

# Elements of Convex Analysis and Optimization

Master II Datascience

Pascal Bianchi, Olivier Fercoq, Walid Hachem

September 18, 2019

# Contents

<b>1</b>	<b>Convex analysis</b>	<b>3</b>
1.1	Convex sets, separation and relative interior . . . . .	3
1.1.1	Convex sets . . . . .	3
1.1.2	Separation results . . . . .	4
1.1.3	Relative interior . . . . .	5
1.2	Convex functions . . . . .	6
1.2.1	Definition and properties . . . . .	6
1.2.2	Operations preserving convexity . . . . .	7
1.3	Subdifferential . . . . .	8
1.4	Lower semi-continuity . . . . .	10
1.5	Minimizers, coercivity, strict and strong convexity . . . . .	10
1.6	Exercises . . . . .	11
<b>2</b>	<b>Fenchel-Legendre transform</b>	<b>15</b>
2.1	Definitions and properties . . . . .	15
2.2	The Fenchel-Moreau theorem . . . . .	16
2.3	The Fenchel-Young inequality and its consequences . . . . .	17
2.4	Exercises . . . . .	18
<b>3</b>	<b>Duality</b>	<b>20</b>
3.1	Parametric Duality . . . . .	20
3.1.1	Primal and dual problems . . . . .	20
3.1.2	Lagrangian . . . . .	21
3.2	Fenchel-Rockafellar duality . . . . .	22
3.2.1	Main results . . . . .	23
3.2.2	Affine equality constraints . . . . .	24
3.2.3	Operations on subdifferentials . . . . .	25
3.2.4	The Attouch-Brezis theorem* . . . . .	26
3.3	Lagrangian duality . . . . .	27
3.3.1	Inequality and affine equality constraints . . . . .	27
3.3.2	Inequality constraints only . . . . .	30
3.4	Examples, Exercises and Problems . . . . .	31
<b>4</b>	<b>Fixed Points Algorithms</b>	<b>35</b>
4.1	$\alpha$ -averaged operators . . . . .	35
4.2	The gradient algorithm . . . . .	37
4.3	The proximal point and the proximal gradient algorithms . . . . .	38
4.3.1	The proximity operator . . . . .	38
4.3.2	The proximal point algorithm . . . . .	38

4.3.3	The proximal gradient algorithm . . . . .	39
4.4	Applications . . . . .	39
4.4.1	Projected gradient algorithm . . . . .	39
4.4.2	Iterative soft-thresholding . . . . .	40
<b>5</b>	<b>Monotone operators</b> . . . . .	<b>41</b>
5.1	Basic definitions and facts . . . . .	41
5.1.1	Monotone and maximal monotone operators . . . . .	41
5.1.2	The proximal point algorithm . . . . .	43
5.2	Splitting algorithms . . . . .	46
5.2.1	The Douglas-Rachford splitting algorithm . . . . .	46
5.2.2	The Forward-Backward algorithm . . . . .	49
5.2.3	Discussion . . . . .	54
5.3	Exercises . . . . .	55
<b>6</b>	<b>Dual methods</b> . . . . .	<b>57</b>
6.1	Method of multipliers . . . . .	57
6.1.1	Problem setting . . . . .	57
6.1.2	Algorithm . . . . .	57
6.1.3	Application: a splitting method . . . . .	59
6.2	Augmented Method of Multipliers . . . . .	60
6.3	Alternating Direction Method of Multipliers (ADMM) . . . . .	60
6.4	The Vũ-Condat method . . . . .	61

# Chapter 1

## Elements of convex analysis

Let  $\mathcal{X}$  be a Euclidian space endowed with a scalar product denoted by  $\langle \cdot, \cdot \rangle$  and an associated norm  $\| \cdot \|$ .

### Convex sets, separation and relative interior

#### Convex sets

In this section,  $C$  denotes a subset of  $\mathbf{E}$ .

**Definition 1.1** (Convex set). The set  $C$  is **convex** if

$$\forall (x, y) \in C^2, \forall t \in [0, 1], \quad tx + (1 - t)y \in C.$$

**Proposition 1.1.** 1. Any arbitrary intersection of convex sets is convex.

2. The closure and the interior of a convex set are convex.

*Proof.* The proof is left as an exercise. □

**Lemma 1.2.** Let  $C$  be a convex set. For every  $k \geq 1$ , every  $(a_1, \dots, a_k) \in C^k$  and every  $(\alpha_1, \dots, \alpha_k) \in \mathbb{R}_+^k$  such that  $\sum_{i=1}^k \alpha_i = 1$ ,

$$\sum_{i=1}^k \alpha_i a_i \in C.$$

Such a weighted sum, with non negative weights  $\alpha_1, \dots, \alpha_k$  summing to one, is called a **convex combination** of the points  $a_1, \dots, a_k$ .

*Proof.* By induction. The statement holds for  $k = 1$ . Now consider  $k \geq 2$  and assume that  $a_k > 0$ .

$$\sum_{i=1}^k \alpha_i a_i = \alpha_k a_k + (1 - \alpha_k) \sum_{i=1}^{k-1} \frac{\alpha_i}{1 - \alpha_k} a_i.$$

The sum in the righthand side lies in  $C$  by induction. Thus, the righthand side lies in  $C$  by definition of a convex set. □

## Separation results

**Proposition 1.3** (Projection). *Let  $C \subset \mathcal{X}$  be a non-empty, closed, and convex set. For every  $x \in \mathcal{X}$ , there is a unique point in  $C$ , denoted by  $P_C(x)$ , such that*

$$\text{for all } y \in C, \|y - x\| \geq \|P_C(x) - x\|.$$

*Moreover, the mapping  $P_C : \mathcal{X} \rightarrow \mathcal{X}$  satisfies the following.*

1.  $\forall y \in C, \langle y - P_C(x), x - P_C(x) \rangle \leq 0$ .
2.  $\forall (x, y) \in \mathcal{X}^2, \|P_C(y) - P_C(x)\| \leq \|y - x\|$ .

The point  $P_C(x)$  is called the projection of  $x$  on  $C$ .

*Proof.*

1. Let  $d_C(x) = \inf_{y \in C} \|y - x\|$ . There exists a sequence  $(y_n)_n$  in  $C$  such that  $\|y_n - x\| \rightarrow d_C(x)$ . The sequence is bounded, so extract a subsequence which converges to  $y_0$ . By continuity of  $y \mapsto \|y - x\|$ , we have  $\|y_0 - x\| = d_C(x)$ , as required.

To prove uniqueness, consider a point  $z \in C$  such that  $\|z - x\| = d_C(x)$ . By convexity of  $C$ ,  $w = (y_0 + z)/2 \in C$ , so  $\|w - x\| \geq d_C(x)$ . According to the parallelogram identity<sup>1</sup>,

$$\begin{aligned} 4d_C(x)^2 &= 2\|y_0 - x\|^2 + 2\|z - x\|^2 \\ &= \|y_0 + z - 2x\|^2 + \|y_0 - z\|^2 \\ &= 4\|w - x\|^2 + \|y_0 - z\|^2 \\ &\geq 4d_C(x)^2 + \|y_0 - z\|^2. \end{aligned}$$

Thus,  $\|y_0 - z\| = 0$  and  $y_0 = z$ .

2. Let  $p = P_C(x)$  and let  $y \in C$ . For  $\epsilon \in [0, 1]$ , let  $z_\epsilon = p + \epsilon(y - p)$ . By convexity,  $z_\epsilon \in C$ . Consider the function 'squared distance from  $x$ ':

$$\varphi(\epsilon) = \|z_\epsilon - x\|^2 = \|\epsilon(y - p) + p - x\|^2.$$

For  $0 < \epsilon \leq 1$ ,  $\varphi(\epsilon) \geq d_C(x)^2 = \varphi(0)$ . Furthermore, for  $\epsilon$  sufficiently close to zero,

$$\varphi(\epsilon) = d_C(x)^2 - 2\epsilon \langle y - p, x - p \rangle + o(\epsilon),$$

whence  $\varphi'(0) = -2 \langle y - p, x - p \rangle$ . In the case  $\varphi'(0) < 0$ , we would have, for  $\epsilon$  close to 0,  $\varphi(\epsilon) < \varphi(0) = d_C(x)^2$ , which is impossible. So  $\varphi'(0) \geq 0$  and the result follows.

3. Adding the inequalities

$$\begin{aligned} \langle P_C(y) - P_C(x), x - P_C(x) \rangle &\leq 0, \text{ and} \\ \langle P_C(x) - P_C(y), y - P_C(y) \rangle &\leq 0, \end{aligned}$$

yields  $\langle P_C(y) - P_C(x), y - x \rangle \geq \|P_C(x) - P_C(y)\|^2$ . The conclusion follows using Cauchy-Schwarz inequality.  $\square$

**Proposition 1.4** (Separation of a point and a convex set by an hyperplane). *Let  $C$  be a non empty, closed and convex set. Consider  $x_0 \in \mathcal{X} \setminus C$ . Then, there exists  $a \in \mathbb{R}$  and  $w \in \mathcal{X} \setminus \{0\}$  such that*

$$\begin{aligned} \forall x \in C, \langle w, x \rangle + a &\leq 0, \\ \langle w, x_0 \rangle + a &> 0. \end{aligned}$$

---

<sup>1</sup>  $2\|a\|^2 + 2\|b\|^2 = \|a + b\|^2 + \|a - b\|^2$ .

*Proof.* Set  $w = x_0 - P_C(x_0)$  and  $a = \langle w, P_C(x_0) \rangle$ . Note that  $w \neq 0$  and  $\langle w, x_0 \rangle + a = \|w\|^2 > 0$ . For every  $x \in C$ ,  $\langle w, x - x_0 \rangle \leq 0$  by Proposition 1.3.  $\square$

The set  $H := \{x \in \mathcal{X} : \langle w, x \rangle + a = 0\}$  defines an hyperplane. This hyperplane splits the space  $\mathcal{X}$  into two half spaces. Proposition 1.4 states that  $C$  is included in one of these half spaces, while the point  $x_0$  lies in the (interior of) the other. Otherwise stated,  $x_0$  is *separated* from  $C$  by the hyperplane  $H$ .

We denote by  $\text{cl}(C)$  the closure of a set  $C$  and by  $\text{int}(C)$  its interior. We denote by  $\text{bdry}(C)$  the boundary of a set  $C$ , defined by  $\text{bdry}(C) = \text{cl}(C) \setminus \text{int}(C)$ .

**Theorem 1.5** (Supporting hyperplane). *Let  $C$  be a non-empty convex set and let  $x_0 \in \text{bdry}(C)$ . There exists  $w \in \mathcal{X} \setminus \{0\}$  such that  $\forall x \in C$ ,  $\langle w, x - x_0 \rangle \leq 0$ .*

*Proof.* Let  $C$  and  $x_0$  as in the statement. There is a sequence  $(x_n)$  with  $x_n \in \mathcal{X} \setminus \text{cl}(C)$  and  $x_n \rightarrow x_0$ , otherwise there would be a ball included in  $C$  that would contain  $x_0$ , and  $x_0$  would be in  $\text{int}(C)$ . By Prop. 1.4,

$$\forall n, \exists w_n \in \mathcal{X} \setminus \{0\}, \forall x \in C, \langle w_n, x \rangle < \langle w_n, x_n \rangle. \quad (1.1.1)$$

It can be assumed without restriction that  $\|w_n\| = 1$  (otherwise, just replace  $w_n$  by  $w_n/\|w_n\|$ ) in the above statement. Since the sequence  $(w_n)$  is bounded, we can extract a subsequence  $(w_{k_n})_n$  that converges to some  $w \in \mathcal{X}$ . By continuity of the norm,  $\|w\| = 1$ . Letting  $n \rightarrow \infty$  in (1.1.1), we obtain the result.  $\square$

## Relative interior

**Definition 1.2.** A set  $E \subset \mathcal{X}$  is called an **affine space** if, for all  $(x, y) \in E^2$  and for all  $t \in \mathbb{R}$ ,  $x + t(y - x) \in E$ .

If  $E$  is a set and  $x \in \mathcal{X}$  is a point, then we define the sum  $a + E$  as the set of points of the form  $a + x$  for  $x \in E$ . It is easy to check that if  $A$  is an affine space and if  $x_0 \in E$ , then  $E - x_0$  is a vector space, which does not depend on the choice of  $x_0$  in  $E$ . The **dimension** of the affine space is the dimension of the corresponding vector space.

**Definition 1.3.** The **affine hull**  $\text{aff}(C)$  of a set  $C \subset \mathcal{X}$  is the smallest affine space that contains  $C$ .

**Definition 1.4.** Let  $C \subset \mathcal{X}$ . The **relative interior** of  $C$ , denoted by  $\text{ri}(C)$  is the set of points  $x \in C$  which admit a neighborhood  $V$  such that  $V \cap \text{aff}(C) \subset C$ .

*n.b.:* For the students aware of topological notions,  $\text{ri}(C)$  is the interior of  $C$  in the topology induced by  $\text{aff}(C)$ . Obviously,  $\text{int}(C) \subset \text{ri}(C)$ .

**Theorem 1.6.** *Let  $C$  be a non empty and convex set of  $\mathcal{X}$ , then  $\text{ri}(C) \neq \emptyset$ .*

*Proof.* It is not difficult to prove that for every  $x_0 \in \mathcal{X}$ ,  $\text{ri}(x_0 + C) = x_0 + \text{ri}(C)$ . Hence, it is sufficient to make the proof in the case where  $0 \in C$ . In this case,  $\text{aff}(C)$  is a vector space. Denote by  $d$  its dimension. One can construct an independent family  $(z_i)_{1 \leq i \leq d}$  of elements of  $C$  (construct it by yourself, as an exercise). Let  $\bar{x} = \frac{1}{d+1} \sum_{i=1}^d z_i$  be the barycenter of the points,  $z_1, \dots, z_d$  and  $0$ . As  $0 \in C$ , the point  $\bar{x}$  is a convex combination of elements of  $C$ , it thus lies in  $C$  by Lemma 1.2. Let  $\varepsilon > 0$  and consider the neighborhood  $V$  of  $\bar{x}$  defined as the open ball of radius  $\varepsilon$  and centered at  $\bar{x}$ . Choose  $x \in V \cap \text{aff}(C)$ . As  $x \in \text{aff}(C)$ , the point  $x$  writes as a

linear combination of the vectors  $(z_i)_i$ , say  $x = \sum_{i=1}^d w_i z_i$  for some  $(w_1, \dots, w_d) \in \mathbb{R}^d$ . For every  $i = 1, \dots, d$ ,

$$|w_i - (d+1)^{-1}| \|z_i\| \leq \|x - \bar{x}\| < \varepsilon.$$

Choose  $\varepsilon < \min_i \|z_i\| (d(d+1))^{-1}$ . Then, for every  $i$ ,

$$\frac{1 - 1/d}{d+1} < w_i < \frac{1 + 1/d}{d+1}.$$

Thus,  $w_i > 0$  and  $\sum_i w_i < 1$ . By Lemma 1.2 again,  $x \in C$ . This shows that  $V \cap \text{aff}(C) \subset C$ . Hence,  $\bar{x} \in \text{ri}(C)$ .  $\square$

**Theorem 1.7.** *Let  $M : \mathcal{X} \rightarrow \mathcal{Y}$  be a linear operator and  $C$  be a convex subset of  $\mathcal{X}$ . Then,  $\text{ri}(MC) = M \text{ri}(C)$ .*

*Proof.* To be completed. See (Rockafellar, 2015, Th. 6.6).  $\square$

## Convex functions

### Definition and properties

For all  $f : \mathcal{X} \rightarrow [-\infty, +\infty]$ , the **domain** of  $f$ , denoted by  $\text{dom}(f)$ , is the set of points  $x$  such that  $f(x) < +\infty$ .

A function  $f$  is called **proper** if  $\text{dom}(f) \neq \emptyset$  (i.e.  $f \not\equiv +\infty$ ) and if  $f$  *never* takes the value  $-\infty$ .

**Definition 1.5.** Let  $f : \mathcal{X} \rightarrow [-\infty, +\infty]$ . The **epigraph** of  $f$ , denoted by  $\text{epi } f$ , is the subset of  $\mathcal{X} \times \mathbb{R}$  defined by:

$$\text{epi } f = \{(x, t) \in \mathcal{X} \times \mathbb{R} : t \geq f(x)\}.$$

**Definition 1.6** (Convex function).  $f : \mathcal{X} \rightarrow [-\infty, +\infty]$  is **convex** if its epigraph is convex.

**Proposition 1.8.** *A function  $f : \mathcal{X} \rightarrow [-\infty, +\infty]$  is convex if and only if*

$$\forall (x, y) \in \text{dom}(f)^2, \forall t \in (0, 1), \quad f(tx + (1-t)y) \leq tf(x) + (1-t)f(y).$$

*Proof.* Assume that  $f$  satisfies the inequality. Let  $(x, u)$  and  $(y, v)$  be two points of the epigraph :  $u \geq f(x)$  and  $v \geq f(y)$ . In particular,  $(x, y) \in \text{dom}(f)^2$ . Let  $t \in ]0, 1[$ . The inequality implies that  $f(tx + (1-t)y) \leq tu + (1-t)v$ . Thus,  $t(x, u) + (1-t)(y, v) \in \text{epi}(f)$ , which proves that  $\text{epi}(f)$  is convex.

Conversely, assume that  $\text{epi}(f)$  is convex. Let  $(x, y) \in \text{dom}(f)^2$ . For  $(x, u)$  and  $(y, v)$  two points in  $\text{epi}(f)$ , and  $t \in [0, 1]$ , the point  $t(x, u) + (1-t)(y, v)$  belongs to  $\text{epi}(f)$ . So,  $f(tx + (1-t)y) \leq tu + (1-t)v$ . If  $f(x)$  et  $f(y)$  are  $> -\infty$ , we can choose  $u = f(x)$  and  $v = f(y)$ , which demonstrates the inequality. If  $f(x) = -\infty$ , we can choose  $u$  arbitrary close to  $-\infty$ . Letting  $u$  go to  $-\infty$ , we obtain  $f(tx + (1-t)y) = -\infty$ , which proves the statement.  $\square$

**Lemma 1.9.** *If  $f : \mathcal{X} \rightarrow [-\infty, +\infty]$  is convex, then  $\text{dom}(f)$  is convex.*

*Proof.* To do as an exercise.  $\square$

**Proposition 1.10.** *Let  $f : \mathcal{X} \rightarrow (-\infty, +\infty]$  be a convex function. Then,  $f$  is continuous at every point in  $\text{int}(\text{dom } f)$ .*

*Proof.* To do as an exercise.  $\square$

## Operations preserving convexity

In the sequel, we use the convention that the supremum of the empty set in  $\mathbb{R}$  is equal to  $-\infty$ .

**Definition 1.7.** The **upper hull** of a collection of functions  $(f_\alpha : \alpha \in I)$  on  $\mathcal{X} \rightarrow [-\infty, +\infty]$ , where  $I$  is an arbitrary set, is the function  $x \mapsto \sup_{\alpha \in I} f_\alpha(x)$ .

**Proposition 1.11.** *The upper hull of of familly of l.s.c. functions is l.s.c. The upper hull of of familly of convex functions is convex.*

*Proof.* Let  $f$  denote the upper hull of a collection  $(f_\alpha : \alpha \in I)$ . We remark that  $\text{epi}(f) = \bigcap_{\alpha} \text{epi}(f_\alpha)$ . If every function  $f_\alpha$  is l.s.c (resp. convex), then  $\text{epi}(f_\alpha)$  is closed (resp. convex). Now  $\text{epi}(f)$  is closed (resp. convex) as an arbitrary intersection of closed (resp. convex) sets. Thus,  $f$  is l.s.c. (resp. convex).  $\square$

**Proposition 1.12.** *Let  $F : \mathcal{X} \times \mathcal{Y} \rightarrow [-\infty, \infty]$  be a convex function. Then, the function defined on  $\mathcal{X} \rightarrow [-\infty, +\infty]$  by  $y \mapsto \inf_{x \in \mathcal{X}} F(x, y)$  is convex.*

*Proof.* Consider  $u, v \in \mathcal{X}$  and  $t \in (0, 1)$ . Denote  $f : y \mapsto \inf_{x \in \mathcal{X}} F(x, y)$ . Then,

$$\begin{aligned} tf(u) + (1-t)f(v) &= t \inf_{x \in \mathcal{X}} F(x, u) + (1-t) \inf_{x' \in \mathcal{X}} F(x', v) \\ &= \inf_{x \in \mathcal{X}, x' \in \mathcal{X}} tF(x, u) + (1-t)F(x', v) \\ &\geq \inf_{x \in \mathcal{X}, x' \in \mathcal{X}} F(tx + (1-t)x', tu + (1-t)v) \\ &= \inf_{x \in \mathcal{X}} F(x, tu + (1-t)v) \\ &= f(tu + (1-t)v). \end{aligned}$$

$\square$

**Definition 1.8.** A map  $A : \mathcal{X} \rightarrow \mathcal{Y}$  is said **affine** if there exists a linear operator  $M : \mathcal{X} \rightarrow \mathcal{Y}$  and a vector  $b \in \mathcal{Y}$  such that  $A : x \mapsto Mx + b$ .

**Proposition 1.13.** *Let  $f : \mathcal{Y} \rightarrow [-\infty, +\infty]$  be a convex function, and let  $A : \mathcal{X} \rightarrow \mathcal{Y}$  be an affine map. Then,  $f \circ A$  is convex.*

*Proof.* Let  $(x, y) \in \mathcal{X}^2$  and  $t \in (0, 1)$ . By Prop. 1.8,  $f \circ A(tx + (1-t)y) = f(tA(x) + (1-t)A(y)) \leq tf(A(x) + (1-t)f(A(y))$ , which proves that  $f \circ A$  is convex.  $\square$

**Proposition 1.14.** *Let  $m \geq 1$  be an integer and let  $f_1, \dots, f_m$  be convex functions on  $\mathcal{X} \rightarrow (-\infty, +\infty]$ . Then,  $\sum_i f_i$  is convex.*

*Proof.* The point is easily shown from Prop. 1.8.  $\square$

**Definition 1.9.** The infimal convolution of two functions  $f, g : \mathcal{X} \rightarrow [-\infty, +\infty]$  is the function  $f \square g : \mathcal{X} \rightarrow [-\infty, +\infty]$  defined by

$$f \square g : y \mapsto \inf_{x \in \mathcal{X}} f(x) + g(y - x). \quad (1.2.1)$$

We say that  $f \square g$  is *exact in a point*  $y$  if the infimum in (1.2.1) is reached. We say that  $f \square g$  is *exact* if it is exact in every point. It is clear that  $f \square g = g \square f$ .

**Proposition 1.15.** *Consider two convex functions  $f, g : \mathcal{X} \rightarrow (-\infty, +\infty]$ . Then,  $f \square g$  is convex.*



*Proof.* By Prop. 1.13 and Prop. 1.14, the map  $F$  defined on  $\mathcal{X} \times \mathcal{X}$  by  $F : (x, y) \mapsto f(x) + g(y - x)$  is convex. Hence,  $f \square g$  is convex by Prop. 1.12.  $\square$

**Definition 1.10.** The **infimal postcomposition** of a linear operator  $M : \mathcal{X} \rightarrow \mathcal{Y}$  and a function  $f : \mathcal{X} \rightarrow [-\infty, +\infty]$  is the function  $M \triangleright f : \mathcal{Y} \rightarrow [-\infty, +\infty]$  defined by

$$M \triangleright f : y \mapsto \inf\{f(x) : x \in \mathcal{X}, Mx = y\}.$$

**Proposition 1.16.** Consider  $f : \mathcal{X} \rightarrow (-\infty, +\infty]$  and a linear operator  $M : \mathcal{X} \rightarrow \mathcal{Y}$ . Then,  $\text{dom}(M \triangleright f) = M \text{dom } f$ . Moreover, if  $f$  is convex, then  $M \triangleright f$  is convex.

*Proof.* The first point follows directly from the definition. Assume that  $f$  is convex. Then, the mapping  $F$  defined on  $\mathcal{X} \times \mathcal{Y}$  by  $F : (x, y) \mapsto f(x) + \iota_{\text{gra}(M)}(x, y)$  is convex, where  $\text{gra}(M) = \{(x, Mx) : x \in \mathcal{X}\}$ . The function  $M \triangleright f$  coincides with  $y \mapsto \inf_{x \in \mathcal{X}} F(x, y)$ . It is thus convex by Prop. 1.12.  $\square$

## Subdifferential

**Definition 1.11** (Subdifferential). Let  $f : \mathcal{X} \rightarrow [-\infty, +\infty]$  and  $x \in \text{dom}(f)$ . A vector  $\phi \in \mathcal{X}$  is called a **subgradient** of  $f$  at  $x$  if:

$$\forall y \in \mathcal{X}, \quad f(y) \geq f(x) + \langle \phi, y - x \rangle.$$

The **subdifferential** of  $f$  in  $x$ , denoted by  $\partial f(x)$ , is the set of all the subgradients of  $f$  at  $x$ . By convention,  $\partial f(x) = \emptyset$  if  $x \notin \text{dom}(f)$ . Formally,  $\partial f : \mathcal{X} \rightarrow 2^{\mathcal{X}}$  is a set-valued mapping.

**Theorem 1.17.** Let  $f : \mathcal{X} \rightarrow [-\infty, +\infty]$  be a convex function and  $x \in \text{ri}(\text{dom } f)$ . Then  $\partial f(x)$  is non-empty.

*Proof.* Let  $x_0 \in \text{ri}(\text{dom } f)$ . We assume that  $f(x_0) > -\infty$  (otherwise the proof is trivial). We may restrict ourselves to the case  $x_0 = 0$  and  $f(x_0) = 0$  (up to replacing  $f$  by the function  $x \mapsto f(x + x_0) - f(x_0)$ ).

In this case, for all vector  $\phi \in \mathcal{X}$ ,

$$\phi \in \partial f(0) \quad \Leftrightarrow \quad \forall x \in \text{dom } f, \quad \langle \phi, x \rangle \leq f(x).$$

Let  $\mathcal{A} = \text{aff}(\text{dom } f)$ . As  $\mathcal{A}$  contains point zero, it is an Euclidean vector space.

Let  $C$  be the closure of  $\text{epi } f \cap (\mathcal{A} \times \mathbb{R})$ . The set  $C$  is a convex closed set in  $\mathcal{A} \times \mathbb{R}$ , which is endowed with the scalar product  $\langle (x, u), (x', u') \rangle = \langle x, x' \rangle + uu'$ .

The point  $(0, 0) = (x_0, f(x_0))$  belongs to the boundary of  $C$ . Therefore, Th. 1.5 applies in  $\mathcal{A} \times \mathbb{R}$ . There is a vector  $w \in \mathcal{A} \times \mathbb{R}$ ,  $w \neq 0$ , such that

$$\forall z \in C, \quad \langle w, z \rangle \leq 0$$

Write  $w = (\phi, u) \in \mathcal{A} \times \mathbb{R}$ . For  $z = (x, t) \in C$ , we have

$$\langle \phi, x \rangle + ut \leq 0.$$

Let  $x \in \text{dom}(f)$ . In particular  $f(x) < \infty$  and for all  $t \geq f(x)$ ,  $(x, t) \in C$ . Thus,

$$\forall x \in \text{dom}(f), \quad \forall t \geq f(x), \quad \langle \phi, x \rangle + ut \leq 0. \quad (1.3.1)$$

Letting  $t$  tend to  $+\infty$ , we obtain  $u \leq 0$ .

Let us prove by contradiction that  $u < 0$ . Suppose not (*i.e.*  $u = 0$ ). Then  $\langle \phi, x \rangle \leq 0$  for all  $x \in \text{dom}(f)$ . As  $0 \in \text{ri}(\text{dom}(f))$ , there is a set  $\tilde{V}$ , open in  $\mathcal{A}$ , such that  $0 \in \tilde{V} \subset \text{dom } f$ . Thus for  $x \in \mathcal{A}$ , there is an  $\epsilon > 0$  such that  $\epsilon x \in \tilde{V} \subset \text{dom}(f)$ . According to (1.3.1),  $\langle \phi, \epsilon x \rangle \leq 0$ , so  $\langle \phi, x \rangle \leq 0$ . Similarly,  $\langle \phi, -x \rangle \leq 0$ . Therefore,  $\langle \phi, x \rangle \equiv 0$  on  $\mathcal{A}$ . Since  $\phi \in \mathcal{A}$ ,  $\phi = 0$  as well. Finally  $w = 0$ , which is a contradiction.

As a result,  $u < 0$ . Dividing inequality (1.3.1) by  $-u$ , and taking  $t = f(x)$ , we get

$$\forall x \in \text{dom}(f), \forall t \geq f(x), \quad \left\langle \frac{-1}{u} \phi, x \right\rangle \leq f(x).$$

So  $\frac{-1}{u} \phi \in \partial f(0)$ . □

**Definition 1.12.** We say that a function  $f$  is a **minorant** of a function  $g$  if  $f \leq g$ .

A function  $f : \mathcal{X} \rightarrow \mathbb{R}$  is affine if there exists  $a \in \mathcal{X}$  and  $b \in \mathbb{R}$  such that  $f(x) = \langle a, x \rangle + b$  for every  $x \in \mathcal{X}$ .

**Proposition 1.18.** Let  $f : \mathcal{X} \rightarrow (-\infty, +\infty]$  be a convex function. Then,  $f$  admits an affine minorant.

