

Continuous optimization, an introduction

Antonin Chambolle (and many others, ...)

February 9, 2022

Contents

1	Introduction	2
2	(First order) Descent methods, rates	2
2.1	Gradient descent	3
2.2	What can we achieve?	7
2.3	Second order methods: Newton's method	9
2.4	Multistep first order methods	10
2.4.1	Heavy ball method	10
2.4.2	The conjugate gradient method	12
2.4.3	Accelerated algorithm: Nesterov 83	14
2.5	Nonsmooth problems?	15
2.5.1	Subgradient descent	15
2.5.2	Implicit descent	15
3	Krasnoselskii-Mann's convergence theorem	17
3.1	A "general" convergence theorem	17
3.2	Varying steps	21
3.3	A variant with errors	21
3.4	Examples	22
3.4.1	Gradient descent	22
3.4.2	Composition of averaged operators	22
4	An introduction to convex analysis and monotone operators	23
4.1	Convexity	24
4.1.1	Convex functions	24
4.1.2	Separation of convex sets	25
4.1.3	Subgradient	27
4.1.4	Subdifferential calculus	29
4.1.5	Remark: KKT's theorem	32
4.2	Convex duality	33
4.2.1	Legendre-Fenchel conjugate	33
4.2.2	Examples	35
4.2.3	Relationship between the growth of f and f^*	36
4.2.4	The conjugate of a sum: Inf-convolutions	38
4.3	Example: the proximity operator	39
4.3.1	A useful variant of inf-convolutions	41
4.3.2	Fenchel-Rockafellar duality	41

4.4	Generalization: Elements of monotone operators theory	43
5	Algorithms. Operator splitting	48
5.1	Abstract algorithms for monotone operators	48
5.1.1	Explicit algorithm	49
5.1.2	Proximal point algorithm	49
5.1.3	Forward-Backward splitting	50
5.1.4	Douglas-Rachford splitting	51
5.1.5	Three-operators splitting	51
5.2	Descent algorithms, acceleration, “FISTA”	53
5.2.1	Forward-Backward descent	53
5.2.2	FISTA	54
5.2.3	Convergence rates	56
5.3	ADMM, Douglas-Rachford splitting	61
5.4	Other saddle-point algorithms: Primal-dual algorithm	63
5.4.1	Rate	65
5.4.2	Extensions	65
6	“Large scale” optimization	66
6.1	Coordinate descent and stochastic coordinate descent	67
6.1.1	Does coordinate descent / alternating minimization work?	67
6.1.2	Block coordinate descent	68
6.1.3	Random coordinate descent	68
6.2	Stochastic gradient descent	70
6.3	SGD for learning problems	70
6.3.1	Improvements of SGD	72

1 Introduction

These lectures notes have been prepared by A. Chambolle for the M2 course “Continuous optimization” given between Oct. and Dec. 2016 in Paris 6, in the Master “modélisation mathématique” of Ecole Polytechnique, Université Pierre-et-Marie Curie (Paris 6) and Ecole National des Ponts et Chaussées. Much of the material is taken from [9], or/and inspired by famous textbooks [33, 29, 38, 15, 2]. Updated for the 2017-2018-2019 courses. The 2020 courses will take place in Université Paris-Dauphine PSL.

The notes gather various material mostly on first order optimisation and iterative algorithms for generally convex problems, including operator splitting, acceleration, etc.

2 (First order) Descent methods, rates

Most of what we describe in this section is in finite dimension, although extension to Hilbert spaces is in general easy. We will discuss rates of convergence, in particular, which we try to make independent on the dimension. The complexity of the iterations, on the other hand, are usually very dimension-dependent, and this is the reason for which high order descent methods are not practical for modern high dimensional problems (imaging, data analysis...).

2.1 Gradient descent

The main source for this section is the reference textbook of Polyak [33]. Consider the problem of minimising

$$\min_{x \in X} f(x)$$

with X a finite dimensional vector space (or Hilbert) and f a real valued, C^1 function (or at least differentiable). We denote X^* the dual of X (which can of course be represented by X through the scalar product).

The differential $df(x) \in X^*$ is defined as the linear part of the closest affine function to f at $(x, f(x))$:

$$f(y) = f(x) + df(x) \cdot (y - x) + o(|x - y|).$$

The function f is said to be (Fréchet) differentiable at x if such an affine approximation exists, and it is C^1 in X if

$$\begin{aligned} df : X &\rightarrow X^* \\ x &\mapsto df(x) \end{aligned}$$

is defined everywhere and continuous. The local inversion theorem guarantees that in this case, near points where $df(x) \neq 0$, the level set $\{f = f(x)\}$ is a C^1 hypersurface with tangent space $\text{Ker } df(x) = \{h : df(x) \cdot h = 0\}$.

When X has a Euclidean or Hilbertian structure (a scalar product), then $df(x)$ has the (Riesz) representation $df(x) \cdot h = \langle \nabla f(x), h \rangle_X$ ($\forall h$), where now $\nabla f(x)$ is the gradient of f at x (which depends on the metric structure of X). One has obviously $\text{Ker } df(x) = \nabla f(x)^\perp$ and $\nabla f(x)$ is a normal vector to the level surface $\{f = f(x)\}$ of f at x , pointing towards the larger values.

For this reason, the most simple idea to minimize the function f is to introduce the “gradient descent algorithm” with step τ :

$$x^{k+1} = x^k - \tau \nabla f(x^k) =: T_\tau(x^k).$$

As said above, $-\nabla f(x^k)$ is a descent direction. Near x^k , indeed,

$$f(x) = f(x^k) + \langle \nabla f(x^k), x - x^k \rangle + o(|x - x^k|)$$

so that

$$f(x^{k+1}) = f(x^k) - \tau |\nabla f(x^k)|^2 + o(\tau) < f(x^k)$$

if $\tau > 0$ is small enough and $\nabla f(x^k) \neq 0$. One can use various strategies to choose τ :

- optimal step: $\min_\tau f(x^k - \tau \nabla f(x^k))$ (with a “line search”, such as for instance for the “conjugate gradient method”);
- Armijo-type rule: find $i \geq 0$ such that $f(x^k - \tau \rho^i \nabla f(x^k)) \leq f(x^k) - c \tau \rho^i |\nabla f(x^k)|^2$, $\rho < 1, c < 1$ fixed;
- Gradient with fixed step: $\tau > 0$ is given, and one sees that one can interpret x^{k+1} as the minimizer of a quadratic approximation of f :

$$x^{k+1} = \arg \min_x f(x^k) + \langle \nabla f(x^k), x - x^k \rangle + \frac{1}{2\tau} |x - x^k|^2.$$

Observe that the latter choice has the form

$$\min_x f(x^k) + df(x^k) \cdot (x - x^k) + \frac{1}{2\tau} d(x, x^k)^2$$

with the distance $d(x, x^k) = |x - x^k|$ given by the Euclidean metric structure of X . It can be natural, in some cases, to consider varying the metric (inside, or sometimes even outside of the Euclidean framework). A particular situation (which is therefore not anymore a “gradient descent” method in the above sense, at least in general), is the

- “Frank-Wolfe”-type method¹: $\min_{x \in X^k} f(x^k) + \langle \nabla f(x^k), x - x^k \rangle$ where X^k is appropriately defined.

Here we replace the metric term $|x - x^k|^2$ with a constraint on x depending on the previous iterates. For instance: $X^k := \{x : |x - x^k| \leq \varepsilon\}$, in which case $\tau = \varepsilon/|\nabla f(x^k)|$ (then we recover a gradient descent method). In other instances, once the minimizer x of the above problem is found, one lets $g^k = x^k - x$ and $x^{k+1} = x^k - \tau g^k$, g^k playing the role of a gradient, but with some local metric (hence the name “conditional gradient”).

Convergence analysis: if τ is too large with respect to the Lipschitz constant of ∇f , or ∇f is not Lipschitz, easy to build infinitely oscillating examples (ex: $f(x) = \|x\|$).

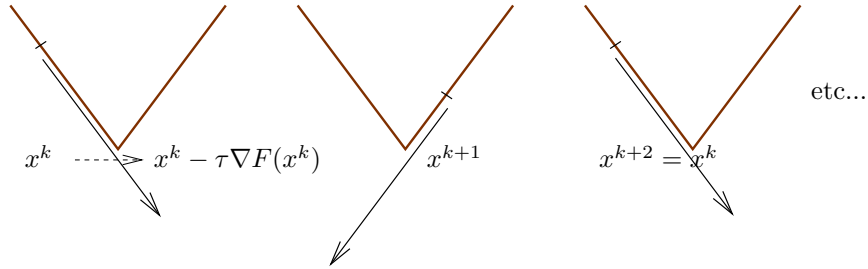


Figure 1: The gradient descent may never converge if the step is too large or the function not smooth enough

If f is C^1 , ∇f is L -Lipschitz, $0 < \tau < 2/L$, $\inf f > -\infty$ then the method converges (in \mathbb{R}^N) in the following sense: $\nabla f(x^k) \rightarrow 0$.

Proof:

$$\begin{aligned} f(x^{k+1}) &= f(x^k) - \int_0^\tau \langle \nabla f(x^k - s\nabla f(x^k)), \nabla f(x^k) \rangle \\ &= f(x^k) - \tau \|\nabla f(x^k)\|^2 + \int_0^\tau \langle \nabla f(x^k) - \nabla f(x^k - s\nabla f(x^k)), \nabla f(x^k) \rangle \\ &\leq f(x^k) - \tau(1 - \frac{L\tau}{2}) \|\nabla f(x^k)\|^2. \end{aligned} \quad (1)$$

(Observe that we just use here that $D^2 f$ is bounded from above by LI (if $f \in C^2$), or more generally, letting $x = x^k$ and $y = x^k - s\nabla f(x^k)$, we use

$$\langle \nabla f(x) - \nabla f(y), x - y \rangle \leq L\|x - y\|^2$$

¹or “conditional gradient”.

which is a consequence of the Lipschitz property of ∇f , but is a weaker property.)

Then letting $\kappa = \tau(1 - \tau L/2) > 0$, one finds that

$$f(x^n) + \kappa \sum_{k=0}^{n-1} \|\nabla f(x^k)\|^2 \leq f(x^0).$$

This shows the claim. If in addition f is “infinite at infinity” (coercive) then x^k has subsequences which converge, therefore to a stationary point.

Remark 2.1. If $\tau = 0$, the iteration does nothing (and hence converges to the initial point...). If $\tau = 2/L$, the iteration might oscillate forever, as shows the example of the function $f(x) := L|x|^2/2$.

Remark 2.2. Taking x^* a minimizer, $\tau = 1/L$, we deduce that

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

The convex case Further information on the second order behaviour of f allows to improve the analysis of the algorithm. The gradient descent method is better analysed assuming that f is convex. One can show (i) that the iteration is a 1-Lipschitz mapping (hence the iterates have to get closer to fixed points, or at least can not move away, during the process), (ii) basic convergence rates (that is, a speed of convergence of $f(x^k)$ towards its minimal value).

First, if f is convex we have the following additional property:

Theorem 2.3 (Baillon-Haddad²). *If f is convex and ∇f is L -Lipschitz, then for all x, y ,*

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

(∇f is said to be “ $(1/L)$ -co-coercive”.)

We will see later a general proof of this result based on convex analysis. In finite dimension, if f is C^2 , then the proof is easy: one has $0 \leq D^2 f \leq LI$ (because f is convex, and because ∇f is L -Lipschitz). Then

$$\nabla f(x) - \nabla f(y) = \int_0^1 D^2 f(y + s(x - y))(x - y) ds =: A(x - y).$$

with $A = \int_0^1 D^2 f(y + s(x - y)) ds$ symmetric with $0 \leq A \leq LI$. Hence:

$$\begin{aligned} \|\nabla f(x) - \nabla f(y)\|^2 &= \|A(x - y)\|^2 = \langle AA^{1/2}(x - y), A^{1/2}(x - y) \rangle \leq \\ &L \langle A^{1/2}(x - y), A^{1/2}(x - y) \rangle \leq L \langle A(x - y), x - y \rangle = L \langle \nabla f(x) - \nabla f(y), x - y \rangle \end{aligned}$$

which is the result. If f is not C^2 , one could smooth f by convolution with a smooth, compactly supported kernel, derive the result and then pass to the limit.

Lemma 2.4. *If f is convex with L -Lipschitz gradient, then the mapping $T_\tau = I - \tau \nabla f$ is a weak contraction when $0 \leq \tau \leq 2/L$ (that is, T_τ is 1-Lipschitz, or “non-expansive”).*

²This is a modest corollary of a much more general result, in arbitrary topological spaces, for operators which satisfy “cyclic monotonicity” conditions, see [1].

Proof:

$$\begin{aligned}\|T_\tau x - T_\tau y\|^2 &= \|x - y\|^2 - 2\tau \langle x - y, \nabla f(x) - \nabla f(y) \rangle + \tau^2 \|\nabla f(x) - \nabla f(y)\|^2 \\ &\leq \|x - y\|^2 - \frac{2\tau}{L} \left(1 - \frac{\tau L}{2}\right) \|\nabla f(x) - \nabla f(y)\|^2.\end{aligned}$$

Remark 2.5. T_τ is “averaged” for $0 < \tau < 2/L$, that is one can write

$$T_\tau = \theta(I - \frac{2}{L}\nabla f) + (1 - \theta)I$$

for $\theta = \tau L/2 \in]0, 1[$. The convergence of the iterates of this class of operators will be proved later on, see Section 3.1.

Convergence rate in the convex case. Additionally, we Then, using that, for x^* a minimizer,

$$f(x^*) \geq f(x^k) + \langle \nabla f(x^k), x^* - x^k \rangle$$

(we will see this is a general property of convex functions), we find

$$\frac{f(x^k) - f(x^*)}{\|x^* - x^k\|} \leq \|\nabla f(x^k)\| \quad (2)$$

Assuming still that $0 < \tau L < 2$, and using Lemma 2.4 which implies that $\|x^k - x^*\| \leq \|x^0 - x^*\|$, it follows $(f(x^k) - f(x^*))/\|x^0 - x^k\| \leq \|\nabla f(x^k)\|$. Hence from (1) we derive, letting $\Delta_k = f(x^k) - f(x^*)$, $\kappa = \tau(1 - \tau L/2) \in]0, 1/(2L)[$, that

$$\Delta_{k+1} \leq \Delta_k - \frac{\kappa}{\|x^0 - x^*\|^2} \Delta_k^2 \quad (3)$$

We can show the following:

Lemma 2.6. *Let $(a_k)_k$ be a sequence of nonnegative numbers satisfying for $k \geq 0$:*

$$a_{k+1} \leq a_k - c^{-1} a_k^2$$

Then, for all $k \geq 0$,

$$a_k \leq \frac{c}{k+1}$$

Proof: First observe that if we replace a_k with a_k/c , the property becomes $a_{k+1} \leq a_k - a_k^2$: hence it is enough to prove it for $c = 1$. Then, as $a_k(1 - a_k) \geq a_{k+1} \geq 0$, one has $0 \leq a_k \leq 1$ for all $k \geq 0$. We show the inequality by induction: for $k = 0$, $a_0 \leq 1$. If $k \geq 1$ and if $ka_{k-1} \leq 1$, then we write that

$$\begin{aligned}(k+1)a_k &\leq (k+1)(a_{k-1} - a_{k-1}^2) \\ &= (k+1)a_{k-1} - (k+1)a_{k-1}^2 = ka_{k-1} + a_{k-1}(1 - (k+1)a_{k-1}) \\ &\leq 1 + a_{k-1}(1 - (k+1)a_k)\end{aligned}$$

since $0 \leq a_k \leq a_{k-1}$? Hence $(1 - (k+1)a_k)(1 + a_{k-1}) \geq 0$. It follows that $(k+1)a_k \leq 1$. Applying this Lemma to the recursion (3) we deduce:

Theorem 2.7. *The gradient descent with fixed step satisfies*

$$\Delta_k \leq \frac{\|x^0 - x^*\|^2}{\kappa(k+1)}$$

Observe that this rate is not very good and also a bit pessimistic (it should improve if $x^k \rightarrow x^*$ because (2) improves: but without further knowledge of f it is impossible to guess how much). On the other hand, it does not prove, a priori, anything on the sequence (x^k) itself. Observe also, to conclude that $\kappa = \tau(1 - \tau L/2) = (2/L)[(\tau L/2)(1 - \tau L/2)]$ is maximal for $\tau L/2 = 1/2$, that is, $\tau = 1/L$. In that case, $\kappa = 1/(2L)$ and the rate is bounded by

$$\Delta_k \leq 2L \frac{\|x^0 - x^*\|^2}{k+1}.$$

Strongly convex case. A function f is γ -strongly convex if and only if $f(x) - \gamma\|x\|^2/2$ is convex: if f is C^2 , it is equivalent to $D^2f \geq \gamma I$. We will discuss more precisely this definition in Section 4.1. In this, case if x^* is the minimizer (which in this case always exists and is unique)

$$x^{k+1} - x^* = x^k - x^* - \tau(\nabla f(x^k) - \nabla f(x^*)) = \int_0^1 (I - \tau D^2 f(x^* + s(x^k - x^*))(x^k - x^*)) ds$$

hence (using that $(1 - \tau L)I \leq I - \tau D^2 f \leq (1 - \tau\gamma)I$)

$$\|x^{k+1} - x^*\| \leq \max\{1 - \tau\gamma, \tau L - 1\} \|x^k - x^*\|.$$

If f is not C^2 one can still show this by smoothing. The best constant is for $\tau = 2/(L + \gamma)$ and gives, for $q = (L - \gamma)/(L + \gamma) \in [0, 1]$

$$\|x^k - x^*\| \leq q^k \|x^0 - x^*\|.$$

One can easily deduce the following (apparently) more general result:

Theorem 2.8. *Let f be C^2 , x^* be a strict local minimum of f where D^2f is definite positive. Then if x^0 is close enough to x^* , the gradient descent method with optimal step (obtained with a line search) will converge linearly. (Or with fixed step small enough.)*

2.2 What can we achieve?

This paragraph contains a very elementary introduction to lower bounds and complexity. We follow the description in [9], where we essentially give elementary variants of deeper results found in [26, 29].)

Idea: consider a “hard problem”, for instance, for $x \in \mathbb{R}^n$, $L > 0$, $\gamma \geq 0$, $1 \leq p \leq n$, functions of the form:

$$f(x) = \frac{L - \gamma}{8} \left((x_1 - 1)^2 + \sum_{i=2}^p (x_i - x_{i-1})^2 \right) + \frac{\gamma}{2} \|x\|^2, \quad (4)$$

which is tackled by a “first order method”, which is such that the iterates x^k are restricted to the subspace spanned by the gradients of already computed iterates, i.e. for $k \geq 0$

$$x^k \in x^0 + \{\nabla f(x^0), \nabla f(x^1), \dots, \nabla f(x^{k-1})\}, \quad (5)$$

where x^0 is an arbitrary starting point.

Starting from an initial point $x^0 = 0$, any first order method of the considered class can transmit the information of the data term only at the speed of one index per iteration. This makes such problems very hard to solve by any first order methods in the

considered class of algorithms. Indeed if one starts from $x^0 = 0$ in the above problem (whose solution, for $\gamma = 0$, is given by $x_l^* = 1$, $k = 1, \dots, p$, and 0 for $l > p$), then at the first iteration, only the first component x_1^1 will be updated (since $\partial_i f(x^0) = 0$ for $i \geq 2$), and by induction one can check that at iteration k , $x_l^k = 0$ for $l \geq k + 1$.

The solution satisfies $\nabla f = 0$, therefore is characterized by

$$x_i = \frac{L - \gamma}{L + \gamma} \frac{x_{i+1} + x_{i-1}}{2}, \quad i \leq p - 1,$$

with $x_0 = 1$ and $x_p = (L - \gamma)/(L + 3\gamma)x_{p-1}$. The best possible point at iteration k satisfies this equation for $i \leq k$, and $x_{k+1} = 0$. In case $\gamma = 0$ we find that this point x is affine: $x_i = (1 - i/(k + 1))^+$, and $x_i - x_{i-1} = -1/(k + 1)$ for $i \leq k + 1$. Hence

$$f(x) = \frac{L}{8} \sum_{i=1}^{k+1} \frac{1}{(k+1)^2} = \frac{L}{8} \frac{1}{k+1}$$

is the best possible value which can be reached at step k .

If one looks for a bound independent on the dimension with (here for homogeneity reasons) $f(x^k) \sim L \|x^0 - x^*\|^2 a_k$ (for a sequence (a_k)), using here that $x_i^* = 1$ for $i \leq p$ and 0 for $i > p$, $x^0 = 0$, and $f(x^*) = 0$, one obtains

$$f(x^k) - f(x^*) \geq \frac{L}{8p(k+1)} \|x^0 - x^*\|^2$$

($k < p$) (while if $k = p$, $x^k = x^*$). For $k = p - 1$ one finds

$$f(x^k) - f(x^*) \geq \frac{L}{8} \frac{\|x^0 - x^k\|^2}{(k+1)^2}$$

hence no first order method can reach a bound of the considered form which is better than this. (*It does not contradict* a bound of the form $f(x^k) - f(x^*) = o(1/k^2)$, for instance!)

It follows a variant of the results in [29] (where a slightly different function is used), see Theorems 2.1.7 and 2.1.13.

Theorem 2.9. *For any $n \geq 2$, any $x^0 \in \mathbb{R}^n$, $L > 0$, and $k < n$, there exists a convex, one times continuously differentiable function f with L -Lipschitz continuous gradient, such that for any first-order algorithm satisfying (5), it holds that*

$$f(x^k) - f(x^*) \geq \frac{L \|x^0 - x^*\|^2}{8(k+1)^2}, \quad (6)$$

where x^* denotes a minimiser of f .

Observe that the above lower bound is valid only if number of iterates k is less than the problem size. We can not improve this with a quadratic function, as the conjugate gradient method (which is a first-order method) is then known to find the global minimiser here after at most p steps.

But practical problems are often so large that it is not possible to perform as many iterations as the dimension of the problem, and will always fulfill similar assumptions.

If choosing $\gamma > 0$ so that the function (4) becomes γ -strongly convex, a lower bound for first order methods is given Theorem 2.1.13 in [29]. It is hard to derive precisely for

p finite, however in $\mathbb{R}^{\mathbb{N}} \simeq \ell^2(\mathbb{N})$, for $p = +\infty$, one finds that the solution is given by $x = q^i$, $q = (\sqrt{Q} - 1)/(\sqrt{Q} + 1)$ where $Q = L/\gamma$ is the condition number of the problem (q satisfies $2 = (L - \gamma)/(L + \gamma)(q + 1/q)$). If $x^0 = 0$,

$$\|x^0 - x^*\|^2 = \sum_{i=1}^{\infty} q^{2i} = \frac{q^2}{1 - q^2}, \text{ while}$$

$$\|x^k - x^*\|^2 \geq \sum_{i=k+1}^{\infty} q^{2i} = q^{2k} \|x^0 - x^*\|^2.$$

The strong convexity of f shows that

$$f(x^k) \geq f(x^*) + \frac{\gamma}{2} q^{2k} \|x^0 - x^*\|^2$$

and it follows:

Theorem 2.10. *For any $x^0 \in \mathbb{R}^{\infty} \simeq \ell_2(\mathbb{N})$ and $\gamma, L > 0$ there exists a γ -strongly convex, one times continuously differentiable function f with L -Lipschitz continuous gradient, such that for any algorithm in the class of first order algorithms defined through (5) it holds that for all k ,*

$$f(x^k) - f(x^*) \geq \frac{\gamma}{2} \left(\frac{\sqrt{Q} - 1}{\sqrt{Q} + 1} \right)^{2k} \|x^0 - x^*\|^2 \quad (7)$$

where $Q = L/\gamma \geq 1$ is the condition number, and x^* a minimiser of f .

In finite dimension, a similar result will hold for k small enough (with respect to n).

The meaning of the two results above is the following: given a first order method, one will never be able to beat in general the rates in the theorem without additional assumptions or properties of the function f or the space (dimension, etc).

2.3 Second order methods: Newton's method

The idea of Newton's method relies on using second order information to improve the precision of the approximation of the function at step k . (In practice, one solves the equation $\nabla f(x) = 0$ using Newton's standard method.) We have

$$f(x) = f(x^k) + \langle \nabla f(x^k), x - x^k \rangle + \frac{1}{2} \langle D^2 f(x^k)(x - x^k), x - x^k \rangle + o(\|x - x^k\|^2).$$

If we are near a minimizer, we can assume $D^2 f(x^k) > 0$ (hopefully), and hence find x^{k+1} by solving

$$\min_x f(x^k) + \langle \nabla f(x^k), x - x^k \rangle + \frac{1}{2} \langle D^2 f(x^k)(x - x^k), x - x^k \rangle$$

Compare with the Gradient descent with step τ in a metric defined by a symmetric positive definite matrix $A > 0$, which would be:

$$\min_x f(x^k) + \langle \nabla f(x^k), x - x^k \rangle + \frac{1}{2\tau} \langle A(x - x^k), x - x^k \rangle$$

hence we can see Newton's method as a gradient descent in the metric which best approximates the function. We find that x^{k+1} is given by

$$\nabla f(x^k) + D^2 f(x^k)(x^{k+1} - x^k) = 0 \Leftrightarrow x^{k+1} = x^k - D^2 f(x^k)^{-1} \nabla f(x^k).$$

We have the following "quadratic" convergence rate.

Theorem 2.11. *Assume f is C^2 , D^2f is M -Lipschitz, and $D^2f \geq \gamma$ (strong convexity). Let $q = M/(2\gamma^2)\|\nabla f(x^0)\|$ and assume x^0 is close enough to the minimizer x^* , so that $q < 1$. Then $\|x^k - x^*\| \leq (2\gamma/M)q^{2^k}$.*

This is extremely fast (the precision is doubled at each iteration, this is called a quadratic rate), but there are strong conditions, and the algorithm can be hard to implement.

Proof: first see that

$$\begin{aligned}\nabla f(x+h) &= \nabla f(x) + \int_0^1 D^2f(x+sh)h ds \\ &= \nabla f(x) + D^2f(x)h + \int_0^1 (D^2f(x+sh) - D^2f(x))h ds\end{aligned}$$

so that

$$\|\nabla f(x+h) - \nabla f(x) - D^2f(x)h\| \leq \frac{M}{2}\|h\|^2.$$

Hence

$$\begin{aligned}\|\nabla f(x^{k+1}) - \overbrace{\nabla f(x^k) - D^2f(x^k)(x^{k+1} - x^k)}^0\| &\leq \frac{M}{2}\|x^{k+1} - x^k\|^2 \\ \Rightarrow \|\nabla f(x^{k+1})\| &\leq \frac{M}{2}\|D^2f(x^k)^{-1}\|^2\|\nabla f(x^k)\|^2 \leq \frac{M}{2\gamma^2}\|\nabla f(x^k)\|^2\end{aligned}$$

Hence letting $g_k = \|\nabla f(x^k)\|$, for all k one has

$$\log g_{k+1} \leq 2 \log g_k + \log \frac{M}{2\gamma^2} \Rightarrow \log g_k \leq 2^k \log g_0 + (2^k - 1) \log \frac{M}{2\gamma^2} = 2^k \log q - \log \frac{M}{2\gamma^2}$$

so that

$$\|\nabla f(x^k)\| \leq \frac{2\gamma^2}{M}q^{2^k}.$$

As f is strongly convex, $\langle \nabla f(x^k), x^k - x^* \rangle \geq \gamma\|x^k - x^*\|^2$, and we can conclude.

The main issue with this is that it is very important to have $q < 1$, otherwise the method could not work. The (very) “good” rate of convergence is obtained only if the starting point is good enough.

There are quite a few very important variants of Newton’s method, which are designed so that one does not have to explicitly evaluate $D^2f(x^k)^{-1}$, usually called “Quasi-Newton” type methods: one replaces $D^2f(x^k)$ with a metric H_k which is improved at each iteration, hoping that $H_k \rightarrow D^2f(x^*)$ in the limit. The most famous (and very efficient) variant is known as the “BFGS” method (after Broyden-Fletcher-Goldfarb-Shanno, detailed in 4 papers of 1970) and its improvements (limited memory “L-BFGS”) [8, 25]. This topic is covered extensively for instance in [30, Chap. 6] and various toolboxes exist which implement this method.

2.4 Multistep first order methods

2.4.1 Heavy ball method

This description follows Polyak’s book [32] where the method is introduced. The idea is to iterate:

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

$\alpha, \beta \geq 0$. This mimicks the equation $\ddot{x} = -\nabla f(x) - \dot{x}$ of a heavy ball in a potential $f(x)$ with friction, which can be discretized as (for instance):

$$\frac{x^{k+1} - 2x^k + x^{k-1}}{(\delta t)^2} + \frac{x^{k+1} - x^k}{\delta t} = -\nabla f(x^k).$$

The method requires that f is C^2 , γ -convex, with L -Lipschitz gradient (at least near a solution x^*), that is:

$$\gamma I \leq D^2 f \leq LI.$$

Then (see [33])

Theorem 2.12. *Let x^* be a (local) minimizer of f such that $\gamma I \leq D^2 f(x^*) \leq LI$, and choose α, β with $0 \leq \beta < 1$, $0 < \alpha < 2(1 + \beta)/L$. There exists $q < 1$ such that if $q < q' < 1$ and if x^0, x^1 are close enough to x^* , one has*

$$\|x^k - x^*\| \leq c(q')q'^k.$$

Moreover, this is almost optimal in the sense of Theorem 7: if

$$\alpha = \frac{4}{(\sqrt{L} + \sqrt{\gamma})^2}, \beta = \left(\frac{\sqrt{L} - \sqrt{\gamma}}{\sqrt{L} + \sqrt{\gamma}} \right)^2 \quad \text{then } q = \frac{\sqrt{L} - \sqrt{\gamma}}{\sqrt{L} + \sqrt{\gamma}}.$$

Proof: this is an example of a proof where one analyses the iteration of a linearized system near the optimum. Close enough to x^* , one has

$$x^{k+1} = x^k - \alpha D^2 f(x^*)(x^k - x^*) + o(\|x^k - x^*\|) + \beta(x^k - x^{k-1}),$$

and one can write that $z^k = (x^k - x^*, x^{k-1} - x^*)^T$ satisfies, for $B = D^2 f(x^*)$,

$$z^{k+1} = \begin{pmatrix} (1 + \beta)I - \alpha B & -\beta I \\ I & 0 \end{pmatrix} z^k + o(z^k).$$

We study the eigenvalues of the matrix A which appears in this iteration: We have

$$A \begin{pmatrix} x \\ y \end{pmatrix} = \begin{pmatrix} (1 + \beta)I - \alpha B & -\beta I \\ I & 0 \end{pmatrix} \begin{pmatrix} x \\ y \end{pmatrix} = \rho \begin{pmatrix} x \\ y \end{pmatrix}$$

if and only if

$$(1 + \beta)x - \alpha Bx - \beta y = \rho x, \quad x = \rho y$$

(and $x, y \neq 0$) hence if $(1 + \beta)x - \alpha Bx - \beta/\rho x = \rho x$. We find that

$$Bx = \frac{1}{\alpha} \left(1 + \beta - \rho - \frac{\beta}{\rho} \right) x$$

hence $\frac{1}{\alpha} \left(1 + \beta - \rho - \frac{\beta}{\rho} \right) = \mu \in [\gamma, L]$ is an eigenvalue of B . We derive the equation

$$\rho^2 - (1 + \beta - \alpha\mu)\rho + \beta = 0$$

which gives two eigenvalues with product β and sum $1 + \beta - \alpha\mu$. If $\beta \in [0, 1]$ and $-(1 + \beta) < 1 + \beta - \alpha\mu < (1 + \beta)$ (extreme cases where $\pm(1, \beta)$ are solutions) then $|\rho| < 1$, that is, if $0 < \alpha < (2 + \beta)/\mu$. Since $\mu < L$ one deduces that if $0 \leq \beta < 1$, $0 < \alpha < (2 + \beta)/L$, the eigenvalues of A are all in $(-1, 1)$ (incidentally, it has $2n$ eigenvalues).

