



Optimisation pour l'apprentissage automatique

Clément W. Royer

Notes de cours - M2 Big Data - 2021/2022

- La version la plus récente de ces notes se trouve à l'adresse :
<https://www.lamsade.dauphine.fr/~croyer/ensdocs/TUN/PolyTUN.pdf>.
- Pour tout commentaire, merci d'envoyer un mail à `clement.royer@dauphine.psl.eu`.
Merci à Radhia Bessi pour ses retours et suggestions concernant ce document.

- **Historique du document**

- 2021.12.08 : Première version du chapitre 4.
- 2021.12.07 : Corrections mineures dans le chapitre 3.
- 2021.12.02 : Seconde version du chapitre 3.
- 2021.11.29 : Première version du chapitre 3, mise à jour du chapitre 2.
- 2021.11.24 : Première version du chapitre 2.
- 2021.11.23 : Mise à jour du chapitre 1.
- 2021.11.19 : Première version avec chapitre 1 et annexe.

- **Objectifs d'apprentissage**

- Savoir formaliser un problème d'optimisation et reconnaître des classes particulières de problèmes;
- Pouvoir proposer un algorithme dédié à une classe de problèmes spécifique;
- Identifier les caractéristiques des problèmes d'optimisation issus de la science des données.

Sommaire

1	Introduction à l'optimisation	7
1.1	Définition et contexte	7
1.1.1	Problème et processus d'optimisation	7
1.1.2	Le contexte actuel	8
1.2	Bases de l'optimisation mathématique	10
1.2.1	Formulation générale et premières définitions	10
1.2.2	Solutions et minima	11
1.2.3	Conditions d'optimalité	12
1.2.4	Cas des fonctions convexes	13
1.2.5	Cas des fonctions fortement convexes	14
1.3	Conclusion du chapitre 1	14
2	Optimisation différentiable et méthodes de gradient	15
2.1	Moindres carrés linéaires et régression	15
2.1.1	Compléments d'algèbre linéaire	15
2.1.2	Contexte de la régression linéaire	17
2.1.3	Problèmes aux moindres carrés linéaires	17
2.1.4	Lien avec la régression linéaire	18
2.2	Minimisation de fonctions quadratiques	19
2.2.1	Formulation et conditions d'optimalité	19
2.2.2	Descente de gradient pour les problèmes quadratiques	20
2.2.3	Algorithme du gradient conjugué	21
2.3	Descente de gradient	22
2.3.1	Algorithme	23
2.3.2	Choix de la taille de pas	24
2.3.3	Analyse de complexité pour la descente de gradient	25
2.3.4	Cas convexe et fortement convexe	27
2.4	Accélération	30
2.4.1	Introduction au concept de momentum	30
2.4.2	Méthode du gradient accéléré	30
2.4.3	Autres algorithmes accélérés	32
2.5	Conclusion du chapitre 2	34

3 Méthodes de gradient stochastique	35
3.1 Motivation	35
3.2 Méthode du gradient stochastique	36
3.2.1 Algorithme	36
3.2.2 Analyse	37
3.3 Réduction de variance	40
3.3.1 Variantes à lots (batch)	40
3.3.2 Autres variantes basées sur la réduction de variance	41
3.4 Méthodes de gradient stochastique pour l'apprentissage profond	42
3.4.1 Gradient stochastique avec momentum	42
3.4.2 AdaGrad	43
3.4.3 RMSProp	43
3.4.4 Adam	44
3.5 Conclusion	45
4 Optimisation non lisse et régularisation	46
4.1 Introduction : algorithme du perceptron	46
4.2 Optimisation non lisse	47
4.2.1 Des fonctions aux problèmes non lisses	47
4.2.2 Méthodes de sous-gradient	48
4.3 Régularisation	49
4.3.1 Problèmes régularisés	49
4.3.2 Régularisation et parcimonie	50
4.3.3 Méthodes proximales	51
4.4 Conclusion	53
Annexe A Notations et bases mathématiques	55
A.1 Notations	55
A.1.1 Conventions de notation générique	55
A.1.2 Notations scalaires et vectorielles	55
A.1.3 Notations matricielles	56
A.2 Éléments de mathématiques	56
A.2.1 Algèbre linéaire vectorielle	57
A.2.2 Algèbre linéaire matricielle	59
A.2.3 Calcul différentiel	62

Avant-propos

Le but de ce cours est de présenter les algorithmes d'optimisation les plus utilisés en apprentissage, et ce en lien avec des formulations classiques dans ce domaine. Il se veut découpé en plusieurs parties selon le découpage des séances de cours.

Ce cours n'est pas un cours de mathématiques, mais il repose sur plusieurs concepts élémentaires d'algèbre linéaire et de calcul différentiel. Ceux-ci ne seront pas rappelés en détail en cours, mais sont en revanche détaillés en annexe de ce polycopié. De même, les notations cruciales du polycopié sont décrites ci-après. L'ensemble des notations utilisées est disponible en Annexe [A.1](#).

Notations

- Les scalaires seront représentés par des lettres minuscules : $a, b, c, \alpha, \beta, \gamma$.
- Les vecteurs seront représentés par des lettres minuscules en **gras** : $\mathbf{a}, \mathbf{b}, \mathbf{c}, \boldsymbol{\alpha}, \boldsymbol{\beta}, \boldsymbol{\gamma}$.
- Les lettres majuscules en gras seront utilisées pour les matrices : $\mathbf{A}, \mathbf{B}, \mathbf{C}$.
- Les lettres majuscules cursives seront utilisées pour les ensembles : $\mathcal{A}, \mathcal{B}, \mathcal{C}$.
- La définition d'une nouvelle quantité ou d'un nouvel opérateur sera indiquée par $:=$.
- L'ensemble des entiers naturels sera noté \mathbb{N} , l'ensemble des entiers relatifs sera noté \mathbb{Z} .
- L'ensemble des réels sera noté \mathbb{R} . L'ensemble des réels positifs sera noté \mathbb{R}_+ et l'ensemble des réels strictement positifs sera noté \mathbb{R}_{++} .
- On notera \mathbb{R}^d l'ensemble des vecteurs à d composantes réelles, et on considérera toujours que d est un entier supérieur ou égal à 1.
- Un vecteur $\mathbf{w} \in \mathbb{R}^d$ sera pensé (par convention) comme un vecteur colonne. On notera $w_i \in \mathbb{R}$ sa i -ème coordonnée dans la base canonique de \mathbb{R}^d . On aura ainsi $\mathbf{w} = \begin{bmatrix} w_1 \\ \vdots \\ w_d \end{bmatrix}$. Étant donné un vecteur (colonne) $\mathbf{w} \in \mathbb{R}^d$, le vecteur ligne correspondant sera noté \mathbf{w}^T . On aura donc $\mathbf{w}^T = [w_1 \ \cdots \ w_n]$ et $[\mathbf{w}^T]^T = \mathbf{w}$.
- Le produit scalaire entre deux vecteurs de \mathbb{R}^d est défini par $\mathbf{w}^T \mathbf{v} = \mathbf{v}^T \mathbf{w} = \sum_{i=1}^d w_i v_i$. On définit ensuite la norme d'un vecteur $\mathbf{w} \in \mathbb{R}^d$ par $\|\mathbf{w}\| = \sqrt{\mathbf{w}^T \mathbf{w}}$.
- On notera $\mathbb{R}^{n \times d}$ l'ensemble des matrices à n lignes et d colonnes à coefficients réels, où n et d seront des entiers supérieurs ou égaux à 1. Une matrice $\mathbf{A} \in \mathbb{R}^{n \times n}$ est dite carrée (dans le cas général, on parlera de matrice rectangulaire).
- On identifiera les vecteurs de \mathbb{R}^n avec les matrices de $\mathbb{R}^{n \times 1}$.
- Étant donnée une matrice $\mathbf{A} \in \mathbb{R}^{d \times d}$, on notera A_{ij} le coefficient en ligne i et colonne j de la matrice. La diagonale de \mathbf{A} est formée par l'ensembles des coefficients A_{ii} pour $i = 1, \dots, \min\{n, d\}$. La notation $[A_{ij}]_{\substack{1 \leq i \leq n \\ 1 \leq j \leq d}}$ sera donc équivalente à \mathbf{A} . Sans ambiguïté sur la taille de la matrice, on notera simplement $[A_{ij}]$.
- Selon les besoins, on utilisera \mathbf{a}_i^T pour la i -ème ligne de \mathbf{A} ou \mathbf{a}_j pour la j -ème colonne de \mathbf{A} .
Selon le cas, on aura donc $\mathbf{A} = \begin{bmatrix} \mathbf{a}_1^T \\ \vdots \\ \mathbf{a}_n^T \end{bmatrix}$ ou $\mathbf{A} = [\mathbf{a}_1 \ \cdots \ \mathbf{a}_d]$.
- Pour une matrice $\mathbf{A} = [A_{ij}] \in \mathbb{R}^{n \times d}$, la matrice transposée de \mathbf{A} , notée \mathbf{A}^T , est la matrice de $\mathbb{R}^{d \times n}$ telle que

$$\forall i = 1 \dots n, \ \forall j = 1 \dots d, \quad A_{ji}^T = A_{ij}.$$
- Pour tout $d \geq 1$, \mathbf{I}_d correspond à la matrice identité $\mathbb{R}^{d \times d}$ (avec 1 sur la diagonale et 0 ailleurs).

Chapitre 1

Introduction à l'optimisation

Le but de ce cours est d'étudier les problèmes d'optimisation issus de la science des données, ainsi que les algorithmes qui leur sont appliqués. Le concept d'optimisation existe bien au-delà de l'apprentissage, mais il est particulièrement important pour ce dernier. De nombreuses tâches d'apprentissage peuvent en effet être formulées comme des problèmes d'optimisation, et ceux-ci peuvent être résolus efficacement par des algorithmes modernes. Ce chapitre présente les principes de base derrière la formulation et l'étude (à la fois mathématique et algorithmique) d'un problème d'optimisation.

1.1 Définition et contexte

1.1.1 Problème et processus d'optimisation

Avant d'être un concept mathématique, un problème d'optimisation peut être verbalisé : on en donne une décomposition ci-dessous.

Définition 1.1.1 *Un problème d'optimisation mathématique est formé par les trois éléments suivants.*

- Une **fondation** : il s'agit du critère qui permet de quantifier la qualité d'une décision. Selon le cas, une meilleure décision peut signifier une valeur de l'objectif plus faible (auquel cas on aura affaire à un problème de minimisation) ou plus élevée (il s'agira alors d'un problème de maximisation).
- Des **variables de décision** : ce sont les différents choix qu'il est possible d'effectuer, dont la valeur influence celle de la fonction objectif. On cherche les meilleures valeurs de ces variables relativement à l'objectif.
- Des **contraintes** : celles-ci spécifient des conditions qui doivent être satisfaites par les variables de décision pour que les valeurs correspondantes soient acceptables.

Tout processus d'optimisation passe nécessairement par une phase de **modélisation**, durant laquelle on transforme un problème concret en objet mathématique que l'on peut essayer de résoudre. Pour cela, on se base sur des **algorithmes d'optimisation** qui, bien que pouvant être formulés à la main, ont généralement vocation à être implémentés sur un ordinateur. Lorsque l'algorithme appliqué

au problème renvoie une réponse, celle-ci doit être examinée pour savoir si elle fournit une solution satisfaisante au problème : au besoin, le modèle pourra être modifié, et l'algorithme ré-exécuté. En pratique, le processus d'optimisation peut donc boucler après une phase d'**interprétation** ou de **discussion** avec des experts du domaine d'application dont le problème est issu.

Remarque 1.1.1 *L'optimisation numérique obéit généralement aux principes suivants :*

- ***Il n'y a pas d'algorithme universel*** : certaines méthodes seront très efficaces sur certains types de problèmes, et peu efficaces sur d'autres. L'étude de la structure du problème facilite le choix d'un algorithme.
- ***Il peut y avoir un fort écart entre la théorie et la pratique*** : le calcul en précision finie peut induire des erreurs d'arrondis qui conduisent à une performance en-deçà de celle prévue par la théorie mathématique.
- ***La théorie guide la pratique, et vice-versa*** : pour la plupart des problèmes, il est possible de définir des expressions mathématiques qui permettent de vérifier si l'on a trouvé une solution du problème. Celles-ci sont à la base de nombreux algorithmes que nous étudierons. De manière générale, l'étude théorique d'un problème permet des avancées garanties qui surpassent souvent les heuristiques. Réciproquement, la performance de certaines méthodes a amené les chercheurs en optimisation à en proposer une étude théorique, permettant ainsi d'expliquer le comportement pratique.

1.1.2 Le contexte actuel

Le domaine de l'optimisation, et notamment son volet numérique, a pris son essor durant la seconde moitié du XXe siècle; depuis les années 80, la théorie mathématique de l'optimisation s'est développée de manière significative, en même temps que des implémentations très efficaces étaient développées pour tirer partie des machines de calcul de l'époque. Cette tendance se poursuit encore aujourd'hui, mais elle s'accompagne également d'un changement de paradigme.

En sciences des données, les problèmes d'optimisation impliquent souvent l'utilisation d'un ensemble massif de données dans la définition de la fonction objectif. On cherche ainsi à construire un modèle ou une prédiction sur la base d'un échantillon, qui n'est pas nécessairement un reflet exact de la distribution sous-jacente des données. Dans un contexte de données potentiellement distribuées, dont le coût d'utilisation peut être drastiquement élevé, les algorithmes d'optimisation ne sont pas tous égaux : de fait, certains algorithmes ayant fait leurs preuves dans des contextes basés sur des modèles s'avèrent peu performants dans ce contexte guidé par les données (ou *data-driven*).

Exemple introductif Considérons par exemple le problème de la figure 1.1. Sur les trois figures, les ronds rouges et carrés bleus sont placés aux mêmes endroits : ils représentent des données (images, mots, etc...) que l'on souhaite séparer via une droite. On parle de classification binaire, car il s'agit d'affecter une couleur/une forme à chacun des points de l'espace : la droite est appelée un classificateur linéaire. Sur chacune figure, le classificateur sépare effectivement les ronds rouges des carrés bleus, et en ce sens réussit la mission fixée au départ. Cependant, si l'on considère que ces points ne sont que des échantillons d'un grand jeu de données, distribué selon les deux "blobs" rouges et bleus des deux dernières figures, il est clair que le meilleur classificateur est celui de la figure de droite.

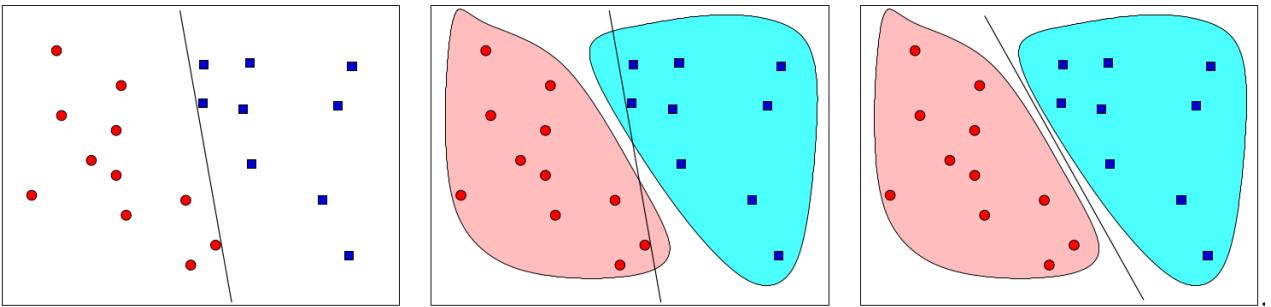


Figure 1.1: Différents classificateurs linéaires pour des échantillons (ronds rouges/points bleus) et leurs distributions sous-jacentes. Source : S. J. Wright, *Optimization Algorithms for Data Analysis* [6].

Il est possible de définir un cadre mathématique qui modélise ces différents aspects. On suppose que l'on dispose de n points \mathbf{x}_i en dimension d auxquels on affecte le label 1 s'ils correspondent à un rond rouge et -1 s'ils correspondent à un carré bleu : on notera $y(\mathbf{x}_i)$ le label du vecteur \mathbf{x}_i . Une modélisation possible du problème de classification est alors

$$\underset{\mathbf{w} \in \mathbb{R}^d}{\text{minimiser}} \frac{1}{n} \sum_{i=1}^n \max \{1 - y(\mathbf{x}_i) \mathbf{x}_i^T \mathbf{w}, 0\} \quad (1.1.1)$$

Il s'agit d'un problème d'optimisation (ici de minimisation), que l'on peut reformuler de manière convexe et résoudre efficacement via des algorithmes (par exemple de points intérieurs).

En résolvant ce problème, nous pouvons obtenir le classificateur de la partie gauche de la figure 1.1. Celui-ci est extrêmement sensible à une perturbation sur les données, et ne sépare pas bien les distributions. De fait, le problème que l'on souhaiterait résoudre porte sur une distribution inconnue \mathcal{D} des (\mathbf{x}_i, y_i) . Celui-ci se formulerait comme suit :

$$\underset{\mathbf{w} \in \mathbb{R}^d}{\text{minimiser}} \mathbb{E}_{\mathbf{x} \sim \mathcal{D}} [\max \{1 - y(\mathbf{x}) \mathbf{x}^T \mathbf{w}, 0\}] \quad (1.1.2)$$

où la moyenne empirique de (1.1.1) est remplacée dans (1.1.2) par une espérance mathématique.

Dans cet exemple, on peut identifier les différents défis posés par les problèmes d'optimisation de ce type :

- le problème que l'on souhaiterait résoudre est le problème sur les distributions (1.1.2), mais on ne peut accéder qu'à des versions échantillonées telles que (1.1.1);
- si l'on veut se rapprocher de la distribution, il est nécessaire de considérer un grand nombre d'échantillons, ce qui rendra le calcul de l'objectif du problème (1.1.1) très coûteux;
- le nombre de paramètres du modèle d peut être extrêmement important;
- pour certaines applications, l'utilisation d'un modèle linéaire ne sera pas suffisante, et l'on utilisera plutôt des formulations non linéaires.

Les points ci-dessus forment autant de difficultés mathématiques et numériques contre lesquelles doivent se prévaloir les algorithmes d'optimisation. Dans ce cours, nous effectuerons un tour d'horizon des techniques utilisées en optimisation qui ont fait leurs preuves dans des contextes industriels mais aussi modernes, comme l'apprentissage.

1.2 Bases de l'optimisation mathématique

1.2.1 Formulation générale et premières définitions

La transformation d'une description formelle d'un problème d'optimisation en objet mathématique conduit à l'écriture suivante :

$$\underset{\mathbf{w} \in \mathbb{R}^d}{\text{minimiser}} f(\mathbf{w}) \quad \text{s. c.} \quad \mathbf{w} \in \mathcal{F}, \quad (1.2.1)$$

où minimiser représente ce que l'on cherche à faire (on peut vouloir minimiser ou maximiser), $\mathbf{w} \in \mathbb{R}^d$ est un vecteur regroupant les variables de décision du problème, $f : \mathbb{R}^d \rightarrow \mathbb{R}$ est la fonction objectif qui mesure la qualité des décisions, et $\mathcal{F} \subset \mathbb{R}^d$ est l'ensemble des points réalisables ou admissibles, qui vérifient les contraintes posées sur les variables de décision.¹

Définition 1.2.1 i) Un point $\mathbf{w} \in \mathbb{R}^d$ tel que $\mathbf{w} \in \mathcal{F}$ est dit **admissible**, ou **réalisable**.

ii) Un point $\mathbf{w} \in \mathbb{R}^d$ tel que $\mathbf{w} \notin \mathcal{F}$ est dit **irréalisable**.

iii) Si $\mathcal{F} = \emptyset$, alors le problème (1.2.1) est dit **irréalisable**.