*Proof.* The result is trivial if  $f$  is identically equal to  $+\infty$ . In the other case ( $f$  is proper),  $\text{dom}(f) \neq \emptyset$ , it holds that  $\text{ri}(\text{dom } f) \neq \emptyset$  by Th. 1.6. Consider  $x_0 \in \text{ri}(\text{dom } f)$ . By Th. 1.17, there exists  $\varphi \in \mathcal{X}$  such that for every  $x \in \mathcal{X}$ ,  $f(x) \geq f(x_0) + \langle \varphi, x - x_0 \rangle$ . The mapping  $x \mapsto f(x_0) + \langle \varphi, x - x_0 \rangle$  is an affine minorant of  $f$ . □

When  $f$  is differentiable at  $x \in \text{dom } f$ , we denote by  $\nabla f(x)$  its gradient at  $x$ . The link between differentiation and subdifferential is given by the following proposition :

**Proposition 1.19.** Let  $f : \mathcal{X} \rightarrow (-\infty, \infty]$  be a convex function, differentiable in  $x$ . Then  $\partial f(x) = \{\nabla f(x)\}$ .

*Proof.* If  $f$  is differentiable at  $x$ , the point  $x$  necessarily belongs to  $\text{int}(\text{dom}(f))$ . Let  $\phi \in \partial f(x)$  and  $t \neq 0$ . Then for all  $y \in \text{dom}(f)$ ,  $f(y) - f(x) \geq \langle \phi, y - x \rangle$ . Applying this inequality to  $y = x + t(\phi - \nabla f(x))$  (which belongs to  $\text{dom}(f)$  for  $t$  small enough) leads to :

$$\frac{f(x + t(\phi - \nabla f(x))) - f(x)}{t} \geq \langle \phi, \phi - \nabla f(x) \rangle.$$

The left term converges to  $\langle \nabla f(x), \phi - \nabla f(x) \rangle$ . Finally,

$$\langle \nabla f(x) - \phi, \phi - \nabla f(x) \rangle \geq 0,$$

*i.e.*  $\phi = \nabla f(x)$ . □

**Example 1.1.** The absolute-value function  $x \mapsto |x|$  defined on  $\mathbb{R} \rightarrow \mathbb{R}$  admits as a subdifferential the sign application, defined by :

$$\text{sign}(x) = \begin{cases} \{1\} & \text{si } x > 0 \\ [-1, 1] & \text{si } x = 0 \\ \{-1\} & \text{si } x < 0. \end{cases}$$

## Lower semi-continuity

**Definition 1.13** (Reminder : **lim inf** : **limit inferior**).

The **limit inferior** of a sequence  $(u_n)_{n \in \mathbb{N}}$ , where  $u_n \in [-\infty, \infty]$ , is

$$\liminf(u_n) = \sup_{n \geq 0} \left( \inf_{k \geq n} u_k \right).$$

Since the sequence  $V_n = \inf_{k \geq n} u_k$  is non decreasing, an equivalent definition is

$$\liminf(u_n) = \lim_{n \rightarrow \infty} \left( \inf_{k \geq n} u_k \right).$$

**Definition 1.14** (Lower semicontinuous function). A function  $f : \mathcal{X} \rightarrow [-\infty, \infty]$  is called **lower semicontinuous (l.s.c.)** at  $x \in \mathcal{X}$  if for all sequence  $(x_n)$  which converges to  $x$ ,

$$\liminf f(x_n) \geq f(x).$$

The function  $f$  is said to be **lower semicontinuous**, if it is l.s.c. at  $x$ , for all  $x \in \mathcal{X}$ . The function  $f$  is said to be **closed**, if  $\text{epi}(f)$  is closed.

**Proposition 1.20** (epigraphical characterization). *Let  $f : \mathcal{X} \rightarrow [-\infty, +\infty]$ . Then  $f$  is l.s.c. if and only if it is closed.*

*Proof.* If  $f$  is l.s.c., and if  $(x_n, t_n) \in \text{epi } f \rightarrow (\bar{x}, \bar{t})$ , then,  $\forall n, t_n \geq f(x_n)$ . Consequently,

$$\bar{t} = \liminf t_n \geq \liminf f(x_n) \geq f(\bar{x}).$$

Thus,  $(\bar{x}, \bar{t}) \in \text{epi } f$ , and  $\text{epi } f$  is closed.

Conversely, if  $f$  is *not* l.s.c., there exists an  $x \in \mathcal{X}$ , and a sequence  $(x_n) \rightarrow x$ , such that  $f(x) > \liminf f(x_n)$ , i.e., there is an  $\epsilon > 0$  such that  $\forall n \geq 0, \inf_{k \geq n} f(x_k) \leq f(x) - \epsilon$ . Thus, for all  $n, \exists k_n \geq k_{n-1}, f(x_{k_n}) \leq f(x) - \epsilon$ . We have built a sequence  $(w_n) = (x_{k_n}, f(x) - \epsilon)$ , each term of which belongs to  $\text{epi } f$ , and which converges to a limit  $\bar{w} = (x, f(x) - \epsilon)$  which is outside the epigraph. Consequently,  $\text{epi } f$  is not closed.  $\square$

**Proposition 1.21.** *The sum of two l.s.c. functions is l.s.c.*

*Proof.* Let  $f, g : \mathcal{X} \rightarrow [-\infty, +\infty]$  be two l.s.c. functions, and let  $(x_k)$  be a sequence converging to  $x \in \mathcal{X}$ . Then,  $\liminf f(x_k) + g(x_k) \geq \liminf f(x_k) + \liminf g(x_k) \geq f(x) + g(x)$ .  $\square$

**Notation.** We denote by  $\Gamma_0(\mathcal{X})$  the set of all closed proper convex functions on  $\mathcal{X} \rightarrow (-\infty, +\infty]$ .

## Minimizers, coercivity, strict and strong convexity

A point  $x$  is called a **minimizer** of  $f$  if  $f(x) \leq f(y)$  for all  $y \in \mathcal{X}$ . The set of minimizers of  $f$  is denoted  $\arg \min(f)$ .

**Proposition 1.22** (Fermat's rule).  $x \in \arg \min f \Leftrightarrow 0 \in \partial f(x)$ .

*Proof.* This follows from the definition:  $0 \in \partial f(x) \Leftrightarrow \forall y, f(y) \geq f(x) + \langle 0, y - x \rangle$ . This means that  $x \in \arg \min f$ .  $\square$

**Definition 1.15.** A function  $f : \mathcal{X} \rightarrow (-\infty, +\infty]$  is said **coercive** if  $\lim_{\|x\| \rightarrow +\infty} f(x) = +\infty$ .

Otherwise stated,  $f$  is coercive if, for every sequence  $(x_n)$  such that  $\|x_n\| \rightarrow +\infty$ , it holds that  $f(x_n) \rightarrow +\infty$ .

**Proposition 1.23.** *Let  $f : \mathcal{X} \rightarrow (-\infty, +\infty]$  be coercive and l.s.c. Then,  $f$  admits a minimizer.*

*Proof.* We consider the case where  $f$  is not identically equal to  $+\infty$  (otherwise the proof is trivial). Consider a sequence  $(x_n)$  such that  $f(x_n) \rightarrow \inf f(\mathcal{X})$ . Such a sequence exists by definition of the infimum. The sequence is bounded. Indeed, if it were unbounded, one would be able to extract a subsequence  $(x_{\varphi_n})$  converging to  $+\infty$  in norm. The coercivity assumption would imply that  $f(x_{\varphi_n}) \rightarrow +\infty$  which would contradict the fact that  $f(x_n) \rightarrow \inf f(\mathcal{X})$ . As  $(x_n)$  is bounded, one can extract a converging subsequence, say  $x_{\psi_n} \rightarrow x^*$  for some  $x^* \in \mathcal{X}$ . As  $f$  is l.s.c.,  $\inf f(\mathcal{X}) = \lim f(x_{\psi_n}) = \liminf f(x_{\psi_n}) \geq f(x^*)$ . Thus,  $f(x^*) = \inf f(\mathcal{X})$  and  $x^*$  is a minimizer.  $\square$

**Proposition 1.24.** *A convex function  $f : \mathcal{X} \rightarrow (-\infty, \infty]$  is coercive iff there exists  $\gamma > 0$  and  $\alpha \in \mathbb{R}$  such that  $f(x) \geq \gamma\|x\| + \alpha$ .*

*Proof.* To do.  $\square$

**Definition 1.16.** A function  $f : \mathcal{X} \rightarrow (-\infty, +\infty]$  is said **strictly convex** if for every  $x, y \in \text{dom } f$  s.t.  $x \neq y$ , and every  $t \in (0, 1)$ ,  $f(tx + (1-t)y) < tf(x) + (1-t)f(y)$ .

**Proposition 1.25.** *A strictly convex function  $f : \mathcal{X} \rightarrow (-\infty, +\infty]$  with non empty domain admits at most one minimizer.*

*Proof.* By contradiction, consider two distinct minimizers  $x, y$ . Then,  $f((x+y)/2) < (f(x) + f(y))/2 = \inf f(\mathcal{X})$ , which is impossible.  $\square$

*Example 1.2.* The squared Euclidean norm  $x \mapsto \|x\|^2$  is strictly convex.

**Definition 1.17.** Let  $\mu > 0$ . A function  $f : \mathcal{X} \rightarrow (-\infty, +\infty]$  is said  **$\mu$ -strongly convex** if  $f - \frac{\mu}{2}\|\cdot\|^2$  is convex. It is said **strongly convex** if it is  $\mu$ -strongly convex for some  $\mu > 0$ .

**Proposition 1.26.** *Any strongly convex function  $f : \mathcal{X} \rightarrow (-\infty, +\infty]$  is strictly convex and coercive. As a consequence, a strongly convex function in  $\Gamma_0(\mathcal{X})$  admits a unique minimizer.*

*Proof.* We only consider the case where  $f \not\equiv +\infty$ , otherwise the first result is trivial. Set  $g = f - \frac{\mu}{2}\|\cdot\|^2$ . By Prop. 1.18, the convex function  $g$  admits an affine minorant, say  $h$ . Thus,  $f \geq h + \frac{\mu}{2}\|\cdot\|^2$ . As  $h$  is affine, it holds that  $h(x) + \frac{\mu}{2}\|x\|^2$  tends to  $+\infty$  as  $\|x\| \rightarrow +\infty$ . Hence,  $f$  is coercive. Consider  $x \neq y$  in  $\text{dom } f$ , and  $t \in (0, 1)$ . By convexity of  $g$ ,  $g(tx + (1-t)y) \leq tg(x) + (1-t)g(y)$ . Therefore,

$$f(tx + (1-t)y) \leq tf(x) + (1-t)f(y) + \frac{\mu}{1} (\|tx + (1-t)y\|^2 - t\|x\|^2 - (1-t)\|y\|^2) .$$

The term enclosed in the parenthesis is strictly negative by the strict convexity of the squared norm (see Example 1.2). Thus,  $f$  is strictly convex. The second result follows from Propositions 1.23 and 1.25.  $\square$

## Exercises

**Exercise 1.1.** Let  $f : \mathcal{X} \rightarrow [-\infty, \infty]$ , and assume that  $\partial f(x)$  and  $\partial f(y)$  are non empty for some  $x, y \in \mathcal{X}$ . Show that

$$\forall u \in \partial f(x), \forall v \in \partial f(y), \langle x - y, u - v \rangle \geq 0 .$$

**Exercise 1.2.** Consider  $m$  Euclidean spaces  $\mathcal{X}_1, \dots, \mathcal{X}_m$ , where  $m \geq 1$  is an integer. For every  $i \in \{1, \dots, m\}$ , let  $f_i : \mathcal{X}_i \rightarrow (-\infty, +\infty]$  be a function and consider the function  $f : \prod_{i=1}^m \mathcal{X}_i \rightarrow (-\infty, +\infty]$  such that

$$f : (x_1, \dots, x_m) \mapsto \sum_{i=1}^m f_i(x_i).$$

1. Prove that

$$\inf f = \sum_{i=1}^m \inf f_i.$$

2. Prove that

$$\arg \min f = \prod_{i=1}^m \arg \min f_i.$$

3. Prove that

$$\partial f = \prod_{i=1}^m \partial f_i.$$

**Exercise 1.3.** Prove that the set  $\{x \in \mathcal{X} : Ax = b\}$  is affine, where  $A$  is a linear operator on  $\mathcal{X} \rightarrow \mathcal{Y}$  and  $b \in \mathcal{Y}$ . Conversely, prove that any affine space can be written under this form.

**Exercise 1.4.** Consider a function  $f : \mathcal{X} \rightarrow (-\infty, +\infty]$  and let  $\phi \in \mathcal{X}$ . Prove that the subdifferential of the function  $x \mapsto f(x) + \langle \phi, x \rangle$  coincides, at a given point  $x \in \mathcal{X}$ , with  $\partial f(x) + \phi$ .

**Exercise 1.5** (normal cone). For every set  $C \subset \mathcal{X}$  and every  $x \in C$ , we define  $N_C(x) = \{\varphi \in \mathcal{X} : \forall y \in C, \langle \varphi, y - x \rangle \leq 0\}$ . We set  $N_C(x) = \emptyset$  for every  $x \notin C$ .

1. Prove the identity  $N_C = \partial \iota_C$ .
2. Show that for all  $x \in C$ ,  $N_C(x)$  is a cone, i.e.  $\forall \varphi \in N_C(x), \forall \lambda \geq 0, \lambda \varphi \in N_C(x)$ . It is called the **normal cone** of  $C$  at  $x$ .
3. Prove that  $N_C(x) = \{0\}$  whenever  $x \in \text{int}(C)$ .
4. Derive  $N_C(x)$  in the following cases:  $C = (-\infty, 0]$ ,  $C = (-\infty, 1]$ ,  $C = \mathbb{R}_- \times \mathbb{R}_-$ ,  $C = \{(u, v) \in \mathbb{R}^2 : u^2 + v^2 \leq 1\}$ .

**Exercise 1.6.** Let  $C$  and  $D$  be two convex subsets of  $\mathcal{X}$ . Prove that  $\text{ri}(C - D) = \text{ri}(C) - \text{ri}(D)$ .  
*Hint: apply Th. 1.7 to the convex set  $C \times D$ .*

**Exercise 1.7.** Show that if  $f$  is convex, proper, with  $\text{dom } f = \mathcal{X}$ , and if  $f$  is bounded, then  $f$  is constant.

**Exercise 1.8.** \* Let  $f$  be a convex function and  $x, y$  in  $\text{dom } f$ ,  $t \in (0, 1)$  and  $z = tx + (1 - t)y$ . Assume that the three points  $(x, f(x))$ ,  $(z, f(z))$  and  $(y, f(y))$  are aligned. Show that for all  $u \in (0, 1)$ ,  $f(ux + (1 - u)y) = uf(x) + (1 - u)f(y)$ .

**Exercise 1.9** (the question of  $-\infty$  values).

Let  $f : \mathcal{X} \rightarrow [-\infty, +\infty]$  be convex and assume that  $\text{ri dom } f$  contains a point  $x$  such that  $f(x) > -\infty$ . Prove that  $f$  *never* takes the value  $-\infty$ . Thus,  $f$  is proper.  
(Hint: use Prop. 1.17).

**Exercise 1.10.** Determine the subdifferentials of the following functions, at the considered points:

1. In  $\mathcal{X} = \mathbb{R}$ ,  $f(x) = \iota_{[0,1]}$ , at  $x = 0, x = 1$  and  $0 < x < 1$ .
2. In  $\mathcal{X} = \mathbb{R}^2$ ,  $f(x_1, x_2) = \iota_{x_1 < 0}$ , at  $x$  such that  $x_1 = 0, x_1 < 0$ .
3.  $\mathcal{X} = \mathbb{R}$ ,

$$f(x) = \begin{cases} +\infty & \text{si } x < 0 \\ -\sqrt{x} & \text{si } x \geq 0 \end{cases}$$

at  $x = 0$ , and  $x > 0$ .

4.  $\mathcal{X} = \mathbb{R}^n$ ,  $f(x) = \|x\|$ , determine  $\partial f(x)$ , for any  $x \in \mathbb{R}^n$ .
5.  $\mathcal{X} = \mathbb{R}$ ,  $f(x) = x^3$ . Show that  $\partial f(x) = \emptyset, \forall x \in \mathbb{R}$ . Explain this result.
6.  $\mathcal{X} = \mathbb{R}^n$ ,  $C = \{y : \|y\| \leq 1\}$ ,  $f(x) = \iota_C(x)$ . Give the subdifferential of  $f$  at  $x$  such that  $\|x\| < 1$  and at  $x$  such that  $\|x\| = 1$ .

*Hint:* For  $\|x\| = 1$ :

- Show that  $\partial f(x) = \{\phi : \forall y \in C, \langle \phi, y - x \rangle \leq 0\}$ .
- Show that  $x \in \partial f(x)$  using Cauchy-Schwarz inequality. Deduce that the cone  $\mathbb{R}^+x = \{tx : t \geq 0\} \subset \partial f(x)$ .
- To show the converse inclusion : Fix  $\phi \in \partial f$  and pick  $u \in \{x\}^\perp$  (i.e.,  $u$  s.t.  $\langle u, x \rangle = 0$ ). Consider the sequence  $y_n = \|x + t_n u\|^{-1}(x + t_n u)$ , for some sequence  $(t_n)_n, t_n > 0, t_n \rightarrow 0$ . What is the limit of  $y_n$  ?

Consider now  $u_n = t_n^{-1}(y_n - x)$ . What is the limit of  $u_n$  ? Conclude about the sign of  $\langle \phi, u \rangle$ .

Do the same with  $-u$ , conclude about  $\langle \phi, u \rangle$ . Conclude.

7. Let  $f : \mathbb{R} \rightarrow \mathbb{R}$ , differentiable. Show that:  $f$  is convex, if and only if

$$\forall (x, y) \in \mathbb{R}^2, \langle \nabla f(y) - \nabla f(x), y - x \rangle \geq 0.$$

**Exercise 1.11.** Show that a function  $f$  is l.s.c. if and only if its level sets :

$$L_{\leq \alpha} = \{x \in \mathcal{X} : f(x) \leq \alpha\}$$

are closed.

(see, e.g., [Rockafellar et al. \(1998\)](#), theorem 1.6.)

**Exercise 1.12.** Let  $f : \mathcal{X} \rightarrow [-\infty, +\infty]$  be l.s.c. and convex. Assume that  $-\infty \in f(\mathcal{X})$ . Then  $f(\mathcal{X}) \subset \{-\infty, +\infty\}$ .

**Exercise 1.13** (proximity operator). We define the proximity operator of a function  $f \in \Gamma_0(\mathcal{X})$  as the mapping defined for every  $x \in \mathcal{X}$  by

$$\text{prox}_f(x) = \arg \min_{y \in \mathcal{X}} f(y) + \frac{\|y - x\|^2}{2}. \quad (1.6.1)$$

1. Show that the mapping  $\text{prox}_f : \mathcal{X} \rightarrow \mathcal{X}$  is well defined i.e., that the arg min is always a singleton.

*Hint:* Use Prop. 1.26

2. Let  $\gamma > 0$ . Show that  $p = \text{prox}_{\gamma f}(x)$  iff  $p \in x + \gamma \partial f(x)$ .
3. Deduce the identity  $\text{prox}_{\gamma f} = (I + \gamma \partial f)^{-1}$ .
4. Prove that the fixed points of  $\text{prox}_{\gamma f}$  (*i.e.*, the points  $x \in \mathcal{X}$  s.t.  $x = \text{prox}_{\gamma f}(x)$ ) coincide with the minimizers of  $f$ .
5. Evaluate  $\text{prox}_{\gamma f}$  when  $f(x)$  coincides with  $\iota_C(x)$  ( $C$  closed convex and non empty),  $|x|$ ,  $\|x\|_1$ ,  $\|x\|^2$ ,  $\|x\|$ .

## Chapter 2

# Fenchel-Legendre transform

### Definitions and properties

**Definition 2.1.** Let  $f : \mathcal{X} \rightarrow [-\infty, +\infty]$ . The **Fenchel-Legendre transform**, or Fenchel-Legendre conjugate, of  $f$  is the function  $f^* : \mathcal{X} \rightarrow [-\infty, \infty]$ , defined by

$$f^*(\phi) = \sup_{x \in \mathcal{X}} \langle \phi, x \rangle - f(x), \quad \phi \in \mathcal{X}.$$

The following lemma is straightforward.

**Lemma 2.1.** *Let  $f, g : \mathcal{X} \rightarrow [-\infty, +\infty]$ . If  $f \leq g$  then  $f^* \geq g^*$ .*

**Proposition 2.2.** *Let  $f : \mathcal{X} \rightarrow [-\infty, +\infty]$ . Then  $f^*$  is convex and l.s.c.*

*Proof.* Remark that  $f^*$  is the upper hull of the collection of functions  $(\phi \mapsto \langle \phi, x \rangle - f(x) : x \in \mathcal{X})$ . Each of these functions is affine, it is thus l.s.c. and convex. The result follows from Prop.1.11.  $\square$

**Remark 2.1.** By definition of  $f^*$ , it is clear that  $-\infty \in f^*(\phi)$  iff  $f \equiv +\infty$ . In other words,  $f^*$  never takes the value  $-\infty$  unless  $f$  is identically  $+\infty$ .

For every  $\phi \in \mathcal{X}$ , we denote by  $\mathcal{A}_\phi(f)$  the set of affine minorants of gradient  $\phi$ , that is:

$$\mathcal{A}_\phi(f) := \{\alpha : \alpha : \mathcal{X} \rightarrow \mathbb{R} \text{ is an affine minorant of } f \text{ and } \nabla \alpha \equiv \phi\}.$$

Note that every  $\alpha \in \mathcal{A}_\phi(f)$  has the form  $\alpha(x) = \langle \phi, x \rangle + \alpha(0)$  for every  $x \in \mathcal{X}$ . If there exists  $\alpha \in \mathcal{A}_\phi(f)$  such that for every  $\beta \in \mathcal{A}_\phi(f)$ ,  $\alpha \geq \beta$ , we say that  $\alpha$  is the **maximal element** of  $\mathcal{A}_\phi(f)$ . If such a maximal element exists, it is unique.

**Proposition 2.3.** *Let  $f : \mathcal{X} \rightarrow (-\infty, +\infty]$  be proper and let  $\phi \in \mathcal{X}$ . Then,*

1. *If  $f^*(\phi) = +\infty$ , then  $f$  admits no affine minorant of gradient  $\phi$  (i.e.,  $\mathcal{A}_\phi(f) = \emptyset$ ).*
2. *If  $f^*(\phi) \in \mathbb{R}$ , the function  $x \mapsto \langle \phi, x \rangle - f^*(\phi)$  is the maximal element in  $\mathcal{A}_\phi(f)$ .*

*Proof.* For every  $\alpha \in \mathcal{A}_\phi(f)$ ,  $\alpha(x) = \langle \phi, x \rangle + \alpha(0)$  for every  $x \in \mathcal{X}$ . As  $\alpha \leq f$ , it holds that  $-\alpha(0) \geq \langle \phi, x \rangle - f(x)$  for every  $x \in \mathcal{X}$ . Hence,  $-\alpha(0) \geq f^*(\phi)$ . This cannot happen if  $f^*(\phi) = +\infty$ , and the first point is proved. For every  $x \in \mathcal{X}$ , the former inequality implies that  $\langle \phi, x \rangle - f^*(\phi) \geq \langle \phi, x \rangle + \alpha(0)$ . This proves that  $\langle \phi, \cdot \rangle - f^*(\phi) \geq \alpha$  for every  $\alpha \in \mathcal{A}_\phi(f)$ , and the second point is proved.  $\square$



The set

$$\mathcal{A}(f) = \bigcup_{\phi \in \mathcal{X}} \mathcal{A}_\phi(f)$$

is the set of affine minorants of  $f$ . Recall that this set is empty for every function  $f : \mathcal{X} \rightarrow (-\infty, +\infty]$  that is convex. The upper hull of  $\mathcal{A}(f)$  is called the **affine envelope** of  $f$  *i.e.*, it is the function

$$x \mapsto \sup\{\alpha(x) : \alpha \in \mathcal{A}\},$$

where we recall the convention  $\sup \emptyset = -\infty$ . We use the notation  $f^{**}$  to designate  $(f^*)^*$ . The function  $f^{**}$  is called the **biconjugate**, and satisfies for every  $x \in \mathcal{X}$ ,

$$f^{**}(x) = \sup_{\phi \in \mathcal{X}} \langle \phi, x \rangle - f^*(\phi).$$

**Proposition 2.4.** *For every  $f : \mathcal{X} \rightarrow [-\infty, +\infty]$ ,  $f^{**}$  is the affine envelope of  $f$ . In particular,  $f \geq f^{**}$ .*

*Proof.* Denote by  $h$  the affine envelope of  $f$ . Assume that  $f$  is proper. Consider  $\alpha \in \mathcal{A}(f)$ . There exists  $\phi$  such that  $\alpha \in \mathcal{A}_\phi(f)$ . By Prop. 2.3,  $\alpha(x) \leq \langle \phi, x \rangle - f^*(\phi)$  for all  $x \in \mathcal{X}$ . Hence,  $\alpha(x) \leq f^{**}(x)$ . Taking the supremum w.r.t.  $\alpha \in \mathcal{A}$ ,  $h(x) \leq f^{**}(x)$ . To show the other inequality, note that for every  $\phi \in \text{dom}(f^*)$ ,  $\langle \phi, \cdot \rangle - f^*(\phi) \in \mathcal{A}(f)$  by Prop. 2.3. Thus, for every  $x \in \mathcal{X}$ ,

$$f^{**}(x) = \sup\{\langle \phi, x \rangle - f^*(\phi) : \phi \in \text{dom}(f^*)\} \leq h(x).$$

This shows that  $f^{**} = h$ .

When  $f \equiv +\infty$ , it is clear that  $h \equiv +\infty$ . We have  $f^* \equiv -\infty$  which in turn yields  $f^{**} \equiv +\infty$ . Thus,  $f^{**} = h$ .

When  $-\infty \in f(\mathcal{X})$ ,  $f$  has no affine minorant and  $h \equiv -\infty$ . We have  $f^* \equiv +\infty$  which implies  $f^{**} \equiv -\infty$ . Thus,  $f^{**} = h$  as well.

The fact that  $f \geq f^{**}$  is straightforward from the definition of the affine envelope.  $\square$

## The Fenchel-Moreau theorem

**Theorem 2.5** (Fenchel-Moreau). *Consider  $f : \mathcal{X} \rightarrow (-\infty, +\infty]$ . Then  $f = f^{**}$  if and only if  $f$  is convex and l.s.c.*

*Proof.* By Prop. 2.2,  $f^{**}$  is always convex and l.s.c. Therefore,  $f = f^{**}$  implies that  $f$  is convex and l.s.c. We prove the converse. By definition, the affine envelope minorates  $f$ . Thus, by Prop. 2.3,  $f^{**} \leq f$ . We now show that  $f^{**} \geq f$ . The result is trivial when  $f \equiv +\infty$ , we thus consider that  $f$  is proper. By contradiction, assume that there exists  $x_0$  such that  $f^{**}(x_0) < f(x_0)$ . Clearly, the point  $z_0 := (x_0, f^{**}(x_0))$  is such that  $z_0 \notin \text{epi}(f)$ . As  $f$  is proper, l.s.c. and convex,  $\text{epi}(f)$  is a non-empty, closed and convex set. By Prop. 1.4, there exists a non-zero vector  $(\phi, u) \in \mathcal{X} \times \mathbb{R}$  and  $a \in \mathbb{R}$  such that  $\langle w, z \rangle + a \leq 0$  for every  $z \in \text{epi}(f)$  and  $\langle w, z_0 \rangle + a > 0$ . This reads:

$$\forall x \in \text{dom}(f), \forall t \geq f(x), \langle \phi, x \rangle + ut + a \leq 0, \tag{2.2.1}$$

$$\text{and } \langle \phi, x_0 \rangle + uf^{**}(x_0) + a > 0, \tag{2.2.2}$$

Letting  $t \rightarrow +\infty$  in (2.2.1), it follows that  $u \leq 0$ .

First consider the case where  $u < 0$ . Putting  $t = f(x)$  in (2.2.1) and dividing both sides of the inequality by  $-u$ , it follows that for every  $x \in \text{dom}(f)$ ,

$$\left\langle -\frac{\phi}{u}, x \right\rangle + a \leq f(x).$$

Therefore, the affine map  $\langle -\phi/u, x \rangle + a$  lies in  $\mathcal{A}(f)$ . By Prop. 2.4,  $f^{**}(x_0) \geq \langle -\phi/u, x_0 \rangle + a$  for every  $x$ . This contradicts (2.2.2).

Now consider the case where  $u = 0$ . Assume that  $f \geq 0$ . Eq. (2.2.1) reads  $\langle \phi, x \rangle + a \leq 0$  for every  $x \in \text{dom}(f)$ . Choose  $\varepsilon > 0$ . In particular,  $\langle \phi, x \rangle + a \leq \varepsilon f(x)$ , because of the assumption that  $f \geq 0$ . Therefore,  $\langle \phi/\varepsilon, \cdot \rangle + a/\varepsilon \in \mathcal{A}(f)$ . By Prop. 2.4,  $f^{**}(x_0) \geq \langle \phi/\varepsilon, x_0 \rangle + \frac{a}{\varepsilon}$ . As a consequence,  $\varepsilon f^{**}(x_0) \geq \langle \phi, x_0 \rangle + a$ . Letting  $\varepsilon \rightarrow 0$ , we obtain that  $0 \geq \langle \phi, x_0 \rangle + a$ , which contradicts (2.2.2).

We have thus shown that if  $f$  is convex, l.s.c. and non negative,  $f = f^{**}$ . Now consider the case where  $f$  is convex, l.s.c., but may take negative values. By Prop. 1.18,  $f$  admits an affine minorant, say  $\alpha$ . The function  $f - \alpha$  is convex, l.s.c. and non negative, and thus satisfies  $f - \alpha = (f - \alpha)^{**}$ . The proof is concluded upon noting that  $(f - \alpha)^{**} = f^{**} - \alpha$  which follows from straightforward algebra.  $\square$

**Definition 2.2.** Let  $f : \mathcal{X} \rightarrow [-\infty, +\infty]$ . The function  $\check{f} : \mathcal{X} \rightarrow [-\infty, +\infty]$  defined for every  $x \in \mathcal{X}$  by

$$\check{f}(x) = \sup\{g(x) : g : \mathcal{X} \rightarrow [-\infty, +\infty] \text{ is a convex and l.s.c. minorant of } f\},$$

is called the l.s.c. convex envelope of  $f$ .