We use here the following fundamental classical lemma [21]:

Lemma 2.13. *Let A be a $N \times N$ matrix and assume that all its eigenvalues (complex or real) have modulus $\leq \rho$. Then for any $\rho' > \rho$, there exists a norm $\|\cdot\|_*$ in \mathbb{C}^N such that $\|A\|_* := \sup_{\|\xi\|_* \leq 1} \|A\xi\|_* < \rho'$.*

This is an important result of linear algebra. The proof is as follows: up to a change of a basis, A is triangular: there exists P such that

$$P^{-1}AP = T$$

with $T = (t_{i,j})_{i,j}$, $t_{i,i} = \lambda_i$, an eigenvalue, and $t_{i,j} = 0$ if $i > j$. Then, if $D_s = \text{diag}(s, s^2, s^3, \dots, s^N) = (s^i \delta_{i,j})_{i,j}$, $D_s P^{-1} A P D_s^{-1} = (x_{i,j}^s)$ with

$$x_{i,j}^s = \sum_{k,l} s^i \delta_{i,k} t_{k,l} s^{-l} \delta_{l,j} = s^{i-j} t_{i,j}$$

and (since $t_{i,j} = 0$ for $i > j$), $x_{i,j}^s \rightarrow \lambda_i \delta_{i,j}$ as $s \rightarrow +\infty$.

Hence, if s is large enough, denoting $\|\xi\|_\infty = \max_i |\xi_i|$ the ∞ -norm,

$$\max_{\|\xi\|_\infty \leq 1} \|D_s P^{-1} A P D_s^{-1} \xi\|_\infty \leq \max_i (|\lambda_i| + (\rho' - \rho)) \leq \rho'$$

if s is large. Hence, if $\|\xi\|_* := \|D_s P^{-1} \xi\|_\infty$, one has

$$\|A\|_* = \sup_{\|\xi\|_* \leq 1} \|A\xi\|_* \leq \rho'.$$

It follows, in particular, that if $\rho' < 1$, $\|A^k\|_* \leq \|A\|_*^k \leq \rho'^k \rightarrow 0$ as $k \rightarrow \infty$. Applying this to our problem, we see that (choosing $\rho' < 1$)

$$\|z^{k+1}\|_* = \|Az^k + o(z^k)\|_* \leq (\rho' + \varepsilon) \|z^k\|_*$$

if $\|z^k\|_*$ is small enough. Starting from z^0 such that this holds for ε with $\rho' + \varepsilon < 1$, we find that it holds for all $k \geq 0$ and that $\|z^{k+1}\|_* \leq (\rho' + \varepsilon)^k \|z^0\|_*$, showing the linear convergence.

2.4.2 The conjugate gradient method

(For this section we refer again to Polyak [33].)

The conjugate gradient is “the best” two-steps method, in the sense that it can be defined as follows: given x^k, x^{k-1} , we let $x^{k+1} = x^k - \alpha_k \nabla f(x^k) + \beta_k (x^k - x^{k-1})$ where α_k, β_k are minimizing

$$\min_{\alpha, \beta} f(x^k - \alpha \nabla f(x^k) + \beta (x^k - x^{k-1})).$$

In particular, we deduce that

$$\langle \nabla f(x^{k+1}), \nabla f(x^k) \rangle = 0 \quad \text{and} \quad \langle \nabla f(x^{k+1}), x^k - x^{k-1} \rangle = 0 \quad (8)$$

and it also follows

$$\langle \nabla f(x^{k+1}), x^{k+1} - x^k \rangle = 0. \quad (9)$$

Notice moreover that

$$\begin{aligned} \nabla f(x^{k+1}) &= \nabla f(x^k) - \alpha_k D^2 f(x^k + s(x^{k+1} - x^k)) \nabla f(x^k) \\ &\quad + \beta_k D^2 f(x^k + s(x^{k+1} - x^k)) (x^k - x^{k-1}) \end{aligned} \quad (10)$$

for some $s \in [0, 1]$.

However this method is in general “conceptual”, meaning that one cannot hope to efficiently evaluate the values α_k, β_k and hence the new point x_k , except when f is quadratic: $f(x) = (1/2) \langle Ax, x \rangle - \langle b, x \rangle + c$ (A symmetric). Denoting then the gradients $p^k = Ax^k - b$ and the residuals $r^k = x^k - x^{k-1}$, we find that (cf (10))

$$p^{k+1} = p^k - \alpha_k Ap^k + \beta_k Ar^k \quad (11)$$

and using the orthogonality formulas (8),

$$0 = \|p^k\|^2 - \alpha_k \langle Ap^k, p^k \rangle + \beta_k \langle Ar^k, p^k \rangle, \quad 0 = \langle p^k, r^k \rangle - \alpha_k \langle Ap^k, r^k \rangle + \beta_k \langle Ar^k, r^k \rangle$$

we can compute explicitly the values of α_k, β_k (exercise).

Lemma 2.14. *The gradients (p^i) are all orthogonal.*

Proof: we start from $x^{k+1} = x^k - \alpha_k p^k + \beta_k (x^k - x^{k-1})$ and deduce (since ∇f is affine, or simply from (11))

$$p^{k+1} = p^k - \alpha_k Ap^k + \beta_k (p^k - p^{k-1}).$$

Assume that (p^0, \dots, p^i) are orthogonal, and that $\alpha_l, l = 0, \dots, i-1$, do not vanish (or we have found the solution, why?). Then

$$\langle Ap^k, p^l \rangle = \frac{1}{\alpha_k} \langle p^k - p^{k+1} + \beta_k (p^k - p^{k-1}), p^l \rangle = 0$$

if $l \leq k-2, k \leq i-1$ or if $i \geq l \geq k+2$. In particular, $\langle Ap^k, p^i \rangle = 0$ if $k \leq i-2$. Hence:

$$\langle p^{i+1}, p^k \rangle = \langle p^i, p^k \rangle - \alpha_k \langle Ap^i, p^k \rangle + \beta_k \langle p^i - p^{i-1}, p^k \rangle = 0 \quad (12)$$

if $k \leq i-2$. It remains therefore to check that $\langle p^{i-1}, p^{i+1} \rangle = 0$ and $\langle p^i, p^{i+1} \rangle = 0$. The latter is already known (8), hence we are left with the case $k = i-1$. If $k = i-1$: we use again $x^{k+1} = x^k - \alpha_k p^k + \beta_k (x^k - x^{k-1})$ to derive (with $r^0 = 0$)

$$r^{k+1} = -\alpha_k p^k + \beta_k r^k$$

so that $\forall k, r^k \in \text{vect} \{p^0, \dots, p^{k-1}\}$. Knowing (8) that $\langle p^{i+1}, r^i \rangle = 0$, one obtains from the previous (for $k = i-1$):

$$0 = -\alpha_{i-1} \langle p^{i+1}, p^{i-1} \rangle + \beta_{i-1} \langle p^{i+1}, r^{i-1} \rangle = -\alpha_{i-1} \langle p^{i+1}, p^{i-1} \rangle,$$

where we have used that $p^{i+1} \perp \text{vect} \{p^0, \dots, p^{i-2}\} \ni r^{i-1}$. This shows that $\langle p^{i+1}, p^{i-1} \rangle = 0$. Hence (p^0, \dots, p^{i+1}) are orthogonal. This holds as long as x^{i+1} is not a solution (then $p^{i+1} = 0$).

Corollary 2.15. *The solution is found in $k = \text{rk } A$ iterations.*

Indeed, if $p^{k+1} \neq 0$ then $p^i = Ax^i - b, i = 0, \dots, k+1$ are $k+2$ orthogonal vectors in $\text{Im}A - b$ which is an affine space of dimension k and contains at most $k+1$ independent points. One remarkable point is that also the directions r_i satisfy an orthogonality conditions: they are A -orthogonal: $\langle Ar_i, r_j \rangle = 0$ for all $i \neq j$, hence the name “conjugate directions”.

Variants One can show that the following rules defines the same points (for quadratic functions)

$$\begin{cases} p^k = \nabla f(x^k) \\ \beta_k = \frac{\|p^k\|^2}{\|p^{k-1}\|^2} \\ r^k = -p^k + \beta_k r^{k-1} \\ \alpha_k = \arg \min_{\alpha \geq 0} f(x^k + \alpha p^k), x^{k+1} = x^k + \alpha_k p^k \end{cases} \quad (\beta_0 = 0)$$

A variant replaces the 2nd line with $\beta_k = \langle p^k, p^k - p^{k-1} \rangle / \|p^{k-1}\|^2$. If f not quadratic, these variants can be implemented.

Optimality The conjugate gradient computes x^k as the minimum of f in the space generated by the orthogonal gradients (p^0, \dots, p^k) . It is then possible to prove that for a strongly convex quadratic function, that is if $\gamma I \leq A \leq LI$, then

$$\|x^k - x^*\| \leq 2\sqrt{Q}q^k \|x^0 - x^*\|$$

with $q = (\sqrt{Q} - 1)/(\sqrt{Q} + 1)$, $Q = L/\gamma$ the condition number. This is the same rate as the Heavy-Ball.

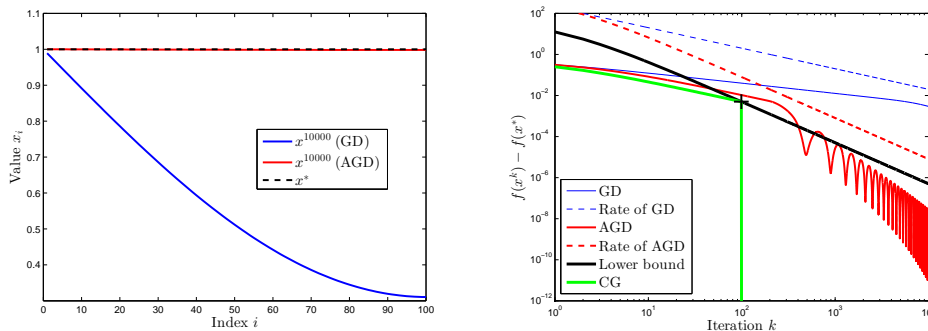


Figure 2: Comparison between accelerated vs non-accelerated gradient schemes. Top: Comparisons of the solutions x of GD and AGD after 10000(!) iterations. Bottom: Rate of convergence for GD, AGD together with their theoretical worst case rates, and the lower bound for smooth optimization. For comparison we also provide the rate of convergence for CG. Note that CG exactly touches the lower bound at $k = 99$ (problem (4) with $\gamma = 0$, $p = n = 100$)

2.4.3 Accelerated algorithm: Nesterov 83

We rapidly mention the “Accelerated Gradient Descent” (AGD) Algorithm by Yu. Nesterov [28].

Algorithm: $x^0 = x^{-1}$ given, x^{k+1} defined by:

$$\begin{cases} y^k = x^k + \frac{t_k - 1}{t_{k+1}}(x^k - x^{k-1}) \\ x^{k+1} = y^k - \tau \nabla f(y^k) \end{cases}$$

where $\tau = 1/L$ and for instance $t_k = 1 + k/2$. Then,

$$f(x^k) - f(x^*) \leq \frac{2L}{(k+1)^2} \|x^0 - x^*\|^2$$

We will prove this later in these notes. For strongly convex problems, a variant exists with again “optimal” rate of convergence.

2.5 Nonsmooth problems?

2.5.1 Subgradient descent

The first basic approach to tackle nonsmooth problems (or more generally problems where the (local) Lipschitz constant of the gradient is unknown and possibly rapidly varying) is called a “subgradient descent”. The idea, given f convex, is to iterate:

$$x^{k+1} = x^k - h_k \frac{\nabla f(x^k)}{\|\nabla f(x^k)\|}.$$

In practice, the gradient here can be replaced with any selection of the subgradient if f is not differentiable at x^k , see Section 4.1 for the technical details.

Then if x^* is a solution,

$$\begin{aligned} \|x^{k+1} - x^*\|^2 &= \|x^k - x^*\|^2 - 2 \frac{h_k}{\|\nabla f(x^k)\|} \langle \nabla f(x^k), x^k - x^* \rangle + h_k^2 \\ &\leq \|x^k - x^*\|^2 - 2 \frac{h_k}{\|\nabla f(x^k)\|} (f(x^k) - f(x^*)) + h_k^2. \end{aligned}$$

We have used here a basic property of convex functions, which is that they are above their affine approximations, so that $f(x^*) \geq f(x^k) + \langle \nabla f(x^k), x^* - x^k \rangle$.

Hence, assuming in addition f is M -Lipschitz (near x^* at least)

$$\min_{0 \leq i \leq k} f(x^i) - f(x^*) \leq M \frac{\|x^0 - x^*\|^2 + \sum_{i=0}^k h_i^2}{2 \sum_{i=0}^k h_i}$$

and choosing $h_i = C/\sqrt{k+1}$ for k iterations, we obtain

$$\min_{0 \leq i \leq k} f(x^i) - f(x^*) \leq M \frac{C^2 + \|x^0 - x^*\|^2}{2C\sqrt{k+1}}$$

(the best choice is $C \sim \|x^0 - x^*\|$ but this is of course unknown).

In general, one chooses steps such that $\sum_i h_i^2 < +\infty$, $\sum_i h_i = +\infty$, such as $h_i = 1/i$. It results in a very slowly converging algorithm which should be used only when there is no other obvious choice.

2.5.2 Implicit descent

Consider a gradient descent where instead of using the gradient at x^k , one is able to evaluate the gradient in x^{k+1} :

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

This is of course often “conceptual”, however we will see that in many instances it can be computed or approximated. It says that x^{k+1} is a critical point of (and one can ask that it minimises)

$$f(x) + \frac{1}{2\tau} \|x - x^k\|^2.$$

Observe that if one lets

$$f_\tau(x) := \min_y f(y) + \frac{1}{2\tau} \|y - x\|^2 \quad (13)$$

(this defines an “inf-convolution”) which is well-defined if f is bounded from below (or $\geq -\alpha\|x\|^2$ and $\tau < 1/\alpha$) and lower-semicontinuous (if not, the min has to be replaced with an inf), then one can show that f_τ is semi-concave and when differentiable, $\nabla f_\tau(x) = (x - y_x)/\tau$ where y_x solves (13) (and is thus, in this case, unique).

Proof: equivalently, one may observe that:

$$\begin{aligned} & f_\tau(x - h) - 2f_\tau(x) + f_\tau(x + h) \\ & \leq f(y_x) + \frac{1}{2\tau} \|x - h - y_x\|^2 - 2f(y_x) - \frac{1}{\tau} \|x - y_x\|^2 + f(y_x) + \frac{1}{2\tau} \|x + h - y_x\|^2 \leq \frac{1}{\tau} \|h\|^2 \end{aligned}$$

showing that $f_\tau(x) - \|x\|^2/(2\tau)$ is concave; or more directly one observes that:

$$f_\tau(x) - \frac{1}{2\tau} \|x\|^2 = \min_y f(y) + \frac{1}{2\tau} \|y\|^2 - \langle x, y \rangle$$

is a concave function as an inf of linear functions (of the variable x). This shows that f_τ is $(1/\tau)$ -“semi-concave”. Hence f_τ is differentiable a.e. (even twice, Aleksandrov’s theorem [16]), and if $\nabla f_\tau(x)$ exists, one has

$$\begin{aligned} f_\tau(x + h) & \leq f(y_x) + \frac{1}{2\tau} \|x + h - y_x\|^2, \text{ hence} \\ f_\tau(x + h) - f_\tau(x) & \leq \frac{1}{\tau} \langle x - y_x, h \rangle + \frac{\|h\|^2}{2\tau}, \end{aligned}$$

so that for all h ,

$$\nabla f_\tau(x) \cdot h \leq \frac{1}{\tau} \langle x - y_x, h \rangle$$

showing the claim. Then, $y_x = x - \tau \nabla f_\tau(x)$.

Conversely, if y_x is unique, then $\nabla f_\tau(x)$ exists and is $(x - y_x)/\tau$. This follows from the observation that if $x_n \rightarrow x$ and y_{x_n} is a minimizer for x_n , as

$$f(y_{x_n}) + \frac{1}{2\tau} \|x_n - y_{x_n}\|^2 \leq f(y_x) + \frac{1}{2\tau} \|x_n - y_x\|^2$$

showing that (f being bounded from below) (y_{x_n}) is a bounded sequence. If $(y_{x_{n_k}})$ is a subsequence which converges to some \bar{y} passing to the limit in

$$f(y_{x_{n_k}}) + \frac{1}{2\tau} \|x_{n_k} - y_{x_{n_k}}\|^2 \leq f(y) + \frac{1}{2\tau} \|x_{n_k} - y\|^2$$

and using the semi-continuity of f , we find that \bar{y} is a minimizer for x , hence $\bar{y} = y_x$ and $y_{x_{n_k}} \rightarrow y_x$: the multivalued mapping $x \mapsto y_x$ is thus continuous at points where the argument is unique. Now, we can write that

$$f_\tau(x + h) \leq f_\tau(x) + \frac{1}{\tau} \langle x - y_x, h \rangle + \frac{\|h\|^2}{2\tau}$$

and in the same way (exchanging x and $x + h$)

$$\begin{aligned} f_\tau(x) &\leq f_\tau(x + h) - \frac{1}{\tau} \langle x + h - y_{x+h}, h \rangle + \frac{\|h\|^2}{2\tau} \\ &= f_\tau(x + h) - \frac{1}{\tau} \langle x - y_{x+h}, h \rangle - \frac{\|h\|^2}{2\tau} \end{aligned}$$

hence for $t > 0$, small:

$$\frac{1}{\tau} \langle x - y_{x+th}, h \rangle \leq \frac{f_\tau(x + th) - f_\tau(x)}{t} \leq \frac{1}{\tau} \langle x - y_x, h \rangle + O(t)$$

and in the limit $t \rightarrow 0$ we recover the claim.

This proof is finite-dimensional, we will however see later on for convex functions in Hilbert spaces that the same result is true.

We find that

$$x^{k+1} = x^k - \tau \nabla f(x^{k+1}) \Leftrightarrow x^{k+1} = x^k - \tau \nabla f_\tau(x^k)$$

hence the implicit descent is an explicit descent on f_τ ! Which has the same minimisers. It converges to critical points of f_τ (as $D^2 f_\tau \leq I/\tau$), as before (and under the same assumptions). These are local minimizers of $f(\cdot) + \|\cdot - x\|^2/(2\tau)$.

Example 2.16 (Lasso problem). Consider:

$$\min_x \|x\|_1 + \frac{1}{2} \|Ax - b\|^2$$

If $\|x\|_M^2 = \langle Mx, x \rangle$ and $M = I/\tau - A^*A$, $\tau < 1/\|A\|^2$, then

$$\min_x \frac{1}{2} \|x - x^k\|_M^2 + \|x\|_1 + \frac{1}{2} \|Ax - b\|^2$$

is solved by

$$x^{k+1} = S_\tau(x^k - \tau A^*(Ax^k - b))$$

where $S_\tau \xi$ is the unique minimizer of

$$\min_x \|x\|_1 + \frac{1}{2\tau} \|x - \xi\|^2,$$

called the ‘‘shrinkage’’ operator. This converges with rate $O(1/k)$ to a solution.

3 Krasnoselskii-Mann’s convergence theorem

3.1 A ‘‘general’’ convergence theorem

We show here a general form of a convergence theorem of Krasnoselskii and Mann for the iterates of weak contractions (or nonexpansive operators) (it is found in all convex optimisation books, cf for instance [4, 2]). We state first a simple form. Consider (a priori, in a Banach space \mathcal{X}) an operator $T : \mathcal{X} \rightarrow \mathcal{X}$ which is 1-Lipschitz:

$$\|Tx - Ty\| \leq \|x - y\| \quad \forall x, y \in \mathcal{X}.$$

If in addition it is ρ -Lipschitz with $\rho < 1$, then Picard’s classical fixed point theorem shows that the iterates $x^k = T^k x^0$, $k \geq 1$, form a Cauchy sequence and therefore

converge to a fixed point, necessarily unique. This relies on the fact that the space is complete.

However, for $\rho = 1$, this does not always work ($\rho = 1$ does not provide much relevant information, as when $T = I$). For instance, if $Tx = -x$, there is only one fixed point but the iterates never converge, unless $x^0 = 0$. The simplest statement of Krasnoselskii-Mann's theorem shows that if T is *averaged* and has fixed points, then the iterates weakly converge to a fixed point. The statement is true in Hilbert (or finite-dimensional Euclidean) spaces, as well as in some class of reflexive Banach space. We assume in what follows that X is Hilbert, and will mention the changes and properties needed for the property to hold for more generality.

For $\theta \in]0, 1[$, we define the (θ -)averaged operator T_θ by letting

$$T_\theta x = (1 - \theta)x + \theta Tx.$$

We also let $T_0 = I$, $T_1 = T$, and $F = \{x \in \mathcal{X} : Tx = x\}$. Observe that for any $\theta \in]0, 1[$, F is the set of fixed point of T_θ .

Theorem 3.1. *Let $x \in \mathcal{X}$, $0 < \theta < 1$, and assume $F \neq \emptyset$. Then $(T_\theta^k x)_{k \geq 1}$ weakly converges to some point $x^* \in F$.*

The proof consists in four simple steps. We denote $x^0 = x$, $x^k = T^k x$, $k \geq 1$.

Step 1 First, since T_θ is 1-Lipschitz, then for any $x^* \in F$, $\|T_\theta x^k - x^*\| \leq \|x^k - x^*\|$ and the sequence $(\|x^k - x^*\|)_k$ is nonincreasing. $(x_k)_k$ is said to be “Fejér-monotone” with respect to F , see [2, Chap. 5] for details and interesting properties.

It follows that one can define $m(x^*) = \inf_k \|x^k - x^*\| = \lim_k \|x^k - x^*\|$. If there exists $x^* \in F$ such that $m(x^*) = 0$ then the theorem is proved (with strong convergence), otherwise we proceed to the next step. We will see later on what happens if the sequence is “quasi-Fejér-monotone”, which happens for instance if T is computed with errors. Hence we assume that $m(x^*) > 0$ for all $x^* \in F$.

Step 2 We now show that $x^{k+1} - x^k \rightarrow 0$ strongly. The operator T_τ is said to be “asymptotically regular”.

First, for \mathcal{X} a Hilbert space, the proof is a straightforward application of the parallelogram identity, to:

$$x^{k+1} - x^* = (1 - \theta)(x^k - x^*) + \theta(T_1 x^k - x^*).$$

We find that for all k :

$$\begin{aligned} \|x^{k+1} - x^*\|^2 &= (1 - \theta)\|x^k - x^*\|^2 + \theta\|T_1 x^k - x^*\|^2 - \theta(1 - \theta)\|T_1 x^k - x^k\|^2 \\ &\leq \|x^k - x^*\|^2 - \frac{1 - \theta}{\theta}\|x^{k+1} - x^k\|^2 \end{aligned}$$

from which one deduces that $\sum_k \|x^{k+1} - x^k\|^2 < \infty$, hence the result. In addition, one observes that the sequence $(1 - \theta)/\theta\|x^{k+1} - x^k\|^2$ (which is nonincreasing) is controlled in the following way:

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

As $x^{k+1} - x^k = \theta(T_1 x^k - x^k)$ we obtain a rate for the error $T_1 x^k - x^k$, in the Hilbertian setting, given by:

$$\|T_1 x^k - x^k\| \leq \frac{\|x^0 - x^*\|}{\sqrt{\theta(1-\theta)}\sqrt{k+1}}. \quad (14)$$

Now, we have to mention that the result also holds in more general spaces. It is easy to extend in *uniformly convex* spaces, meaning that the unit ball satisfies the following property:

Uniformly convex unit ball: $\forall \varepsilon > 0, \theta \in (0, 1), \exists \delta > 0$ such that for all $x, y \in \mathcal{X}$ with $\|x\| \leq 1, \|y\| \leq 1$ and $\|x - y\| \geq \varepsilon$,

$$\|\theta x + (1 - \theta)y\| \leq (1 - \delta) \max\{\|x\|, \|y\|\}$$

Of course, the following holds:

Lemma 3.2. *If \mathcal{X} is a Hilbert space then it is uniformly convex.*

Indeed the parallelogram identity yields:

$$\begin{aligned} \|\theta x + (1 - \theta)y\|^2 &= \theta^2 \|x\|^2 + (1 - \theta)^2 \|y\|^2 + 2\theta(1 - \theta) \langle x, y \rangle \\ &= \theta \|x\|^2 + (1 - \theta) \|y\|^2 - \theta(1 - \theta) \|x - y\|^2 \\ &\leq \max\{\|x\|^2, \|y\|^2\} - \theta(1 - \theta)\varepsilon^2 \leq (1 - \delta)^2 \max\{\|x\|^2, \|y\|^2\} \end{aligned}$$

for $\delta = 1 - \sqrt{1 - \theta(1 - \theta)\varepsilon^2}$, where we have used, of course, that $\|x\|, \|y\| \leq 1$.

Yet the same property also holds in many Banach spaces (such as L^p spaces, $0 < p < 1$, etc). It is well known (as ‘‘Milman-Pettis’’ theorem) that such a space is reflexive (while the converse is not true). We can prove the asymptotic regularity, that is, that $x^{k+1} - x^k \rightarrow 0$, only relying on this property, as follows: We recall that we assume $m(x^*) > 0$ for all $x^* \in F$. Assume that along a subsequence, one has $\|x^{k_l+1} - x^{k_l}\| \geq \varepsilon > 0$. Observe that

$$x^{k_l+1} - x^* = (1 - \theta)(x^{k_l} - x^*) + \theta(T_1 x^{k_l} - x^*)$$

and that

$$(x^{k_l} - x^*) - (T_1 x^{k_l} - x^*) = x^{k_l} - T_1 x^{k_l} = -\frac{1}{\theta}(x^{k_l+1} - x^{k_l})$$

so that $\|(x^{k_l} - x^*) - (T_1 x^{k_l} - x^*)\| \geq \varepsilon/\theta > 0$. Hence thanks to the uniform convexity of the ball (remember that $(x^k - x^*)_k$ is globally bounded since its norm is nonincreasing), we obtain that for some $\delta > 0$,

$$m(x^*) \leq \|x^{k_l+1} - x^*\| \leq (1 - \delta) \max\{\|x^{k_l} - x^*\|, \|T_1 x^{k_l} - x^*\|\}$$

but since $\|T_1 x^{k_l} - x^*\| \leq \|x^{k_l} - x^*\|$, it follows

$$m(x^*) \leq (1 - \delta) \|x^{k_l} - x^*\|.$$

As $k_l \rightarrow \infty$, we get a contradiction if $m(x^*) > 0$.

The result in [11] shows that asymptotic regularity holds in any normed space.

Step 3. Assume now that \bar{x} is the weak limit of some subsequence $(x^{k_l})_l$. Then, we claim it is a fixed point. We use Opial's lemma:

Lemma 3.3 ([31, Lem. 1]). *If in a Hilbert space \mathcal{X} the sequence $(x_n)_n$ is weakly convergent to x_0 then for any $x \neq x_0$,*

$$\liminf_n \|x_n - x\| > \liminf_n \|x_n - x_0\|$$

Proof of Opial's lemma (obvious): one has

$$\|x_n - x\|^2 = \|x_n - x_0\|^2 + 2\langle x_n - x_0, x_0 - x \rangle + \|x_0 - x\|^2.$$

Since $\langle x_n - x_0, x_0 - x \rangle \rightarrow 0$ by weak convergence, we deduce

$$\liminf_n \|x_n - x\|^2 = \liminf_n (\|x_n - x_0\|^2 + \|x_0 - x\|^2) = \|x_0 - x\|^2 + \liminf_n \|x_n - x_0\|^2$$

and the claim follows.

Proof that \bar{x} is a fixed point: since T_θ is a contraction, we observe that for each k ,

$$\begin{aligned} \|x^k - \bar{x}\| &\geq \|T_\theta x^k - T_\theta \bar{x}\| \\ &= \|x^{k+1} - x^k + x^k - T_\theta \bar{x}\| \geq \|x^k - T_\theta \bar{x}\| - \|x^{k+1} - x^k\| \end{aligned}$$

and we deduce (thanks to the previous Step 2):

$$\liminf_l \|x^{k_l} - \bar{x}\| \geq \liminf_l \|x^{k_l} - T_\theta \bar{x}\|.$$

Opial's lemma implies that $T_\theta \bar{x} = \bar{x}$.

One advantage of this approach is that it can be extended to Banach spaces [31] where "Opial's property" (the statement of the Lemma) holds (in the norm for which T is a contraction). On the other hand, not all spaces satisfy this property (it is shown that in separable Banach spaces, there is an equivalent norm for which the property is true [44], but this is useless if T is not nonexpansive for this norm...)

Remark 3.4. Another classical approach in Hilbert spaces to prove this claim is to use "Minty's trick" to study the limit of "monotone" operators: Since T_θ is a contraction, for each $y \in \mathcal{X}$ we have (thanks to Cauchy-Schwarz's inequality)

$$\langle (I - T_\theta)x_{n_k} - (I - T_\theta)y, x_{n_k} - y \rangle \geq 0$$

and as we have just proved that $(I - T_\theta)x_{n_k} \rightarrow 0$ (strongly), then

$$\langle -(I - T_\theta)y, \bar{x} - y \rangle \geq 0.$$

Choose $y = \bar{x} + \varepsilon z$ for $z \in \mathcal{X}$ and $\varepsilon > 0$: it follows after dividing by ε that

$$\langle (I - T_\theta)(\bar{x} + \varepsilon z), z \rangle \geq 0.$$

and since T_θ is Lipschitz, sending $\varepsilon \rightarrow 0$ we recover $\langle (I - T_\theta)\bar{x}, z \rangle \geq 0$ for any z , which shows that $\bar{x} \in F$.

Step 4. To conclude, assume that a subsequence $(x^{m_l})_l$ of $(x^k)_k$ converges weakly to another fixed point \bar{y} . Then it must be that $\bar{y} = \bar{x}$, otherwise Opial's lemma 3.3 again would imply both that $m(\bar{x}) < m(\bar{y})$ and $m(\bar{y}) < m(\bar{x})$:

$$m(\bar{y}) = \liminf_l \|x^{m_l} - \bar{y}\| < \liminf_l \|x^{m_l} - \bar{x}\| = m(\bar{x}).$$

It follows that the whole sequence (x^k) must weakly converge to \bar{x} .

3.2 Varying steps

One can consider more generally iterations of the form

$$x^{k+1} = x^k + \tau_k(T_1 x^k - x^k)$$

with varying steps τ_k . Then, if $0 < \underline{\tau} \leq \tau_k \leq \bar{\tau} < 1$, the convergence still holds, with almost the same proof. (This is obvious in the Hilbertian setting, cf Step. 2.)

Remark 3.5. A sufficient condition is that $\sum_k \tau_k(1 - \tau_k) = \infty$, see [34]. In addition, a slight improvement to the proof in Step. 2 shows that

$$\sum_{i=0}^k (1 - \tau_i) \tau_i \|T_1 x^i - x^i\|^2 \leq \|x^0 - x^*\|^2 - \|x^{k+1} - x^*\|^2$$

so that $\min_{0 \leq i \leq k} \|T_1 x^i - x^i\| \leq \|x^0 - x^*\| / \sqrt{\sum_{i=0}^k (1 - \tau_i) \tau_i}$. In fact, in a general normed space one has the estimate

$$\|T_1 x^k - x^k\| \leq \frac{1}{\sqrt{\pi}} \frac{\|x^0 - x^*\|}{\sqrt{\sum_{i=0}^k \tau_i (1 - \tau_i)}}$$

which improves (14), see [11].

3.3 A variant with errors

Assume now the sequence (x_k) is an inexact iteration of T_θ :

$$\|x^{k+1} - T_\theta x^k\| \leq \varepsilon_k.$$

Then one has the following result:

Theorem 3.6 (Variant of Thm 3.1). *If $\sum_k \varepsilon_k < \infty$, then $x^k \rightarrow \bar{x}$ a fixed point of T (if one exists).*

Proof: now, x^k is “quasi-Fejér monotone”: denoting $e_k = x^{k+1} - T_\theta x^k$ so that $\|e_k\| \leq \varepsilon_k$,

$$\|x^{k+1} - x^*\| = \|T_\theta x^k - T_\theta x^* + e_k\| \leq \|x^k - x^*\| + \varepsilon_k$$

for all k , and any $x^* \in F$. Hence, $\|x^{k+1} - x^*\| \leq \|x^0 - x^*\| + \sum_{i=0}^k \varepsilon_i$ is bounded. Letting $a_k = \sum_{i=k}^{\infty} \varepsilon_i$ which is finite and goes to 0 as $k \rightarrow \infty$, this can be rewritten

$$\|x^{k+1} - x^*\| + a_{k+1} \leq \|x^k - x^*\| + a_k$$

so that once more one can define

$$m(x^*) := \lim_{k \rightarrow \infty} \|x^k - x^*\| = \inf_{k \geq 0} \|x^k - x^*\| + a_k$$

Again, if $m(x^*) = 0$ the theorem is proved, otherwise, one can continue the proof as before: now,

$$x^{k_l+1} - x^* = (1 - \theta)(x^{k_l} - x^* + e_{k_l}) + \theta(T_1 x^{k_l} - x^* + e_{k_l})$$

while

$$(x^{k_l} - x^* + e_{k_l}) - (T_1 x^{k_l} - x^* + e_{k_l}) = x^{k_l} - T_1 x^{k_l} = -\frac{1}{\theta}(x^{k_l+1} - x^{k_l} - e_{k_l})$$

so that $\|(x^{k_l} - x^*) - (T_1 x^{k_l} - x^*)\| \geq (\varepsilon - \varepsilon_{k_l})/\theta > \varepsilon/(2\theta) > 0$ if l is large enough, and one can invoke again Lemma 3.2 to find that

$$\begin{aligned} m(x^*) \leq \|x^{k_l+1} - x^*\| &\leq (1 - \delta) \max\{\|x^{k_l} - x^* + e_{k_l}\|, \|T_1 x^{k_l} - x^* + e_{k_l}\|\} \\ &\leq (1 - \delta) (\|x^{k_l} - x^*\| + \varepsilon_{k_l}) \end{aligned}$$

and again sending $l \rightarrow \infty$ we obtain that $m(x^*) \leq (1 - \delta)m(x^*)$, a contradiction if $m(x^*) > 0$. The rest of the proof (steps 3, 4) is almost identical.