Définition 1.2.2 i) L'ensemble des solutions du problème (1.2.1) est noté

$$\underset{\mathbf{w} \in \mathbb{R}^d}{\text{argmin}} \{f(\mathbf{w}) \text{ s. c. } \mathbf{w} \in \mathcal{F}\}, \quad (1.2.2)$$

qu'on appelle **argument minimal** (ou **argmin**) du problème. C'est un sous-ensemble de \mathbb{R}^d (et de \mathcal{F}) : il s'agit de l'ensemble des points de \mathcal{F} qui donnent la valeur minimale de f .

ii) La **valeur minimale** (ou **le minimum**) du problème (1.2.1) est notée :

$$\min_{\mathbf{w} \in \mathbb{R}^d} \{f(\mathbf{w}) \text{ s. c. } \mathbf{w} \in \mathcal{F}\}. \quad (1.2.3)$$

Lorsque le problème possède une solution, la valeur minimale est la valeur de f en toute solution.

Remarque 1.2.1 Le problème (1.2.1) n'admet pas forcément de solution (prendre par exemple $d = 1$, $\mathcal{F} = \mathbb{R}$ et $f(w) = w$) : dans ce cas, l'argument minimal est l'ensemble vide, et la valeur minimale est $-\infty$.

On définit de la même manière que précédemment les problèmes de maximisation ainsi que les opérateurs argmax et max. Dans ce cours, on se concentrera sur les problèmes de minimisation, utilisant en cela la propriété ci-dessous.

Proposition 1.2.1 Pour toute fonction $f : \mathbb{R}^d \rightarrow \mathbb{R}$ et tout ensemble $\mathcal{F} \subset \mathbb{R}^d$, résoudre le problème de maximisation

$$\underset{\mathbf{w} \in \mathbb{R}^d}{\text{maximiser}} f(\mathbf{w}) \quad \text{s. c.} \quad \mathbf{w} \in \mathcal{F}$$

équivaut à résoudre le problème de minimisation

$$\underset{\mathbf{w} \in \mathbb{R}^d}{\text{minimiser}} -f(\mathbf{w}) \quad \text{s. c.} \quad \mathbf{w} \in \mathcal{F},$$

¹L'abréviation *s.c.* signifie "sous la/les contrainte(s)"; en anglais, on utilise *s.t.*, pour *subject to*.

dans la mesure où

$$\underset{\mathbf{w} \in \mathbb{R}^d}{\operatorname{argmax}} \{f(\mathbf{w}) s. c. \mathbf{w} \in \mathcal{F}\} = \underset{\mathbf{w} \in \mathbb{R}^d}{\operatorname{argmin}} \{-f(\mathbf{w}) s. c. \mathbf{w} \in \mathcal{F}\}$$

et

$$\max_{\mathbf{w} \in \mathbb{R}^d} \{f(\mathbf{w}) s. c. \mathbf{w} \in \mathcal{F}\} = - \min_{\mathbf{w} \in \mathbb{R}^d} \{-f(\mathbf{w}) s. c. \mathbf{w} \in \mathcal{F}\}.$$

Tout problème de maximisation se ramène donc à un problème de minimisation via une reformulation. Nous verrons d'autres exemples de reformulation, qui est une technique très utilisée pour transformer un problème en un autre problème équivalent mais que l'on saura mieux traiter théoriquement et/ou algorithmiquement.

Dans la majorité de ce cours, et notamment dans le reste de ce chapitre, on se concentrera sur des problèmes d'optimisation sans contraintes pour en caractériser plus précisément les solutions.

1.2.2 Solutions et minima

Soit le problème de minimisation sans contraintes

$$\underset{\mathbf{w} \in \mathbb{R}^d}{\operatorname{minimiser}} f(\mathbf{w}), \quad (1.2.4)$$

où $f : \mathbb{R}^d \rightarrow \mathbb{R}$.

Une solution du problème d'optimisation (1.2.4) est un vecteur qui conduit à la meilleure valeur possible de l'objectif. Pour un problème de minimisation sans contraintes, cela correspond au concept de minimum global.

Définition 1.2.3 (Minimum global) Un point $\mathbf{w}^* \in \mathbb{R}^d$ est un **minimum global** du problème (1.2.4), ou plus simplement de la fonction objectif f , lorsque la propriété suivante est vérifiée :

$$\forall \mathbf{w} \in \mathbb{R}^d, \quad f(\mathbf{w}^*) \leq f(\mathbf{w}). \quad (1.2.5)$$

Lorsque l'inégalité est stricte pour $\mathbf{w} \neq \mathbf{w}^*$, ce minimum est unique.

Déterminer le minimum global d'une fonction est un problème difficile en général : il est donc courant d'utiliser une autre notion (plus faible) de solution, basée sur l'optimalité d'un point relativement à son voisinage.

Définition 1.2.4 (Minimum local) Un point $\mathbf{w}^* \in \mathbb{R}^d$ est un **minimum local** du problème (1.2.4), ou plus simplement de la fonction objectif f , lorsque la propriété suivante est vérifiée :

$$\exists \epsilon > 0, \quad \forall \mathbf{w} \in \mathbb{R}^d, \quad \|\mathbf{w} - \mathbf{w}^*\| \leq \epsilon \Rightarrow f(\mathbf{w}^*) \leq f(\mathbf{w}). \quad (1.2.6)$$

Tout minimum global est ainsi un minimum local, mais l'inverse n'est pas vrai.

1.2.3 Conditions d'optimalité

Lorsqu'on cherche à minimiser une fonction f dérivable, il est possible de fournir des caractérisations pratiques de minima locaux (voire globaux) à l'aide des dérivées de f , qui fournissent une information locale sur la variation de la fonction. Ces caractérisations sont appelées **conditions d'optimalité**.

On s'intéresse d'abord aux fonctions dérivables, pour lesquelles on peut fournir une condition d'optimalité basée sur le gradient. Conformément à la suite du cours, nous donnerons le résultat dans le cas particulier des fonctions de classe \mathcal{C}^1 .

Théorème 1.2.1 (Condition nécessaire d'optimalité à l'ordre 1) Soient $f \in \mathcal{C}^1(\mathbb{R}^d)$ et $\mathbf{w}^* \in \mathbb{R}^d$. Alors

$$[\mathbf{w}^* \text{ minimum local de } f] \implies \nabla f(\mathbf{w}^*) = \mathbf{0}. \quad (1.2.7)$$

Tout minimum local de f est donc un point en lequel le gradient de f s'annule, mais l'inverse n'est pas forcément vraie ! Comme son nom l'indique, cette condition est nécessaire mais pas suffisante : il peut en effet exister des points en lesquels le gradient s'annule sans qu'il s'agisse de minima. On parle alors de maxima locaux (ex : 0 pour la fonction $x \mapsto -x^2$) ou de points selles (ex : 0 pour $x \mapsto x^3$) : cette dernière catégorie pose notamment des problèmes en apprentissage profond, où les problèmes à résoudre possèdent souvent de nombreux points selle.

Plus globalement, si f est une fonction \mathcal{C}^1 , un point \mathbf{w} tel que $\nabla f(\mathbf{w}) = \mathbf{0}$ s'appelle un **point stationnaire d'ordre 1**.

Si l'on suppose que f est deux fois dérivable avec dérivée seconde continue, on peut alors établir des conditions d'optimalité plus fortes que celles à l'ordre un.

Théorème 1.2.2 (Condition nécessaire d'optimalité à l'ordre 2) Soient $f \in \mathcal{C}^2(\mathbb{R}^d)$ et $\mathbf{w}^* \in \mathbb{R}^d$. Alors

$$[\mathbf{w}^* \text{ minimum local de } f] \implies \begin{cases} \nabla f(\mathbf{w}^*) = \mathbf{0}, \\ \nabla^2 f(\mathbf{w}^*) \succeq \mathbf{0}. \end{cases} \quad (1.2.8)$$

La nouvelle propriété fait intervenir la dérivée à l'ordre 2, qui est une matrice : pour un minimum local, la matrice hessienne est nécessairement *semi-définie positive*.²

Comme pour l'ordre 1, il est possible qu'un point \mathbf{w}^* vérifie $\nabla f(\mathbf{w}^*) = \mathbf{0}$ et $\nabla^2 f(\mathbf{w}^*) \succeq \mathbf{0}$ sans être un minimum local (voir 0 pour $w \mapsto w^3$ en dimension 1). Un point \mathbf{w} que $\nabla f(\mathbf{w}) = \mathbf{0}$ et $\nabla^2 f(\mathbf{w}) \succeq \mathbf{0}$ s'appelle un **point stationnaire à l'ordre 2** : l'ensemble des points stationnaires d'ordre 2 contient l'ensemble des minima locaux, mais n'y est pas forcément égal.

Contrairement à l'ordre 1, l'ordre 2 permet de définir une version suffisante des conditions d'optimalité, qui permet de reconnaître un minimum local.

Théorème 1.2.3 (Condition suffisante d'optimalité à l'ordre 2) Soient $f \in \mathcal{C}^2(\mathbb{R}^d)$ et $\mathbf{w}^* \in \mathbb{R}^d$. Alors

$$\left. \begin{array}{l} \nabla f(\mathbf{w}^*) = \mathbf{0}, \\ \nabla^2 f(\mathbf{w}^*) \succ \mathbf{0}. \end{array} \right\} \implies [\mathbf{w}^* \text{ minimum local de } f]. \quad (1.2.9)$$

Il suffit donc que le gradient en un point soit nul et que la matrice hessienne en ce point soit définie positive³ pour que ce point soit un minimum local.

²Une matrice symétrique $\mathbf{H} \in \mathbb{R}^{d \times d}$ est dite semi-définie positive si $\mathbf{v}^T \mathbf{H} \mathbf{v} \geq 0$ pour tout $\mathbf{v} \in \mathbb{R}^d$, ce que l'on note $\mathbf{H} \succeq \mathbf{0}$.

³Une matrice symétrique $\mathbf{H} \in \mathbb{R}^{d \times d}$ est dite définie positive si $\mathbf{v}^T \mathbf{H} \mathbf{v} > 0$ pour tout $\mathbf{v} \in \mathbb{R}^d$ non nul, ce que l'on note $\mathbf{H} \succ \mathbf{0}$.

C^1 : ensemble des fonctions continues et dérivable

1.2.4 Cas des fonctions convexes

Nous considérons maintenant une classe de fonctions particulièrement intéressantes en optimisation : les fonctions convexes. Comme on le verra dans ce cours, il est possible de développer des algorithmes très efficaces pour minimiser des fonctions convexes. Dans la présente section, on s'intéresse aux propriétés des problèmes de minimisation impliquant une fonction convexe.

La convexité est une propriété au départ géométrique, qui s'applique à des ensembles.

Définition 1.2.5 (Ensemble convexe) *Un ensemble $\mathcal{F} \subseteq \mathbb{R}^n$ est dit **convexe** si pour tout couple de points \mathcal{F} , le segment reliant ces deux points est inclus dans \mathcal{F} . Mathématiquement, cela s'écrit :*

$$\forall (\mathbf{w}, \mathbf{v}) \in \mathcal{F}^2, \quad \forall \alpha \in [0, 1], \quad \alpha \mathbf{w} + (1 - \alpha) \mathbf{v} \in \mathcal{F}. \quad (1.2.10)$$

Il est aussi possible de définir la convexité pour des fonctions.

Définition 1.2.6 (Fonction convexe) *Une fonction $f : \mathbb{R}^d \rightarrow \mathbb{R}$ est une **fonction convexe** sur un ensemble convexe $\mathcal{F} \subseteq \mathbb{R}^d$ si et seulement si*

$$\forall (\mathbf{w}, \mathbf{v}) \in \mathcal{F}^2, \quad \forall \alpha \in [0, 1], \quad f(\alpha \mathbf{w} + (1 - \alpha) \mathbf{v}) \leq \alpha f(\mathbf{w}) + (1 - \alpha) f(\mathbf{v}). \quad (1.2.11)$$

Exemple 1.2.1 (Fonctions convexes d'une variable)

1. La fonction $w \mapsto w^2$ est convexe sur \mathbb{R} .
2. La fonction $w \mapsto |w|$ est convexe sur \mathbb{R} .

Exemple 1.2.2 (Fonctions convexes en dimension d)

1. Si f est une norme sur \mathbb{R}^d , alors il s'agit d'une fonction convexe sur \mathbb{R}^d .
2. Une fonction quadratique $\mathbf{w} \mapsto \frac{1}{2} \mathbf{w}^T \mathbf{A} \mathbf{w} + \mathbf{b}^T \mathbf{w}$ avec $\mathbf{b} \in \mathbb{R}^d$ et $\mathbf{A} \in \mathbb{R}^{d \times d}$ semi-définie positive est convexe sur \mathbb{R}^d .

La notion de convexité est préservée par certaines opérations : en particulier, toute somme (et plus généralement toute combinaison linéaire à coefficients positifs) de fonctions convexes est une fonction convexe.

En plus de la définition (1.2.11), on peut caractériser la convexité au moyen des dérivées à l'ordre un et deux lorsque celles-ci existent.

Théorème 1.2.4 *Soit $f : \mathbb{R}^d \rightarrow \mathbb{R}$ une fonction de classe \mathcal{C}^1 et $\mathcal{F} \subset \mathbb{R}^d$ un ensemble convexe. Alors, f est convexe sur \mathcal{F} si et seulement si*

$$\forall (\mathbf{w}, \mathbf{v}) \in \mathcal{F}^2, \quad f(\mathbf{v}) \geq f(\mathbf{w}) + \nabla f(\mathbf{w})^T (\mathbf{v} - \mathbf{w}). \quad (1.2.12)$$

La propriété (1.2.12) joue un rôle fondamental en optimisation convexe.

Théorème 1.2.5 *Soit $f : \mathbb{R}^d \rightarrow \mathbb{R}$ une fonction de classe \mathcal{C}^2 et $\mathcal{F} \subset \mathbb{R}^d$ un ensemble convexe. Alors, f est convexe sur \mathbb{R}^d si et seulement si*

$$\forall \mathbf{w} \in \mathcal{F}, \quad \nabla^2 f(\mathbf{w}) \succeq \mathbf{0}, \quad (1.2.13)$$

càd que la matrice hessienne en tout point de \mathbb{R}^d est toujours semi-définie positive.

En termes d'optimisation, les fonctions convexes possèdent la propriété suivante, extrêmement intéressante du point de vue de l'optimisation.

Théorème 1.2.6 Soit $f : \mathbb{R}^d \rightarrow \mathbb{R}$ une fonction convexe sur \mathbb{R}^d . Alors, on a les propriétés suivantes :

- i) Tout minimum local de f est un minimum global.
- ii) Si f est de classe C^1 , alors tout point \mathbf{w}^* tel que $\nabla f(\mathbf{w}^*) = \mathbf{0}$ est un minimum global du problème.

1.2.5 Cas des fonctions fortement convexes

On peut définir une propriété encore plus puissante que la convexité, dont on verra qu'elle a un impact majeur sur les résultats d'optimisation. Il s'agit de la convexité forte, définie ci-dessous de trois manières, avec et sans dérivées.

Définition 1.2.7 (Fonction fortement convexe) Soit $\mathcal{F} \subset \mathbb{R}^d$ un ensemble convexe et $f : \mathbb{R}^d \rightarrow \mathbb{R}$. On dit que f est une fonction **fortement convexe** de paramètre $\mu > 0$, ou **μ -fortement convexe** sur \mathcal{F} , si et seulement si

$$\forall (\mathbf{w}, \mathbf{v}) \in \mathcal{F}^2, \quad \forall \alpha \in [0, 1], \quad f(\alpha\mathbf{w} + (1 - \alpha)\mathbf{v}) \leq \alpha f(\mathbf{w}) + (1 - \alpha)f(\mathbf{v}) - \mu \frac{\alpha(1 - \alpha)}{2} \|\mathbf{v} - \mathbf{w}\|^2.$$

Théorème 1.2.7 Soit $f : \mathbb{R}^d \rightarrow \mathbb{R}$ une fonction de classe C^1 , $\mathcal{F} \subset \mathbb{R}^d$ un ensemble convexe et $\mu > 0$. La fonction f est μ -fortement convexe sur \mathcal{F} si et seulement si

$$\forall (\mathbf{w}, \mathbf{v}) \in \mathcal{F}^2, \quad f(\mathbf{v}) \geq f(\mathbf{w}) + \nabla f(\mathbf{w})^\top (\mathbf{v} - \mathbf{w}) + \frac{\mu}{2} \|\mathbf{v} - \mathbf{w}\|^2. \quad (1.2.14)$$

Théorème 1.2.8 Soit $f : \mathbb{R}^d \rightarrow \mathbb{R}$ une fonction de classe C^2 , $\mathcal{F} \subset \mathbb{R}^d$ un ensemble convexe et $\mu > 0$. La fonction f est μ -fortement convexe sur \mathcal{F} si et seulement si

$$\forall \mathbf{w} \in \mathcal{F}, \quad \nabla^2 f(\mathbf{w}) \succeq \mu \mathbf{I}_d.$$

Les fonctions fortement convexes offrent un contexte encore plus favorable pour l'optimisation que les fonctions convexes; la raison en incombe au résultat ci-dessous.

Théorème 1.2.9 Si $f : \mathbb{R}^d \rightarrow \mathbb{R}$ est fortement convexe, elle possède un unique minimum global.

De plus, si f est de classe C^1 , ce minimum est l'unique solution de $\nabla f(\mathbf{w}) = \mathbf{0}$.

1.3 Conclusion du chapitre 1

L'optimisation est un outil mathématique qui permet de modéliser de nombreux problèmes en sciences des données et au-delà. L'optimisation comporte toujours une phase de modélisation, qui consiste à formuler mathématiquement le problème en définissant une fonction objectif, des variables de décisions et des contraintes. Cela permet de placer le problème dans une classe, et ainsi de caractériser ce qu'est une solution de ce problème. Dans un contexte d'optimisation continue, les concepts de dérivées et de convexité sont extrêmement utiles pour caractériser les solutions d'un problème donné.

Chapitre 2

Optimisation différentiable et méthodes de gradient

Le but de ce chapitre est de présenter les techniques classiques en optimisation différentiable, qui ont conduit à la plupart des méthodes utilisées en sciences des données. L'algorithme de base, dit de *descente de gradient* repose sur l'information donnée par les dérivées, ce qui permet d'obtenir de meilleurs points tant qu'un point stationnaire n'est pas atteint.

Nous commençons par introduire les problèmes quadratiques et aux moindres carrés, classiques en optimisation non linéaire et qui serviront de base à notre étude. Nous détaillons ensuite l'algorithme de descente de gradient, dont le champ d'application va au-delà des problèmes quadratiques.

2.1 Moindres carrés linéaires et régression

Dans cette partie, nous nous intéressons à une classe de problèmes d'optimisation qui joue un rôle essentiel en analyse de données : les problèmes aux moindres carrés linéaires.