**Lemma 2.6.** For every  $f : \mathcal{X} \rightarrow [-\infty, +\infty]$ ,  $\check{f}$  is l.s.c. and convex.

*Proof.* As an exercise.  $\square$

**Lemma 2.7.** Consider a convex function  $f : \mathcal{X} \rightarrow [-\infty, +\infty]$  and a point  $x \in \mathcal{X}$ . Then,  $f$  is l.s.c. at  $x$  if and only if  $f(x) = \check{f}(x)$ .

*Proof.* As an exercise.  $\square$

**Theorem 2.8.** Consider a convex function  $f : \mathcal{X} \rightarrow (-\infty, +\infty]$  and a point  $x \in \mathcal{X}$ . Then,  $f(x) = f^{**}(x)$  if and only if  $f$  is l.s.c. at  $x$ .

*Proof.* We prove that  $\check{f} = f^{**}$  and the result follows from lem. 2.7.

Clearly the set of convex and l.s.c. minorants of  $f$  contains  $\mathcal{A}(f)$ . This implies that  $\check{f} \geq f^{**}$  i.e., the l.s.c. convex envelope majorizes the affine envelope. We prove the other inequality. By definition of the l.s.c. convex envelope,  $\check{f} \leq f$ . By Lemma 2.1,  $\check{f}^* \geq f^*$  and thus  $\check{f}^{**} \leq f^{**}$ . By Prop. 1.18,  $f$  admits an affine minorant, and thus admits a real valued l.s.c. and convex minorant. Therefore,  $-\infty \notin \check{f}(\mathcal{X})$ . Lem. 2.6 and Th. 2.5 together imply that  $\check{f} = \check{f}^{**}$ . Thus,  $\check{f} \leq f^{**}$ . Hence,  $\check{f} = f^{**}$ .  $\square$

## The Fenchel-Young inequality and its consequences

**Proposition 2.9** (Fenchel - Young). Let  $f : \mathcal{X} \rightarrow [-\infty, \infty]$ . For all  $(x, \phi) \in \mathcal{X}^2$ , the following inequality holds:

$$f(x) + f^*(\phi) \geq \langle \phi, x \rangle,$$

with equality if and only if  $\phi \in \partial f(x)$ .

*Proof.* The inequality is an immediate consequence of the definition of  $f^*$ . The condition for equality to hold (i.e., for the converse inequality to be valid), is obtained with the equivalence

$$f(x) + f^*(\phi) \leq \langle \phi, x \rangle \Leftrightarrow \forall y, f(x) + \langle \phi, y \rangle - f(y) \leq \langle \phi, x \rangle \Leftrightarrow \phi \in \partial f(x).$$

$\square$

**Proposition 2.10.** *Let  $f : \mathcal{X} \rightarrow (-\infty, \infty]$  be proper, convex, and l.s.c. at some point  $x \in \mathcal{X}$ . Then,*

$$\phi \in \partial f(x) \Leftrightarrow x \in \partial f^*(\phi).$$

*Proof.* By the equality case in Prop. 2.9,  $\phi \in \partial f(x)$  iff  $f(x) + f^*(\phi) = \langle \phi, x \rangle$ . Since  $f$  is l.s.c. at  $x \in \mathcal{X}$ , Th. 2.8 implies that  $f(x) = f^{**}(x)$ . Thus,  $\phi \in \partial f(x)$  iff  $f^{**}(x) + f^*(\phi) = \langle \phi, x \rangle$ . By Prop. 2.9 again, this is equivalent to  $x \in \partial f^*(\phi)$ .  $\square$

If  $A : \mathcal{X} \rightarrow 2^{\mathcal{X}}$  is a set-valued mapping, we define  $A^{-1} : \mathcal{X} \rightarrow 2^{\mathcal{X}}$  the **inverse** of  $A$  as

$$A^{-1}(x) = \{y : x \in A(y)\},$$

namely,  $y \in A^{-1}(x) \Leftrightarrow x \in A(y)$ .

**Corollary 2.11.** *Let  $f \in \Gamma_0(\mathcal{X})$ . Then,  $\partial f^{-1} = \partial f^*$ .*

*Proof.* The result follows from Prop. 2.10 and the Fenchel-Moreau theorem Th. 2.5.  $\square$

**Theorem 2.12.** *Consider  $f : \mathcal{X} \rightarrow (-\infty, +\infty]$  a proper and convex function. Let  $x \in \mathcal{X}$ . Consider the following assertions:*

- i)  $x \in \text{ri}(\text{dom } f)$ ,
- ii)  $\partial f(x) \neq \emptyset$ ,
- iii)  $f$  is l.s.c. at  $x$ ,
- iv)  $f(x) = f^{**}(x)$ .
- v)  $\partial f(x) = \partial f^{**}(x)$ .

*Then,  $i) \Rightarrow ii) \Rightarrow iii) \Leftrightarrow iv) \Rightarrow v)$ .*

*Proof.*  $i) \Rightarrow ii)$  is given by Th. 1.17.  $iii) \Leftrightarrow iv)$  is given by Th. 2.8.

We prove that  $ii) \Rightarrow iv)$ . Choose  $\phi \in \partial f(x)$ . By the equality case in the Fenchel-Young inequality (Prop. 2.9),  $f(x) = \langle \phi, x \rangle - f^*(\phi)$ . Thus,  $f(x) \leq \sup_{\psi \in \mathcal{X}} \langle \psi, x \rangle - f^*(\psi)$  and the righthand side coincides with  $f^{**}(x)$ . Thus,  $f(x) \leq f^{**}(x)$ . On the otherhand  $f^{**} \leq f$  since  $f^{**}$  is the affine envelope of  $f$  (Prop. 2.4). Thus  $f(x) = f^{**}(x)$ .

Finally, we prove that  $iv) \Rightarrow v)$ . Consider  $\phi \in \mathcal{X}$ . By Prop. 2.9,  $\phi \in \partial f(x) \Leftrightarrow f(x) + f^*(\phi) = \langle \phi, x \rangle$ . By the standing hypothesis, this rewrites as  $f^{**}(x) + f^*(\phi) = \langle \phi, x \rangle$ . Since  $f^*$  is convex, l.s.c., and does not take the value  $-\infty$  (as  $f$  is proper), Th. 2.5 implies that  $f^* = f^{***}$ . Therefore, the following equivalence holds:  $\phi \in \partial f(x) \Leftrightarrow f^{**}(x) + f^{***}(\phi) = \langle \phi, x \rangle$ . By Prop. 2.9, this is equivalent to  $\phi \in \partial f^{**}(x)$ .  $\square$

## Exercises

**Exercise 2.1** (invariance by Fenchel-Legendre conjugation). Let  $f : \mathcal{X} \rightarrow [-\infty, +\infty]$ . Show that  $f = f^* \Leftrightarrow f(x) = \frac{1}{2}\|x\|^2$ .

**Exercise 2.2.** Set  $\mathcal{X} := \mathbb{R}^m$  where  $m \geq 1$  is an integer. Prove that the Fenchel conjugate of  $\iota_{(-\infty, 0]^m}$  is equal to  $\iota_{[0, +\infty)^m}$ .

**Exercise 2.3.** Let  $E$  be a linear subspace of  $\mathcal{X}$  and denote by  $E^\perp$  its supplementary set. Prove that the Fenchel conjugate of  $\iota_E$  is equal to  $\iota_{E^\perp}$ .

**Exercise 2.4.** 1. On  $\mathcal{X} := \mathbb{R}$ , define

$$f(x) = \begin{cases} 1/x & \text{if } x > 0; \\ +\infty & \text{otherwise .} \end{cases}$$

Prove that

$$f^*(\phi) = \begin{cases} -2\sqrt{-\phi} & \text{if } \phi \leq 0; \\ +\infty & \text{otherwise .} \end{cases}$$

2. On  $\mathcal{X} := \mathbb{R}$ , define  $f(x) = \exp(x)$ . Prove that

$$f^*(\phi) = \begin{cases} \phi \ln(\phi) - \phi & \text{if } \phi > 0; \\ 0 & \text{if } \phi = 0; \\ +\infty & \text{if } \phi < 0. \end{cases}$$

**Exercise 2.5** (Fenchel-Legendre transform of an infimal post composition). Consider  $f : \mathcal{X} \rightarrow (-\infty, +\infty]$  and let  $M : \mathcal{X} \rightarrow \mathcal{Y}$  be a linear operator. Prove that for every  $\phi \in \mathcal{Y}$ ,  $(M \triangleright f)^*(\phi) = f^*(M^*\phi)$ .

**Exercise 2.6** (Fenchel-Legendre transform of a convex coercive function). Show that a function  $f : \mathcal{X} \rightarrow (-\infty, +\infty]$  satisfies  $f(x) \geq \gamma\|x\| + \alpha$  for some  $\gamma > 0$  and  $\alpha \in \mathbb{R}$  iff  $0 \in \text{int dom } f^*$ .

*Hint:*  $0 \in \text{int dom } f^* \Leftrightarrow f^*$  is bounded on a small ball around zero.

Deduce that when  $f$  is convex, then  $f$  is coercive iff  $0 \in \text{int dom } f^*$ .

**Exercise 2.7** (Fenchel-Legendre transform of an infimal convolution). Let  $f, g : \mathcal{X} \rightarrow (-\infty, +\infty]$  be proper convex functions.

1. Assume that  $f$  and  $g$  both admit an affine minorant with gradient  $\phi$ . Show that  $f \square g$  has an affine minorant of gradient  $\phi$ .
2. Prove that  $(f \square g)^* = f^* + g^*$ .

**Exercise 2.8** (Moreau's identity). Consider  $f \in \Gamma_0(\mathcal{X})$ . The goal of this exercise is to show Moreau's identity: for every  $x \in \mathbb{R}^n$ ,

$$\text{prox}_f(x) + \text{prox}_{f^*}(x) = x.$$

1. Let  $x \in \mathbb{R}^n$  and  $p = \text{prox}_f(x)$ . Show that  $x - p \in \partial f(p)$ .
2. Using the result of Exercise ??, show that  $p \in \partial f^*(x - p)$ .
3. Prove Moreau's identity.
4. Show that Moreau's formula generalizes the famous identity  $\Pi_E + \Pi_{E^\perp} = I$ , where  $\Pi_E$  and  $\Pi_{E^\perp}$  are the orthogonal projectors onto some linear subspace  $E \subset \mathbb{R}^n$  and its supplementary space  $E^\perp$  respectively.

*Hint: choose  $f$  as the indicator function of a properly chosen set.*

5. Generalization: for  $\gamma > 0$ , prove that for all  $x \in \mathcal{X}$ ,

$$\text{prox}_{\gamma f}(x) + \gamma \text{prox}_{\gamma^{-1} f^*}\left(\frac{x}{\gamma}\right) = x. \quad (2.4.1)$$

# Chapter 3

## Duality

### Parametric Duality

In this paragraph, we consider a function  $F : \mathcal{X} \times \mathcal{Y} \rightarrow (-\infty, +\infty]$  where  $\mathcal{X}$  and  $\mathcal{Y}$  are two Euclidean spaces. Without risk of ambiguity, we can denote with the same symbol  $\langle \cdot, \cdot \rangle$  the scalar products in each of these spaces. We equip the product space  $\mathcal{X} \times \mathcal{Y}$  with the scalar product  $\langle (\nu, \phi), (x, y) \rangle := \langle \nu, x \rangle + \langle \phi, y \rangle$ .

#### Primal and dual problems

The **primal function** and the **dual function** are the mappings  $\mathcal{P} : \mathcal{X} \rightarrow [-\infty, +\infty]$  and  $\mathcal{D} : \mathcal{Y} \rightarrow [-\infty, +\infty]$  respectively defined by

$$\begin{aligned}\mathcal{P}(x) &= F(x, 0) \quad (\forall x \in \mathcal{X}) \\ \mathcal{D}(\phi) &= -F^*(0, -\phi) \quad (\forall \phi \in \mathcal{Y}).\end{aligned}$$

We respectively refer to  $p := \inf \mathcal{P}$  and  $d := \sup \mathcal{D}$  as the **primal (resp. dual) value** and to  $\arg \min \mathcal{P}$  and  $\arg \max \mathcal{D}$  as the set of **primal (resp. dual) solutions**. Finally, the **value function** is the mapping  $\vartheta : \mathcal{Y} \rightarrow [-\infty, +\infty]$  defined for all  $y \in \mathcal{Y}$  by

$$\vartheta(y) := \inf F(\mathcal{X}, y).$$

**Lemma 3.1.** *The following holds:*

1.  $p = \vartheta(0)$  and  $d = \vartheta^{**}(0)$ .
2. For every  $\phi \in \mathcal{Y}$ ,  $\mathcal{D}(\phi) = -\vartheta^*(-\phi)$ .
3. If  $\vartheta^*$  is proper,  $\arg \max \mathcal{D} = -\partial \vartheta^{**}(0)$ . Otherwise,  $\arg \max \mathcal{D} = \mathcal{Y}$  and  $d \in \{-\infty, +\infty\}$ .

*Proof.* The equality  $p = \vartheta(0)$  is obvious from the definition. For every  $\phi \in \mathcal{Y}$ ,

$$\vartheta^*(\phi) = \sup_{(x,y) \in \mathcal{X} \times \mathcal{Y}} \langle \phi, y \rangle - F(x, y) = F^*(0, \phi),$$

and the second point is proved. The fact that  $d = \vartheta^{**}(0)$ , follows from the definition of  $\vartheta^{**}(0) = \sup_{\phi \in \mathcal{Y}} -\vartheta^*(\phi)$ . The first point is proved. This also shows that the dual solutions are the

opposite of the minimizers of  $\vartheta^*$ :  $\arg \max \mathcal{D} = -\arg \min \vartheta^*$ . If  $\vartheta$  is proper, then

$$\begin{aligned} \phi \in \arg \min \vartheta^* &\stackrel{(a)}{\Leftrightarrow} 0 \in \arg \min \vartheta^*(\phi) \\ &\stackrel{(b)}{\Leftrightarrow} \vartheta^*(\phi) + \vartheta^{**}(0) = 0 \\ &\stackrel{(c)}{\Leftrightarrow} \vartheta^{***}(\phi) + \vartheta^{**}(0) = 0 \\ &\stackrel{(d)}{\Leftrightarrow} \phi \in \partial \vartheta^{**}(0), \end{aligned}$$

where the equivalence (a) is due to Prop. 1.22, (b) is due to Prop. 2.9, (c) uses the fact that  $\vartheta^* = \vartheta^{***}$  which follows from Th. 2.5 and (d) is due to Prop. 2.9 again. This completes the proof of the third point when  $\vartheta$  is proper. If  $\vartheta$  is not proper, then  $\vartheta^*$  is everywhere  $+\infty$  or everywhere  $-\infty$ . In both cases,  $\arg \min \vartheta^* = \mathcal{Y}$ .  $\square$

The quantity  $(p - d)$  is called the **duality gap**.

**Theorem 3.2.** *It holds that  $p \geq d$ . Moreover, if  $\vartheta$  is convex and if*

$$0 \in \text{ri dom } \vartheta, \quad (3.1.1)$$

*then the following holds:*

1.  $d = p < +\infty$ .
2. The set of dual solutions  $\arg \max \mathcal{D}$  is non empty, and coincides with  $-\partial \vartheta(0)$ .

*Proof.* By Prop. 2.4,  $p \geq d$ . We prove that  $p < +\infty$ . By Eq. (3.1.1),  $0 \in \text{dom } \vartheta$ , thus  $\vartheta(0) < +\infty$ , which reads  $p < +\infty$  by Lemma 3.1.

First assume that  $\vartheta$  is proper. Th. 2.10 along with Eq. (3.1.1) imply that  $\vartheta(0) = \vartheta^{**}(0)$  which reads  $p = d$  by Lemma 3.1. The first point is proved. This also implies that  $\partial \vartheta(0) = \partial \vartheta^{**}(0)$ . By Lemma 3.1, we obtain that  $\arg \max \mathcal{D} = -\partial \vartheta(0)$  (and this set is non empty). Hence, the second point is proved.

Finally, consider the case where  $\vartheta$  is not proper. Since,  $\vartheta(0) < +\infty$ , this means that  $-\infty \in \vartheta(\mathcal{Y})$ . By contradiction, assume that  $p \in \mathbb{R}$ , which reads  $\vartheta(0) \in \mathbb{R}$  by Lemma 3.1. By Th. 1.17, there exists  $\phi \in \partial \vartheta(0)$ . Thus,  $\vartheta$  admits  $\vartheta(0) + \langle \phi, \cdot \rangle$  as an affine minorant. This contradicts the fact that  $-\infty \in \vartheta(\mathcal{Y})$ . We have shown that  $p = -\infty$ , which also implies that  $d = -\infty$ : the first point is proved. Since  $\vartheta(0) = -\infty$ , it holds that  $\partial \vartheta(0) = \mathcal{Y}$ . Since  $\arg \max \mathcal{D} = \mathcal{Y}$  by Lemma 3.1, the second point holds as well.  $\square$

Eq. (3.1.1) is sometimes referred to as a **qualification condition**.

**Remark 3.1.** If  $F$  is convex, then  $\vartheta$  is convex by Prop. 1.11.

## Lagrangian

**Definition 3.1.** The **Lagrangian** associated with the function  $F$  is the mapping on  $\mathcal{X} \times \mathcal{Y} \rightarrow [-\infty, +\infty]$  defined by:

$$\mathcal{L} : (x, \phi) \mapsto \inf_{y \in \mathcal{Y}} (F(x, y) + \langle \phi, y \rangle).$$

*Assumption 3.1.* For all  $x \in \mathcal{X}$ ,  $F(x, \cdot)$  is convex and l.s.c. at zero.

**Proposition 3.3.** *For every  $\phi \in \mathcal{Y}$ ,*

$$\mathcal{D}(\phi) = \inf \mathcal{L}(\mathcal{X}, \phi).$$

*Moreover, under Assumption 3.1, for every  $x \in \mathcal{X}$ ,*

$$\mathcal{P}(x) = \sup \mathcal{L}(x, \mathcal{Y}).$$

*Proof.* After straightforward algebra, we obtain for every  $(x, \phi) \in \mathcal{X} \times \mathcal{Y}$ ,

$$\mathcal{L}(x, \phi) = -\sup_{y \in \mathcal{Y}} \langle -\phi, y \rangle - F(x, y). \quad (3.1.2)$$

Thus,  $\inf \mathcal{L}(\mathcal{X}, \phi) = -\sup_{(x, y)} \langle -\phi, y \rangle - F(x, y) = -F^*(0, -\phi) = \mathcal{D}(\phi)$ . To prove the second point, set  $x \in \mathcal{X}$  and consider the mapping  $F_x : y \mapsto F(x, y)$ . By Eq. (3.1.2),  $\mathcal{L}(x, \phi) = -F_x^*(-\phi)$ . Thus,  $\sup_{\phi} \mathcal{L}(x, \phi) = F_x^{**}(0)$ . By Th. 2.8,  $F_x^{**}(0) = F_x(0)$  and since  $\mathcal{P}(x) = F_x(0)$ , the result is proved.  $\square$

As an immediate consequence of Prop. 3.3, the following equalities hold under Assumption 3.1:

$$p = \inf_{x \in \mathcal{X}} \sup_{\phi \in \mathcal{Y}} \mathcal{L}(x, \phi) \quad (3.1.3)$$

$$d = \sup_{\phi \in \mathcal{Y}} \inf_{x \in \mathcal{X}} \mathcal{L}(x, \phi). \quad (3.1.4)$$

**Definition 3.2.** We say that  $(x, \phi)$  is a **saddle point** of  $\mathcal{L}$  if

$$\inf \mathcal{L}(\mathcal{X}, \phi) = \mathcal{L}(x, \phi) = \sup \mathcal{L}(x, \mathcal{Y}).$$

**Theorem 3.4.** Under Assumption 3.1, the following properties are equivalent:

- i)  $x$  is primal optimal,  $\phi$  is dual optimal and  $p = d$ .
- ii)  $(x, \phi)$  is a saddle point of the Lagrangian.

*Proof.* i)  $\Rightarrow$  ii)  $d = \mathcal{D}(\phi) = \inf \mathcal{L}(\mathcal{X}, \phi) \leq \mathcal{L}(x, \phi) \leq \sup \mathcal{L}(x, \mathcal{Y}) = \mathcal{P}(x) = p$ .

ii)  $\Rightarrow$  i) As  $(x, \phi)$  is a saddle point, we have on the one hand  $\mathcal{L}(x, \phi) = \sup \mathcal{L}(x, \mathcal{Y}) = \mathcal{P}(x)$  and on the other hand  $\mathcal{L}(x, \phi) = \inf \mathcal{L}(\mathcal{X}, \phi) = \mathcal{D}(\phi)$ . Thus,  $\mathcal{P}(x) = \mathcal{D}(\phi)$ . Hence  $p \leq \mathcal{P}(x) = \mathcal{D}(\phi) \leq d$ . We conclude using  $p \geq d$  (Th. 3.2).  $\square$

**Corollary 3.5.** Assume that  $\vartheta$  is convex and let Assumption 3.1 hold true. If

$$0 \in \text{ri dom } \vartheta,$$

then  $d = p < +\infty$  and the two following properties are equivalent:

- i)  $x$  is primal optimal.
- ii) There exists  $\phi \in \mathcal{Y}$  such that  $(x, \phi)$  is a saddle point of the Lagrangian.

In that case, every  $\phi$  satisfying ii) is a dual solution.

*Proof.* i)  $\Rightarrow$  ii). By Th. 3.2,  $d = p < +\infty$  and there exists  $\phi \in \arg \max \mathcal{D}$ . By Th. 3.4,  $(x, \phi)$  is a saddle point of  $\mathcal{L}$ .

ii)  $\Rightarrow$  i) is an immediate consequence of Th. 3.4. The same holds for the last statement of the theorem.  $\square$

## Fenchel-Rockafellar duality

In this section, we introduce two functions  $f : \mathcal{X} \rightarrow (-\infty, +\infty]$ ,  $g : \mathcal{Y} \rightarrow (-\infty, +\infty]$  and a linear operator  $M : \mathcal{X} \rightarrow \mathcal{Y}$ . We address the special case where  $F : \mathcal{X} \times \mathcal{Y} \rightarrow (-\infty, +\infty]$  is given by

$$F : (x, y) \mapsto f(x) + g(Mx - y). \quad (3.2.1)$$

Eq. (3.2.1) is a standing assumption throughout Section 3.2.

## Main results

**Proposition 3.6.** *Let  $f : \mathcal{X} \rightarrow (-\infty, +\infty]$  and  $g : \mathcal{Y} \rightarrow (-\infty, +\infty]$  be two functions and let  $M : \mathcal{X} \rightarrow \mathcal{Y}$  be a linear operator. The following holds.*

i) *The primal function  $\mathcal{P} : \mathcal{X} \rightarrow [-\infty, +\infty]$  is given by*

$$\mathcal{P} : x \mapsto f(x) + g(Mx).$$

ii) *The Lagrangian  $\mathcal{L} : \mathcal{X} \times \mathcal{Y} \rightarrow [-\infty, +\infty]$  is given by*

$$\mathcal{L} : (x, \phi) \mapsto \begin{cases} f(x) + \langle \phi, Mx \rangle - g^*(\phi) & \text{if } x \in \text{dom } f, \\ +\infty & \text{otherwise.} \end{cases} \quad (3.2.2)$$

iii) *The dual function  $\mathcal{D}$  on  $\mathcal{Y} \rightarrow [-\infty, +\infty]$  is given by*

$$\mathcal{D} : \phi \mapsto -f^*(-M^*\phi) - g^*(\phi).$$

*Proof.* Point (i) is immediate. For every  $(x, \phi) \in \mathcal{X} \times \mathcal{Y}$ ,  $\mathcal{L}(x, \phi) = \inf_{y \in \mathcal{Y}} f(x) + g(Mx - y) + \langle \phi, y \rangle$ . In particular,  $\mathcal{L}(x, \phi) = +\infty$  if  $x \notin \text{dom } f$ . If  $x \in \text{dom } f$ ,

$$\begin{aligned} \mathcal{L}(x, \phi) &= f(x) + \langle \phi, Mx \rangle + \inf_{y \in \mathcal{Y}} g(Mx - y) - \langle \phi, Mx - y \rangle \\ &= f(x) + \langle \phi, Mx \rangle + \inf_{w \in \mathcal{Y}} g(w) - \langle \phi, w \rangle \\ &= f(x) + \langle \phi, Mx \rangle - \sup_{w \in \mathcal{Y}} \langle \phi, w \rangle - g(w), \end{aligned}$$

which proves (ii). By Prop. 3.3,  $\mathcal{D}(\phi) = \inf \mathcal{L}(\mathcal{X}, \phi)$  for all  $\phi \in \mathcal{Y}$ . Hence,

$$\begin{aligned} \mathcal{D}(\phi) &= \inf_{x \in \text{dom } f} (f(x) + \langle \phi, Mx \rangle - g^*(\phi)) \\ &= - \sup_{x \in \text{dom } f} (\langle -\phi, Mx \rangle - f(x)) - g^*(\phi) \\ &= -f^*(-M^*\phi) - g^*(\phi). \end{aligned}$$

This proves (iii). □

**Theorem 3.7** (Fenchel-Rockafellar). *Let  $f : \mathcal{X} \rightarrow (-\infty, +\infty]$  and  $g : \mathcal{Y} \rightarrow (-\infty, +\infty]$  be two convex functions and let  $M : \mathcal{X} \rightarrow \mathcal{Y}$  be a linear operator. Assume that*

$$0 \in \text{ri}(M \text{ dom } f - \text{dom } g). \quad (3.2.3)$$

*Then,*

$$\inf_{x \in \mathcal{X}} f(x) + g(Mx) = - \min_{\phi \in \mathcal{Y}} f^*(-M^*\phi) + g^*(\phi). \quad (3.2.4)$$

*Proof.* The value function associated with the function  $F$  defined in (3.2.1) is the function  $\vartheta : \mathcal{Y} \rightarrow [-\infty, +\infty]$  given by  $\vartheta : y \mapsto \inf_{x \in \mathcal{X}} f(x) + g(Mx - y)$ . As  $f, g$  are convex,  $F$  is convex by Prop. 1.13 and Prop. 1.14. Thus,  $\vartheta$  is convex by Prop. 1.12. A point  $y$  belongs to  $\text{dom } \vartheta$  iff there exists  $x \in \mathcal{X}$  such that  $f(x) + g(Mx - y) < +\infty$ . This is equivalent to the existence of  $x \in \text{dom } f$ , such that  $Mx - y \in \text{dom } g$ . Otherwise stated,  $y \in Mx - \text{dom } g$ . This shows that  $\text{dom } \vartheta = M \text{ dom } f - \text{dom } g$ . Thus condition (3.2.3) reads  $0 \in \text{ri dom } \vartheta$ . By Th. 3.2, we obtain that  $\inf \mathcal{P} = \max \mathcal{D}$ . This is equivalent to (3.2.4) by Prop. 3.6. □



**Proposition 3.8.** Consider a function  $f : \mathcal{X} \rightarrow (-\infty, +\infty]$ , a function  $g \in \Gamma_0(\mathcal{Y})$  and a linear map  $M : \mathcal{X} \rightarrow \mathcal{Y}$ . Let  $(x, \phi)$  be a point in  $\mathcal{X} \times \mathcal{Y}$ . The following statements are equivalent:

i)  $(x, \phi)$  satisfies

$$0 \in \partial f(x) + M^* \phi \quad (3.2.5a)$$

$$0 \in -Mx + \partial g^*(\phi). \quad (3.2.5b)$$

ii)  $(x, \phi)$  is a saddle point of the Lagrangian.

iii)  $x$  is a primal solution,  $\phi$  is a dual solution, and  $p = d$ .

*Proof.* As  $g \in \Gamma_0(\mathcal{Y})$ , the Assumption 3.1 is satisfied for  $F$  given by (3.2.1). By Th. 3.4, this shows that ii)  $\Leftrightarrow$  iii).

By Prop. 3.6,  $(x, \phi)$  is a saddle point of the Lagrangian (3.2.2) iff

$$\begin{cases} x \in \text{dom } f \\ x \in \arg \min_{w \in \mathcal{X}} (f(w) + \langle \phi, Mw \rangle - g^*(\phi)) \\ \phi \in \arg \min_{\psi \in \mathcal{Y}} (f(x) - \langle \psi, Mx \rangle + g^*(\psi)) \end{cases} \quad (3.2.6)$$

By Fermat's rule (Prop. 1.22) and Exercice 1.4, the last two inclusions are equivalent to Eq. (3.2.5). Since  $\partial f$  takes the value  $\emptyset$  outside  $\text{dom } f$ , the condition  $x \in \text{dom } f$  can be discarded from 3.2.6. Thus, i)  $\Leftrightarrow$  ii).  $\square$

**Proposition 3.9.** Consider a convex function  $f : \mathcal{X} \rightarrow (-\infty, +\infty]$ , a function  $g \in \Gamma_0(\mathcal{Y})$  and a linear map  $M : \mathcal{X} \rightarrow \mathcal{Y}$ . Assume that condition (3.2.3) holds. Then, a point  $x \in \mathcal{X}$  is primal optimal iff there exists  $\phi \in \mathcal{Y}$  s.t. the condition (3.2.5) holds (in which case, any such  $\phi$  is dual optimal).