Remark 3.7. In practice, what do you think about the condition $\sum_k \varepsilon_k < \infty$?

3.4 Examples

3.4.1 Gradient descent

It follows the convergence for the explicit and implicit gradient descent for convex functions. Consider indeed the iteration $x^{k+1} = T_\tau(x^k) := x^k - \tau \nabla f(x^k)$, for f convex with L -Lipschitz gradient. Then, Lemma 2.4 claims that

$$T_{2/L}(x) = x - \frac{2}{L} \nabla f(x)$$

is a weak contraction (1-Lipschitz or “nonexpansive” operator).

We observe that if $0 < \tau < 2/L$, one has

$$T_\tau(x) = x - \frac{\tau L}{2} \frac{2}{L} \nabla f(x) = \frac{\tau L}{2} T_{2/L}(x) + \left(1 - \frac{\tau L}{2}\right) x$$

is an averaged operator (with here $\theta = L\tau/2 \in]0, 1[$). Theorem 3.1 yields the convergence of the iterates. Moreover, one still has convergence if one uses varying steps τ_k with $0 < \inf_k \tau_k \leq \sup_k \tau_k < 2/L$. One can also consider (summable) errors. Eventually, thanks to 14, one has the rate

$$\left\| \frac{2}{L} \nabla f(x^k) \right\| \leq \frac{\|x^0 - x^*\|}{\sqrt{(1 - L\tau/2)L\tau/2\sqrt{k+1}}}.$$

(Compare this with (2), Theorem 2.7, Remark 2.2.)

For the implicit descent, we can use the fact that it is an explicit descent on the function f_τ , which has $1/\tau$ -Lipschitz gradient, to get a similar result: Let $x^{k+1} = x^k - \lambda \nabla f_\tau(x^k) = x^k + (\lambda/\tau)(y_{x^k} - x^k)$ (where y_x solves (13)) for $0 < \lambda < 2\tau$, then x^k converges (weakly) to a minimizer of f_τ (which is also a minimizer of f)...

3.4.2 Composition of averaged operators

An important remark is the following: Let T_θ, S_λ be averaged operators: $T_\theta = (1 - \theta)I + \theta T_1$, $S_\lambda = (1 - \lambda)I + \lambda S_1$. Then $T_\theta \circ S_\lambda$ is also averaged: letting $\mu = \theta + \lambda(1 - \theta) \in]0, 1[$, one has

$$T_\theta \circ S_\lambda = (1 - \mu)I + \mu \frac{(1 - \theta)\lambda S_1 + \theta T_1 \circ ((1 - \lambda)I + \lambda S_1)}{\theta + (1 - \theta)\lambda}.$$

An important application is the following: consider the problem

$$\min_x f(x) + g(x) \tag{15}$$

where f, g convex, lsc, f has L -Lipschitz gradient and g is such that one knows how to compute, for all y and all $\tau > 0$:

$$g_\tau(x) := \min_y g(y) + \frac{1}{2\tau} \|x - y\|^2. \tag{16}$$

Then one can compose the averaged operators

$$T_\tau x := x - \tau \nabla f(x),$$

$0 < \tau < 2/L$, and

$$S_\tau x := y_x$$

which solves (16) (and is $(1/2)$ -averaged, as it is $x - \tau \nabla g_\tau(x)$ where ∇g_τ is $1/\tau$ -Lipschitz). Hence, if one defines the iterates $x^{k+1} := S_\tau \circ T_\tau x^k$, $k \geq 0$, then $x^k \rightharpoonup x^*$ (weakly) where x^* is a fixed point if $S_\tau \circ T_\tau$. As $S_\tau x$ satisfies

$$\nabla g(S_\tau x) + \frac{1}{\tau}(S_\tau x - x) = 0,$$

one has

$$\begin{aligned} 0 &= \nabla g(S_\tau(T_\tau x^*)) + \frac{1}{\tau}(S_\tau(T_\tau x^*) - T_\tau x^*) \\ &= \nabla g(x^*) + \frac{1}{\tau}(x^* - (x^* - \tau \nabla f(x^*))) = \nabla g(x^*) + \nabla f(x^*) \end{aligned}$$

so that x^* is a minimizer of (15). We deduce the following:

Theorem 3.8. *The iterates of the “forward-backward” algorithm $x^{k+1} := S_\tau \circ T_\tau x^k$ weakly converge to a minimizer of (15).*

We will see later on that one can say much more about this approach. Compare this with Example 2.16.

Remark 3.9. What about the “explicit-explicit” (“forward-forward”) iteration

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

with $\tau < \min\{2/L_f, 2/L_g\}$ where L_f, L_g are the Lipschitz constants of the gradients of f, g , respectively?

We will see later on other useful examples of composition of averaged operators.

4 An introduction to convex analysis and monotone operators

Most of this section is in Hilbert spaces, though many results are also valid in more general vector spaces, but often with more involved proofs.

4.1 Convexity

See for instance: [38, 15] for a general introduction. We discuss here the following notions: Convex function; Subgradients; Inf-convolution; Sums of subgradients; Convex Conjugate (Legendre-Fenchel); Fenchel-Rockafellar duality; Moreau-Yosida's regularization (inf-convolution); Moreau's identity.

4.1.1 Convex functions

An extended-valued function $f : \mathcal{X} \rightarrow [-\infty, +\infty]$ is said to be *convex* if and only if its *epigraph*

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

is a convex set, that is, if when $\lambda \geq f(x)$, $\mu \geq f(y)$, and $t \in [0, 1]$, one has $t\lambda + (1-t)\mu \geq f(tx + (1-t)y)$.³ It is *proper* if it is not identically $+\infty$ and nowhere $-\infty$: in this case, it is convex if and only if for all $x, y \in \mathcal{X}$ and $t \in [0, 1]$,

$$f(tx + (1-t)y) \leq tf(x) + (1-t)f(y).$$

It is *strictly convex* if the above inequality is strict whenever $x \neq y$ and $0 < t < 1$. It is *strongly convex* (or μ -convex) if in addition, there exists $\mu > 0$ such that for all $x, y \in \mathcal{X}$ and $t \in [0, 1]$,

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

Thanks to the parallelogram identity, in the Hilbertian setting, one easily checks that this is equivalent to require that $x \mapsto f(x) - \mu/2 \|x\|^2$ is still convex. The function is also said to be, in this case, " μ -convex". The archetypical example of a μ -convex function is a quadratic plus affine function $\mu \|x\|^2/2 + \langle b, x \rangle + c$.

The *domain* of a proper convex function f is the set $\text{dom } f = \{x \in \mathcal{X} : f(x) < +\infty\}$. It is obviously a convex set.

We say that f is *lower semi-continuous* (l.s.c.) if for all $x \in \mathcal{X}$, if $x_n \rightarrow x$, then

$$f(x) \leq \liminf_{n \rightarrow \infty} f(x_n).$$

It is easy to see that f is l.s.c. if and only if $\text{epi } f$ is closed.

A trivial but important example is the *characteristic function* or *indicator function* of a set C :

$$\delta_C(x) = \begin{cases} 0 & \text{if } x \in C, \\ +\infty & \text{else,} \end{cases}$$

which is convex, l.s.c., and proper as soon as C is convex, closed and nonempty. The minimisation of such functions will allow to easily model convex constraints in our problems.

One can show the following result:

Lemma 4.1. *If there exists $B \subset \text{dom } f$ an open ball where the proper convex function f is bounded from above, then f is locally Lipschitz in the interior of $\text{dom } f$. In finite dimension, a proper convex function f is locally Lipschitz in the relative interior of $\text{dom } f$, $\text{ri dom } f$.*

³This definition avoids the embarrassing expression $(+\infty) + (-\infty)$, see for instance [38, Sec. 4].

In finite dimension, the relative interior is defined as the interior of $\text{dom } f$ in the space $x + \text{vect}(\text{dom } f - x)$ for any $x \in \text{dom } f$; this is never empty (but may have, in extreme cases, dimension zero).

Proof of the lemma: we assume that $B = B(0, \delta)$, $\delta > 0$, and let $M = \sup_B f < \infty$. Observe also that for $x \in B$, by convexity $f(x) \geq 2f(0) - f(-x) \geq 2f(0) - M$ so that $|f| \leq M + 2|f(0)|$. We prove that f is Lipschitz in $B(0, \delta/2)$: indeed, if $x, y \in B(0, \delta/2)$, there is $z \in B(0, \delta)$ such that $y = (1-t)x + tz$ for some $t \in [0, 1]$, and $\|z - x\| \geq \delta/2$. In particular by convexity, $f(y) - f(x) \leq t(f(z) - f(x)) \leq t2(M - f(0))$. Now, $t(z - x) = y - x$ so that $t \leq \|y - x\|/\|z - x\| \leq 2\|y - x\|/\delta$ hence: $f(y) - f(x) \leq (4(M - f(0))/\delta)\|y - x\|$ which shows the claim (one could show in fact in the same way that f is Lipschitz in any ball contained in $B(0, \delta)$).

Now, let \bar{x} in the interior of $\text{dom } f$. Observe that for some $\lambda > 1$, $\lambda\bar{x} \in \text{dom } f$ and as a consequence $B' = 1/\lambda(\lambda\bar{x}) + (1 - 1/\lambda)B(0, \delta) = B(\bar{x}, \delta(1 - 1/\lambda)) \subset \text{dom } f$; moreover, if $x \in B'$, $x = 1/\lambda(\lambda\bar{x}) + (1 - 1/\lambda)z$ for some z with $f(z) \leq M$ hence $f(x) \leq 1/\lambda f(\lambda\bar{x}) + (1 - 1/\lambda)M$, so that $\sup_{B'} f < \infty$. Hence as before f is Lipschitz in a smaller ball.

In finite dimension, assume $0 \in \text{dom } f$ and let d be the dimension of $\text{vect dom } f$. It means there exist x_1, \dots, x_d independent points in $\text{dom } f$. Now, the d -dimensional set $\{\sum_i t_i x_i : t_i > 0, \sum_i t_i \leq 1\}$ (the interior of the convex envelope of $\{0, x_1, \dots, x_d\}$) is an open set in $\text{vect dom } f$, moreover if $x = \sum_i t_i x_i$, $f(x) \leq \sum_i t_i f(x_i) + (1 - \sum_i t_i)f(0) \leq M := \max\{f(0), f(x_1), \dots, f(x_d)\}$. Hence we can apply the first part of the theorem, and f is locally Lipschitz in the relative interior of the domain.

Remark 4.2. Note that in infinite dimension one can possibly find noncontinuous linear forms⁴ hence noncontinuous convex functions. However, one can show that a convex proper *lower semi-continuous* function is always locally bounded in the interior of its domain, and therefore locally Lipschitz (as if 0 is an interior point and one considers the convex closed set $C = \{x : f(x) \leq 1 + f(0)\}$, one can check that $\cup_{n \geq 1} nC = \mathcal{X}$, as if $x \in \mathcal{X}$, $t \mapsto f(tx)$ is locally Lipschitz near $t = 0$. Hence $\overset{\circ}{C} \neq \emptyset$ by Baire's property: it follows that there is an open ball where f is bounded, as requested), cf [15, Cor. 2.5].

4.1.2 Separation of convex sets

In this section we establish two important “separation” theorems for convex sets, which are geometric variants of Hahn-Banach's theorem, in the particular setting of Hilbert spaces. In this setting, unlike in the general case, these are quite obvious results.

Theorem 4.3. *Let \mathcal{X} be a (real) Hilbert space, $C \subset \mathcal{X}$ a closed, convex set and $x \notin C$. Then there exists a closed hyperplane which “separates” strictly x and C : precisely, in the Hilbertian setting, one can find $v \in X, \alpha \in \mathbb{R}$ such that*

$$\langle v, x \rangle > \alpha \geq \langle v, y \rangle \quad \forall y \in C$$

Proof: introduce the projection $z = \Pi_C(x)$ defined by $\|x - z\| = \min_{y \in C} \|x - y\|$ (existence is classically shown by proving that any minimizing sequence is a Cauchy sequence, thanks to the parallelogram identity [or strong convexity of $\|x - \cdot\|^2$]). The first order optimality condition for z is found by writing that for any $y \in C$, $\|x - z\|^2 \leq \|x - (z + t(y - z))\|^2$ for $t \in (0, 1]$ and then sending $t \rightarrow 0$. We find

$$\langle x - z, y - z \rangle \leq 0 \quad \forall y \in C.$$

⁴the typical example is a linear function defined by $f(e_n) = n$ where $(e_n)_{n \geq 1}$ is an independent family, which is then completed into a basis \mathcal{B} , then, one lets $f(e) = 0$ if $e \in \mathcal{B} \setminus \{e_n : n \geq 1\}$.

It follows that if $v = x - z \neq 0$, $y \in C$,

$$\langle v, x \rangle = \langle x - z, x \rangle = \|x - z\|^2 + \langle x - z, z \rangle \geq \|x - z\|^2 + \langle x - z, y \rangle = \|v\|^2 + \langle v, y \rangle.$$

The result follows (letting for instance $\alpha = \|v\|^2/2 + \sup_{y \in C} \langle v, y \rangle$). The proof can easily be extended to the situation where $\{x\}$ is replaced with a compact convex set not intersecting C .

Corollary 4.4. *In a real Hilbert space \mathcal{X} , a closed convex set C is weakly closed.*

Indeed, if $x \notin C$, one finds v, α with $\langle v, x \rangle > \alpha \geq \langle v, y \rangle$ for all $y \in C$ and this defines a neighborhood $\{\langle v, \cdot \rangle > \alpha\}$ of x for the weak topology which does not intersect C .

Theorem 4.5. *Let \mathcal{X} be a (real) Hilbert space, $C \subset \mathcal{X}$ an open convex set and $C' \subset \mathcal{X}$ a convex set with $C' \cap C = \emptyset$. Then there exists a closed hyperplane which “separates” C and C' : precisely, in the Hilbertian setting, one can find $v \in X, \alpha \in \mathbb{R}$, $v \neq 0$, such that*

$$\langle v, x \rangle \geq \alpha \geq \langle v, y \rangle \quad \forall x \in C, y \in C'$$

Proof: first assume that $C' = \{\bar{x}\}$ is a singleton. The difficult case is whenever $\bar{x} \in \overline{C} \setminus C$, otherwise we can apply Theorem 4.3 to separate (strictly) \bar{x} and \overline{C} . By assumption, there exists a ball $B = B(y, \delta) \subset C$. Let $x_n = y + (1 + 1/n)(\bar{x} - y)$, which is such that $x_n \rightarrow \bar{x}$ as $n \rightarrow \infty$. Since

$$\bar{x} = \frac{n}{n+1}x_n + \frac{1}{n+1}y,$$

one has $x_n \notin \overline{C}$, otherwise by convexity one would deduce that $B(\bar{x}, \delta/(n+1)) \subset \overline{C}$ so that $\bar{x} \in C$, a contradiction.

By Theorem 4.3 there exists v_n such that for all $x \in \overline{C}$,

$$\langle v_n, x_n \rangle \leq \langle v_n, x \rangle$$

and we can assume $\|v_n\| = 1$. Up to a subsequence, we may then assume that $v_n \rightarrow v$ weakly in \mathcal{X} . In the limit, (using that $x_n \rightarrow \bar{x}$ strongly) we obtain $\langle v, \bar{x} \rangle \leq \langle v, x \rangle \forall x \in C$, which is our claim if $v \neq 0$.

Using again the ball $B(y, \delta) \subset C$, one has for any $\|z\| \leq 1$

$$\langle v_n, x_n \rangle \leq \langle v_n, y - \delta z \rangle$$

so that $\langle v_n, y - x_n \rangle \geq \delta \langle v_n, z \rangle$: and taking the supremum over all possible z we find $\langle v_n, y - x_n \rangle \geq \delta$. In the limit we deduce $\langle v, y - \bar{x} \rangle \geq \delta$ which shows that $v \neq 0$.

Now, to show the general case, one lets $A = C' - C = \{y - x : y \in C', x \in C\}$: this is an open convex set and by assumption, $0 \notin A$. Hence by the previous part, there exists $v \neq 0$ such that $\langle v, y - x \rangle \leq \langle v, 0 \rangle = 0$ for all $y \in C', x \in C$, which is the thesis of the Theorem.

These simple examples of separation theorems are geometric versions of the Hahn-Banach theorem and are valid in fact in a much more general setting, see [7, 15].

4.1.3 Subgradient

Given a convex, extended valued, $f : \mathcal{X} \rightarrow]-\infty, +\infty]$, its subgradient at a point x is defined as the set

$$\partial f(x) := \{p \in \mathcal{X} : f(y) \geq f(x) + \langle p, y - x \rangle \forall y \in \mathcal{X}\}.$$

This is a closed, convex set.

If f is (Fréchet-)differentiable at x , then it is easy to see that $\partial f(x) = \{\nabla f(x)\}$: one has

$$f(y) = f(x) + \langle \nabla f(x), y - x \rangle + o(|y - x|)$$

so that if $p \in \partial f(x)$,

$$\langle \nabla f(x) - p, y - x \rangle + o(|y - x|) \geq 0.$$

Taking $y = x + th$, for $h \in \mathcal{X}$ and $t > 0$ small, we find after dividing by t and sending $t \rightarrow 0$ that $\langle \nabla f(x) - p, h \rangle \geq 0$. Hence $p = \nabla f(x)$. We leave to the reader the proof that $\nabla f(x) \in \partial f(x)$ (hence $\nabla f(x) \neq \emptyset$), which follows from convexity.

A converse? Now we want to understand better the structure of the subgradient in relationship to the behaviour of f at a point x , and in particular what can be said whenever $\partial f(x) = \{p\}$ is a singleton. First, we observe that if $x \in \text{dom } f$, $v \in \mathcal{X}$, $t > s > 0$ one has

$$f(x + sv) = f((s/t)(x + tv) + (1 - s/t)x) \leq \frac{s}{t}f(x + tv) + (1 - \frac{s}{t})f(x)$$

so that

$$\frac{f(x + sv) - f(x)}{s} \leq \frac{f(x + tv) - f(x)}{t}.$$

It follows that

$$f'(x; v) := \lim_{t \downarrow 0^+} \frac{f(x + tv) - f(x)}{t} = \inf_{t > 0} \frac{f(x + tv) - f(x)}{t} \quad (17)$$

is well defined (in $[-\infty, \infty]$), and $< +\infty$ as soon as $\{x + tv : t > 0\} \cap \text{dom } f \neq \emptyset$. If $x \in \overbrace{\text{dom } f}$, then $f'(x; v) < \infty$ for all v , moreover as $f'(x; 0) = 0 \leq f'(x; v) + f'(x; -v)$ it is not $-\infty$ either. In fact, $f'(x; \cdot)$ is a limit of convex functions, and hence convex, moreover, it is clearly positively 1-homogeneous: $f'(x; \lambda v) = \lambda f'(x; v)$ for all $\lambda \geq 0$ and all v .

If this quantity is finite, then the function has a Gateaux derivative in the direction v (however, usual definitions of Gateaux differentiability require that this derivative be a continuous linear form of v).

By definition, one easily sees that $f'(x; v) \geq \langle p, v \rangle$ if and only if $p \in \partial f(x)$. ($f'(x; v) \geq \langle p, v \rangle \Rightarrow f(x + tv) - f(x) \geq t \langle p, v \rangle$ for all $t > 0$, $v \in \mathcal{X}$.) This means that

$$\partial f'(x; \cdot)(0) = \partial f(x). \quad (18)$$

If in addition, f is locally bounded near $x \in \overbrace{\text{dom } f}$ (for this, as we have seen in Lemma 4.1, it is enough that f be locally bounded near one point of the domain, or that f be lsc, cf Remark 4.2), then one can easily deduce that also $f'(x; \cdot)$ is, and in particular it is Lipschitz (globally, as it is 1-homogeneous).

(In addition, in finite dimension, the convergence in (17) is uniform for $\|v\| \leq 1$ because of Ascoli-Arzelà's theorem: In fact, if $t \leq t_0$ small enough and $\|v\| \leq 2$,

$$h_x^t(v) := \frac{f(x+tv) - f(x)}{t} \leq \frac{f(x+t_0v) - f(x)}{t_0} \leq M$$

for some M and the proof of Lemma 4.1 shows that the h_x^t are uniformly Lipschitz in $B(0, 1)$.)

We will see later on (Sections 4.2, 4.2.2) that since in these cases, $f'(x; \cdot)$ is continuous, $f'(x; v) = \sup_{p \in \partial f(x; \cdot)(0)} \langle p, x \rangle$, so that $f'(x; \cdot)$ is the *support function* of $\partial f(x)$ which in particular cannot be empty.

Moreover, we deduce that if $\partial f(x) = \{p\}$ is a singleton, then $f'(x; v) = \langle p, x \rangle$. In finite dimension (as the convergence $h_x^t \rightarrow f'(x; \cdot)$ is uniform) we deduce that f is differentiable at x . In infinite dimension, we deduce that f is Gateaux-differentiable. It is not necessarily Fréchet-differentiable: for instance in $\ell^2(\mathbb{N})$, the convex function

$$f(x) = \sup_{i \geq 0} \left(\sqrt{\frac{1}{i+1} + x_i^2} - \sqrt{\frac{1}{i+1}} \right)$$

which is bounded near 0 ($\|x\| \leq 1 \Leftrightarrow \sum_i |x_i|^2 \leq 1 \Rightarrow |x_i| \leq 1 \forall i \geq 0$) satisfies $\partial f(0) = \{0\}$, however if $v = e_i = (\delta_{i,j})_{j \geq 0}$, then

$$\frac{f(0+tv) - f(0)}{t} = \frac{1}{t} \left(\sqrt{\frac{1}{i+1} + t^2} - \sqrt{\frac{1}{i+1}} \right) = \sqrt{2} - 1$$

if $t = 1/\sqrt{i+1}$, showing that the differentiability is only Gateaux. (It is, as for any v and $t > 0$,

$$\frac{f(tv) - f(0)}{t} = \sup_{i \geq 0} \frac{1}{t} \left(\sqrt{\frac{1}{i+1} + t^2 v_i^2} - \sqrt{\frac{1}{i+1}} \right)$$

and for each i , the quantity in the sup is less than $|v_i|$. Given ε , one can find i_0 such that $|v_i| \leq \varepsilon$ for $i > i_0$, while for $i = 0, \dots, i_0$, if t is small enough one can make the quantity below the sup less than ε . Hence the Gateaux derivative exists and is zero.)

Using Lemma 4.1, we can deduce the following two results:

Lemma 4.6. *Let f be proper, convex. Assume it is lsc, or continuous in one point. Then, in the interior of the domain, $\partial f(x) \neq \emptyset$. In finite dimension, f has a nonempty subdifferential everywhere in $\text{ri dom } f$.*

Lemma 4.7. *Let f be proper, convex. Then if f is Gateaux-differentiable at x , $\partial f(x) = \{\nabla f(x)\}$. Conversely if x is in the interior of $\text{dom } f$ and f is continuous at some point⁵, then if $\partial f(x)$ is a singleton, f is Gateaux-differentiable at x .*

In finite dimension, ∂f is a singleton if and only if f is differentiable at x .

Minimality condition An obvious remark which stems from the definition of a subgradient is that this notion allows to generalise the Euler-Lagrange stationary conditions ($\nabla f(x) = 0$ if x is a minimiser of f) to nonsmooth convex functions: we have indeed

$$\boxed{x \in \mathcal{X} \text{ is a global minimiser of } f \text{ if and only if } 0 \in \partial f(x).} \quad (19)$$

In the same way, one has:

⁵for instance if it is lsc, cf Rem. 4.2 and Lemma 4.1.

Lemma 4.8. *if $x \in \text{dom } f$ is a local minimiser of $f + g$, f convex, g C^1 near x , then for all $y \in \mathcal{X}$,*

$$f(y) \geq f(x) - \langle \nabla g(x), y - x \rangle$$

and $-\nabla g(x) \in \partial f(x)$.

Indeed, one just writes that for $t > 0$ small enough,

$$f(x) + g(x) \leq f(x + t(y - x)) + g(x + t(y - x)) \leq f(x) + t(f(y) - f(x)) + g(x + t(y - x))$$

so that

$$\frac{g(x) - g(x + t(y - x))}{t} \leq f(y) - f(x)$$

and we recover the claim in the limit $t \rightarrow 0$.

Subgradient of a strongly convex function If the function f is strongly convex or “ μ -convex” and $p \in \partial f(x)$, then x is by definition a minimiser of $y \mapsto f(y) - \langle p, y - x \rangle$ which is also μ -convex. In particular, letting $h(y) = f(y) - \langle p, y - x \rangle - \mu \|y - x\|^2/2$, one has that x is a minimizer of $h(y) + \mu \|y - x\|^2/2$ and h is convex. Hence, by Lemma 4.8,

$$0 = -\nabla \left(\frac{\mu}{2} \|\cdot - x\|^2 \right) (x) \in \partial h(x).$$

Hence, $h(y) \geq h(x)$ for all $y \in \mathcal{X}$, that is

$$f(y) - \langle p, y - x \rangle - \mu \|y - x\|^2/2 \geq f(x).$$

We deduce that for any $x, y \in \mathcal{X}$ and $p \in \partial f(x)$:

$$f(y) \geq f(x) + \langle p, y - x \rangle + \frac{\mu}{2} \|y - x\|^2 \quad (20)$$

An equivalent (but important) remark is that if f is strongly convex and x is a minimiser, then one has (since $0 \in \partial f(x)$)

$$f(y) \geq f(x) + \frac{\mu}{2} \|y - x\|^2 \quad (21)$$

for all $y \in \mathcal{X}$.

Domain The domain of ∂f is the set $\text{dom } \partial f = \{x \in \mathcal{X} : \partial f(x) \neq \emptyset\}$. Clearly, $\text{dom } \partial f \subset \text{dom } f$, in fact if f is convex, l.s.c. and proper, we will see later on (see Prop 4.22 or [15]) that $\text{dom } \partial f$ is dense in $\text{dom } f$ (even when $\text{dom } f$ has empty interior, as for instance when $f(u) = \int_{\Omega} |\nabla u|^2 dx$ for $u \in L^2(\Omega)$). The fact it is not empty will also follow.

In finite dimension, one has seen that for a proper convex function, $\text{dom } \partial f$ contains at least the relative interior of $\text{dom } f$ (that is, the interior in the vector subspace which is generated by $\text{dom } f$).

4.1.4 Subdifferential calculus

Theorem 4.9. *Assume f, g are convex, proper. Then for all x , $\partial f(x) + \partial g(x) \subset \partial(f + g)(x)$. Moreover if there exists $\bar{x} \in \text{dom } f$ where g is continuous, then $\partial f(x) + \partial g(x) = \partial(f + g)(x)$. In finite dimension, if $\text{ri } \text{dom } g \cap \text{ri } \text{dom } f \neq \emptyset$, this is also true.*

Proof: the inclusion is obvious from the definition. For the reverse inclusion, we assume $p \in \partial(f+g)(x)$ and want to show that it can be decomposed as $q+r$ with $q \in \partial f(x)$ and $r \in \partial g(x)$. By definition, we have that $f(y)+g(y) \geq f(x)+g(x)+\langle p, y-x \rangle$.

Thanks to the assumption that g is continuous at \bar{x} , $\text{epi}(g(\cdot) - \langle p, \cdot \rangle)$ contains a ball B centered at $(\bar{x}, g(\bar{x}) - \langle p, \bar{x} \rangle + 1)$ and has non empty interior. Denote E this interior, and F the following translation/flip of $\text{epi } f$:

$$F = \{(y, t) : -t \geq f(y) - [f(x) + g(x) - \langle p, x \rangle]\},$$

which is convex. For $(y, t) \in F$, one has $-t \geq f(y) - [f(x) + g(x) - \langle p, x \rangle] \geq -[g(y) - \langle p, y \rangle]$, that is $t \leq [g(y) - \langle p, y \rangle]$ so that $(y, t) \notin E$. Hence by Theorem 4.5 there exists $(q, \lambda) \neq (0, 0)$, such that for all $(y, t) \in E$, $(y', t') \in F$,

$$\langle q, y \rangle + \lambda t \geq \langle q, y' \rangle + \lambda t'.$$

As t' can be sent to $-\infty$ (or t to $+\infty$), $\lambda \geq 0$. Moreover since \bar{x} is in $\text{dom } f$, if $\lambda = 0$ one finds that $\langle q, y - \bar{x} \rangle \leq 0$ for all $y \in \text{dom } g$ which contains a ball centered in \bar{x} , so that $q = 0$, which is a contradiction. Hence $\lambda > 0$ so that without loss of generality we can assume $\lambda = 1$.

In particular choosing $t' = f(x) + g(x) - \langle p, x \rangle - f(y')$,

$$\langle q, y \rangle + t \geq \langle q, y' \rangle + f(x) + g(x) - \langle p, x \rangle - f(y').$$

for all $(y, t) \in E$. The closure of E contains $\text{epi}(g(\cdot) - \langle p, \cdot \rangle)$: indeed any $(y, t) \in \text{epi}(g(\cdot) - \langle p, \cdot \rangle)$ is on the boundary of the set $\{ty + (1-t)B : 0 < t < 1\} \subset \text{epi}(g(\cdot) - \langle p, \cdot \rangle)$. Hence it follows that for all y, y' ,

$$\begin{aligned} \langle q, y \rangle + g(y) - \langle p, y \rangle &\geq \langle q, y' \rangle + f(x) + g(x) - \langle p, x \rangle - f(y') \\ &\Leftrightarrow f(y') + g(y) \geq f(x) + g(x) + \langle p, y - x \rangle + \langle q, y' - y \rangle \\ &= f(x) + g(x) + \langle p - q, y - x \rangle + \langle q, y' - x \rangle \end{aligned}$$

showing that $q \in \partial f(x)$ and $r = p - q \in \partial g(x)$, as requested.

In finite dimension, the proof relies on the previous result and the fact that subgradients, for points in the relative interior of a convex function, are the sum of a subgradient of a Lipschitz function and a vector orthogonal to the domain, which is a consequence of the following easy fact (actually valid for \mathcal{X} Hilbert):

Lemma 4.10. *Let $f : \mathcal{X} \rightarrow \mathbb{R} \cup \{+\infty\}$ be convex, proper and let $W \subset \mathcal{X}$ be an affine, closed subspace with $\text{dom } f \subset W$. Then for any $x \in W$,*

$$\partial f(x) = \partial(f|_W)(x) + W^\perp.$$

We denote W_0 the vector space $\{x - y : (x, y) \in W^2\}$. If $p \in \partial f(x)$ and $y \in W$, one has

$$f(y) - f(x) \geq \langle p, y - x \rangle = \langle \Pi_{W_0}(p), y - x \rangle$$

since $y - x \in W_0$, so that $\Pi_{W_0}(p) \in \partial(f|_W)(x)$. Conversely, if $\tilde{p} \in W_0$ is an element of $\partial(f|_W)(x)$, obviously for any $y \in \mathcal{X}$ and $q \in W^\perp$,

$$f(y) \geq f(x) + \langle \tilde{p} + q, y - x \rangle$$

since either $y \notin W$ and $f(y) = +\infty$, or $y \in W$ and $\langle g, y - x \rangle = 0$. This shows the lemma.