2.1.1 Compléments d'algèbre linéaire

Afin d'étudier les problèmes aux moindres carrés, nous introduisons une décomposition classique de l'algèbre linéaire; celle-ci nous permettra de caractériser les solutions d'une certaine classe de problèmes d'optimisation.

Théorème 2.1.1 (Décomposition en valeurs singulières) *Toute matrice $X \in \mathbb{R}^{n \times d}$ admet une décomposition en valeurs singulières (SVD¹) de la forme*

$$X = U\Sigma V^T,$$

où $U \in \mathbb{R}^{n \times n}$ est orthogonale ($U^T U = I_n$), $V \in \mathbb{R}^{d \times d}$ est orthogonale ($V^T V = I_d$) et $\Sigma \in \mathbb{R}^{n \times d}$ est telle que $\Sigma_{ij} = 0$ si $i \neq j$ et les coefficients diagonaux, notés $\sigma_i = \Sigma_{ii}$, vérifient $\sigma_1 \geq \dots \geq$

¹Dans la suite, on utilisera fréquemment l'algorithme anglo-saxon SVD pour faire référence à la décomposition en valeurs singulières.

$\sigma_{\min\{n,d\}} \geq 0$, soit

$$\Sigma = \left[\begin{array}{cccc|c} \sigma_1 & 0 & \cdots & 0 & 0 \\ 0 & \ddots & & 0 & \vdots \\ 0 & \cdots & 0 & \sigma_{\min\{n,d\}} & 0 \\ \hline 0 & \cdots & \cdots & & 0 \end{array} \right]$$

L'ensemble des valeurs $\{\sigma_1, \dots, \sigma_{\min\{n,d\}}\}$ est appelé l'ensemble des valeurs singulières de la matrice X . Le plus grand entier r tel que $\sigma_r > 0$ et $\sigma_{r+1} = 0$ est appelé le rang de X (il vaut 0 si X est la matrice nulle).

Exemple 2.1.1 La décomposition en valeurs singulières d'une matrice de $\mathbb{R}^{3 \times 2}$ est de la forme

$$A = \underbrace{[u_1 \ u_2 \ u_3]}_U \overbrace{\left[\begin{array}{cc} \sigma_1 & 0 \\ 0 & \sigma_2 \\ 0 & 0 \end{array} \right]}^{\Sigma} \underbrace{[v_1^T \ v_2^T]}_{V^T}$$

où $\sigma_1 \geq \sigma_2 \geq 0$, les u_i forment une base orthonormée de \mathbb{R}^3 et les v_i forment une base orthonormée de \mathbb{R}^2 . Cette matrice est de rang 2 si $\sigma_2 > 0$, de rang 1 si $\sigma_1 > 0 = \sigma_2$ et de rang 0 sinon.

Tout comme la décomposition en valeurs singulières généralise la décomposition en valeurs propres, la pseudo-inverse généralise la notion d'inverse.

Théorème 2.1.2 (Formule de pseudo-inverse) Soit une matrice $X \in \mathbb{R}^{n \times d}$ et $U\Sigma V^T$ une décomposition en valeurs singulières de cette matrice, où $\Sigma \in \mathbb{R}^{n \times d}$ est de la forme :

$$\left[\begin{array}{cccc|c} \sigma_1 & 0 & \cdots & 0 & 0 \\ 0 & \ddots & & 0 & \vdots \\ 0 & \cdots & 0 & \sigma_r & 0 \\ \hline 0 & \cdots & \cdots & & 0 \end{array} \right]$$

avec $\sigma_1 \geq \cdots \geq \sigma_r > 0$ et $r \leq \min\{m, n\}$. Alors, la pseudo-inverse de X est donnée par

$$X^\dagger = V \Sigma^\dagger U^T, \tag{2.1.1}$$

où $\Sigma^\dagger \in \mathbb{R}^{n \times m}$ est la pseudo-inverse de Σ définie par

$$\Sigma^\dagger = \left[\begin{array}{cccc|c} \frac{1}{\sigma_1} & 0 & \cdots & 0 & 0 \\ 0 & \ddots & & 0 & \vdots \\ 0 & \cdots & 0 & \frac{1}{\sigma_r} & 0 \\ \hline 0 & \cdots & \cdots & & 0 \end{array} \right].$$

On étend cette définition aux matrices nulles $X = \mathbf{0}_{n \times d}$ en posant $X^\dagger = \mathbf{0}_{d \times n}$.

2.1.2 Contexte de la régression linéaire

On considère un jeu de données ayant n éléments ou individus, et on associe à chaque individu d caractéristiques² sous la forme d'un vecteur de \mathbb{R}^d . Soient $\mathbf{x}_1, \dots, \mathbf{x}_n$ ces vecteurs : on les regroupe alors sous la forme d'une matrice de données

$$\mathbf{X} = \begin{bmatrix} \mathbf{x}_1^T \\ \vdots \\ \mathbf{x}_n^T \end{bmatrix} \in \mathbb{R}^{n \times d}. \quad (2.1.2)$$

Exemple 2.1.2 • Chaque ligne de \mathbf{X} représente un individu, et les d composantes de \mathbf{x}_i sont des données médicales (âge, poids, taux de cholestérol, etc).

- Chaque ligne de \mathbf{X} est une “vectorisation” d'une image 2D, et les valeurs de \mathbf{x}_i sont celles des pixels, en niveau de gris. Ainsi, une image de taille 480*640 serait transformée (en mettant les lignes bout à bout, par exemple) en un vecteur de taille $d = 480 * 640 = 307200$.

Dans un contexte d'apprentissage supervisé, on associe chaque vecteur de caractéristiques \mathbf{x}_i à un label $y_i \in \mathbb{R}$, qui peut représenter une classe à laquelle l'individu appartient (malade/non malade, image de chien ou de chat, etc). Ces labels sont concaténés pour former un vecteur de labels $\mathbf{y} \in \mathbb{R}^n$.

Par conséquent, notre but n'est plus seulement d'analyser l'information de la matrice \mathbf{X} , mais bien de trouver une relation entre les caractéristiques \mathbf{X} et les labels \mathbf{y} . Si l'on postule que cette relation est linéaire, on va donc chercher une fonction $h : \mathbb{R}^d \rightarrow \mathbb{R}$ de la forme $h(\mathbf{x}) = \mathbf{x}^T \mathbf{w}$ ³. On souhaite que h permette d'obtenir les y_i à partir des \mathbf{x}_i , c'est-à-dire que l'on voudrait avoir

$$h(\mathbf{x}_i) = \mathbf{x}_i^T \mathbf{w} = y_i \quad \forall i = 1, \dots, n,$$

que l'on peut ré-écrire sous forme matricielle comme

$$\mathbf{X} \mathbf{w} = \mathbf{y}. \quad (2.1.3)$$

On se trouve donc en présence d'un système linéaire que l'on va vouloir résoudre. Rien ne garantit a priori que ce système possède une solution, ou que cette solution (si elle existe) est unique : de fait, en présence de bruit dans les données, il est peu probable d'avoir une unique solution. La notion de solution à ce système n'est donc pas la plus pertinente. On va donc utiliser l'optimisation pour définir des notions plus appropriées.

2.1.3 Problèmes aux moindres carrés linéaires

Comme vu en section 2.1.2, on cherche à déterminer un modèle linéaire qui puisse expliquer nos données; lorsqu'il n'est pas possible de les expliquer de manière exacte, on va chercher le modèle linéaire qui explique *au mieux* nos données. Comme on cherche au départ à satisfaire la condition $\mathbf{x}_i^T \boldsymbol{\beta} = y_i$, on va maintenant chercher à avoir $|\mathbf{x}_i^T \boldsymbol{\beta} - y_i| \approx 0$. Ces considérations conduisent à la classe de problèmes d'optimisation suivante.

²Ou *features* en anglais.

³Pour simplifier, on considérera un modèle linéaire et non affine, c'est-à-dire sans terme constant.

Définition 2.1.1 (Moindres carrés linéaires) *Un problème aux moindres carrés linéaires est de la forme*

$$\underset{\mathbf{w} \in \mathbb{R}^d}{\text{minimiser}} \frac{1}{2} \|\mathbf{X}\mathbf{w} - \mathbf{y}\|^2 = \frac{1}{2} \sum_{i=1}^n (\mathbf{x}_i^\top \mathbf{w} - y_i)^2, \quad (2.1.4)$$

où $\mathbf{X} \in \mathbb{R}^{n \times d}$ et $\mathbf{y} \in \mathbb{R}^n$.

Ce problème modélise bien notre objectif, qui est de déterminer un vecteur \mathbf{w} qui permette d'expliquer au mieux nos données.

Théorème 2.1.3 (Solution des moindres carrés linéaires) *On considère le problème (2.1.4) et le vecteur $\mathbf{w}^* = \mathbf{X}^\dagger \mathbf{y}$. On a alors les propriétés suivantes :*

- 1) *Si $\text{rang}(\mathbf{X}) = d \leq n$, alors \mathbf{w}^* est l'unique solution du problème aux moindres carrés. C'est l'unique solution du système linéaire (2.1.3) lorsque $d = n$, et c'est une solution de ce système lorsque $d < n$ et $\mathbf{y} \in \text{Im}(\mathbf{X})$.*
- 2) *Si $\text{rang}(\mathbf{X}) = n < d$, alors \mathbf{w}^* est à la fois une solution de (2.1.3) et de norme minimale pour (2.1.4), ce qui signifie que*

$$\mathbf{w}^* \in \underset{\mathbf{w} \in \mathbb{R}^d}{\text{argmin}} \{ \|\mathbf{w}\| \text{ s. c. } \mathbf{X}\mathbf{w} = \mathbf{y} \}.$$

Les problèmes (2.1.3) et (2.1.4) admettent chacun une infinité de solutions.

- 3) *Si $\text{rang}(\mathbf{X}) < \min\{m, n\}$ et $\mathbf{y} \in \text{Im}(\mathbf{X})$, alors \mathbf{w}^* est une solution du système linéaire et de norme minimale au sens des moindres carrés. Les problèmes (2.1.3) et (2.1.4) admettent chacun une infinité de solutions.*
- 4) *Si $\text{rang}(\mathbf{X}) < \min\{m, n\}$ et $\mathbf{y} \notin \text{Im}(\mathbf{X})$, alors il n'existe pas de solution au système linéaire (2.1.3); en revanche, \mathbf{w}^* est la solution de norme minimale au sens des moindres carrés.*

2.1.4 Lien avec la régression linéaire

La régression linéaire est une technique classique d'analyse de données. Elle considère un jeu de données $\{(\mathbf{x}_i, y_i)\}_{i=1}^n$, où $\mathbf{x}_i \in \mathbb{R}^d$ et $y_i \in \mathbb{R}$, et vise à construire un modèle linéaire qui explique au mieux les données. La régression linéaire avec une fonction de perte ℓ_2 correspond au problème :

$$\underset{\mathbf{w} \in \mathbb{R}^d}{\text{minimiser}} \frac{1}{2n} \|\mathbf{X}\mathbf{w} - \mathbf{y}\|^2 = \frac{1}{2n} \sum_{i=1}^n (\mathbf{x}_i^\top \mathbf{w} - y_i)^2.$$

Ce problème est identique au problème de moindres carrés linéaires (2.1.4), et peut donc être résolu par les mêmes outils.

Dans le contexte classique de la régression linéaire, on suppose que les données proviennent réellement d'un modèle linéaire mais entaché d'un bruit :

$$\mathbf{y} = \mathbf{X}\mathbf{w}^* + \boldsymbol{\epsilon},$$

où $\boldsymbol{\epsilon} \in \mathcal{N}(\mathbf{0}, \mathbf{I})$ est un vecteur avec des composantes i.i.d. (indépendantes, identiquement distribuées) suivant une loi normale de moyenne 0 et de variance 1. La figure 2.1 propose un tel échantillon.

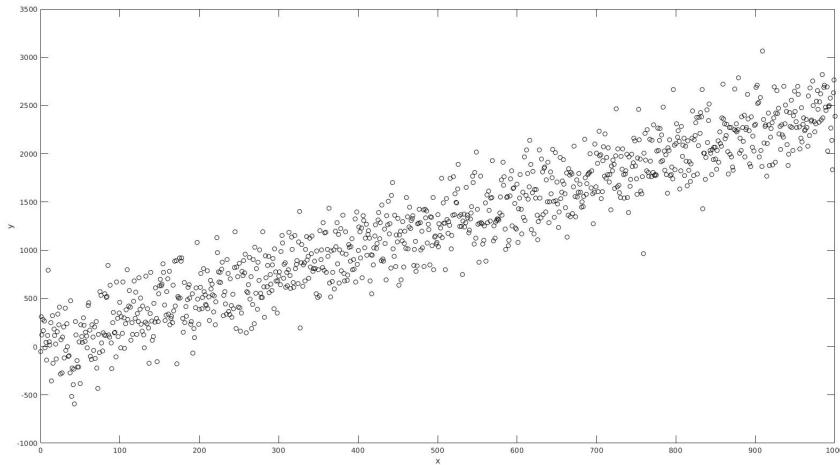


Figure 2.1: Données générées à partir d'un modèle linéaire entaché d'un bruit gaussien.

Comme pour cette figure, on supposera que $n \gg d$ et que $\mathbf{X}^T \mathbf{X} \in \mathbb{R}^{d \times d}$ est inversible. Dans ce contexte, on cherche à obtenir la valeur la plus vraisemblable de \mathbf{w} étant données les observations, ce qui est obtenu en résolvant :

$$\underset{\mathbf{w} \in \mathbb{R}^d}{\text{maximiser}} L(y_1, \dots, y_n; \mathbf{w}) := \left[\frac{1}{\sqrt{2\pi}} \right]^m \exp \left(-\frac{1}{2} \sum_{i=1}^m (\mathbf{x}_i^T \mathbf{w} - y_i)^2 \right). \quad (2.1.5)$$

Ce problème est un problème de maximisation, qui possède (sous les hypothèses énoncées plus haut) une unique solution appelée l'**estimateur du maximum de vraisemblance**, ou *maximum likelihood estimator* en anglais. On peut montrer que les problèmes (2.1.5) et (2.1.4) ont la même solution, ce qui signifie ici que l'estimateur du maximum de vraisemblance est $\mathbf{X}^\dagger \mathbf{y} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}$. Cet estimateur possède de nombreuses propriétés statistiques (par exemple, sa moyenne par rapport aux données est égale à \mathbf{w}^*). On voit ainsi que la formulation des moindres carrés, qui provient purement de l'optimisation, répond à un problème bien posé dans le langage statistique.

2.2 Minimisation de fonctions quadratiques

P2 : quadratique

Le problème aux moindres carrés étudié plus haut est un cas particulier d'un problème de minimisation dit **quadratique**, car il implique une fonction qui est un polynôme de degré 2 en les variables de décision. Il s'agit de la forme la plus simple des problèmes d'optimisation non linéaires : les problèmes quadratiques sont en fait fortement liés à la résolution de systèmes linéaires, comme on le verra dans les sections ci-dessous.

2.2.1 Formulation et conditions d'optimalité

Définition 2.2.1 (Problème d'optimisation quadratique) Un problème d'optimisation quadratique est de la forme

$$\underset{\mathbf{w} \in \mathbb{R}^d}{\text{minimiser}} q(\mathbf{w}) := \frac{1}{2} \mathbf{w}^T A \mathbf{w} + \mathbf{b}^T \mathbf{w} + c, \quad (2.2.1)$$

① La fct cout :

$$q(w) = \frac{1}{2} w^T A w + b^T w + c \Rightarrow \text{fct cout.}$$

fct objective.

1) le terme $\frac{1}{2} w^T A w$ est quadratique en w .

$$w^T A w = \text{produit scalaire}$$

A = matrice symétrique positive.

ou semi-definie positive

2) le terme $b^T w$ est linéaire en w .

3) c est une constante indépendante de w

Propriété mathématique

1- la fct $q(w)$ est **convexe** si la matrice A est positive définie ou semi-positive définie.
 \Rightarrow convexe garantie l'existence et l'unicité d'un **minimum global**.

② Convexe Vs convex

Convexe

1) fct est convexe $\Rightarrow A$ est semi-positive définie \Rightarrow Il existe des minima locaux

2) fct est fortement convexe $\Rightarrow A$ est positive définie \Rightarrow ! minimum = minimum global

où $A \in \mathbb{R}^{d \times d}$ est une matrice symétrique, $b \in \mathbb{R}^d$ et $c \in \mathbb{R}$.

Proposition 2.2.1 Soit le problème (2.2.1). Alors :

- L'ensemble des solutions du problème ne change pas avec la valeur de c ;
- La fonction objectif q est de classe C^∞ ;
- La fonction objectif q est minorée si la matrice A est semi-définie positive;
- La fonction objectif q est convexe si et seulement si A est semi-définie positive; elle est fortement convexe si et seulement si A est définie positive.

Dans la suite, on travaillera avec des fonctions quadratiques convexes, pour que le problème de minimisation ait un intérêt en soi (les fonctions quadratiques non convexes apparaissent en optimisation non linéaire, et ont un intérêt dans ce contexte).

Hypothèse 2.2.1 La fonction objectif du problème (2.2.1) est convexe.

L'application des conditions d'optimalité vues en section 1.2 conduit au résultat suivant.

Théorème 2.2.1 (Condition d'optimalité à l'ordre un) On considère le problème (2.2.1), sous l'hypothèse 2.2.1 et un point $w^* \in \mathbb{R}^d$. Alors w^* est un minimum global de q si et seulement si

$$Aw^* = -b. \quad (2.2.2)$$

On voit donc que la recherche d'un minimum d'un problème quadratique convexe se ramène à la résolution d'un système linéaire. On peut alors utiliser les techniques décrites au début de ce chapitre pour déterminer de manière explicite la solution de ce problème.

Corollaire 2.2.1 Sous les hypothèses du théorème 2.2.1, on considère le vecteur $w^* = -A^\dagger b$.

- Si la fonction objectif est fortement convexe, $w^* = -A^{-1}b$, est l'unique minimum global du problème (2.2.1).
- Si la fonction objectif est convexe, w^* est un minimum global; de plus, parmi toutes les solutions du problème, il s'agit de celle avec la plus petite norme euclidienne :

$$\forall \hat{w} \in \underset{w}{\operatorname{argmin}} q(w), \quad \|\hat{w}\| \leq \|w^*\|.$$

2.2.2 Descente de gradient pour les problèmes quadratiques

En pratique, il peut être trop coûteux de résoudre le problème quadratique (2.2.1) de manière exacte, c'est-à-dire en résolvant le système linéaire (2.2.2). On remplace donc cette résolution dite directe par une résolution itérative, en construisant une suite de points qui converge vers une solution.

L'une des manières classiques de procéder, qui sera décrite de manière plus générale en Section 2.3, consiste à effectuer des pas dans la direction opposée au gradient de la quadratique. En effet, la fonction

$$q : w \mapsto \frac{1}{2} w^T A w + b^T w$$



est minimale en une solution du système linéaire $\mathbf{A}\mathbf{w} = -\mathbf{b}$, et une telle solution vérifie $\nabla q(\mathbf{w}) = \mathbf{A}\mathbf{w} + \mathbf{b} = \mathbf{0}$.

Par conséquent, si \mathbf{w} est tel que $\nabla q(\mathbf{w}) \neq \mathbf{0}$, alors ce vecteur ne peut pas être un minimum global (ni local), et il doit être possible de trouver un meilleur point. On peut montrer qu'une direction de décroissance de la fonction est donnée par $-\nabla q(\mathbf{w})$: la méthode de descente de gradient effectue alors un déplacement dans cette direction.

