

Optimisation convexe et méthodes de descente

Introduction	11
1.1 Convexité	12
1.1.1 Définitions	12
1.1.2 Existence d'un minimiseur	13
1.1.3 Conditions d'optimalité	14
1.1.4 Forte convexité	16
1.2 Dualité	18
1.2.1 Conjuguée convexe ou conjuguée de Legendre-Fenchel	18
1.2.2 Point-selle	19
1.3 Opérateur proximal	20
1.3.1 Définition et caractérisation	20
1.3.2 Identité de Moreau	22
1.4 Optimisation convexe et méthodes proximales	23
1.4.1 Méthodes de gradient	23
1.4.2 Méthodes d'éclatement	27
1.4.3 Itérations de Bregman	29

Introduction

De nombreux problèmes rencontrés en traitement d'images peuvent être abordés en introduisant une fonctionnelle d'énergie qui traduit le modèle considéré, en mesurant l'écart de toute fonction donnée à ce modèle. Si le modèle est suffisamment réaliste, la solution recherchée est logiquement celle qui en est le plus proche, c'est-à-dire celle qui minimise la fonctionnelle associée. Les fonctionnelles sont généralement composées de plusieurs termes séparés, chacun correspondant à une composante du modèle considéré. Ces termes peuvent être différentiables ou non, mais sont généralement convexes¹. On se concentre dans ce chapitre sur le cas des fonctionnelles convexes, car il offre un cadre de travail naturel pour la minimisation. La convexité assure en effet l'existence d'une solution (donnée par le minimum global), mais surtout la non-existence de minima locaux. Elle permet donc d'envisager des stratégies basiques de *descente* pour rechercher

1. Ce n'est pas toujours le cas, cf. chapitre 2

un minimum : elles consistent à trouver la direction de descente de l'énergie, qui est donnée par le gradient lorsque la fonctionnelle est différentiable. On verra qu'il est possible d'étendre ce genre de méthodes dans le cas non différentiable.

La régularité (au sens large) des fonctionnelles joue un rôle prépondérant dans le choix et la conception des algorithmes de minimisation. La différentiabilité permet par exemple des calculs explicites dans les schémas de descente de gradient. Néanmoins, la différentiabilité n'est pas toujours acquise. Si la fonctionnelle n'est que partiellement différentiable, on montre qu'il est avantageux d'exploiter la régularité partielle de la fonctionnelle, ce qui conduit à des méthodes dites par *éclatement*. Dans le cas plus général, on introduit un opérateur dit *proximal*, qui est à l'origine d'une classe d'algorithmes dit *proximaux*. Ces derniers comprennent en particulier les méthodes classiques de gradient implicite, gradient explicite ou encore gradient projeté. Une autre sorte de régularité fournit également des résultats intéressants : il s'agit de la *forte convexité*. Cette propriété permet en outre de proposer des algorithmes accélérés, comme on le verra dans le chapitre 5. Enfin, un dernier aspect important en optimisation convexe reste la complexité des calculs. Pour qu'un algorithme soit utilisable sur des données réelles, il est nécessaire d'en assurer la convergence dans un temps raisonnable. On verra dans le chapitre 6 qu'une stratégie d'éclatement peut en réduire la complexité, en permettant notamment les calculs parallèles. Néanmoins, ce genre d'approches implique des résolutions approchées.

L'objectif de ce chapitre est de donc de rappeler et d'établir certains résultats classiques en optimisation convexe continue, en vue de les appliquer dans les deux prochains chapitres. Nous commencerons par des rappels sur le cadre de la convexité (section 1.1). On verra ensuite le cas plus connu des fonctions différentiables, puis on introduira la notion de sous-différentiabilité. Cette notion est au coeur de la théorie des opérateurs proximaux (section 1.3) qui offre une classe très large d'algorithmes, appelés algorithmes proximaux, qui généralisent les méthodes de gradient (section 1.4). Enfin, on présentera des stratégies classiques qui permettent d'exploiter les propriétés de régularité d'une seule partie de la fonctionnelle (méthodes d'éclatement).

1.1 Convexité

Dans cette section, on fait quelques rappels sur les fonctions convexes, en considérant le cas général des fonctions à valeurs dans $\mathbb{R} \cup \{\pm\infty\}$. Ce choix nous permettra en particulier de considérer des problèmes de minimisation sous contraintes, mais sous une forme non contrainte. On rappellera ensuite les résultats d'existence de minimum, puis les conditions d'optimalité du premier ordre, d'abord dans le cas familier des fonctions différentiables, puis dans le cas plus général des fonctions sous-différentiables. Enfin, on introduira les fonctions fortement convexes, pour lesquelles on verra plus tard qu'il est possible de proposer des algorithmes accélérés.

1.1.1 Définitions

Commençons par quelques définitions classiques.

Domaine, fonctions propres Soit X un espace hilbertien, de dual noté X^* . On va considérer dans ce chapitre des fonctions à valeurs dans la ligne réelle étendue $X \rightarrow \mathbb{R} \cup \{\pm\infty\}$. On appellera alors *domaine* de la fonction F , noté $\text{dom}(F)$ l'ensemble des points $x \in X$ tels que $f(x) < +\infty$.

Une fonction $F : E \rightarrow \mathbb{R} \cup \{\pm\infty\}$ est dite *propre* si elle ne prend pas la valeur $-\infty$ et si elle n'est pas identiquement égale à $+\infty$. Autrement dit, le domaine d'une fonction propre n'est pas vide.

Inégalité de Jensen Un ensemble $C \subset X$ est dit *convexe* si, pour tous éléments x_1 et x_2 de C , le segment $[x_1; x_2]$ défini par $\{\lambda x_1 + (1 - \lambda)x_2 \mid \lambda \in [0; 1]\}$ est contenu dans C . Une fonction $F : E \rightarrow \mathbb{R} \cup \{+\infty\}$ est dite *convexe* si son épigraphe (c'est-à-dire l'ensemble des points situés au-dessus de son graphe) est un ensemble convexe. Elle est dite *concave* si $-F$ est convexe.

On montre que F est convexe si et seulement si elle vérifie l'inégalité de JENSEN, qui pour tout couple $(x_1, x_2) \in (\text{dom}(F))^2$ s'écrit

$$\forall \lambda \in]0; 1[, \quad F(\lambda x_1 + (1 - \lambda)x_2) \leq \lambda F(x_1) + (1 - \lambda)F(x_2).$$

Si cette inégalité est stricte pour tout $x_1 \neq x_2$, alors F est dite *strictement convexe*.

1.1.2 Existence d'un minimiseur

Un des intérêts de l'optimisation convexe repose sur le fait qu'il n'existe pas de minimum local dans lequel les algorithmes de recherche de minimum par descente pourraient être piégés, comme l'atteste le résultat suivant :

Proposition 1 *Soit F une fonction convexe sur X . Si F admet un minimum local en x^* , alors il admet un minimum global en x^* .*

Un autre intérêt des fonctions convexes est l'existence de minimiseurs dans des cas simples à caractériser. Voici dans ce qui suit quelques résultats connus d'existence (et/ou d'unicité) de minimiseurs. Pour plus de détails, le lecteur pourra se reporter par exemple à [10].

Cas général On s'intéresse tout d'abord aux fonctions convexes (non strictement convexes). Commençons par le cas des fonctions continues sur un compact, pour lesquelles l'existence d'un minimiseur est assuré :

Proposition 2 *Soit F une fonction convexe sur un compact $K \subset X$ non vide. On suppose que F est continue sur K . Alors F admet au moins un minimum dans K .*

On enchaîne ensuite avec le cas non borné ; ce cas nécessite de considérer des fonctions dites *coercives*, c'est-à-dire telles que

$$F(x) \rightarrow +\infty \quad \text{si} \quad \|x\| \rightarrow +\infty.$$

On peut alors montrer que, dans ce cas, l'existence d'un minimiseur est également assurée :

Proposition 3 *Soit F une fonction convexe sur X . On suppose que F est continue et coercive sur X . Alors F admet au moins un minimum dans X .*

Ce résultat se généralise à la minimisation sur un fermé quelconque non vide.

Cas de la stricte convexité Lorsqu'on ajoute une hypothèse de stricte convexité, les résultats qui précèdent incluent un résultat d'unicité. En effet, on peut montrer que

Proposition 4 *Soit F une fonction strictement convexe sur X . Alors F admet au plus un minimum sur X .*

Cette proposition nous permet donc d'énoncer un résultat important sur la minimisation d'une fonction strictement convexe, lorsque celle-ci est coercive :

Théorème 1 *Soit F une fonction strictement convexe et coercive sur X . Soit A un fermé non vide. On suppose que F est continue sur A . Alors F admet exactement un minimum dans A .*

Dans tout ce qui suit, on supposera toujours (sans le préciser) l'existence d'au moins un minimiseur.