In particular, one sees that for any $x \in W$, $\partial\delta_W(x) = W^\perp$ (where we recall δ_W is the indicator or characteristic function of W); also, in finite dimension, if one chooses $W = \text{vect dom } f$, we remark that for a.e. point (for the Lebesgue measure in W) in $\text{ri dom } f$, then $\partial f(x) = \{\nabla f|_W(x)\} + W^\perp$.

We now show that, when \mathcal{X} is finite dimensional, then if $\text{ri dom } f \cap \text{ri dom } g \neq \emptyset$, $\partial(f + g) = \partial f + \partial g$. We deduce this from the previous result (which we refer as Theorem 4.9) and from Lemma 4.10; a direct proof is found in [38, Thm 23.8], it is actually not much simpler.

First, up to a translation, we may assume $0 \in \text{ri dom } g \cap \text{ri dom } f$ so that $\text{ri dom } g$ is the interior of $\text{dom } g$ in $V = \text{vect dom } g$ and $\text{ri dom } f$ is the interior of $\text{dom } f$ in $W = \text{vect dom } f$.

If $x \notin W \cap V$, $f(x) + g(x) = +\infty$ and $\partial f(x) + \partial g(x) = \partial(f + g)(x) = \emptyset$, so we assume $x \in W \cap V$. From the Lemma we have that

$$\partial(f + g)(x) = \partial((f + g)|_W)(x) + W^\perp$$

since $\text{dom}(f + g) \subset W$. Now since $f|_W$ is continuous at all points of $\text{ri dom } f$, and by assumption one of such points is in $\text{dom } g|_W$, we deduce from Theorem 4.9 that

$$\partial((f + g)|_W)(x) + W^\perp = \partial(f|_W)(x) + \partial(g|_W)(x) + W^\perp = \partial f(x) + \partial(g|_W)(x) + W^\perp$$

where in the last equality we have used again that $\partial f = \partial(f|_W) + W^\perp$. On the other hand still because of Lemma 4.10,

$$\partial(g|_W)(x) + W^\perp = \partial(g + \delta_W)(x) = \partial((g + \delta_W)|_V)(x) + V^\perp = \partial(g|_V + \delta_{W \cap V})(x) + V^\perp.$$

Now, using the fact that $g|_V$ is continuous at some point of W (again, from the assumption $\text{ri dom } f \cap \text{ri dom } g \neq \emptyset$), we can use Theorem 4.9 again to deduce that

$$\partial(g|_V + \delta_{W \cap V})(x) = \partial(g|_V)(x) + \partial\delta_{W \cap V}(x)$$

Since we have assumed $x \in W \cap V$, one has⁶ $\partial\delta_{W \cap V}(x) = (W \cap V)^\perp = W^\perp + V^\perp$ so that, using Lemma 4.10 once more:

$$\partial(g|_W)(x) + W^\perp = \partial(g|_V)(x) + \partial\delta_{W \cap V}(x) + V^\perp = \partial g(x) + W^\perp.$$

We deduce that $\partial(f + g)(x) = \partial f(x) + \partial g(x)$.

Theorem 4.11. *Let $A : \mathcal{X} \rightarrow \mathcal{Y}$ be a continuous operator between two Hilbert spaces and f a proper, convex function on \mathcal{Y} . Let $g = f(Ax)$, then if there is \bar{x} such that f is continuous at $A\bar{x}$, $\partial g(x) = A^*\partial f(Ax)$. In finite dimension, one can just require that $A\bar{x} \in \text{ri dom } f$.*

Proof: $A^*\partial f(Ax) \subset \partial g(x)$ is easy. If $p \in \partial g(x)$, one has for all z ,

$$f(Az) \geq f(Ax) + \langle p, z - x \rangle. \quad (22)$$

Hence $\overset{\circ}{\text{epi } f}$ (which is non empty because f is continuous at some point) and

$$E = \{(Az, f(Ax) + \langle p, z - x \rangle) : z \in \mathcal{X}\} \subset \mathcal{Y} \times \mathbb{R}$$

⁶We use here that $(V \cap W)^\perp = V^\perp + W^\perp$ which easily follows from the obvious relationship $(A + B)^\perp = A^\perp \cap B^\perp$ and $(V^\perp)^\perp = V$ — which is elementary duality.

have no common point: if $(y, t) \in E$, then by (22) $t \leq f(y)$. Then there exists by Theorem 4.5 (q, λ) such that

$$-\langle q, y \rangle + \lambda t \geq -\langle q, y' \rangle + \lambda t'$$

for all $(y, t) \in \text{epi } f$ and $(y', t') \in E$. Again, $\lambda \geq 0$, and if $\lambda = 0$ one can find a contradiction as in the previous proof. Then, assuming $\lambda = 1$, one obtains for all $z \in \mathcal{X}$, $y \in \mathcal{Y}$,

$$-\langle q, y \rangle + f(y) \geq -\langle q, Az \rangle + f(Ax) + \langle p, z - x \rangle = f(Ax) + \langle p - A^*q, z \rangle - \langle p, x \rangle.$$

This is possible only if $p = A^*q$, otherwise one can send the right-hand side to $+\infty$. Hence, $p = A^*q$, $\langle p, x \rangle = \langle q, Ax \rangle$ and

$$f(y) \geq f(Ax) + \langle q, y - Ax \rangle$$

for all y , so that $q \in \partial f(Ax)$.

In finite dimension, we leave the proof to the reader (see also [38, Thm 23.9]).

4.1.5 Remark: KKT's theorem

Theorem 4.12 (Karush-Kuhn-Tucker). *Let f, g_i , $i = 1, \dots, m$ be C^1 , convex and assume*

$$\exists \bar{x}, (g_i(\bar{x}) < 0 \forall i = 1, \dots, m) \quad (\text{Slater's condition})$$

Then x^ is a solution of*

$$\min_{g_i(x) \leq 0, i=1, \dots, m} f(x)$$

if and only if there exists $(\lambda_i)_{i=1}^m$, $\lambda_i \geq 0$ such that

$$\nabla f(x^*) + \sum_{i=1}^m \lambda_i \nabla g_i(x^*) = 0 \quad (23)$$

and for all $i = 1, \dots, m$:

$$\lambda_i g_i(x^*) = 0 \quad (\text{complementary slackness condition})$$

Proof: first, if (23) holds together with the complementary slackness condition, then it is easy to show that x^* , which is a minimizer of the convex function $f + \sum_i \lambda_i g_i$, is a solution of the constrained problem: if x satisfies the constraints, then

$$f(x) \geq f(x) + \sum_i \lambda_i g_i(x) \geq f(x^*) + \sum_i \lambda_i g_i(x^*) = f(x^*).$$

Conversely, consider for all i the function

$$\delta_i(x) = \begin{cases} 0 & \text{if } g_i(x) \leq 0, \\ +\infty & \text{else,} \end{cases}$$

then the problem is equivalent to $\min_x f(x) + \sum_i \delta_i(x)$. By Slater's condition, we know that there exists \bar{x} where all functions f, δ_i are continuous. Hence by Thm. 4.9,

$$0 \in \partial(f + \sum_i \delta_i)(x^*) = \nabla f(x^*) + \sum_{i=1}^m \partial \delta_i(x^*).$$

It remains to characterize $\partial\delta_i(x^*)$: if $g_i(x^*) < 0$ then it is negative in a neighborhood of x^* and $\partial\delta_i(x^*) = \{0\}$. If $g_i(x^*) = 0$, then we need to characterize the vectors p such that for all y with $g_i(y) \leq 0$,

$$0 \geq \langle p, y - x^* \rangle.$$

Let $v \perp \nabla g_i(x^*)$, and consider $y = x^* - t(\nabla g_i(x^*) + v)$: then

$$g_i(y) = -t \langle \nabla g_i(x^*), \nabla g_i(x^*) + v \rangle + o(t) = -t \|\nabla g_i(x^*)\|^2 + o(t) < 0$$

if $t > 0$ is small enough, hence

$$0 \leq \langle p, \nabla g_i(x^*) + v \rangle.$$

We easily deduce that we must have $p = \lambda_i \nabla g_i(x^*)$, for some $\lambda_i \geq 0$ (in other words, $\partial\delta_i(x^*) = \mathbb{R}_+ \nabla g_i(x^*)$). The theorem follows.

Remark 4.13. The standard KKT theorem suggests also the possibility of some affine equality constraints $h_i(x) = 0$, $i = 1, \dots, m'$, and the Slater condition just assumes $h_i(\bar{x}) = 0$. The proof above needs to be tuned a little to address this case. In practice, one can observe that when solving, for some $i \in 1, \dots, m'$ the problem with either the constraint $+h_i \leq 0$ or $-h_i \leq 0$, one finds two solutions x^\pm with value \mathbf{m}^\pm and: either $\mathbf{m}^+ < \mathbf{m}^-$, in which case one easily shows that $-h_i(x^-) = 0$ (otherwise one could find a better value for \mathbf{m}^- in the interval $[x^+, x^-]$), or $\mathbf{m}^+ > \mathbf{m}^-$ and $h_i(x^+) = 0$, or $\mathbf{m}^+ = \mathbf{m}^-$ and the problem is equivalent when removing the constraint $h_i = 0$. As a result, the initial problem is shown to be equivalent to

$$\min_x \{f(x) : g_i(x) \leq 0, i = 1, \dots, m; \epsilon_i h_i(x) \leq 0, i = 1, \dots, m'\}$$

where $\epsilon_i \in \{-1, 0, 1\}$, and the standard KKT conditions follow by applying the Theorem to this new problem, observing that one can perturb slightly \bar{x} to find a new point \bar{x}' with $\epsilon_i h_i(\bar{x}') < 0$ for all i with $\epsilon_i \neq 0$.

4.2 Convex duality

4.2.1 Legendre-Fenchel conjugate

Given a function $f : \mathcal{X} \rightarrow \mathbb{R} \cup \{+\infty\}$, we introduce the *Legendre-Fenchel conjugate*

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

which is defined for all $p \in \mathcal{X}$, as a supremum of continuous linear forms: in particular, it is obviously a convex, lsc function. Observe that here we rely on the Riesz theorem to define the conjugate, in a more general vector space E , the proper definition should be as a function defined in a dual space E' , see for instance [15].

Obviously for all x, y ,

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

and in particular $f(x) \geq \langle y, x \rangle - f^*(y)$. Thus, the *biconjugate* f^{**} , defined as f^* by $f^{**}(y) = \sup_{x \in \mathcal{X}} \langle y, x \rangle - f^*(y)$, clearly satisfies

$$f^{**} \leq f.$$

The following is the most important result about the Legendre-Fenchel conjugate (it is also elementary in our Hilbertian setting):

Theorem 4.14. *If f has no affine minorant, $f^* \equiv +\infty$ and $f^{**} \equiv -\infty$. Otherwise, f^{**} is the largest convex lsc function below f , called the convex lsc envelope of f (sometimes also the Γ -regularization, or the convex relaxation). In this case then either $f \equiv +\infty$, or f^* , f^{**} are proper.*

This is a consequence of the separation theorem. Observe that the convex lsc envelope of a function f is always well defined as the sup of all the convex lsc functions below f , or $-\infty$ if there is none. Observe also that it is the function whose epigraph is the closed convex envelope of $\text{epi } f$. The special case of a function with no affine minorant is very specific: for instance, a function f equal to $-\infty$ in $B(0, 1)$ and $+\infty$ else, despite being convex lsc, is such that $f^* \equiv +\infty$ and $f^{**} \equiv -\infty$.

Proof: if $f \equiv +\infty$ then $f^* \equiv -\infty$ and $f^{**} \equiv +\infty$: the theorem is trivial. So we assume there exists x with $f(x) < +\infty$. As we have seen, $f^{**} \leq f$ is a convex lsc below f . Either it is $-\infty$ everywhere, or it is proper and there exists an affine function a such that $f \geq f^{**} \geq a$. Indeed, choosing (x, t) with $t < f^{**}(x) \leq f(x) < +\infty$, the separation Theorem 4.3 applied to the closed convex set $\text{epi } f^{**}$ and $(x, t) \notin \text{epi } f^{**}$ shows the existence of (p, λ, α) with

$$-\langle p, x \rangle + \lambda t < \alpha \leq -\langle p, y \rangle + \lambda s$$

for any $(y, s) \in \text{epi } f^{**}$. As usual $\lambda \geq 0$ (sending $s \rightarrow \infty$), moreover $\lambda \neq 0$ otherwise choosing $y = x$ yields a contradiction. Hence one may assume $\lambda = 1$ and one obtains $f(y) \geq f^{**}(y) \geq t + \langle p, y - x \rangle =: a(y)$, which shows the claim. By definition, one has of course in this case that $f^*(p) \leq \langle p, x \rangle - t < \infty$ and $f^{**}(y) \geq \langle p, y \rangle - f^*(p) \forall y$.

One sees that if f has no affine minorant, then $f^* \equiv +\infty$ and $f^{**} \equiv -\infty$; while in the other case f^* and f^{**} are proper as soon as $f \not\equiv +\infty$.

Assuming that we are in the latter case, let g be convex, lsc with $g \leq f$. To show that f^{**} is maximal among such functions, we must show that $g \leq f^{**}$. Since $g \leq f$, then $f^* \leq g^*$, so that $g^{**} \leq f^{**}$. Hence it is enough to show that $g^{**} = g$. As before, considering p with $f^*(p) < +\infty$ one has $f^{**} \geq \langle p, \cdot \rangle - f^*(p)$, so that it is enough to consider only functions g with $g(x) \geq \langle p, x \rangle - f^*(p) \forall x$ (otherwise replace g with $x \mapsto \max\{g(x), \langle p, x \rangle - f^*(p)\}$).

The next (not essential) simplification consists in replacing f with $f'(x) = f(x) - \langle p, x \rangle + f^*(p) \geq 0$. Indeed,

$$\begin{aligned} (f')^*(y) &= \sup_x \langle y, x \rangle - f(x) + \langle p, x \rangle - f^*(p) \\ &= -f^*(p) + \sup_y \langle y + p, x \rangle - f(x) = f^*(y + p) - f^*(p), \end{aligned}$$

so that

$$\begin{aligned} (f')^{**}(x) &= \sup_y \langle y, x \rangle - f^*(y + p) + f^*(p) \\ &= f^*(p) - \langle p, x \rangle + \sup_y \langle y + p, x \rangle - f^*(y + p) = f^{**}(x) - \langle p, x \rangle + f^*(p). \end{aligned}$$

Hence $f = f^{**} \Leftrightarrow f' = (f')^{**}$ and it is enough to show the result for nonnegative functions.

Assume therefore that f is convex, lsc, with $0 \leq f \not\equiv +\infty$. If $f^{**} \neq f$, then there exists x with $f^{**}(x) < f(x)$. That is, $(x, f^{**}(x)) \notin \text{epi } f$ and from Theorem 4.3, there exists p, λ, α with

$$\langle p, x \rangle - \lambda f^{**}(x) > \alpha \geq \langle p, y \rangle - \lambda s$$

for all $y \in \text{dom } f$ and $s \geq f(y)$. In particular, as $\text{dom } f \neq \emptyset$, letting $s \rightarrow +\infty$ we see that $\lambda \geq 0$.

Case 1: $\lambda > 0$: then we can divide the inequality and assume that $\lambda = 1$. It follows that $f^{**}(x) < -\alpha + \langle p, x \rangle$, while $\alpha \geq \langle p, y \rangle - f(y)$ for all y , hence taking the sup over y , $\alpha \geq f^*(p)$. Hence, $f^{**}(x) < \langle p, x \rangle - f^*(p)$, a contradiction.

Case 2: $\lambda = 0$: then $\langle p, x \rangle > \alpha \geq \langle p, y \rangle$ for all $y \in \text{dom } f$. Observe then that for $t > 0$, using that $f \geq 0$ in $\text{dom } f$ and $f = +\infty$ outside,

$$f^*(tp) = \sup_y t \langle p, y \rangle - f(y) \leq t \sup_{y \in \text{dom } f} \langle p, y \rangle \leq t\alpha,$$

so that

$$f^{**}(x) = \sup_q \langle q, x \rangle - f^*(q) \geq \sup_{t>0} \langle tp, x \rangle - f^*(tp) \geq \sup_{t>0} t(\langle p, x \rangle - \alpha) = +\infty$$

which is again a contradiction.

Remark 4.15 (Legendre-Fenchel Identity). If x realizes the sup in $f^*(y) = \sup_x \langle y, x \rangle - f(x)$ then for all z ,

$$\langle y, x \rangle - f(x) \geq \langle y, z \rangle - f(z) \Leftrightarrow f(z) \geq f(x) + \langle y, z - x \rangle$$

which means that $y \in \partial f(x)$. Conversely if $y \in \partial f(x)$, by definition one easily deduces that $f^*(y) \leq \langle y, x \rangle - f(x)$, and moreover that $f^{**}(x) = f(x)$, $y \in \partial f^{**}(x)$, and f is lsc at x . In particular we see that $\partial f^{**}(x) \supseteq \partial f(x)$ for all x .

One derives the celebrated *Legendre-Fenchel identity*:

$$y \in \partial f(x) \Leftrightarrow \langle x, y \rangle = f(x) + f^*(y) \Rightarrow x \in \partial f^*(y), \quad (24)$$

the latter being also an equivalence if f is lsc, convex (if $f = f^{**}$).

One also can check that conversely, if the convex function f is lsc at x , then $f^{**}(x) = f(x)$. This is true because f^{**} is the lsc envelope of f (since it is convex), which can be defined by $z \mapsto \inf_{z_n \rightarrow z} \liminf_n f(z_n)$.

One can derive as a corollary the following variant of Theorem 4.14, which may be useful (see Sec. 4.3.2).

Corollary 4.16. *Let $f : \mathcal{X} \rightarrow \mathbb{R} \cup \{+\infty\}$ be convex, proper and assume that f is lower-semicontinuous at $x \in \mathcal{X}$. Then $f^{**}(x) = f(x)$.*

To prove this, observe that the lower semicontinuity assumption implies that for any $t < f(x)$, there exists $\delta > 0$ such that $f(y) \geq t$ for all $y \in B(x, \delta)$, the open ball of center x and radius δ . In other words,

$$\text{epi } f \cap B(x, \delta) \times (-\infty, t) = \emptyset$$

Since the second set is open, also $\overline{\text{epi } f}$ does not intersect it. Since f is convex, $\overline{\text{epi } f} = \text{epi } f^{**}$ (thanks to Theorem 4.14) and one deduces that $t \leq f^{**}(x)$, which proves the claim.

4.2.2 Examples

1. $f(x) = \|x\|^2/(2\alpha)$, $\alpha > 0$: $f^*(y) = \alpha\|y\|^2/2$;

2. $f(x) = |x|^p/p$: $f^*(y) = |y|^{p'}/p'$, $1/p + 1/p' = 1$;
3. $F(f) = \|f\|_{L^p}^p/p$: $F^*(g) = \|g\|_{L^{p'}}^{p'}/p'$ (the duality is in L^2 , however this is also true in the $(L^p, L^{p'})$ duality, see [15]);
4. $f(x) = \delta_{B(0,1)}(x) = 0$ if $x \in B(0, 1)$, $+\infty$ else: $f^*(p) = |p|$.

The last example is a particular case of the following situation: if f is convex, 1-homogeneous, then

$$f^*(y) = \sup_x \langle y, x \rangle - f(x) = \sup_{t>0} \sup_x \langle y, tx \rangle - f(tx) = \sup_{t>0} t f^*(y) \in \{0, +\infty\}$$

and precisely

$$f^*(y) = \begin{cases} 0 & \text{if } \langle y, x \rangle \leq f(x) \ \forall x \in \mathcal{X}, \\ +\infty & \text{if } \exists x \in \mathcal{X}, \langle y, x \rangle > f(x). \end{cases}$$

Letting $C = \{y : \langle y, x \rangle \leq f(x) \ \forall x \in \mathcal{X}\} = \partial f(0)$, one has $f^* = \delta_C$ (C is clearly closed and convex, and f^* convex lsc). Eventually, observe that if f is lsc, then $f^{**} = f$ which shows that in this case

$$f(x) = \sup_{y \in \partial f(0)} \langle y, x \rangle.$$

Observe in particular that $\partial f(x) = \{y \in \partial f(0) : \langle y, x \rangle = f(x)\}$.

This example, in turn, is a particular case of the following: if f is β -homogeneous, $\beta > 1$, then

$$f^*(ty) = \sup_x \langle ty, x \rangle - f(x) = t^\alpha \sup_x \langle y, t^{1-\alpha} x \rangle - f(t^{-\alpha/\beta} x) = t^\alpha f^*(y)$$

if $1 - \alpha = -\alpha/\beta$, hence if $1/\alpha + 1/\beta = 1$.

4.2.3 Relationship between the growth of f and f^*

Lemma 4.17. *If f is finite everywhere, then $f^*(tp)/t \rightarrow +\infty$ as $t \rightarrow +\infty$ for all $p \in \mathcal{X}$ (f^* is superlinear). The converse is true in finite dimension if f is convex, lsc.*

Proof: if f^* is not superlinear, there exists $p, c < \infty$, such that $f^*(tp) \leq ct$ for all $t > 0$: hence $f^{**}(x) \geq \sup_{t>0} t \langle p, x \rangle - f^*(tp) \geq \sup_{t>0} t(\langle p, x \rangle - c) = +\infty$ as soon as x is such that $\langle p, x \rangle > c$. Of course then, $f(x) \geq f^{**}(x) = +\infty$.

Conversely, in finite dimension, let f be convex, lsc and assume that there is x with $f(x) = +\infty$. We can assume without loss of generality that $f \geq 0$ (cf proof of Thm 4.14).

Then, since $\overline{\text{dom } f} \neq \mathcal{X}$ (in finite dimension only, in infinite dimension $\text{dom } f$ could be dense, for instance think of $f(u) = \int |\nabla u|^2 dx$ for $u \in L^2$) one can consider $x \notin \overline{\text{dom } f}$. Then, there exists by Theorem 4.3 p, α with $\langle p, x \rangle > \alpha \geq \langle p, y \rangle \ \forall y \in \text{dom } f$. We have

$$f^*(tp) = \sup_y \langle tp, y \rangle - f^*(y) \leq \sup_{y \in \text{dom } f} t \langle p, y \rangle \leq t\alpha$$

for $t > 0$, so that $f^*(tp)/t \leq \alpha$ and f^* is not superlinear.

Remark 4.18. In infinite dimension, one needs to strengthen a bit the assumption, for instance if $f \geq g(|p|)$ with g superlinear then f^* is finite everywhere.

Proposition 4.19. *Let f a convex, lsc function: then f is μ -convex if and only if f^* has $(1/\mu)$ -Lipschitz gradient.*

Proof: observe that if f is μ -convex one has in particular, given $x \in \text{dom } \partial f$, for $p \in \partial f(x)$, that (20) holds:

$$f(y) \geq f(x) + \langle p, y - x \rangle + \frac{\mu}{2} \|y - x\|^2 \quad (20)$$

for all y , hence taking the conjugate (cf Example 1 in the previous Section), we find for all q :

$$\begin{aligned} f^*(q) &\leq \sup_y \langle q, y \rangle - f(x) - \langle p, y - x \rangle - \frac{\mu}{2} \|y - x\|^2 \\ &= \langle q, x \rangle - f(x) + \sup_y \langle q - p, y - x \rangle - \frac{\mu}{2} \|y - x\|^2 = \langle q, x \rangle - f(x) + \frac{1}{2\mu} \|q - p\|^2 \\ &= \langle p, x \rangle - f(x) + \langle q - p, x \rangle + \frac{1}{2\mu} \|q - p\|^2 = f^*(p) + \langle q - p, x \rangle + \frac{1}{2\mu} \|q - p\|^2. \end{aligned} \quad (25)$$

We have used that $\langle p, x \rangle - f(x) = f^*(p)$ which follows from (24). In particular we see that f^* has at most a quadratic growth, and we deduce that it is locally Lipschitz (Lemma 4.1), and its subgradient is not empty everywhere. Moreover, we deduce from (25) that when $x \in \partial f^*(p) \Leftrightarrow p \in \partial f(x)$ (cf (24)),

$$f^*(p) + \langle q - p, x \rangle \leq f^*(q) \leq f^*(p) + \langle q - p, x \rangle + \frac{1}{2\mu} \|q - p\|^2,$$

in other words, $f^*(q) = f^*(p) + \langle q - p, x \rangle + o(\|q - p\|)$ which shows that f^* is (Fréchet)-differentiable and $x = \nabla f^*(p)$.

Eventually, given $p, q \in \mathcal{X}$ and $x = \nabla f^*(p)$, $y = \nabla f^*(q)$, one has by (24) that $p \in \partial f(x)$, $q \in \partial f(y)$ and by strong convexity, using (20) and the same with x, y switched and p replaced with q , and summing, we find

$$\langle q - p, y - x \rangle \geq \mu \|y - x\|^2$$

so that in particular, $\|\nabla f^*(q) - \nabla f^*(p)\| \leq (1/\mu) \|q - p\|$: ∇f^* is $(1/\mu)$ -Lipschitz. In fact, we see that

$$\langle q - p, \nabla f^*(q) - \nabla f^*(p) \rangle \geq \mu \|\nabla f^*(q) - \nabla f^*(p)\|^2,$$

which expresses that ∇f^* is “ μ -co-coercive”, a property which is stronger than being $(1/\mu)$ -Lipschitz.

Conversely, if f^* has $(1/\mu)$ -Lipschitz gradient, let us show that f is μ -convex. Observe that

$$\begin{aligned} f^*(q) &= f^*(p) + \int_0^1 \langle \nabla f^*(p + s(q - p)), q - p \rangle ds \\ &= f^*(p) + \langle \nabla f^*(p), q - p \rangle + \int_0^1 \langle \nabla f^*(p + s(q - p)) - \nabla f^*(p), q - p \rangle ds \\ &\leq f^*(p) + \langle \nabla f^*(p), q - p \rangle + \frac{1}{\mu} \|q - p\|^2 \int_0^1 s ds. \end{aligned}$$

If $p \in \partial f(x)$, so that $x = \nabla f^*(p)$, we deduce

$$f^*(q) \leq f^*(p) + \langle q - p, x \rangle + \frac{1}{2\mu} \|q - p\|^2.$$

Hence taking the conjugate:

$$\begin{aligned} f(y) = f^{**}(y) &\geq \sup_q \langle q, y \rangle - \left(f^*(p) + \langle q - p, x \rangle + \frac{1}{2\mu} \|q - p\|^2 \right) \\ &= \langle p, x \rangle - f^*(p) + \sup_q \langle p - q, x - y \rangle - \frac{1}{2\mu} \|q - p\|^2 = \langle p, y \rangle - f^*(p) + \frac{\mu}{2} \|x - y\|^2. \end{aligned}$$

By (24), $\langle p, x \rangle - f^*(p) = f(x)$ (as f is convex lsc), and we find

$$f(y) \geq f(x) + \langle p, y - x \rangle + \frac{\mu}{2} \|x - y\|^2,$$

showing that f is strongly convex. Notice in particular that we have found another proof of Theorem 2.3, valid also in Hilbert spaces for convex lsc functions.

4.2.4 The conjugate of a sum: Inf-convolutions

A natural question, given two convex functions f and g , is whether one can derive an expression for the conjugate $(f + g)^*$. The answer is given by a particular “convolution” formula, called the “inf-convolution”. Letting f, g be convex, lsc functions it is defined as follows:

$$f \square g(x) = \inf_y f(x - y) + g(y). \quad (26)$$

It is easy to show that this defines a convex function (more generally, given $G(x, y)$ convex in (x, y) , we let the reader show that $x \mapsto \inf_y G(x, y)$ is also convex. One can show in addition the following result:

Lemma 4.20. *We assume f, g are convex, lsc. If there is $p \in \mathcal{X}$ where f^* is continuous and g^* is finite, then the inf is reached in (26) and $f \square g$ is convex, lsc. In finite dimension, it is enough to have $p \in \text{ri dom } f^* \cap \text{ri dom } g^*$.*

Proof: consider indeed $x_n \rightarrow x$ and y_n such that

$$f \square g(x_n) \geq f(x_n - y_n) + g(y_n) - \frac{1}{n}.$$

Consider a subsequence with

$$\lim_k f(x_{n_k} - y_{n_k}) + g(y_{n_k}) = \lim_n \inf f(x_n - y_n) + g(y_n) \leq \lim_n \inf f \square g(x_n)$$

Observe that if f^* is continuous at p , then it means that there is a constant c such that

$$f^*(q) \leq c + \delta_{B(0, \varepsilon)}(q - p)$$

(where δ_C is the characteristic function of C which is zero in C and $+\infty$ elsewhere) while $g^*(p) < +\infty$: so that for all z

$$f(z) = f^{**}(z) \geq \langle p, z \rangle - c + \varepsilon \|z\|, \quad g(z) \geq \langle p, z \rangle - g^*(p).$$

Hence,

$$\begin{aligned} f(x_{n_k} - y_{n_k}) + g(y_{n_k}) &\geq \langle p, x_{n_k} - y_{n_k} \rangle - c + \varepsilon \|x_{n_k} - y_{n_k}\| + \langle p, y_{n_k} \rangle - g^*(p) \\ &= \langle p, x_{n_k} \rangle + \varepsilon \|x_{n_k} - y_{n_k}\| - (c + g^*(p)) \end{aligned}$$

so that $(x_{n_k} - y_{n_k})_k$ is a bounded sequence, hence there exists y and a subsequence of (y_{n_k}) (not relabelled) with $y_{n_k} \rightharpoonup y$. In the limit (as, f, g are weakly lsc),

$$f \square g(x) \leq \liminf_k f(x_{n_k} - y_{n_k}) + g(y_{n_k}) \leq \liminf_n f \square g(x_n).$$

Eventually, we observe that if the sequence $x_n \equiv x$, then this proves that there is a minimizer y in (26). We can derive a second, more precise variant of Theorem 4.9:

Corollary 4.21. *Let f, g be convex, lsc: if there exists $x \in \text{dom } f \cap \text{dom } g$ such that f is continuous at x (in finite dimension, $x \in \text{ri dom } f \cap \text{ri dom } g$), then*

- $(f + g)^* = f^* \square g^*$,
- $\partial(f + g) = \partial f + \partial g$.

The first point is clear: as by our assumption, $f^* \square g^*$ is lsc, and:

$$\begin{aligned} (f^* \square g^*)^*(x) &= \sup_{p, q} \langle x, p \rangle - f^*(q) - g^*(p - q) \\ &= \sup_{p, q} \langle x, q \rangle - f^*(q) + \langle x, p - q \rangle - g^*(p - q) = f(x) + g(x). \end{aligned}$$

The second point is because if $p \in \partial(f + g)(x)$, using that $x \in \partial(f^* \square g^*)(p)$ and

$$f^* \square g^*(p) = f^*(q) + g^*(p - q)$$

for some q , one obtains letting $p - q = r$:

$$\begin{aligned} f^*(s) + g^*(t) &\geq f^* \square g^*(s + t) \geq f^* \square g^*(p) + \langle x, s + t - p \rangle \\ &\geq f^*(q) + \langle x, s - q \rangle + g^*(r) + \langle x, t - r \rangle \end{aligned}$$

for all s, t . Hence $x \in \partial f^*(q) \cap \partial g^*(r)$, which shows that $p = q + r \in \partial f(x) + \partial g(x)$.