Algorithme 1: Descente de gradient pour une fonction quadratique

Entrées : $\mathbf{A} \in \mathbb{R}^{d \times d}$ symétrique semi-définie positive, $\mathbf{b} \in \mathbb{R}^d$.

Initialisation : Choisir $\mathbf{w}_0 \in \mathbb{R}^d$.

Pour $k = 0, 1, \dots$

1. Si $\mathbf{A}\mathbf{w}_k + \mathbf{b} = \mathbf{0}$ terminer.
2. Choisir $\alpha_k > 0$.
3. Calculer $\mathbf{w}_{k+1} = \mathbf{w}_k - \alpha_k(\mathbf{A}\mathbf{w}_k + \mathbf{b})$.

Fin

L'algorithme 1 donne une illustration de la descente de gradient sur un problème quadratique (supposé convexe même si cela n'est pas requis par l'algorithme). Lorsque le pas α_k est bien choisi, on peut montrer que cette méthode converge vers une solution du problème, et on peut même quantifier combien d'itérations sont nécessaires pour atteindre une solution approchée.

2.2.3 Algorithme du gradient conjugué

La méthode du gradient conjugué est l'une des méthodes les plus utilisées pour résoudre un problème de la forme (2.2.1) lorsque \mathbf{A} est symétrique définie positive. On considérera ainsi des problèmes de la forme

$$\underset{\mathbf{w} \in \mathbb{R}^d}{\text{minimiser}} \frac{1}{2} \mathbf{w}^T \mathbf{A} \mathbf{w} + \mathbf{b}^T \mathbf{w}, \quad (2.2.3)$$

où $\mathbf{A} = \mathbf{A}^T \succ 0$ et $\mathbf{b} \in \mathbb{R}^d$.

Le principe de la méthode du gradient conjugué consiste à construire une famille de vecteurs formant une base de \mathbb{R}^d \mathbf{p}_i telle que la solution du problème (2.2.3) s'obtienne par combinaison linéaire de ces directions. Mathématiquement, l'algorithme du gradient conjugué résout (2.2.3) dans des sous-espaces vectoriels emboîtés de \mathbb{R}^n , que l'on appelle espaces de Krylov. À l'itération j de l'algorithme, on calcule ainsi

$$\mathbf{w}_j \in \underset{\mathbf{w} \in \mathcal{K}_j}{\text{argmin}} \frac{1}{2} \mathbf{w}^T \mathbf{A} \mathbf{w} + \mathbf{b}^T \mathbf{w},$$

où \mathcal{K}_j est le sous-espace engendré par les vecteurs $\mathbf{b}, \mathbf{A}\mathbf{b}, \dots, \mathbf{A}^{j-1}\mathbf{b}$. On peut montrer que $\mathcal{K}_d = \mathbb{R}^d$, ce qui conduit au résultat suivant.

Théorème 2.2.2 *Si l'algorithme du gradient conjugué est appliqué au problème (2.2.3) avec \mathbf{A} symétrique définie positive, alors il termine en au plus n itérations.*

Algorithme 2: Méthode du gradient conjugué

Entrées : $A \in \mathbb{R}^{d \times d}$ symétrique réelle, $\mathbf{b} \in \mathbb{R}^d$, $\zeta \in [0, 1]$.

Initialisation : Poser $\mathbf{w}_0 = \mathbf{0}$, $\mathbf{r}_0 = A\mathbf{w}_0 + \mathbf{b} = \mathbf{b}$, $\mathbf{p}_0 = -\mathbf{r}_0 = -\mathbf{b}$, $j = 0$.

Tant que $\|\mathbf{r}_j\| > \zeta \|\mathbf{r}_0\|$

1. $\alpha_j = \frac{\|\mathbf{r}_j\|^2}{\mathbf{p}_j^T A \mathbf{p}_j};$
2. $\mathbf{w}_{j+1} = \mathbf{w}_j + \alpha_j \mathbf{p}_j;$
3. $\mathbf{r}_{j+1} = \mathbf{r}_j + \alpha_j A \mathbf{p}_j;$
4. $\beta_{j+1} = \frac{\|\mathbf{r}_{j+1}\|^2}{\|\mathbf{r}_j\|^2};$
5. $\mathbf{p}_{j+1} = -\mathbf{r}_{j+1} + \beta_{j+1} \mathbf{p}_j;$
6. $j \leftarrow j + 1.$

Fin

En pratique, une implémentation du gradient conjugué suivrait le schéma décrit par l'algorithme 2. On se donne une précision voulu sur la norme du résidu (qui correspond au gradient de la fonction objectif), que l'on souhaite réduire à une certaine fraction, potentiellement nulle, de la valeur initiale. Tant que cette valeur n'est pas atteinte (ce qui nous garantit que le gradient est non nul, et donc qu'il ne s'agit pas d'un minimum local), on effectue un nouveau pas dans la direction pré-calculée, puis on calcule une nouvelle direction. Il est possible d'obtenir des bornes de convergence qui déterminent le nombre maximum d'itérations nécessaires pour atteindre une précision donnée : lorsque $\zeta = 0$, ce nombre est égal à n d'après le théorème 2.2.2.

De manière remarquable (surtout en comparaison avec les algorithmes que nous verrons par la suite), la méthode du gradient conjugué ne requiert pas d'information liée à la matrice A , comme ses valeurs propres ou sa norme : elle acquiert cette information au fur et à mesure de l'algorithme, via des produits de la matrice A avec des vecteurs (en ce sens, le stockage de la matrice A n'est pas requis). Cette propriété a notamment contribué au succès de l'algorithme du gradient conjugué, très utilisé en calcul scientifique.

2.3 Descente de gradient

Dans cette partie, nous quittons les problèmes sous forme explicite pour nous intéresser à des problèmes fondamentalement **non linéaires** de la forme :

$$\underset{\mathbf{w} \in \mathbb{R}^d}{\text{minimiser}} f(\mathbf{w}). \quad (2.3.1)$$

Le but sera de se ramener à l'étude de fonctions quadratiques construites à partir de f (et de ses dérivées). On travaillera dans la suite sous l'hypothèse suivante.

Hypothèse 2.3.1 La fonction objectif f du problème (2.3.1) est de classe $C_L^{1,1}(\mathbb{R}^d)$ pour $L > 0$, et minorée sur \mathbb{R}^d par $f_{low} \in \mathbb{R}$ (on a donc $f(\mathbf{x}) \geq f_{low} \forall \mathbf{x} \in \mathbb{R}^d$).

Notre but est donc de construire et d'analyser des algorithmes basés sur l'information fournie par le gradient.

L'algorithme de descente de gradient est la méthode la plus classique en optimisation dérivable. Elle se base sur le principe élémentaire suivant, tiré de la condition d'optimalité (1.2.7) : pour tout point $\mathbf{w} \in \mathbb{R}^d$,

1. Soit $\nabla f(\mathbf{w}) = 0$, auquel cas \mathbf{w} est potentiellement un minimum local (c'en est un si f est convexe, et il est même global);
2. Soit $\nabla f(\mathbf{w}) \neq 0$, et on peut alors montrer que la fonction f décroît *localement* dans la direction de l'opposé du gradient $-\nabla f(\mathbf{w})$.

La seconde propriété, que nous démontrerons ci-dessous, constitue l'essence même de l'algorithme de descente de gradient.

2.3.1 Algorithme

L'algorithme de descente de gradient est un processus itératif qui se base sur l'itération suivante :

$$\mathbf{w} \leftarrow \mathbf{w} - \alpha \nabla f(\mathbf{w}), \quad (2.3.2)$$

où $\alpha > 0$ est un paramètre appelé **taille de pas** ou longueur de pas⁴. Lorsque $\nabla f(\mathbf{w}) = 0$, on remarque que le vecteur \mathbf{w} ne change pas lors de la mise à jour : cette propriété est logique puisque dans une telle situation, il est impossible d'utiliser le gradient pour déterminer un meilleur point. En revanche, dès lors que $\nabla f(\mathbf{w}) \neq 0$, on s'attend à ce qu'il existe des valeurs de α pour lesquelles une telle mise à jour permette d'obtenir un meilleur point (avec une valeur de fonction objectif plus faible).

En utilisant la règle (2.3.2) au sein d'un processus itératif, on peut construire un algorithme dont le but consiste à minimiser la fonction objectif f : il s'agit de l'algorithme de **descente de gradient**, parfois également appelé *méthode de la plus forte pente*, dont l'énoncé complet est donné par l'algorithme 3.

Algorithme 3: Descente de gradient pour la minimisation d'une fonction f .

Initialisation: Choisir $\mathbf{w}_0 \in \mathbb{R}^d$.

Pour $k = 0, 1, \dots$

- 1. Calculer le gradient $\nabla f(\mathbf{w}_k)$.
- 2. Définir une longueur de pas $\alpha_k > 0$.
- 3. Poser $\mathbf{w}_{k+1} = \mathbf{w}_k - \alpha_k \nabla f(\mathbf{w}_k)$.

FinPour

Tel qu'écrit ici, l'algorithme ne possède pas de critère d'arrêt; il s'agit en fait d'un schéma générique, dont il existe de nombreuses variantes qui correspondent à des choix spécifiques de critère d'arrêt, de taille de pas, voire de point initial. Nous passons en revue ces différents aspects ci-dessous.

⁴On parle en anglais de *stepsize* ou *steplength*.

Critère d'arrêt Dans la pratique, un algorithme est soumis à des contraintes de budget (en termes d'opérations arithmétiques, de temps d'exécution, d'itérations). Des critères d'arrêt de la méthode sont donc nécessaires afin de forcer l'algorithme à terminer si ces limites sont atteintes. Pour l'algorithme 3, il pourrait s'agir de terminer la méthode après avoir effectué k_{\max} itérations, auquel cas $\mathbf{w}_{k_{\max}}$ représenterait la meilleure solution obtenue.

Par ailleurs, pour mesurer la capacité d'un algorithme à converger vers une solution du problème, on introduit généralement un critère d'arrêt basé sur les conditions d'optimalité. Par exemple, le critère suivant est très fréquemment utilisé pour les algorithmes basés sur le gradient :

$$\|\nabla f(\mathbf{w}_k)\| < \epsilon, \quad (2.3.3)$$

où $\epsilon > 0$ représente une précision donnée, de sorte que ce critère est plus difficile à vérifier lorsque ϵ est très faible.

Enfin, il est toujours possible (et souvent recommandé) d'ajouter des critères d'arrêt de secours, qui permettent de ne pas continuer à faire tourner un algorithme inutilement. Par exemple, lorsque la différence entre deux itérés successifs est du niveau de la précision machine, cela signifie que l'algorithme ne progresse plus, et on peut donc en arrêter l'exécution.

Choix du point initial La performance d'un algorithme peut être grandement améliorée lorsque le point initial est bien choisi. En règle générale, il est cependant difficile de déterminer un tel point : des stratégies de démarrage multiple, qui consistent à effectuer quelques itérations en partant de différents points choisis au hasard, peuvent permettre de déterminer un bon point initial. Mieux encore, dans de nombreuses applications, il existe déjà une valeur de référence ou une idée de ce que pourrait être la solution : dans ce cas, il est souvent avantageux de partir d'un tel point, et de chercher à l'améliorer via le processus d'optimisation.

2.3.2 Choix de la taille de pas

Nous présentons ici les principales techniques de choix de taille de pas; là encore, toute connaissance complémentaire sur le problème peut conduire à de meilleurs choix.

Taille de pas constante L'une des stratégies les plus courantes consiste à utiliser une taille de pas constante pour toutes les itérations de l'algorithme, soit $\alpha_k = \alpha > 0$ pour tout k . Selon le budget disponible, plusieurs valeurs peuvent être testées afin de ne pas se limiter à une seule possibilité. En apprentissage, il est ainsi fréquent de tester une grille de valeurs. Si f vérifie l'hypothèse 2.3.1, on sait qu'il existe des valeurs constantes qui conduiront à un algorithme convergent. En particulier, le choix

$$\alpha_k = \alpha = \frac{1}{L}, \quad (2.3.4)$$

où L représente la constante de Lipschitz pour le gradient, est adapté au problème. Cependant, ce choix nécessite de connaître cette constante de Lipschitz, et une telle information n'est pas forcément aisée à obtenir en pratique.

Taille de pas décroissante Une autre technique classique pour choisir le pas consiste à utiliser une suite de tailles de pas décroissante, de sorte que $\alpha_k \rightarrow 0$ lorsque $k \rightarrow \infty$. Un tel choix peut permettre d'établir des garanties de convergence; en revanche, il peut aussi conduire à un arrêt prématuré de l'algorithme en pratique, si les tailles de pas diminuent trop rapidement. La vitesse avec laquelle la taille de pas α_k tend vers 0 est un aspect critique de ces stratégies.

Choix adaptatif via une recherche linéaire Les techniques de recherche linéaire sont très populaires en optimisation continue (elles sont moins employées en sciences des données, pour des raisons que nous détaillerons plus loin). Pour une itération k donnée, le but d'une recherche linéaire est de calculer la longueur de pas α_k qui conduit à une décroissance de la fonction objectif dans la direction choisie (dans le cas de l'algorithme 3, il s'agit de la direction $-\nabla f(\mathbf{w}_k)$). Une recherche linéaire exacte recherche la taille de pas conduisant à la décroissance maximale, ce qui peut être coûteux à effectuer en pratique. On lui préfèrera plutôt des approches *inexactes*, dont la plus répandue est la technique de **retour arrière** (*backtracking* en anglais) décrite par l'algorithme 3.

Algorithme 4: Recherche linéaire avec retour arrière dans la direction d .

Entrées : $\mathbf{w} \in \mathbb{R}^d$, $\mathbf{d} \in \mathbb{R}^d$, $\alpha_0 \in \mathbb{R}^d$.

Initialisation : Définir $\alpha = \alpha_0$.

Tant que $f(\mathbf{w} + \alpha\mathbf{d}) \geq f(\mathbf{w})$

| $\alpha \rightarrow \frac{\alpha}{2}$.

Fin

Sortie : α .

Une telle recherche linéaire peut être utilisée à l'étape 2 de l'algorithme 3 en appelant l'algorithme 4 avec $\mathbf{w} = \mathbf{w}_k$, $\mathbf{d} = -\nabla f(\mathbf{w}_k)$ et (par exemple) $\alpha_0 = 1$. Dans ce cas, on cherche généralement une longueur de pas vérifiant la condition dite d'Armijo, définie par

$$f(\mathbf{w}_k - \alpha \nabla f(\mathbf{w}_k)) < f(\mathbf{w}_k) - c\alpha \|\nabla f(\mathbf{w}_k)\|^2, \quad (2.3.5)$$

avec $c \in (0, 1/2)$. Il existe de très nombreuses extensions de ce schéma, qui ont l'avantage de garantir l'obtention d'un meilleur point au sens de la valeur de la fonction objectif. Cependant, ces techniques requièrent des évaluations supplémentaires de la fonction objectif, ce qui représente un coût parfois non négligeable.

2.3.3 Analyse de complexité pour la descente de gradient

Nous allons maintenant établir des garanties dites de complexité pour l'algorithme 3. Étant donné un critère de convergence dépendant d'une précision ϵ , il s'agira de borner le nombre d'itérations nécessaires pour satisfaire ce critère comme une fonction de la précision ϵ . Nous obtiendrons des résultats différents selon que nous nous intéresserons aux fonctions non convexes, convexes ou fortement convexes.

Nous nous plaçons dans le cadre de l'hypothèse 2.3.1, qui nous permet de mesurer l'évolution des valeurs de la fonction objectif entre deux pas de descente de gradient.

Proposition 2.3.1 *On considère la k -ième itération de l'algorithme 3 appliquée à une fonction f vérifiant l'hypothèse 2.3.1. Supposons que $\nabla f(\mathbf{w}_k) \neq 0$; alors, si la taille de pas est choisie de sorte que $0 < \alpha_k < \frac{2}{L}$, on a*

$$f(\mathbf{w}_k - \alpha_k \nabla f(\mathbf{w}_k)) < f(\mathbf{w}_k).$$

De plus, si $\alpha_k = \frac{1}{L}$, on obtient :

$$f(\mathbf{w}_k - \frac{1}{L} \nabla f(\mathbf{w}_k)) < f(\mathbf{w}_k) - \frac{1}{2L} \|\nabla f(\mathbf{w}_k)\|^2. \quad (2.3.6)$$

Démonstration. On utilise l'inégalité (A.2.1) appliquée aux vecteurs $(\mathbf{w}_k, \mathbf{w}_k - \alpha_k \nabla f(\mathbf{w}_k))$:

$$\begin{aligned} f(\mathbf{w}_k - \alpha_l \nabla f(\mathbf{w}_k)) &\leq f(\mathbf{w}_k) + \nabla f(\mathbf{w}_k)^T [-\alpha_k \nabla f(\mathbf{w}_k)] + \frac{L}{2} \| -\alpha_k \nabla f(\mathbf{w}_k) \|^2 \\ &= f(\mathbf{w}_k) - \alpha_k \nabla f(\mathbf{w}_k)^T \nabla f(\mathbf{w}_k) + \frac{L}{2} \alpha_k^2 \|\nabla f(\mathbf{w}_k)\|^2 \\ &= f(\mathbf{w}_k) + \left(-\alpha_k + \frac{L}{2} \alpha_k^2 \right) \|\nabla f(\mathbf{w}_k)\|^2. \end{aligned}$$

Si $-\alpha_k + \frac{L}{2} \alpha_k^2 < 0$, le terme $(-\alpha_k + \frac{L}{2} \alpha_k^2) \|\nabla f(\mathbf{w}_k)\|^2$ sera négatif, ce qui impliquera par l'inégalité ci-dessus que $f(\mathbf{w}_k - \alpha_l \nabla f(\mathbf{w}_k)) < f(\mathbf{w}_k)$. Comme $-\alpha_k + \frac{L}{2} \alpha_k^2 < 0 \Leftrightarrow \alpha_k < \frac{2}{L}$ et $\alpha_k > 0$ par définition de l'algorithme, on a donc prouvé la première partie de la proposition.

Pour établir (2.3.6), il suffit de considérer les équations ci-dessus avec la valeur $\alpha_k = \frac{1}{L}$. \square

Le résultat de la proposition 2.3.1 permet de garantir que la fonction objectif peut décroître à pas constant. Nous allons voir comment exploiter ce résultat selon la nature de la fonction objectif.

Cas non convexe En optimisation non convexe, un résultat de complexité correspond à une borne sur le nombre d'itérations requis pour obtenir $\|\nabla f(\mathbf{w}_k)\| \leq \epsilon$. Pour l'algorithme de descente de gradient, on peut ainsi obtenir le résultat suivant.

Théorème 2.3.1 (Complexité de la descente de gradient pour les fonctions non convexes) Soit f une fonction non convexe vérifiant l'hypothèse 2.3.1. On suppose que l'algorithme 3 est appliqué à f avec $\alpha_k = \frac{1}{L}$; alors, pour tout $\epsilon > 0$, l'algorithme atteint un itéré \mathbf{w}_k tel que $\|\nabla f(\mathbf{w}_k)\| \leq \epsilon$ en au plus $\mathcal{O}(\epsilon^{-2})$ itérations.