1.1.3 Conditions d'optimalité

On présente dans ce paragraphe des résultats utiles permettant de caractériser, lorsqu'ils existent, les minima des fonctions considérées. On se focalise plus particulièrement sur les conditions dites *du premier ordre*, qui concernent les fonctions différentiables ou sous-différentiables, et qui donnent un critère sur le gradient ou le sous-gradient.

Cas différentiable Commençons par traiter le cas plus familier des fonctions différentiables. Les conditions nécessaires d'optimalité sont connus sous le nom d'équation ou d'inégalité d'EULER.

Théorème 2 (Équation d'Euler) *Soit F une fonction convexe sur X . On suppose que F est différentiable sur X . Alors F admet un minimum en x^* si et seulement si x^* vérifie l'équation d'EULER*

$$\nabla F(x^*) = 0.$$

Ce résultat est également valable sur tout ouvert convexe $\Omega \subset X$. Il **n'est pas valable** sur des ensembles fermés (où le minimum, s'il existe, peut être atteint sur le bord). C'est l'objet du résultat suivant, généralisable à tout convexe Ω :

Proposition 5 (Inégalité d'Euler) *Soit F une fonction convexe sur X . On suppose que F est différentiable sur X . Alors F admet un minimum en x^* si et seulement si x^* vérifie l'inégalité d'EULER*

$$\forall x \in X, \quad \langle x - x^*, \nabla F(x^*) \rangle \geq 0.$$

Sous-différentiabilité On quitte maintenant le cadre des fonctions différentiables. Commençons par introduire la notion de sous-différentiabilité, qui généralise celle de la différentiabilité dans le cas des fonctions convexes. Soit $x_0 \in X$ tel que $x_0 \in \text{dom}(F)$ avec F une fonction convexe. On définit le *sous-différentiel* de F en x_0 , noté $\partial F(x_0)$, comme étant l'ensemble des points $p \in X^*$ vérifiant

$$\forall x \in X, \quad \langle x - x_0, p \rangle + F(x_0) \leq F(x)$$

appelés, quand ils existent, *sous-gradients de F en x_0* . On dit alors que F est *sous-différentiable en x_0* si son sous-différentiel en x_0 est non vide. Par convention, on définit le sous-différentiel de F en x_0 comme étant l'ensemble vide si $F(x_0) = +\infty$. On peut montrer par ailleurs que le sous-différentiel en x_0 d'une fonction convexe F différentiable en x_0 est donné par le singleton $\{\nabla F(x_0)\}$. Réciproquement, on établit que, si le sous-différentiel de F en x_0 est réduit à un vecteur p , alors F est différentiable en x_0 , de gradient p .

Calcul de sous-différentiel Établissons ici quelques règles de calcul de sous-différentiel qui nous seront utiles par la suite. Commençons par remarquer que, pour tout $\alpha > 0$, on a

$$\forall x_0 \in \text{dom}(F), \quad \partial(\alpha F)(x_0) = \alpha \partial F(x_0).$$

En effet, on a par définition du sous-différentiel

$$\begin{aligned} x \in \partial(\alpha F)(x_0) &\iff \forall x \in X, \quad \langle x - x_0, p \rangle + \alpha F(x_0) \leq \alpha F(x) \\ &\iff \forall p \in X, \quad \left\langle x - x_0, \frac{x}{\alpha} \right\rangle + F(x_0) \leq F(x) \\ x \in \partial(\alpha F)(x_0) &\iff \frac{x}{\alpha} \in \partial F(x_0). \end{aligned}$$

Supposons à présent que f est une fonction convexe différentiable et F convexe. Posons $G = F + f$. Calculons $\partial G(x_0)$ pour tout $x_0 \in \text{dom}(F) \cap \text{dom}(f)$. Si $p \in \partial F(x_0)$, alors

$$\forall x \in X, \quad \langle x - x_0, p \rangle + F(x_0) \leq F(x).$$

On a par ailleurs, puisque $\partial f(x_0) = \{\nabla f(x_0)\}$,

$$\forall x \in X, \quad \langle x - x_0, \nabla f(x_0) \rangle + f(x_0) \leq f(x).$$

En additionnant les deux, on montre que $p + \nabla f(x_0) \in \partial G(x_0)$. Ainsi, on prouve que $\partial F(x_0) + \nabla f(x_0) \subset \partial G(x_0)$ ². Supposons maintenant que $x \in \partial G(x_0)$ et démontrons l'inclusion inverse. On a par définition

$$\forall x \in X, \forall \lambda \in]0; 1[, \quad \langle [\lambda x + (1 - \lambda)x_0] - x_0, p \rangle + G(x_0) \leq G([\lambda x + (1 - \lambda)x_0])$$

soit $\lambda \langle x - x_0, p \rangle + F(x_0) + f(x_0) \leq \lambda F(x) + (1 - \lambda) F(x_0) + f(x_0 + \lambda(x - x_0))$.

En simplifiant et en utilisant la formule de TAYLOR-YOUNG au premier ordre pour f , on obtient que

$$\begin{aligned} \lambda \langle x - x_0, p \rangle + f(x_0) &\leq \lambda F(x) - \lambda F(x_0) + f(x_0) + \lambda \langle x - x_0, \nabla f(x_0) \rangle \\ &\quad + \lambda \|x - x_0\| \varepsilon(\lambda(x - x_0)). \end{aligned}$$

2. En réalité, on démontre aussi que, de manière générale, $\partial F_1(x_0) + \partial F_2(x_0) \subset \partial(F_1 + F_2)(x_0)$ pour tout $x_0 \in X$.

En divisant à nouveau par λ puis en le faisant tendre vers 0, il s'ensuit que

$$\forall x \in X, \quad \langle x - x_0, p - \nabla f(x_0) \rangle + F(x_0) \leq F(x).$$

On en déduit que $p - \nabla f(x_0) \in \partial F(x_0)$, donc $\partial G(x_0) \subset \partial F(x_0) + \nabla f(x_0)$. Finalement, si f est différentiable, alors

$$\forall x_0 \in X, \quad \partial G(x_0) = \partial F(x_0) + \nabla f(x_0).$$

Cas sous-différentiable Généralisons à présent les conditions d'optimalité de premier ordre au cas des fonctions sous-différentiables :

Théorème 3 (Équation d'Euler (2)) Soit F une fonction convexe sur X . On suppose que F est sous-différentiable sur X . Alors F admet un minimum en x^* si et seulement si x^* vérifie l'équation d'EULER

$$0 \in \partial F(x^*).$$

DÉMONSTRATION : Puisque $\langle y - x_0, 0 \rangle$ est nul pour tout $y \in X$, on en déduit que x_0 minimise F si et seulement si $F(x_0) < +\infty$ et que

$$\forall y \in X, \quad \langle y - x_0, 0 \rangle + F(x_0) \leq F(y)$$

c'est-à-dire si et seulement si $0 \in \partial F(x_0)$. ■

1.1.4 Forte convexité

Introduisons enfin la notion de forte convexité, qui apparaîtra en particulier dans le chapitre suivant.

Définition Une fonction $F : E \rightarrow \mathbb{R} \cup \{+\infty\}$ est dite *fortement convexe*, de module $\alpha > 0$, si pour tout couple $(x_1, x_2) \in (\text{dom}(F))^2$

$$\forall \lambda \in]0; 1[, \quad F(\lambda x_1 + (1 - \lambda)x_2) \leq \lambda F(x_1) + (1 - \lambda)F(x_2) - \frac{\alpha}{2} \lambda(1 - \lambda) \|x_1 - x_2\|^2.$$

Une fonction fortement convexe est en particulier strictement convexe.

Exemples La fonction $x \mapsto \alpha \|x + a\|^2/2$ est fortement convexe de module α pour tout $a \in X$ et tout $\alpha > 0$. On peut également montrer que la somme d'une fonction convexe et d'une fonction fortement convexe est fortement convexe.

On peut enfin établir que, si la fonction F est fortement convexe de module α , alors, pour tout $x \in X$, la fonction

$$y \mapsto F(y) - \frac{\alpha}{2} \|x - y\|^2$$

est convexe.

Caractérisation Établissons à présent une caractérisation bien utile de la forte convexité :

Proposition 6 Soit F une fonction définie sur l'espace euclidien X . Supposons que F est sous-différentiable. Alors F est fortement convexe de module α si et seulement si

$$\forall (x_1, x_2) \in (\text{dom}(F))^2, \quad F(x_2) \geq F(x_1) + \langle p, x_2 - x_1 \rangle + \frac{\alpha}{2} \|x_2 - x_1\|^2$$

avec $p \in \partial F(x_2)$.