*Proof.* As in the proof of Th. 3.7, the stated hypotheses imply that  $\vartheta$  is convex and  $0 \in \text{ri dom } \vartheta$ . The result is an application of Cor. 3.5 along with the characterization of the saddle points of  $\mathcal{L}$  by Prop. 3.8.  $\square$

## Affine equality constraints

Here, we make the hypothesis that  $g := \iota_{\{b\}}$  where  $b \in \mathcal{Y}$ . This will be a standing assumption in Section 3.2.2.

In this case, the primal optimal points coincide with the solutions to the following minimization problem:

**Minimization under affine equality constraints.**

$$\text{Minimize } f(x) \text{ w.r.t. } x \in \mathcal{X} \text{ s.t. } Mx = b. \quad (3.2.7)$$

The set  $\{x \in \text{dom } f : Mx = b\}$  is called the **feasible set**. A point of  $\mathcal{X}$  is said **feasible** if it belongs to the feasible set.

**Proposition 3.10.** Consider a function  $f : \mathcal{X} \rightarrow (-\infty, +\infty]$ , a linear map  $M : \mathcal{X} \rightarrow \mathcal{Y}$  and a point  $b \in \mathcal{Y}$ . Then,

i) The primal function is

$$\mathcal{P} : x \mapsto \begin{cases} f(x) & \text{if } x \text{ is feasible} \\ +\infty & \text{otherwise.} \end{cases}$$

The primal optimal points are the solutions to Problem (3.2.7).

ii) The lagrangian  $\mathcal{L} : \mathcal{X} \times \mathcal{Y} \rightarrow [-\infty, +\infty]$  is given by

$$\mathcal{L} : (x, \phi) \mapsto \begin{cases} f(x) + \langle \phi, Mx - b \rangle & \text{if } x \in \text{dom } f \\ +\infty & \text{otherwise.} \end{cases}$$

iii) The dual function is  $\mathcal{D} : \phi \mapsto -f^*(-M^*\phi) - \langle b, \phi \rangle$ .

*Proof.* Apply Prop. 3.6 with  $g := \iota_{\{b\}}$ . In this case,  $g \circ M = \iota_{M^{-1}(b)}$ . This proves point (i). Points (ii-iii) follows from the identity  $g^* = \langle \cdot, b \rangle$ .  $\square$

**Proposition 3.11.** Consider a function  $f : \mathcal{X} \rightarrow (-\infty, +\infty]$ , a linear map  $M : \mathcal{X} \rightarrow \mathcal{Y}$  and a point  $b \in \mathcal{Y}$ . Let  $(x, \phi)$  be a point in  $\mathcal{X} \times \mathcal{Y}$ . The following statements are equivalent.

i)  $(x, \phi)$  satisfies

$$\begin{cases} x \text{ is feasible} \\ 0 \in \partial f(x) + M^*\phi. \end{cases} \quad (3.2.8)$$

ii)  $(x, \phi)$  is a saddle point of the Lagrangian.

iii)  $x$  is a primal solution,  $\phi$  is a dual solution, and  $p = d$ .

*Proof.* Use Prop. 3.8 with  $g = \iota_{\{b\}}$ .  $\square$

**Proposition 3.12.** Consider a convex function  $f : \mathcal{X} \rightarrow (-\infty, +\infty]$ , a linear map  $M : \mathcal{X} \rightarrow \mathcal{Y}$  and a point  $b \in \mathcal{Y}$ . Assume that

$$\exists x_0 \in \text{ri}(\text{dom } f), Mx_0 = b. \quad (3.2.9)$$

Then, a point  $x \in \mathcal{X}$  is primal optimal iff there exists  $\phi \in \mathcal{Y}$  s.t. the conditions (3.2.8) hold (in which case, any such  $\phi$  is dual optimal).

*Proof.* Set  $g = \iota_{\{b\}}$ . Condition (3.2.3) is equivalent to  $b \in \text{ri}(M \text{ dom } f)$ . By Th. 1.7, this is again equivalent to  $b \in M \text{ ri dom } f$ , hence the condition (3.2.9). The result follows from Prop. 3.9.  $\square$

## Operations on subdifferentials

**Proposition 3.13.** Let  $f : \mathcal{X} \rightarrow (-\infty, +\infty]$  and  $g : \mathcal{Y} \rightarrow (-\infty, +\infty]$  be two convex functions and  $M : \mathcal{X} \rightarrow \mathcal{Y}$  be a linear operator. Consider  $x \in \mathcal{X}$ . Then,

$$\partial f(x) + M^*\partial g(Mx) \subset \partial(f + g \circ M)(x). \quad (3.2.10)$$

Moreover, if the following condition holds:

$$0 \in \text{ri}(M \text{ dom } f - \text{dom } g),$$

then

$$\partial(f + g \circ M)(x) = \partial f(x) + M^*\partial g(Mx).$$

*Proof.* In order to prove (3.2.10), consider  $\phi \in \partial f(x)$  and  $\psi \in \partial g(Mx)$ . For every  $y \in \mathcal{X}$ , it holds that

$$\begin{aligned} f(y) &\geq f(x) + \langle \phi, y - x \rangle \\ g(My) &\geq g(Mx) + \langle \psi, M(y - x) \rangle . \end{aligned}$$

Summing the inequalities, we obtain  $(f + g \circ M)(y) \geq (f + g \circ M)(x) + \langle \phi + M^*\psi, y - x \rangle$ . This proves that  $\phi + M^*\psi \in \partial(f + g \circ M)(x)$  and Eq. (3.2.10) follows.

Let  $x \in \mathcal{X}$  and  $\phi \in \mathcal{X}$ . Applying Theorem 3.7 on  $(f - \langle \phi, \cdot \rangle)$  and  $g$ , we can see that there exists  $\psi \in \mathcal{Y}$  such that  $f(x) - \langle \phi, x \rangle + g(Mx) = -(f - \langle \phi, \cdot \rangle)^*(-M^*\psi) - g^*(\psi)$ . One can easily check that  $(f - \langle \phi, \cdot \rangle)^* = f^*(\cdot + \phi)$ . Hence,  $f(x) - \langle \phi, x \rangle + g(Mx) = -f(-M^*\psi + \phi) - g^*(\psi)$ . Said otherwise,

$$f(x) + f(-M^*\psi + \phi) - \langle \phi, x \rangle + g(Mx) + g^*(\psi) = 0 ,$$

and so

$$[f(x) + f(-M^*\psi + \phi) - \langle -M^*\psi + \phi, x \rangle] + [g(Mx) + g^*(\psi) - \langle \psi, Mx \rangle] = 0 .$$

Each of the two terms in the brackets is nonnegative according to Fenchel-Young inequality. Thus both are null. The equality case in Fenchel-Young inequality allows us to conclude that  $\psi \in \partial g(Mx)$  et  $-M^*\psi + \phi \in \partial f(x)$ . The last inclusion can be read also as  $\phi \in \partial f(x) + M^*\partial g(Mx)$  which is .  $\square$

**Corollary 3.14.** *Let  $m \geq 2$  be an integer and consider functions  $f_1, \dots, f_m : \mathcal{X} \rightarrow (-\infty, +\infty]$ . Then,  $\partial f_1 + \dots + \partial f_m \subset \partial(f_1 + \dots + f_m)$ .*

*Moreover, if the functions are convex and satisfy the condition*

$$0 \in \bigcap_{i=2}^m \text{ri} \left( \text{dom } f_i - \bigcap_{j=1}^{i-1} \text{dom } f_j \right)$$

*or the stronger condition*

$$0 \in \bigcap_{i=1}^m \text{ri } \text{dom } f_i .$$

*Then,  $\partial(f_1 + \dots + f_m) = \partial f_1 + \dots + \partial f_m$ .*

*Proof.* By induction using Prop. 3.13.  $\square$

### The Attouch-Brezis theorem\*

Let  $f$  and  $g$  be two functions on  $\mathcal{X} \rightarrow (-\infty, +\infty]$ . We recall that  $(f \square g)^* = f^* + g^*$  (see Exercise 2.7).

**Theorem 3.15** (Attouch-Brezis). *Let  $f, g : \mathcal{X} \rightarrow (-\infty, +\infty]$  be two convex functions such that*

$$0 \in \text{ri}(\text{dom } f - \text{dom } g) . \tag{3.2.11}$$

*Then  $(f + g)^* = f^* \square g^*$ , and this function belongs to  $\Gamma_0(\mathcal{X})$ . Moreover, the infimal convolution is exact.*

*Proof.* Let  $\psi \in \mathcal{X}$ . Note that  $(f + g)^*(\psi) = -\inf_{x \in \mathcal{X}} (f - \langle \psi, \cdot \rangle)(x) + g(x)$ . Clearly,  $\text{dom}(f - \langle \psi, \cdot \rangle) = \text{dom } f$ . Thus, Eq. (3.2.11) and Th. 3.7 imply that  $(f + g)^*(\psi) = \min_{\phi \in \mathcal{X}} (f - \langle \psi, \cdot \rangle)^*(-\phi) + g^*(\phi)$ . It is straightforward to show that  $(f - \langle \psi, \cdot \rangle)^*(-\phi) = f^*(\psi - \phi)$ , which proves the first part of the result.

Finally, the condition (3.2.11) implies that  $0 \in \text{dom } f - \text{dom } g$  which also reads  $\text{dom } f \cap \text{dom } g \neq \emptyset$ . By Remark 2.1,  $(f + g)^*$  does not take the value  $-\infty$ . Also Prop. 2.3 and Prop. 1.18,  $(f + g)^*$  is not identically  $+\infty$ . Hence,  $(f + g)^*$  is proper.  $\square$

## Lagrangian duality

In this paragraph, we consider an integer  $m \geq 1$ . For every  $(a, b) \in \mathbb{R}^m \times \mathbb{R}^m$ , the notations  $a \leq b$  (or  $b \geq a$ ) is used as an equivalent to  $a - b \in (-\infty, 0]^m$ . Similarly, we write  $a < b$  as a short for  $a - b \in (-\infty, 0)^m$ . Finally, we denote by  $a_1, \dots, a_m$  the real components of any vector  $a \in \mathbb{R}^m$ .

### Inequality and affine equality constraints

Consider a function  $f : \mathcal{X} \rightarrow (-\infty, +\infty]$ , a linear operator  $M : \mathcal{X} \rightarrow \mathcal{Y}$ , a point  $b \in \mathcal{Y}$ , and a map  $g : \mathcal{X} \rightarrow \mathbb{R}^m$ . We address the case where the function  $F : \mathcal{X} \times (\mathcal{Y} \times \mathbb{R}^m) \rightarrow [-\infty, +\infty]$  is defined for every  $x \in \mathcal{X}$ ,  $u \in \mathcal{Y}$ ,  $v \in \mathbb{R}^m$ , by

$$F(x, (u, v)) = f(x) + \iota_{\{b\}}(Mx - u) + \iota_{(-\infty, 0]^m}(g(x) - v). \quad (3.3.1)$$

We consider Eq. (3.3.1) as a standing assumption throughout Section 3.3.

In this case, the primal optimal points coincide with the solutions to the following minimization problem:

**Minimization under inequality and affine equality constraints.**

$$\text{minimize } f(x) \text{ w.r.t. } x \in \mathcal{X} \text{ s.t. } Mx = b \text{ and } g(x) \leq 0. \quad (3.3.2)$$

The set  $\{x \in \text{dom } f : Mx = b, g(x) \leq 0\}$  is called the **feasible set**. A point of  $\mathcal{X}$  is said **feasible** if it belongs to the feasible set.

**Proposition 3.16.** *Consider a function  $f : \mathcal{X} \rightarrow (-\infty, +\infty]$ , a linear map  $M : \mathcal{X} \rightarrow \mathcal{Y}$ , a point  $b \in \mathcal{Y}$  and a map  $g : \mathcal{X} \rightarrow \mathbb{R}^m$ . Then,*

i) *The primal function is*

$$\mathcal{P} : x \mapsto \begin{cases} f(x) & \text{if } x \text{ is feasible} \\ +\infty & \text{otherwise.} \end{cases}$$

*The primal optimal points are the solutions to Problem (3.3.2).*

ii) *The lagrangian  $\mathcal{L} : \mathcal{X} \times (\mathcal{Y} \times \mathbb{R}^m) \rightarrow [-\infty, +\infty]$  is defined by*

$$\mathcal{L} : (x, (\lambda, \nu)) \mapsto \begin{cases} f(x) + \langle \lambda, Mx - b \rangle + \langle \nu, g(x) \rangle - \iota_{[0, +\infty)^m}(\nu) & \text{if } x \in \text{dom } f \\ +\infty & \text{otherwise.} \end{cases} \quad (3.3.3)$$

*Proof.* The first point is immediate. We prove the second. Recall that  $\mathcal{L}(x, (\lambda, \nu)) = \inf_{u, v} F(x, (u, v)) + \langle \lambda, u \rangle + \langle \nu, v \rangle$ . By the definition of  $F$  in Eq. (3.3.1),  $\mathcal{L}(x, (\lambda, \nu)) = +\infty$  when  $x \notin \text{dom } f$ . If  $x \in \text{dom } f$ ,

$$\begin{aligned} \mathcal{L}(x, (\lambda, \nu)) &= f(x) + \inf_{u \in \mathcal{Y}} (\iota_{\{b\}}(Mx - u) + \langle \lambda, u \rangle) + \inf_{v \in \mathbb{R}^m} (\iota_{(-\infty, 0]^m}(g(x) - v) + \langle \nu, v \rangle) \\ &= f(x) + \langle \lambda, Mx - b \rangle + \langle \nu, g(x) \rangle \\ &\quad + \inf_{u \in \mathcal{Y}} (\iota_{\{0\}}(Mx - b - u) - \langle \lambda, Mx - b - u \rangle) \\ &\quad + \inf_{v \in \mathbb{R}^m} (\iota_{(-\infty, 0]^m}(g(x) - v) - \langle \nu, g(x) - v \rangle) \\ &= f(x) + \langle \lambda, Mx - b \rangle + \langle \nu, g(x) \rangle + \inf_{w \in \mathbb{R}^m} (\iota_{(-\infty, 0]^m}(w) - \langle \nu, w \rangle). \end{aligned}$$

where we use the change of variables  $w = g(x) - v$  in the last equation. The proof of point (ii) is completed upon noting that  $-\inf_{w \in \mathbb{R}^m} (\iota_{(-\infty, 0]^m}(w) - \langle \nu, w \rangle) = \iota_{(-\infty, 0]^m}^*(\nu) = \iota_{[0, +\infty)^m}(\nu)$  (see Exercise 2.2).  $\square$

We denote by  $g_1, \dots, g_m : \mathcal{X} \rightarrow \mathbb{R}$  the components of  $g$  i.e.,  $g : x \mapsto (g_1(x), \dots, g_m(x))$ .

**Theorem 3.17.** *Consider a proper function  $f : \mathcal{X} \rightarrow (-\infty, +\infty]$ , a linear map  $M : \mathcal{X} \rightarrow \mathcal{Y}$ , a point  $b \in \mathcal{Y}$  and a map  $g : \mathcal{X} \rightarrow \mathbb{R}^m$ . Let  $(x, \lambda, \nu)$  be a point in  $\mathcal{X} \times \mathcal{Y} \times \mathbb{R}^m$ . Consider the following statements.*

i)  $(x, \lambda, \nu)$  satisfies

$$\begin{cases} x \text{ is feasible} & (3.3.4a) \\ \nu \geq 0 & (3.3.4b) \\ 0 \in \partial f(x) + M^* \lambda + \sum_{i=1}^m \nu_i \partial g_i(x) & (3.3.4c) \\ \forall i \in \{1, \dots, m\}, \nu_i g_i(x) = 0. & (3.3.4d) \end{cases}$$

ii)  $(x, (\lambda, \nu))$  is a saddle point of the Lagrangian.

iii)  $x$  is a primal solution,  $(\lambda, \nu)$  is a dual solution, and  $p = d$ .

Then, i)  $\Rightarrow$  ii)  $\Leftrightarrow$  iii). When  $f, g_1, \dots, g_m$  are convex, the three points are equivalent.

*Proof.* Note that Assumption 3.1 holds when  $F$  is defined as in (3.3.1). By Th. 3.4, ii)  $\Leftrightarrow$  iii). We now prove that i)  $\Rightarrow$  ii). As  $f$  is proper, a point  $(x, (\lambda, \nu))$  can be a saddle point of  $\mathcal{L}$  only if  $x \in \text{dom } f$  and  $\nu \geq 0$ . Therefore,  $(x, (\lambda, \nu))$  is a saddle point of  $\mathcal{L}$  iff the following four conditions hold together:

$$\begin{aligned} & x \in \text{dom } f \\ & \nu \geq 0 \\ & x \in \arg \min_{x' \in \mathcal{X}} f(x') + \langle \lambda, Mx' \rangle + \langle \nu, g(x') \rangle \\ & (\lambda, \nu) \in \arg \max_{(\lambda', \nu') \in \mathcal{Y} \times \mathbb{R}^m} \langle \lambda, Mx - b \rangle + \langle \nu', g(x) \rangle - \iota_{[0, +\infty)^m}(\nu'). \end{aligned} \quad (3.3.5)$$

Using Exercise 1.2, Eq. (3.3.5) is equivalent to

$$\begin{cases} \lambda \in \arg \min_{\lambda' \in \mathcal{Y}} \langle \lambda', -Mx + b \rangle \\ \nu_i \in \arg \min_{\nu' \in \mathbb{R}} \iota_{[0, +\infty)}(\nu') - \nu' g_i(x) \quad (\forall i \in \{1, \dots, m\}), \end{cases}$$

By Fermat's rule (Prop. 1.22) and the fact Exercise 1.5, the above condition is equivalent to

$$\begin{cases} 0 = -Mx + b \\ 0 \in N_{[0, +\infty)}(\nu_i) - g_i(x) \quad (\forall i \in \{1, \dots, m\}), \end{cases}$$

where we recall that  $N_{[0, +\infty)}(a)$  is the normal cone of  $[0, +\infty)$ , equal to  $\{0\}$  if  $a > 0$ , to  $(-\infty, 0]$  if  $a = 0$ , and to the empty set if  $a < 0$ . Hence, for a given  $i \in \{1, \dots, m\}$ , condition  $g_i(x) \in N_{[0, +\infty)}(\nu_i)$  reduces to  $\nu_i \geq 0$ ,  $g_i(x) \leq 0$  and  $\nu_i g_i(x) = 0$ . Putting all pieces together,  $(x, (\lambda, \nu))$  is a saddle point of  $\mathcal{L}$  iff

$$\begin{aligned} & x \in \text{dom } f, Mx = b, g(x) \leq 0 \\ & \nu \geq 0 \\ & 0 \in \partial (f + \langle \lambda, M \cdot \rangle + \langle \nu, g(\cdot) \rangle)(x) \\ & \nu_i g_i(x) = 0 \quad (\forall i \in \{1, \dots, m\}). \end{aligned} \quad (3.3.6)$$

where we applied Fermat's rule (Prop. 1.22) to obtain the condition (3.3.6). Using Cor. 3.14, it holds that

$$\partial f(x) + M^* \lambda + \sum_{i=1}^m \nu_i \partial g_i(x) \subset \partial(f + \langle \lambda, M \cdot \rangle + \langle \nu, g(\cdot) \rangle)(x). \quad (3.3.7)$$

Thus, the condition  $0 \in \partial f(x) + M^* \lambda + \sum_{i=1}^m \nu_i \partial g_i(x)$  implies Eq. (3.3.6). Thus,  $i) \Rightarrow ii)$ . As  $f$  is proper and  $\text{dom } g_i = \mathcal{X}$  for all  $i \in \{1, \dots, m\}$ , the inclusion (3.3.7) holds with equality if  $f, g_1, \dots, g_m$  are convex (use again Cor. 3.14). In the latter case,  $i) \Leftrightarrow ii)$ .  $\square$

The set of conditions in Eq. (3.3.4) are often referred to as the **Karush-Kuhn-Tucker (KKT) conditions**. Eq. (3.3.4a) is often referred to as the primal feasibility condition. Eq. (3.3.4b) is often referred to as the dual feasibility condition, because it is seen from the expression of the Lagrangian in (3.3.3) that the dual function  $-\inf \mathcal{L}(\mathcal{X}, \cdot)$  is greater than  $-\infty$  only if  $\nu \geq 0$ . Finally, the  $m$  conditions in Eq. (3.3.4d) are called the **complementary slackness conditions**.

**Theorem 3.18.** *Consider a proper function  $f : \mathcal{X} \rightarrow (-\infty, +\infty]$ , a linear map  $M : \mathcal{X} \rightarrow \mathcal{Y}$ , a point  $b \in \mathcal{Y}$  and a map  $g : \mathcal{X} \rightarrow \mathbb{R}^m$ . Assume that  $f, g_1, \dots, g_m$  are convex and that*

$$\exists x_0 \in \text{ri}(\text{dom } f), Mx_0 = b \text{ and } g(x_0) < 0. \quad (3.3.8)$$

*Then, a point  $x \in \mathcal{X}$  is primal optimal iff there exists  $(\lambda, \nu) \in \mathcal{Y} \times \mathbb{R}^m$  s.t. the conditions (3.3.4) hold (in which case, any such  $(\lambda, \nu)$  is dual optimal).*

*Proof.* The result is an application of Cor. 3.5 along with the characterization of the saddle points of  $\mathcal{L}$  by Th. 3.17. It is thus sufficient to check the hypotheses of Cor. 3.5 when  $F$  is defined as in (3.3.1). Assumption 3.1 clearly holds. The value function  $\vartheta$  is defined on  $\mathcal{Y} \times \mathbb{R}^m \rightarrow [-\infty, +\infty]$  by  $\vartheta : (u, v) \mapsto f(x) + \iota_{\{b\}}(Mx - u) + \iota_{(-\infty, 0]^m}(g(x) - v)$ . As  $f, g_1, \dots, g_m$  are convex, so is  $\vartheta$ . We must check that  $0 \in \text{ri dom } \vartheta$  and the proof is completed. A point  $(u, v)$  is in  $\text{dom } \vartheta$  iff there exists  $x \in \text{dom } f$  s.t.  $u = Mx - b$  and  $v \geq g(x)$ . Define  $C := \{(x, v) : x \in \text{dom } f, v \in \mathbb{R}^m, v \geq g(x)\}$ . It holds that

$$\text{dom } \vartheta = \mathcal{M}C - (b, 0),$$

where  $\mathcal{M} : \mathcal{X} \times \mathbb{R}^m \rightarrow \mathcal{Y} \times \mathbb{R}^m$  is the linear operator defined by  $\mathcal{M} : (x, v) \mapsto (Mx, v)$ . Note that  $C$  is convex. By Th. 1.7,

$$\text{ri dom } \vartheta = \mathcal{M} \text{ri}(C) - (b, 0).$$

We shall prove below that  $(x_0, 0) \in \text{ri}(C)$ . As  $\mathcal{M} \cdot (x_0, 0) - (b, 0) = (0, 0)$ , this implies that  $(0, 0) \in \text{ri dom } \vartheta$ , and concludes the proof.

We now show that  $(x_0, 0) \in \text{ri}(C)$ . As  $x_0 \in \text{ri dom } f$ , there exists  $\varepsilon > 0$  s.t.  $B_{\mathcal{X}}(x_0, \varepsilon) \cap \text{aff}(\text{dom } f) \subset \text{dom } f$  where  $B_{\mathcal{X}}(x_0, \varepsilon)$  represents the open Euclidean ball of center  $x_0$  and radius  $\varepsilon$ . Moreover, for every  $i \in \{1, \dots, m\}$ ,  $g_i(x_0) < 0$ . Note that  $g_i$  is continuous by Prop. 1.10. Thus, there exists  $\varepsilon' > 0$  s.t. for every  $i$ ,  $\sup\{g_i(x) : x \in B_{\mathcal{X}}(x_0, \varepsilon')\} < g_i(x_0)/2$ . Define

$$V := B_{\mathcal{X}}(x_0, \varepsilon \wedge \varepsilon') \times \prod_{i=1}^m (g_i(x_0)/2, -g_i(x_0)/2).$$

The set  $V$  is a neighborhood of  $(x_0, 0)$ . We prove that  $V \cap \text{aff}(C) \subset C$ , which implies that  $(x_0, 0) \in \text{ri}(C)$ . Note that  $C$  is a subset of the affine space  $\text{aff}(\text{dom } f) \times \mathbb{R}^m$ . Therefore,  $\text{aff}(C) \subset \text{aff}(\text{dom } f) \times \mathbb{R}^m$ . Thus,

$$V \cap \text{aff}(C) \subset (B_{\mathcal{X}}(x_0, \varepsilon \wedge \varepsilon') \cap \text{aff}(\text{dom } f)) \times \prod_{i=1}^m (g_i(x_0)/2, -g_i(x_0)/2)$$

Let  $(x, v) \in V \cap \text{aff}(C)$ . It holds that  $x \in B_{\mathcal{X}}(x_0, \varepsilon) \cap \text{aff}(\text{dom } f)$ , thus  $x \in \text{dom } f$ . Moreover, for every  $i \in \{1, \dots, m\}$ ,  $v_i > g_i(x_0)/2$ . Thus,

$$\nu_i > \sup\{g_i(y) : y \in B_{\mathcal{X}}(x_0, \varepsilon')\},$$

which implies that  $v_i > g_i(x)$ , because  $x \in B_{\mathcal{X}}(x_0, \varepsilon')$ . We have shown that  $v \geq g(x)$  and that  $x \in \text{dom } f$ . Stated otherwise, we have shown that  $V \cap \text{aff}(C) \subset C$ . As a conclusion,  $(x_0, 0) \in \text{ri}(C)$ .  $\square$

Eq. (3.3.8) is referred to as **Slater's condition**.

### Inequality constraints only

In this paragraph, we re-state the results of Section 3.3.1 in the special case where only inequality constraints are present. This is obtained by setting  $\mathcal{Y} = \{0\}$ ,  $M = 0$  and  $b = 0$ . The primal problem (3.3.2) reduces to the following problem:

**Minimization under inequality constraints.**

$$\text{Minimize } f(x) \text{ w.r.t. } x \in \mathcal{X} \text{ s.t. } g(x) \leq 0. \quad (3.3.9)$$

The primal optimal points are the solutions to Problem (3.3.9). The Lagrangian  $\mathcal{L} : \mathcal{X} \times \mathbb{R}^m \rightarrow [-\infty, +\infty]$  is

$$\mathcal{L} : (x, \nu) \mapsto \begin{cases} f(x) + \langle \nu, g(x) \rangle - \iota_{[0, +\infty)^m}(\nu) & \text{if } x \in \text{dom } f \\ +\infty & \text{otherwise.} \end{cases}$$

A point  $(x, \nu)$  is a saddle point of  $\mathcal{L}$  iff it satisfies the KKT conditions (3.3.4) only replacing the third condition (3.3.4c) by the simplest condition

$$0 \in \partial f(x) + \sum_{i=1}^m \nu_i \partial g_i(x), \quad (3.3.10)$$

and the statement of Th. 3.17 holds under this substitution. Although Th. 3.18 continue to hold, its statement can be simplified in this case. We re-state it as follows.

**Theorem 3.19.** *Consider a proper function  $f : \mathcal{X} \rightarrow (-\infty, +\infty]$  and a map  $g : \mathcal{X} \rightarrow \mathbb{R}^m$ . Assume that  $f, g_1, \dots, g_m$  are convex and that*

$$\exists x_0 \in \text{dom } f, g(x_0) < 0. \quad (3.3.11)$$

*Then, a point  $x \in \mathcal{X}$  is primal optimal iff there exists  $\nu \in \mathbb{R}^m$  s.t. the conditions (3.3.4a), (3.3.4b), (3.3.4d) and (3.3.10) hold (in which case, any such  $\nu$  is dual optimal).*

*Proof.* As in the proof of Th. 3.18, the main point is to prove that Eq. (3.3.11) implies that  $0 \in \text{ri dom } \vartheta$  where  $\vartheta$  is the value function. When  $\mathcal{Y} = \{0\}$ ,  $M = 0$  and  $b = 0$ , the latter reduces to the function  $\vartheta : \mathbb{R}^m \rightarrow [-\infty, +\infty]$  s.t.  $\vartheta : v \mapsto f(x) + \iota_{(-\infty, 0]^m}(g(x) - v)$  and its domain is  $\text{dom } \vartheta = \{v : \exists x \in \text{dom } f, v \geq g(x)\}$ . By the continuity of  $g$ , Condition (3.3.11) implies that  $0 \in \text{int dom } \vartheta$ . The proof is concluded upon noting that  $\text{int dom } \vartheta \subset \text{ri dom } \vartheta$ .  $\square$

## Examples, Exercises and Problems

In addition to the following exercises, a large number of feasible and instructive exercises can be found in [Boyd and Vandenberghe \(2009\)](#), chapter 5, pp 273-287.

**Exercise 3.1.** The goal of this exercise is to study the following semi-definite program where the variables are semi-definite matrices (the set of semi-definite matrices is denoted  $S_+^n$ )

$$\begin{aligned} & \inf_{x \in \mathbb{R}^{3 \times 3}} X_{3,3} + \iota_{S_+^3}(X) \\ & \text{under the constraints: } X_{1,2} + X_{2,1} + X_{3,3} = 1 \\ & \quad X_{2,2} = 0 \end{aligned}$$

We will denote  $E_{3,3} = \begin{bmatrix} 0 & 0 & 0 \\ 0 & 0 & 0 \\ 0 & 0 & 1 \end{bmatrix}$ ,  $f(x) = \langle E_{3,3}, X \rangle + \iota_{S_+^3}(X)$ ,  $G(X) = [X_{1,2} + X_{2,1} + X_{3,3}; X_{2,2}]$ ,  $e_1 = [1; 0]$  and  $A(X) = G(X) - e_1$ .