Démonstration. Soit K un indice tel que pour tout $k = 0, \dots, K-1$, on ait $\|\nabla f(\mathbf{w}_k)\| > \epsilon$. D'après la proposition 2.3.1, on a

$$\forall k = 0, \dots, K-1, \quad f(\mathbf{w}_{k+1}) \leq f(\mathbf{w}_k) - \frac{1}{2L} \|\nabla f(\mathbf{w}_k)\|^2 < f(\mathbf{w}_k) - \frac{1}{2L} \epsilon^2.$$

En sommant cette relation pour toutes les itérations d'indices 0 à $K-1$, on obtient :

$$\sum_{k=0}^{K-1} f(\mathbf{w}_{k+1}) \leq \sum_{k=0}^{K-1} f(\mathbf{w}_k) - \frac{K}{2L} \epsilon^2,$$

ce qui donne après simplification :

$$f(\mathbf{w}_K) \leq f(\mathbf{w}_0) - \frac{K}{2L} \epsilon^2.$$

Comme $f(\mathbf{w}_K) \geq f_{low}$ d'après l'hypothèse 2.3.1, on aboutit à

$$K \leq 2L(f(\mathbf{w}_0) - f_{low})\epsilon^{-2}.$$

Par conséquent, on vient de montrer que le nombre d'itérations pour lesquelles $\|\nabla f(\mathbf{w}_k)\| > \epsilon$ est majoré par

$$\lceil 2L(f(\mathbf{w}_0) - f_{low})\epsilon^{-2} \rceil = \mathcal{O}(\epsilon^{-2}).$$

\square

Remarque 2.3.1 La notion de borne de complexité est parfois remplacée par celle, équivalente, de vitesse de convergence, qui ne fait pas directement intervenir de précision. Dans le cas présent, par exemple, on dira que l'algorithme de descente de gradient converge en $\mathcal{O}\left(\frac{1}{\sqrt{K}}\right)$, dans le sens où on peut établir (via un raisonnement similaire à celui utilisé pour prouver le théorème 2.3.1) que

$$\min_{0 \leq k \leq K-1} \|\nabla f(\mathbf{w}_k)\| \leq \frac{C}{\sqrt{K}},$$

où la constante C dépend de $f(\mathbf{w}_0)$, f_{low} et L .

On notera que le résultat du théorème 2.3.1 ne garantit pas que la méthode converge vers un minimum local, seulement vers un point stationnaire d'ordre 1 (qui peut donc être un point selle ou un maximum local). En pratique, s'il est facile de trouver des exemples "pathologiques" où l'algorithme de descente de gradient converge vers un point selle ou reste bloqué en un maximum local, on observe plutôt que la méthode converge vers de "bons points". Un résultat récent a permis de préciser ce comportement en présence d'une dérivée seconde.

Théorème 2.3.2 ([?]) Sous les hypothèses du théorème 2.3.1, on suppose que f est de classe \mathcal{C}^2 et que le point initial \mathbf{w}_0 est choisi aléatoirement dans \mathbb{R}^d . Alors, la méthode de descente de gradient converge presque sûrement vers un point stationnaire d'ordre deux.

Pour des fonctions e théorème 2.3.1 garantit donc (avec probabilité 1) que la descente de gradient ne converge pas vers un point en lequel la matrice hessienne n'est pas semi-définie positive.

2.3.4 Cas convexe et fortement convexe

On suppose maintenant (en plus de l'hypothèse 2.3.1) que la fonction objectif est convexe : pour le même algorithme que ci-dessus, on peut alors démontrer que des propriétés plus fortes sont vérifiées à un coût plus faible. Ces résultats illustrent l'intérêt des fonctions convexes en optimisation.

Pour le reste de cette section, on considérera l'hypothèse suivante.

Hypothèse 2.3.2 La fonction f est convexe, et admet un minimum $\mathbf{w}^* \in \mathbb{R}^d$. On notera la valeur minimale par $f^* := f(\mathbf{w}^*) = \min_{\mathbf{w} \in \mathbb{R}^d} f(\mathbf{w})$.

Pour une précision $\epsilon > 0$, on va maintenant s'intéresser au nombre d'itérations nécessaires pour atteindre \mathbf{w}_k tel que $f(\mathbf{w}_k) - f^* \leq \epsilon$. En pratique, f^* n'est pas toujours connue, mais un tel critère a son importance (et peut toujours être estimé en évaluant la valeur de l'objectif).

Théorème 2.3.3 (Complexité de la descente de gradient pour les problèmes convexes) On considère une fonction f vérifiant les hypothèses 2.3.1 et 2.3.2. On considère l'algorithme 3 appliqué à f avec $\alpha_k = \frac{1}{L}$; alors, pour tout $\epsilon > 0$, l'algorithme obtient \mathbf{w}_k tel que $f(\mathbf{w}_k) - f^* \leq \epsilon$ en au plus $\mathcal{O}(\epsilon^{-1})$ itérations.

Démonstration. Soit K un indice tel que pour tout $k = 0, \dots, K-1$, on ait $f(\mathbf{w}_k) - f^* > \epsilon$.

Pour tout $k = 0, \dots, K-1$, la caractérisation de la convexité via le gradient (1.2.12) appliquée à \mathbf{w}_k et \mathbf{w}^* donne

$$f(\mathbf{w}^*) \geq f(\mathbf{w}_k) + \nabla f(\mathbf{w}_k)^T (\mathbf{w}^* - \mathbf{w}_k).$$

En combinant cette propriété avec (2.3.6), il vient :

$$\begin{aligned} f(\mathbf{w}_{k+1}) &\leq f(\mathbf{w}_k) - \frac{1}{2L} \|\nabla f(\mathbf{w}_k)\|^2 \\ &\leq f(\mathbf{w}^*) + \nabla f(\mathbf{w}_k)^T (\mathbf{w}_k - \mathbf{w}^*) - \frac{1}{2L} \|\nabla f(\mathbf{w}_k)\|^2. \end{aligned}$$

On remarque ensuite que

$$\nabla f(\mathbf{w}_k)^T (\mathbf{w}_k - \mathbf{w}^*) - \frac{1}{2L} \|\nabla f(\mathbf{w}_k)\|^2 = \frac{L}{2} \left(\|\mathbf{w}_k - \mathbf{w}^*\|^2 - \|\mathbf{w}_k - \mathbf{w}^* - \frac{1}{L} \nabla f(\mathbf{w}_k)\|^2 \right).$$

Par conséquent, si l'on se rappelle que $\mathbf{w}_{k+1} = \mathbf{w}_k - \frac{1}{L} \nabla f(\mathbf{w}_k)$, on aboutit à :

$$\begin{aligned} f(\mathbf{w}_{k+1}) &\leq f(\mathbf{w}^*) + \frac{L}{2} \left(\|\mathbf{w}_k - \mathbf{w}^*\|^2 - \|\mathbf{w}_k - \mathbf{w}^* - \frac{1}{L} \nabla f(\mathbf{w}_k)\|^2 \right) \\ &= f(\mathbf{w}^*) + \frac{L}{2} \left(\|\mathbf{w}_k - \mathbf{w}^*\|^2 - \|\mathbf{w}_{k+1} - \mathbf{w}^*\|^2 \right), \end{aligned}$$

soit encore

$$f(\mathbf{w}_{k+1}) - f(\mathbf{w}^*) \leq \frac{L}{2} \left(\|\mathbf{w}_k - \mathbf{w}^*\|^2 - \|\mathbf{w}_{k+1} - \mathbf{w}^*\|^2 \right). \quad (2.3.7)$$

En sommant les inégalités (2.3.7) pour les indices k entre 0 et $K - 1$, on aboutit à

$$\sum_{k=0}^{K-1} f(\mathbf{w}_{k+1}) - f(\mathbf{w}^*) \leq \frac{L}{2} \left(\|\mathbf{w}_0 - \mathbf{w}^*\|^2 - \|\mathbf{w}_K - \mathbf{w}^*\|^2 \right) \leq \frac{L}{2} \|\mathbf{w}_0 - \mathbf{w}^*\|^2.$$

Enfin, en utilisant notre hypothèse $f(\mathbf{w}_k) - f(\mathbf{w}^*) > \epsilon$, le membre de gauche de l'équation peut être minoré, et on obtient ainsi

$$K\epsilon \leq \frac{L}{2} \|\mathbf{w}_0 - \mathbf{w}^*\|^2.$$

On a donc prouvé que le nombre d'itérations pour lesquelles $f(\mathbf{w}_k) - f(\mathbf{w}^*) > \epsilon$ est majoré par

$$\left\lceil \frac{L \|\mathbf{w}_0 - \mathbf{w}^*\|^2}{2\epsilon} \right\rceil = \mathcal{O}(\epsilon^{-1}).$$

□

Remarque 2.3.2 Par le même raisonnement, on peut montrer que la vitesse de convergence de la descente de gradient dans le cas convexe est $\mathcal{O}(\frac{1}{K})$, c'est-à-dire qu'il existe $C > 0$ (dont la valeur dépend de $\|\mathbf{w}_0 - \mathbf{w}^*\|$ et L) telle que

$$f(\mathbf{w}_K) - f^* \leq \frac{C}{K}.$$

Nous nous tournons finalement vers le cas fortement convexe.

Théorème 2.3.4 (Complexité de la descente de gradient pour les fonctions fortement convexes)

Soit f une fonction vérifiant les hypothèses 2.3.1 et 2.3.2, dont on suppose de surcroît qu'elle est μ -fortement convexe de paramètre $\mu \in (0, L]$. On considère l'algorithme 3 appliqué à f avec $\alpha_k = \frac{1}{L}$; alors, pour tout $\epsilon > 0$, la méthode obtient \mathbf{w}_k tel que $f(\mathbf{w}_k) - f^* \leq \epsilon$ en au plus $\mathcal{O}(\frac{L}{\mu} \ln(\frac{1}{\epsilon}))$ itérations.

Démonstration. On utilise la propriété de convexité forte (1.2.14), que l'on rappelle ci-dessous : pour tous $(\mathbf{u}, \mathbf{v}) \in (\mathbb{R}^n)^2$, on a :

$$f(\mathbf{v}) \geq f(\mathbf{u}) + \nabla f(\mathbf{u})^T(\mathbf{v} - \mathbf{u}) + \frac{\mu}{2}\|\mathbf{v} - \mathbf{u}\|^2.$$

En minimisant des deux côtés de l'inégalité par rapport à \mathbf{v} , on obtient $\mathbf{v} = \mathbf{w}^*$ dans le membre de gauche, et $\mathbf{v} = \mathbf{u} - \frac{1}{\mu}\nabla f(\mathbf{u})$ dans celui de droite.⁵ On obtient ainsi

$$\begin{aligned} f^* &\geq f(\mathbf{u}) + \nabla f(\mathbf{u})^T \left[-\frac{1}{\mu} \nabla f(\mathbf{u}) \right] + \frac{\mu}{2} \left\| -\frac{1}{\mu} \nabla f(\mathbf{u}) \right\|^2 \\ f^* &\geq f(\mathbf{u}) - \frac{1}{2\mu} \|\nabla f(\mathbf{u})\|^2. \end{aligned}$$

Un ré-arrangement des termes conduit à

$$\|\nabla f(\mathbf{u})\|^2 \geq 2\mu [f(\mathbf{u}) - f^*], \quad (2.3.8)$$

qui est donc valable pour tout vecteur $\mathbf{u} \in \mathbb{R}^n$. En combinant (2.3.8) (appliqué avec $\mathbf{u} = \mathbf{w}_k$) avec (2.3.6), il vient

$$f(\mathbf{w}_{k+1}) \leq f(\mathbf{w}_k) - \frac{1}{2L} \|\nabla f(\mathbf{w}_k)\|^2 \leq f(\mathbf{w}_k) - \frac{\mu}{L} (f(\mathbf{w}_k) - f^*).$$

On a ainsi,

$$f(\mathbf{w}_{k+1}) - f^* \leq \left(1 - \frac{\mu}{L}\right) (f(\mathbf{w}_k) - f^*),$$

d'où l'on tire par récursion que

$$f(\mathbf{w}_K) - f^* \leq \left(1 - \frac{\mu}{L}\right)^K (f(\mathbf{w}_0) - f^*)$$

pour tout $K \geq 0$ (on notera que l'on obtient ainsi une vitesse de convergence). Si l'on suppose que $f(\mathbf{w}_K) - f^* \geq \epsilon$, cela conduit à :

$$\begin{aligned} \epsilon &\leq \left(1 - \frac{\mu}{L}\right)^K (f(\mathbf{w}_0) - f^*) \\ \ln(\epsilon) &\leq K \ln \left(1 - \frac{\mu}{L}\right) + \ln(f(\mathbf{w}_0) - f^*) \\ K &\leq \frac{\ln(\epsilon(f(\mathbf{w}_0) - f^*)^{-1})}{\ln(1 - \frac{\mu}{L})} \\ K &\leq \frac{\ln((f(\mathbf{w}_0) - f^*) \epsilon^{-1})}{\ln(L/(L - \mu))}. \end{aligned}$$

En appliquant

$$\ln \left(\frac{L}{L-\mu} \right) = \ln \left(1 + \frac{\mu}{L-\mu} \right) \geq \frac{1}{\frac{L-\mu}{\mu} + 1/2} = \frac{1}{L/\mu - 1/2} \geq \frac{\mu}{L},$$

⁵Pour un raisonnement complet, voir l'exercice 2 du TD 2.

on arrive ainsi à

$$K \leq \frac{L}{\mu} \ln((f(\mathbf{w}_0) - f^*) \epsilon^{-1}),$$

ce qui montre alors que le nombre d'itérations pour lesquelles $f(\mathbf{w}_k) - f(\mathbf{w}^*) > \epsilon$ doit être majoré par

$$\left\lceil \frac{L}{\mu} \ln((f(\mathbf{w}_0) - f^*) \epsilon^{-1}) \right\rceil = \mathcal{O}\left(\frac{L}{\mu} \ln\left(\frac{1}{\epsilon}\right)\right).$$

□

- Remarque 2.3.3**
- Par une démonstration similaire, on peut montrer (et on dira donc) que la vitesse de convergence de la descente de gradient dans le cas fortement convexe est en $\mathcal{O}\left((1 - \frac{\mu}{L})^k\right)$;
 - Dans le cas fortement convexe, on peut également établir ces résultats pour le critère de convergence des itérés $\|\mathbf{w}_k - \mathbf{w}^*\| \leq \epsilon$: la distance entre l'itéré courant et l'optimum (qui est unique) décroît donc à un taux $\mathcal{O}\left((1 - \frac{\mu}{L})^k\right)$.

Pour terminer, on notera que les preuves de vitesses de convergence dans le cas convexe (et fortement convexe) sont souvent plus techniques que celles du cas non convexe. Cette technicité est nécessaire pour certifier que les vitesses de convergence de ces algorithmes sont meilleures. Dans ce cours, nous nous intéressons principalement aux résultats, et à ce qu'ils impliquent sur la facilité de résolution de ces problèmes.

2.4 Accélération

2.4.1 Introduction au concept de momentum

Dans la partie précédente, nous avons obtenu des bornes de complexité dites *au pire cas* : celles-ci quantifient la performance d'un algorithme sur une classe de problèmes donnée, mais elles n'indiquent pas la meilleure performance possible pour *une classe d'algorithmes donnée*. On parle ainsi de bornes supérieures, par opposition aux bornes inférieures dont nous allons discuter ici.

Pour l'optimisation d'une fonction $\mathcal{C}_L^{1,1}$ non convexe par une méthode de gradient, on sait que la borne en $\mathcal{O}(\epsilon^{-2})$, que nous avons obtenue dans le théorème 2.3.1 ne peut pas être améliorée sans que l'algorithme utilise plus d'information. En revanche, pour le cas convexe, il existe des algorithmes qui peuvent atteindre une borne de complexité $\mathcal{O}(\epsilon^{-1/2})$, ce qui représente une amélioration notable par rapport à la borne $\mathcal{O}(\epsilon^{-1})$ que nous avons obtenu pour la descente de gradient (cf théorème 2.3.3). Ces méthodes ayant une meilleure complexité reposent de manière plus ou moins explicite sur une technique dite d'**accélération**.

Le principe général de l'accélération (que l'on appelle parfois *momentum*) est d'utiliser de l'information des itérations précédentes (généralement la dernière), ce qui permet potentiellement de profiter du momentum de l'itération précédente pour effectuer un meilleur pas.

2.4.2 Méthode du gradient accéléré

L'utilisation la plus connue des techniques d'accélération est due au mathématicien Yurii Nesterov [4] : l'algorithme associé est d'ailleurs souvent appelé "méthode de Nesterov". Une description générique de cette méthode est donnée par l'algorithme 5.

Algorithme 5: Méthode du gradient accéléré

Initialisation : $\mathbf{w}_0 \in \mathbb{R}^d$, $\mathbf{w}_{-1} = \mathbf{w}_0$.

Pour $k = 0, 1, \dots$

1. Calculer une taille de pas $\alpha_k > 0$ et un paramètre $\beta_k > 0$.
2. Définir le nouveau point comme :

$$\mathbf{w}_{k+1} = \mathbf{w}_k - \alpha_k \nabla f(\mathbf{w}_k + \beta_k(\mathbf{w}_k - \mathbf{w}_{k-1})) + \beta_k(\mathbf{w}_k - \mathbf{w}_{k-1}). \quad (2.4.1)$$

Fin

Comme l'algorithme de descente de gradient vu en section 2.3, l'algorithme du gradient accéléré requiert une seule évaluation de gradient à chaque itération; en revanche, contrairement à l'algorithme de descente de gradient, dans l'algorithme 5, le gradient n'est pas évalué en le point courant \mathbf{w}_k , mais en une combinaison linéaire de ce point avec le pas précédent $\mathbf{w}_k - \mathbf{w}_{k-1}$: ce second terme, appelé **terme de momentum**, est à l'origine de la performance améliorée des techniques accélérées.

L'algorithme 5 peut s'écrire d'une autre manière en utilisant deux suites de vecteurs démarrant respectivement à \mathbf{w}_0 et $\mathbf{z}_0 = \mathbf{w}_0$; on réécrit alors la mise à jour (2.4.1) comme suit :

$$\begin{cases} \mathbf{w}_{k+1} = \mathbf{z}_k - \alpha_k \nabla f(\mathbf{z}_k) \\ \mathbf{z}_{k+1} = \mathbf{w}_{k+1} + \beta_{k+1}(\mathbf{w}_{k+1} - \mathbf{w}_k). \end{cases} \quad (2.4.2)$$

Grâce à cette nouvelle formulation, on distingue bien le processus en deux étapes contenu dans la méthode du gradient accéléré : un pas de gradient sur \mathbf{z}_k et un pas de type "momentum" sur \mathbf{w}_{k+1} .

Choix des paramètres En plus du choix de la *taille de pas*, déjà présent dans l'algorithme de descente de gradient, l'algorithme 5 possède un autre paramètre β_k , dit de momentum. Le paramètre α_k peut être choisi selon les mêmes techniques que celles décrites en section 2.3.2 : le choix classique est celui d'une taille de pas constante $\alpha_k = \frac{1}{L}$, qui permet d'obtenir les résultats de complexité attendus.

Le choix de β_k est également crucial dans l'obtention de ces garanties. Ainsi, les valeurs de β_k proposées originellement par Nesterov dépendent de la fonction objectif :

- Si f est μ -fortement convexe, on pose :

$$\beta_k = \beta = \frac{\sqrt{L} - \sqrt{\mu}}{\sqrt{L} + \sqrt{\mu}} \quad (2.4.3)$$

pour tout k . Cela requiert de connaître à la fois la constante de Lipschitz pour le gradient et le paramètre de convexité forte.