DÉMONSTRATION : • Commençons par prouver le sens direct. Supposons que F est fortement convexe de module α . Par définition, on a pour tout $(x_1, x_2) \in (\text{dom}(F))^2$ et $\lambda \in]0; 1[$

$$F(\lambda x_2 + (1 - \lambda)x_1) \leq \lambda F(x_2) + (1 - \lambda) F(x_1) - \frac{\alpha}{2} \lambda(1 - \lambda) \|x_2 - x_1\|^2.$$

que l'on peut réécrire

$$F(x_1 + \lambda(x_2 - x_1)) - F(x_1) \leq \lambda [F(x_2) - F(x_1)] - \frac{\alpha}{2} \lambda(1 - \lambda) \|x_2 - x_1\|^2.$$

Soit maintenant $p \in \nabla F(x_1)$. Par définition, p vérifie

$$\forall z \in X, \quad \langle z - x_1, p \rangle + F(x_1) \leq F(z).$$

En particulier, pour $z = x_1 + \lambda(x_2 - x_1)$,

$$\langle x_1 + \lambda(x_2 - x_1) - x_1, p \rangle + F(x_1) \leq F(x_1 + \lambda(x_2 - x_1))$$

soit

$$\lambda \langle x_2 - x_1, p \rangle \leq F(x_1 + \lambda(x_2 - x_1)) - F(x_1).$$

Ainsi, on obtient que

$$\lambda \langle x_2 - x_1, p \rangle \leq \lambda [F(x_2) - F(x_1)] - \frac{\alpha}{2} \lambda(1 - \lambda) \|x_2 - x_1\|^2$$

et puisque λ est strictement positif, on peut simplifier par λ , puis faire tendre λ vers 0 et obtenir l'inégalité recherchée.

• Montrons maintenant l'autre sens. Supposons que F vérifie pour tout $(x, x') \in (\text{dom}(F))^2$

$$F(x') \geq F(x) + \langle p, x' - x \rangle + \frac{\alpha}{2} \|x' - x\|^2$$

avec $p \in \nabla F(x)$. Soit $(x_1, x_2) \in (\text{dom}(F))^2$ et $\lambda \in]0; 1[$. Appliquons cette relation à $x = \lambda x_1 + (1 - \lambda)x_2$ et $x' = x_1$:

$$F(\lambda x_1 + (1 - \lambda)x_2) \geq F(x_1) + \langle p, \lambda x_1 + (1 - \lambda)x_2 - x_1 \rangle + \frac{\alpha}{2} \|\lambda x_1 + (1 - \lambda)x_2 - x_1\|^2$$

puis à $x' = x_2$:

$$F(\lambda x_1 + (1 - \lambda)x_2) \geq F(x_2) + \langle p, \lambda x_1 + (1 - \lambda)x_2 - x_2 \rangle + \frac{\alpha}{2} \|\lambda x_1 + (1 - \lambda)x_2 - x_2\|^2$$

ce qui donne respectivement, en simplifiant,

$$F(\lambda x_1 + (1 - \lambda)x_2) \geq F(x_1) + (1 - \lambda) \langle p, x_2 - x_1 \rangle + \frac{\alpha}{2} (1 - \lambda)^2 \|x_1 - x_2\|^2$$

et

$$F(\lambda x_1 + (1 - \lambda)x_2) \geq F(x_2) + \lambda \langle p, x_1 - x_2 \rangle + \frac{\alpha}{2} \lambda^2 \|x_1 - x_2\|^2.$$

En multipliant la première inégalité par λ , la seconde par $(1 - \lambda)$, puis en les ajoutant, on obtient finalement la relation de forte convexité souhaitée, ce qui achève la preuve. ■

Minimisation Commençons par remarquer que la caractérisation de la forte convexité assure que, si x^* est le minimiseur (unique) de F , alors on a

$$\forall x \in \text{dom}(F), \quad F(x) \geq F(x^*) + \frac{\alpha}{2} \|x - x^*\|^2$$

puisque x^* vérifie l'équation d'EULER. On peut par ailleurs signaler le résultat suivant [1] :

Proposition 7 Soit F une fonction fortement convexe de module α définie sur l'espace euclidien X . Alors F admet un unique minimiseur x^* , et toute suite minimisante (c'est-à-dire toute suite $(x_n)_n$ telle que $(f(x_n))_n$ converge vers $f(x^*)$) converge vers x^* . Par ailleurs, on a pour tout $x \in X$

$$\|x - x^*\|^2 \leq \frac{4}{\alpha} (F(x) - F(x^*)).$$

1.2 Dualité

1.2.1 Conjuguée convexe ou conjuguée de Legendre-Fenchel

On va introduire dans cette section deux notions importantes, qui sont la semi-continuité inférieure et la conjuguée convexe (ou conjuguée de LEGENDRE-FENCHEL, déjà rencontrée dans le chapitre 3).

Semi-continuité inférieure Une fonction $F : X \rightarrow \mathbb{R} \cup \{+\infty\}$ est dite *semi-continue inférieurement* (abrégé en s.c.i.³) si

$$\forall x^* \in X, \quad \liminf_{x \rightarrow x^*} F(x) \geq F(x^*).$$

On peut montrer [10] que F est s.c.i. si et seulement si son épigraphe est fermé.

Conjuguée convexe Soit F une fonction convexe définie sur X . On définit sa conjuguée convexe, ou encore conjuguée de LEGENDRE-FENCHEL, en posant

$$\forall y \in X^*, \quad F^*(y) = \sup_{x \in X} \{ \langle x, y \rangle - F(x) \}.$$

Commençons par remarquer que F^* est à valeurs dans $\mathbb{R} \cup \{+\infty\}$ car F est propre. On peut par ailleurs montrer que F^* est convexe et s.c.i. Signalons également le résultat suivant :

Théorème 4 Si F est s.c.i., alors sa biconjuguée $F^{**} = (F^*)^*$, qui est la conjuguée convexe de sa conjuguée convexe, est F elle-même. Autrement dit, on a la relation suivante

$$\forall x \in X, \quad F(x) = \sup_{y \in X^*} \{ \langle x, y \rangle - F^*(y) \}.$$

Une manière d'interpréter ce résultat est de montrer que la biconjuguée F^{**} de F est en réalité la plus grande fonction convexe, propre et s.c.i. située en-dessous de F . On a donc l'égalité lorsque F possède déjà ces propriétés.

3. En anglais, l.s.c. pour *lower semicontinuous*.

Cas de la forte convexité Lorsque la fonction F est fortement convexe, cela induit une régularité remarquable sur sa conjuguée convexe, comme en témoigne le résultat suivant [2, 8] :

Théorème 5 Soit F une fonction convexe.

- Si F est différentiable et ∇F est lipschitzienne, de constante L , alors F^* est fortement convexe de module $(1/L)$.
- Si F est fortement convexe, de module α , alors F^* est différentiable, de gradient lipschitzien, de constante de LIPSCHITZ $1/\alpha$.

En particulier, si F est convexe, propre et s.c.i., alors $F^{**} = F$, ce qui implique que F est fortement convexe de module α si et seulement si F^* est différentiable, de gradient $(1/\alpha)$ -lipschitzien. C'est pourquoi la forte convexité peut être interprétée comme une forme de régularité.

1.2.2 Point-selle

Enfin, rappelons quelques notions utiles à propos des points-selles.

Définition et propriétés Un couple $(\bar{x}, \bar{y}) \in X \times Y$ est un point-selle de la fonction \mathcal{L} sur $X \times Y$ si

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

On peut montrer qu'une fonction \mathcal{L} à valeurs réelles définie sur $X \times Y$ possède un point-selle (\bar{x}, \bar{y}) sur $X \times Y$ si et seulement si

$$\max_{y \in Y} \inf_{x \in X} \mathcal{L}(x, y) = \min_{x \in X} \sup_{y \in Y} \mathcal{L}(x, y)$$

Ce nombre est alors égal à $\mathcal{L}(\bar{x}, \bar{y})$ (appelée *valeur-selle*).

Cas convexe-concave On suppose que pour tout $y \in Y$, la fonction $f_y : x \mapsto \mathcal{L}(x, y)$ est convexe, et que, pour tout $x \in X$, la fonction $g_x : y \mapsto \mathcal{L}(x, y)$ est concave. On dit alors que la fonction \mathcal{L} est *convexe-concave*. Dans ce cas, l'existence d'un point-selle est assurée par le théorème suivant :

Théorème 6 Soit $\mathcal{L} : X \times Y \rightarrow \mathbb{R}$ une fonction convexe-concave. Posons pour tout $y \in Y$ et pour tout $x \in X$:

$$f_y : x \mapsto \mathcal{L}(x, y) \quad \text{et} \quad g_x : y \mapsto \mathcal{L}(x, y).$$

On suppose que tout $y \in Y$, la fonction g_y est s.c.i. et que pour tout $x \in X$, la fonction f_x est s.c.i. Alors \mathcal{L} possède un point-selle sur $X \times Y$.