4.3 Example: the proximity operator

(Also known as *Proximal map*.) Given f convex lsc, proper, observe that for any $\tau > 0$, $x \in \mathcal{X}$, $y \mapsto f(y) + \|y - x\|^2/(2\tau)$ is strongly convex and hence has a unique minimizer. We define

$$f_\tau(x) := \min_{y \in \mathcal{X}} f(y) + \frac{1}{2\tau} \|y - x\|^2 \quad (27)$$

as the inf-convolution of f and $\|\cdot\|^2/(2\tau)$. It is clearly a convex, lsc function thanks to Lemma 4.20 (and the “min” is reached, but this is also because we are minimizing a strongly convex, lsc function in a Hilbert or Euclidean space). As we have seen before (Lemma 4.8), one has at the minimizer y_x

$$\partial f(y_x) + \frac{1}{\tau}(y_x - x) \ni 0. \quad (28)$$

This characterizes the unique minimizer of (27) and in particular it means that the following operator is uniquely defined:

$$y_x = (I + \tau \partial f)^{-1}(x) =: \text{prox}_{\tau f}(x).$$

As already shown, $(x - y_x)/\tau = \nabla f_\tau(x)$. Actually, in the convex case, there is a direct proof: one has, letting $\eta = \text{prox}_{\tau f}(y)$ and $\xi = \text{prox}_{\tau f}(x)$,

$$\begin{aligned} f_\tau(y) &= f(\eta) + \frac{\|\eta - y\|^2}{2\tau} = f(\eta) + \frac{\|(\eta - x) + (x - y)\|^2}{2\tau} \\ &= f(\eta) + \frac{\|\eta - x\|^2}{2\tau} + \left\langle \frac{x - \eta}{\tau}, y - x \right\rangle + \frac{\|x - y\|^2}{2\tau} \\ &\geq f(\xi) + \frac{\|\xi - x\|^2}{2\tau} + \frac{\|\eta - \xi\|^2}{2\tau} + \left\langle \frac{x - \xi}{\tau}, y - x \right\rangle + \left\langle \frac{\xi - \eta}{\tau}, y - x \right\rangle + \frac{\|x - y\|^2}{2\tau} \\ &= f_\tau(x) + \left\langle \frac{x - \xi}{\tau}, y - x \right\rangle + \frac{\tau}{2} \left\| \frac{y - \eta}{\tau} - \frac{x - \xi}{\tau} \right\|^2. \end{aligned}$$

In the third line, we have used the fact that ξ is the minimiser of a $(1/\tau)$ -strongly convex problem, so that $f(\eta) + \|\eta - x\|^2/(2\tau) \geq f(\xi) + \|\eta - x\|^2/(2\tau) + \|\eta - \xi\|^2/(2\tau)$ for all η . We deduce from the inequality

$$f_\tau(y) \geq f_\tau(x) + \left\langle \frac{x - \xi}{\tau}, y - x \right\rangle + \frac{\tau}{2} \left\| \frac{y - \eta}{\tau} - \frac{x - \xi}{\tau} \right\|^2$$

both that $(x - \xi)/\tau$ is a subgradient of f_τ at x , and that the map $x \mapsto (x - \text{prox}_{\tau f}(x))/\tau$ is τ -co-coercive, hence $(1/\tau)$ -Lipschitz: indeed, writing the same inequality after having swapped x and y , and summing the two inequalities, we obtain

$$\left\langle \frac{y - \eta}{\tau} - \frac{x - \xi}{\tau}, y - x \right\rangle \geq \tau \left\| \frac{y - \eta}{\tau} - \frac{x - \xi}{\tau} \right\|^2.$$

In particular, f_τ is C^1 . Also, we find that

$$\text{prox}_{\tau f}(x) = x - \tau \nabla f_\tau(x)$$

is a $(1/2)$ -averaged operator (it is $(1/2)I + (1/2)(x - 2\tau \nabla f_\tau(x))$, see Lemma 2.4).

Moreau's identity Thanks to (24), (28) yields

$$y_x \in \partial f^*\left(\frac{x - y_x}{\tau}\right) \Leftrightarrow \frac{x - y_x}{\tau} + \frac{1}{\tau} \partial f^*\left(\frac{x - y_x}{\tau}\right) \ni \frac{x}{\tau} \Leftrightarrow \frac{x - y_x}{\tau} = (I + \frac{1}{\tau} \partial f^*)^{-1}\left(\frac{x}{\tau}\right).$$

We deduce *Moreau's Identity*, valid for any convex, lsc, proper function f :

$$x = \text{prox}_{\tau f}(x) + \tau \text{prox}_{\frac{1}{\tau} f^*}\left(\frac{x}{\tau}\right) \quad (29)$$

One also can show the following:

Proposition 4.22. *Let f be proper, convex, lsc: then $\text{dom } \partial f$ is dense in $\text{dom } F$.*

Indeed, let $x \in \text{dom } f$: then $f_\tau(x) \leq f(x)$. In particular, denoting $x_\tau = \text{prox}_{\tau f}(x)$,

$$f_\tau(x) = f(x_\tau) + \frac{1}{2\tau} \|x - x_\tau\|^2 \leq f(x).$$

We use again that f , being proper, is larger than some affine function: hence there is p, c such that $\langle p, x_\tau \rangle + c + \frac{1}{2\tau} \|x - x_\tau\|^2 \leq f(x)$ from which it follows that $\|x_\tau - x\| \leq c' \sqrt{\tau}$ for some constant $c' > 0$. Hence $x_\tau \rightarrow x$. Now, $\partial f(x_\tau) \ni (x - x_\tau)/\tau \neq \emptyset$ hence $x_\tau \in \text{dom } f$, which shows the proposition. As a by-product of the proof, one sees that:

Proposition 4.23. *Let f be proper, lsc, convex and f_τ defined by (27). Then for all x , $f_\tau(x) \rightarrow f(x)$ as $\tau \rightarrow 0$.*

(We leave to the reader the proof that if $f(x) = +\infty$, $f_\tau(x) \rightarrow +\infty$, which is easy using that f is lsc.)

Examples: $f(x) = \|x\|_1 = \sum_i |x_i|$, $x \in \mathbb{R}^d$:

$$\text{prox}_{\tau f}(x) = ((|x_i| - \tau)^+ \text{sign}(x_i))_{i=1}^d.$$

If $f(x) = \delta_{|x_i| \leq 1}$, $\text{prox}_{\tau f}(x) = (\max\{-1, \min\{1, x_i\}\})_{i=1}^d$.
 If $f(x) = \|x\|^2/2$, $\text{prox}_{\tau f}(x) = x/(1 + \tau)$.

4.3.1 A useful variant of inf-convolutions

Consider now the modified inf-convolution problem

$$h(x) = \inf_{y \in \mathcal{Y}} f(x - Ky) + g(y)$$

where $K : \mathcal{Y} \rightarrow \mathcal{X}$ is a continuous operator and f, g are convex, lsc, proper. Then one can show similarly that if there exists p such that $f^*(p) < +\infty$ and g^* is continuous at K^*p , h is lsc and since

$$\begin{aligned} h^*(q) &= \sup_{x \in \mathcal{X}, y \in \mathcal{Y}} \langle q, x \rangle - f(x - Ky) - g(y) \\ &= \sup_{x \in \mathcal{X}, y \in \mathcal{Y}} \langle q, x - Ky \rangle + \langle K^*q, y \rangle - f(x - Ky) - g(y) = f^*(q) + g^*(K^*q) \end{aligned}$$

it follows that $h = [f^*(\cdot) + g^*(K^*\cdot)]^*$.

The proof is exactly as the proof of Lemma 4.20, but now one uses that $g^* \leq a + \delta_{B(K^*p, \varepsilon)}$ for some $a \in \mathbb{R}$ and $\varepsilon > 0$, so that $g(y) \geq -a + \langle p, Ky \rangle + \varepsilon\|y\|$ and $f^*(p) \in \mathbb{R}$ so that $f(x) \geq \langle p, x \rangle - f^*(p)$.

Then, if $x_n \rightarrow x$ and y_n is such that $f(x_n - Ky_n) + g(y_n) \leq h(x_n) + 1/n$, and if $\liminf_n h(x_n) < +\infty$, one find that along a subsequence $\|y_n\|$ is bounded, hence we may assume it converges weakly to some y (and as a consequence Ky_n converges weakly to Ky). Hence

$$h(x) \leq f(x - Ky) + g(y) \leq \liminf_n f(x_n - Ky_n) + g(y_n) \leq \liminf_n h(x_n)$$

and the semicontinuity follows. In addition, we deduce that the “inf” is in fact a “min”.

A useful application is the following: let g be convex, lsc and proper and K a continuous operator, and define

$$g^K(x) := \inf_{y: Ky=x} g(y).$$

Then, if there exists p where g^* is continuous at K^*p , g^K is lsc and $g^K = [g^*(K^*\cdot)]^*$. It is enough to apply the previous result with $f = \delta_{\{0\}}$, so that $f^* \equiv 0$ and $p \in \text{dom } f^*$.

4.3.2 Fenchel-Rockafellar duality

Consider now a minimization problem of the form

$$\min_{x \in \mathcal{X}} f(Kx) + g(x) \tag{30}$$

where $K : \mathcal{X} \rightarrow \mathcal{Y}$ is a continuous linear map and f, g are convex, lsc. Then, clearly

$$\begin{aligned} (\mathcal{P}) &= \min_x f(Kx) + g(x) = \min_x \sup_y \langle y, Kx \rangle - f^*(y) + g(x) \\ &\geq \sup_y \inf_x \langle K^*y, x \rangle + g(x) - f^*(y) = \sup_y - (g^*(-K^*y) + f^*(y)) = (\mathcal{D}) \end{aligned}$$

A natural question is when there is equality: this is true under various criteria: we will give a simple example below.

The problem “(P)” is usually called the *primal problem* and “(D)” the *dual problem* (observe though that there is a symmetry between these problems...) Notice that the *primal-dual gap*

$$\mathcal{G}(x, y) = f(Kx) + g(x) + g^*(-K^*y) + f^*(y)$$

is a measure of optimality. If it vanishes at (x^*, y^*) , then $(P) = (D)$, and (x^*, y^*) is a *saddle point* of the *Lagrangian*

$$\mathcal{L}(x, y) = \langle y, Kx \rangle - f^*(y) + g(x), \quad (31)$$

as one has

$$\mathcal{L}(x^*, y) \leq \mathcal{L}(x^*, y^*) \leq \mathcal{L}(x, y^*) \quad (32)$$

for all $x \in \mathcal{X}$, $y \in \mathcal{Y}$. [Indeed, for all y, x , $\mathcal{L}(x^*, y) \leq f(Kx^*) + g(x^*) = -f^*(y^*) - g^*(-K^*y^*) \leq \mathcal{L}(x, y^*)$.]

Theorem 4.24. *If there exists $\bar{x} \in \text{dom } g$ with f continuous at $K\bar{x}$, then $(P) = (D)$. Moreover under these assumptions, (D) has a solution.*

We show the result following a classical approach, see [15, (4.21)] for more general variants. In finite dimension, it is shown in [38, Cor 31.2.1] that equality holds if there exists $x \in \text{ri dom } g$ with $Kx \in \text{ri dom } f$, or even more generally that $0 \in \text{ri}(\text{dom } f - K \text{dom } g)$ (the proof works as below).

Proof: the method is called the “perturbation method”: We introduce, for $z \in \mathcal{Y}$,

$$\Phi(z) := \inf_{x \in \mathcal{X}} f(Kx + z) + g(x).$$

Assume $\Phi(0) > -\infty$ (otherwise there is nothing to prove), then by assumption, one can find M and ε such that for $|z| < \varepsilon$, $\Phi(z) \leq f(K\bar{x} + z) + g(\bar{x}) \leq M < +\infty$. Being Φ convex, we deduce that it is locally Lipschitz near 0 and in particular thanks to Corollary 4.16, $\Phi(0) = \Phi^{**}(0) = \sup_y -\Phi^*(y)$. We compute:

$$\begin{aligned} \Phi^*(y) &= \sup_{z \in \mathcal{Y}} \langle y, z \rangle - \inf_{x \in \mathcal{X}} (f(Kx + z) + g(x)) \\ &= \sup_{x, z} \langle y, z + Kx \rangle - \langle K^*y, x \rangle - f(Kx + z) - g(x) = f^*(y) + g^*(-K^*y). \end{aligned}$$

The claim follows. Moreover, since Φ is Lipschitz near 0 it is also subdifferentiable: there exists $y \in \partial\Phi(0)$. This subdifferential provides a solution to the “dual” problem $\max_y -\Phi^*(y)$.

Exercise: show the result in finite dimension if $0 \in \text{ri}(\text{dom } f - K \text{dom } g)$ (one needs to show again that Φ is lsc at 0).

Observe that one has by optimality in (32) that $Kx^* - \partial f^*(y^*) \ni 0$, $K^*y^* + \partial g(x^*) \ni 0$, which may be written

$$0 \in \begin{pmatrix} \partial g(x) \\ \partial f^*(y) \end{pmatrix} + \begin{pmatrix} 0 & K^* \\ -K & 0 \end{pmatrix} \begin{pmatrix} x \\ y \end{pmatrix} \quad (33)$$

meaning the solution is found by finding the “zero” of the sum of two *monotone operators* (see Section 4.4).

Example Consider the problem

$$\min_x \lambda \|Dx\|_1 + \frac{1}{2} \|x - x^0\|^2$$

where $D : \mathbb{R}^n \rightarrow \mathbb{R}^m$ is a continuous operator, $x^0 \in \mathbb{R}^n$, $\|\cdot\|_1$ is the ℓ^1 -norm. One has

$$f = \lambda \|\cdot\|_1, \quad K = D, \quad g = \frac{1}{2} \|\cdot - x^0\|^2.$$

Then the Lagrangian is

$$\mathcal{L}(x, y) = \langle y, Dx \rangle - f^*(y) + g(x)$$

where $f^*(y) = 0$ if $|y_i| \leq \lambda$ for $i = 1, \dots, n$, and $+\infty$ else. To find the dual problem, we compute $g^*(z) = \langle z, x^0 \rangle + \|z\|^2/2$, and we obtain

$$\max \left\{ \langle D^*y, x^0 \rangle - \frac{1}{2} \|D^*y\|^2 : |y_i| \leq \lambda, i = 1, \dots, n \right\}.$$

This can be rewritten as a projection problem:

$$\min_{|y_i| \leq \lambda} \|D^*y - x^0\|^2.$$

4.4 Generalization: Elements of monotone operators theory

For more results, see [6]. We mostly mention the main properties, which extend the properties shown so far for subgradients.

Observe that if f is convex, one has for all $x, y, p \in \partial f(x), q \in \partial f(y)$

$$f(y) \geq f(x) + \langle p, y - x \rangle, \quad f(x) \geq f(y) + \langle q, x - y \rangle$$

so that, summing,

$$\langle p - q, x - y \rangle \geq 0.$$

This leads to introduce the class of operators which satisfy such an inequality, which share many properties with subgradients. Consider in the Hilbert space \mathcal{X} a multi-valued operator $A : \mathcal{X} \rightarrow \mathcal{P}(\mathcal{X})$. By a slight abuse of notation, we will also denote A the graph $\{(x, y) : x \in \mathcal{X}, y \in Ax\}$.

We introduce the following definitions:

Definition 1. *The operator A is said monotone if for all $x, y \in \mathcal{X}, p \in Ax, q \in Ay$,*

$$\langle p - q, x - y \rangle \geq 0.$$

It is (μ) -strongly monotone if

$$\langle p - q, x - y \rangle \geq \mu \|x - y\|^2.$$

It is (μ) -co-coercive if

$$\langle p - q, x - y \rangle \geq \mu \|p - q\|^2.$$

It is maximal if the graph $\{(x, p) : p \in Ax\} \subset \mathcal{X} \times \mathcal{X}$ is maximal with respect to inclusion, among all the graphs of monotone operators.

In dimension 1, monotone graphs are graphs of nondecreasing functions. Obviously then, they also coincide with (sub)gradients of convex functions. In higher dimension, this is not true anymore (example: an antisymmetric linear mapping in \mathbb{R}^d , $d \geq 2$).

One sees that the subgradient of a convex function f is monotone, strongly monotone if f is strongly convex, co-coercive if ∇f is Lipschitz (cf Theorem 2.3).

A subgradient is maximal if and only if it is the subgradient of a lower-semicontinuous function. A simple proof is due to Rockafellar: if f is lsc, to show that ∂f is maximal we must show that if $x \in \mathcal{X}$ and $p \notin \partial f(x)$ then one can find y and $q \in \partial f(y)$ with $\langle p - q, x - y \rangle < 0$. Replacing f with $f(x) - \langle p, x \rangle$ we can assume that $p = 0$. Saying that $0 \notin \partial f(x)$ is precisely saying that x is not a minimizer, that is, there exists $y \in \mathcal{X}$ with $f(y) < f(x)$.

Consider now $y = \text{prox}_f(x)$, the minimizer of $f(y) + \|y - x\|^2/2$. As we have seen, $q = x - y \in \partial f(y)$. One has

$$\langle p - q, x - y \rangle = \langle -q, x - y \rangle = -\|x - y\|^2 < 0,$$

unless $y = x$. But $y = x$ would imply that $q = 0 \in \partial f(x)$, a contradiction. Hence ∂f is maximal. (The proof can be extended to non-Hilbert spaces, see [39].)

Conversely if ∂f is maximal, since $\partial f^{**} \supset \partial f$, then this operator is also the subgradient of the convex, lsc function f^{**} . We are *not* proving here that $f = f^{**}$, only that ∂f is *also* the subgradient of the convex, lsc function f^{**} .

A monotone operator is not necessarily a subgradient: for instance, in \mathbb{R}^2 , the linear operator

$$A = \begin{pmatrix} 0 & 1 \\ -1 & 0 \end{pmatrix}$$

is monotone but not the subgradient of a convex function. In order for a monotone operator to be (included in) the subgradient of a convex function, it needs to be *cyclically monotone* [37, 1]: for any $x_0, x_1, \dots, x_n = x_0$ and $p_i \in Ax_i$, $p_0 = p_n$,

$$\sum_{i=0}^{n-1} \langle p_i, x_{i+1} - x_i \rangle \leq 0.$$

An important case of monotone operator is obtained from nonexpansive (1-Lipschitz mappings) T , as in Section 3. Indeed, it is obvious to check that $I - T$ is maximal monotone:

$$\langle (x - Tx) - (x - Ty), x - y \rangle = \|x - y\|^2 - \langle Tx - Ty, x - y \rangle \geq 0$$

thanks to Cauchy-Schwartz inequality and the fact T is 1-Lipschitz.

Given A a monotone operator, its inverse is simply $A^{-1} : p \mapsto \{x : Ax \ni p\}$, with graph $\{(p, x) : p \in Ax\}$. It is therefore maximal if A is maximal, co-coercive if A is strongly monotone (cf Prop. 4.19). Clearly, $(\partial f)^{-1} = \partial f^*$ (see (24)).

Theorem 4.25 (Minty [24]). *The resolvent of a maximal-monotone operator A , defined by*

$$x \mapsto y = (I + A)^{-1}x =: J_A x \Leftrightarrow y + Ay \ni x$$

is a well (everywhere) defined single-valued nonexpansive mapping. (Conversely, for a monotone operator A if $(I + A)$ is surjective then A is maximal.)

One will see that the resolvent is also a $(1/2)$ -averaged operator (and any $(1/2)$ -averaged operator has this form).

Proof: Let us introduce the graph $G = \{(y+x, y-x) : x \in \mathcal{X}, y \in Ax\}$. If $(a, b), (a', b') \in G$, with $a = y+x, b = y-x$ and $a' = y'+x', b = y'-x'$, then

$$\begin{aligned} \|b-b'\|^2 &= \|y-y'\|^2 - 2\langle y-y', x-x' \rangle + \|y+y'\|^2 \\ &= \|a-a'\|^2 - 4\langle y-y', x-x' \rangle \leq \|a-a'\|^2 \end{aligned}$$

showing that G is the graph of a 1-Lipschitz function. Moreover, if $G' \supseteq G$ is also the graph of a 1-Lipschitz function, then defining $A' = \{((a-b)/2, (a+b)/2) : (a, b) \in G'\}$ the same computation shows that $A' \supseteq A$ is the graph of a monotone operator, hence $A' = A$ if A is maximal. (Conversely, if G is defined for all a then clearly G and therefore A are maximal, as being 1-Lipschitz G is necessarily single-valued.)

So the theorem is equivalent to the question whether a 1-Lipschitz function which is not defined in the whole of \mathcal{X} can be extended. This result (which is true only in Hilbert spaces) is known as Kirszbraun-Valentine's theorem (1935). [The proof we give is derived from [17, 2.10.43].]

The basic brick is the following extension from n to $n+1$ points:

Lemma 4.26. *If $(x_i)_{i=1}^n, (y_i)_{i=1}^n$ are points in Hilbert spaces respectively \mathcal{X}, \mathcal{Y} such that $\forall i, j, \|y_i - y_j\| \leq \|x_i - x_j\|$, then for any $x \in \mathcal{X}$ there exists $y \in \mathcal{Y}$ with $\|y_i - y\| \leq \|x_i - x\|$ for all $i = 1, \dots, n$.*

It is enough to prove this for $x = 0$: we need to find a common point to $\bar{B}(y_i, \|x_i\|)$. There is nothing to prove if $x = x_i$ for some i , so we assume $x_i \neq 0, i = 1, \dots, n$. We define

$$\bar{c} = \min \left\{ c \geq 0 : \bigcap_{i=1}^n \bar{B}(y_i, c\|x_i\|) \neq \emptyset \right\} > 0$$

(if the y_i are distinct, which we may also assume). This is a min because the closed balls are weakly compact, and we can consider y such that $\|y - y_i\| \leq \bar{c}\|x_i\|, i = 1, \dots, n$. Then we observe that y must be a convex combination of the points $(y_i)_{i \in I}$ such that $\|y - y_i\| = \bar{c}\|x_i\|$. Indeed, if not, let y' be the projection of y onto $\overline{\text{co}}\{y_i : i \in I\}$. As for any $i \in I, \langle y_i - y', y - y' \rangle \leq 0$ one has, letting $y_t = (1-t)y + ty'$, that for any $i \in I$:

$$\begin{aligned} \|y_i - y_t\|^2 &= \|y_i - y + t(y - y')\|^2 = \|y_i - y\|^2 + 2t\langle y_i - y, y - y' \rangle + t^2\|y - y'\|^2 \\ &= \|y_i - y\|^2 + 2t\langle y_i - y', y - y' \rangle - 2t\|y - y'\|^2 + t^2\|y - y'\|^2 \\ &\leq \|y_i - y\|^2 - t(2-t)\|y - y'\|^2 < \|y_i - y\|^2 \end{aligned}$$

if $t \in (0, 2)$. Hence if $t > 0$ is small enough, one sees that $\|y_i - y_t\| < \|y_i - y\| = \bar{c}\|x_i\|$ for $i \in I$, while since for $i \notin I, \|y_i - y\| < \bar{c}\|x_i\|$, one can still guarantee the same strict inequality for y_t if t is small enough. But this contradicts the definition of \bar{c} , since then there would exist $c < \bar{c}$ such that $y_t \in \bigcap_{i=1}^n \bar{B}(y_i, c\|x_i\|)$.

We therefore can write $y = \sum_{i \in I} \theta_i y_i$ as a convex combination ($\theta_i \in [0, 1], \sum_{i \in I} \theta_i =$

1). Then since $2\langle a, b \rangle = \|a\|^2 + \|b\|^2 - \|a - b\|^2$,

$$\begin{aligned}
0 &= \left\| \sum_{i \in I} \theta_i y_i - y \right\|^2 = \sum_{i, j \in I} \theta_i \theta_j \langle y_i - y, y_j - y \rangle \\
&= \frac{1}{2} \sum_{i, j \in I} \theta_i \theta_j (\|y_i - y\|^2 + \|y_j - y\|^2 - \|y_i - y_j\|^2) \\
&\geq \frac{1}{2} \sum_{i, j \in I} \theta_i \theta_j (\bar{c}^2 \|x_i\|^2 + \bar{c}^2 \|x_j\|^2 - \|x_i - x_j\|^2) \\
&= \bar{c}^2 \sum_{i, j \in I} \theta_i \theta_j \langle x_i, x_j \rangle - \frac{1 - \bar{c}^2}{2} \|x_i - x_j\|^2
\end{aligned}$$

which shows that

$$(1 - \bar{c}^2) \sum_{i, j \in I} \theta_i \theta_j \|x_i - x_j\|^2 \geq 2\bar{c}^2 \sum_{i \in I} \theta_i \|x_i\|^2$$

so that $\bar{c} \leq 1$. Hence, y satisfies $\|y - y_i\| \leq \|x_i\|$, as requested, which shows Lemma 4.26.

We can conclude the proof of Theorem 4.25: if there exists $x \in \mathcal{X}$ such that $\{x\} \times \mathcal{X} \cap G = \emptyset$, consider the set

$$K = \bigcap_{(a, b) \in G} \bar{B}(b, \|x - a\|)$$

which is an intersection of weakly compact sets.

We show that because the compact sets defining K have the “finite intersection property”, K can not be empty: Choosing $(a_0, b_0) \in G$, if $\bar{B}_0 = \bar{B}(b_0, \|x - a_0\|)$, we see that

$$K = \bar{B}_0 \cap \left(\bigcap_{(a, b) \in G} \bar{B}(b, \|x - a\|) \right)$$

hence $\bar{B}_0 \setminus K = \bar{B}_0 \cap \bigcup_{(a, b) \in G} \bar{B}(b, \|x - a\|)^c$. If this is \bar{B}_0 , by compactness one can extract a finite covering $\bigcup_{i=1}^n \bar{B}(b_i, \|x - a_i\|)^c$ for $(a_i, b_i) \in G$, $i = 1, \dots, n$. We find that

$$\bar{B}_0 \cap \bigcup_{i=1}^n \bar{B}(b_i, \|x - a_i\|)^c = \bar{B}_0$$

or equivalently that

$$\bar{B}_0 \cap \bigcap_{i=1}^n \bar{B}(b_i, \|x - a_i\|) = \emptyset$$

which contradicts Lemma 4.26. Hence, $\bar{B}_0 \setminus K \neq \bar{B}_0$ which means that $K \neq \emptyset$. Choosing $y \in K$, we find that $G \cup \{(x, y)\}$ is the graph of a 1-Lipschitz function and is strictly larger than G , which contradicts the maximality of A .

The non-expansiveness of $(I + A)^{-1}$ follows from, if $y + Ay \ni x$, $y' + Ay' \ni x'$, $p = x - y \in Ay$, $p' = x' - y' \in Ay'$:

$$\|x - x'\|^2 = \|y - y'\|^2 + 2\langle p - p', y - y' \rangle + \|p - p'\|^2 \geq \|y - y'\|^2 + \|p - p'\|^2,$$

that is, for $T = (I + A)^{-1}$:

$$\|Tx - Tx'\|^2 + \|(I - T)x - (I - T)x'\|^2 \leq \|x - x'\|^2. \quad (34)$$

An operator which satisfies (34) is *firmly non-expansive*.

Let us now consider the “*reflexion operator*”

$$R_A = 2J_A - I = 2(I + A)^{-1} - I \quad (35)$$

Lemma 4.27. *R_A is nonexpansive, and in particular, $J_A = I/2 + R_A/2$ is (1/2)-averaged.*

More generally we prove the following: *An operator T is firmly non-expansive if and only if it is 1/2-averaged, that is, $R = 2T - I$ is non-expansive (so that indeed $T = I/2 + R/2$ is 1/2-averaged).*

It follows in an obvious way from the parallelogram identity, since for any x, x' ,

$$\begin{aligned} \|Rx - Rx'\|^2 &= \|(Tx - x) - (Tx' - x') + Tx - Tx'\|^2 \\ &= 2\|(I - T)x - (I - T)x'\|^2 + 2\|Tx - Tx'\|^2 - \|x - x'\|^2 \leq \|x - x'\|^2 \\ &\Leftrightarrow \|(I - T)(x) - (I - T)(x')\|^2 + \|Tx - Tx'\|^2 \leq \|x - x'\|^2. \end{aligned}$$

We have shown that if A is maximal monotone, then $J_A = (I + A)^{-1}$ is defined everywhere and single-valued, then that it is firmly non-expansive, and eventually that an operator is firmly non-expansive if and only if it is (1/2)-averaged. We conclude by showing that if an operator $T = I/2 + R/2$ is (1/2)-averaged (R is non-expansive), then there exists a maximal monotone operator A such that $T = J_A$.

The proof follows by the same (or reverse) construction as in the beginning of the proof of Minty’s theorem: we consider the graph

$$G = \{(x + y)/2, (x - y)/2) : x \in \mathcal{X}, y = Rx\} = \{(Tx, (I - T)x) : x \in \mathcal{X}\}$$

and denote by A the corresponding operator ($y \in Ax \Leftrightarrow (x, y) \in G$). Then A is monotone: if $(\xi, \eta), (\xi', \eta') \in G$, then for some $x, x' \in \mathcal{X}$, $\xi = (x + Rx)/2$, $\eta = (x - Rx)/2$, etc., and we find:

$$\begin{aligned} \langle \xi - \xi', \eta - \eta' \rangle &= \frac{1}{4} \langle x + Rx - x' - Rx', x - Rx - x' + Rx' \rangle \\ &= \frac{1}{4} (\|x - x'\|^2 - \|Rx - Rx'\|^2) \geq 0. \end{aligned}$$

Moreover, A is maximal, if not, one could build as before from $A' \supset A$ a non-expansive graph $\{(\xi + \eta, \xi - \eta) : \eta \in A'\xi\}$ strictly larger than the graph $\{(x, Rx) : x \in \mathcal{X}\}$, which is of course impossible. By construction, $ATx \ni (I - T)x$ for all x , hence $(I + A)Tx \ni x \Leftrightarrow Tx = (I + A)^{-1}x$.

To sum up, we have shown the following result:

Theorem 4.28. *Let T be an operator, then the following are equivalent:*

- $T = (I + A)^{-1}$ for some maximal operator A ;
- T is firmly non-expansive;
- T is (1/2)-averaged ($2T - I$ is nonexpansive).

A consequence is that if $x^0 \in \mathcal{X}$ and $x^{k+1} = (I + A)^{-1}x^k$, $k \geq 0$, then $x^k \rightharpoonup x$ where x is a fixed point of $(I + A)^{-1}$, that is, $Ax = 0$, if such a point exists (Theorem 3.1). We will return soon to these iterations.

Another way to interpret Theorem 4.25 is to observe that it says that a strongly monotone maximal operator has a well-defined single-valued inverse everywhere. Indeed, if A is maximal μ -monotone, then $A' = A/\mu - I$ is maximal monotone hence $I + A'$ is surjective with single-valued inverse, and so is A . From

$$\langle p - q, x - y \rangle \geq \mu \|x - y\|^2, \quad p \in Ax, q \in Ay$$

we deduce if $B = A^{-1}$ that

$$\langle p - q, Bp - Bq \rangle \geq \mu \|Bp - Bq\|^2,$$

showing that B is co-coercive and $(1/\mu)$ -Lipschitz.

The maximal monotone operator $A_\tau = [x - (I + \tau A)^{-1}x]/\tau$ is called a *Yosida* approximation of A : it is a $(1/\tau)$ -Lipschitz-continuous mapping, with full domain (in case $A = \partial f$, $A_\tau = \nabla f_\tau$). τA_τ is firmly non-expansive, since $I - \tau A_\tau$ is. It has very important properties, see in particular Brézis' book [6]. We mention in particular Theorems 2.2, Prop. 2.5, and Cor. 2.7 in that book: the first two say that for a maximal monotone operator A , $C = \overline{\text{dom } A}$ is convex and $\lim_{\tau \rightarrow 0} J_{\tau A} x$ is the orthogonal projection of x onto C , in addition if $x \in \text{dom } A$, $A_\tau x \rightarrow A^0 x$, the element of Ax with minimal norm, while if not, $|A_\tau x| \rightarrow \infty$. The last shows that if for A, B two maximal monotone operators $\overbrace{\text{dom } A \cap \text{dom } B} \neq \emptyset$, then also $A + B$ is maximal monotone. The Yosida approximation is used in [6] to show the existence of solutions to $\dot{x} + Ax \ni 0$ for A maximal-monotone, by showing it is obtained as the limit of the solutions of $\dot{x} + A_\tau x \ni 0$ (which trivially exist because of Cauchy-Lipschitz's theorem). This allows to define properly the “gradient flow” of a convex lsc function, which is the time-continuous equivalent of the gradient descent algorithms. An exhaustive study of maximal monotone operators in Hilbert spaces is found in [2].

We will use the generalization of Moreau's identity (29):

$$x = (I + \tau A)^{-1}(x) + \tau(I + \frac{1}{\tau}A^{-1})^{-1}(\frac{x}{\tau}). \quad (36)$$

which is proved exactly in the same way as (29).

5 Algorithms. Operator splitting

We introduce here the “Forward-Backward splitting” technique. We discuss convergence rates and introduce acceleration, in particular the famous “FISTA / Nesterov acceleration”.

We also introduce other splitting: Douglas-Rachford (DR), Alternating directions method of multipliers (ADMM), Primal-Dual.

5.1 Abstract algorithms for monotone operators

In this section, we describe rapidly general algorithms for solving the equations

$$0 \in Ax \quad \text{or} \quad 0 \in Ax + Bx$$

where A, B are maximal monotone operators (sometimes subgradients, sometimes not). The idea is to generalise algorithms already seen, and then to have at hand general results which will be useful for studying more concrete algorithms.