1. Show that if  $X \in S_+^3$  and  $X_{2,2} = 0$ , then  $X_{2,1} = X_{2,3} = 0$
2. Give an optimal solution to this problem.  
*Hint*: determine the set of feasible points.
3. Compute the dual of this problem and solve it. What do you observe?
4. What does  $0 \in A(\text{dom } f)$  mean for the feasibility of the optimization problem?
5. Show that  $[0; \epsilon] \in A(\text{dom } f)$  if and only if  $\epsilon \geq 0$ . Deduce that the constraints are not qualified.

**Exercise 3.2** (Examples of duals, [Borwein and Lewis \(2006\)](#), chap.4). Compute the dual of the following problems. In other words, calculate the dual function  $\mathcal{D}$  and write the problem of maximizing the latter as a convex minimization problem.

1. Linear program

$$\begin{aligned} & \inf_{x \in \mathbb{R}^n} \langle c, x \rangle \\ & \text{under constraint } Gx \preceq b \end{aligned}$$

where  $c \in \mathbb{R}^n$ ,  $b \in \mathbb{R}^p$  and  $G \in \mathbb{R}^{p \times n}$ .

*Hint* : you should find that the dual problem is again a linear program, with equality constraints.

2. Linear program on the non negative orthant

$$\begin{aligned} & \inf_{x \in \mathbb{R}^n} \langle c, x \rangle + \iota_{x \succeq 0} \\ & \text{under constraint } Gx \preceq b \end{aligned}$$

*Hint* : you should obtain a linear program with inequality constraints again.

3. Quadratic program

$$\begin{aligned} & \inf_{x \in \mathbb{R}^n} \frac{1}{2} \langle x, Cx \rangle \\ & \text{under constraint } Gx \preceq b \end{aligned}$$

where  $C$  is symmetric, positive, definite.

*Hint* : you should obtain an unconstrained quadratic problem.



- Assume in addition that the constraints are linearly independent, *i.e.*  $\text{rang}(G) = p$ , *i.e.*  $G = \begin{pmatrix} w_1^\top \\ \vdots \\ w_p^\top \end{pmatrix}$ , where  $(w_1, \dots, w_p)$  are linearly independent. Compute then the dual value.

**Exercise 3.3** (dual gap). Consider the three examples in exercise 3.2, and assume, as in example 3., that the constraints are linearly independent. Show the duality gap is zero under the respective following conditions:

1. Show that there is zero duality gap in examples 1 and 3 (linear and quadratic programs).

*Hint* : Slater.

2. For example 2, Assume that  $\exists \bar{x} > 0 : G\bar{x} = b$ . Show again that the duality gap is zero.

*Hint* (spoiler) : Show that  $0 \in \text{int dom } \mathcal{V}$ . In other words, show that for all  $y \in \mathbb{R}^p$  close enough to 0, there is some small  $\bar{u} \in \mathbb{R}^n$ , such that  $x = \hat{x} + \bar{u}$  is admissible, and  $Gx \leq b + y$ . To do so, exhibit some  $u \in \mathbb{R}^n$  such that  $Gu = -\mathbf{1}_p$  (why does it exist ?) Pick  $t$  such that  $\hat{x} + tu > 0$ . Finally, consider the ‘threshold’  $Y = -t\mathbf{1}_p \prec 0$  and show that, if  $y \succ Y$ , then  $\mathcal{V}(y) < \infty$ . Conclude.

**Exercise 3.4** (Gaussian Channel, Water filling.). In signal processing, a *Gaussian channel* refers to a transmitter-receiver framework with Gaussian noise: the transmitter sends an information  $X$  (real valued), the receiver observes  $Y = X + \epsilon$ , where  $\epsilon$  is a noise.

A Channel is defined by the joint distribution of  $(X, Y)$ . If it is Gaussian, the channel is called *Gaussian*. In other words, if  $X$  and  $\epsilon$  are Gaussian, we have a Gaussian channel.

Say the transmitter wants to send a word of size  $p$  to the receiver. He does so by encoding each possible word  $w$  of size  $p$  by a certain vector of size  $n$ ,  $\mathbf{x}_n^w = (x_1^w, \dots, x_n^w)$ . To stick with the Gaussian channel setting, we assume that the  $x_i^w$ 's are chosen as *i.i.d.* replicates of a Gaussian, centered random variable, with variance  $x$ .

The receiver knows the code (the dictionary of all  $2^p$  possible  $\mathbf{x}_n^w$ 's) and he observes  $\mathbf{y}_n = \mathbf{x}_n^w + \boldsymbol{\varepsilon}$ , where  $\boldsymbol{\varepsilon} \sim \mathcal{N}(\mathbf{0}, \sigma^2 \mathbf{I}_n)$ . We want to recover  $w$ .

The *capacity* of the channel, in information theory, is (roughly speaking) the maximum ratio  $C = n/p$ , such that it is possible (when  $n$  and  $p$  tend to  $\infty$  while  $n/p \equiv C$ ), to recover a word  $w$  of size  $p$  using a code  $\mathbf{x}_n^w$  of length  $n$ .

For a Gaussian Channel,  $C = \log(1 + x/\sigma^2)$ . ( $x/\sigma^2$  is the ratio signal/noise). For  $n$  Gaussian channels in parallel, with  $\alpha_i = 1/\sigma_i^2$ , then

$$C = \sum_{i=1}^n \log(1 + \alpha_i x_i).$$

The variance  $x_i$  represents a *power* affected to channel  $i$ . The aim of the transmitter is to maximize  $C$  under a *total power constraint* :  $\sum_{i=1}^n x_i \leq P$ . In other words, the problem is

$$\max_{x \in \mathbb{R}^n} \sum_{i=1}^n \log(1 + \alpha_i x_i) \quad \text{under constraints : } \forall i, x_i \geq 0, \quad \sum_{i=1}^n x_i \leq P. \quad (3.4.1)$$

1. Write problem (3.4.1) as a minimization problem under constraint  $g(x) \leq 0$ . Show that this is a convex problem (objective and constraints both convex).

2. Show that the constraints are qualified. (hint: Slater).
3. Write the Lagrangian function
4. Using the KKT theorem, show that a primal optimal  $x^*$  exists and satisfies :
  - $\exists K > 0$  such that  $x_i = \max(0, K - 1/\alpha_i)$ .
  - $K$  is given by

$$\sum_{i=1}^n \max(K - 1/\alpha_i, 0) = P$$

5. Justify the expression *water filling*

**Exercise 3.5** (Max-entropy). Let  $p = (p_1, \dots, p_n)$ ,  $p_i > 0$ ,  $\sum_i p_i = 1$  a probability distribution over a finite set. If  $x = (x_1, \dots, x_n)$  is another probability distribution ( $x_i \geq 0$ ), and if we use the convention  $0 \log 0 = 0$ , the entropy of  $x$  with respect to  $p$  is

$$H_p(x) = - \sum_{i=1}^n x_i \log \frac{x_i}{p_i}.$$

To deal with the case  $x_i < 0$ , introduce the function  $\psi : \mathbb{R} \rightarrow (-\infty, \infty]$ :

$$\psi(u) = \begin{cases} u \log(u) & \text{if } u > 0 \\ 0 & \text{if } u = 0 \\ +\infty & \text{otherwise .} \end{cases}$$

If  $g : \mathbb{R}^n \rightarrow \mathbb{R}^p$ , the general formulation of the max-entropy problem under constraint  $g(x) \preceq 0$  is

$$\begin{aligned} & \text{maximize over } \mathbb{R}^n && \sum_i (-\psi(x_i) + x_i \log(p_i)) \\ & \text{under constraints} && \sum x_i = 1; g(x) \preceq 0. \end{aligned}$$

In terms of minimization, the problem writes

$$\inf_{x \in \mathbb{R}^n} \sum_{i=1}^n \psi(x_i) - \langle x, c \rangle + \iota_{\langle \mathbf{1}_n, x \rangle = 1} + \iota_{g(x) \preceq 0}. \quad (3.4.2)$$

with  $c = \log(p) = (\log(p_1), \dots, \log(p_n))$  and  $\mathbf{1}_n = (1, \dots, 1)$  (the vector of size  $n$  which coordinates are equal to 1).

### A : preliminary questions

1. Show that

$$\partial \iota_{\langle \mathbf{1}_n, x \rangle} = \begin{cases} \{ \lambda_0 \mathbf{1}_n : \lambda_0 \in \mathbb{R} \} := \mathbb{R} \mathbf{1}_n & \text{if } \sum_i x_i = 1 \\ \emptyset & \text{otherwise.} \end{cases}$$

2. Show that  $\psi$  is convex

*hint* : compute first the Fenchel conjugate of the function  $\exp$ , then use proposition ??.

Compute  $\partial \psi(u)$  for  $u \in \mathbb{R}$ .

3. Show that

$$\partial\left(\sum_i \psi(x_i)\right) = \begin{cases} \sum_i (\log(x_i) + 1) \mathbf{e}_i & \text{if } x \succ 0 \\ \emptyset & \text{otherwise,} \end{cases}$$

where  $(\mathbf{e}_1, \dots, \mathbf{e}_n)$  is the canonical basis of  $\mathbb{R}^n$ .

4. Check that, for any set  $A$ ,  $A + \emptyset = \emptyset$ .

5. Consider the unconstrained optimization problem, (3.4.2) where the term  $\iota_{g(x) \leq 0}$  has been removed. Show that there exists a unique primal optimal solution, which is  $x^* = p$ .

*Hint:* Do not use Lagrange duality, apply Fermat's rule (section 1.5) instead. Then, check that the conditions for subdifferential calculus rules (proposition ??) apply.

**B : Linear inequality constraints** In the sequel, we assume that the constraints are linear, independent, and independent from  $\mathbf{1}_n$ , i.e.:  $g(x) = Gx - b$ , where  $b \in \mathbb{R}^p$ , and  $G$  is a  $p \times n$  matrix,

$$G = \begin{pmatrix} (\mathbf{w}^1)^\top \\ \vdots \\ (\mathbf{w}^p)^\top \end{pmatrix},$$

where  $\mathbf{w}^j \in \mathbb{R}^n$ , and the vectors  $(\mathbf{w}^1, \dots, \mathbf{w}^p, \mathbf{1}_n)$  are linearly independent. We also assume the existence of some point  $\hat{x} \in \mathbb{R}^n$ , such that

$$\forall i, \hat{x}_i > 0, \sum_i \hat{x}_i = 1, G\hat{x} = b. \quad (3.4.3)$$

1. Show that the constraints are qualified, in the Lagrangian sense.

*Hint (spoiler) :* proceed as in exercise 3.3, (2). This time, you need to introduce a vector  $u \in \mathbb{R}^n$ , such that  $Gu = -\mathbf{1}_p$  and  $\sum u_i = 0$  (again, why does it exist ?). The remaining of the argument is similar to that of exercise 3.3, (2).

2. Using the KKT conditions, show that any primal optimal point  $x^*$  must satisfy :  
 $\exists Z > 0, \exists \lambda \in \mathbb{R}^{+p} :$

$$x_i^* = \frac{1}{Z} p_i \exp \left[ - \sum_{j=1}^p \lambda_j \mathbf{w}_i^j \right] \quad (i \in \{1, \dots, n\})$$

(this is a Gibbs-type distribution).

## Chapter 4

# Fixed Points Algorithms

We wish to obtain numerically a minimizer of a convex function  $f : \mathcal{X} \rightarrow (-\infty, +\infty]$ . By Fermat's rule, this amounts to finding a point  $x$  such that  $0 \in \partial f(x)$ . In the case where  $f$  is differentiable, this is equivalent to finding a point  $x$  such that  $\nabla f(x) = 0$ . In this case, a celebrated algorithm is the gradient algorithm, which consists in generating a sequence  $(x^k : k = 0, 1, \dots)$  that is defined recursively by the equation

$$x^{k+1} = x^k - \gamma \nabla f(x^k),$$

where  $\gamma$  is a positive step size. Formally, we can write this algorithm as  $x^{k+1} = T(x^k)$ , where  $T(x) = x - \gamma \nabla f(x)$ . The minimizers of  $f$  coincide with the fixed points of  $T$ .

More generally, many optimization algorithms take the form  $x^{k+1} = T(x^k)$ , where the mapping  $T$  is chosen in such a way that its fixed points are solutions of the problem. We thus need to devise conditions on  $T$  that guarantee the convergence of the algorithm towards a fixed point.

### $\alpha$ -averaged operators

In all this chapter, an **operator** is a mapping from a nonempty subset  $D$  of the Euclidean space  $\mathcal{X}$  to  $\mathcal{X}$ . The image of  $x$  by an operator  $T$  will be denoted arbitrarily as  $T(x)$  or  $Tx$ . The composition  $T \circ R$  of the operator  $T$  with the operator  $R$  will be denoted compactly as  $TR$  when defined. The set of fixed points of  $T$  will be denoted  $\text{Fix}(T)$ . The identity operator on  $D$  is denoted as  $I$ .

**Definition 4.1.** Given a real number  $L \geq 0$ , the operator  $R$  defined on  $D$  is said  $L$ -Lipschitz if for all  $x, y \in D$ ,  $\|Rx - Ry\| \leq L\|x - y\|$ .

If  $L < 1$ , we say that  $R$  is a **contraction**. If  $L = 1$ ,  $R$  is said **non-expansive**.

**Remark 4.1.** Banach's fixed point theorem states that a contraction  $R$  defined on  $\mathcal{X}$  has a unique fixed point, and that each sequence that is defined by the recursion  $x^{k+1} = R(x^k)$  converges to this fixed point as  $k \rightarrow \infty$ . However, the contraction is a strong assumption, and Banach's fixed point theorem is often unapplicable in the field of optimization. Regarding the non expansiveness, this assumption is not sufficient, as the counter-example  $R = -I$  shows.

**Definition 4.2.** Let  $\alpha \in (0, 1]$ . The operator  $T$  on  $D$  is said  $\alpha$ -**averaged** if there exists a non-expansive operator  $R$  on  $D$  such that  $T = \alpha R + (1 - \alpha)I$ .

A  $1/2$ -averaged operator is said **firmly non-expansive**.

**Proposition 4.1.** Let  $\alpha \in (0, 1]$ . The following statements are equivalent.

1.  $T$  is  $\alpha$ -averaged.

2. For all  $x, y \in D$ ,  $\|Tx - Ty\|^2 \leq \|x - y\|^2 - \frac{1-\alpha}{\alpha} \|(I - T)x - (I - T)y\|^2$ .

*Proof.* Write  $T = \alpha R + (1 - \alpha)I$  where  $R$  is non-expansive. Equivalently,  $R = I - \frac{1}{\alpha}(I - T)$ . Writing  $\lambda = \frac{1}{\alpha}$  and  $Q = I - T$ , we get  $R = I - \lambda Q$ . We now develop:

$$\|Rx - Ry\|^2 = \lambda^2 \|Qx - Qy\|^2 + \|x - y\|^2 - 2\lambda \langle Qx - Qy, x - y \rangle.$$

Since  $R$  is non-expansive,  $\|x - y\|^2 \geq \|Rx - Ry\|^2$ , thus,

$$\begin{aligned} 0 &\geq \lambda \|Qx - Qy\|^2 - 2 \langle Qx - Qy, x - y \rangle \\ &= \lambda \|Qx - Qy\|^2 - 2 \|x - y\|^2 + 2 \langle Tx - Ty, x - y \rangle. \end{aligned}$$

We also have

$$\|Qx - Qy\|^2 = \|x - y\|^2 + \|Tx - Ty\|^2 - 2 \langle Tx - Ty, x - y \rangle.$$

By substituting the scalar product in the previous expression, we get

$$0 \geq (\lambda - 1) \|Qx - Qy\|^2 - \|x - y\|^2 + \|Tx - Ty\|^2$$

which is the required inequality.  $\square$

**Proposition 4.2.** *An operator  $T$  is firmly non-expansive iff  $\langle Tx - Ty, x - y \rangle \geq \|Tx - Ty\|^2$  for all  $x, y \in D$ .*

*Proof.* Write  $\|(I - T)x - (I - T)y\|^2 = \|x - y - (Tx - Ty)\|^2 = \|x - y\|^2 + \|Tx - Ty\|^2 - 2 \langle Tx - Ty, x - y \rangle$  and use Proposition 4.1-2. with  $\alpha = 1/2$ .  $\square$

**Theorem 4.3** (Krasnosel'skii Mann). *Let  $D \subset \mathcal{X}$  be a closed set. Let  $0 < \alpha < 1$ , and let  $T : D \rightarrow D$  be an  $\alpha$ -averaged operator such that  $\text{Fix}(T) \neq \emptyset$ . Then, each sequence  $(x^k)$  verifying the recursion  $x^{k+1} = T(x^k)$  with  $x^0 \in D$  converges to a point of  $\text{Fix}(T)$ .*

*Proof.* Let  $x^* \in \text{Fix}(T)$ . Since  $T$  is  $\alpha$ -averaged, Proposition 4.1-2. implies that for each  $k$ ,

$$\begin{aligned} \|x^{k+1} - x^*\|^2 &= \|Tx^k - Tx^*\|^2 \\ &\leq \|x^k - x^*\|^2 - \frac{1-\alpha}{\alpha} \|(I - T)x^k\|^2, \end{aligned} \tag{4.1.1}$$

where we used the fact that  $(I - T)x^* = 0$ . By iterating this inequality, we get that

$$0 \leq \|x^{k+1} - x^*\|^2 \leq \|x^0 - x^*\|^2 - \frac{1-\alpha}{\alpha} \sum_{i=0}^k \|(I - T)x^i\|^2$$

which implies that the series  $\sum_i \|(I - T)x^i\|^2$  is convergent, thus,  $(I - T)x^k \rightarrow 0$  as  $k \rightarrow \infty$ . Since the sequence of iterates  $(x^k)$  belongs to  $D$ , each accumulation point  $\bar{x}$  of this sequence belongs to  $D$ , and since  $T$  is  $\alpha$ -averaged, it is continuous, thus,  $(I - T)\bar{x} = 0$ , in other words,  $\bar{x} \in \text{Fix}(T)$ .

The inequality (4.1.1) implies moreover that the sequence  $(\|x^k - x^*\|^2)$  is decreasing. In particular, the sequence  $(x^k)$  is bounded. It admits an accumulation point that is an element of  $\text{Fix}(T)$  as we just showed. As a consequence, the inequality (4.1.1) remains true after replacing the fixed point  $x^*$  (which was arbitrarily chosen) with  $\bar{x}$ . Therefore, the sequence  $(\|x^k - \bar{x}\|^2)$  is decreasing and has zero as an accumulation point. Thus, it converges to zero. We showed that  $x^k \rightarrow \bar{x}$ , where  $\bar{x} \in \text{Fix}(T)$ .  $\square$

**Lemma 4.4** (Composition). *Let  $T$  and  $S$  be two operators from  $D$  to  $D$  such that  $T$  is  $\alpha$ -averaged and  $S$  is  $\beta$ -averaged, where  $0 < \alpha, \beta < 1$ . Then, there exists  $0 < \delta < 1$  such that the composed operator  $TS$  is  $\delta$ -averaged.*

*Proof.* For all  $x, y \in D$ ,

$$\begin{aligned}\|(I - TS)x - (I - TS)y\|^2 &= \|(I - S)x - (I - S)y + Sx - Sy - (TSx - TSy)\|^2 \\ &= \|(I - S)x - (I - S)y + (I - T)Sx - (I - T)Sy\|^2 \\ &\leq 2(\|(I - S)x - (I - S)y\|^2 + \|(I - T)Sx - (I - T)Sy\|^2)\end{aligned}$$

where we used the well known inequality  $\|a + b\|^2 \leq 2(\|a\|^2 + \|b\|^2)$ . By Proposition 4.1,  $\|(I - S)x - (I - S)y\|^2 \leq \frac{\beta}{1-\beta}(\|x - y\|^2 - \|Sx - Sy\|^2)$ , and similarly for  $T$ . Setting  $\kappa = \max(\beta/(1 - \beta), \alpha/(1 - \alpha))$ , we have

$$\|(I - TS)x - (I - TS)y\|^2 \leq 2\kappa(\|x - y\|^2 - \|Sx - Sy\|^2 + \|Sx - Sy\|^2 - \|TSx - TSy\|^2),$$

and finally,  $\|(I - TS)x - (I - TS)y\|^2 \leq 2\kappa(\|x - y\|^2 - \|TSx - TSy\|^2)$ . Setting  $\delta = (1 + (2\kappa)^{-1})^{-1}$ , we get that  $2\kappa = \delta/(1 - \delta)$ , and  $TS$  is  $\delta$ -averaged.  $\square$

## The gradient algorithm

Let  $f : \mathcal{X} \rightarrow (-\infty, +\infty]$ . We make the following assumption:

*Assumption 4.1.*  $f$  is convex and differentiable on  $\mathcal{X}$ , and  $\nabla f$  is  $L$ -Lipschitz.

Such functions are called smooth functions in the field of optimization theory.

**Theorem 4.5** (Baillon-Haddad). *Under Assumption 4.1,  $\forall x, y \in \mathcal{X}$ ,*

$$\langle \nabla f(x) - \nabla f(y), x - y \rangle \geq \frac{1}{L} \|\nabla f(x) - \nabla f(y)\|^2.$$

*In other words,  $L^{-1}\nabla f$  is firmly non-expansive.*

*Proof.* We first establish the following inequality. For all  $x, y$ ,

$$f(y) \leq f(x) + \langle \nabla f(x), y - x \rangle + \frac{L}{2} \|y - x\|^2. \quad (4.2.1)$$

To that end, we define the function  $\varphi(t) = f(x + t(y - x))$  on  $\mathbb{R}$ , and we observe that  $f(x) = \varphi(0)$  and  $f(y) = \varphi(1)$ . The function  $\varphi$  is differentiable with  $\varphi'(t) = \langle \nabla f(x + t(y - x)), y - x \rangle$ . Consequently,  $f(y) = f(x) + \int_0^1 \langle \nabla f(x + t(y - x)), y - x \rangle dt$ . Thus,  $f(y) = f(x) + \langle \nabla f(x), y - x \rangle + \delta$ , where  $\delta = \int_0^1 \langle \nabla f(x + t(y - x)) - \nabla f(x), y - x \rangle dt$ . Using the fact that  $\nabla f$  is  $L$ -Lipschitz, Inequality (4.2.1) is straightforward.

Second, we show that

$$f(y) \geq f(x) + \langle \nabla f(x), y - x \rangle + \frac{1}{2L} \|\nabla f(y) - \nabla f(x)\|^2. \quad (4.2.2)$$

For this, we fix  $x$  and we set  $\psi(y) = f(y) - \langle \nabla f(x), y - x \rangle$ . We easily verify that  $\psi$  is convex, that its derivative  $\nabla \psi$  is  $L$ -Lipschitz, and that  $\nabla \psi(x) = 0$ , in other words,  $x$  is a minimizer for  $\psi$ . In particular,  $f(x) = \psi(x) \leq \psi(y - \frac{1}{L}\nabla \psi(y))$ . We now use the inequality (4.2.1) at the points  $y - \frac{1}{L}\nabla \psi(y)$  and  $y$  after replacing  $f$  with  $\psi$ . We get

$$f(x) \leq \psi(y) - \frac{1}{L} \langle \nabla \psi(y), \nabla \psi(y) \rangle + \frac{1}{2L} \|\nabla \psi(y)\|^2.$$

Inequality (4.2.2) is established by noticing that  $\nabla \psi(y) = \nabla f(y) - \nabla f(x)$ . The proof of the theorem is completed by adding (4.2.2) to the inequality that we obtain by exchanging  $x$  and  $y$  in (4.2.2). The last statement of the theorem is a consequence of Proposition 4.2.  $\square$

**Lemma 4.6.** Let  $0 < \gamma < 2/L$ . Under Assumption 4.1,  $I - \gamma \nabla f$  is  $\gamma L/2$ -averaged.

*Proof.* By the previous theorem, there exists a non-expansive operator  $R$  such that  $L^{-1} \nabla f = (I + R)/2$ . Thus,  $I - \gamma \nabla f = (1 - \alpha)I + \alpha(-R)$  where  $\alpha = \gamma L/2$ , and the proof is completed by noticing that  $-R$  is non-expansive.  $\square$

**Theorem 4.7** (gradient algorithm). Let Assumption 4.1 hold true, and assume in addition that  $\arg \min f \neq \emptyset$ . Let  $0 < \gamma < 2/L$ . Then, each sequence  $(x^k)$  satisfying the recursion  $x^{k+1} = x^k - \gamma \nabla f(x^k)$  converges to a minimizer of  $f$ .

*Proof.* The assumption  $\arg \min f \neq \emptyset$  ensures that  $\text{Fix}(I - \gamma \nabla f) \neq \emptyset$ . The result is an immediate consequence of Theorem 4.3 and Lemma 4.6.  $\square$

## The proximal point and the proximal gradient algorithms

### The proximity operator

The proximity operator associated with a function  $f$  is the  $\mathcal{X} \rightarrow \mathcal{X}$  mapping defined as

$$\text{prox}_f(x) = \arg \min_{y \in \mathcal{X}} f(y) + \frac{1}{2} \|y - x\|^2 \quad (4.3.1)$$

provided this definition makes sense.

**Proposition 4.8.** Let  $f \in \Gamma_0(\mathcal{X})$ . Then,

1.  $\text{prox}_f$  is a well-defined mapping on  $\mathcal{X}$ .
2.  $\text{prox}_f$  is firmly non-expansive.
3.  $p = \text{prox}_f(x) \Leftrightarrow x \in p + \partial f(p)$ .

*Proof.* Let us fix  $x \in \mathcal{X}$ . The convexity is preserved if we subtract  $\frac{1}{2} \|\cdot\|^2$  from the function  $f + \frac{1}{2} \|\cdot - x\|^2$ . Thus,  $f + \frac{1}{2} \|\cdot - x\|^2$  is a strongly convex  $\Gamma_0(\mathcal{X})$  function. By Proposition 1.26, its argmin is a singleton, which establishes the first point. Let  $p$  be the minimizer of this function. By Fermat's rule,  $0 \in \partial(f + \frac{1}{2} \|\cdot - x\|^2)(p)$ , thus,  $0 \in \partial f(p) + p - x$  by Proposition 3.13. Consequently,  $x - p \in \partial f(p)$ , and the last point is established. Let  $y \in \mathcal{X}$  and  $q = \text{prox}_f(y)$ . Then,  $y - q \in \partial f(q)$ . This implies that  $\langle (x - p) - (y - q), p - q \rangle \geq 0$  (see Exercise 1.1), or equivalently, that  $\langle x - y, \text{prox}_f(x) - \text{prox}_f(y) \rangle \geq \|\text{prox}_f(x) - \text{prox}_f(y)\|^2$ . By Proposition 4.2, the operator  $\text{prox}_f$  is firmly non-expansive.  $\square$

The proof of the following proposition is an easy exercise.

**Proposition 4.9.** Let  $n \in \mathbb{N}_*$ , and let  $f_1, \dots, f_n$  be functions in  $\Gamma_0(\mathcal{X})$ . For all  $x = (x_1, \dots, x_n) \in \mathcal{X}^N$ , set  $f(x) = f_1(x_1) + \dots + f_n(x_n)$ . Then,  $f \in \Gamma_0(\mathcal{X}^N)$ , and  $\text{prox}_f(x) = (\text{prox}_{f_1}(x_1), \dots, \text{prox}_{f_n}(x_n))$ .

### The proximal point algorithm

Our purpose is to find a minimizer of a function  $g \in \Gamma_0(\mathcal{X})$ . By Fermat's rule (Proposition 1.22), this amounts to finding a point  $\bar{x}$  such that  $0 \in \partial g(\bar{x})$ . By Proposition 4.8-3., this inclusion reads  $\bar{x} \in \text{Fix prox}_g$ . The set of minimizers is not altered if we replace  $g$  with  $\gamma g$  for any  $\gamma > 0$ . The **proximal point algorithm** reads:

$$x^{k+1} = \text{prox}_{\gamma g}(x^k). \quad (4.3.2)$$

The following theorem follows from Proposition 4.8 and Theorem 4.3.

**Theorem 4.10** (proximal point algorithm). *Let  $g \in \Gamma_0(\mathcal{X})$  be such that  $\arg \min g \neq \emptyset$ . Then, given  $\gamma > 0$ , every sequence  $(x^k)$  verifying the recursion (4.3.2) converges towards a minimizer of  $g$ .*

Contrary to the gradient algorithm, the proximal point algorithm does not require that the function under study be smooth, and does not put any constraint on the step size  $\gamma > 0$ . On the other hand, the proximal point algorithm requires solving Problem (4.3.1) at each iteration, which can be computationally demanding.

## The proximal gradient algorithm

In this paragraph, we consider two functions  $f$  and  $g$  such that  $f$  satisfies Assumption 4.1, and  $g \in \Gamma_0(\mathcal{X})$ . Our purpose is to find a minimizer of  $f + g$ . Observing that  $\text{dom } f = \mathcal{X}$  and using Proposition 3.13, this amounts to finding a point  $\bar{x}$  such that  $0 \in \nabla f(\bar{x}) + \partial g(\bar{x})$ . Equivalently,

$$\bar{x} - \nabla f(\bar{x}) \in \bar{x} + \partial g(\bar{x}).$$

By Proposition 4.8, this inclusion can be read as  $\bar{x} = \text{prox}_g(\bar{x} - \nabla f(\bar{x}))$ . We can extend this remark by observing that there is an identity between the minimizers of  $f + g$  and those of  $\gamma f + \gamma g$  for all  $\gamma > 0$ . In other words, we established the following property:

**Proposition 4.11.** *Let  $f, g$  two functions such that  $f$  satisfies Assumption 4.1, and such that  $g \in \Gamma_0(\mathcal{X})$ . Assume that  $\arg \min(f + g) \neq \emptyset$ . Then,  $\bar{x} \in \arg \min(f + g)$  iff  $\bar{x} = \text{prox}_{\gamma g}(\bar{x} - \gamma \nabla f(\bar{x}))$ .*

This proposition suggests the following algorithm, termed the **proximal gradient algorithm**:

$$x^{k+1} = \text{prox}_{\gamma g}(x^k - \gamma \nabla f(x^k)). \quad (4.3.3)$$

**Theorem 4.12** (Proximal gradient algorithm). *Let  $f, g$  two functions such that  $f$  satisfies Assumption 4.1, and such that  $g \in \Gamma_0(\mathcal{X})$ . Assume that  $\arg \min(f + g) \neq \emptyset$ . Let  $0 < \gamma < 2/L$ . Then, every sequence  $(x^k)$  verifying the recursion (4.3.3) converges towards a minimizer of  $f + g$ .*

*Proof.* The operator  $I - \gamma \nabla f$  is  $\gamma L/2$ -averaged by Lemma 4.6. The operator  $\text{prox}_{\gamma g}$  is firmly non-expansive, in other words  $(1/2)$ -averaged, by Proposition 4.8. Thus, the composition  $\text{prox}_{\gamma g}(I - \gamma \nabla f)$  is  $\delta$ -averaged for some  $\delta \in (0, 1)$  by Lemma 4.4. Theorem 4.3 allows to conclude.  $\square$

The proximal gradient algorithm is often used when the computation of  $\text{prox}_{\gamma g}$  is easy. Examples are given here.