Conditions d'optimalité Voyons comment on peut caractériser les points selle d'une fonction \mathcal{L} dans le cas convexe-concave :

Théorème 7 (Équation d'Euler) Soit $\mathcal{L} : X \times Y \rightarrow \mathbb{R}$ une fonction convexe-concave et sous-différentiable. Posons pour tout $y \in Y$ et pour tout $x \in X$:

$$f_y : x \mapsto \mathcal{L}(x,y) \quad \text{et} \quad g_x : y \mapsto \mathcal{L}(x,y).$$

Alors (\bar{x}, \bar{y}) est un point-selle de \mathcal{L} si et seulement si

$$0 \in \partial f_{\bar{y}}(\bar{x}) \quad \text{et} \quad 0 \in \partial g_{\bar{x}}(\bar{y}).$$

1.3 Opérateur proximal

On introduit dans cette section un opérateur, appelé *opérateur proximal*, introduit par Jean-Jacques MOREAU [12]. Il permet de concevoir une classe de méthodes d'optimisation convexe applicables à des fonctions non différentiables (mais sous-différentiables), qui généralisent les méthodes de descente de gradient.

1.3.1 Définition et caractérisation

Point et opérateur proximal Soit F une fonction convexe, s.c.i. et propre, définie sur l'espace euclidien X (avec $d \in \mathbb{N}^*$). On définit pour tout $x \in X$:

$$\text{prox}_F(x) = \underset{y \in X}{\text{argmin}} \left\{ \frac{1}{2} \|x - y\|^2 + F(y) \right\}.$$

appelé *point proximal* de x relativement à la fonction F . L'opérateur qui à tout x associe son point proximal relativement à F est appelé *opérateur proximal* (ou opérateur de proximité) associé à F .

Commençons par vérifier que, pour tout $x \in X$, le point proximal $\text{prox}_F(x)$ est bien défini. Soit $x \in X$. On remarque tout d'abord que la fonction

$$G : y \mapsto \frac{1}{2} \|x - y\|^2 + F(y)$$

est strictement convexe et s.c.i. car c'est la somme d'une fonction fortement convexe et d'une fonction convexe, toutes deux s.c.i. On en déduit que la fonction G admet un unique minimum, ce qui assure la bonne définition de $\text{prox}_F(x)$.

Caractérisation du point proximal On cherche à présent à donner une caractérisation plus manipulable de $\text{prox}_F(x)$. Comme il s'agit du minimiseur d'une fonction convexe, on peut le caractériser à l'aide de l'équation d'EULER. C'est l'objet du résultat suivant :

Proposition 8 Soit F une fonction convexe, s.c.i. et propre. Alors pour tout $x \in X$,

$$p = \text{prox}_F(x) \quad \iff \quad x - p \in \partial F(p)$$

DÉMONSTRATION : • Soit $x \in X$. Supposons que $p \in X$ vérifie $F(p) < +\infty$. Par définition du sous-gradient, $x - p \in \partial F(p)$ est défini par

$$\forall y \in X, \quad \langle y - p, x - p \rangle + F(p) \leq F(y).$$

Or, puisque $\langle y - p, x - p \rangle = \|y - p\|^2/2 + \|x - p\|^2/2 - \|y - x\|^2/2$, cette définition s'écrit

$$\frac{1}{2} \|y - p\|^2 + \frac{1}{2} \|x - p\|^2 - \frac{1}{2} \|x - y\|^2 + F(p) \leq F(y)$$

$$\text{soit} \quad \forall y \in X, \quad \frac{1}{2} \|y - p\|^2 + \frac{1}{2} \|x - p\|^2 + F(p) \leq \frac{1}{2} \|x - y\|^2 + F(y).$$

On en déduit en particulier que, si $x - p \in \partial F(p)$, alors

$$\forall y \in X, \quad \frac{1}{2} \|x - p\|^2 + F(p) \leq \frac{1}{2} \|x - y\|^2 + F(y)$$

ce qui, par définition de l'optimalité, assure que $p = \text{prox}_F(x)$.

• Supposons à présent que $p = \text{prox}_F(x)$. On remarque déjà que, dans ce cas, on a nécessairement $F(p) < +\infty$, car sinon, on aurait $\|p - y\|^2/2 + F(p) = +\infty$ ce qui est absurde pour le minimum d'une fonction propre. On peut donc définir le sous-différentiel de F en p . Par définition de l'optimalité, on a pour tout $y \in X$ et pour tout $\lambda \in]0; 1[$

$$\frac{1}{2} \|x - p\|^2 + F(p) \leq \frac{1}{2} \|x - [\lambda p + (1 - \lambda)y]\|^2 + F(\lambda p + (1 - \lambda)y).$$

On a alors d'une part

$$\begin{aligned} \|x - [\lambda p + (1 - \lambda)y]\|^2 &= \|(x - p) + (1 - \lambda)(p - y)\|^2 \\ &= \|x - p\|^2 + (1 - \lambda)^2 \|p - y\|^2 - 2(1 - \lambda) \langle y - p, x - p \rangle \end{aligned}$$

et d'autre part, par convexité,

$$F(\lambda p + (1 - \lambda)y) \leq \lambda F(p) + (1 - \lambda)F(y).$$

Il s'ensuit que, après simplification, on a pour tout $y \in X$ et pour tout $\lambda \in]0; 1[$,

$$(1 - \lambda) \langle y - p, x - p \rangle + (1 - \lambda)F(p) \leq \frac{1}{2} (1 - \lambda)^2 \|p - y\|^2 + (1 - \lambda)F(y).$$

En divisant par $1 - \lambda$ puis en faisant tendre λ vers 1, on montre alors que $x - p$ est un sous-gradient de F en p . ■

Puisque $p = \text{prox}_F(x)$ peut être caractérisé par la relation $x - p \in \partial F(p)$, soit, en d'autres termes, $x \in p + \partial F(p)$ (l'addition s'entendant de manière ensembliste), on peut formellement définir l'opérateur

$$\mathbf{I} + \partial F : \begin{cases} E & \rightarrow & \mathcal{P}(E) \\ p & \mapsto & p + \partial F(p). \end{cases}$$

Dans ce cas, $p = \text{prox}_F(x)$ équivaut alors à $x \in (\mathbf{I} + \partial F)(p)$, et comme p est unique, cela nous permet de noter l'opérateur proximal

$$\text{prox}_F = (\mathbf{I} + \partial F)^{-1}.$$

Projection sur un convexe fermé L'opérateur proximal peut être vu comme la généralisation de la projection sur un convexe fermé. Commençons par rappeler le résultat suivant :

Théorème 8 Soit $C \subset E$ un convexe fermé. Il existe une unique application proj_C de E dans C , appelée projection sur le convexe C , qui à tout $x \in E$ associe le point $\text{proj}_C(x)$ de C , telle que la distance de x à C soit égale à celle de x à $\text{proj}_C(x)$. Le vecteur $\text{proj}_C(x)$ est l'unique point de C vérifiant les deux propositions équivalentes suivantes :

$$\begin{aligned} \forall y \in C, \quad \|x - \text{proj}_C(x)\| &\leq \|x - y\| \\ \forall y \in C, \quad \langle x - \text{proj}_C(x), y - \text{proj}_C(x) \rangle &\leq 0. \end{aligned} \tag{1.1}$$

La relation (1.1) permet en particulier de caractériser le point $\text{proj}_C(x)$:

$$\text{proj}_C(x) = \underset{y \in C}{\text{argmin}} \left\{ \frac{1}{2} \|x - y\|^2 \right\}.$$

Or, si C est un convexe fermé de X , alors la fonction caractéristique χ_C définie par

$$\forall x \in X, \quad \chi_C(x) = \begin{cases} 0 & \text{si } x \in C \\ +\infty & \text{sinon} \end{cases}$$

est convexe. Montrons à présent que $\text{prox}_{\chi_C} = \text{proj}_C$. En effet, pour tout $x \in X$, on a par définition

$$\text{prox}_{\chi_C}(x) = \underset{y \in X}{\text{argmin}} \left\{ \frac{1}{2} \|x - y\|^2 + \chi_C(y) \right\} = \underset{y \in C}{\text{argmin}} \left\{ \frac{1}{2} \|x - y\|^2 \right\} = \text{proj}_C(x).$$