5.1.1 Explicit algorithm

Let us first consider the equivalent of the “gradient descent”:

$$x^{k+1} = x^k - \tau p^k, p^k \in Ax^k.$$

Even if A is single-valued and continuous, then this might not converge. For instance, if $A = \begin{pmatrix} 0 & -1 \\ 1 & 0 \end{pmatrix}$ then

$$x^k = \begin{pmatrix} 1 & -\tau \\ \tau & 1 \end{pmatrix}^k x^0.$$

But the eigenvalues of this matrix are $1 + \pm\tau i$ and have modulus $\sqrt{1 + \tau^2}$, so that the iteration always diverges.

So one needs to require a further condition on A . We recall (Baillon-Haddad) that the gradient descent works for convex functions with Lipschitz gradient, whose gradient is a co-coercive monotone operator. We can show here the same:

Theorem 5.1. *Let A maximal monotone be μ -co-coercive (in particular, single-valued):*

$$\langle Ax - Ay, x - y \rangle \geq \mu \|Ax - Ay\|^2.$$

Assume there exists a solution to $Ax = 0$. Then the iteration $x^{k+1} = x^k - \tau Ax^k$ converges to x^ with $Ax^* = 0$ if $0 < \tau < 2\mu$.*

For the proof we just show that $I - \tau A$ is an averaged operator. Let us compute

$$\begin{aligned} \|(I - \tau A)x - (I - \tau A)y\|^2 + \|\tau Ax - \tau Ay\|^2 \\ = \|x - y\|^2 - 2\tau \langle x - y, Ax - Ay \rangle + 2\tau^2 \|Ax - Ay\|^2 \\ \leq \|x - y\|^2 - 2\tau(\mu - \tau) \|Ax - Ay\|^2. \end{aligned}$$

This shows that if $0 \leq \tau \leq \mu$, τA and $(I - \tau A)$ are firmly non-expansive hence $(1/2)$ -averaged. It follows that for $0 \leq \tau < 2\mu$, $(I - \tau A)$ is averaged. Hence by Theorem 3.1 the iterates weakly converge, as $k \rightarrow \infty$, to a fixed point of $(I - \tau A)$ (if it exists). If $\tau = 0$ this is not interesting, if $0 < \tau < 2\mu$, then it is a zero of A , which exists by assumption.

5.1.2 Proximal point algorithm

Then we consider the “implicit descent” $x^{k+1} \in x^k - \tau Ax^{k+1}$. This is precisely which is solved by $x^{k+1} = (I + \tau A)^{-1} x^k$, which is well-posed if A is maximal monotone (Th. 4.25). The corresponding iteration

$$x^{k+1} = (I + \tau A)^{-1} x^k$$

is known as the *proximal point algorithm*. It obviously converges to a fixed point as the operator is $(1/2)$ -averaged (if the fixed point, that is a point with $Ax = 0$, exists). Moreover, as we have seen, one can consider more generally, if $R_{\tau A} = 2(I + \tau A)^{-1} - I$,

$$x^{k+1} = (1 - \theta_k)x^k + \theta_k R_{\tau A} x^k = x^k + 2\theta_k ((I + \tau A)^{-1} x^k - x^k) = x^k - 2\theta_k \tau A x^k,$$

for $0 < \underline{\theta} \leq \theta_k \leq \bar{\theta} < 1$ and still get convergence. More generally, we prove:

Theorem 5.2 (PPA Algorithm). *Let $x^0 \in \mathcal{X}$, $\tau_k \geq \underline{\tau} > 0$, $0 \leq \underline{\lambda} \leq \lambda_k \leq \bar{\lambda} \leq 2$, and let*

$$x^{k+1} = x^k + \lambda_k((I + \tau_k A)^{-1}x^k - x^k). \quad (37)$$

If there exists x with $Ax \ni 0$, then x^k weakly converges to a zero of A .

Proof. The proof follows the lines of the proof of Thm 3.1.

A first remark is that one obviously has $\|x^{k+1} - x\|^2 \leq \|x^k - x\|^2$ for each x with $Ax \ni 0$, which is a fixed point of $J_{\tau A}$ for any τ , as in that case (37) is iterating an averaged operator with same fixed point. But we can be more precise. We have:

$$\begin{aligned} \|x^{k+1} - x\|^2 &= \|x^k - x\|^2 + \lambda_k^2 \|J_{\tau_k A} x^k - x^k\|^2 + 2\lambda_k \langle x^k - x, J_{\tau_k A} x^k - x^k \rangle \\ &= \|x^k - x\|^2 + \lambda_k^2 \|J_{\tau_k A} x^k - x^k\|^2 \\ &\quad + \lambda_k (\|J_{\tau_k A} x^k - x\|^2 - \|x^k - x\|^2 - \|J_{\tau_k A} x^k - x^k\|^2). \end{aligned}$$

Now, as $J_{\tau_k A}$ is firmly non-expansive,

$$\|J_{\tau_k A} x^k - x\|^2 + \|(I - J_{\tau_k A})x^k - (I - J_{\tau_k A})x\|^2 \leq \|x^k - x\|^2$$

where in addition $(I - J_{\tau_k A})x = 0$. Hence:

$$\begin{aligned} \|x^{k+1} - x\|^2 &\leq \|x^k - x\|^2 + \lambda_k^2 \|J_{\tau_k A} x^k - x^k\|^2 - 2\lambda_k \|J_{\tau_k A} x^k - x^k\|^2 \\ &= \|x^k - x\|^2 - \lambda_k(2 - \lambda_k) \|J_{\tau_k A} x^k - x^k\|^2. \end{aligned}$$

Letting $c = \underline{\lambda}(2 - \bar{\lambda}) > 0$, we deduce that $(x^k)_k$ is Fejér-monotone with respect to $\{x : Ax \ni 0\}$ and that

$$c \sum_{k=0}^n \|J_{\tau_k A} x^k - x^k\|^2 + \|x^{n+1} - x\|^2 \leq \|x^0 - x\|^2$$

for all $n \geq 0$, in particular $\|J_{\tau_k A} x^k - x^k\| \rightarrow 0$, as well as, by the scheme, $x^{k+1} - x^k$. We want to conclude as in the proof of Theorem 3.1. However with varying τ_k , it is not obvious that a limit point \bar{x} of a subsequence x^{k_i} is such that $A\bar{x} \ni 0$. To see this one can use the maximal-monotonicity of A . If $x' \ni \mathcal{X}$, $y' \in Ax'$, denoting $e_k := J_{\tau_k A} x^k - x^k \rightarrow 0$ we have:

$$A(x^k + e^k) \ni \frac{e^k}{\tau_k},$$

so that

$$\left\langle y' - \frac{e^k}{\tau_k}, x' - x^k - e^k \right\rangle \geq 0.$$

In the limit along the subsequence x^{k_i} , we find $\langle y', x' - \bar{x} \rangle \geq 0$, so that $A\bar{x} \ni 0$. The rest of the proof relies on Opial's lemma and is as in the proof of Theorem 3.1.

We could also consider (summable) errors. See [2] for variants, [14] for a similar proof with errors.

5.1.3 Forward-Backward splitting

We now consider a mixture of the two previous, namely the ‘‘forward-backward’’ splitting

$$x^{k+1} = (I + \tau A)^{-1}(I - \tau B)x^k \quad (38)$$

where A is maximal monotone and B μ -co-coercive. Then, as before, if $0 < \tau < 2\mu$, the algorithm is the composition of two averaged operator and converges weakly to a fixed point if it exists. We see that

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

As B is continuous, this is equivalent to $(A + B)x \ni 0$. Hence, if $A + B$ has a zero, this algorithm converges to a zero of $A + B$.

5.1.4 Douglas-Rachford splitting

This method was introduced under the following form in a paper of Lions and Mercier (79):

$$x^{k+1} = J_{\tau A}(2J_{\tau B} - I)x^k + (I - J_{\tau B})x^k \quad (39)$$

Theorem 5.3. *Let $x^0 \in \mathcal{X}$. Then if x^k defined by (39), $x^k \rightharpoonup x$ such that $w = J_{\tau B}x$ is a solution of $Aw + Bw \ni 0$ (if it exists).*

Proof: we use

$$J_{\tau A} = \frac{1}{2}I + \frac{1}{2}R_{\tau A}, \quad J_{\tau B} = \frac{1}{2}I + \frac{1}{2}R_{\tau B}.$$

Hence the operator in the algorithm is

$$\frac{1}{2}R_{\tau B} + \frac{1}{2}R_{\tau A} \circ R_{\tau B} + \left(\frac{1}{2}I - \frac{1}{2}R_{\tau B}\right) = \frac{1}{2}I + \frac{1}{2}R_{\tau A} \circ R_{\tau B}$$

so that it is $(1/2)$ -averaged (and hence a resolvent). We deduce from Thm 3.1 that the iterations converge to a fixed point, if it exists, of $R_{\tau A} \circ R_{\tau B}$. One has

$$\begin{aligned} R_{\tau A} \circ R_{\tau B}x = x &\Leftrightarrow 2J_{\tau A}(2J_{\tau B}x - x) - (2J_{\tau B}x - x) = x \Leftrightarrow J_{\tau A}(2J_{\tau B}x - x) = J_{\tau B}x \\ &\Leftrightarrow 2J_{\tau B}x - x \in J_{\tau B}x + \tau A(J_{\tau B}x) \Leftrightarrow J_{\tau B}x \in x + \tau A(J_{\tau B}x). \end{aligned}$$

Letting $w = J_{\tau B}x$, we see that w satisfies

$$w \in w + \tau Bw + \tau Aw$$

hence $Aw + Bw = 0$. Conversely, if w satisfies this equation and $x = w + Bw = w - Aw$, we see that x is a fixed point. We know, then, by Theorem 3.1, that $x^k \rightharpoonup x$. Then $w = J_{\tau B}x$ is a solution of $Aw + Bw \ni 0$. Further conditions on A, B ensuring that $J_{\tau B}x^k$ converges to a solution are found in [23], variants with errors in [14].

The iterations $x^{k+1} = R_{\tau A}R_{\tau B}x^k$ are known as the *Peaceman-Rachford* splitting algorithm and converge under some conditions to the same point.

5.1.5 Three-operators splitting

This approach, introduced in [12], generalizes the two previous methods. Given A, B, C three maximal-monotone operators with C co-coercive: for all $x, y \in \mathcal{X}$:

$$\langle Cx - Cy, x - y \rangle \geq \gamma \|x - y\|^2,$$

one wants to find $\xi \in \mathcal{X}$ such that $A\xi + B\xi + C\xi \ni 0$, and we assume there is at least a solution. One introduces for $\tau > 0$:

$$T_{\tau} := I - J_{\tau B} + J_{\tau A} \circ (2J_{\tau B} - I - \tau C \circ J_{\tau B}).$$

We observe that if A or B is 0, T_{τ} is similar to a forward-backward algorithm, while if $C = 0$, it reduces to the previous Douglas-Rachford operator.

The following is easy:

Lemma 5.4. *A point x is a fixed point of T_τ if and only if $\xi = J_{\tau B}x$ satisfies $A\xi + B\xi + C\xi \ni 0$.*

Hence, given ξ which solves $A + B + C \ni 0$, any point $x \in \xi + \tau B\xi$ is a fixed point of T_τ . The main result in [12] is then the following:

Theorem 5.5. *For $0 < \tau < 2\gamma$, T_τ is averaged.*

As a consequence, for such values of τ , the algorithm given by $x^0 \in \mathcal{X}$, $x^{k+1} = T_\tau x^k$, $k \geq 0$, produces a sequence which weakly converges to a fixed point x , such that $J_\tau x$ solves the problem.

Proof of the theorem: First, we have seen already that if $\tau < 2\gamma$, $(I - \gamma C)$ is averaged, and more precisely there exists S nonexpansive such that

$$I - \gamma C = (1 - \theta)I + \theta S =: S_\theta$$

for $\theta = \tau/(2\gamma)$. In addition, one can write $J_{\tau B} = (I + R_{\tau B})/2$ and $J_{\tau A} = (I + R_{\tau A})/2$. Hence,

$$\begin{aligned} T_\tau &= I - J_{\tau B} + J_{\tau A} \circ (J_{\tau B} - I + S_\theta \circ J_{\tau B}) \\ &= \frac{1}{2}(I - J_{\tau B} + S_\theta J_{\tau B}) + \frac{1}{2}R_{\tau A} \circ (J_{\tau B} - I + S_\theta J_{\tau B}). \end{aligned}$$

This can be written:

$$\begin{aligned} T_\tau &= \frac{1-\theta}{2}I + \frac{\theta}{2}(I - J_{\tau B} + S J_{\tau B}) + \frac{1}{2}R_{\tau A}((1 - \theta)R_{\tau B} + \theta(J_{\tau B} - I + S J_{\tau B})) \\ &= (1 - \frac{1+\theta}{2})I + \frac{1+\theta}{2}\tilde{T} \end{aligned}$$

with

$$\tilde{T} = \frac{\theta}{1+\theta}(I - J_{\tau B} + S J_{\tau B}) + \frac{1}{1+\theta}R_{\tau A}((1 - \theta)R_{\tau B} + \theta(J_{\tau B} - I + S J_{\tau B})).$$

Then, for $x, y \in \mathcal{X}$ we have:

$$\begin{aligned} \|\tilde{T}x - \tilde{T}y\|^2 &\leq \frac{\theta}{1+\theta}\|(I - J_{\tau B} + S J_{\tau B})x - (I - J_{\tau B} + S J_{\tau B})y\|^2 \\ + \frac{1}{1+\theta}\|R_{\tau A}((1 - \theta)R_{\tau B} + \theta(J_{\tau B} - I + S J_{\tau B}))x - R_{\tau A}((1 - \theta)R_{\tau B} + \theta(J_{\tau B} - I + S J_{\tau B}))y\|^2 \\ &\leq \frac{\theta}{1+\theta}\|(I - J_{\tau B} + S J_{\tau B})x - (I - J_{\tau B} + S J_{\tau B})y\|^2 \\ + \frac{1}{1+\theta}\|((1 - \theta)R_{\tau B} + \theta(J_{\tau B} - I + S J_{\tau B}))x - ((1 - \theta)R_{\tau B} + \theta(J_{\tau B} - I + S J_{\tau B}))y\|^2 \end{aligned}$$

where we have used that $R_{\tau A}$ is 1-Lipschitz. In addition,

$$\begin{aligned} &\|((1 - \theta)R_{\tau B} + \theta(J_{\tau B} - I + S J_{\tau B}))x - ((1 - \theta)R_{\tau B} + \theta(J_{\tau B} - I + S J_{\tau B}))y\|^2 \\ &\leq (1 - \theta)\|R_{\tau B}x - R_{\tau B}y\|^2 + \theta\|(J_{\tau B} - I + S J_{\tau B})x - (J_{\tau B} - I + S J_{\tau B})y\|^2 \\ &\leq (1 - \theta)\|x - y\|^2 + \theta\|(J_{\tau B} - I + S J_{\tau B})x - (J_{\tau B} - I + S J_{\tau B})y\|^2 \end{aligned}$$

using now that $R_{\tau B}$ is 1-Lipschitz. In the end we obtain:

$$\begin{aligned} \|\tilde{T}x - \tilde{T}y\|^2 &\leq \frac{\theta}{1+\theta}\|(I - J_{\tau B} + S J_{\tau B})x - (I - J_{\tau B} + S J_{\tau B})y\|^2 \\ &\quad + \frac{\theta}{1+\theta}\|(J_{\tau B} - I + S J_{\tau B})x - (J_{\tau B} - I + S J_{\tau B})y\|^2 + \frac{1-\theta}{1+\theta}\|x - y\|^2 \end{aligned}$$

We conclude with the parallelogram identity which shows that

$$\begin{aligned} & \|(I - J_{\tau B} + SJ_{\tau B})x - (I - J_{\tau B} + SJ_{\tau B})y\|^2 \\ & \quad + \|(J_{\tau B} - I + SJ_{\tau B})x - (J_{\tau B} - I + SJ_{\tau B})y\|^2 \\ & = 2(\|(I - J_{\tau B})x - (I - J_{\tau B})y\|^2 + \|SJ_{\tau B}x - SJ_{\tau B}y\|^2) \\ & \leq 2(\|(I - J_{\tau B})x - (I - J_{\tau B})y\|^2 + \|J_{\tau B}x - J_{\tau B}y\|^2) \leq 2\|x - y\|^2 \end{aligned}$$

since S is 1-Lipschitz and since $J_{\tau B}$ is firmly non expansive. Hence,

$$\|\tilde{T}x - \tilde{T}y\|^2 \leq \frac{2\theta}{1+\theta}\|x - y\|^2 + \frac{1-\theta}{1+\theta}\|x - y\|^2 = \|x - y\|^2$$

showing that T_τ is $(1 + \theta)/2$ -averaged.

Remark 5.6. The averaging here is not as good as the one found in [12], which is $1/(2 - \theta)$.

5.2 Descent algorithms, acceleration, “FISTA”

5.2.1 Forward-Backward descent

In case $A = \partial g$ and $B = \nabla f$, algorithm (38), which aims at finding a point x where $\partial g(x) + \nabla f(x) \ni 0$, or equivalently a minimizer of

$$\min_{x \in \mathcal{X}} F(x) := f(x) + g(x) \quad (40)$$

where g is, a “simple” convex lsc function and f is a convex function with Lipschitz gradient. The basic idea of the Forward-Backward splitting scheme (FBS) is to combine an explicit step of descent in the smooth part f with a implicit step of descent in g . It iterates the operator:

$$\bar{x} \mapsto \hat{x} = T_\tau \bar{x} := \text{prox}_{\tau g}(\bar{x} - \tau \nabla f(\bar{x})) = (I + \tau \partial g)^{-1}(\bar{x} - \tau \nabla f(\bar{x})). \quad (41)$$

Another name found in the literature [27] is the “composite gradient” descent, as one may see here $(\hat{x} - \bar{x})/\tau$ as a generalised gradient for F at \bar{x} . The essential reason why all this is reasonable is that clearly, a fixed point $\hat{x} = \bar{x}$ will satisfy the Euler Lagrange equations $\nabla f(\bar{x}) + \partial g(\bar{x}) \ni 0$ of (40). Observe that in the particular case where $g = \delta_C$ is the characteristic function of a closed, convex set C , then $\text{prox}_{\tau g}(x)$ reduces to $\Pi_C(x)$ (the orthogonal projection onto C) and the mapping T_τ defines a projected gradient descent method.

Algorithm 1 Forward-Backward descent with fixed step

Choose $x_0 \in \mathcal{X}$
for all $k \geq 0$ **do**

$$x^{k+1} = T_\tau x^k = \text{prox}_{\tau g}(x^k - \tau \nabla f(x^k)). \quad (42)$$

end for

The theoretical convergence rate of the plain FBS descent is not very good, as one can merely show the same as for the gradient descent:

Theorem 5.7. Let $x^0 \in \mathcal{X}$ and x^k be recursively defined by (42), with $\tau \leq 1/L$. Then not only x^k converges to a minimiser, but one has the rates

$$F(x^k) - F(x^*) \leq \frac{1}{2^k} \|x^* - x^0\|^2 \quad (43)$$

where x^* is any minimiser of f . If in addition f or g are strongly convex with parameters μ_f, μ_g (with $\mu = \mu_f + \mu_g > 0$), one has

$$F(x^k) - F(x^*) + \frac{1+\tau\mu_g}{2\tau} \|x^k - x^*\|^2 \leq \omega^k \frac{1+\tau\mu_g}{2\tau} \|x^0 - x^*\|^2. \quad (44)$$

where $\omega = (1 - \tau\mu_f)/(1 + \tau\mu_g)$.

However, its behaviour is improved if the objective is smoother than actually known, moreover, it is quite robust to perturbations and can be overrelaxed, see in particular [10].

5.2.2 FISTA

An “optimal” accelerated version is also available for this method, cf Section 2.4.3. This is well described in [29], [27], although a somewhat simpler proof is found in [3], where the algorithm, in the cases where $\mu = \mu_f + \mu_g = 0$, is called “FISTA”. The general iteration takes the form:

Algorithm 2 FISTA with fixed step

Choose $x^0 = x^{-1} \in \mathcal{X}$ and $t_0 \geq 0$
for all $k \geq 0$ **do**

$$y^k = x^k + \beta_k(x^k - x^{k-1}) \quad (45)$$

$$x^{k+1} = T_\tau y^k = \text{prox}_{\tau g}(y^k - \tau \nabla f(y^k)) \quad (46)$$

where, in case $\mu = 0$,

$$t_{k+1} = \frac{1 + \sqrt{1 + 4t_k^2}}{2} \geq \frac{k+1}{2}, \quad (47)$$

$$\beta_k = \frac{t_k - 1}{t_{k+1}}, \quad (48)$$

and if $\mu = \mu_f + \mu_g > 0$,

$$t_{k+1} = \frac{1 - qt_k^2 + \sqrt{(1 - qt_k^2)^2 + 4t_k^2}}{2}, \quad (49)$$

$$\beta_k = \frac{t_k - 1}{t_{k+1}} \frac{1 + \tau\mu_g - t_{k+1}\tau\mu}{1 - \tau\mu_f}, \quad (50)$$

where $q = \tau\mu/(1 + \tau\mu_g) < 1$.

end for

In the latter case, we assume $L > \mu_f$, otherwise f is quadratic and the problem is trivial. The following result is then true:

Theorem 5.8. Assume $t_0 = 0$ and let x^k be generated by the algorithm, in either case $\mu = 0$ or $\mu > 0$. Then, one has the decay rate

$$F(x^k) - F(x^*) \leq \min \left\{ (1 + \sqrt{q})(1 - \sqrt{q})^k, \frac{4}{(k+1)^2} \right\} \frac{1 + \tau\mu_g}{2\tau} \|x^0 - x^*\|^2.$$

It must be mentioned that in the case $\mu = 0$, a classical choice for t_k is also $t_k = (k + 1)/2$, which gives essentially the same rate. An important issue is the stability of these rates when the proximal operators can be only evaluated inexactly — the situation here is worse than for the nonaccelerated algorithm, which has been addressed in several papers.

The proof of both Theorems 5.7 and 5.8 rely on the following essential but straightforward descent rule: let $\hat{x} = T_\tau \bar{x}$, then for all $x \in \mathcal{X}$,

$$F(x) + (1 - \tau\mu_f) \frac{\|x - \bar{x}\|^2}{2\tau} \geq \frac{1 - \tau L}{\tau} \frac{\|\hat{x} - \bar{x}\|^2}{2} + F(\hat{x}) + (1 + \tau\mu_g) \frac{\|x - \hat{x}\|^2}{2\tau}. \quad (51)$$

In particular, if $\tau L \leq 1$,

$$F(x) + (1 - \tau\mu_f) \frac{\|x - \bar{x}\|^2}{2\tau} \geq F(\hat{x}) + (1 + \tau\mu_g) \frac{\|x - \hat{x}\|^2}{2\tau}. \quad (52)$$

The proof is elementary: by definition, \hat{x} is the minimiser of the $(\mu_g + (1/\tau))$ -strongly convex function

$$x \mapsto g(x) + f(\bar{x}) + \langle \nabla f(\bar{x}), x - \bar{x} \rangle + \frac{\|x - \bar{x}\|^2}{2\tau}.$$

It follows that for all x (cf (20)):

$$\begin{aligned} F(x) + (1 - \tau\mu_f) \frac{\|x - \bar{x}\|^2}{2\tau} &\geq g(x) + f(\bar{x}) + \langle \nabla f(\bar{x}), x - \bar{x} \rangle + \frac{\|x - \bar{x}\|^2}{2\tau} \\ &\geq g(\hat{x}) + f(\bar{x}) + \langle \nabla f(\bar{x}), \hat{x} - \bar{x} \rangle + \frac{\|\hat{x} - \bar{x}\|^2}{2\tau} + (1 + \tau\mu_g) \frac{\|x - \hat{x}\|^2}{2\tau}. \end{aligned}$$

But since ∇f is L -Lipschitz, $f(\bar{x}) + \langle \nabla f(\bar{x}), \hat{x} - \bar{x} \rangle \geq f(\hat{x}) - (L/2)\|\hat{x} - \bar{x}\|^2$ so that equation (51) follows.

Remark 5.9. One can more precisely deduce from this computation that

$$F(x) + (1 - \tau\mu_f) \frac{\|x - \bar{x}\|^2}{2\tau} \geq F(\hat{x}) + (1 + \tau\mu_g) \frac{\|x - \hat{x}\|^2}{2\tau} + \left(\frac{\|\hat{x} - \bar{x}\|^2}{2\tau} - D_f(\hat{x}, \bar{x}) \right). \quad (53)$$

where $D_f(x, y) := f(x) - f(y) - \langle \nabla f(y), x - y \rangle \leq (L/2)\|x - y\|^2$ is the ‘‘Bregman f -distance’’ from y to x [5]. In particular, (52) holds as soon as

$$D_f(\hat{x}, \bar{x}) \leq \frac{\|\hat{x} - \bar{x}\|^2}{2\tau}$$

which is always true if $\tau \leq 1/L$ but might also occur in other situations, and in particular, be tested ‘‘on the fly’’ during the iterations. This allows to implement efficient backtracking strategies ‘à la’ Armijo for the algorithms described in this section when the Lipschitz constant of f is not a priori known.

Remark 5.10. Observe that if $X \subset \mathcal{X}$ is a closed convex set containing the domain of F , and on which the projection Π_X can be computed, then the same inequality (52) holds if $\hat{x} = T_\tau \Pi_X \bar{x}$ (requiring only that ∇f is Lipschitz on X). This means that the same rates are valid if one replaces (45) with

$$y^k = \Pi_X(x^k + \beta_k(x^k - x^{k-1}))$$

which is feasible if X is the domain of F .

5.2.3 Convergence rates

Unaccelerated scheme We start with the rates of the unaccelerated FB descent scheme and prove Theorem 5.7.

First, if $\mu_f = \mu_g = 0$: we start from inequality (52), letting, for $k \geq 0$, $\bar{x} = x^k$ and $\hat{x} = x^{k+1}$. It follows that for any x :

$$F(x) + \frac{\|x - x^k\|^2}{2\tau} \geq F(x^{k+1}) + \frac{\|x - x^{k+1}\|^2}{2\tau}.$$

Choosing $x = x^k$ shows that $F(x^k)$ is nonincreasing. Summing then these inequalities from $k = 0$ to $n - 1$, $n \geq 1$ yields

$$\sum_{k=1}^n (F(x^k) - F(x)) + \sum_{k=1}^n \frac{1}{2\tau} \|x - x^k\|^2 \leq \sum_{k=0}^{n-1} \frac{1}{2\tau} \|x - x^k\|^2.$$

After cancellations and using $F(x^k) \geq F(x^n)$ for $k = 0, \dots, n$, it remains just

$$n(F(x^n) - F(x)) + \frac{1}{2\tau} \|x - x^n\|^2 \leq \frac{1}{2\tau} \|x - x^0\|^2$$

so that, in particular $F(x^n) - F(x^*) \leq \|x^* - x^0\|^2 / (2n\tau)$.

Now, if $\mu_f > 0$ or $\mu_g > 0$ we can improve this computation: we now have for any x :

$$F(x) + (1 - \tau\mu_f) \frac{\|x - x^k\|^2}{2\tau} \geq F(x^{k+1}) + (1 + \tau\mu_g) \frac{\|x - x^{k+1}\|^2}{2\tau}.$$

Choosing $x = x^k$ shows that $F(x^k)$ is nonincreasing. Letting

$$\omega = \frac{1 - \tau\mu_f}{1 + \tau\mu_g} \leq 1, \tag{54}$$

and summing these inequalities from $k = 0$ to $n - 1$, $n \geq 1$, after multiplication by ω^{-k-1} , yields

$$\sum_{k=1}^n \omega^{-k} (F(x^k) - F(x)) + \sum_{k=1}^n \omega^{-k} \frac{1 + \tau\mu_g}{2\tau} \|x - x^k\|^2 \leq \sum_{k=0}^{n-1} \omega^{-k-1} \frac{1 - \tau\mu_f}{2\tau} \|x - x^k\|^2.$$

After cancellations and using $F(x^k) \geq F(x^n)$ for $k = 0, \dots, n$, we get

$$\omega^{-n} \left(\sum_{k=0}^{n-1} \omega^k \right) (F(x^n) - F(x)) + \omega^{-n} \frac{1 + \tau\mu_g}{2\tau} \|x - x^n\|^2 \leq \frac{1 - \tau\mu_f}{2\tau} \|x - x^0\|^2.$$

We deduce, in case $\mu = \mu_f + \mu_g > 0$ so that $\omega < 1$,

$$F(x^k) - F(x^*) + \frac{1 + \tau\mu_g}{2\tau} \|x^k - x^*\|^2 \leq \omega^k \frac{1 + \tau\mu_g}{2\tau} \|x^0 - x^*\|^2. \tag{55}$$

which is a “linear convergence rate” (however we will see that one can do better).

Convergence rates for FISTA Now we show the accelerated convergence rates. The basic idea consists in first choosing in (52) a generic point of the form $((t - 1)x^k + x)/t$, $t \geq 1$, which is a convex combination of the iterate x^k and another generic point (in

practice a minimizer) x . We find after some calculation (systematically using the strong convexity inequalities when possible)

$$\begin{aligned} t(t-1)(F(x^k) - F(x)) - \mu \frac{t-1}{2} \|x - x^k\|^2 \\ + (1 - \tau\mu_f) \frac{\|(t-1)x^k + x - ty^k\|^2}{2\tau} \\ \geq t^2(F(x^{k+1}) - F(x)) + (1 + \tau\mu_g) \frac{\|(t-1)x^k + x - tx^{k+1}\|^2}{2\tau}. \end{aligned} \quad (56)$$

Consider first the case where $\mu = \mu_f + \mu_g = 0$. Then we have

$$\begin{aligned} t^2(F(x^{k+1}) - F(x)) + \frac{\|(t-1)x^k + x - tx^{k+1}\|^2}{2\tau} \\ \leq t(t-1)(F(x^k) - F(x)) + \frac{\|(t-1)x^k + x - ty^k\|^2}{2\tau}. \end{aligned}$$

We see that the term $F(x^k) - F(x)$ is “shrunk” at each step by a factor $(t-1)/t < 1$, while the other term is not. How can we exploit this?

The basic idea in the proof is to use a *variable* parameter $t = t_{k+1}$, and choose y^k to ensure that the term $(t_{k+1} - 1)x^k + x - t_{k+1}y^k$ in the right hand side is the same as the term $(t_{k+1} - 1)x^k + x - t_{k+1}x^{k+1}$ of the left hand side at the *previous* iterate, that is,

$$(t_{k+1} - 1)x^k + x - t_{k+1}y^k = (t_k - 1)x^{k-1} + x - t_kx^k$$

so that if we sum the inequalities for $k = 0, \dots, n$ the norms will cancel. Hence, we choose:

- $t_{k+1}(t_{k+1} - 1) = t_k^2$;
- $y^k = x^k + \beta_k(x^k - x^{k-1})$ with $\beta_k = (t_k - 1)/t_{k+1}$;

we obtain the recursion

$$\begin{aligned} t_{k+1}^2(F(x^{k+1}) - F(x)) + \frac{\|(t_{k+1} - 1)x^k + x - t_{k+1}x^{k+1}\|^2}{2\tau} \\ \leq t_k^2(F(x^k) - F(x)) + \frac{\|(t_k - 1)x^{k-1} + x - t_kx^k\|^2}{2\tau}. \end{aligned}$$

which we can sum from $k = 0, \dots, n-1$ to obtain

$$F(x^n) - F(x) + \frac{1}{2t_n^2\tau} \|(t_{k+1} - 1)x^k + x - t_{k+1}x^{k+1}\|^2 \leq \frac{1}{2t_n^2\tau} \|x^0 - x\|^2.$$

Observe that $t_{k+1}^2 - t_{k+1} - t_k^2 = 0$ yields $t_{k+1} = (1 + \sqrt{1 + 4t_k^2})/2$ (one can choose $t_0 = 0, t_1 = 1$), and in particular $t_{k+1} \geq 1/2 + t_k \geq (k+1)/2$ for $k \geq 1$, by induction. Therefore, choosing $x = x^*$,

$$F(x^n) - F(x^*) \leq \frac{2}{2(n+1)^2\tau} \|x^0 - x\|^2. \quad (57)$$

An important remark is that, if one takes $x = x^*$, $F(x^k) - F(x^*) \geq 0$ so that in fact one can get the same inequalities if one only ensures $t_{k+1}(t_{k+1} - 1) \leq t_k^2$, and not =.

For instance, the sequence $t_0 = 0, t_k = (k+1)/2$ for $k \geq 1$ is admissible and yields the same rate.

