all components of g are linear, and that the set X is a polytope (that is, a bounded polyhedron). Show that the dual objective function is also in this case piece-wise linear. How many linear pieces can it be built from, at most?

192

where

**Exercise 6.9** (Lagrangian relaxation) Consider the problem to

Lagrangian relax the first constraint. Describe the Lagrangian function and the dual problem. Calculate the Lagrangian dual function at these four points:  $\mu = 0, 1, 2, 3$ . Give the best lower and upper bounds on the optimal value of the original problem that you have found.

**Exercise 6.10 (surrogate relaxation)** Consider an optimization problem of the form

minimize 
$$f(\boldsymbol{x})$$
,  
subject to  $g_i(\boldsymbol{x}) \leq 0$ ,  $i = 1, \dots, m$ ,  $(P)$   
 $\boldsymbol{x} \in X$ ,

where the functions  $f, g_i : \mathbb{R}^n \to \mathbb{R}$  are continuous and the set  $X \subset \mathbb{R}^n$  is closed and bounded. The problem is assumed to have an optimal solution,  $x^*$ . Introduce parameters  $\mu_i \ge 0, i = 1, \ldots, m$ , and define

$$s(\boldsymbol{\mu}) := \min f(\boldsymbol{x}),$$
  
subject to  $\boldsymbol{\mu}^{\mathrm{T}} \boldsymbol{g}(\boldsymbol{x}) \leq 0,$  (S)  
 $\boldsymbol{x} \in X.$ 

This problem therefore has exactly one explicit constraint.

(a) [weak duality] Show that  $x^*$  is a feasible solution to the problem (S) and that  $s(\mu) \leq f^*$  therefore always holds, that is, the problem (S) is a relaxation of the original one. Motivate also why maximum<sub> $\mu \geq 0^m$ </sub>  $s(\mu) \leq f^*$  must hold. Explain the potential usefulness of this result!

(b) [example] Consider the linear 0/1 problem

$$z^* = \underset{\text{subject to}}{\text{maximum }} z = 5x_1 + 8x_2 + 7x_3 + 9x_4,$$
  

$$subject \text{ to} \qquad 3x_1 + 2x_2 + 3x_3 + 3x_4 \leq 6, \qquad (1)$$
  

$$2x_1 + 3x_2 + 3x_3 + 4x_4 \leq 5, \qquad (2)$$
  

$$2x_1 + x_2 + 2x_3 + x_4 = 3,$$
  

$$x_1 , x_2 , x_3 , x_4 \in 0/1.$$

Surrogate relax the constraints (1) and (2) with multipliers  $\mu_1, \mu_2 \ge 0$ and formulate the problem (S). Let  $\bar{\mu} = (1,2)^{\mathrm{T}}$ . Calculate  $s(\bar{\mu})$ .

Consider again the original problem and Lagrangian relax the constraints (1) and (2) with multipliers  $\mu_1, \mu_2 \ge 0$ . Calculate the Lagrangian dual objective value at  $\mu = \bar{\mu}$ .

Compare the two results! (c) [comparison with Lagrangian duality] Let  $\mu \ge 0^m$  and

$$q(\boldsymbol{\mu}) := \min_{\boldsymbol{x} \in X} \{f(\boldsymbol{x}) + \boldsymbol{\mu}^{\mathrm{T}} \boldsymbol{g}(\boldsymbol{x})\}.$$

Show that  $q(\boldsymbol{\mu}) \leq s(\boldsymbol{\mu})$ , and that

$$\underset{\boldsymbol{\mu} \ge \mathbf{0}^m}{\operatorname{maximum}} q(\boldsymbol{\mu}) \le \underset{\boldsymbol{\mu} \ge \mathbf{0}^m}{\operatorname{maximum}} s(\boldsymbol{\mu}) \le f^*$$

holds.