1.3.2 Identité de Moreau

Décomposition de Moreau Si $X \subset X$ est un sous-espace vectoriel de X , alors on a classiquement la décomposition orthogonale suivante⁴ :

$$E = E \oplus E^\perp$$

où X^\perp est l'ensemble des points qui sont orthogonaux aux points de X , c'est-à-dire

$$E^\perp = \{y \in X \mid \forall x \in X, \langle x, y \rangle = 0\}.$$

On a alors la relation suivante :

$$\forall x \in X, \quad x = \text{proj}_E(x) + \text{proj}_{E^\perp}(x).$$

D'après la remarque du paragraphe précédent, puisque X et X^\perp sont convexes, on peut réécrire la formule précédente en utilisant des opérateurs proximaux :

$$\forall x \in X, \quad x = \text{prox}_{\chi_E}(x) + \text{prox}_{\chi_{E^\perp}}(x).$$

Il est à présent utile de remarquer que χ_{E^\perp} est la conjuguée convexe de χ_E . En effet,

$$\forall y \in X, \quad (\chi_E)^*(y) = \sup_{x \in X} \{ \langle x, y \rangle - \chi_E(x) \} = \sup_{x \in X} \{ \langle x, y \rangle \} = \begin{cases} 0 & \text{si } y \in X^\perp \\ +\infty & \text{si } y \notin X^\perp \end{cases}$$

ce qui implique que $\forall x \in X, \quad x = \text{prox}_{\chi_E}(x) + \text{prox}_{(\chi_E)^*}(x)$.

On peut généraliser ce résultat à n'importe quel opérateur proximal : c'est l'identité de MOREAU.

4. que l'on peut généraliser à un sous-espace fermé d'un espace de HILBERT.

Proposition 9 (Identité de Moreau) Soit F une fonction convexe, s.c.i. et propre. Alors on a

$$\forall x \in X, \quad x = \text{prox}_F(x) + \text{prox}_{F^*}(x).$$

DÉMONSTRATION : Soit $x \in X$. Posons $u = \text{prox}_F(x)$ et $v = x - u$, et montrons que $v = \text{prox}_{F^*}(x)$. La caractérisation de u assure que $x - u \in \partial F(u)$, soit $v \in \partial F(u)$. Par définition du sous-différentiel, on en déduit que

$$\forall y \in X, \quad \langle y - u, v \rangle + F(u) \leq F(y)$$

soit
$$\forall y \in X, \quad \langle u, v \rangle - F(u) \geq \langle y, v \rangle - F(y)$$

Cette dernière relation étant valable pour tout $y \in X$, on peut passer à la borne supérieure, ce qui entraîne que

$$\langle u, v \rangle - F(u) \geq \sup_{y \in X} \{ \langle y, v \rangle - F(y) \} = F^*(v).$$

Par ailleurs, pour tout $y \in X$, on a par optimalité que $F^*(y) \geq \langle u, y \rangle - F(u)$, ce qui implique que

$$\forall y \in X, \quad \langle u, v - y \rangle - F^*(y) \geq F^*(v) \quad \text{soit} \quad \langle y - v, u \rangle + F^*(v) \leq F^*(y)$$

ce qui assure que $x - v \in \partial(F^*)(v)$, soit $v = \text{prox}_{F^*}(x)$. ■

On peut encore généraliser ce résultat, et montrer que pour tout $\gamma > 0$,

$$\forall x \in X, \quad x = \text{prox}_{\gamma F}(x) + \gamma \text{prox}_{F^*/\gamma}(x/\gamma).$$

Ce résultat est à nouveau appelé *identité de MOREAU*.

1.4 Optimisation convexe et méthodes proximales

Soit F une fonction convexe s.c.i. propre. On cherche à résoudre le problème d'optimisation suivant :

$$\min_{x \in X} F(x) \tag{1.2}$$

où F prend généralement la forme d'une somme de fonctions convexes, et est coercive. Elle admet alors un minimum. On s'attachera ici à étudier principalement les cas où F est une fonction convexe ou la somme de deux fonctions convexes (étant entendu qu'on peut très souvent se ramener à l'un de ces deux cas).

Étudions dans cette section différentes méthodes d'optimisation convexe connues à la lumière des opérateurs proximaux. On supposera que le problème étudié admet au moins une solution. Pour une revue plus complète sur ce sujet, le lecteur pourra se reporter à [6].

1.4.1 Méthodes de gradient

Ces méthodes nécessitent que les fonctions en jeu soient régulières (au moins sous-différentiables).

Gradient explicite Commençons par considérer le cas où F est une fonction convexe différentiable, de gradient continu sur X et lipschitzienne, de constante de LIPSCHITZ L . Tout minimum x^* satisfait la condition d'optimalité de premier ordre ($\nabla F(x^*) = 0$). On peut en particulier écrire, pour tout $\tau > 0$, la relation de point fixe

$$x^* - \tau \nabla F(x^*) = x^*.$$

On peut alors considérer l'algorithme itératif de recherche de point fixe :

$$x_0 \in X \quad \text{et} \quad \forall k \in \mathbb{N}, \quad x_{k+1} = x_k - \tau \nabla F(x_k) \quad (1.3)$$

connu sous le nom de *méthode (de descente) de gradient explicite* (explicite car on évalue le gradient de F au point précédent déjà connu). Il faut bien sûr s'assurer que ce schéma itératif converge. Par exemple, si F est supposée convexe et ∇F lipschitzienne de constante L , alors un développement de TAYLOR assure la convergence de ce schéma dès que $\tau < 2/L$.

La méthode du gradient explicite est une méthode bien connue pour résoudre le problème (1.2). Elle s'interprète de la manière suivante : la convexité de la fonction F assurant la non-existence de minima locaux, il suffit pour trouver x^* à partir d'un point x_0 de réaliser des pas de descentes, c'est-à-dire de s'approcher de x^* en suivant une direction dans laquelle F décroît. Pour cela, si F est différentiable, la meilleure direction (localement) est celle donnée par l'opposé du gradient. L'itération (1.3) revient donc à effectuer un pas (fixe) dans cette direction. On voit alors que le choix du pas τ est crucial, en particulier lorsqu'on s'approche de x^* : s'il est choisi trop grand, on dépasse x^* lorsqu'on s'en approche, tandis que, s'il est choisi trop petit, la convergence est trop longue. Notons enfin que le choix de τ dans les deux cas particuliers abordés repose sur la connaissance de certaines constantes (de LIPSCHITZ et éventuellement de la forte convexité de F) qui ne sont pas toujours accessibles.

Gradient implicite Relâchons l'hypothèse de régularité sur ∇F . On peut alors écrire, pour tout $\tau > 0$, la relation

$$x^* + \tau \nabla F(x^*) = x^* \quad \text{soit} \quad x^* + \partial(\tau F)(x^*) = \{x^*\}.$$

Ainsi, $(I + \partial(\tau F))(x^*) = x^*$, qu'on peut encore écrire $x^* = (I + \partial(\tau F))^{-1}(x^*)$. On en déduit que le minimum est caractérisé par la nouvelle relation de point fixe

$$x^* = \text{prox}_{\tau F}(x^*).$$

Si on sait calculer l'opérateur proximal de F , cette recherche de point fixe incite donc à proposer cette fois le schéma itératif suivant :

$$x_0 \in X \quad \text{et} \quad \forall k \in \mathbb{N}, \quad x_{k+1} = \text{prox}_{\tau F}(x_k) \quad (1.4)$$

Cet algorithme est appelé *algorithme du point proximal*. Il a été proposé pour la première fois en 1970 par Bernard MARTINET [11]. Réécrivons l'algorithme du point proximal autrement : par définition de l'opérateur proximal,

$$\forall n \in \mathbb{N}, \quad x_{k+1} = (I + \partial(\tau F))^{-1}(x_k)$$

qu'on peut écrire $(I + \partial(\tau F))(x_{k+1}) = x_{k+1} + \tau \nabla F(x_{k+1}) = x_k$

car ici, F est différentiable. L'algorithme proposé devient alors :

$$x_0 \in X \quad \text{et} \quad \forall n \in \mathbb{N}, \quad x_{k+1} = x_k - \tau \nabla F(x_{k+1})$$

plus connu sous le nom de *méthode (de descente) de gradient implicite*, car, écrite sous cette forme, pour trouver le point x_{k+1} , on descend dans la direction du gradient de F au point x_{k+1} que l'on est en train de calculer.