It can be interesting to take slightly smaller t_k , such as $(k+1)/\alpha$ for $\alpha > 2$. One can show in particular the convergence of the iterates (x^k) to a solution in this case, while it is still an open problem for $\alpha = 2$. It has been even observed by Charles Dossal (U. Bordeaux) that in that case, one can show that

$$F(x^n) - F(x^*) = o\left(\frac{1}{n^2}\right)$$

which *does not* contradict the lower bound (6).

Convergence rates for FISTA, strongly convex case We start again from (56) but now we assume that $\mu = \mu_f + \mu_g > 0$. Then, we observe that

$$\begin{aligned} & -\mu \frac{t-1}{2} \|x - x^k\|^2 + (1 - \tau\mu_f) \frac{\|x - x^k + t(x^k - y^k)\|^2}{2\tau} \\ &= (1 - \tau\mu_f - \mu\tau(t-1)) \frac{\|x - x^k\|^2}{2\tau} + \frac{1 - \tau\mu_f}{\tau} t \langle x - x^k, x^k - y^k \rangle + t^2(1 - \tau\mu_f) \frac{\|x^k - y^k\|^2}{2\tau} \\ &= \frac{(1 + \tau\mu_g - t\mu\tau)}{2\tau} \|x - x^k + t \frac{1 - \tau\mu_f}{1 + \tau\mu_g - t\mu\tau} (x^k - y^k)\|^2 + t^2(1 - \tau\mu_f) \left(1 - \frac{1 - \tau\mu_f}{1 + \tau\mu_g - t\mu\tau}\right) \frac{\|x^k - y^k\|^2}{2\tau} \\ &= \frac{(1 + \tau\mu_g - t\mu\tau)}{2\tau} \|x - x^k + t \frac{1 - \tau\mu_f}{1 + \tau\mu_g - t\mu\tau} (x^k - y^k)\|^2 - t^2(t-1) \frac{\tau\mu(1 - \tau\mu_f)}{1 + \tau\mu_g - t\mu\tau} \frac{\|x^k - y^k\|^2}{2\tau}. \end{aligned}$$

It follows that for any $x \in \mathcal{X}$,

$$\begin{aligned} & t(t-1)(F(x^k) - F(x)) + (1 + \tau\mu_g - t\mu\tau) \frac{\|x - x^k - t \frac{1 - \tau\mu_f}{1 + \tau\mu_g - t\mu\tau} (y^k - x^k)\|^2}{2\tau} \\ & \geq t^2(F(x^{k+1}) - F(x)) + (1 + \tau\mu_g) \frac{\|x - x^{k+1} - (t-1)(x^{k+1} - x^k)\|^2}{2\tau} \\ & \quad + t^2(t-1) \frac{\tau\mu(1 - \tau\mu_f)}{1 + \tau\mu_g - t\mu\tau} \frac{\|x^k - y^k\|^2}{2\tau}. \end{aligned} \quad (58)$$

We let $t = t_{k+1}$ above, then we can get a useful recursion if we let

$$\omega_k = \frac{1 + \tau\mu_g - t_{k+1}\mu\tau}{1 + \tau\mu_g} = 1 - t_{k+1} \frac{\mu\tau}{1 + \tau\mu_g} \in [0, 1] \quad (59)$$

$$t_{k+1}(t_{k+1} - 1) \leq \omega_k t_k^2, \quad (60)$$

$$\beta_k = \frac{t_k - 1}{t_{k+1}} \frac{1 + \tau\mu_g - t_{k+1}\mu\tau}{1 - \tau\mu_f} = \omega_k \frac{t_k - 1}{t_{k+1}} \frac{1 + \tau\mu_g}{1 - \tau\mu_f}, \quad (61)$$

$$y^k = x^k + \beta_k (x^k - x^{k-1}) \quad (62)$$

Denoting $\alpha_k = 1/t_k$ and

$$q = \frac{\tau\mu}{1 + \tau\mu_g} = \frac{\tau\mu_f + \tau\mu_g}{1 + \tau\mu_g} < 1, \quad (63)$$

one finds the same rules as in formula (2.2.9), p. 80 in [29] (with the minor difference that here we may chose $t_0 = 0, t_1 = 1$, and we have shifted the numbering of the

sequences $(x^k), (y^k)$. In this case, we find

$$\begin{aligned} & t_{k+1}^2(F(x^{k+1}) - F(x)) + \frac{1 + \tau\mu g}{2\tau} \|x - x^{k+1} - (t_{k+1} - 1)(x^{k+1} - x^k)\|^2 \\ & \leq \omega_k \left(t_k^2(F(x^k) - F(x)) + \frac{1 + \tau\mu g}{2\tau} \|x - x^k - (t_k - 1)(x^k - x^{k-1})\|^2 \right) \end{aligned}$$

so that

$$\begin{aligned} & t_k^2(F(x^k) - F(x)) + \frac{1 + \tau\mu g}{2\tau} \|x - x^k - (t_k - 1)(x^k - x^{k-1})\|^2 \\ & \leq \left(\prod_{n=0}^{k-1} \omega_n \right) \left[t_0^2(F(x^0) - F(x)) + \frac{1 + \tau\mu g}{2\tau} \|x - x^0\|^2 \right] \end{aligned} \quad (64)$$

The update rule for t_k reads

$$t_{k+1}(t_{k+1} - 1) = (1 - qt_{k+1})t_k^2, \quad (65)$$

so that,

$$t_{k+1} = \frac{1 - qt_k^2 + \sqrt{(1 - qt_k^2)^2 + 4t_k^2}}{2}. \quad (66)$$

We need to make sure that $qt_{k+1} \leq 1$ so that (59) holds. This is proved exactly as in the proof of Lemma 2.2.4 in [29]. Assuming (as in [29]) that $\sqrt{q}t_k \leq 1$, we observe that (65) yields

$$qt_{k+1}^2 = qt_{k+1} + (1 - qt_{k+1})qt_k^2.$$

If $qt_{k+1} \geq 1$, it yields $qt_{k+1}^2 \leq qt_{k+1}$ hence $qt_{k+1} \leq q < 1$, a contradiction. Hence $qt_{k+1} < 1$ and we obtain that qt_{k+1}^2 is a convex combination of 1 and qt_k^2 , so that $\sqrt{q}t_{k+1} \leq 1$. We have shown that as soon as $\sqrt{q}t_0 \leq 1$ (which we will now assume), $\sqrt{q}t_k \leq 1$ for all k . Eventually, we also observe that

$$t_{k+1}^2 = (1 - qt_k^2)t_{k+1} + t_k^2$$

showing that t_k is an increasing sequence. It remains to estimate the factor

$$\theta_k = t_k^{-2} \prod_{n=0}^{k-1} \omega_n \quad (k \geq 1).$$

From (60) (with an equality) we find

$$1 - \frac{1}{t_{k+1}} = \omega_k \frac{t_k^2}{t_{k+1}^2}$$

so that

$$t_0^2 \theta_k = \frac{t_0^2}{t_k^2} \prod_{n=0}^{k-1} \omega_n = \prod_{n=1}^k \left(1 - \frac{1}{t_n} \right) \leq (1 - \sqrt{q})^k$$

since $1/t_k \geq \sqrt{q}$. If $t_0 \geq 1$ it shows $\theta_k \leq (1 - \sqrt{q})^k / t_0^2$. If $t_0 \in [0, 1[$, we rather write

$$\theta_k = \frac{\omega_0}{t_k^2} \prod_{n=1}^{k-1} \omega_n = \frac{\omega_0}{t_1^2} \prod_{n=2}^k \left(1 - \frac{1}{t_n} \right)$$

and observe that (66) yields (using $2 - q \geq 1 \geq q$)

$$t_1 = \frac{1 - qt_0^2 + \sqrt{1 + 2(2 - q)t_0^2 + q^2t_0^4}}{2} \geq 1.$$

Also, $\omega_0 \leq 1 - q$ (from (59)), so that

$$\theta_k \leq (1 + \sqrt{q})(1 - \sqrt{q})^k.$$

The next step is to bound also θ_k by $O(1/k^2)$. We exactly follow Lemma 2.2.4 in [29]. In our notation, it reads

$$\frac{1}{\sqrt{\theta_{k+1}}} - \frac{1}{\sqrt{\theta_k}} = \frac{\theta_k - \theta_{k+1}}{\sqrt{\theta_k\theta_{k+1}}(\sqrt{\theta_k} + \sqrt{\theta_{k+1}})} \geq \frac{\theta_k(1 - (1 - 1/t_{k+1}))}{2\theta_k\sqrt{\theta_{k+1}}}$$

since θ_k is nonincreasing. It follows

$$\frac{1}{\sqrt{\theta_{k+1}}} - \frac{1}{\sqrt{\theta_k}} \geq \frac{1}{2t_{k+1}\sqrt{\theta_{k+1}}} = \frac{1}{2} \frac{1}{\sqrt{\prod_{n=0}^k \omega_n}} \geq \frac{1}{2},$$

showing that $1/\sqrt{\theta_k} \geq (k - 1)/2 + t_1/\sqrt{\omega_0} \geq (k + 1)/2$. Hence, provided $\sqrt{q}t_0 \leq 1$, we also find:

$$\theta_k \leq \frac{4}{(k + 1)^2}. \quad (67)$$

We have shown the following result, due to Nesterov and stated, with a different proof, in [29]:

Theorem 5.11. *If $\sqrt{q}t_0 \leq 1$, $t_0 \geq 0$, then the sequence (x^k) produced by iterations $x^k = T_\tau y^k$ with (66), (59), (61), (62) satisfies*

$$F(x^k) - F(x^*) \leq \min \left\{ \frac{(1 - \sqrt{q})^k}{t_0^2}, \frac{4}{(k + 1)^2} \right\} \left(t_0^2(F(x^0) - F(x^*)) + \frac{1 + \tau\mu_g}{2\tau} \|x^0 - x^*\|^2 \right) \quad (68)$$

if $t_0 \geq 1$, and

$$F(x^k) - F(x^*) \leq \min \left\{ (1 + \sqrt{q})(1 - \sqrt{q})^k, \frac{4}{(k + 1)^2} \right\} \left(t_0^2(F(x^0) - F(x^*)) + \frac{1 + \tau\mu_g}{2\tau} \|x^0 - x^*\|^2 \right) \quad (69)$$

if $t_0 \in [0, 1]$, where x^* is a minimiser of F .

Theorem 5.8 is a particular case of this result, for $t_0 = 0$.

Remark 5.12. (Constant steps.) In case $\mu > 0$ (which is $q > 0$), then an admissible choice which satisfies (59),(60), (61), is to take $t = 1/\sqrt{q}$, $\omega = 1 - \sqrt{q}$,

$$\beta = \omega^2 \frac{1 + \tau\mu_g}{1 - \tau\mu_f} = \frac{\sqrt{1 + \tau\mu_g} - \sqrt{\tau\mu}}{\sqrt{1 + \tau\mu_g} + \sqrt{\tau\mu}}.$$

Then, (68) becomes

$$F(x^k) - F(x^*) \leq (1 - \sqrt{q})^k \left(F(x^0) - F(x^*) + \mu \frac{\|x^0 - x^*\|^2}{2} \right).$$

Remark 5.13. (Monotone ‘‘FISTA’’, monotone algorithms.) The algorithms studied here are not necessarily ‘‘monotone’’ in the sense that the objective F is not always nonincreasing. A workaround implemented in various papers [43, 3] consists in choosing for x^{k+1} any point with $F(x^{k+1}) \leq F(T_\tau y^k)$ ⁷, which will not change much (56) except

⁷this makes sense only if the evaluation of F is easy and does not take too much time.

that in the last term, x^{k+1} should be replaced with $T_\tau y^k$. Then, the same computations carry on, and it is enough to replace the update rule (62) for y^k with

$$\begin{aligned} y^k &= x^k + \beta_k(x^k - x^{k-1}) + \omega_k \frac{t_k}{t_{k+1}} \frac{1+\tau\mu_g}{1-\tau\mu_f} (T_\tau y^{k-1} - x^k) \\ &= x^k + \beta_k \left((x^k - x^{k-1}) + \frac{t_k}{t_{k+1}} (T_\tau y^{k-1} - x^k) \right) \end{aligned} \quad (62')$$

to obtain the same rates of convergence. The most sensible choice for x^{k+1} is to take $T_\tau y^k$ if $F(T_\tau y^k) \leq F(x^k)$, and x^k else, in which case one of the two terms ($x^k - x^{k-1}$ or $T_\tau y^{k-1} - x^k$) vanishes in (62').

Conclusion: compare the geometric rate (54) with $\omega = 1 - \sqrt{q}$ for q given by (63), what do we observe?

5.3 ADMM, Douglas-Rachford splitting

We now consider a class of method which solves another kind of problem, namely of the form

$$\min_{Ax+By=\zeta} f(x) + g(y) \quad (70)$$

where in practice one will ask that the convex, lsc functions f, g are “simple” (and even more than this).

Observe that if f^* is continuous at some point A^*p and if g^* is continuous at some point B^*q , cf Section. 4.3.1 (or, in finite dimension, if $A^*p \in \text{ri dom } f^*$, $B^*q \in \text{ri dom } g^*$), we can define

$$\tilde{f}(\xi) = \min_{Ax=\xi} f(x), \quad \tilde{g}(\eta) = \min_{By=\eta} g(y),$$

moreover the min is reached in both cases.

Then, one has $\tilde{f}^*(p) = f(A^*p)$, $\tilde{g}^*(q) = g(B^*q)$ and the problem reads

$$\min_{\xi} \tilde{f}(\xi) + \tilde{g}(\zeta - \xi);$$

it can be seen as an inf-convolution problem. Moreover Corollary 4.21 shows that the value of (70) is also

$$\sup_p \langle \zeta, p \rangle - f^*(A^*p) - g^*(B^*p) \quad (71)$$

which gives a dual form for (70).

An “augmented Lagrangian” approach for (70) consists in introducing the constraint in the form

$$\min_{x,y} \sup_z f(x) + g(y) - \langle z, Ax + By - \zeta \rangle + \frac{\gamma}{2} \|Ax + By - \zeta\|^2$$

which we observe is equivalent (as the sup is $+\infty$ if $Ax + By \neq \zeta$).

If we introduce the function

$$\mathcal{D}(z) = \inf_{x,y} f(x) + g(y) - \langle z, Ax + By - \zeta \rangle + \frac{\gamma}{2} \|Ax + By - \zeta\|^2$$

we find that, denoting \bar{x}, \bar{y} the solution of the problem for z and \bar{x}_h, \bar{y}_h the solution for $z + h$ (the min is reached, why?),

$$\begin{aligned} \mathcal{D}(z) &= f(\bar{x}) + g(\bar{y}) - \langle z + h, A\bar{x} + B\bar{y} - \zeta \rangle + \frac{\gamma}{2} \|A\bar{x} + B\bar{y} - \zeta\|^2 + \langle h, A\bar{x} + B\bar{y} - \zeta \rangle \\ &\geq f(\bar{x}_h) + g(\bar{y}_h) - \langle z + h, A\bar{x}_h + B\bar{y}_h - \zeta \rangle + \frac{\gamma}{2} \|A\bar{x}_h + B\bar{y}_h - \zeta\|^2 \\ &\quad + \frac{\gamma}{2} \|A(\bar{x} - \bar{x}_h) + B(\bar{y} - \bar{y}_h)\|^2 + \langle h, A\bar{x} + B\bar{y} - \zeta \rangle \end{aligned}$$

where we have used the strong convexity of the norm with respect to $Ax + By$. We find

$$\mathcal{D}(z) - \langle h, A\bar{x} + B\bar{y} - \zeta \rangle \geq \mathcal{D}(z + h) + \frac{\gamma}{2} \|(A\bar{x} + B\bar{y} - \zeta) - (A\bar{x}_h + B\bar{y}_h - \zeta)\|^2$$

which shows that $\zeta - A\bar{x} - B\bar{y} \in \partial^+ \mathcal{D}(z)$ (the super-gradient of the concave function \mathcal{D} at z) and that $z \mapsto A\bar{x} + B\bar{y} - \zeta$ is γ -co-coercive, and $(1/\gamma)$ -Lipschitz. Hence a natural algorithm, known as ‘‘augmented Lagrangian’’, consists in iteratively solving

$$\begin{cases} (x^{k+1}, y^{k+1}) = \arg \min_{x,y} f(x) + g(y) - \langle z^k, Ax + By - \zeta \rangle + \frac{\gamma}{2} \|Ax + By - \zeta\|^2, \\ z^{k+1} = z^k + \gamma(\zeta - Ax^{k+1} - By^{k+1}) : \end{cases} \quad (72)$$

it is precisely a gradient ascent with fixed step for the concave function \mathcal{D} , and will converge (it should be shown then that also x^k, y^k converge to a solution).

Unfortunately, this algorithm is usually not implementable, as the joint minimization step cannot in general be performed. This is why it was proposed [20, 19] to perform these minimizations alternatively instead than simultaneously, see Algorithm 3

Algorithm 3 ADMM

Choose $\gamma > 0, y^0, z^0$.

for all $k \geq 0$ **do**

Find x^{k+1} by minimising $x \mapsto f(x) - \langle z^k, Ax \rangle + \frac{\gamma}{2} \|\zeta - Ax - By^k\|^2$,

Find y^{k+1} by minimising $y \mapsto g(y) - \langle z^k, By \rangle + \frac{\gamma}{2} \|\zeta - Ax^{k+1} - By\|^2$,

Update $z^{k+1} = z^k + \gamma(\zeta - Ax^{k+1} - By^{k+1})$.

end for

We will relate this approach to other known converging algorithms. Then in a next section, we will show how we can derive rates of convergence for this approach. A classical reference for the convergence is [14], see also <http://stanford.edu/~boyd/admm.html>.

Let us observe that in terms of the functions \tilde{f}, \tilde{g} , the algorithm computes, letting $\xi^k = Ax^k, \eta^k = By^k$:

$$\xi^{k+1} = \arg \min_{\xi} \tilde{f}(\xi) - \langle z^k, \xi \rangle + \frac{\gamma}{2} \|\zeta - \xi - \eta^k\|^2 = \text{prox}_{\tilde{f}/\gamma}(\zeta + \frac{1}{\gamma} z^k - \eta^k); \quad (73)$$

$$\eta^{k+1} = \arg \min_{\eta} \tilde{g}(\eta) - \langle z^k, \eta \rangle + \frac{\gamma}{2} \|\zeta - \xi^{k+1} - \eta\|^2 = \text{prox}_{\tilde{g}/\gamma}(\zeta + \frac{1}{\gamma} z^k - \xi^{k+1}). \quad (74)$$

Thanks to Moreau’s identity (29),

$$\text{prox}_{\gamma \tilde{f}^*}(z^k + \gamma(\zeta - \eta^k)) = z^k + \gamma(\zeta - \eta^k) - \gamma \xi^{k+1}, \quad (75)$$

$$\text{prox}_{\gamma\tilde{g}^*}(z^k + \gamma(\zeta - \xi^{k+1})) = z^k + \gamma(\zeta - \xi^{k+1}) - \gamma\eta^{k+1} = z^{k+1}. \quad (76)$$

Letting $\tilde{f}_\zeta^*(p) := \tilde{f}^*(p) - \langle \zeta, p \rangle = f^*(A^*p) - \langle \zeta, p \rangle$, the first line can also be rewritten

$$\gamma(\xi^{k+1} - \zeta) = z^k - \gamma\eta^k - \text{prox}_{\gamma\tilde{f}_\zeta^*}(z^k - \gamma\eta^k). \quad (77)$$

If we let $u^k = z^k - \gamma\eta^k$, $v^{k+1} = z^k + \gamma(\zeta - \xi^{k+1})$, we find that

$$\gamma\eta^{k+1} = z^k + \gamma(\zeta - \xi^{k+1}) - \text{prox}_{\gamma\tilde{g}^*}(z^k + \gamma(\zeta - \xi^{k+1})) = v^{k+1} - \text{prox}_{\gamma\tilde{g}^*}(v^{k+1}).$$

and

$$u^{k+1} = z^{k+1} - \gamma\eta^{k+1} = 2\text{prox}_{\gamma\tilde{g}^*}(v^{k+1}) - v^{k+1}.$$

On the other hand, (77) gives

$$\text{prox}_{\gamma\tilde{f}_\zeta^*}(u^k) + \gamma\eta^k = z^k + \gamma(\zeta - \xi^{k+1}) = v^{k+1}.$$

Hence the iteration reads

$$v^{k+1} = \text{prox}_{\gamma\tilde{f}_\zeta^*}(2\text{prox}_{\gamma\tilde{g}^*}(v^k) - v^k) + v^k - \text{prox}_{\gamma\tilde{g}^*}(v^k),$$

which is precisely a Douglas-Rachford iteration for the problem

$$0 \in \partial\tilde{g}^* + \partial\tilde{f}_\zeta^*$$

which is the equation for (71).

The theory seen up to now shows that $v^k \rightarrow v$ a fixed point of the iteration, which is such that $\text{prox}_{\gamma\tilde{g}^*}(v)$ is a solution of the dual problem. In practice, z^k will converge to a Lagrange Multiplier for (72), and x^k, y^k to a solution, as soon as there is enough coercivity (in particular, in finite dimension).

5.4 Other saddle-point algorithms: Primal-dual algorithm

We remark that thanks to (76) and (73), one has

$$\frac{z^k - z^{k-1}}{\gamma} = \zeta - \xi^k - \eta^k$$

hence

$$\xi^{k+1} = \text{prox}_{\tilde{f}/\gamma}(\xi^k + \frac{1}{\gamma}(2z^k - z^{k-1}))$$

while as before

$$z^{k+1} = \text{prox}_{\gamma\tilde{g}^*}(z^k - \gamma(\xi^{k+1} - \zeta)).$$

This is the form of a *primal-dual* algorithm (of ‘‘Arrow-Hurwicz’’ type) which aims at solving a fixed point problem of the form (letting $\tau = 1/\gamma$):

$$\xi + \tau\partial\tilde{f}(\xi) \ni \xi + \tau z, \quad z + \gamma\partial\tilde{g}^*(z) \ni z - \gamma(\xi - \zeta).$$

More generally, for a problem in the standard form

$$\min_x f(Kx) + g(x) = \min_x \sup_y \langle Kx, y \rangle + g(x) - f^*(y),$$

one can implement the Algorithm 4 described below.

Algorithm 4 PDHG

Input: initial pair of primal and dual points (x^0, y^0) , steps $\tau, \sigma > 0$.

for all $k \geq 0$ **do**

 find (x^{k+1}, y^{k+1}) by solving

$$x^{k+1} = \text{prox}_{\tau g}(x^k - \tau K^* y^k) \quad (78)$$

$$y^{k+1} = \text{prox}_{\sigma f^*}(y^k + \sigma K(2x^{k+1} - x^k)). \quad (79)$$

end for

Let us write now the iterates as follows:

$$\begin{cases} \frac{x^{k+1} - x^k}{\tau} + \partial g(x^{k+1}) \ni -K^* y^k = K^*(y^{k+1} - y^k) - K^* y^{k+1} \\ \frac{y^{k+1} - y^k}{\sigma} + \partial f^*(y^{k+1}) \ni K(x^{k+1} - x^k) + Kx^{k+1}, \end{cases}$$

that is

$$\begin{pmatrix} \frac{1}{\tau}I & -K^* \\ -K & \frac{1}{\sigma}I \end{pmatrix} \begin{pmatrix} x^{k+1} - x^k \\ y^{k+1} - y^k \end{pmatrix} + \begin{pmatrix} \partial g(x^{k+1}) \\ \partial f^*(y^{k+1}) \end{pmatrix} + \begin{pmatrix} 0 & K^* \\ -K & 0 \end{pmatrix} \begin{pmatrix} x^{k+1} \\ y^{k+1} \end{pmatrix} \ni 0. \quad (80)$$

We find that this algorithm is a proximal point algorithm for the variable $z^k = (x^k, y^k)^T$, the monotone operator which is the sum of the subgradient of the convex function $(x, y) \mapsto (g(x) + f^*(y))$ and the antisymmetric linear operator $\begin{pmatrix} 0 & K^* \\ -K & 0 \end{pmatrix}$, in the metric

$$M_{\tau, \sigma} := \begin{pmatrix} \frac{1}{\tau}I & -K^* \\ -K & \frac{1}{\sigma}I \end{pmatrix}$$

if this metric is positive definite. To see this we observe that if A is a monotone operator and M a symmetric positive definite operator, $M^{-1}A$ defines a monotone operator in the scalar product $\langle \cdot, \cdot \rangle_M = \langle M \cdot, \cdot \rangle$: if $p \in M^{-1}Ax$, $q \in M^{-1}Ay$,

$$\langle p - q, x - y \rangle_M = \langle M(p - q), x - y \rangle \geq 0$$

as $Mp \in Ax$, $Mq \in Ay$. Hence, in this metric, the resolvent J_A^M is given by $y = (I + M^{-1}A)^{-1}x$, which satisfies the equation $y + M^{-1}Ay \ni x$, that is, $M(y - x) + Ay \ni 0$.

When is the matrix $M_{\tau, \sigma}$ positive definite? We have

$$\left\langle M_{\tau, \sigma} \begin{pmatrix} \xi \\ \eta \end{pmatrix}, \begin{pmatrix} \xi \\ \eta \end{pmatrix} \right\rangle = \frac{1}{\tau} \|\xi\|^2 + \frac{1}{\sigma} \|\eta\|^2 - 2 \langle K\xi, \eta \rangle$$

which is positive if and only if for any $X, Y \geq 0$

$$\sup_{\|\xi\| \leq X, \|\eta\| \leq Y} 2 \langle K\xi, \eta \rangle = 2 \|K\| XY < \frac{X^2}{\tau} + \frac{Y^2}{\sigma}$$

if and only if

$$2 \|K\| < \min_{X \geq 0, Y \geq 0} \frac{X}{\tau Y} + \frac{Y}{\sigma X} = \frac{2}{\sqrt{\tau \sigma}}$$

if and only if

$$\tau \sigma \|K\|^2 < 1. \quad (81)$$

We deduce:

Theorem 5.14. *If (81) is satisfied, then $z^k = (x^k, y^k)^T$ defined by Algorithm 4 converges to a fixed point $(x, y)^T$ of the operator, that is, a solution of (33) (if one exists).*

5.4.1 Rate

To find a rate, we do as follows. Taking the scalar product of (80) with $z^{k+1} - z$ where z is an arbitrary point, we find

$$\begin{aligned} \langle z^{k+1} - z^k, z^{k+1} - z \rangle_{M_{\tau,\sigma}} + \left\langle \begin{pmatrix} 0 & K^* \\ -K & 0 \end{pmatrix} \begin{pmatrix} x^{k+1} \\ y^{k+1} \end{pmatrix}, \begin{pmatrix} x^{k+1} - x \\ y^{k+1} - y \end{pmatrix} \right\rangle \\ + g(x^{k+1}) + f^*(y^{k+1}) \leq g(x) + f^*(y) \end{aligned}$$

The scalar product is

$$-\langle K^* y^{k+1}, x \rangle + \langle K x^{k+1}, y \rangle$$

while

$$\langle z^{k+1} - z^k, z^{k+1} - z \rangle_{M_{\tau,\sigma}} = \frac{1}{2} \|z^{k+1} - z^k\|_{M_{\tau,\sigma}}^2 + \frac{1}{2} \|z^{k+1} - z\|_{M_{\tau,\sigma}}^2 - \frac{1}{2} \|z^k - z\|_{M_{\tau,\sigma}}^2.$$

Therefore, introducing the Lagrangian of (31), as

$$\mathcal{L}(x^{k+1}, y) - \mathcal{L}(x, y^{k+1}) = g(x^{k+1}) + \langle y, K x^{k+1} \rangle - f^*(y) - g(x) - \langle y^{k+1}, K x \rangle + f^*(y^{k+1})$$

we obtain for any $z = (x, y)^T$:

$$\mathcal{L}(x^{k+1}, y) - \mathcal{L}(x, y^{k+1}) + \frac{1}{2} \|z^{k+1} - z^k\|_{M_{\tau,\sigma}}^2 + \frac{1}{2} \|z^{k+1} - z\|_{M_{\tau,\sigma}}^2 \leq \frac{1}{2} \|z^k - z\|_{M_{\tau,\sigma}}^2.$$

Summing from $k = 0$ to $n - 1$ and using the convexity of $(\xi, \eta)^T \mapsto \mathcal{L}(\xi, y) - \mathcal{L}(x, \eta)$, we find if we let $Z^n = (X^n, Y^n)^T = (\sum_{k=1}^n z^n)/n$ that

$$\mathcal{L}(X^n, y) - \mathcal{L}(x, Y^n) \leq \frac{1}{2n} \|z^0 - z\|_{M_{\tau,\sigma}}^2. \quad (82)$$

This is a weak form of a rate (as it depends on (x, y)), and there is still some work to convert it into a true rate for the energy. The simplest case is when $\text{dom } f^*$, $\text{dom } g$ are bounded, then one can take the sup in x, y to find that

$$\mathcal{G}(X^n, Y^n) \leq \frac{C}{2n}$$

where $C = \sup\{\|z^0 - z\|_{M_{\tau,\sigma}}^2 : z = (x, y), x \in \text{dom } g, y \in \text{dom } f^*\}$.

5.4.2 Extensions

We present here an extension of Algorithm 4 due to Condat and in a generalized form to Vu (referred usually as Condat-Vu's primal-dual algorithm). A first observation (cf Vu, Bot) is that one can replace ∂g and ∂f^* with monotone operators, and get similar results.

A second observation, due to Condat, is that one can iterate the operator with an explicit step of a co-coercive operator. Typically, if h is a convex function with L_h -Lipschitz gradient, one can replace (80) with

$$\begin{pmatrix} \frac{1}{\tau} I & -K^* \\ -K & \frac{1}{\sigma} I \end{pmatrix} \begin{pmatrix} x^{k+1} - x^k \\ y^{k+1} - y^k \end{pmatrix} + \begin{pmatrix} \partial g(x^{k+1}) \\ \partial f^*(y^{k+1}) \end{pmatrix} + \begin{pmatrix} 0 & K^* \\ -K & 0 \end{pmatrix} \begin{pmatrix} x^{k+1} \\ y^{k+1} \end{pmatrix} \ni \begin{pmatrix} -\nabla h(x^k) \\ 0 \end{pmatrix}.$$

This iteration is of the form (38) and will converge if the operator

$$C = M_{\tau,\sigma}^{-1} \begin{pmatrix} \nabla h(x) \\ 0 \end{pmatrix}$$

is μ -co-coercive with $\mu > 1/2$, in the metric $M_{\tau,\sigma}$. That is, if for all z, z' :

$$\langle M_{\tau,\sigma}(z - z'), Cz - Cz' \rangle \geq \mu \|Cz - Cz'\|_{M_{\tau,\sigma}}^2.$$

Note that

$$\|Cz - Cz'\|_{M_{\tau,\sigma}}^2 = \left\langle M_{\tau,\sigma}^{-1} \begin{pmatrix} \nabla h(x) - \nabla h(x') \\ 0 \end{pmatrix}, \begin{pmatrix} \nabla h(x) - \nabla h(x') \\ 0 \end{pmatrix} \right\rangle$$

and that

$$M_{\tau,\sigma} \begin{pmatrix} \xi \\ \eta \end{pmatrix} = \begin{pmatrix} \nabla h(x) - \nabla h(x') \\ 0 \end{pmatrix} \Rightarrow \xi = (I - \sigma\tau K^*K)^{-1}(\tau(\nabla h(x) - \nabla h(x'))),$$

hence (using also the $1/L_h$ -co-coercivity of ∇h):

$$\begin{aligned} \|Cz - Cz'\|_{M_{\tau,\sigma}}^2 &= \langle \xi, \nabla h(x) - \nabla h(x') \rangle \leq \frac{\tau}{1 - \sigma\tau L^2} \|\nabla h(x) - \nabla h(x')\|^2 \\ &\leq \frac{\tau L_h}{1 - \sigma\tau L^2} \langle x - x', \nabla h(x) - \nabla h(x') \rangle \\ &= \frac{\tau L_h}{1 - \sigma\tau L^2} \langle M_{\tau,\sigma}(z - z'), Cz - Cz' \rangle. \end{aligned}$$

Here, $L = \|K\|$ (the operator norm). Hence C is μ -co-coercive for $\mu = (1 - \sigma\tau L^2)/(\tau L_h)$ and one deduces the algorithm converges provided

$$\frac{1}{\sigma} \left(\frac{1}{\tau} - \frac{L_h}{2} \right) > L^2.$$

In this case again we get the convergence of the Vu-Condat algorithm, which reads:

Algorithm 5 PDHG with explicit step

Input: initial pair of primal and dual points (x^0, y^0) , steps $\tau, \sigma > 0$.

for all $k \geq 0$ **do**

 find (x^{k+1}, y^{k+1}) by solving

$$x^{k+1} = \text{prox}_{\tau g}(x^k - \tau(K^*y^k + \nabla h(x^k))) \tag{83}$$

$$y^{k+1} = \text{prox}_{\sigma f^*}(y^k + \sigma K(2x^{k+1} - x^k)). \tag{84}$$

end for

Exercise: Show that a fixed point of these iterations solves

$$\min_x f(Kx) + g(x) + h(x) = \min_x \sup_y \langle y, Kx \rangle - f^*(y) + g(x) + h(x).$$

6 “Large scale” optimization

In this lecture, we only mention rapidly two techniques currently used to avoid computing full gradients. Such approaches are useful for solving very large dimensional problems.