- Pour une fonction f convexe quelconque, β_k est calculé de manière adaptative en utilisant deux suites :

$$t_{k+1} = \frac{1}{2}(1 + \sqrt{1 + 4t_k^2}), t_0 = 0, \quad \beta_k = \frac{t_k - 1}{t_{k+1}}. \quad (2.4.4)$$

Le résultat suivant résume les résultats de complexité qui peuvent être obtenus pour l'algorithme 5.

Théorème 2.4.1 Soit l'algorithme 5 appliqué à une fonction f vérifiant les hypothèses 2.3.1 et 2.3.2, avec $\alpha_k = \frac{1}{L}$ et β_k choisi selon la règle (2.4.4). Alors, pour tout $\epsilon > 0$, l'algorithme atteint \mathbf{w}_k tel que $f(\mathbf{w}_k) - f^* \leq \epsilon$ en au plus

- i) $\mathcal{O}(\epsilon^{-1/2})$ itérations pour une fonction convexe;
- ii) $\mathcal{O}\left(\sqrt{\frac{L}{\mu}} \ln\left(\frac{1}{\epsilon}\right)\right)$ itérations pour une fonction μ -fortement convexe.

On peut également obtenir des vitesses de convergence pour le gradient accéléré, qui reflètent également cette amélioration. Ainsi, pour des fonctions μ -fortement convexes, on peut démontrer que $f(\mathbf{w}_k) - f^* = \mathcal{O}\left(\left(1 - \sqrt{\frac{\mu}{L}}\right)^k\right)$, ce qui est donc plus rapide que la vitesse en $\mathcal{O}\left((1 - \frac{\mu}{L})^k\right)$ de la descente de gradient.

2.4.3 Autres algorithmes accélérés

Boule lestée La méthode de la boule lestée (heavy ball) est un précurseur de l'algorithme du gradient accéléré, développée par Boris T. Polyak in 1964, dédiée à la minimisation des fonctions fortement convexes. Sa k -ième itération s'écrit :

$$\mathbf{w}_{k+1} = \mathbf{w}_k - \alpha \nabla f(\mathbf{w}_k) + \beta(\mathbf{w}_k - \mathbf{w}_{k+1}),$$

où la taille de pas et le paramètre de momentum sont typiquement choisis constants en fonction des constantes de Lipschitz et paramètre de convexité forte du problème. La différence entre ce schéma et celui de Nesterov repose sur le point en lequel le gradient est évalué : en ce sens, la méthode de la boule lestée effectue d'abord un pas de gradient, puis un pas de momentum, tandis que l'algorithme du gradient accéléré adopte l'approche inverse. Il est possible de montrer que l'algorithme de la boule lestée possède la même complexité que le gradient accéléré sur des fonctions quadratiques strictement convexes, en revanche le résultat est faux sur des fonctions strictement convexes générales.

Gradient conjugué Nous avons déjà vu l'algorithme du gradient conjugué en section 2.2.3. Cette méthode est dédiée à la minimisation de fonctions quadratiques; contrairement à la méthode de la boule lestée, le gradient conjugué ne requiert pas de choix des paramètres L ou μ , car il exploite plus d'information des itérations précédentes. On peut résumer la k -ième itération en la double récursion suivante :

$$\mathbf{w}_{k+1} = \mathbf{w}_k + \alpha_k p_k, \quad p_k = -\nabla f(x_k) + \beta_k p_{k-1},$$

où α_k and β_k sont calculés automatiquement par l'algorithme. On reconnaît ainsi un pas dans une direction mêlant le gradient et la direction prise à l'itération précédente. Il est possible de démontrer une vitesse de convergence similaire à celle du gradient accéléré sur ce problème; en revanche, l'algorithme du gradient conjugué est garanti de terminer en au plus d itérations en dimension d . Lorsque d est très grand, la borne de complexité correspond à celle des autres algorithmes.

Exemple 2.4.1 (Minimisation d'une fonction quadratique fortement convexe) On revient sur les problèmes de minimisation quadratique de la forme

$$\underset{\mathbf{w} \in \mathbb{R}^d}{\text{minimiser}} q(\mathbf{w}) := \frac{1}{2} \mathbf{w}^T \mathbf{A} \mathbf{w} + \mathbf{b}^T \mathbf{w}$$

où $A \in \mathbb{R}^{d \times d}$ et $b \in \mathbb{R}^d$. On suppose ici que A est une matrice symétrique définie positive, de sorte que le problème est fortement convexe. Ce problème est dérivable (car la fonction objectif est polynomiale en chacune des variables) et μ -fortement convexe, où μ représente la plus petite valeur propre de A . Ce problème admet un unique minimum global donné par la solution de l'équation $\nabla q(\mathbf{w}) = Aw + b = 0$. Comme nous l'avons vu en section 2.2, on peut calculer une solution de ce système via l'inverse; cependant, le coût d'une telle opération peut être trop élevé en grande dimension. On remplace alors le calcul exact de la solution par un appel à une méthode basée sur le gradient, par exemple l'algorithme 3 ou 5. Comme $q \in \mathcal{C}_{\|A\|}^{1,1}(\mathbb{R}^d)$, le choix de taille de pas 2.3.4 est approprié.

Si la descente de gradient est utilisée, on sait que l'on peut obtenir $q(\mathbf{w}_k) - q^* \leq \epsilon$ (avec $q^* = \min_{\mathbf{w} \in \mathbb{R}^d} q(\mathbf{w})$) en au plus $\mathcal{O}\left(\frac{L}{\mu} \ln\left(\frac{1}{\epsilon}\right)\right)$ itérations, tandis qu'appliquer le gradient accéléré ou la méthode de la boule lestée conduirait à une borne en $\mathcal{O}\left(\frac{L}{\mu} \ln\left(\frac{1}{\epsilon}\right)\right)$. Enfin, avec le gradient conjugué, on aurait une borne de complexité en $\mathcal{O}\left(\min\{d, \frac{L}{\mu} \ln\left(\frac{1}{\epsilon}\right)\}\right)$.

Pour terminer cette partie, nous mentionnons brièvement que les méthodes de quasi-Newton peuvent être vues comme des techniques accélérées; nous reviendrons sur ces techniques dans un chapitre dédié.

2.5 Conclusion du chapitre 2

Si l'on considère un problème d'optimisation générique, il n'est en général pas possible d'obtenir une formule explicite pour ses solutions : c'est pourtant le cas pour certaines classes de problèmes telles que les moindres carrés linéaires, ou plus généralement les problèmes quadratiques. Cependant, même dans cette situation, le calcul direct de la solution peut se révéler coûteux, et être remplacé par une procédure itérative qui converge vers une solution.

L'algorithme de descente de gradient est l'exemple classique d'une méthode itérative pour l'optimisation différentiable, dont de nombreuses variantes ont été proposées. Ces variantes diffèrent notamment dans leur stratégie de choix de longueur de pas. Pour analyser le comportement de la descente de gradient, on se base sur des vitesses de convergence (ou des bornes de complexité) qui sont spécifiques à une classe de fonctions donnée. Ainsi, la vitesse de convergence de la descente de gradient est meilleure sur un problème fortement convexe que sur un problème convexe, et meilleure sur un problème convexe que sur un problème non convexe.

On peut alors se demander quelles sont les bornes de complexité optimales que l'on peut obtenir avec des méthodes basées sur le gradient. Dans le cas non convexe, il n'existe pas de méthode basée sur le gradient avec une meilleure complexité. En revanche, lorsque la fonction objectif est convexe ou fortement convexe, on peut construire des algorithmes dits "accélérés", pour lesquels les meilleures vitesses de convergence possibles sont atteintes. Ces méthodes se basent essentiellement sur le concept de momentum, qui consiste à combiner l'information du gradient avec le déplacement effectué à l'itération précédente : ce paradigme est à la base de nombreux algorithmes d'optimisation très efficaces en pratique, même sur des problèmes non convexes.

Chapitre 3

Méthodes de gradient stochastique

3.1 Motivation

Dans ce chapitre, nous allons véritablement prendre en compte la structure des problèmes liés aux sciences des données. On suppose ainsi que l'on dispose d'un échantillon de n exemples sous la forme $\{(\mathbf{x}_i, y_i)\}_{i=1}^n$ où $\mathbf{x}_i \in \mathbb{R}^d$ et $y_i \in \mathbb{R}$ sont obtenus à partir d'une distribution de données. Comme dans l'exemple de la régression linéaire, on recherche un modèle h tel que $h(\mathbf{x}_i) \approx y_i$ pour tout $i = 1, \dots, n$. On supposera ici qu'un modèle est paramétré par un vecteur $\mathbf{w} \in \mathbb{R}^d$ (c'est-à-dire $h(\mathbf{x}_i) = h(\mathbf{w}; \mathbf{x}_i)$), de sorte qu'il suffit de déterminer le vecteur \mathbf{w} pour définir le modèle. Afin de quantifier la capacité du modèle à représenter nos données, on définit une fonction de coût (ou de perte) de la forme $\ell : (h, y) \mapsto \ell(h, y)$. Le but de cette fonction est de pénaliser des valeurs (h, y) telles que $h \neq y$. La fonction $(h, y) \mapsto (h - y)^2$, que nous avons déjà rencontrée dans le contexte des moindres carrés, est un exemple d'une fonction de coût pour les modèles linéaires. Avec cette fonction de coût, la perte en un échantillon donné est $\ell(h(\mathbf{w}; \mathbf{x}_i), y_i)$: on souhaite que notre modèle soit le meilleur relativement à l'ensemble des exemples. On considère ainsi la moyenne des pertes, ce qui conduit au problème d'optimisation suivant.

Définition 3.1.1 (Problème d'optimisation en somme finie) Soit un jeu de données $\{(\mathbf{x}_i, y_i)\}_{i=1}^n$ où $\mathbf{x}_i \in \mathbb{R}^d$ et $y_i \in \mathbb{R}$, une classe de modèles $\{h(\mathbf{w}; \cdot)\}_{\mathbf{w} \in \mathbb{R}^d}$ et une fonction de coût ℓ . On définit le problème d'optimisation en somme finie suivant :

$$\underset{\mathbf{w} \in \mathbb{R}^d}{\text{minimiser}} f(\mathbf{w}) = \frac{1}{n} \sum_{i=1}^n \ell(h(\mathbf{w}; \mathbf{x}_i), y_i) = \frac{1}{n} \sum_{i=1}^n f_i(\mathbf{w}). \quad (3.1.1)$$

Supposons que l'on applique l'algorithme 3 de descente de gradient à ce problème, en supposant que tous les f_i sont de classe \mathcal{C}^1 (donc que f l'est). La k -ième itération de cet algorithme s'écrira

$$\mathbf{w}_{k+1} = \mathbf{w}_k - \alpha_k \nabla f(\mathbf{w}_k) = \mathbf{w}_k - \frac{\alpha_k}{n} \sum_{i=1}^n \nabla f_i(\mathbf{w}_k). \quad (3.1.2)$$

On voit ainsi que chaque itération de descente de gradient requiert l'accès à **l'ensemble des données** pour effectuer un calcul de gradient. Dans un contexte de données massives, le nombre d'exemples n peut être extrêmement large, et cet algorithme peut être trop coûteux pour être utilisé en pratique.

Remarque 3.1.1 Dans un contexte stochastique ou “online”, les exemples peuvent être directement générés de la distribution à la volée. Dans ce contexte, il peut être impossible d’effectuer une moyenne discrète sur les échantillons, et donc d’obtenir un problème sous la forme d’une somme finie. On peut néanmoins formuler le problème en utilisant une espérance mathématique, et ainsi chercher à résoudre :

$$\min_{\mathbf{w} \in \mathbb{R}^d} \mathbb{E}_{(\mathbf{x},y)} [f_{(\mathbf{x},y)}(\mathbf{w})]. \quad (3.1.3)$$

Ce problème est bien souvent le véritable but de l’apprentissage automatique, mais est souvent remplacé par (3.1.1) pour tenir compte de l’échantillon fini de données à disposition. Par ailleurs, le gradient de cette fonction objectif peut être compliqué voire impossible à calculer, ce qui exclut d’emblée l’utilisation d’algorithmes tels que la descente de gradient. La méthode du gradient stochastique sera au contraire applicable dans ce cas.

3.2 Méthode du gradient stochastique

3.2.1 Algorithme

L’idée derrière l’algorithme du gradient stochastique est remarquablement simple. Si l’on considère le problème minimiser $\min_{\mathbf{w} \in \mathbb{R}^d} \frac{1}{n} \sum_{i=1}^n f_i(\mathbf{w})$ sous l’hypothèse que chaque fonction f_i est dérivable, chaque itération consiste à choisir un indice i_k au hasard et à faire un pas dans la direction opposée au gradient de la fonction f_{i_k} . L’algorithme 6 décrit ce processus.

Algorithme 6: Méthode du gradient stochastique.

Initialisation: $\mathbf{w}_0 \in \mathbb{R}^d$.

for $k = 0, 1, \dots$ **do**

- 1. Calculer un pas $\alpha_k > 0$.
- 2. Tirer un indice $i_k \in \{1, \dots, n\}$.
- 3. Calculer le nouvel itéré :

$$\mathbf{w}_{k+1} = \mathbf{w}_k - \alpha_k \nabla f_{i_k}(\mathbf{w}_k). \quad (3.2.1)$$

end

Le vecteur $\nabla f_{i_k}(\mathbf{w}_k)$ s’appelle un **gradient stochastique** au point \mathbf{w}_k . La propriété principale de l’itération (3.2.1) est qu’elle n’accède qu’à un seul exemple du jeu de données : par conséquent, son coût en termes d’accès aux données est **n fois inférieur à celui d’une itération de descente de gradient** (3.1.2).

Remarque 3.2.1 En règle générale, et même si le problème d’optimisation s’écrit sous la forme d’une somme finie, l’algorithme 6 peut diverger. À titre d’exemple, considérons le problème

$$\underset{w \in \mathbb{R}}{\text{minimiser}} \frac{1}{2} (f_1(w) + f_2(w))$$

avec $f_1(w) = 2w^2$ et $f_2 = -w^2$. Si $w_0 > 0$ et $i_k = 2$ pour tout k , alors l’algorithme divergera.

Il est ais  de construire des exemples pour lesquels la m thode du gradient stochastique diverge ainsi. En pratique, cependant, pour les probl mes de somme finie issus de l'apprentissage automatique, les donn es sont suffisamment **corr l es** pour qu'une mise   jour relativement   un exemple affecte la pr dition sur d'autres exemples. Cette observation explique en partie le succ s des m thodes de gradient stochastique dans ce contexte.

Remarque 3.2.2 L'algorithme 6 est souvent d sign  par l'acronyme **SGD**, pour **Stochastic Gradient Descent**, ou descente de gradient stochastique. Cela  tant, il n'est pas possible de garantir que l'algorithme du gradient stochastique est une m thode de descente (qui d croit la fonction objectif   chaque it ration). Pour cette raison, dans ce document, nous utiliserons la terminologie (par ailleurs adopt e chez d'autres auteurs) de **gradient stochastique** (mais pourrons nous autoriser l'utilisation de l'acronyme **SGD**, qui se retrouve dans de nombreuses impl mentations).

3.2.2 Analyse

Dans cette section, on d crit les  tapes principales de l'obtention de vitesses de convergence pour l'algorithme du gradient stochastique, sous une version l g rement modifi e de l'hypoth se 2.3.1.

Hypoth se 3.2.1 La fonction objectif $f = \frac{1}{n} \sum_{i=1}^n f_i$ du probl me (3.1.1) est de classe $\mathcal{C}_L^{1,1}(\mathbb{R}^d)$ pour $L > 0$, et minor e sur \mathbb{R}^d par $f_{low} \in \mathbb{R}$. Par ailleurs, chaque fonction f_i est de classe \mathcal{C}^1 .

L'argument principal de l'analyse de la descente de gradient est le r sultat de la proposition 2.3.1, que l'on rappelle ci-dessous :

$$f(\mathbf{w}_{k+1}) \leq f(\mathbf{w}_k) + \nabla f(\mathbf{w}_k)^T (\mathbf{w}_{k+1} - \mathbf{w}_k) + \frac{L}{2} \|\mathbf{w}_{k+1} - \mathbf{w}_k\|^2.$$

Il est possible d'en obtenir une version  quivalente pour l'algorithme du gradient stochastique, sous des hypoth ses sur le tirage al atoire des indices, et donc des gradients "stochastiques" associ s. Ces conditions sont r sum es dans l'hypoth se suivante.

Hypoth se 3.2.2 (Hypoth se sur les gradients stochastiques)   chaque it ration de l'algorithme 6 d'indice k , l'indice i_k est tir  tel que :

i) $\mathbb{E}_{i_k} [\nabla f_{i_k}(\mathbf{w}_k)] = \nabla f(\mathbf{w}_k); \mathbb{E}_{i_k} [\|\nabla f_{i_k}(\mathbf{w}_k)\|^2] \leq \sigma^2 + \|\nabla f(\mathbf{w}_k)\|^2$ avec $\sigma^2 > 0$.

La propri t  i) de l'hypoth se 3.2.2 impose que le gradient stochastique $\nabla f_{i_k}(\mathbf{w}_k)$ soit un estimateur sans biais du v ritable gradient $\nabla f(\mathbf{w}_k)$. La seconde propri t  controle la variance de la norme de ce gradient stochastique, de sorte   garantir que la variation al atoire de celui-ci soit contrl e e. Il existe plusieurs mani res de g n rer al atoirement des indices v rifiant ces propri t s, au premier rang desquelles le tirage uniforme.

Exemple 3.2.1 ( chantillonnage uniforme) Si   l'it ration k , l'indice i_k est tir  uniform m nt dans $\{1, \dots, n\}$, alors l'algorithme 6 v rifie l'hypoth se 3.2.2.

Proposition 3.2.1 Sous les hypoth ses 2.3.1 et 3.2.2, on consid re la k -i me it ration de l'algorithme 6. Alors, on a

$$\mathbb{E}_{i_k} [f(\mathbf{w}_{k+1})] - f(\mathbf{w}_k) \leq \nabla f(\mathbf{w}_k)^T \mathbb{E}_{i_k} [\mathbf{w}_{k+1} - \mathbf{w}_k] + \frac{L}{2} \mathbb{E}_{i_k} [\|\mathbf{w}_{k+1} - \mathbf{w}_k\|^2]. \quad (3.2.2)$$

Comme le montre l'inégalité ci-dessus, on ne peut garantir la décroissance d'une itération à l'autre qu'**en moyenne** (sous réserve que le membre de droite de (3.2.2) soit négatif). Cependant, une telle propriété suffit pour obtenir des vitesses de convergence (ou des bornes de complexité) pour le gradient stochastique appliqué aux problèmes non convexes, convexes ou fortement convexes. Ces résultats sont fortement dépendants du choix des tailles de pas $\{\alpha_k\}_k$: ce choix est de fait un enjeu majeur en apprentissage, qui correspond à la calibration du taux d'apprentissage (ou *learning rate*).

On présente ci-dessous les différentes possibilités pour un tel choix, et leur impact sur les garanties théoriques, dans le cadre des fonctions fortement convexes. Nous ferons donc l'hypothèse suivante.