La méthode du gradient implicite est *a priori* plus difficile à mettre en place que celle du gradient explicite, car il s'agit d'évaluer le gradient en un point non connu. Si cette opération est réalisable, elle peut alors présenter un intérêt pour le cas où ∇F n'est pas lipschitzien (et pour lequel la convergence de la méthode de gradient explicite n'est pas assurée). On peut en effet considérer la fonction auxiliaire γF , définie pour tout $\gamma > 0$ par

$$\forall x \in X, \quad \gamma F(x) = \min_{y \in X} \left\{ \frac{1}{2\gamma} \|x - y\|^2 + F(y) \right\}$$

appelée *enveloppe de MOREAU* d'indice γ de la fonction F . Cette fonction est convexe, s.c.i. et propre, et est différentiable, de gradient

$$\forall x \in X, \quad \nabla \gamma F(x) = \frac{1}{\gamma} (x - \text{prox}_{\gamma F}(x)) = \text{prox}_{F^*/\gamma}(x/\gamma).$$

On montre alors en particulier que l'itération (1.4) s'écrit

$$x_{k+1} = x_k - \tau \nabla^T F(x_k)$$

qui s'interprète comme une itération de gradient explicite pour la fonction auxiliaire γF . On peut alors choisir un pas de descente τ assurant la convergence de cet algorithme si $\nabla^T F$ est lipschitzienne par exemple. C'est bien le cas ici, car on peut montrer que la constante de LIPSCHITZ de $\nabla^T F$ vaut $1/\tau$. La méthode de gradient implicite est donc plus stable que celle de gradient explicite. Un autre intérêt de cette approche réside dans le fait qu'elle est facilement généralisable au cas sous-différentiable, comme on le verra au paragraphe suivant.

Sous-gradient implicite Lorsque F n'est pas différentiable, mais sous-différentiable, la condition d'optimalité de premier ordre devient $0 \in \partial F(x^*)$, ce qui permet à nouveau d'écrire la caractérisation du minimum par la relation de point fixe suivante, valable pour tout $\tau > 0$,

$$x^* = \text{prox}_{\tau F}(x^*).$$

Cela conduit donc au même algorithme itératif :

$$x_0 \in X \quad \text{et} \quad \forall k \in \mathbb{N}, \quad x_{k+1} = \text{prox}_{\tau F}(x_k)$$

Cette méthode peut à nouveau être interprétée comme une méthode de *sous-gradient implicite*, car l'algorithme considéré peut s'écrire de manière équivalente :

$$x_0 \in X \quad \text{et} \quad \forall n \in \mathbb{N}, \quad x_{k+1} = x_k - \tau g_{k+1} \quad \text{où} \quad g_{k+1} \in \partial F(x_{k+1}).$$

Gradient explicite-implicite On se place maintenant dans le cas où F est de la forme $F = f + g$, avec f différentiable et g sous-différentiable ; on supposera de plus que ∇f est continue et L -lipschitzienne. On s'intéresse donc au problème suivant

$$\min_{x \in X} \left\{ f(x) + g(x) \right\}.$$

Il est possible d'appliquer la méthode de sous-gradient implicite à la fonction F , mais on choisit ici de présenter une méthode qui permet de tirer parti de la différentiabilité d'une partie de la fonction F .

La condition d'optimalité de premier ordre assure que, si on note x^* l'optimum, alors on a

$$0 \in \tau \partial F(x^*) \quad \text{soit} \quad 0 \in \tau \nabla f(x^*) + \tau \partial g(x^*)$$

i.e.
$$x^* - \tau \nabla f(x^*) - x^* \in \partial(\tau g)(x^*)$$

où τ est un réel strictement positif quelconque. On en déduit que le problème d'optimisation est équivalent au problème de recherche de point fixe

$$x^* = \text{prox}_{\tau g}(x^* - \tau \nabla f(x^*)).$$

On peut alors proposer l'algorithme itératif de recherche de point fixe

$$x_0 \in X \quad \text{et} \quad \forall k \in \mathbb{N}, \quad x_{k+1} = \text{prox}_{\tau g}(x_k - \tau \nabla f(x_k))$$

en choisissant correctement le paramètre τ afin d'en assurer la convergence. On décompose généralement cet algorithme en deux étapes :

1. on commence par évaluer un point intermédiaire $x_{k+1/2} = x_k - \tau \nabla f(x_k)$; cette étape ne met en jeu que la fonction f , et s'interprète comme une descente de gradient *explicite*;
2. on calcule ensuite $x_{k+1} = \text{prox}_{\tau g}(x_{k+1/2})$: ce calcul ne dépend que de la fonction g , et s'interprète d'après ce qui précède comme une descente de gradient *implicite*.

Cet algorithme est également connu sous le nom de *forward-backward splitting*. Tout comme pour la méthode de gradient explicite, le pas τ doit être choisi en fonction de la constante de LIPSCHITZ de la partie différentiable de F , afin d'assurer la convergence du schéma proposé.

Dans toutes ces méthodes de gradient, des stratégies de pas τ_n variables peuvent alors être envisagées pour accélérer la convergence. Il est également possible d'ajouter des pas de relaxation (qui consiste à utiliser des points intermédiaires, situés sur la droite reliant le point précédant et le point courant).

Application : méthode du gradient projeté On présente maintenant une application très classique des algorithmes présentés ci-dessus. Considérons une minimisation sur un convexe $C \subset X$, c'est-à-dire le problème

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

où f est différentiable, de gradient lipschitzien. En remarquant que la contrainte d'appartenant à un convexe peut être intégrée dans une fonction caractéristique χ_C , on peut se ramener au cas précédent en écrivant le problème comme le problème sans contrainte suivant

$$\min_{x \in X} F(x) \quad \text{avec} \quad F(x) = f(x) + g(x) \quad \text{et} \quad g = \chi_C.$$

La fonction g est ici sous-différentiable. Ainsi, on cherche à trouver le point fixe x^*

$$x^* = \text{prox}_{\tau g}(x^* - \tau \nabla f(x^*)) \quad \text{soit} \quad x^* = \text{prox}_g(x^* - \tau \nabla f(x^*))$$

Or, on rappelle que l'opérateur proximal associé à une fonction caractéristique d'un convexe est la projection sur ce même convexe; on cherche donc à résoudre

$$x^* = \text{proj}_C(x^* - \tau \nabla f(x^*)).$$

La méthode du gradient explicite-implicite s'écrit alors

$$x_0 \in X \quad \text{et} \quad \forall k \in \mathbb{N}, \quad x_{k+1} = \text{proj}_C(x_k - \tau \nabla f(x_k))$$

ce qui revient à effectuer dans un premier temps une descente de gradient (explicite), puis à projeter le point ainsi trouvé sur le convexe C , d'où l'appellation de *méthode du gradient projeté*.

1.4.2 Méthodes d'éclatement

On a vu avec la méthode de gradient implicite-explicite qu'il est parfois possible d'exploiter séparément les propriétés des termes composant la fonction F . Dans le cas cité, il s'agit de profiter de la différentiabilité d'un des deux termes, qui est une propriété de régularité plus forte que la sous-différentiabilité. Ce genre d'approches est connu sous le nom de méthode d'*éclatement* (*splitting* en anglais). On propose dans cette section deux autres méthodes d'éclatement classiques.

Méthode d'éclatement de Dykstra La méthode d'éclatement de DYKSTRA s'applique aux problèmes de la forme

$$\min_{x \in X} \left\{ F(x) + \frac{1}{2} \|x - u\|^2 \right\} \quad \text{avec} \quad F(x) = f(x) + g(x) \quad (1.5)$$

où u est un vecteur de X donné, f et g deux fonctions convexes sous-différentiables. Dans le cas où les fonctions f et g sont les fonctions caractéristiques d'ensembles convexes, on voit que ce problème s'interprète comme la projection sur l'intersection des deux convexes du vecteur u . C'est dans ce cadre que cette méthode a été initialement proposée (c'est pourquoi elle est également connue sous le nom de *méthode de projection de DYKSTRA*). Cette méthode peut être utilisée lorsque l'opérateur proximal associé à F n'est pas calculable (ou difficilement), mais que ceux associés à f et g respectivement le sont. L'idée est donc d'exploiter la calculabilité de ces deux opérateurs.