## Applications

### Projected gradient algorithm

Let  $C \subset \mathcal{X}$  be a non empty closed convex set. We want to solve the problem

$$\inf_{x \in C} f(x), \quad (4.4.1)$$

where the function  $f$  satisfies Assumption 4.1. We define the *indicator function*  $\iota_C$  of the set  $C$  as

$$\iota_C(x) = \begin{cases} 0 & \text{if } x \in C \\ +\infty & \text{otherwise.} \end{cases}$$



Problem (4.4.1) is equivalent to the problem

$$\inf_{x \in \mathcal{X}} f(x) + \iota_C(x).$$

One can immediately verify that  $\text{prox}_{\iota_C} = P_C$ . Consequently, the proximal gradient algorithm is

$$x^{k+1} = P_C(x^k - \gamma \nabla f(x^k)).$$

Under the assumptions of Theorem 4.12, this algorithm converges to a minimizer of  $f + \iota_C$ , i.e., a minimizer of  $f$  on  $C$ .

### Iterative soft-thresholding

Setting  $\mathcal{X} = \mathbb{R}^n$ , we want to solve the problem

$$\inf_{x \in \mathcal{X}} f(x) + \eta \|x\|_1 \tag{4.4.2}$$

where  $f$  satisfies Assumption 4.1, and where  $\|x\|_1$  is the  $\ell_1$  of the vector  $x$ , defined as  $\|x\|_1 = |x_1| + \dots + |x_n|$  by writing  $x = (x_1, \dots, x_n)$ .

**Proposition 4.13.** *The function  $\text{prox}_{\eta \|\cdot\|_1}$  coincides with the so-called soft thresholding function, which is defined for all  $x \in \mathbb{R}$  as:*

$$S_\eta(x) = \begin{cases} x - \eta & \text{if } x > \eta \\ 0 & \text{if } x \in [-\eta, \eta] \\ x + \eta & \text{if } x < -\eta. \end{cases}$$

*Proof.* Put  $p = \text{prox}_{\eta \|\cdot\|_1}(x)$ . By Proposition 4.8,  $x \in p + \partial_{\eta \|\cdot\|_1}(p)$ . In case  $p > 0$ , this implies after Example 1.1 that  $x = p + \eta$ , or equivalently  $p = x - \eta > 0$ . In case  $p < 0$ , we have  $p = x + \eta < 0$ . Finally, if  $p = 0$ ,  $x \in [-\eta, \eta]$ . Thus,  $p = S_\eta(x)$ .  $\square$

From Proposition 4.9, we get that for all  $x = (x_1, \dots, x_n)$ ,

$$\text{prox}_{\eta \|\cdot\|_1}(x) = (S_\eta(x_1), \dots, S_\eta(x_n)).$$

In this situation, the proximal gradient algorithm takes the form

$$\begin{aligned} y^k &= x^k - \gamma \nabla f(x^k) \\ x_i^k &= S_{\gamma\eta}(y_i^k) \quad (\forall i = 1, \dots, n). \end{aligned}$$

Under the assumptions of Theorem 4.12, the iterates  $x^k$  converge to a minimizer of (4.4.2).

## Chapter 5

# An introduction to monotone operator theory

Monotone operators can be seen as generalizations of the subdifferentials. They provide a theoretical framework for designing and studying many optimization algorithms. In the previous chapter, the convergence of some optimization algorithms was given a geometrical interpretation. This interpretation extends to the framework of monotone operators, providing powerful and elegant convergence proofs for a wider panel of algorithms.

### Basic definitions and facts

#### Monotone and maximal monotone operators

In this chapter, an **operator** on the Euclidean space  $\mathcal{X}$  is a mapping from  $\mathcal{X}$  to  $2^{\mathcal{X}}$  (sometimes we shall talk about a “set-valued operator”). The subdifferential of a function (see Definition 1.11) is a typical example of an operator. A single-valued operator, *i.e.*, a mapping  $A$  from  $D \subset \mathcal{X}$  to  $\mathcal{X}$  as in the previous chapter, is seen here as a particular case of a set-valued operator if we put  $Ax = \emptyset$  when  $x \notin D$ . The image of  $x$  by an operator  $A$  will be arbitrarily denoted as  $A(x)$  or  $Ax$ . The **domain** of  $A$  is the set  $\text{dom}(A) = \{x \in \mathcal{X} : Ax \neq \emptyset\}$ . The **image** of  $A$  is the set  $\text{Im}(A) = \bigcup_{x \in \mathcal{X}} Ax$ . The **graph** of  $A$  is the subset of  $\mathcal{X} \times \mathcal{X}$  defined as  $\text{gra}(A) = \{(x, y) \in \mathcal{X} \times \mathcal{X} : y \in Ax\}$ . It is obvious that an operator can be identified by its graph. As in Page 18, the inverse of  $A$ , denoted as  $A^{-1}$ , is the operator whose graph is  $\text{gra}(A^{-1}) = \{(x, y) \in \mathcal{X} \times \mathcal{X} : (y, x) \in \text{gra}(A)\}$ . Observe that  $\text{dom}(A^{-1}) = \text{Im}(A)$ . The **set of zeros** of an operator  $A$  is the set  $\mathcal{Z}(A) = \{x \in \text{dom}(A) : 0 \in Ax\}$ . The **composition**  $AB$  of the operator  $A$  with the operator  $B$  is the operator whose graph is  $\text{gra}(AB) = \{(x, z) \in \mathcal{X} \times \mathcal{X} : \exists y \in \mathcal{X}, y \in Bx, z \in Ay\}$ . Given a real number  $\gamma$ , we denote as  $\gamma A$  the operator  $(\gamma I)A$ , where we recall that  $I$  is the identity operator. Finally, the **sum**  $A + B$  of  $A$  and  $B$  is the operator  $x \mapsto (A + B)(x) = \{y + z : y \in Ax, z \in Bx\}$  when  $x \in \text{dom}(A) \cap \text{dom}(B)$ , and  $\emptyset$  elsewhere.

**Definition 5.1.** A set-valued operator  $A$  on the Euclidean space  $\mathcal{X}$  is said **monotone** if for all  $x, y \in \text{dom}(A)$  and for all  $u \in Ax$  and  $v \in Ay$ , it holds that  $\langle x - y, u - v \rangle \geq 0$ .

In Exercise 1.1, we showed that the subdifferential of a function  $f : \mathcal{X} \rightarrow [-\infty, \infty]$  is an example of a monotone operator.

**Exercise 5.1.** Show that if  $\mathcal{X} = \mathbb{R}$  and if  $A$  is single-valued, then  $A$  is monotone *iff*  $A$  is non-decreasing.

The **resolvent** of an operator  $A$  is the operator  $Q_A = (I + A)^{-1}$ . Resolvents of monotone operators will be of prime importance in the remainder.

**Proposition 5.1.** *An operator  $A$  is monotone iff for all  $\gamma > 0$ , the operator  $Q_{\gamma A}$  is a non-expansive operator from  $\text{Im}(I + \gamma A)$  to  $\mathcal{X}$ .*

This proposition shows in particular that when  $A$  is monotone,  $Q_{\gamma A}$  is single-valued. In other words, given  $u \in \mathcal{X}$ , the inclusion  $u \in x + \gamma Ax$  admits at most one solution.

*Proof.* Observe that  $x \in Q_{\gamma A}(u) \Leftrightarrow u = x + \gamma p$  with  $p \in Ax$ . Thus, the statement of the proposition can be rephrased as follows:  $A$  is monotone iff  $\forall \gamma > 0, \forall x, y \in \text{dom } A, \forall p \in Ax, \forall q \in Ay, \|x - y\|^2 \leq \|x + \gamma p - y - \gamma q\|^2$ .

We have

$$\|x + \gamma p - y - \gamma q\|^2 - \|x - y\|^2 = 2\gamma \langle x - y, p - q \rangle + \gamma^2 \|p - q\|^2.$$

If  $A$  is monotone, then the right hand side is obviously non negative. Conversely, if this term is non negative for all  $\gamma > 0$ , then dividing it by  $\gamma$  and making  $\gamma \downarrow 0$ , we get that  $\langle x - y, p - q \rangle \geq 0$ .  $\square$

**Proposition 5.2.** *The operator  $A$  is monotone iff for all  $\gamma > 0$ , the operator  $Q_{\gamma A}$  is a firmly non-expansive operator from  $\text{Im}(I + \gamma A)$  to  $\mathcal{X}$ .*

*Proof.* We only need to prove that if  $A$  is monotone, then  $Q_{\gamma A}$  is firmly non expansive. Given  $u, v \in \text{dom } Q_{\gamma A}$ , let  $x = Q_{\gamma A}(u)$  and  $y = Q_{\gamma A}(v)$ . Since  $u - x \in \gamma Ax$  and  $v - y \in \gamma Ay$ , we get that  $0 \leq \langle x - y, u - x - (v - y) \rangle = \langle x - y, u - v \rangle - \|x - y\|^2$  by the monotonicity of  $A$ . The result follows from Proposition 4.2.  $\square$

In this chapter, we shall study iterative algorithms for finding an element of  $\mathcal{Z}(A)$  when  $A$  is a monotone operator. One motivation for doing this study is provided by the monotonicity of the subdifferential (Exercise 1.1) and by Fermat's rule (Proposition 1.22).

To this end, we first observe that  $\mathcal{Z}(A) = \text{Fix } Q_{\gamma A}$  for any  $\gamma > 0$ . Second, the previous proposition shows that  $Q_{\gamma A}$  is a firmly non expansive operator. Consequently, if the domain of  $Q_{\gamma A}$  is the whole space  $\mathcal{X}$ , then the theorem of Krasnosel'skii Mann (Theorem 4.3) shows that for every initial value  $x_0$ , the sequence  $(x^k)$  produced by the algorithm  $x^{k+1} = Q_{\gamma A}(x^k)$  converges to a point of  $\mathcal{Z}(A)$ . We thus need to guarantee that  $\text{dom } Q_{\gamma A} = \mathcal{X}$ .

Noting that the graph inclusion provides a partial ordering for the monotone operators, we make the following definition:

**Definition 5.2.** A monotone operator  $A$  on  $\mathcal{X}$  is said **maximal** if its graph is a maximal element in the graph inclusion ordering.

Figure 5.1 is an illustration for a non-maximal and a maximal monotone operator on  $\mathbb{R}$ .

**Exercise 5.2.** Show that the maximal monotonicity of  $A$ ,  $A^{-1}$ , and  $\gamma A$  for  $\gamma > 0$  are equivalent.

**Exercise 5.3.** Let  $A$  be a single-valued monotone operator with  $\text{dom } A = \mathcal{X}$ . Show that if  $A$  is continuous, then it is maximal monotone.

*Hint:* Let  $x, u \in \mathcal{X}$  be such that  $\langle x - y, u - Ay \rangle \geq 0$  for all  $y \in \mathcal{X}$ . To show that  $u = Ax$ , work on the scalar product  $\langle u - Ax, u - Az_\lambda \rangle$ , where  $z_\lambda = x + \lambda(u - Ax)$ , and  $\lambda > 0$ .

**Exercise 5.4.** Show that a single-valued linear operator  $M$  is maximal monotone iff the symmetric operator  $M + M^*$  is nonnegative (in the sense that  $\langle x, (M + M^*)x \rangle \geq 0$  for all  $x \in \mathcal{X}$ ).

**Theorem 5.3 (Minty).** *A monotone operator  $A$  on  $\mathcal{X}$  is maximal iff  $\text{Im}(I + A) = \mathcal{X}$ .*

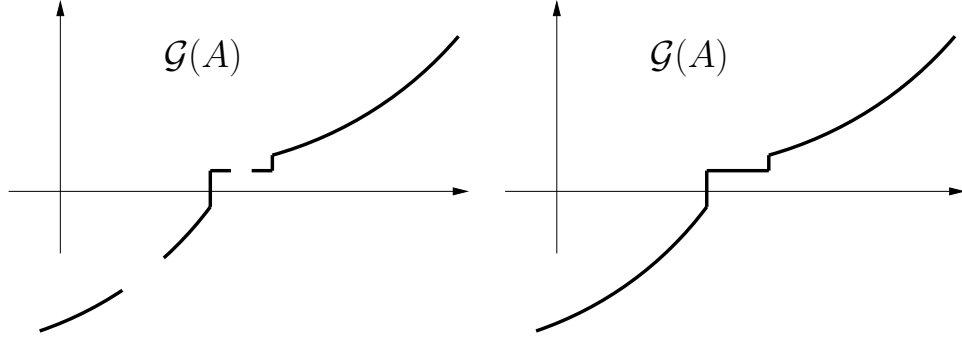


Figure 5.1: Non-maximal (left) *vs.* maximal (right) monotone operator on  $\mathbb{R}$ .

Equivalently,  $A$  is maximal *iff*  $\text{dom } Q_A = \mathcal{X}$ .

*Proof.* Assume that  $\text{Im}(I + A) = \mathcal{X}$ . Let  $B$  be a monotone operator such that  $\text{gra}(A) \subset \text{gra}(B)$ . Given  $(x, u) \in \text{gra}(B)$ , it holds by assumption that there exists  $y \in \text{dom } A$  such that  $x + u \in y + Ay$ , thus,  $x + u \in y + By$ . On the other hand,  $x + u \in x + Bx$ . Consequently,  $x = (I + B)^{-1}(x + u) = y$  by Proposition 5.1, and  $u \in Ax$ . Thus,  $(x, u) \in \text{gra}(A)$ , and  $A$  is maximal.

The converse will be proven in Exercise 5.7. □

Using Propositions 5.1 and 5.2, we immediately get the following corollary of Theorem 5.3.

**Corollary 5.4.** *Given a set-valued operator  $A$  on  $\mathcal{X}$ , the following assertions are equivalent:*

- i)  $A$  is a maximal monotone operator.
- ii) For all  $\gamma > 0$ ,  $Q_{\gamma A}$  is a non expansive operator with domain  $\mathcal{X}$ .
- iii) For all  $\gamma > 0$ ,  $Q_{\gamma A}$  is a firmly non expansive operator with domain  $\mathcal{X}$ .

**Corollary 5.5.** *If  $g \in \Gamma_0(\mathcal{X})$ , then  $\partial g$  is a maximal monotone operator, and  $Q_{\partial g} = \text{prox}_g$ .*

*Proof.* Combine Proposition 4.8 with the previous corollary. □

## The proximal point algorithm

Along with Theorem 4.3, Corollary 5.4 leads to the following result.

**Theorem 5.6** (proximal point algorithm). *Let  $A$  be a maximal monotone operator on  $\mathcal{X}$  such that  $\mathcal{Z}(A) \neq \emptyset$ . Then, for each  $\gamma > 0$  and each initial value  $x_0 \in \mathcal{X}$ , the sequence  $(x^k)$  produced by the algorithm*

$$x^{k+1} = Q_{\gamma A}(x^k)$$

*converges to a point of  $\mathcal{Z}(A)$ .*

Corollary 5.5 shows that Theorem 4.10 is a particular case of Theorem 5.6.

We provide herein another instance of Theorem 5.6: the so-called *method of multipliers*.

## The method of multipliers

**Lemma 5.7.** *Let  $\mathcal{X}$  and  $\mathcal{Y}$  be two Euclidean spaces, let  $f \in \Gamma_0(\mathcal{X})$ , and let  $M \in \mathbf{L}(\mathcal{X}, \mathcal{Y})$ .*

*i) If  $f$  is coercive, then  $M \triangleright f \in \Gamma_0(\mathcal{Y})$ . Moreover, for each  $y \in \mathcal{Y}$ , the set*

$$\arg \min_{u \in \mathcal{X}} f(u) + \frac{1}{2} \|Mu - y\|^2$$

*is nonempty, and  $\text{prox}_{M \triangleright f}(y) = Mx$ , where  $x$  is an arbitrary element of this set.*

*ii) If  $M$  is injective (which supposes  $\dim(\mathcal{X}) \leq \dim(\mathcal{Y})$ ), then  $M \triangleright f \in \Gamma_0(\mathcal{Y})$ . Moreover, for each  $y \in \mathcal{Y}$ , the set*

$$\arg \min_{u \in \mathcal{X}} f(u) + \frac{1}{2} \|Mu - y\|^2$$

*is reduced to a singleton  $\{x\}$ , and  $\text{prox}_{M \triangleright f}(y) = Mx$ .*

*Proof.* We first establish (i). We first recall from Proposition 1.16 that  $M \triangleright f$  is a convex function with the nonempty domain  $\text{dom}(M \triangleright f) = M \text{dom}(f)$ . Since  $f$  is coercive,  $\inf f(\mathcal{X}) > -\infty$  by Proposition 1.23. Thus,  $M \triangleright f$  is proper since  $\inf(M \triangleright f)(\mathcal{Y}) \geq \inf f(\mathcal{X})$ . Let us show that  $M \triangleright f$  is l.s.c. We first show that if  $y \in \text{dom}(M \triangleright f)$ , then  $(M \triangleright f)(y)$  is attained by some  $x$  such that  $y = Mx$ . Indeed,  $(M \triangleright f)(y) = \inf f + \iota_y(M \cdot)$ . The function  $f + \iota_y(M \cdot)$  is convex and proper, and by Proposition 1.21, it is l.s.c. Moreover, it is coercive, being larger than  $f$ . The result follows from Proposition 1.23. Consider now a sequence  $y_k \rightarrow y \in \mathcal{Y}$ , and let  $c_k = (M \triangleright f)(y_k)$ . To establish the lower semi continuity of  $M \triangleright f$ , we only need to consider the case where the sequence  $(c_k)$  is bounded. In this case, there exists  $x_k$  such that  $Mx_k = y_k$  and  $f(x_k) = c_k$ . Consider an arbitrary sequence of such  $x_k$ . By the boundedness of  $(c_k)$  and the coercivity of  $f$ , the sequence  $(x_k)$  belongs to a compact (see, e.g., the proof of Proposition 1.23), thus, it has accumulation points. Given an arbitrary subsequence  $(\psi(k))$  such that  $x_{\psi(k)} \rightarrow x^*$ , we have  $Mx^* = y$ , and by the lower semi continuity of  $f$ ,  $\liminf f(x_{\psi(k)}) \geq f(x^*) \geq (M \triangleright f)(y)$ . Consequently,  $M \triangleright f \in \Gamma_0(\mathcal{Y})$ .

The function

$$u \mapsto f(u) + \frac{1}{2} \|Mu - y\|^2 = f(u) + \frac{1}{2} \langle u, M^*Mu \rangle - \langle Mu, y \rangle + \frac{\|y\|^2}{2} \quad (5.1.1)$$

is clearly a coercive function in  $\Gamma_0(\mathcal{X})$ . Consequently, its argmin is nonempty by Proposition 1.23. The last result can be checked by direct calculation.

We now establish (ii) by succinctly pointing out the differences with (i). Consider  $y \in \text{dom}(M \triangleright f)$ . Since  $M$  is injective,  $M \triangleright f(y) = f(M^{-1}y) > -\infty$ , thus,  $M \triangleright f$  is proper. To show that  $M \triangleright f(y)$  is l.s.c., we similarly take a sequence  $y_k \rightarrow y \in \mathcal{Y}$ , and assume that the  $c_k = (M \triangleright f)(y_k)$  are bounded. Here we simply have  $c_k = f(x_k)$  where  $x_k = M^{-1}y_k$ . We now have  $x_k \rightarrow M^{-1}y$ , and the lower semicontinuity of  $M \triangleright f$  follows from the lower semicontinuity of  $f$ . The function (5.1.1) is now the sum of a function in  $\Gamma_0(\mathcal{X})$  and the strongly convex function  $\langle u, M^*Mu \rangle / 2$ . Thus, it is a strongly convex function in  $\Gamma_0(\mathcal{X})$ . Consequently, its argmin is reduced to a singleton by Proposition 1.26. The last result can also be checked by direct calculation.  $\square$

Given two Euclidean spaces  $\mathcal{X}$  and  $\mathcal{Y}$ , let  $f \in \Gamma_0(\mathcal{X})$ ,  $b \in \mathcal{Y}$ , and  $M \in \mathbf{L}(\mathcal{X}, \mathcal{Y})$ . We consider the minimization problem under affine equality constraints which is detailed in Section 3.2.2. In the framework of the Fenchel-Rockafellar duality theory, the primal problem is written

$$p = \inf_{\mathcal{X}} f + g \circ M,$$

where  $g = \iota_{\{b\}} \in \Gamma_0(\mathcal{Y})$ . Observing that  $g^* = \langle b, \cdot \rangle$ , the dual problem is

$$d = -\inf_{\mathcal{Y}} f^* \circ (-M^*) + \langle b, \cdot \rangle .$$

The set  $\mathcal{S}$  of saddle points of the Lagrangian, whether it exists, is the set of points  $(x, \phi) \in \mathcal{X} \times \mathcal{Y}$  that satisfy the system of inclusions (3.2.8). If the qualification condition (3.2.9) holds true, then by Theorem 3.7,  $p = d$ , and the set of minimizers of the dual problem is not empty. We now observe that the primal problem can be written as  $p = (M \triangleright f)(b)$ . If  $f$  is coercive, then  $M \triangleright f \in \Gamma_0(\mathcal{Y})$  by Lemma 5.7-(i). Obviously, the qualification condition (3.2.9) ensures that  $b \in \text{dom}(M \triangleright f)$ . Thus, the set of minimizers of the primal problem is not empty (*details to be provided*). This implies that  $\mathcal{S} \neq \emptyset$  by Proposition 3.8-(iii). The method of multipliers is a technique for finding a point of  $\mathcal{S}$ . Let  $\gamma > 0$ . Starting with an arbitrary  $\phi^0 \in \mathcal{Y}$ , the algorithm consists in following iterations:

$$x^{k+1} \in \arg \min_{v \in \mathcal{X}} \left\{ f(v) + \langle Mv, \phi^k \rangle + \frac{\|Mv - b\|^2}{2\gamma} \right\} , \quad (5.1.2a)$$

$$\phi^{k+1} = \phi^k + \gamma^{-1} (Mx^{k+1} - b) . \quad (5.1.2b)$$

**Theorem 5.8** (method of multipliers). *Given two Euclidean spaces  $\mathcal{X}$  and  $\mathcal{Y}$ , let  $f$  be a coercive function in  $\Gamma_0(\mathcal{X})$ , let  $b \in \mathcal{Y}$ , and let  $M \in \mathbf{L}(\mathcal{X}, \mathcal{Y})$ . Consider the minimization problem with affine equality constraints (3.2.7). Assume that the set  $\mathcal{S}$  of saddle points for this problem, i.e., the points  $(x, \phi) \in \mathcal{X} \times \mathcal{Y}$  verifying (3.2.8), is not empty. Let  $\gamma > 0$  be arbitrary. Consider the sequence of iterates  $((x^k, \phi^k))$  produced by the algorithm (5.1.2). Then the sequence  $(\phi^k)$  converges to  $\phi^* \in \mathcal{Y}$ , and the sequence  $(x^k)$  belongs to a compact set. Given any accumulation point  $x^*$  of  $(x^k)$ , it holds that  $(x^*, \phi^*) \in \mathcal{S}$ .*

To prove this theorem, we show that the method of multipliers results from applying the proximal point algorithm to the dual problem  $\inf_{\mathcal{Y}} f^* \circ (-M^*) + \langle b, \cdot \rangle$ .

*Proof.* We show that the iterates

$$\phi^{k+1} = \text{prox}_{\gamma^{-1}(f^* \circ (-M^*) + \langle b, \cdot \rangle)}(\phi^k) \quad (5.1.3)$$

are those which are provided by Algorithm (5.1.2). Consider the function  $F(y) = ((-M) \triangleright f)(y - b)$ . We know from Exercise 2.5 that  $((-M) \triangleright f)^* = f^* \circ (-M^*)$ . Thus,  $F^* = f^* \circ (-M^*) + \langle b, \cdot \rangle$ . Moreover,  $F \in \Gamma_0(\mathcal{Y})$  by Lemma 5.7-(i). Consequently, we can apply Moreau's identity (2.4.1) to obtain

$$\phi^{k+1} = \phi^k - \gamma^{-1} \text{prox}_{\gamma F}(\gamma \phi^k) . \quad (5.1.4)$$

Developing the prox operator at the right hand side, we get that

$$\begin{aligned} \text{prox}_{\gamma F}(\phi^k) &= \arg \min_{y \in \mathcal{Y}} ((-M) \triangleright f)(y - b) + \frac{1}{2\gamma} \|y - \gamma \phi^k\|^2 \\ &= b + \arg \min_{y \in \mathcal{Y}} ((-M) \triangleright f)(y) + \frac{1}{2\gamma} \|y + b - \gamma \phi^k\|^2 \\ &= b + \text{prox}_{\gamma((-M) \triangleright f)}(\gamma \phi^k - b) . \end{aligned}$$

Using Lemma 5.7-(i), we get that  $\text{prox}_{\gamma((-M) \triangleright f)}(\gamma \phi^k - b) = -Mx^{k+1}$ , where  $x^{k+1}$  satisfied the inclusion (5.1.2a), and  $\phi^{k+1}$  is given by (5.1.2b).

Since  $\mathcal{S} \neq \emptyset$ , the dual problem  $\inf f^* \circ (-M)^* + \langle b, \cdot \rangle$  has a minimizer, thus, the iterates  $\phi^k$  which are given by Equation (5.1.3) converge to a dual solution  $\phi^\infty$  by Theorem 4.10 (or Theorem 5.6). Equation (5.1.2a) can be rewritten as

$$x^{k+1} \in \arg \min_v \left\{ f(v) + \frac{\|Mv + \gamma\phi^k - b\|^2}{2\gamma} \right\}.$$

By the boundedness of  $(\phi^k)$  and the coercivity of  $f$ , we obtain that the sequence  $(x^k)$  belongs to a compact. Moreover, considering Equation (5.1.4) and taking  $k \rightarrow \infty$ , we get by the continuity of the prox that  $\text{prox}_{\gamma F}(\gamma\phi^\infty) = 0$ . This shows that  $Mx^k \rightarrow b$ , thus, any accumulation point  $x^\infty$  of  $(x^k)$  is a primal solution.  $\square$

## Splitting algorithms

Let  $A$  and  $B$  be two maximal monotone operators such that  $\mathcal{Z}(A+B) \neq \emptyset$ . A splitting algorithm is an iterative algorithm for finding an element of  $\mathcal{Z}(A+B)$  by performing operations that involve  $A$  and  $B$  separately. This requirement is usually unavoidable for obtaining implementable optimization algorithms. The remainder of this chapter is devoted to such algorithms.

### The Douglas-Rachford splitting algorithm

The **Cayley transform** of a maximal monotone operator  $A$  is the mapping

$$C_A : \mathcal{X} \longrightarrow \mathcal{X}, \quad x \longmapsto 2Q_A - I.$$

The proof of the following proposition is left to the reader as an exercise:

**Proposition 5.9.**  *$C_A$  is a non expansive operator defined on  $\mathcal{X}$ .*

Given two maximal monotone operators  $A$  and  $B$ , the **Douglas-Rachford** (or Lions-Mercier) operator is the single-valued operator defined on  $\mathcal{X}$  as

$$T_{A,B} = \frac{I + C_A C_B}{2}.$$

More explicitly,

$$T_{A,B}(x) = Q_A(2Q_B(x) - x) - Q_B(x) + x. \quad (5.2.1)$$

It is obvious from the definition of  $T_{A,B}$  and from Proposition 5.9 that  $T_{A,B}$  is firmly non-expansive.

**Proposition 5.10.** *Let  $A$  and  $B$  be two maximal monotone operators. Then,  $\mathcal{Z}(A+B) \neq \emptyset \Leftrightarrow \text{Fix } T_{A,B} \neq \emptyset$ . If  $\mathcal{Z}(A+B) \neq \emptyset$ , then  $Q_B(\text{Fix } T_{A,B}) = \mathcal{Z}(A+B)$ .*

*Proof.* We have the following equivalences:

$$\begin{aligned} x \in \text{Fix } T_{A,B}, \ y = Q_B(x) &\Leftrightarrow Q_A(2y - x) = y, \ y = Q_B(x) \\ &\Leftrightarrow 2y - x \in y + Ay, \ x \in y + By \\ &\Leftrightarrow 2y \in 2y + Ay + By, \ x \in y + By \\ &\Leftrightarrow y \in \mathcal{Z}(A+B), \ x \in y + By. \end{aligned}$$

$\square$

This proposition, the firm non-expansiveness of  $T_{A,B}$ , and Theorem 4.3 lead to the following result:

**Theorem 5.11** (Douglas-Rachford algorithm). *Let  $A$  and  $B$  be two maximal monotone operators on  $\mathcal{X}$  such that  $\mathcal{Z}(A+B) \neq \emptyset$ . Let  $\gamma > 0$  and  $x_0 \in \mathcal{X}$  be arbitrary, and consider the sequence  $(x^k)$  produced by the algorithm*

$$x^{k+1} = T_{\gamma A, \gamma B}(x^k).$$

*Then,  $\text{Fix } T_{\gamma A, \gamma B} \neq \emptyset$ , the sequence  $(x^k)$  converges to a point of  $\text{Fix } T_{\gamma A, \gamma B}$ , and the sequence  $(Q_{\gamma B}(x^k))$  converges to a point of  $\mathcal{Z}(A+B)$ .*

We shall describe two instances of this algorithm.