6.1 Coordinate descent and stochastic coordinate descent

6.1.1 Does coordinate descent / alternating minimization work?

Assume one wants to solve

$$\min_{x_1, \dots, x_n} f(x_1, \dots, x_n)$$

and one knows how to solve, for any $i = 1, \dots, n$ and given $(x_j)_{j \neq i}$

$$\min_{\xi} f(x_1, \dots, x_{i-1}, \xi, x_{i+1}, \dots, x_n).$$

Then, it is natural to consider the following algorithm: (x^0) being given, one computes for $k \geq 0$, $i = 1, \dots, n$:

$$x_i^{k+1} \in \arg \min_{\xi} f(x_1^{k+1}, \dots, x_{i-1}^{k+1}, \xi, x_{i+1}^k, \dots, x_n^k). \quad (85)$$

Denoting $x = (x_1, \dots, x_n)$, does this converge? It depends. The following straightforward (classical) example shows that it is easily not the case.

Consider, for $x = (x_1, x_2) \in \mathbb{R}^2$, $f(x_1, x_2) = x_1^2/2 + |x_1 - x_2|$, which is minimal for $(x_1, x_2) = (0, 0)$. From (x_1^k, x_2^k) , the algorithm will first produce $x_1^{k+1} = \max\{-1, \min\{x_2^k, 1\}\}$ and then $x_2^{k+1} = x_1^{k+1}$. Hence, one has $x_1^k = x_2^k = x_2^1$ for any $k \geq 1$ and unless $x_2^0 = 0$, one never converges to the minimizer.

On the other hand, assume f is C^1 , bounded from below, coercive (infinite at infinity), so that the sequence (x^k) is bounded, and we are in finite dimension. A first remark is that by construction, $f(x^k)$ is decreasing, and converges to some value f^* . In addition, one has (with obvious notation, and without assuming particularly that the x_i are one-dimensional scalars):

$$\frac{\partial_i f(x_1^{k+1}, x_i^{k+1}, x_{i+1}^k, \dots, x_n^k)}{\partial x_i} = 0.$$

If $(x^k)_k$ converges, then one easily deduces that $\nabla f(x^k) = 0$, hence x^k is a critical point with value $f(x^k) = f^*$. But (x^k) could have subsequences converging to different limits.

In case f is convex, one can show that these limits are minimizers. Indeed, assume $\lim_i x_i^{k_i} = x_i$ and let also $x'_i = \lim_i x_i^{k_i+1}$ (possibly passing to another subsequence). Clearly, one has $f(x) = f(x') = f^* = f(x'_1, \dots, x'_{i-1}, x_i, \dots, x_n)$ for any i , and one also easily finds that x'_i is a minimizer of $f(x'_1, \dots, x'_{i-1}, \bullet, x_{i+1}, \dots, x_n)$, as well as x_i since f has the same value at all these points. In particular,

$$\frac{\partial f(x'_1, \dots, x'_i, x_{i+1}, \dots, x_n)}{\partial x_i} = \frac{\partial f(x'_1, \dots, x'_{i-1}, x_i, \dots, x_n)}{\partial x_i} = 0. \quad (86)$$

We show by induction that $\nabla f(x') = 0$. To start with, from (86) for $i = 1, 2$ we deduce that (x'_1, x_2) is a minimizer the convex function $f(\bullet, \bullet, x_3, \dots, x_n)$. But since $f(x'_1, x'_2, x_3, \dots, x_n) = f(x'_1, x_2, \dots, x_n) = f^*$, also (x'_1, x'_2) is a minimizer and in particular, the gradients of f with respect to x_1 and x_2 vanish at this point. By induction, if we assume that the gradient of f with respect to x_j , $j = 1, \dots, m$ vanishes in $(x'_1, \dots, x'_m, x_{m+1}, \dots, x_n)$, using (86) for $i = m$ and $i = m+1$ we find that (x'_m, x_{m+1}) is a minimizer of the convex function $f(x'_1, \dots, x'_{m-1}, \bullet, \bullet, x_{m+2}, \dots, x_n)$, and using the induction assumption, we have that $(x'_1, \dots, x'_m, x_{m+1})$ is a minimizer of the convex function $f(\bullet, \dots, \bullet, x_{m+2}, \dots, x_n)$. As the value is f^* , also (x'_1, \dots, x'_{m+1}) is a minimizer and the $(m+1)$ first gradients of f vanish at this point. By induction we deduce that $\nabla f(x') = 0$ and that x' (hence also x) is a minimizer. A similar proof in a more complex situation (with a convex, separable nonsmooth term) is found in [42] (Tseng).

6.1.2 Block coordinate descent

Instead of finding the minimizer of f with respect to one variable, one could perform a step of gradient descent. In particular, if f has Lipschitz gradients, one could rather replace (85) with

$$x_i^{k+1} = x_i^k - \tau_i \nabla_i f(x_1^{k+1}, \dots, x_{i-1}^{k+1}, x_i^k, \dots, x_n^k). \quad (87)$$

Here, $\nabla_i := \partial/\partial x_i$. Assume that $\partial_i f$ is L_i -Lipschitz (uniformly). Then one has as usual (see (1))

$$\begin{aligned} f(x_1^{k+1}, \dots, x_i^{k+1}, x_{i+1}^k, \dots, x_n^k) &\leq f(x_1^{k+1}, \dots, x_{i-1}^{k+1}, x_i^k, \dots, x_n^k) \\ &\quad - \tau_i \left(1 - \frac{L_i \tau_i}{2}\right) \|\nabla_i f(x_1^{k+1}, \dots, x_{i-1}^{k+1}, x_i^k, \dots, x_n^k)\|^2 \end{aligned}$$

Choosing for instance $\tau_i = \frac{1}{L_i}$, we deduce

$$\begin{aligned} f(x_1^{k+1}, \dots, x_i^{k+1}, x_{i+1}^k, \dots, x_n^k) &+ \frac{1}{2L_i} \|\nabla_i f(x_1^{k+1}, \dots, x_{i-1}^{k+1}, x_i^k, \dots, x_n^k)\|^2 \\ &\leq f(x_1^{k+1}, \dots, x_{i-1}^{k+1}, x_i^k, \dots, x_n^k) \end{aligned}$$

and in particular, one can reproduce the same proof as before in the convex case and show that any limit point is a minimizer. One interesting point here is that in general, the Lipschitz constant with respect to one variable is smaller than with respect to all the variables (think for instance to $(x_1, x_2) \mapsto (x_1 + x_2)^2$: its gradient is $\sqrt{2}$ -Lipschitz, while its partial gradients are 1-Lipschitz), so that the steps performed in the coordinate descent method are longer than for a gradient descent.

Up to now, we have considered alternating minimizations or block coordinate descent with a *cyclic* rule, where each coordinate is optimized in ascending order. This is a bit arbitrary. Let us now show that (on average) one can obtain good performances with a random update. This is our first example of a stochastic algorithm.

6.1.3 Random coordinate descent

We consider the following algorithm, with a notation slightly differing from the previous sections: we pick x^0 . At iteration $k \geq 0$, we consider x^k . We pick randomly a coordinate i with some probability p_i ($\sum_{i=1}^n p_i = 1$, $p_i > 0$), and let $i_k := i$. Then we let $x_j^{k+1} = x_j^k$ for $j \neq i_k$, and

$$x_{i_k}^{k+1} = x_{i_k}^k - \tau_{i_k} \nabla_{i_k} f(x^k). \quad (88)$$

As before, we have

$$f(x^{k+1}) \leq f(x^k) - \tau_{i_k} \left(1 - \frac{L_{i_k} \tau_{i_k}}{2}\right) \|\nabla_{i_k} f(x^k)\|^2 \quad (89)$$

As a consequence, knowing the point x^k , the expectation $\mathbb{E}(f(x^{k+1})|x^k)$ satisfies

$$\mathbb{E}(f(x^{k+1})|x^k) \leq f(x^k) - \sum_{i=1}^n p_i \tau_i \left(1 - \frac{L_i \tau_i}{2}\right) \|\nabla_i f(x^k)\|^2.$$

We can pick for instance $\tau_i = 1/L_i$ and $p_i = L_i/(\sum_j L_j)$, meaning that we pick more often the coordinates with larger Lipschitz constants. In this case, the previous estimate becomes

$$\mathbb{E}(f(x^{k+1})|x^k) \leq f(x^k) - \frac{1}{2 \sum_j L_j} \sum_{i=1}^n \|\nabla_i f(x^k)\|^2 = f(x^k) - \frac{1}{2 \sum_j L_j} \|\nabla f(x^k)\|^2. \quad (90)$$

Computing then the expectation with respect to x^k , we obtain

$$\mathbb{E}(f(x^{k+1})) \leq \mathbb{E}(f(x^k)) - \frac{1}{2 \sum_j L_j} \mathbb{E}(\|\nabla f(x^k)\|^2). \quad (91)$$

In particular, this is a decreasing sequence, and one has

$$\frac{1}{2 \sum_j L_j} \sum_{k=0}^{\infty} \mathbb{E}(\|\nabla f(x^k)\|^2) \leq f(x^0) < \infty$$

which shows that $\mathbb{E}(\|\nabla f(x^k)\|^2) \rightarrow 0$ ($\nabla f(x^k) \rightarrow 0$ almost surely, up to subsequences).

More generally, we pick $\tau_i = \theta/L_i$ for $\theta \in]0, 2[$ and introduce the norm $\|g\|_M^2 := \sum_{i=1}^n m_i |g_i|^2$, for $m_i := p_i/L_i$. Then the same computation as above yields

$$\mathbb{E}(f(x^{k+1})|x^k) \leq f(x^k) - \sum_{i=1}^n \frac{\theta(2-\theta)p_i}{L_i} \|\nabla_i f(x^k)\|^2 = f(x^k) - \frac{\theta(2-\theta)}{2} \|\nabla f(x^k)\|_M^2.$$

We assume that there exists a minimizer x^* and let $\Delta_k := f(x^k) - f(x^*)$, and show the following result:

Lemma 6.1. *Assume $\{f \leq f(x^0)\}$ is bounded. Then*

$$\mathbb{E}(\Delta_k) \leq \frac{2C^2}{\theta(2-\theta)} \frac{1}{k+1} \quad (92)$$

where $C \geq \sup_{f(x) \leq f(x^0)} \|x - x^*\|_{M^{-1}}$.

Proof. By convexity, we observe that

$$f(x) - f(x^*) \leq \langle \nabla f(x), x - x^* \rangle \leq \|\nabla f(x)\|_M \|x - x^*\|_{M^{-1}},$$

which is $\leq C \|\nabla f(x)\|$ if $f(x) \leq f(x^0)$ and C is as in the statement. Hence using (90), we find that

$$\mathbb{E}(f(x^{k+1}) - f(x^*)|x^k) \leq f(x^k) - f(x^*) - \frac{\theta(2-\theta)}{2} \frac{(f(x^k) - f(x^*))^2}{C^2}.$$

Now, by convexity (from Jensen's inequality), we know that $\mathbb{E}(\Delta_k)^2 \leq \mathbb{E}(\Delta_k^2)$ so that

$$\mathbb{E}(\Delta_{k+1}) \leq \mathbb{E}(\Delta_k) - \frac{\theta(2-\theta)}{2C^2} \mathbb{E}(\Delta_k^2) \leq \mathbb{E}(\Delta_k) - \frac{\theta(2-\theta)}{2C^2} \mathbb{E}(\Delta_k)^2.$$

Inequality (92) follows then from Lemma 2.6.

One sees here that it might be interesting to use non-uniform probabilities to improve the process, however it is not obvious how (one should minimize the ‘‘diameter’’ C , which is given by $C^2 = \sup_{f(x) \leq f(x^0)} \sum_i L_i |x_i - x_i^*|^2 / p_i$).

To compare with a standard gradient descent, one can use the choice already mentioned above, $\theta = 1$ and $p_i = L_i / \sum_j L_j$, for which $m_i = 1 / \sum_j L_j$. The rate becomes

$$\mathbb{E}(\Delta_{nk}) \leq \left(\frac{2}{n} \sum_{j=1}^n L_j \right) \frac{\sup_{f(x) \leq f(x^0)} \|x - x^*\|^2}{k + 1/n}$$

after k “epochs” (i.e., passes over all the data, at least on average: we consider that it requires n iterations to approximate one step of a full gradient descent). This is to be compared to the rate in Theorem 2.7:

$$\Delta_k \leq 2L \frac{\|x^0 - x^*\|^2}{k+1}$$

at the k th iteration of a gradient descent, where now L is the global Lipschitz constant of f .

So the relevant question here is: which is smallest of L and $\frac{1}{n} \sum_j L_j$? One always have

$$\max_j L_j \leq L \leq \sqrt{\sum_{j=1}^n L_j^2}, \quad (93)$$

hence

$$\frac{1}{n} \sum_j L_j \leq L,$$

whereas the upper bound in (93) satisfies

$$\frac{1}{n} \sum_j L_j \leq \frac{1}{\sqrt{n}} \sqrt{\sum_{j=1}^n L_j^2}.$$

Hence, in the worst case, the complexity of the random coordinate descent is similar to the gradient descent, while if L is closer to the upper bound in (93), the complexity is smaller by a factor $1/\sqrt{n}$, where n is the number of coordinates.

This approach has of course many extensions. It was first extended to the (separable) non-smooth case in [35]: it is shown that for an objective of the form $f(x) + \sum_{i=1}^n \psi_i(x_i)$ one can replace the k th iteration (88) with the proximal iteration

$$x_{i_k}^{k+1} = (I + \tau_{i_k} \partial \psi_i)^{-1}(x_{i_k} - \tau_{i_k} \nabla_{i_k} f(x^k))$$

with $\tau_{i_k} = 1/L_{i_k}$, and obtain essentially the same rate. Acceleration has been proposed shortly after, for a very complete variant (including non differentiable separable terms, parallel updates, and Nesterov-type acceleration...) see in particular [18].

6.2 Stochastic gradient descent

6.3 SGD for learning problems

We now consider a different problem, arising for instance in statistical learning, when one has to minimize (for large $n \geq 1$) a sum of convex functions of the form

$$\min_x \frac{1}{n} \sum_i f_i(x) + \psi(x) \quad (94)$$

Note that if ψ is strongly convex, one can derive a dual problem

$$\max_{y_1, \dots, y_n} -\frac{1}{n} \sum_i f_i^*(y_i) - \psi^*\left(-\frac{1}{n} \sum_i y_i\right)$$

where now ψ^* has Lipschitz gradient, and tackle the problem by a proximal variant of the (random) coordinate descent algorithm (such as in [42, 35, 18]), as mentioned in the last section. See also the variant termed “stochastic dual coordinate ascent” [40, 41].

We will focus on a direct gradient descent approach for the objective function $f(x) := (1/n) \sum_i f_i(x)$ (and hence the case $\psi = 0$, to simplify), considering however that if n is too large, it might not be a good idea to evaluate ∇f at each iteration. We assume here that each f_i is convex with L_i -Lipschitz gradient. We study the following “stochastic gradient” algorithm: starting from x^0 , for each $k \geq 1$, we

- pick $i_k = i \in \{1, \dots, n\}$ with probability $1/n$;
- let $x^{k+1} = x^k - \tau \nabla f_{i_k}(x^k)$, for some $\tau > 0$.

We observe immediately that $\mathbb{E}(x^{k+1}|x^k) = x^k - \tau \sum_i \frac{1}{n} \nabla f_i(x^k) = x^k - \tau \nabla f(x^k)$, so that this corresponds to a stochastic gradient descent where the gradient of f is replaced by a random variable with expectation $\nabla f(x)$ (and one could indeed consider this more general situation).

As usual, one can write that for $j = 1, \dots, n$, if $i_k = i$,

$$f_j(x^{k+1}) \leq f_j(x^k) - \tau \langle \nabla f_j(x^k), \nabla f_i(x^k) \rangle + \frac{L_j \tau^2}{2} \|\nabla f_i(x^k)\|^2$$

and summing, we find that

$$f(x^{k+1}) \leq f(x^k) - \tau \langle \nabla f(x^k), \nabla f_i(x^k) \rangle + \frac{\tau^2}{2} \left(\frac{1}{n} \sum_{j=1}^n L_j \right) \|\nabla f_i(x^k)\|^2.$$

We denote $\bar{L} := (\sum_j L_j)/n$ the average Lipschitz constant. Hence, knowing x^k , one has (using that each i appears with probability $1/n$)

$$\begin{aligned} \mathbb{E}(f(x^{k+1})|x^k) &\leq f(x^k) - \tau \|\nabla f(x^k)\|^2 + \frac{\tau^2}{2} \bar{L} \left(\frac{1}{n} \sum_{i=1}^n \|\nabla f_i(x^k)\|^2 \right) \\ &\leq f(x^k) - \tau \left(1 - \frac{\tau \bar{L}}{2}\right) \|\nabla f(x^k)\|^2 + \frac{\tau^2}{2} \bar{L} \left(\frac{1}{n} \sum_{i=1}^n \|\nabla f_i(x^k) - \nabla f(x^k)\|^2 \right) \end{aligned}$$

One sees that now, there is a problem: for $\tau < 2/\bar{L}$, one can expect that $\mathbb{E}(f(x^k))$ will decrease, until $\mathbb{E}(\|\nabla f(x^k)\|^2)$ (which is of the order of $\|x^k - x^{k+1}\|^2$) becomes comparable to $\mathbb{E}(\frac{1}{n} \sum_{i=1}^n \|\nabla f_i(x^k) - \nabla f(x^k)\|^2)$, which is the variance of the random gradient ∇f_i , averaged on the random point x^k .

Hence, with constant step size, one cannot expect this to converge. The only hope is that the “bad” variance term is of second order in τ . So that the standard solution is to replace τ in the iteration with a variable τ_k , with $\tau_k \rightarrow 0$. To simplify, we make also the assumption that the “variance” is globally bounded

$$\frac{1}{n} \sum_{i=1}^n \|\nabla f_i(x) - \nabla f(x)\|^2 \leq \sigma^2$$

for all x (or all x in some set, provided we can show that the iterates x^k will remain not too far from x^* : this is the case for instance if we assume that all the gradients $-\nabla f_i(x)$ point roughly towards x^* (or the origin) for large $|x|$, in the sense $\langle -\nabla f_i(x), x - x^* \rangle \geq$

$\theta|\nabla f_i(x)||x - x^*|$ for some $\theta \in (0, 1)$, for all i and for $|x|$ large enough). Assuming also $\tau_k \leq 1/\bar{L}$, one has then for $n \geq 1$

$$\left(\sum_{k=0}^{n-1} \tau_k \right) \min_{k=0, \dots, n-1} \mathbb{E}(\|\nabla f(x^k)\|^2) \leq f(x^0) + \frac{\bar{L}}{2} \sigma^2 \sum_{k=0}^{n-1} \tau_k^2$$

so that

$$\min_{k=0, \dots, n-1} \mathbb{E}(\|\nabla f(x^k)\|^2) \leq \frac{f(x^0) + \frac{\bar{L}}{2} \sigma^2 \sum_{k=0}^{n-1} \tau_k^2}{\sum_{k=0}^{n-1} \tau_k}.$$

One obtains a rate which is governed by the ratio

$$\frac{\sum_{k=0}^{n-1} \tau_k^2}{\sum_{k=0}^{n-1} \tau_k}.$$

for instance for $\tau_k \sim 1/k$, this is like $C/\log n$, while for $1/\sqrt{k}$, it is like $\log n/\sqrt{n}$.

The latter choice is nearly optimal, indeed, if one knows all the parameters of the problem and fixes the number of iterations in advance, one can use a fixed step τ : in this case, the best choice is to let $\bar{L}\sigma^2 n\tau^2/2 = f(x^0)$, yielding

$$\min_{k=0, \dots, n-1} \mathbb{E}(\|\nabla f(x^k)\|^2) \leq \frac{f(x^0) + \frac{\bar{L}}{2} \sigma^2 n\tau^2}{n\tau} = \frac{\sqrt{2\bar{L}f(x^0)}}{\sqrt{n}} \sigma$$

This approach is originally due Robbins and Monro [36].

6.3.1 Improvements of SGD

We only refer here to some recent improvements developed in the machine learning literature. The basic idea is to “reduce” along the iterations the variance of the gradient estimate, so that one does not have to send the step τ to zero to compensate. Starting from the early 2010’s, a few variants have been proposed, called for instance “SVRG” (stochastic variances-reduced gradient algorithm) [45], “SAG” (stochastic average gradient) [22], or “SAGA” [13].

For instance, the latter addresses problems of the form (94), where all f_i are supposed to have L -Lipschitz gradient, in the following way: assuming at iteration k one knows x^k and the values $\nabla f_i(y_i^k)$, $i = 1, \dots, n$, at $k + 1$ one does:

1. pick randomly an index $i \in \{1, \dots, n\}$ with uniform probability distribution;
2. let $y_i^{k+1} = x_i$, and for $j \neq i$, let $y_j^{k+1} = y_j^k$. Store $\nabla f_i(y_i^{k+1}) = \nabla f_i(x_i^k)$ in memory (points y need not be recorded, only the gradients are needed);
3. set

$$z^{k+1} = z^k - \tau \left(\nabla f_i(y_i^{k+1}) - \nabla f_i(y_i^k) + \frac{1}{n} \sum_{j=1}^n \nabla f_j(y_j^k) \right)$$

$$x^{k+1} = \text{prox}_{\tau\psi}(z^{k+1}).$$

Then, the results reported in [13] show that:

- If the f_i 's are μ -convex with L -Lipschitz gradient ($L > \mu > 0$), then if $\tau = 1/(2\mu n + L)$ one has

$$\mathbb{E}(\|x^k - x^*\|^2) \leq \left(1 - \frac{\mu}{2(\mu n + L)}\right)^k \left[\|x^* - x^0\|^2 + \frac{n}{\mu n + L} D_f(x^0, x^*)\right]$$

while for $\tau = 1/(3L)$ (not depending on μ), one has:

$$\mathbb{E}(\|x^k - x^*\|^2) \leq \left(1 - \min\left\{\frac{1}{4n}, \frac{\mu}{3L}\right\}\right)^k \left[\|x^* - x^0\|^2 + \frac{2n}{3L} D_f(x^0, x^*)\right]$$

- If the f_i 's have L -Lipschitz gradient, then again for $\tau = 1/(3L)$ one has, introducing the averages $\bar{x}^k := (1/k) \sum_{t=1}^k x^t$,

$$\mathbb{E}(F(\bar{x}^k) - F(x^*)) \leq \frac{4n}{k} \left[\frac{2L}{n} \|x^0 - x^*\|^2 + D_f(x^0, x^*)\right]$$

where F is the global objective in (94). Here, $D_f(x, y) := f(x) - f(y) - \langle \nabla f(y), x - y \rangle$ is the f -“Bregman distance” of x to y cf Remark 5.9 in Section 5.2.2.

References

- [1] Jean-Bernard Baillon and Georges Haddad. Quelques propriétés des opérateurs angle-bornés et n -cycliquement monotones. *Israel J. Math.*, 26(2):137–150, 1977.
- [2] H. H. Bauschke and P. L. Combettes. *Convex analysis and monotone operator theory in Hilbert spaces*. CMS Books in Mathematics/Ouvrages de Mathématiques de la SMC. Springer, New York, 2011. With a foreword by Hedy Attouch.
- [3] Amir Beck and Marc Teboulle. A fast iterative shrinkage-thresholding algorithm for linear inverse problems. *SIAM J. Imaging Sci.*, 2(1):183–202, 2009.
- [4] Dimitri P. Bertsekas. *Convex Optimization Algorithms*. Athena Scientific, 2015.
- [5] L. M. Brègman. A relaxation method of finding a common point of convex sets and its application to the solution of problems in convex programming. *Ž. Vyčisl. Mat. i Mat. Fiz.*, 7:620–631, 1967.
- [6] H. Brézis. *Opérateurs maximaux monotones et semi-groupes de contractions dans les espaces de Hilbert*. North-Holland Publishing Co., Amsterdam, 1973. North-Holland Mathematics Studies, No. 5. Notas de Matemática (50).
- [7] Haïm Brézis. *Analyse fonctionnelle*. Collection Mathématiques Appliquées pour la Maîtrise. [Collection of Applied Mathematics for the Master’s Degree]. Masson, Paris, 1983. Théorie et applications. [Theory and applications].
- [8] Richard H. Byrd, Peihuang Lu, Jorge Nocedal, and Ci You Zhu. A limited memory algorithm for bound constrained optimization. *SIAM J. Sci. Comput.*, 16(5):1190–1208, 1995.
- [9] Antonin Chambolle and Thomas Pock. An introduction to continuous optimization for imaging. *Acta Numer.*, 25:161–319, 2016.

- [10] P. L. Combettes and V. R. Wajs. Signal recovery by proximal forward-backward splitting. *Multiscale Model. Simul.*, 4(4):1168–1200, 2005.
- [11] R. Cominetti, J. A. Soto, and J. Vaisman. On the rate of convergence of Krasnosel’skiĭ-Mann iterations and their connection with sums of Bernoullis. *Israel J. Math.*, 199(2):757–772, 2014.
- [12] Damek Davis and Wotao Yin. A three-operator splitting scheme and its optimization applications. *Set-Valued Var. Anal.*, 25(4):829–858, 2017.
- [13] Aaron Defazio, Francis Bach, and Simon Lacoste-Julien. Saga: A fast incremental gradient method with support for non-strongly convex composite objectives. In *Proceedings of the 27th International Conference on Neural Information Processing Systems - Volume 1*, NIPS’14, pages 1646–1654, Cambridge, MA, USA, 2014. MIT Press.
- [14] Jonathan Eckstein and Dimitri P. Bertsekas. On the Douglas-Rachford splitting method and the proximal point algorithm for maximal monotone operators. *Math. Programming*, 55(3, Ser. A):293–318, 1992.
- [15] I. Ekeland and R. Témam. *Convex analysis and variational problems*, volume 28 of *Classics in Applied Mathematics*. Society for Industrial and Applied Mathematics (SIAM), Philadelphia, PA, english edition, 1999. Translated from the French.
- [16] L. C. Evans and R. F. Gariepy. *Measure theory and fine properties of functions*. CRC Press, Boca Raton, FL, 1992.
- [17] H. Federer. *Geometric measure theory*. Springer-Verlag New York Inc., New York, 1969.
- [18] Olivier Fercoq and Peter Richtárik. Optimization in high dimensions via accelerated, parallel, and proximal coordinate descent. *SIAM Rev.*, 58(4):739–771, 2016.
- [19] Daniel Gabay and Bertrand Mercier. A dual algorithm for the solution of nonlinear variational problems via finite element approximation. *Computers & Mathematics with Applications*, 2(1):17 – 40, 1976.
- [20] R. Glowinski and A. Marroco. Sur l’approximation, par éléments finis d’ordre un, et la résolution, par pénalisation-dualité, d’une classe de problèmes de Dirichlet non linéaires. *Rev. Française Automat. Informat. Recherche Opérationnelle Sér. Rouge Anal. Numér.*, 9(R-2):41–76, 1975.
- [21] Roger A. Horn and Charles R. Johnson. *Matrix analysis*. Cambridge University Press, Cambridge, second edition, 2013.
- [22] Nicolas Le Roux, Mark Schmidt, and Francis Bach. A stochastic gradient method with an exponential convergence rate for finite training sets. In *Proceedings of the 25th International Conference on Neural Information Processing Systems - Volume 2*, NIPS’12, pages 2663–2671, USA, 2012. Curran Associates Inc.
- [23] P. L. Lions and B. Mercier. Splitting algorithms for the sum of two nonlinear operators. *SIAM Journal on Numerical Analysis*, 16(6):964–979, 1979.
- [24] George J. Minty. Monotone (nonlinear) operators in Hilbert space. *Duke Math. J.*, 29:341–346, 1962.

- [25] José Luis Morales and Jorge Nocedal. Remark on “Algorithm 778: L-BFGS-B: Fortran subroutines for large-scale bound constrained optimization” [mr1671706]. *ACM Trans. Math. Software*, 38(1):Art. 7, 4, 2011.
- [26] Arkadi S. Nemirovski and D. Yudin. *Problem complexity and method efficiency in optimization*. John Wiley & Sons Inc., New York, 1983. Translated from the Russian and with a preface by E. R. Dawson, Wiley-Interscience Series in Discrete Mathematics.
- [27] Yu. Nesterov. Gradient methods for minimizing composite functions. *Math. Program.*, 140(1, Ser. B):125–161, 2013.
- [28] Yurii Nesterov. A method for solving the convex programming problem with convergence rate $O(1/k^2)$. *Dokl. Akad. Nauk SSSR*, 269(3):543–547, 1983.
- [29] Yurii Nesterov. *Introductory lectures on convex optimization*, volume 87 of *Applied Optimization*. Kluwer Academic Publishers, Boston, MA, 2004. A basic course.
- [30] Jorge Nocedal and Stephen J. Wright. *Numerical optimization*. Springer Series in Operations Research and Financial Engineering. Springer, New York, second edition, 2006.
- [31] Zdzisław Opial. Weak convergence of the sequence of successive approximations for nonexpansive mappings. *Bull. Amer. Math. Soc.*, 73:591–597, 1967.
- [32] B. T. Polyak. Some methods of speeding up the convergence of iteration methods. *U.S.S.R. Comput. Math. Math. Phys.*, 4(5):1–17, 1964.
- [33] Boris T. Polyak. *Introduction to optimization*. Translations Series in Mathematics and Engineering. Optimization Software, Inc., Publications Division, New York, 1987. Translated from the Russian, With a foreword by Dimitri P. Bertsekas.
- [34] Simeon Reich. Weak convergence theorems for nonexpansive mappings in Banach spaces. *J. Math. Anal. Appl.*, 67(2):274–276, 1979.
- [35] Peter Richtárik and Martin Takáč. Iteration complexity of randomized block-coordinate descent methods for minimizing a composite function. *Math. Program.*, 144(1-2, Ser. A):1–38, 2014.
- [36] Herbert Robbins and Sutton Monro. A stochastic approximation method. *Ann. Math. Statistics*, 22:400–407, 1951.
- [37] R. T. Rockafellar. Characterization of the subdifferentials of convex functions. *Pacific J. Math.*, 17:497–510, 1966.
- [38] R. T. Rockafellar. *Convex analysis*. Princeton Landmarks in Mathematics. Princeton University Press, Princeton, NJ, 1997. Reprint of the 1970 original, Princeton Paperbacks.
- [39] R. Tyrrell Rockafellar. Convex functions, monotone operators and variational inequalities. In *Theory and Applications of Monotone Operators (Proc. NATO Advanced Study Inst., Venice, 1968)*, pages 35–65. Edizioni “Oderisi”, Gubbio, 1969.
- [40] Shai Shalev-Shwartz and Tong Zhang. Stochastic dual coordinate ascent methods for regularized loss minimization. *J. Mach. Learn. Res.*, 14:567–599, 2013.

- [41] Shai Shalev-Shwartz and Tong Zhang. Accelerated proximal stochastic dual coordinate ascent for regularized loss minimization. *Math. Program.*, 155(1-2, Ser. A):105–145, 2016.
- [42] P. Tseng. Convergence of a block coordinate descent method for nondifferentiable minimization. *J. Optim. Theory Appl.*, 109(3):475–494, 2001.
- [43] Paul Tseng. On accelerated proximal gradient methods for convex-concave optimization, 2008. *Submitted to SIAM J. Optim.* / available at <http://www.csie.ntu.edu.tw/~b97058/tseng/papers/apgm.pdf>.
- [44] D. van Dulst. Equivalent norms and the fixed point property for nonexpansive mappings. *J. London Math. Soc. (2)*, 25(1):139–144, 1982.
- [45] Lin Xiao and Tong Zhang. A proximal stochastic gradient method with progressive variance reduction. *SIAM J. Optim.*, 24(4):2057–2075, 2014.