On peut réécrire le problème (1.5) en utilisant les conjuguées convexes : on commence par écrire que, pour tout $x \in X$,

$$F(x) + \frac{1}{2} \|x - u\|^2 = \sup_{x_1, x_2 \in X} \left\{ -f^*(x_1) - g^*(x_2) + \frac{1}{2} \|x - u\|^2 + \langle x, x_1 + x_2 \rangle \right\}.$$

Puisque

$$\frac{1}{2} \|x - u\|^2 + \langle x, x_1 + x_2 \rangle = \frac{1}{2} \|x + x_1 + x_2 - u\|^2 - \frac{1}{2} \|x_1 + x_2 - u\|^2 + \frac{1}{2} \|u\|^2,$$

on en déduit que le problème (1.5) est équivalent à

$$\min_{x \in X} \sup_{x_1, x_2 \in X} \left\{ -f^*(x_1) - g^*(x_2) + \frac{1}{2} \|x + x_1 + x_2 - u\|^2 - \frac{1}{2} \|x_1 + x_2 - u\|^2 + \frac{1}{2} \|u\|^2 \right\}$$

et donc également à

$$\max_{x_1, x_2 \in X} \inf_{x \in X} \left\{ -f^*(x_1) - g^*(x_2) + \frac{1}{2} \|x + x_1 + x_2 - u\|^2 - \frac{1}{2} \|x_1 + x_2 - u\|^2 + \frac{1}{2} \|u\|^2 \right\}.$$

Or, pour tout $x_1, x_2 \in X$, on remarque que

$$\begin{aligned} \inf_{x \in X} \left\{ -f^*(x_1) - g^*(x_2) + \frac{1}{2} \|x + x_1 + x_2 - u\|^2 - \frac{1}{2} \|x_1 + x_2 - u\|^2 + \frac{1}{2} \|u\|^2 \right\} \\ = -f^*(x_1) - g^*(x_2) - \frac{1}{2} \|x_1 + x_2 - u\|^2 + \frac{1}{2} \|u\|^2 \end{aligned}$$

où le minimum est atteint pour $x^* = u - x_1 - x_2$. On peut donc montrer que pour résoudre le problème (1.5), il suffit de résoudre le problème dual

$$\min_{x_1, x_2 \in X} \left\{ f^*(x_1) + g^*(x_2) + \frac{1}{2} \|x_1 + x_2 - u\|^2 \right\}. \quad (1.6)$$

Si on note (x_1^*, x_2^*) la solution de ce problème fortement convexe, la solution x^* du problème initial (primal) est alors donnée par

$$x^* = u - x_1^* - x_2^*.$$

Le problème dual (1.6) peut être résolu par minimisation alternée : on minimise le lagrangien dual par rapport (par exemple) à y_1 pour y_2 fixé, puis l'inverse (avec ou non une mise-à-jour de la première variable duale entre les deux minimisations). Chacune de ces deux minimisations partielles peut être interprétée comme l'évaluation des opérateurs proximaux associés à f et g respectivement.

Méthode de Douglas-Rachford On relâche dans ce paragraphe l'hypothèse de régularité sur les fonctions à minimiser, mais on suppose que les fonctions sont convexes, s.c.i. et propres, ce qui nous permet de manipuler leur conjuguée convexe. L'idée est à nouveau d'exploiter la possibilité d'évaluer les opérateurs proximaux de chacun des termes composant le lagrangien, alors que celui de la somme n'est pas calculable.

On cherche à minimiser le problème de la forme⁵

$$\min_{x \in X} \left\{ f(x) + g(x) \right\} \quad (1.7)$$

où les fonctions f et g sont toutes les deux convexes, s.c.i. et propres. On suppose par ailleurs que leur somme définit une fonction coercive. Ce problème admet donc au moins une solution. Pour la caractériser, on commence par noter que, puisque g est s.c.i. et propre, elle est égale à sa biconjuguée convexe, et le problème étudié devient donc

$$\min_{x \in X} \left\{ f(x) + g^{**}(x) \right\} = \min_{x \in X} \left\{ f(x) + \sup_{y \in X} \left\{ \langle x, y \rangle - g^*(y) \right\} \right\}$$

ce qui nous amène à considérer le problème de recherche de point-selle

$$\min_{x \in X} \sup_{y \in X} \left\{ f(x) + \langle x, y \rangle - g^*(y) \right\}.$$

Les solutions (x^*, y^*) de ce problème satisfont les équations d'EULER

$$-y^* \in \partial f(x^*) \quad \text{et} \quad x^* \in \partial g^*(y^*)$$

ce qui implique que, pour tout $\tau > 0$,

$$x^* - \tau y^* \in x^* + \tau \partial f(x^*) \quad \text{et} \quad y^* + \tau^{-1} x^* \in y^* + \tau^{-1} \partial g^*(y^*)$$

soit

$$x^* = (\text{I} + \tau \partial f)^{-1}(x^* - \tau y^*) = \text{prox}_{\tau f}(x^* - \tau y^*)$$

et

$$y^* = (\text{I} + \tau^{-1} \partial g^*)^{-1}(y^* + \tau^{-1} x^*) = \text{prox}_{\tau^{-1} g^*}(y^* + \tau^{-1} x^*).$$

L'identité de MOREAU assure alors que

$$\text{prox}_{\tau^{-1} g^*}(y^* + \tau^{-1} x^*) = y^* + \tau^{-1} x^* - \tau^{-1} \text{prox}_{\tau g}(\tau(y^* + \tau^{-1} x^*))$$

ce qui entraîne donc

$$y^* = y^* + \tau^{-1} x^* - \tau^{-1} \text{prox}_{\tau g}(\tau(y^* + \tau^{-1} x^*))$$

5. La méthode proposée par Jim DOUGLAS et Henri RACHFORD dans [9] en 1956 visait originellement à résoudre des problèmes linéaires de la forme $u = Ax + Bx$, avec A et B des matrices définies positives.

soit $\text{prox}_{\tau g}(x^* + \tau y^*) - x^* = 0$. Ainsi, les solutions du problème (1.7) sont caractérisées par (pour $\lambda > 0$)

$$\begin{cases} x^* = \text{prox}_{\tau f}(x^* - \tau y^*) \\ 0 = \lambda [\text{prox}_{\tau g}(2x^* - x^* + \tau y^*) - x^*] \end{cases}$$

soit, en posant $y = x^* - \tau y^*$ et en ajoutant y dans la dernière équation,

$$\begin{cases} x^* = \text{prox}_{\tau f}(y) \\ y = y + \lambda [\text{prox}_{\tau g}(2x^* - y) - x^*]. \end{cases} \quad (1.8)$$

Voici donc un algorithme proposé pour résoudre le problème (1.7) basée sur la relation (1.8) :

$$x_0 \in X, \quad \text{et} \quad \forall k \in \mathbb{N}, \quad \begin{cases} x_{k+1} = \text{prox}_{\tau f}(y_k) \\ y_{k+1} = y_k + \lambda_k [\text{prox}_{\tau g}(2x_{k+1} - y_k) - x_{k+1}] \end{cases}$$

où λ_k est choisi dans l'intervalle $[\varepsilon; 2 - \varepsilon]$, avec $\varepsilon > 0$.

1.4.3 Itérations de Bregman

Enfin, signalons la possibilité d'utiliser des distances de BREGMAN [3] pour proposer une variante des algorithmes proximaux. On peut en effet voir dans la définition du point proximal la norme au carré comme une distance ; les auteurs de [5] ont proposé de remplacer cette distance par une classe de pseudo-distances, étudiée par BREGMAN [4], on peut appliquer les méthodes proximales en remplaçant chaque calcul de point proximal par une itération dite de BREGMAN, mais en gagnant en vitesse de convergence.

Distance de Bregman Soit H une fonction strictement convexe et différentiable. On définit la pseudo-distance suivante, qu'on appellera désormais *distance de BREGMAN associée à la fonction H* :

$$\forall (x, y) \in E^2, \quad D_x^H(x, y) = H(y) - H(x) - \langle \nabla H(x), y - x \rangle.$$

Notons que cette distance n'est pas symétrique par rapport à x et y . Puisque la fonction H est supposée strictement convexe, on montre que la quantité $D_x^H(x, y)$ est toujours positive quels que soient x et y . Par ailleurs, il est immédiat que $D_x^H(x, x) = 0$.

On peut aisément vérifier que, si $H = \|\cdot\|^2$, alors $D_x^H(x, y) = \|x - y\|^2/2$.