Hypothèse 3.2.3 *La fonction objectif f est μ -fortement convexe et possède un unique minimum global \mathbf{w}^* ; on notera $f^* = f(\mathbf{w}^*)$.*

On considère d'abord le cas d'une taille de pas constante, pour lequel on peut établir le résultat suivant.

Théorème 3.2.1 (Gradient stochastique à taille de pas constante) *Sous les hypothèses 2.3.1, 3.2.2 et 3.2.3, supposons que l'on applique l'algorithme 6 avec une taille de pas constante*

$$\alpha_k = \alpha \in (0, \frac{1}{2\mu}) \forall k.$$

Alors,

$$\mathbb{E}[f(\mathbf{w}_k) - f^*] \leq \frac{\alpha L \sigma^2}{2\mu} + (1 - 2\alpha\mu)^k \left[f(\mathbf{w}_0) - f^* - \frac{\alpha L \sigma^2}{2\mu} \right]. \quad (3.2.3)$$

Comme on le voit, l'algorithme du gradient stochastique converge en espérance en vitesse linéaire avec une taille de pas constante, tout comme l'algorithme de descente de gradient : cette vitesse signifie que pour garantir $\mathbb{E}[f(\mathbf{w}_k) - f^*] \leq \epsilon$, la méthode requiert au plus $\mathcal{O}(\ln(1/\epsilon))$ itérations. Cependant, ce résultat n'est pas valide pour tout $\epsilon > 0$, contrairement au cas de la descente de gradient. En effet, le résultat (3.2.3) inclut un terme constant dans son second membre, qui représente un biais de $\frac{\alpha L \sigma^2}{4\mu}$. Cela signifie que le théorème 3.2.1 ne peut garantir la convergence que vers un **voisinage** de l'optimum f^* . Dans le même temps, l'utilisation d'une taille de pas constante permet de profiter de pas prometteurs.

Dans sa version originelle (proposée par Robbins et Monro en 1951), la suite des longueurs de pas devait vérifier

$$\sum_{k=0}^{\infty} \alpha_k = \infty \quad \text{et} \quad \sum_{k=0}^{\infty} \alpha_k^2 < \infty.$$

Pour vérifier ces hypothèses, il est nécessaire que α_k tende vers 0. On s'intéresse donc maintenant aux variantes du gradient stochastique basées sur une taille de pas décroissante.

Théorème 3.2.2 (Gradient stochastique à taille de pas décroissante) *Sous les hypothèses ??, 3.2.2 et 3.2.3, on considère l'algorithme 6 appliquée avec une taille de pas décroissante de la forme*

$$\alpha_k = \frac{\beta}{k + \gamma},$$

où $\beta > \frac{1}{\mu}$ et $\gamma > 0$ est choisi de sorte que $\alpha_0 = \frac{\beta}{\gamma} \leq \frac{\mu}{L}$. Dans ce cas, on a

$$\mathbb{E}[f(\mathbf{w}_k) - f^*] \leq \frac{\nu}{\gamma + k}, \quad (3.2.4)$$

avec

$$\nu = \max \left\{ \frac{\beta^2 L \sigma^2}{2(\beta\mu - 1)}, (\gamma + 1)(f(\mathbf{w}_0) - f^*) \right\}.$$

Comme dans le cas de l'algorithme de descente de gradient, le choix d'une taille de pas décroissante peut poser problème, dans la mesure où les pas deviennent nécessairement de plus en plus petits. Par ailleurs, la vitesse de convergence établie par (3.2.4) est sous-linéaire (en $\mathcal{O}(\frac{1}{k})$), ce qui est plus lent que la vitesse linéaire (en $\mathcal{O}((1-t)^k)$) obtenue par une taille de pas constante; en revanche, avec une taille de pas décroissante, la méthode peut satisfaire $\mathbb{E}[f(\mathbf{w}_k) - f^*] \leq \epsilon$ pour toute valeur $\epsilon > 0$.

Remarque 3.2.3 (Une approche hybride) *En apprentissage, une stratégie hybride entre la taille de pas constante et la taille de pas décroissante est souvent utilisée. Celle-ci consiste à lancer l'algorithme avec une taille de pas constante égale à α jusqu'à décroître $f(\mathbf{w}_k) - f^*$ en dessous d'un certain seuil (typiquement égal à $\frac{\alpha L \sigma^2}{2\mu}$). On choisit ensuite $\alpha' < \alpha$, et on continue l'algorithme avec ce nouveau pas, en visant la précision $\frac{\alpha' L \sigma^2}{2\mu}$. En procédant par réductions successives de α , on peut ainsi atteindre n'importe quelle précision, tout en conservant un pas constant durant plusieurs itérations. On peut montrer que la vitesse de convergence d'un tel processus est sous-linéaire, en ce sens que pour tout $\epsilon > 0$,*

$$\mathbb{E}[f(\mathbf{w}_k) - f^*] \leq \epsilon \quad \text{after } \mathcal{O}(1/\epsilon) \text{ itérations.}$$

Pour appliquer cette stratégie de manière adaptive, il est nécessaire de déterminer quand le voisinage est atteint; en pratique, on regarde si l'algorithme stagne, c'est-à-dire ne progresse plus (cela peut être quantifié en termes d'itérés, de norme du gradient stochastique, etc).

Taille de pas dans le cas nonconvexe L'algorithme du gradient stochastique (et ses variantes) sont les méthodes les plus fréquemment utilisées pour l'entraînement des réseaux de neurones : ce problème est fortement non convexe, et l'analyse ci-dessus ne peut donc y être appliquée. Cependant, il est possible de déterminer des vitesses de convergence pour le cas convexe, en étudiant les quantités suivantes :

- $\mathbb{E} \left[\frac{1}{K} \sum_{i=1}^K \|\nabla f(\mathbf{w}_k)\|^2 \right]$ pour les variantes de l'algorithme 6 utilisant une taille de pas constante;
- $\mathbb{E} \left[\frac{1}{\sum_{i=1}^K \alpha_k} \sum_{i=1}^K \alpha_k \|\nabla f(\mathbf{w}_k)\|^2 \right]$ pour les variantes utilisant une taille de pas décroissante.

De par l'utilisation d'approximations du gradient, on obtiendra des bornes qui seront moins bonnes que celles du cas déterministe. À titre d'exemple, l'algorithme 6 appliqué avec une longueur de pas constante vérifiera

$$\mathbb{E} \left[\frac{1}{K} \sum_{i=1}^K \|\nabla f(\mathbf{w}_k)\|^2 \right] \leq \epsilon$$

en au plus $\mathcal{O}(\epsilon^{-4})$ itérations, ce qui est une plus mauvaise borne que dans le cas déterministe, en $\mathcal{O}(\epsilon^{-2})$. De plus, ce résultat ne s'appliquera que pour une valeur de ϵ suffisamment large, à cause du biais introduit par l'utilisation de quantités stochastiques.

Remarque 3.2.4 (Utilisation du momentum) Les variantes du gradient stochastique les plus populaires en pratique (ADAM, ADAGRAD, RMSPROP) incorporent généralement un terme de momentum dans l'itération du gradient stochastique. L'idée sous-jacente est que l'ajout de momentum permet d'accumuler les pas dans des bonnes directions (qui sont souvent de descente pour l'algorithme), tandis que les mauvaises directions (correspondant par exemple aux outliers, ou données aberrantes) se compenseront au fil du temps. L'analyse de ces méthodes accélérées est cependant beaucoup plus complexe que celle du cas déterministe, et à l'heure actuelle la théorie et la pratique sont encore décorrélées, contrairement au cas déterministe.

3.3 Réduction de variance

Comme nous l'avons vu dans la partie précédente, la théorie du gradient stochastique repose sur l'hypothèse 3.2.2, et plus particulièrement sur un contrôle de la variance de la norme du gradient stochastique (à travers la quantité σ^2). En observant la dépendance en σ des vitesses de convergence telles que (3.2.3), on voit qu'une valeur élevé de σ conduit à de plus mauvaises bornes. Cela se traduit numériquement par le fait qu'une méthode avec une forte variance risque de converger lentement.

Les techniques dites de **réduction de variance** ont été précisément développées dans le but de diminuer la variance des estimations de gradient utilisées par l'algorithme du gradient stochastique. On peut classer ces techniques en deux groupes, selon qu'elles utilisent plusieurs exemples à chaque itération ou qu'elles se basent sur les itérations précédentes. Nous nous concentrerons ici sur les stratégies relevant de la première catégorie.

3.3.1 Variantes à lots (batch)

Comme on l'a vu dans la partie précédente, l'itération de l'algorithme du gradient stochastique se ramène à

$$\mathbf{w}_{k+1} = \mathbf{w}_k - \alpha_k \nabla f_{i_k}(\mathbf{w}_k),$$

où i_k est un indice tiré au hasard dans $\{1, \dots, n\}$. Cette méthode utilise ainsi un **seul** exemple pour construire son gradient stochastique, et la variance σ^2 associée à ce gradient (voir hypothèse 3.2.2) provient directement de cet estimateur. A contrario, l'algorithme de descente de gradient utilise **tous** les exemples et le résultat possède une variance égale à zéro (car il s'agit d'une quantité déterministe).

Si l'on cherche à améliorer cette variance, il semble donc naturel de considérer des estimateurs stochastiques du gradient basés sur **plusieurs** exemples à la fois : c'est le principe des méthodes de lots¹ de gradients stochastiques.

Le principe d'une telle méthode est de tirer un ensemble aléatoire d'indices $S_k \subset \{1, \dots, n\}$ puis d'effectuer l'itération suivante :

$$\mathbf{w}_{k+1} = \mathbf{w}_k - \alpha_k \frac{1}{|S_k|} \sum_{i \in S_k} \nabla f_i(\mathbf{w}_k) \quad (3.3.1)$$

Lorsque $|S_k| = 1$, on retrouve la formule du gradient stochastique classique. En considérant $|S_k| = n$, on aurait nécessairement $S_k = \{1, \dots, n\}$, auquel cas l'itération (3.3.1) est équivalente à celle de la descente de gradient.

Plus globalement, on peut identifier deux classes de tailles de lots :

¹Ou *batch*, en anglais.

- $|S_k| \approx n$: une telle variante possède un coût par itération proche de celui de la descente de gradient, et par conséquent est soumise aux mêmes problématiques de coût;
- $|S_k| = n_b << n$, que l'on appelle mini-lot (ou **mini-batching**), qui peut être un choix avantageux sur le plan théorique, raisonnable en pratique et permet de réduire la variance. Cette méthode est communément appelée l'algorithme de mini-lot de gradient stochastique, ou **mini-batch SG**.

Si l'on suppose que la taille du lot est constante, c'est-à-dire que $|S_k| = n_b \forall k$, on peut alors montrer que, pour la même longueur de pas, la variante avec mini-lot requiert n_b fois moins d'itérations que le gradient stochastique. Ces dernières sont n_b fois plus coûteuses, mais l'approche par mini-lots permet d'exploiter le calcul parallèle, en calculant les n_b gradients stochastiques sur différents processeurs. Par ailleurs, l'estimé de gradient stochastique vérifie la propriété suivante.

Proposition 3.3.1 *Sous les hypothèses 2.3.1 et 3.2.2, la variance du gradient stochastique en “mini-lot” vérifie :*

$$\mathbb{E}_{S_k} \left[\left\| \frac{1}{|S_k|} \sum_{i \in S_k} \nabla f_i(\mathbf{w}_k) \right\|^2 \right] \leq \frac{\sigma^2}{n_b}.$$

Pour terminer cette partie, nous mentionnons que les techniques par (mini)-lots restent plus coûteuses (au niveau d'une itération), tout en étant plus sensibles aux redondances dans les données; par ailleurs, la taille du lot (qui peut être fixe ou adaptative) introduit un hyperparamètre supplémentaire à choisir. Ces considérations expliquent en partie que l'approche classique du gradient stochastique reste la plus fréquente en pratique.

3.3.2 Autres variantes basées sur la réduction de variance

Les méthodes d'agrégation (**aggregation gradient**) sont d'autres stratégies de réduction de variance dont les propriétés théoriques (et notamment leurs vitesses de convergence linéaires) ont fait le succès dans la communauté de l'optimisation et de la théorie de l'apprentissage. Elles consistent en substance à effectuer un **pas complet de descente de gradient** à intervalles réguliers durant l'exécution de l'algorithme : cela permet de corriger les composantes du gradient stochastique qui possèderaient une variance trop élevée. En dépit de leur analyse théorique très fouillée, les variantes de ces méthodes ne sont pas encore déployées en pratique, principalement en raison du coût prohibitif que peut représenter un calcul de gradient, même s'il n'est effectué qu'une fois toutes les K itérations (avec $K \gg 1$).

Faire la moyenne des itérés est une autre manière de réduire la variance, moins coûteuse à implémenter que la précédente. L'idée sous-jacente consiste à étudier les propriétés de l'itéré moyen $\frac{1}{K} \sum_{k=0}^{K-1} \mathbf{w}_k$: dans certains cas (par exemple lorsque $\alpha = \frac{1}{\mu(k+1)}$ et f est μ -fortement convexe), on peut montrer de bonnes propriétés pour cet itéré, notamment qu'il s'agit d'une solution plus robuste que le dernier itéré obtenu. Cependant, afin de renvoyer cet itéré moyen, il est nécessaire soit de stocker l'historique des itérés (ce qui peut être coûteux), soit de maintenir un itéré moyen, et cela peut générer des erreurs numériques.

3.4 Méthodes de gradient stochastique pour l'apprentissage profond

Dans cette partie, on s'intéresse aux techniques de gradient stochastique utilisées pour entraîner des modèles d'apprentissage profond. On considère toujours un problème de la forme (3.1.1), sous l'hypothèse 3.2.1. Notre but est d'analyser différentes variations du schéma

$$\mathbf{w}_{k+1} = \mathbf{w}_k - \alpha \mathbf{g}_k, \quad (3.4.1)$$

où $\alpha > 0$ est une longueur de pas/un taux d'apprentissage (*learning rate*) et \mathbf{g}_k est un estimateur stochastique du gradient correspondant à un seul indice (comme le gradient stochastique) ou à un paquet (*batch*) d'indices.

On s'intéressera par la suite à un schéma générique qui couvre toutes les variantes que nous allons étudier. Ce schéma est de la forme :

$$\mathbf{w}_{k+1} = \mathbf{w}_k - \alpha \mathbf{m}_k \oslash \mathbf{v}_k, \quad (3.4.2)$$

où $\alpha > 0$, $\mathbf{m}_k, \mathbf{v}_k \in \mathbb{R}^d$ et \oslash symbolise la division composante à composante :

$$\mathbf{m}_k \oslash \mathbf{v}_k := \left[\frac{[\mathbf{m}_k]_i}{[\mathbf{v}_k]_i} \right]_{i=1,\dots,d}.$$

Cette itération généralise (3.4.1) : en effet, en posant $\mathbf{m}_k = \mathbf{g}_k$ et $\mathbf{v}_k = \mathbf{1}_{\mathbb{R}^d}$, on retrouve l'itération 3.4.1. Ce formalisme permet d'exprimer de façon unifiée les méthodes de gradient stochastique les plus populaires en apprentissage profond, et ainsi de mieux les comparer.

3.4.1 Gradient stochastique avec momentum

Dans la lignée des techniques accélérées que nous avons vues au chapitre 2, la plupart des implémentations de gradient stochastique considèrent l'addition de momentum au pas basique (3.4.1). Ainsi, une itération de gradient stochastique avec momentum s'écrira

$$\mathbf{w}_{k+1} = \mathbf{w}_k - \alpha(1 - \beta)\mathbf{g}_k + \alpha\beta(\mathbf{w}_k - \mathbf{w}_{k-1}), \quad (3.4.3)$$

où $\beta \in (0, 1)$ est un paramètre constant (choisir $\beta = 0$ reviendrait à faire du gradient stochastique). Il s'agit d'une version de l'algorithme de Polyak pour laquelle la direction de type gradient est combinée avec la direction précédente : on retrouve l'idée selon laquelle l'utilisation du momentum permet d'accumuler l'information de l'itération précédente. En pratique, on observe que l'itération (3.4.3) tend à accumuler les bonnes directions (en termes d'optimisation), tandis que les "mauvaises" directions (et donc les "mauvais" pas) pour l'optimisation ont tendance à se compenser.

L'itération (3.4.3) est un cas particulier de (3.4.2), correspondant à $\mathbf{v}_k = \mathbf{1}_{\mathbb{R}^d}$ et \mathbf{m}_k défini de manière récursive par $\mathbf{m}_{-1} = \mathbf{0}_{\mathbb{R}^d}$ et

$$\mathbf{m}_k = (1 - \beta)\mathbf{g}_k - \beta\mathbf{m}_{k-1} \quad \forall k \in \mathbb{N}.$$

où $\beta \in (0, 1)$.

La méthode de gradient stochastique avec momentum est implémentée dans les bibliothèques d'apprentissage profond telles que PyTorch. Elle est particulièrement dans l'entraînement de réseaux de neurones profonds sur des problèmes de vision par ordinateur, et est à l'origine de la montée du "deep learning" au début des années 2010.

Remarque 3.4.1 Les garanties de l'algorithme (3.4.3) sont plus difficiles à établir que dans le cas de la descente de gradient accélérée : les approches par momentum ont cependant rencontré un certain succès pratique, même sur des problèmes non convexes tels que l'entraînement de réseaux de neurones.

3.4.2 AdaGrad

La méthode de gradient adaptatif, ou ADAGRAD, a été proposée en 2011 pour répondre à la difficulté du choix de α dans le gradient stochastique tout en évitant d'avoir recours à des procédures adaptatives coûteuses telles que la recherche linéaire. L'approche d'ADAAGRAD consiste à normaliser chaque composante du gradient stochastique au moyen d'une accumulation des valeurs de chaque composante au cours des itérations. L'algorithme maintient donc une suite $\{\mathbf{r}_k\}_k$ définie par

$$\forall i = 1, \dots, d, \quad \begin{cases} [\mathbf{r}_{-1}]_i = 0 \\ [\mathbf{r}_k]_i = [\mathbf{r}_{k-1}]_i + [\mathbf{g}_k]_i^2 \quad \forall k \geq 0, \end{cases} \quad (3.4.4)$$

L'itération d'ADAAGRAD s'écrit alors

$$\mathbf{w}_{k+1} = \mathbf{w}_k - \alpha \mathbf{g}_k \oslash \sqrt{\mathbf{r}_k}, \quad (3.4.5)$$

où la racine carrée est appliquée à toutes les composantes de \mathbf{r}_k . On reconnaît ainsi l'itération générique (3.4.2) avec $\mathbf{m}_k = \mathbf{g}_k$ et $\mathbf{v}_k = \sqrt{\mathbf{r}_k}$. L'innovation d'ADAAGRAD tient donc non pas dans l'introduction de momentum, mais dans l'utilisation d'une suite de pas différenciée pour chaque coordonnée, de la forme

$$\left\{ \left[\frac{\alpha}{\sqrt{[\mathbf{r}_k]_i}} \right]_{i=1}^d \right\}_k.$$

On opère ainsi une normalisation diagonale (*diagonal scaling*) des composantes du gradient stochastique \mathbf{g}_k , qui permet notamment de traiter le cas de composantes très différentes en amplitude sans avoir à calibrer α très finement. Cependant, de tels pas tendent généralement vers 0 assez rapidement.

Remarque 3.4.2 En pratique, on remplace \mathbf{r}_k par $\mathbf{r}_k + \eta \mathbf{1}_{\mathbb{R}^d}$ où $\eta > 0$ est de valeur faible, pour que l'algorithme soit numériquement plus stable.