### The Douglas-Rachford for minimizing the sum of two functions

Given two functions  $f, g \in \Gamma_0(\mathcal{X})$ , consider the problem  $\inf_{\mathcal{X}} f + g$ . Assume that  $\mathcal{Z}(\partial f + \partial g) \neq \emptyset$ . This is satisfied if the set of minimizers of the problem  $\inf_{\mathcal{X}} f + g$  is non empty and if  $0 \in \text{ri}(\text{dom } f - \text{dom } g)$ , as shown by Fermat's rule and by Proposition 3.13. The Douglas-Rachford algorithm considered here finds a point of  $\mathcal{Z}(\partial f + \partial g)$ . Given a step  $\gamma > 0$ , it consists in the following iterations, starting with an arbitrary  $x^0 \in \mathcal{X}$ .

$$\begin{aligned} z^{k+1} &= \text{prox}_{\gamma g}(x^k), \\ x^{k+1} &= \text{prox}_{\gamma f}(2z^{k+1} - x^k) - z^{k+1} + x^k. \end{aligned} \tag{5.2.2}$$

The following theorem is a straightforward corollary of Theorem 5.11.

**Theorem 5.12** (Douglas-Rachford minimization algorithm). *Let  $f, g \in \Gamma_0(\mathcal{X})$ , and assume that  $\mathcal{Z}(\partial f + \partial g) \neq \emptyset$ . Let  $\gamma > 0$  be arbitrary. Then each sequence  $(z^k)$  produced by Algorithm (5.2.2) converges to a point of  $\mathcal{Z}(\partial f + \partial g)$ .*

### The Douglas-Rachford algorithm in the dual domain: ADMM

Given two Euclidean spaces  $\mathcal{X}$  and  $\mathcal{Y}$  such that  $\dim(\mathcal{X}) \leq \dim(\mathcal{Y})$ , let  $f \in \Gamma_0(\mathcal{X})$ ,  $g \in \Gamma_0(\mathcal{Y})$ , and let  $M$  be an injection in  $\mathbf{L}(\mathcal{X}, \mathcal{Y})$ . Consider the minimization problem

$$p = \inf_{\mathcal{X}} f + g \circ M.$$

Recalling the results of Section 3.2, the dual of this problem in the sense of the Fenchel-Rockafellar theory is

$$d = -\inf_{\mathcal{Y}} f^* \circ (-M^*) + g^*.$$

The set  $\mathcal{S}$  of saddle points of the Lagrangian, whether it exists, is the set of points  $(x, \phi) \in \mathcal{X} \times \mathcal{Y}$  that satisfy the system of inclusions (3.2.5). We recall that if the qualification condition  $0 \in \text{ri}(M \text{dom}(f) - \text{dom}(g))$  holds true, then by Theorem 3.7,  $p = d$ , and the set of minimizers of the dual problem is not empty. If, furthermore, the set of minimizers of the primal problem  $\inf f + g \circ M$  is not empty, then  $\mathcal{S} \neq \emptyset$  by Proposition 3.8-(iii). The Alternating Direction Method of Mutlipliers (ADMM, see e.g., Boyd et al. (2011)) is a popular algorithm for finding a saddle point for this problem. Let  $\gamma > 0$  be some step size. Starting with  $(z^0, \phi^0) \in \mathcal{Y}^2$ , this algorithm is written:



$$x^{k+1} = \arg \min_{v \in \mathcal{X}} f(v) + \left\langle \phi^k, Mv \right\rangle + \frac{1}{2\gamma} \|Mv - z^k\|^2, \quad (5.2.3a)$$

$$z^{k+1} = \arg \min_{w \in \mathcal{Y}} g(w) - \left\langle \phi^k, w \right\rangle + \frac{1}{2\gamma} \|Mx^{k+1} - w\|^2, \quad (5.2.3b)$$

$$\phi^{k+1} = \phi^k + \gamma^{-1} (Mx^{k+1} - z^{k+1}). \quad (5.2.3c)$$

**Theorem 5.13** (ADMM). *Given two Euclidean spaces  $\mathcal{X}$  and  $\mathcal{Y}$  such that  $\dim \mathcal{X} \leq \dim \mathcal{Y}$ , let  $f \in \Gamma_0(\mathcal{X})$ ,  $g \in \Gamma_0(\mathcal{Y})$ , and let  $M$  be an injective element of  $\mathbf{L}(\mathcal{X}, \mathcal{Y})$ . Consider the minimization problem  $\inf f + g \circ M$ . Assume that the set  $\mathcal{S}$  of saddle points for this problem, i.e., the points  $(x, \phi) \in \mathcal{X} \times \mathcal{Y}$  verifying the system of inclusions (3.2.5), is not empty. Let  $\gamma > 0$  be arbitrary. Then, the sequence of iterates  $((x^k, \phi^k))$  produced by the algorithm (5.2.3) converges to a point of  $\mathcal{S}$ .*

ADMM can be seen as an instance of the Douglas-Rachford algorithm applied to the dual problem  $\inf f^* \circ (-M^*) + g^*$ . We take this route to prove Theorem 5.13.

*Proof.* Define the maximal monotone operators  $A = \partial(f^* \circ (-M^*))$  and  $B = \partial g^*$ . Noticing that  $\text{Im}(-M^*) = \mathcal{X}$  and applying Proposition 3.13 to the subdifferential  $\partial(0 + f^* \circ (-M^*))$ , we get that  $A = -M^* \partial f^* \circ (-M^*)$ . Let  $(x, \phi) \in \mathcal{S}$ . By Proposition 2.10, Inclusion (3.2.5a) can be rewritten as  $x \in \partial f^*(-M^* \phi)$ . Plugging this into Inclusion (3.2.5b), we get that  $0 \in A\phi + B\phi$ , thus,  $\mathcal{Z}(A+B) \neq \emptyset$ . Now, let  $(u^k)$  be a sequence defined by the iterations  $u^{k+1} = T_{\gamma^{-1}A, \gamma^{-1}B}(u^k)$ . By Theorem 5.11,  $(u^k)$  converges to a point  $u^\infty$ , and the sequence  $(\phi^k)$ , defined as  $\phi^k = Q_{\gamma^{-1}B}(u^k)$ , converges to  $\phi^\infty \in \mathcal{Z}(A+B)$ . Our first task is to show that  $(\phi^k)$  is the one provided by Algorithm (5.2.3).

Since  $\phi^k = Q_{\gamma^{-1}B}(u^k)$ , we can write  $u^k = \phi^k + \gamma^{-1}z^k$ , where  $z^k \in \partial g^*(\phi^k)$ . Recalling Equation (5.2.1), let  $y^{k+1} = Q_{\gamma^{-1}A}(2\phi^k - u^k) = \text{prox}_{\gamma^{-1}f^* \circ (-M^*)}(\phi^k - \gamma^{-1}z^k)$ . Then the output of the Douglas-Rachford algorithm is  $u^{k+1} = y^{k+1} - \phi^k + u^k = y^{k+1} + \gamma^{-1}z^k$ . Let us make explicit the expression of  $y^{k+1}$ . We know from Exercise 2.5 that  $((-M) \triangleright f)^* = f^* \circ (-M^*)$ . Thus, using Moreau's identity (2.4.1), we get

$$y^{k+1} = \phi^k - \gamma^{-1}z^k - \gamma^{-1} \text{prox}_{\gamma(-M) \triangleright f}(\gamma\phi^k - z^k).$$

Using Lemma 5.7-(ii), we can write

$$\text{prox}_{\gamma(-M) \triangleright f}(\gamma\phi^k - z^k) = -Mx^{k+1}, \quad (5.2.4)$$

where

$$x^{k+1} = \arg \min_{v \in \mathcal{X}} f(v) + \frac{1}{2\gamma} \|Mv + \gamma\phi^k - z^k\|^2,$$

which is the right hand side of (5.2.3a). Replacing the value of  $y^{k+1}$  in the expression of  $u^{k+1}$ , we get  $u^{k+1} = \phi^k - \gamma^{-1}z^k + \gamma^{-1}Mx^{k+1} + \gamma^{-1}z^k = \phi^k + \gamma^{-1}Mx^{k+1}$ . We conclude by using that  $\phi^{k+1} = Q_{\gamma^{-1}B}(u^{k+1}) = \text{prox}_{\gamma^{-1}g^*}(\phi^k + \gamma^{-1}Mx^{k+1})$ . Indeed, by Moreau's identity again,

$$\phi^{k+1} = \phi^k + \gamma^{-1}Mx^{k+1} - \gamma^{-1}z^{k+1},$$

where

$$z^{k+1} = \text{prox}_{\gamma g}(\gamma\phi^k + Mx^{k+1}).$$

This equation and the previous one coincide with (5.2.3b) and (5.2.3c) respectively.

We now show that  $((x^k, \phi^k))$  converges to a point of  $\mathcal{S}$ . Since  $u^k \rightarrow u^\infty$  and  $\phi^k \rightarrow \phi^\infty$ ,  $z^k = \gamma(u^k - \phi^k)$  converges to a point  $z^\infty$ . From Equation (5.2.3c), and the injectivity of  $M$ , the sequence  $(x^k)$  converges to  $x^\infty$ , and  $z^\infty = Mx^\infty$ . Using that  $\phi^k = Q_{\gamma^{-1}\partial g^*}(\phi^k + \gamma^{-1}z^k)$  and taking  $k$  to  $\infty$ , we get that  $\phi^\infty = Q_{\gamma^{-1}\partial g^*}(\phi^\infty + \gamma^{-1}z^\infty)$ , which is equivalent to  $z^\infty \in \partial g^*(\phi^\infty)$ . Since  $z^\infty = Mx^\infty$ , we get Inclusion (3.2.5b) with  $(x, \phi) = (x^\infty, \phi^\infty)$ . Using (5.2.4) and taking  $k$  to  $\infty$ , we obtain that  $Q_{\gamma\partial((-M) \triangleright f)}(\gamma\phi^\infty + Mx^\infty) = -Mx^\infty$  which is equivalent to  $\phi^\infty \in \partial((-M) \triangleright f)(-Mx^\infty)$ , or  $-Mx^\infty \in -M\partial f^*(-M^*\phi^\infty)$ , which can be rewritten as Inclusion (3.2.5a).  $\square$

## The Forward-Backward algorithm

**Definition 5.3.** Given  $\alpha > 0$ , a single-valued operator  $A$  defined on  $\mathcal{X}$  is said  $\alpha$ -cocoercive if  $\alpha A$  is firmly non-expansive.

**Exercise 5.5.** Show that an  $\alpha$ -cocoercive operator is maximal monotone.

Consider a convex and differentiable function  $f$  defined on  $\mathcal{X}$ . If  $\nabla f$  is  $L$ -Lipschitz for some  $L > 0$ , then the Baillon-Haddad theorem (Theorem 4.5) shows that  $\nabla f$  is  $L^{-1}$ -cocoercive.

Given a maximal monotone operator  $A$ , an  $\alpha$ -cocoercive operator  $B$  and a real  $\gamma > 0$ , the **Forward-Backward** operator is the single-valued operator defined on  $\mathcal{X}$  as

$$S_{\gamma A, \gamma B} = Q_{\gamma A}(I - \gamma B).$$

The “forward” operator is  $I - \gamma B$  while the “backward” operator is  $Q_{\gamma A}$ . Applying the latter consists in solving an implicit equation, hence the denomination.

**Proposition 5.14.** Fix  $S_{\gamma A, \gamma B} = \mathcal{Z}(A + B)$ . Furthermore, if  $\gamma \in (0, 2\alpha)$ , then there exists  $\delta \in (0, 1)$  such that  $S_{\gamma A, \gamma B}$  is a  $\delta$ -averaged operator.

*Proof.* To obtain the first result, we write

$$x \in \mathcal{Z}(A + B) \Leftrightarrow 0 \in \gamma Ax + \gamma Bx \Leftrightarrow (I - \gamma B)x \in (I + \gamma A)x \Leftrightarrow x = (I + \gamma A)^{-1}(I - \gamma B)x.$$

Since  $B$  is  $\alpha$ -cocoercive, there exists a non-expansive operator  $R$  such that  $B = (2\alpha)^{-1}R + (2\alpha)^{-1}I$ . Thus, the operator

$$I - \gamma B = \frac{\gamma}{2\alpha}(-R) + \left(1 - \frac{\gamma}{2\alpha}\right)I$$

is  $\gamma/(2\alpha)$ -averaged. Since  $Q_{\gamma A}$  is  $1/2$ -averaged by Corollary 5.4, Lemma 4.4 leads to the second result.  $\square$

With the help of Theorem 4.3, this proposition immediately leads to:

**Theorem 5.15** (Forward-Backward algorithm). *Let  $A$  be a maximal monotone operator, and let  $B$  be a  $\alpha$ -cocoercive operator on  $\mathcal{X}$ . Assume that  $\mathcal{Z}(A + B) \neq \emptyset$ . Let  $\gamma \in (0, 2\alpha)$ . Then, for each  $x_0 \in \mathcal{X}$ , the sequence  $(x^k)$  produced by the algorithm*

$$x^{k+1} = S_{\gamma A, \gamma B}(x^k)$$

*converges to a point of  $\mathcal{Z}(A + B)$ .*

Observe that this theorem generalizes Theorem 4.12: the forward-backward algorithm is a general version of the proximal gradient algorithm.

In the next paragraph, we describe another instance of the forward-backward algorithm: the so-called Vũ-Condat algorithm.

### The Vũ-Condat primal-dual algorithm

Let  $\mathcal{X}$  and  $\mathcal{Y}$  be two Euclidean spaces. Let  $f$  be a convex and differentiable function defined on  $\mathcal{X}$ , and assume that  $\nabla f$  is  $L$ -Lipschitz. Let  $g \in \Gamma_0(\mathcal{X})$  and  $h \in \Gamma_0(\mathcal{Y})$ . Given  $M \in \mathbf{L}(\mathcal{X}, \mathcal{Y})$ , consider the minimization problem

$$\inf_{\mathcal{X}} f + g + h \circ M,$$

The dual problem for this problem is

$$\inf_{\mathcal{Y}} (f + g)^* \circ (-M^*) + h^*.$$

Recalling the inclusions (3.2.5), the set  $\mathcal{S}$  of saddle points for this problem, whether it exists, satisfies the following conditions:

$$0 \in \nabla f(x) + \partial g(x) + M^* \phi \quad (5.2.5a)$$

$$0 \in -Mx + \partial h^*(\phi) \quad (5.2.5b)$$

(notice indeed that  $\partial(f + g) = \nabla f + \partial g$ , since  $f$  and  $g$  satisfy the assumptions of Proposition 3.13). By Theorem 3.7 and Proposition 3.8, the condition

$$0 \in \text{ri}(M \text{ dom}(g) - \text{dom}(h)),$$

along with a non emptiness condition of the set of minimizers of the primal problem, ensures that  $\mathcal{S} \neq \emptyset$ .

Given two step sizes  $\tau, \gamma > 0$ , the Vũ-Condat algorithm is an algorithm for finding a point of  $\mathcal{S}$ . It consists in the following iterations:

$$z^{k+1} = \arg \min_{v \in \mathcal{Y}} h(v) - \langle v, \phi^k \rangle + \frac{1}{2\gamma} \|v - Mx^k\|^2, \quad (5.2.6a)$$

$$\phi^{k+1} = \phi^k + \gamma^{-1}(Mx^k - z^{k+1}), \quad (5.2.6b)$$

$$x^{k+1} = \arg \min_{w \in \mathcal{X}} g(w) + \langle w, \nabla f(x^k) \rangle + \langle w, M^*(2\phi^{k+1} - \phi^k) \rangle + \frac{1}{2\tau} \|w - x^k\|^2. \quad (5.2.6c)$$

**Theorem 5.16** (Vũ-Condat algorithm). *Let  $f$  be a convex and differentiable function defined on  $\mathcal{X}$ . Assume that  $\nabla f$  is  $L$ -Lipschitz for some  $L \geq 0$ . Let  $g \in \Gamma_0(\mathcal{X})$ ,  $h \in \Gamma_0(\mathcal{Y})$ , and  $M \in \mathbf{L}(\mathcal{X}, \mathcal{Y})$ . Consider the minimization problem  $\inf_{\mathcal{X}} f + g + h \circ M$ , and assume that the set  $\mathcal{S}$  of saddle points of this problem, i.e., the set of points  $(x, \phi) \in \mathcal{X} \times \mathcal{Y}$  verifying the inclusions (5.2.5), is not empty. Consider the iterations (5.2.6), and assume that the positive numbers  $\tau$  and  $\gamma$  are such that  $\tau^{-1} - \gamma^{-1}\|M\|^2 > L/2$ , where  $\|\cdot\|$  is the spectral norm. Then the sequence of the iterates  $(x^k, \phi^k)$  produced by the algorithm (5.2.6) converges to a point of  $\mathcal{S}$ .*

The remainder of this section is devoted to the proof of this theorem. We first observe that if  $L = 0$  (i.e.,  $f$  is affine), we can always replace  $L$  with  $\varepsilon > 0$  small enough so that  $\tau^{-1} - \gamma^{-1}\|M\|^2 > \varepsilon/2$ . Thus, we can always assume that  $L > 0$ , as we shall do hereinafter.

We first introduce some new notations. We denote as  $\mathcal{W}$  the Euclidean space  $\mathcal{X} \times \mathcal{Y}$  endowed with the scalar product  $\langle (x, \phi), (y, \psi) \rangle = \langle x, y \rangle + \langle \phi, \psi \rangle$ , where  $x, y \in \mathcal{X}$ ,  $\phi, \psi \in \mathcal{Y}$ , and the scalar products  $\langle x, y \rangle$  and  $\langle \phi, \psi \rangle$  are those of  $\mathcal{X}$  and  $\mathcal{Y}$  respectively (even though we use the

same notations for scalar products and norms in these different spaces, their use will be always clear from the context). In the remainder, when we denote an element  $u \in \mathcal{W}$  as  $u = (x, \phi)$ , we mean that  $x \in \mathcal{X}$  and  $\phi \in \mathcal{Y}$ . The following matrix notation for some operators on  $\mathcal{W}$  will also be convenient: given the set-valued operators  $C_{xx} : \mathcal{X} \rightarrow 2^{\mathcal{X}}$ ,  $C_{xy} : \mathcal{X} \rightarrow 2^{\mathcal{Y}}$ ,  $C_{yx} : \mathcal{Y} \rightarrow 2^{\mathcal{X}}$ , and  $C_{yy} : \mathcal{Y} \rightarrow 2^{\mathcal{Y}}$ , the operator

$$C = \begin{bmatrix} C_{xx} & C_{xy} \\ C_{yx} & C_{yy} \end{bmatrix}$$

is the  $\mathcal{W} \rightarrow 2^{\mathcal{W}}$  operator defined as  $C((x, \phi)) = \{(y, \psi) \in \mathcal{W} : y \in C_{xx}x + C_{xy}\phi, \psi \in C_{yx}x + C_{yy}\psi\}$ .

Now, define on  $\mathcal{W}$  the operators

$$A = \begin{bmatrix} \partial g & M^* \\ -M & \partial h^* \end{bmatrix} \quad \text{and} \quad B = \begin{bmatrix} \nabla f & 0 \\ 0 & 0 \end{bmatrix},$$

where 0 denotes generically an operator whose image is the set  $\{0\}$ . With these notations, the inclusion (5.2.5) is rewritten as  $0 \in Au + Bu$ . Let us inspect the basic properties of  $A$  and  $B$ .

**Exercise 5.6.** Given two maximal monotone operators  $C$  and  $D$  defined respectively on  $\mathcal{X}$  and  $\mathcal{Y}$ , show that  $\begin{bmatrix} C & 0 \\ 0 & D \end{bmatrix}$  is maximal monotone on  $\mathcal{W}$ .

This exercise shows that the operator  $\begin{bmatrix} \partial g & 0 \\ 0 & \partial h^* \end{bmatrix}$  is maximal monotone. We also know by Exercise 5.4 that the linear operator  $\begin{bmatrix} & M^* \\ -M & \end{bmatrix}$  is maximal monotone. In these conditions, it can be shown that the operator  $A$ , which is the sum of these two operators, is maximal monotone (see *e.g.*, (Brézis, 1973, Lemma 2.4). Another proof that is specific to our context will be given as an exercise).

Turning to the operator  $B$ , we notice that

$$B = \nabla F, \quad \text{where} \quad F : \mathcal{W} \rightarrow \mathbb{R}, \quad (x, \phi) \mapsto f(x). \quad (5.2.7)$$

By Baillon-Haddad's theorem,  $B$  is a  $L^{-1}$ -cocoercive operator.

We just showed that  $A$  and  $B$  satisfy the assumptions of Theorem 5.15. Therefore, one first idea for solving our minimization problem is to implement the Forward-Backward algorithm  $u^{k+1} = Q_{\rho A}(I - \rho B)u^k$  for some  $\rho > 0$ . However, due to the structure of  $A$ , implementing the operator  $Q_{\rho A}$  is difficult because the operations involving  $g$  and  $h^*$  would have to be performed jointly (check this). What is needed is an algorithm that performs operations on  $g$  and  $h$  *separately*.

The idea goes as follows. Given  $\tau > 0$  and  $\gamma > 0$ , define the linear operator

$$V = \begin{bmatrix} \tau^{-1}I & M^* \\ M & \gamma I \end{bmatrix}.$$

The inclusion  $0 \in Au + Bu$  can be rewritten as  $(V - B)u \in (V + A)u$ . Observe now that the operator

$$V + A = \begin{bmatrix} \partial g + \tau^{-1}I & 2M^* \\ 0 & \partial h^* + \gamma I \end{bmatrix}$$

has an “upper triangular” structure. Thanks to this structure, we shall obtain an implementable and convergent algorithm if some care is taken in the choice of  $V$ .

The inclusion  $(V - B)u \in (V + A)u$  may suggest the iterative algorithm

$$(V - B)u^k \in (V + A)u^{k+1}. \quad (5.2.8)$$

This inclusion can be written equivalently  $(I - V^{-1}B)u^k \in (I + V^{-1}A)u^{k+1}$ , provided  $V$  is invertible. Below, we shall show that the iteration

$$u^{k+1} = (I + V^{-1}A)^{-1}(I - V^{-1}B)u^k \quad (5.2.9)$$

is well defined for all  $u^k$ , and is equivalent to (5.2.8). Moreover, provided  $\tau^{-1} - \gamma^{-1}\|M\|^2 > L/2$ , the resulting algorithm is an instance of the Forward-Backward algorithm whose convergence is established by Theorem 5.15. The proof will be done in three steps:

Step 1: We show that when  $\gamma\tau^{-1} > \|M\|^2$ , the linear symmetric operator  $V$  is positive definite. This defines a new scalar product  $\langle u, w \rangle_V = \langle u, Vw \rangle$  in the space  $\mathcal{W}$ . We denote as  $\mathcal{W}_V$  the space  $\mathcal{W}$  endowed with this new scalar product, and we denote as  $\|\cdot\|_V$  the associated norm.

Step 2: In the space  $\mathcal{W}_V$ , the operator  $V^{-1}A$  is maximal monotone,  $V^{-1}B$  is cocoercive, and the operator defined by Equation (5.2.9) satisfies the assumptions of Theorem 5.15. Consequently, the sequence  $(u^k)$  is well defined, and it converges to a point of  $\mathcal{Z}(A + B)$ .

Step 3: Getting back to the inclusion (5.2.8), we make use of the upper triangular structure of  $V + A$  to obtain Algorithm (5.2.6). We show that the sequence of iterates  $((x^k, \phi^k))$  converges to a point of  $\mathcal{S}$ .

The first step is achieved by the following lemma:

**Lemma 5.17.** *The symmetric operator  $V$  is positive definite if  $\gamma\tau^{-1} > \|M\|^2$ .*

This condition is obviously satisfied when  $\tau^{-1} - \gamma^{-1}\|M\|^2 > L/2$  as in the statement of Theorem 5.16.

*Proof.* Given an arbitrary non zero vector  $u = (x, \phi) \in \mathcal{W}$ , we have

$$\begin{aligned} \langle u, Vu \rangle &= \tau^{-1}\|x\|^2 + \gamma\|\phi\|^2 + 2\langle \phi, Mx \rangle \\ &\geq \tau^{-1}\|x\|^2 + \gamma\|\phi\|^2 - 2\|M\|\|\phi\|\|x\| \\ &= (\tau^{-1} - \gamma^{-1}\|M\|^2)\|x\|^2 + (\gamma^{1/2}\|\phi\| - \gamma^{-1/2}\|M\|\|x\|)^2. \end{aligned} \quad (5.2.10)$$

Using the inequality  $\tau^{-1} - \gamma^{-1}\|M\|^2 > 0$ , and dealing separately with the cases  $x \neq 0$  and  $x = 0$ , we get that  $\langle u, Vu \rangle > 0$ , hence the result.  $\square$

We now turn to the second step.

**Lemma 5.18.** *The operator  $V^{-1}A$  is maximal monotone in  $\mathcal{W}_V$ .*

*Proof.* Given any  $u, v \in \text{dom}(A)$  and any  $p \in Au$ ,  $q \in Av$ , we have  $\langle u - v, V^{-1}p - V^{-1}q \rangle_V = \langle u - v, p - q \rangle \geq 0$ , hence the monotonicity of  $V^{-1}A$  in  $\mathcal{W}_V$ . The maximality of  $V^{-1}A$  is deduced from the maximality of  $A$ , as it can be checked by the reader as an exercise.  $\square$

**Lemma 5.19.** *Given any vector  $v = (x, 0) \in \mathcal{W}$ , we have*

$$\langle v, V^{-1}v \rangle \leq (\tau^{-1} - \gamma^{-1}\|M\|^2)^{-1}\|x\|^2.$$

*Moreover, for any  $w = (x, \phi) \in \mathcal{W}$ , we have*

$$\left(\tau^{-1} - \gamma^{-1}\|M\|^2\right)\|x\|^2 \leq \langle w, Vw \rangle.$$

*Proof.* It is well known that any invertible matrix  $C$  with a block decomposition  $C = \begin{bmatrix} C_{11} & C_{12} \\ C_{21} & C_{22} \end{bmatrix}$  where  $C_{11}$  is a square matrix satisfies  $[C^{-1}]_{11} = (C_{11} - C_{12}C_{22}^{-1}C_{21})^{-1}$ , where  $[C^{-1}]_{11}$  is the upper left block of  $C^{-1}$  having the same size as  $C_{11}$ . Applying this identity to the linear operator  $V$ , we get that  $\langle v, V^{-1}v \rangle = \langle x, (\tau^{-1} - \gamma^{-1}M^*M)^{-1}x \rangle$ . Since  $M^*M \leq \|M\|^2 I$  in the positive semidefinite ordering of the symmetric matrices, we have  $(\tau^{-1} - \gamma^{-1}M^*M)^{-1} \leq (\tau^{-1} - \gamma^{-1}\|M\|^2)^{-1}I$ , hence the first result.

The second result is deduced immediately from the inequality (5.2.10).  $\square$

**Lemma 5.20.** *The operator  $V^{-1}B$  is  $\kappa/L$ -cocoercive in  $\mathcal{W}_V$ , where  $\kappa = \tau^{-1} - \gamma^{-1}\|M\|^2$ .*

*Proof.* Given  $u = (x, \phi), v = (y, \psi) \in \mathcal{W}$ , we have

$$\begin{aligned} \|V^{-1}Bu - V^{-1}Bv\|_V^2 &= \langle Bu - Bv, V^{-1}(Bu - Bv) \rangle \\ &= \langle (\nabla f(x) - \nabla f(y), 0), V^{-1}(\nabla f(x) - \nabla f(y), 0) \rangle \\ &\leq \kappa^{-1}L^2\|x - y\|^2 \\ &\leq \kappa^{-2}L^2\|u - v\|_V^2 \end{aligned}$$

by the  $L$ -Lipschitz property of  $\nabla f$  and by Lemma 5.19. We now recall from (5.2.7) that  $B = \nabla F$  in  $\mathcal{W}$ . Thus,  $V^{-1}B = \nabla F$  in  $\mathcal{W}_V$ , and the result follows from Baillon-Haddad's theorem.  $\square$

When  $\kappa > L/2$ , all the assumptions of Theorem 5.15 are satisfied, and each sequence of iterates  $(u^k)$  produced by the algorithm (5.2.9) converges to a point of  $\mathcal{Z}(V^{-1}A + V^{-1}B) = \mathcal{Z}(A + B)$  in  $\mathcal{W}_V$ . Since the topologies of  $\mathcal{W}$  and  $\mathcal{W}_V$  are equivalent, this convergence takes place also in  $\mathcal{W}$ , which concludes Step 2.