Convergence des itérations de Bregman L'intérêt des distances de BREGMAN [14] réside dans le fait qu'il est possible de trouver un réel $\alpha > 0$ tel que

$$\forall (x, y) \in E^2, \quad D_x^{\alpha H}(x, y) \geq \frac{1}{2} \|x - y\|^2.$$

On supposera donc désormais que H est telle que l'inégalité précédente soit vraie pour $\alpha = 1$. Une conséquence de cette hypothèse est l'existence d'un unique minimum pour la fonction $y \mapsto F(y) + D_x^H(x, y)$ quelle que soit F une fonction convexe. Il est donc possible de définir une version généralisée de l'opérateur proximal, en remplaçant dans sa définition la norme euclidienne par une distance de BREGMAN :

$$\underset{y \in E}{\text{argmin}} \left\{ F(y) + D_x^H(x, y) \right\}.$$

On a de plus le résultat suivant :

Théorème 9 Soient F une fonction s.c.i., convexe et propre et $x \in X$. Si

$$x^* = \operatorname{argmin}_{y \in X} \left\{ F(y) + D_y^H(y, x) \right\}$$

alors on a

$$\forall y \in X, \quad F(y) + D_y^H(y, x) \geq F(x^*) + D_{x^*}^H(x^*, x) + D_y^H(y, x^*).$$

Cette propriété permet généralement de remplacer dans les algorithmes proximaux l'évaluation du point proximal par la minimisation de la fonction $y \mapsto F(y) + D_y^H(y, x)$ (lorsque celle-ci est calculable), tout en conservant la validité des preuves de convergence. Dans [5] par exemple, les auteurs ont démontré la convergence de l'algorithme du point proximal utilisant une distance de BREGMAN. L'intérêt de les utiliser peut être multiple [13] : il peut permettre de s'affranchir de certaines hypothèses de régularité (sur ∇F par exemple) dans les méthodes de gradient ; l'évaluation de l'opérateur proximal peut être plus simple lorsque la distance de BREGMAN est adaptée au problème (voir le paragraphe suivant).

Exemples de distances de Bregman Un premier exemple simple de distances de BREGMAN est donné par

$$H(x) = \|x\|_M = \sqrt{\langle Mx, x \rangle}$$

où M est une matrice symétrique définie positive. La fonction H définit dans ce cas une norme, et on vérifie que $D_x^H(x, y) = \|x - y\|_M^2 / 2$ est une distance de BREGMAN.

Considérons un second exemple. Commençons par introduire la fonction dite *d'entropie*, définie par

$$\forall x = (x_i) \in \mathbb{R}^d, \quad H(x) = \begin{cases} \sum_{i=0}^{d-1} x_i \ln x_i & \text{si } x \in \Sigma \\ +\infty & \text{sinon} \end{cases}$$

avec $0 \ln 0 = 0$ par convention, avec Σ le simplexe de \mathbb{R}^d , défini par l'ensemble des vecteurs x de \mathbb{R}^d à coefficients positifs et de somme 1. Posons $h(t) = t \ln t$ pour tout réel positif t . La fonction h est strictement convexe (car dérivable sur \mathbb{R}_+^* , de dérivée strictement croissante). On en déduit la stricte convexité de H , ce qui implique que H définit une distance de BREGMAN. Cette version est en particulier utilisée pour minimiser certaines fonctions (supposées différentiables) sur le simplexe, qui consiste à résoudre pour $x \in \mathbb{R}^d$

$$\operatorname{argmin}_{y \in \Sigma} f(y) = \operatorname{argmin}_{y \in \mathbb{R}^d} \left\{ f(y) + \chi_\Sigma(y) \right\}.$$

En posant $F = f + \chi_\Sigma$, écrivons les itérations de la méthode du sous-gradient implicite pour ce problème :

$$x_0 \in \mathbb{R}^d \quad \text{et} \quad \forall k \in \mathbb{N}, \quad x_{k+1} = \operatorname{prox}_{\tau F}(x_k)$$

où la mise-à-jour de x_{k+1} s'écrit explicitement

$$x_{k+1} = \operatorname{argmin}_{y \in \mathbb{R}^d} \left\{ \tau F(y) + \frac{1}{2} \|x_k - y\|^2 \right\}.$$

Si on remplace l'évaluation de cet opérateur proximal par une itération de BREGMAN, alors on est amené à résoudre à la place

$$x_{k+1} = \operatorname{argmin}_{y \in \mathbb{R}^d} \left\{ \tau F(y) + D_{x_k}^H(x_k, y) \right\}.$$

La définition de H assure que ce dernier problème s'écrit

$$x_{k+1} = \operatorname{argmin}_{y \in \mathbb{R}^d} \left\{ \tau f(y) + H(y) - H(x_k) - \langle \nabla H(x_k), y - x_k \rangle \right\}.$$

La fonctionnelle est différentiable sur Σ et on peut résoudre explicitement ce problème à l'aide d'un multiplicateur de LAGRANGE pour la contrainte (égalité) sur la somme des coefficients de x_{k+1} . Plus précisément, x_{k+1} est défini comme la solution du problème de minimisation

$$\min_{y \in \mathbb{R}^d} \left\{ \tau f(y) + \sum_{i=0}^{d-1} y_i \ln y_i - \sum_{i=0}^{d-1} (x_k)_i \ln (x_k)_i - \sum_{i=0}^{d-1} (1 + \ln(x_k)_i)(y_i - (x_k)_i) \right\}$$

sous la contrainte égalité
$$\sum_{i=0}^{d-1} y_i = 1.$$

Les conditions de KUHN-TUCKER s'écrivent pour ce problème

$$\forall i \in \llbracket 0; d-1 \rrbracket, \quad \tau \frac{\partial f}{\partial y_i}((x_{k+1})_i) + 1 + \ln(x_{k+1})_i - (1 + \ln(x_k)_i) + \lambda = 0$$

On en déduit que

$$\forall i \in \llbracket 0; d-1 \rrbracket, \quad \tau \frac{\partial f}{\partial y_i}((x_{k+1})_i) + \ln(x_{k+1})_i = \ln(x_k)_i - \lambda.$$

Ainsi, si les $\tau \frac{\partial f}{\partial y_i} + \ln$ sont inversibles, on a

$$\forall i \in \llbracket 0; d-1 \rrbracket, \quad (x_{k+1})_i = \left(\tau \frac{\partial f}{\partial y_i} + \ln \right)^{-1} \left(\ln(x_k)_i - \lambda \right)$$

et la contrainte égalité assure que λ doit être solution de l'équation

$$\sum_{i=0}^{d-1} \left(\tau \frac{\partial f}{\partial y_i} + \ln \right)^{-1} \left(\ln(x_k)_i - \lambda \right) = 1.$$

Cette méthode peut être envisagée dans le cas où $f(x) = \langle a, x \rangle$ par exemple.

Pour une revue récente sur les algorithmes de projection sur le simplexe, on pourra également se référer à [7].

Références

- [1] Grégoire ALLAIRE. *Analyse numérique et optimisation : Une introduction à la modélisation mathématique et à la simulation numérique*. Éditions École polytechnique, 2005.
- [2] Heinz H. BAUSCHKE and Patrick L. COMBETTES. *Convex analysis and monotone operator theory in HILBERT spaces*. Springer Science & Business Media, 2011.
- [3] Amir BECK and Marc TEBoulLE. Mirror descent and nonlinear projected subgradient methods for convex optimization. *Operations Research Letters*, 31(3) :167–175, 2003.

-
- [4] Lev M. BREGMAN. The relaxation method of finding the common point of convex sets and its application to the solution of problems in convex programming. *USSR Computational Mathematics and Mathematical Physics*, 7(3) :200–217, 1967.
- [5] Yair CENSOR and Stavros Andrea ZENIOS. Proximal minimization algorithm with D-functions. *Journal of Optimization Theory and Applications*, 73(3) :451–464, 1992.
- [6] Patrick L. COMBETTES and Jean-Christophe PESQUET. Proximal splitting methods in signal processing. In *Fixed-point algorithms for inverse problems in science and engineering*, pages 185–212. Springer, 2011.
- [7] Laurent CONDAT. Fast projection onto the simplex and the ℓ_1 ball. *Mathematical Programming*, pages 1–11, 2014.
- [8] Damek DAVIS and Wotao YIN. Faster convergence rates of relaxed Peaceman–Rachford and ADMM under regularity assumptions. *arXiv preprint arXiv :1407.5210*, 2014.
- [9] Jim DOUGLAS and Henry H. RACHFORD. On the numerical solution of heat conduction problems in two and three space variables. *Transactions of the American mathematical Society*, 82(2) :421–439, 1956.
- [10] Ivar EKELAND and Roger TÉMAM. *Convex Analysis and Variational Problems*, volume 28. SIAM, 1999.
- [11] Bernard MARTINET. Brève communication. régularisation d’inéquations variationnelles par approximations successives. *Revue Française d’Informatique et de Recherche Opérationnelle, série rouge*, 4(3) :154–158, 1970.
- [12] Jean-Jacques MOREAU. Proximité et dualité dans un espace hilbertien. *Bulletin de la Société Mathématique de France*, 93 :273–299, 1965.
- [13] Van Quang NGUYEN. *Méthodes d’éclatement basées sur les distances de BREGMAN pour les inclusions monotones composites et l’optimisation*. PhD thesis, Paris 6, 2015.
- [14] Paul TSENG. On accelerated proximal gradient methods for convex-concave optimization. *Submitted to SIAM J. Optim*, 2008.