We now turn to step 3, making the inclusion (5.2.8) explicit. Recall the expression of  $V + A$  above, and observe that

$$V - B = \begin{bmatrix} \tau^{-1}I - \nabla f & M^* \\ M & \gamma I \end{bmatrix}.$$

Writing  $u^k = (x^k, \phi^k)$ , the inclusion (5.2.8) is written

$$\begin{bmatrix} I - \tau \nabla f & \tau M^* \\ \gamma^{-1}M & I \end{bmatrix} \begin{bmatrix} x^k \\ \phi^k \end{bmatrix} \in \begin{bmatrix} \tau \partial g + I & 2\tau M^* \\ 0 & \gamma^{-1} \partial h^* + I \end{bmatrix} \begin{bmatrix} x^{k+1} \\ \phi^{k+1} \end{bmatrix},$$

or equivalently,

$$\begin{aligned} \phi^{k+1} &= \text{prox}_{\gamma^{-1}h^*}(\gamma^{-1}Mx^k + \phi^k), \\ x^{k+1} &= \text{prox}_{\tau g}(x^k - \tau \nabla f(x^k) + \tau M^*(\phi^k - 2\phi^{k+1})). \end{aligned} \tag{5.2.11}$$

The second equation coincides with (5.2.6c). By Moreau's identity (2.4.1), the first equation can be rewritten as

$$\phi^{k+1} = \phi^k + \gamma^{-1}Mx^k - \gamma^{-1} \text{prox}_{\gamma h}(Mx^k + \gamma \phi^k),$$

which amounts to (5.2.6a)–(5.2.6b).

We now show that the limit  $(x^\infty, \phi^\infty)$  of the sequence  $((x^k, \phi^k))$  is a saddle point. Using Equations (5.2.11) and taking  $k$  to infinity, we obtain  $\phi^\infty = Q_{\gamma^{-1}\partial h^*}(\gamma^{-1}Mx^\infty + \phi^\infty)$  and  $x^\infty = Q_{\tau\partial g}(x^\infty - \tau\nabla f(x^\infty) - \tau M^*\phi^\infty)$ . One can see that these equations are equivalent to the inclusions (5.2.5). Theorem 5.16 is proven.

**Remark 5.1.** The Vũ-Condat algorithm described here can be straightforwardly generalized to the case where  $\nabla f$  is replaced with a general cocoercive operator, and where  $\partial g$  and  $\partial h^*$  are replaced with general maximal monotone operators.

## Discussion

We conclude by commenting briefly the algorithms we introduced in this chapter.

Considering the Vũ-Condat algorithm, we first make the following observations:

- If  $h \circ M = 0$ , we recover the proximal gradient algorithm. Indeed, in this case, we can assume without generality loss that  $M = 0$  and  $h = 0$ . The primal and dual problems become decoupled, and so is the case of the inclusions (5.2.5). The iterations (5.2.6c) coincide with the proximal gradient iterations (4.3.3), and Theorem 4.12 is encompassed by Theorem 5.16.
- In the case where  $f = 0$  and  $M = I$ , we recover the Douglas-Rachford algorithm if we set  $\tau = \gamma$ . Indeed, in this case, Equations (5.2.11) become

$$\phi^{k+1} = \text{prox}_{\gamma^{-1}h^*}(\gamma^{-1}x^k + \phi^k), \quad (5.2.12)$$

$$x^{k+1} = \text{prox}_{\gamma g}(x^k + \gamma(\phi^k - 2\phi^{k+1})). \quad (5.2.13)$$

Write  $s^k = x^k + \gamma\phi^k$ , and let

$$y^{k+1} = \text{prox}_{\gamma h}(s^k). \quad (5.2.14)$$

Then, applying Moreau's identity to the right hand side of Equation (5.2.12), we get that  $\gamma\phi^{k+1} = s^k - y^{k+1}$ . Plugging this equation into (5.2.13), we get

$$\begin{aligned} s^{k+1} &= x^{k+1} + \gamma\phi^{k+1} = \text{prox}_{\gamma g}(x^k + \gamma\phi^k - 2(s^k - y^{k+1})) + s^k - y^{k+1} \\ &= \text{prox}_{\gamma g}(2y^{k+1} - s^k) + s^k - y^{k+1}, \end{aligned} \quad (5.2.15)$$

and it remains to compare Equations (5.2.14)–(5.2.15) with Equations (5.2.2) to obtain the result.

Notice that when  $f = 0$  and  $M = I$ , the assumption  $\tau = \gamma$  is not covered by the statement of Theorem 5.16. However, the proof of this theorem can be adapted to include this case. Details are provided in Condat (2013).

We now provide some observations on the computational complexity of the splitting algorithms studied above. An important observation in this respect is that the Douglas-Rachford algorithm and ADMM require the implementation of proximity operators related with both functions involved in the minimization problem. On the other hand, when one of these functions (namely  $f$ ) is smooth, the proximal gradient and the Vũ-Condat algorithms contend themselves with the gradient of this function. The computation of the proximity operator is sometimes demanding. As an example, assume  $f(x) = 0.5\|Ax - b\|^2$ , where  $A$  is a matrix with large dimensions, as it is frequently the case in the fields of statistical learning and large scale optimization. Then, the

computation of  $\nabla f(x) = A^*(Ax - b)$  only requires a matrix-vector multiplication. On the other hand, since  $\text{prox}_{\gamma f}(x)$  is the unique solution of the equation  $\gamma \nabla f(v) + v = x$ , it is written as

$$\text{prox}_{\gamma f}(x) = (\gamma A^* A + I)^{-1}(x + \gamma A^* b).$$

Before starting the algorithm, the inversion of the matrix  $\gamma A^* A + I$  is required. The algorithm can be implemented only if this matrix inversion is affordable. Fast inversion algorithms exist only when  $A$  has a certain structure (Toeplitz, circulant, isometric, etc.). For general matrices, algorithms based on the computation of  $\nabla f$  are often preferred.

Another point concerns the impact of the operator  $M$  on the computational complexity. In the context of ADMM, the iteration (5.2.3a) requires solving an inclusion of the type

$$\gamma \partial f(v) + M^* M v \in \dots$$

If the operator  $M^* M$  has no particular structure, the solution of this inclusion can be computationally demanding. This problem is avoided by the Vũ-Condat algorithm, where the computational impact of  $M$  on the iterations (5.2.6) is limited to matrix-vector multiplications.

## Exercises

**Exercise 5.7** (Fitzpatrick function and proof of Minty's theorem). This exercise is devoted to the long part of the proof of Minty's theorem: if  $A$  is a maximal monotone operator on a Euclidean space  $\mathcal{X}$ , then  $\text{Im}(I + A) = \mathcal{X}$ .

Let  $A$  be a set-valued operator on a Euclidean space  $\mathcal{X}$ , and let  $\mathcal{G}(A)$  be its graph, assumed to be nonempty. The *Brézis-Haraux* function associated with  $A$  is

$$\begin{aligned} \Phi_A : \mathcal{X} \times \mathcal{X} &\longrightarrow (-\infty, \infty] \\ (x, u) &\longmapsto \sup_{(y, v) \in \mathcal{G}(A)} \langle x - y, v - u \rangle, \end{aligned}$$

and the *Fitzpatrick* function associated with this operator is

$$\begin{aligned} \Psi_A : \mathcal{X} \times \mathcal{X} &\longrightarrow (-\infty, \infty] \\ (x, u) &\longmapsto \Phi_A(x, u) + \langle x, u \rangle, \end{aligned}$$

in other words,

$$\Psi_A(x, u) = \sup_{(y, v) \in \mathcal{G}(A)} \langle x, v \rangle - \langle y, v \rangle + \langle y, u \rangle.$$

We start by showing that these functions characterize the maximal monotonicity of the operator  $A$ .

1. Show that  $A$  is monotone *iff*  $\Phi_A = 0$  on  $\mathcal{G}(A)$ .
2. Show that  $\Phi_A(x, u) \leq 0$  if and only if  $\{(x, u)\} \cup \mathcal{G}(A)$  is the graph of a monotone operator.
3. Deduce from the last question that a monotone operator  $A$  is maximal *iff*  $\Phi_A(x, u) > 0$  whenever  $(x, u) \notin \mathcal{G}(A)$ .

We shall assume in the remainder that  $A$  is a maximal monotone operator. The maximality of  $A$  can be characterized by the statement just shown. Equivalently,

$$A \text{ is maximal } \iff \Psi_A(x, u) > \langle x, u \rangle \text{ whenever } (x, u) \notin \mathcal{G}(A). \quad (5.3.1)$$

It will be more convenient to make use of this last result because  $\Psi_A$  is convex, which  $\Phi_A$  is not in general.



4. Show that  $\Psi_A$  is a proper Fenchel-Legendre transform whose expression can be provided.
5. Let  $f$  be a convex and proper function on  $\mathcal{X}$  that satisfies  $f + \|\cdot\|^2/2 \geq 0$ . Show that there exists  $\phi \in \mathcal{X}$  such that

$$\forall x \in \mathcal{X}, \quad \langle x, \phi \rangle - f(x) + \frac{\|\phi\|^2}{2} \leq 0,$$

and deduce that

$$\forall x \in \mathcal{X}, \quad f(x) + \frac{\|x\|^2}{2} \geq \frac{\|\phi + x\|^2}{2}.$$

*Hint:* Use duality.

6. Show that  $\Psi_A + \|\cdot\|^2/2 \geq 0$ , and deduce that

$$\exists (y, v) \in \mathcal{X} \times \mathcal{X}, \quad \forall (x, u) \in \mathcal{X} \times \mathcal{X}, \quad \Psi_A(x, u) + \frac{\|(x, u)\|^2}{2} \geq \frac{\|(x, u) + (y, v)\|^2}{2}.$$

In the remainder,  $(y, v)$  will be a point that satisfies the previous statement.

7. Show that

$$\begin{aligned} \forall (x, u) \in \mathcal{G}(A), \quad \langle x, u \rangle &\geq \frac{\|v\|^2}{2} + \langle x, v \rangle + \frac{\|y\|^2}{2} + \langle y, u \rangle \\ &\geq -\langle y, v \rangle + \langle x, v \rangle + \langle y, u \rangle, \end{aligned}$$

and deduce that  $v \in Ay$ .

8. Replacing  $(x, u)$  with  $(y, v)$  in the first inequality above, show that  $v = -y$ , thus, that  $0 \in \text{Im}(I + A)$ .
9. Given an arbitrary  $z \in \mathcal{X}$ , show that  $z \in \text{Im}(I + A)$ .

*Hint:* Apply the previous result to a properly defined maximal monotone operator.

# Chapter 6

## Dual methods

### Method of multipliers

#### Problem setting

In this section, we seek to solve the following problem:

$$\min_{x:Ax=0} f(x) \tag{6.1.1}$$

where  $f : \mathcal{X} \rightarrow (-\infty, +\infty]$  is proper closed convex function,  $A$  is a linear operator on  $\mathcal{X} \rightarrow \mathcal{Y}$ , where  $\mathcal{X}$  and  $\mathcal{Y}$  are Euclidean sets. The Lagrangian function associated with the above problem is

$$\mathcal{L}(x, \lambda) = f(x) + \langle \lambda, Ax \rangle$$

and the corresponding dual function is given by

$$\begin{aligned} \Phi(\lambda) &= \inf_{x \in \mathcal{X}} \mathcal{L}(x, \lambda) \\ &= - \sup_{x \in \mathcal{X}} \langle -A^T \lambda, x \rangle - f(x) \\ &= -f^*(-A^T \lambda). \end{aligned}$$

Therefore, the dual problem boils down to:

$$\min_{\lambda \in \mathcal{Y}} f^*(-A^T \lambda).$$

In the sequel, we provide two methods for solving this dual problem.

#### Algorithm

In this paragraph, we restrict our attention to the special case where  $f$  is  $\mu$ -strongly convex. This assumption can be quite restrictive in practice, and we shall see later some alternative methods that can be used in a broader setting. We start with the following lemma.

**Lemma 6.1.** *If  $f$  is  $\mu$ -strongly convex, then  $f^*$  is differentiable and  $\nabla f^*$  is  $\mu^{-1}$ -Lipschitz continuous.*

*Proof.* Let  $\lambda \in \mathcal{Y}$ . By the Fenchel-Young property 2.9,  $x \in \partial f^*(\lambda)$  if and only if  $\lambda \in \partial f(x)$ . By Fermat's rule, this is again equivalent to  $x \in \arg \min f - \langle \lambda, \cdot \rangle$ . As  $f$  is strongly convex, the argument of the minimum exists and is unique. As such,  $x$  is well and uniquely defined. Thus,  $\partial f^*(\lambda)$  is a singleton. Otherwise stated,  $f^*$  is differentiable.

We verify the Lipschitz continuity of  $\nabla f^*$ . Let us fix  $\lambda$  and  $\lambda'$ . Set  $x = \nabla f^*(\lambda)$  and  $y = \nabla f^*(\lambda')$ . Since  $\lambda \in \partial f(x)$ , the strong convexity of  $f$  implies

$$f(y) \geq f(x) + \langle \lambda, y - x \rangle + \frac{\mu}{2} \|y - x\|^2$$

and a similar inequality hold by symetry:

$$f(x) \geq f(y) + \langle \lambda', x - y \rangle + \frac{\mu}{2} \|y - x\|^2$$

Summing these inequalities leads to  $0 \geq \langle \lambda - \lambda', y - x \rangle + \mu \|x - y\|^2$ . Hence,  $\|y - x\|^2 \leq \frac{1}{\mu} \|\lambda - \lambda'\| \|x - y\|$  by the Cauchy-Schwarz inequality. Thus  $\|y - x\| \leq 1/\mu \|\lambda - \lambda'\|$  and the proof is complete.  $\square$

**Remark 6.1.** Lemma 6.1 has a converse. We refer to (Hiriart-Urruty and Lemaréchal, 2012, Theorem 4.2.2).

Therefore, if  $f$  is  $\mu$ -strongly convex, the (negative) dual function  $\lambda \mapsto f^*(-A^T \lambda)$  is differentiable and its gradient is as well Lipschitz continuous. In these circumstances, the previous chapter indicates that a gradient descent method can be used in order to solve the dual problem. The gradient descent writes:

$$\begin{aligned} \lambda^{k+1} &= \lambda^k - \gamma \nabla(f^* \circ (-A^T))(\lambda^k) \\ &= \lambda^k + \gamma A \nabla f^*(-A^T \lambda^k) \end{aligned}$$

where  $\gamma > 0$  is the step size of the gradient descent. Define  $x^{k+1} = \nabla f^*(-A^T \lambda^k)$ . By the Fenchel-Young property 2.9, this is equivalent to  $-A^T \lambda^k \in \partial f(x^{k+1})$  or equivalently

$$0 \in \partial f(x^{k+1}) + A^T \lambda^k.$$

By Fermat's rule, the above inclusion is again equivalent to

$$x^{k+1} = \arg \min_{x \in \mathcal{X}} f(x) + \langle A^T \lambda^k, x \rangle$$

(note that the argument of the minimum exists and is unique due to the strong convexity of  $f$ ). We finally obtain the following iterations, called the *method of multipliers*:

$$\begin{aligned} x^{k+1} &= \arg \min_{x \in \mathcal{X}} \mathcal{L}(x, \lambda^k) \\ \lambda^{k+1} &= \lambda^k + \gamma A x^{k+1}. \end{aligned}$$

We have the following convergence result.

**Theorem 6.2.** Assume that  $f$  is  $\mu$ -strongly convex. Assume that the Lagrangian function  $\mathcal{L}$  has a saddle point. Set  $0 < \gamma < \frac{2\mu}{\|A\|^2}$  where  $\|A\|$  is the spectral norm<sup>1</sup> of  $A$ . Then, the sequence  $(x^k, \lambda^k)$  generated by the method of multipliers converges to a saddle point of  $\mathcal{L}$ .

*Proof.* By Lemma 6.1, the dual function  $\Phi$  is differentiable and  $\nabla \Phi$  is Lipschitz continuous. It is not difficult to show that the corresponding Lipschitz constant is upper bounded by  $\frac{\|A\|^2}{\mu}$ . Therefore, the gradient descent on the (negative) dual function yields a sequence  $\lambda^k$  converging

---

<sup>1</sup>the square root of the largest eigenvalue of  $A^T A$

to a dual solution. It remains to show that  $x^k$  converges to a primal solution  $\lambda^*$ . Recall that  $x^{k+1} = \nabla f^*(-A^T \lambda^k)$ . By continuity of  $\nabla f^*$ ,  $x^k$  converges to  $x^* = \nabla f^*(-A^T \lambda^*)$ . Using once again the Fenchel-Young property 2.9 and Fermat's rule, it follows that  $x^* = \arg \min_x \mathcal{L}(x, \lambda^*)$ . Thus,  $x^*$  is primal-optimal and  $(x^*, \lambda^*)$  is a saddle point of  $\mathcal{L}$ .  $\square$

**Remark 6.2.** The method of multipliers can be used under slightly milder assumptions. The assumption that  $f$  is strongly convex can be replaced by the assumption that the function  $y \mapsto \inf\{f(x) : x \text{ s.t. } Ax = y\}$  is strongly convex. The latter condition is indeed necessary and sufficient to ensure that the dual function  $\Phi$  is differentiable with a Lipschitz continuous gradient. In the absence of strong convexity assumption on  $f$ , note that the quantity  $x^k$  may no longer be uniquely defined and the conclusions of the theorem should thus be weakened.

### Application: a splitting method

In this paragraph, we instantiate the Method of Multipliers as a way to solve the following problem:

$$\min_{x \in \mathcal{X}} f(x) + g(Mx) \quad (6.1.2)$$

where  $f : \mathcal{X} \rightarrow (-\infty, +\infty]$ ,  $g : \mathcal{Y} \rightarrow (-\infty, +\infty]$  are proper closed convex functions and  $M : \mathcal{X} \rightarrow \mathcal{Y}$  is a linear operator. The above problem writes equivalently as

$$\min\{F(y) : y \in \mathcal{X} \times \mathcal{Y}, Ay = 0\}$$

where the function  $F$  is defined on  $\mathcal{X} \times \mathcal{Y} \rightarrow (-\infty, +\infty]$  by  $F(x, z) = f(x) + g(z)$  and where  $A : \mathcal{X} \times \mathcal{Y} \rightarrow \mathcal{Y}$  is the linear operator  $A = [M, -I]$  where  $I$  denotes the identity, that is

$$A \begin{pmatrix} x \\ z \end{pmatrix} = Mx - z$$

for every  $(x, z) \in \mathcal{X} \times \mathcal{Y}$ . Provided that the conditions of application of the method of multipliers are in force, the latter writes:

$$\begin{aligned} (x^{k+1}, z^{k+1}) &= \arg \min_{(x, z)} f(x) + g(z) + \langle \lambda^k, Mx - z \rangle \\ \lambda^{k+1} &= \lambda^k + \gamma(Mx^{k+1} - z^{k+1}). \end{aligned}$$

As a remarkable feature, the first update equation above reduces to solving a minimization problem which is separable in  $(x, y)$ . Finally, the algorithm reformulates as

$$\begin{aligned} x^{k+1} &= \arg \min_x f(x) + \langle \lambda^k, Mx \rangle \\ z^{k+1} &= \arg \min_z g(z) - \langle \lambda^k, z \rangle \\ \lambda^{k+1} &= \lambda^k + \gamma(Mx^{k+1} - z^{k+1}). \end{aligned}$$

The algorithm is called a splitting method, which should be understood in the following sense. The first update equation involves only function  $f$  whereas the second one involves only  $g$ . The algorithm has an interest in the where  $f$  and  $g$  are tractable functions and could be separately handled, but the sum  $f + g \circ M$  is difficult to minimize. Examples will be provided in the next chapter.

## Augmented Method of Multipliers

We now apply the proximal point algorithm to the minimization of the (negative) dual function  $f^* \circ (-A^T)$ . This yields

$$\lambda^{k+1} = \text{prox}_{\gamma f^* \circ (-A^T)}(\lambda^k).$$

In the sequel, we denote by  $\mathcal{L}_\gamma$  the *augmented Lagrangian* defined by

$$\mathcal{L}_\gamma(x, \lambda) = f(x) + \langle \lambda, Ax \rangle + \frac{\gamma}{2} \|Ax\|^2.$$

**Theorem 6.3.** *The proximal point iteration  $\lambda^{k+1} = \text{prox}_{\gamma f^* \circ (-A^T)}(\lambda^k)$  writes equivalently*

$$\begin{aligned} \text{Find } x^{k+1} &\in \arg \min_x \mathcal{L}_\gamma(x, \lambda^k) \\ \text{Set } \lambda^{k+1} &= \lambda^k + \gamma Ax^{k+1}. \end{aligned}$$

*Proof.* By the definition of the proximity operator, the above equation reads equivalently

$$0 \in \gamma \partial(f^* \circ (-A^T))(\lambda^{k+1}) + \lambda_{k+1} - \lambda^k$$

or equivalently

$$0 \in -\gamma A \partial f^*(-A^T \lambda^{k+1}) + \lambda_{k+1} - \lambda^k.$$

This means that there exist some  $x^{k+1} \in \partial f^*(-A^T \lambda^{k+1})$  such that  $\lambda^{k+1} = \lambda^k + \gamma Ax_{k+1}$ . By the Fenchel Young property again, the inclusion  $x^{k+1} \in \partial f^*(-A^T \lambda^{k+1})$  reduces to  $-A^T \lambda^{k+1} \in \partial f(x^{k+1})$  or equivalently

$$0 \in \partial f(x^{k+1}) + A^T(\lambda^k + \gamma Ax_{k+1}).$$

By Fermat's rule, this is again equivalent to  $x^{k+1} \in \arg \min_x \mathcal{L}_\gamma(x, \lambda^k)$ .  $\square$

As the augmented method of multipliers coincides with a proximal point algorithm on the (negative) dual function, it follows that the sequence  $\lambda^k$  converges to a solution to the dual problem. We remark that the variable  $x^{k+1}$  is not necessarily uniquely defined because the argmin may not be unique.

**Remark 6.3.** Apply the augmented method of multipliers to the example (6.4.1). We obtain

$$\begin{aligned} (x^{k+1}, z^{k+1}) &= \arg \min_{(x,z)} f(x) + g(z) + \left\langle \lambda^k, Mx - z \right\rangle + \frac{\gamma}{2} \|Mx - z\|^2 \\ \lambda^{k+1} &= \lambda^k + \gamma(Mx^{k+1} - z^{k+1}). \end{aligned}$$

The minimization problem is no longer separable due to the presence of a new quadratic term. In that sense, the augmented method of multipliers cannot be used as a splitting method.

## Alternating Direction Method of Multipliers (ADMM)

The ADMM can be seen as a variant over the augmented method of multipliers, which combines the good features of the later (there is no strong convexity assumption, no restriction on the step size) and the standard method of multipliers (it is a splitting method allowing to “separate”  $f$  and  $g$  in (6.4.1)). The iterations are given by

$$\begin{aligned}
x^{k+1} &\in \arg \min_x f(x) + \langle \lambda^k, Mx \rangle + \frac{\gamma}{2} \|Mx - z^k\|^2 \\
z^{k+1} &= \arg \min_z g(z) - \langle \lambda^k, z \rangle + \frac{\gamma}{2} \|Mx^{k+1} - z\|^2 \\
\lambda^{k+1} &= \lambda^k + \gamma(Mx^{k+1} - z^{k+1}).
\end{aligned} \tag{6.3.1}$$

We remark that in the above iterations, the quantity  $x^k$  may, in certain circumstances, not be uniquely defined. However, when  $M$  is injective or when  $f$  is strongly convex, the arg min in (6.3.1) is unique, and  $x^{k+1}$  is unambiguously defined. Therefore, in Equation (6.3.1), the symbol “ $\in$ ” can be replaced by “ $=$ ”.

The proof of the following result is taken for granted. We refer to [Boyd et al. \(2011\)](#) for a convergence proof.

**Theorem 6.4.** *Consider a sequence  $(x^k, z^k, \lambda^k)$  satisfying the ADMM iterations and assume that a saddle point of the Lagrangian exists. Then,*

- the sequence  $\lambda^k$  converges to a solution to the dual problem,
- any limit point of the sequence  $x^k$  is a primal solution,
- the sequence  $f(x^k) + g(z^k)$  converges to the primal value  $p = \inf f + g \circ M$ ,
- the sequence  $Mx^k - z^k$  tends to zero.

Moreover, if  $M$  is injective, then  $x^k$  converges to a primal solution.

## The Vũ-Condat method

To conclude this chapter on primal-dual algorithms, let us present a method that uses even more structure in the problem in order to reduce even more the per-iteration complexity than the ADMM. The Vũ-Condat method has been designed to solve problems of the type

$$\min_{x \in \mathcal{X}} f_1(x) + f_2(x) + g(Mx) \tag{6.4.1}$$

where  $f_1 : \mathcal{X} \rightarrow (-\infty, +\infty)$  is a convex differentiable function with Lipschitz gradient,  $f_1 : \mathcal{X} \rightarrow (-\infty, +\infty]$  and  $g : \mathcal{Y} \rightarrow (-\infty, +\infty]$  are proper closed convex functions and  $M : \mathcal{X} \rightarrow \mathcal{Y}$  is a linear operator.

The iterations of the algorithm are

$$\begin{aligned}
x^{k+1} &= \arg \min_x f_1(x_k) + \langle \nabla f_1(x_k), x - x_k \rangle + f_2(x) + \langle \lambda^k, Mx \rangle + \frac{\gamma}{2} \|x - x^k\|^2 \\
\lambda^{k+1} &= \arg \min_{\lambda} g^*(\lambda) - \langle \lambda, M(2x_{k+1} - x_k) \rangle + \frac{\sigma}{2} \|\lambda - \lambda_k\|^2
\end{aligned} \tag{6.4.2}$$

The algorithm can also be written as

$$\begin{aligned}
x^{k+1} &= \text{prox}_{\gamma^{-1}f_2} (x_k - \gamma^{-1}(\nabla f_1(x_k) + M^\top \lambda_k)) \\
\lambda^{k+1} &= \text{prox}_{\sigma^{-1}g^*} (\lambda_k + \sigma^{-1}M(2x_{k+1} - x_k))
\end{aligned}$$

This shows that if the gradient of  $f_1$  is available and if the proximity operators of  $f_2$  and  $g^*$  are easy to compute, then each iteration of the algorithm is cheap in computations. In

particular, unlike the augmented method of multipliers or the ADMM, the linear operator  $M$  is never involved in sometimes complex subproblems: it is only involved through matrix-vector multiplications. Hence, the method is particularly suited to large scale optimization problems. The counterpart is that the Vũ-Condat method usually requires more iterations than the ADMM to get an approximate solution of good quality.

**Theorem 6.5.** *If the step-size parameters satisfy*

$$\gamma > \frac{L(\nabla f_1)}{2} + \sigma^{-1} \|M\|^2 \quad (6.4.3)$$

*then the fixed point operator defining the Vũ-Condat algorithm is averaged and thus the sequence  $(x^k, \lambda^k)_k$  defined by (6.4.2) converges to a saddle point of the Lagrangian. In particular, if  $0 \in \text{ri}(\text{dom } g - M \text{ dom } f_2)$ , then*

$$x^k \rightarrow x^* \in \arg \min_x f_1(x) + f_2(x) + g(Mx)$$

# Notations

$\text{cl}(C)$	: Closure of the set $C$
$\dim(\mathcal{V})$	: Dimension of the vector space $\mathcal{V}$
$\text{dom}(A)$	: Domain of the set-valued operator $A$
$\text{dom}(f)$	: Domain of the function $f$
$\text{epi}(f)$	: Epigraph of the function $f$
$f^*$	: Fenchel-Legendre transform of the function $f$
$f \square g$	: Infimal convolution of the functions $f$ and $g$
$\text{Fix}(T)$	: Set of fixed points of the operator $T$
$\text{int}(C)$	: Interior of the set $C$
$\text{l.s.c.}$	: Lower semicontinuous
$M \triangleright f$	: Infimal postcomposition of the the linear operator $M$ and the function $f$
$P_C$	: Projection on the closed convex set $C$
$\text{prox}_f$	: Proximal operator associated with the function $f$
$Q_A$	: Resolvent of the operator $A$
$\text{ri}(C)$	: Relative interior of the set $C$
$\mathcal{W}, \mathcal{X}, \mathcal{Y}$	: Euclidean spaces
$\mathcal{Z}(A)$	: Set of zeros of a set-valued operator $A$ (P. 41)
$\partial f$	: Subdifferential of the function $f$
$\iota_C$	: Indicator function of the set $C$
$\Gamma_0(\mathcal{X})$	: Set of convex, proper and lower semicontinuous functions on the Euclidean space $\mathcal{X}$
$\nabla f$	: Gradient of the function $f$
$\ \cdot\ $	: Euclidean norm of a vector or spectral norm of a matrix



# Bibliography

- Bauschke, H. H. and Combettes, P. L. (2011). *Convex analysis and monotone operator theory in Hilbert spaces*. Springer.
- Borwein, J. and Lewis, A. (2006). *Convex Analysis and Nonlinear Optimization: Theory and Examples*. CMS Books in Mathematics. Springer. 31
- Boyd, S., Parikh, N., Chu, E., Peleato, B., and Eckstein, J. (2011). Distributed optimization and statistical learning via the alternating direction method of multipliers. *Foundations and Trends® in Machine Learning*, 3(1):1–122. 47, 61
- Boyd, S. and Vandenberghe, L. (2009). *Convex optimization*. Cambridge university press. 31
- Brézis, H. (1973). *Opérateurs maximaux monotones et semi-groupes de contractions dans les espaces de Hilbert*. North-Holland mathematics studies. Elsevier Science, Burlington, MA. 51
- Brézis, H. (1987). Analyse fonctionnelle, 2e tirage.
- Condat, L. (2013). A primal-dual splitting method for convex optimization involving Lipschitzian, proximable and linear composite terms. *J. Optim. Theory Appl.*, 158(2):460–479. 54
- Hiriart-Urruty, J.-B. and Lemaréchal, C. (2012). *Fundamentals of convex analysis*. Springer Science & Business Media. 58
- Nesterov, Y. (2004). *Introductory lectures on convex optimization: A basic course*, volume 87. Springer.
- Rockafellar, R. T. (2015). *Convex analysis*. Princeton university press. 6
- Rockafellar, R. T., Wets, R. J.-B., and Wets, M. (1998). *Variational analysis*, volume 317. Springer. 13