La méthode ADAAGRAD est particulièrement pertinente dans le cas de problèmes à gradients *parcimonieux* (ou *sparse* en anglais), pour lesquels les gradients stochastiques ont tendance à avoir beaucoup de composantes nulles. Dans une telle situation, le calcul de \mathbf{r}_k permettra d'ajuster le pas uniquement pour les coordonnées non nulles. Les problèmes issus des systèmes de recommandation possèdent généralement cette propriété, et il s'agit du domaine où ADAAGRAD est généralement considéré comme l'algorithme référence.

3.4.3 RMSProp

L'algorithme RMSPROP (pour *Root Mean Square Propagation*) procède de manière similaire à ADAAGRAD en jouant sur les composantes du gradient. La méthode repose sur une suite $\{\mathbf{r}_k\}_k$ définie par

$$\forall i = 1, \dots, d, \quad \begin{cases} [\mathbf{r}_{-1}]_i = 0 \\ [\mathbf{r}_k]_i = (1 - \lambda)[\mathbf{r}_{k-1}]_i + \lambda[\mathbf{g}_k]_i^2 \quad \forall k \geq 0, \end{cases} \quad (3.4.6)$$

avec $\lambda \in (0, 1)$. On voit donc que la méthode peut donner (en fonction de la valeur de λ) plus de poids aux gradients précédents plutôt qu'au gradient de l'itération : cette idée permet aux longueurs de pas de décroître moins rapidement que celles de ADAGRAD.

Comme pour ADAGRAD, l'itération de RMSPROP s'écrit alors sous la forme (3.4.2) avec $\mathbf{m}_k = \mathbf{g}_k$ et $\mathbf{v}_k = \sqrt{\mathbf{r}_k}$.

Remarque 3.4.3 En pratique, et comme pour ADAGRAD, on remplacera typiquement \mathbf{r}_k par $\mathbf{r}_k + \eta \mathbf{1}_{\mathbb{R}^d}$ pour une petite valeur $\eta > 0$.

La méthode RMSPROP a été utilisée avec succès dans le cadre d'entraînements de réseaux de neurones profonds.

3.4.4 Adam

L'algorithme ADAM, proposé en 2013, peut être vu comme combinant le concept de momentum avec celui d'accumulation d'information sur les gradients précédents utilisés par les algorithmes précédents. Son itération correspond au schéma générique (3.4.2) avec

$$\mathbf{m}_k = \frac{(1 - \beta_1) \sum_{j=0}^k \beta_1^{k-j} \mathbf{g}_j}{1 - \beta_1^{k+1}} \quad (3.4.7)$$

et

$$\mathbf{v}_k = \sqrt{\frac{(1 - \beta_2) \sum_{j=0}^k \beta_2^{k-j} \mathbf{g}_j \odot \mathbf{g}_j}{1 - \beta_2^{k+1}}}. \quad (3.4.8)$$

avec $\beta_1, \beta_2 \in (0, 1)$ et \odot le produit composante à composante, dit de Hadamard :

$$\mathbf{g}_k \odot \mathbf{g}_k = [[\mathbf{g}_k]_i^2]_{i=1}^d.$$

Remarque 3.4.4 En pratique, on utilisera $\mathbf{v}_k + \eta \mathbf{1}_{\mathbb{R}^d}$ pour une petite valeur $\eta > 0$ plutôt que \mathbf{v}_k .

Les formules ci-dessous correspondent respectivement à une combinaison des directions précédemment employées qui met l'accent sur les directions les plus récentes, et à une normalisation des coordonnées des directions obtenues relativement à une moyenne de ces coordonnées donnant aussi plus d'importance aux dernières itérations. Cet aspect important, qui se justifie de manière statistique, semble être à l'origine du succès d'ADAM, dont la performance impressionnante sur de nombreuses tâches d'entraînement de réseaux de neurones a contribué à son utilisation massive. La méthode ADAM (et sa variante ADAMW, basée sur de la régularisation dont nous parlerons plus loin) sont notamment très efficaces sur des tâches de traitement automatique des langues.

3.5 Conclusion

Les méthodes de gradient stochastique reposent sur de l'information partielle du gradient, et ont donc de moins bonnes propriétés de convergence que l'algorithme de descente de gradient : en ce sens, elles ne sont pas forcément intéressantes pour un problème d'optimisation quelconque. En revanche, lorsque le calcul du gradient implique la manipulation d'un jeu de données très volumineux, les stratégies de gradient stochastique prennent tout leur sens, et se révèlent particulièrement adaptés pour de tels problèmes. Cela est dû au très faible coût de ces méthodes, mais aussi à leur efficacité sur des données possiblement aléatoires ou redondantes : dans un tel contexte, il est payant de privilégier les pas aléatoires et peu coûteux, et on observe ainsi une convergence nettement plus rapide en pratique.

Même si l'algorithme du gradient stochastique est très efficace en général, sa performance peut être sensiblement affectée par des gradients stochastiques à forte variance. Les implémentations de cette méthode cherchent généralement à en réduire la variance, en considérant par exemple des paquets d'exemples simultanément. Les variantes les plus populaires en pratique, et notamment en apprentissage profond, sont basées sur le principe d'accélération, ainsi que sur la normalisation diagonale des directions utilisées. Il est à noter que ces techniques peuvent manquer de justification théorique relativement aux problèmes auxquels elles sont appliquées (typiquement l'entraînement d'un réseau de neurones qui donne lieu à un problème non convexe). Cependant, ces méthodes ont été adoptées par la communauté du fait de leur succès pratique, et leur analyse complète reste un problème ouvert.

Chapitre 4

Optimisation non lisse et régularisation

Dans ce chapitre, nous considérons deux caractéristiques fréquentes des problèmes de sciences des données : l'absence potentielle de dérivées, et la volonté d'imposer une certaine structure sur les modèles considérés. Nous illustrons ces deux aspects au moyen d'un exemple classique d'apprentissage, puis étudions les concepts d'optimisation sous-jacents.

4.1 Introduction : algorithme du perceptron

On revient ici sur le problème vu en section 1.1.2, qui possédait la forme suivante :

$$\underset{\mathbf{w} \in \mathbb{R}^d}{\text{minimiser}} \frac{1}{n} \sum_{i=1}^n \max\{1 - y_i \mathbf{x}_i^\top \mathbf{w}, 0\} + \frac{\lambda}{2} \|\mathbf{w}\|_2^2 \quad (4.1.1)$$

où $\{(\mathbf{x}_i, y_i)\}_{i=1}^n$ forme un jeu de données dans $\mathbb{R}^d \times \{-1, 1\}$, et $\lambda > 0$.

L'une des méthodes classiques pour traiter ce problème est l'**algorithme du perceptron**, que l'on présente dans l'algorithme 7.

Algorithme 7: Algorithme du perceptron pour le problème 4.1.2.

Initialisation: $\mathbf{w}_0 \in \mathbb{R}^d$, $\alpha > 0$.

Pour $k = 0, 1, \dots$

 1. Tirer $i_k \in \{1, \dots, n\}$ au hasard.

 2. Calculer

$$\mathbf{w}_{k+1} = \left(1 - \frac{\alpha \lambda}{n}\right) \mathbf{w}_k + \begin{cases} \alpha y_{i_k} \mathbf{x}_{i_k}^\top \mathbf{w}_k & \text{si } 1 - y_{i_k} \mathbf{x}_{i_k}^\top \mathbf{w}_k \\ 0 & \text{sinon,} \end{cases} \quad (4.1.2)$$

Fin

Sous sa forme de base, la méthode du perceptron ressemble fortement à l'algorithme du gradient stochastique avec taille de pas constante : en effet, il s'agit de sélectionner un exemple du jeu de données à chaque itération et d'effectuer une mise à jour sur la base de cet exemple. De fait, il

s'agirait de l'algorithme du gradient stochastique pour le problème

$$\underset{\mathbf{w} \in \mathbb{R}^d}{\text{minimiser}} \frac{1}{n} \sum_{i=1}^n (1 - y_i \mathbf{x}_i^\top \mathbf{w}) + \frac{\lambda}{2} \|\mathbf{w}\|_2^2.$$

Cependant, un tel choix de fonction de coût ne modélisera pas bien notre problème (voir la discussion en section 1.1.2). La fonction de perte que nous avons choisie (*hinge loss*) correspond mieux à nos attentes, mais elle possède l'inconvénient d'être **non lisse**, c'est-à-dire qu'elle n'est pas dérivable en tout point de l'espace. Lorsqu'on est confronté à un tel problème, il est souvent possible de définir des notions de dérivées plus génériques que celle du gradient, et ainsi de développer des algorithmes d'optimisation. C'est ce que fait implicitement l'algorithme du perceptron, et ce que nous verrons en partie 4.2.

La seconde propriété notable du problème (4.1.1) est la présence d'un second terme dans la fonction objectif, qui dépend non pas des données mais uniquement du modèle \mathbf{w} , et est contrôlé par une constante λ . Dans l'algorithme du perceptron, on voit que plus λ sera important, plus l'algorithme tendra à réduire les composantes de w_k . Il s'agit d'un processus de régularisation, que nous étudierons plus en détail et plus généralement en partie 4.3.

4.2 Optimisation non lisse

4.2.1 Des fonctions aux problèmes non lisses

Les problèmes tels que (4.1.1), pour lesquels la fonction objectif n'est pas dérivable, sont appelés des *problèmes non lisses*. Leurs fonctions objectifs sont également qualifiées de non lisses (par opposition aux fonctions lisses). Pour ce chapitre, on définira les fonctions non lisses comme suit.

Définition 4.2.1 (Fonctions non lisses) Une fonction $f : \mathbb{R}^d \rightarrow \mathbb{R}$ est dite **non lisse** si elle n'est pas dérivable partout.

Remarque 4.2.1 Une fonction non lisse peut être continue (voir par exemple la fonction objectif de (4.1.1)).

Exemple 4.2.1 (Exemples de fonctions non lisses)

- $w \mapsto |w|$ de \mathbb{R} dans \mathbb{R} ;
- $w \mapsto \|w\|_1$ de \mathbb{R}^d dans \mathbb{R} ;
- $\text{ReLU} : w \mapsto \max\{w, 0\}$ de \mathbb{R}^d dans \mathbb{R} .

Comme les fonctions non lisses ne sont pas différentiables partout, il n'est pas envisageable de résoudre les problèmes non lisses via des méthodes basées sur le gradient. Il existe cependant plusieurs approches possibles pour traiter ces problèmes.

Une première technique consiste à trouver une formulation lisse équivalente à celle du problème non lisse. Le problème $\underset{w \in \mathbb{R}}{\text{minimiser}} |w|$ est par exemple équivalent à

$$\underset{w, t^+, t^- \in \mathbb{R}}{\text{minimiser}} t^+ + t^- \quad \text{s. c. } w = t^+ - t^-, t^+ \geq 0, t^- \geq 0.$$

Cette reformulation est un programme linéaire, qui fait partie des classes de problèmes classiques, pour lesquelles des algorithmes et codes très efficaces existent.

Lorsque la fonction est non lisse mais lipschitzienne (ce que l'on notera $\mathcal{C}_L^{0,0}$, par analogie avec $\mathcal{C}_L^{1,1}$), on peut en revanche essayer d'appliquer une méthode de gradient, en raison de la propriété suivante.

Théorème 4.2.1 Soit $f : \mathbb{R}^d \rightarrow \mathbb{R}$ une fonction lipschitzienne. Alors, f est dérivable en presque tout point de \mathbb{R}^d .

La fonction ReLU est un exemple classique de fonction lipschitzienne, ce qui fait que les constructions impliquant des fonctions ReLU (comme les réseaux de neurones) ne sont pas dérivables en tout point. En revanche, il est possible de traiter ces problèmes comme dérivables presque partout; cela reste un problème pour certifier qu'un point est un minimum local, car il est fréquent que les fonctions non lisses ne soient pas dérivables en leurs points extrêmaux. La notion de dérivée généralisée, que nous présentons ci-après, est précisément utilisée dans ce but.

4.2.2 Méthodes de sous-gradient

Nous nous concentrons ici sur le cas des fonctions non lisses convexes, cas fréquent en pratique, et nous définissons une notion généralisée de gradient.

Définition 4.2.2 (Sous-gradient et sous-différentiel) Soit $f : \mathbb{R}^d \rightarrow \mathbb{R}$ une fonction convexe. Un vecteur $\mathbf{g} \in \mathbb{R}^d$ est appelé un **sous-gradient** de f en $\mathbf{w} \in \mathbb{R}^d$ si

$$\forall \mathbf{z} \in \mathbb{R}^n, \quad f(\mathbf{z}) \geq f(\mathbf{w}) + \mathbf{g}^T(\mathbf{z} - \mathbf{w}).$$

L'ensemble des sous-gradients de f en \mathbf{w} est appelé le **sous-différentiel** de f en \mathbf{w} ; on le note $\partial f(\mathbf{w})$.

Lorsque la fonction f est dérivable en \mathbf{w} , on a $\partial f(\mathbf{w}) = \{\nabla f(\mathbf{w})\}$, ce qui montre que la notion de sous-différentiel généralise bien celle du gradient.

Les sous-différentiels permettent de caractériser l'optimalité (globale) pour les fonctions convexes, comme le montre le résultat ci-dessous.

Théorème 4.2.2 Soit $f : \mathbb{R}^d \rightarrow \mathbb{R}$ une fonction convexe, et $\mathbf{w} \in \mathbb{R}^d$.

$$\mathbf{0} \in \partial f(\mathbf{w}) \Leftrightarrow \mathbf{w} \text{ minimum of } f$$

Exemple 4.2.2 Soit $f : \mathbb{R} \rightarrow \mathbb{R}$, $f(w) = |w|$.

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

L'ensemble $[-1, 1]$ contient 0, ce qui prouve que $w^* = 0$ est un minimum (en fait l'unique minimum) de f .

Remarque 4.2.2 Il est aussi possible de définir des sous-gradients pour les fonctions convexes: cependant, pour de telles fonctions, le sous-différentiel peut être vide en certains points (par exemple en un maximum local), ce qui en limite la portée. Il existe d'autres notions de dérivée généralisée, qui peuvent être utilisées dans un processus d'optimisation.

Algorithme 8: Méthode de sous-gradient.

Initialisation : $\mathbf{w}_0 \in \mathbb{R}^d$.

Pour $k = 0, 1, \dots$

1. Calculer un sous-gradient $\mathbf{g}_k \in \partial f(\mathbf{w}_k)$.
2. Calculer une taille de pas $\alpha_k > 0$.
3. Poser $\mathbf{w}_{k+1} = \mathbf{w}_k - \alpha_k \mathbf{g}_k$.

Fin

En utilisant la notion de sous-gradient, on peut alors définir l'analogue de la descente de gradient pour des fonctions convexes non lisses : c'est le principe décrit dans l'algorithme 8.

Une telle méthode offre un degré de liberté dans le choix du sous-gradient, ce qui peut poser problème. En effet, un sous-gradient peut ne pas être une direction de descente, et au contraire être une direction de montée *pour tout pas possible*. C'est notamment le cas lorsque l'on se trouve en un minimum pour lequel il existe plusieurs sous-gradients : toute direction autre que $\mathbf{0}$ sera nécessairement une direction où la fonction augmentera localement.

Variantes de la méthode du sous-gradient Les variantes de l'algorithme de descente de gradient ont souvent un équivalent dans le cadre des méthodes de sous-gradient, même si leur analyse est parfois plus délicate. L'une des méthodes les plus populaires est l'algorithme du sous-gradient stochastique, utilisé en optimisation stochastique et (par conséquent) en apprentissage de manière générique.

4.3 Régularisation

4.3.1 Problèmes régularisés

Comme décrit dans l'introduction de ces notes, une pratique courante en apprentissage consiste à favoriser les modèles possédant une structure spécifique. En termes d'optimisation, cela se fait **à travers la fonction objectif**, ce qui donne des problèmes de la forme

$$\underset{\mathbf{w} \in \mathbb{R}^d}{\text{minimiser}} f(\mathbf{w}) + \lambda \Omega(\mathbf{w}).$$

où f est une fonction de coût, Ω est une fonction appelée *terme de régularisation* et $\lambda > 0$ est appelé un paramètre de régularisation.

Exemple 4.3.1 (Régularisation “ridge”) Un problème avec régularisation écrêtée, ou ridge en anglais, est de la forme suivante :

$$\underset{\mathbf{w} \in \mathbb{R}^d}{\text{minimiser}} f(\mathbf{w}) + \frac{\lambda}{2} \|\mathbf{w}\|^2.$$

Le terme de régularisation $\mathbf{w} \mapsto \frac{1}{2} \|\mathbf{w}\|^2$ possède plusieurs interprétations. Il pénalise les vecteurs \mathbf{w} dont les composantes sont larges, et on peut montrer que sa présence est équivalente à imposer une

contrainte sur la norme au carré $\|\mathbf{w}\|^2$. Par ailleurs, la régularisation écrêtée réduit la variance du problème par rapport aux données définissant f , ce qui est important dans notre contexte. Enfin, lorsque $\lambda > 0$ est suffisamment grand et que la fonction de coût f est minorée, on peut montrer que la fonction objectif est fortement convexe, et donc qu'il existe au plus un minimum global au problème.

Remarque 4.3.1 Dans le cadre d'un algorithme de type descente de gradient (ou gradient stochastique classique) utilisé en apprentissage, ajouter un terme de régularisation ℓ_2 correspond au concept de batch normalization.

4.3.2 Régularisation et parcimonie

Lorsque l'on construit un modèle basé sur des données réparties en attributs (*features*) et labels, on peut vouloir un modèle qui explique les données en utilisant peu d'attributs : outre qu'un tel modèle est souvent plus facile à interpréter, il possède l'avantage de sélectionner les attributs les plus significatifs (on parle parfois de *feature selection*). Du point de vue de l'optimisation, si le modèle en question est représenté par un vecteur $\mathbf{w} \in \mathbb{R}^d$, on cherche un vecteur qui soit solution du problème d'optimisation mais possède autant de coordonnées nulles que possible (on dira que l'on cherche un vecteur creux, ou parcimonieux).

Il existe des termes de régularisation qui pénalisent les vecteurs avec des composantes non nulles (et non les composantes globalement larges, comme pour la régularisation écrêtée). La fonction norme ℓ_0 est une expression exacte de cette pénalisation¹. Ainsi, un problème avec régularisation ℓ_0 est de la forme

$$\underset{\mathbf{w}}{\text{minimiser}} f(\mathbf{w}) + \lambda \|\mathbf{w}\|_0, \quad \|\mathbf{v}\|_0 = |\{i | v_i \neq 0\}|.$$

Cependant, cette fonction est non convexe, non lisse et discontinue; elle possède également une nature combinatoire qui la rend complexe à utiliser en optimisation. On lui préfère donc en général la norme ℓ_1 , définie par

$$\|\mathbf{w}\|_1 = \sum_{i=1}^d |w_i|. \tag{4.3.1}$$

Cette fonction est non lisse, mais elle est continue et convexe; il s'agit également d'une norme, ce qui lui confère de nombreuses propriétés intéressantes en optimisation.

Exemple 4.3.2 LASSO (Least Absolute Shrinkage and Selection Operator) On se place dans le cadre de la régression linéaire sur des données $\mathbf{X} \in \mathbb{R}^{n \times d}$ et $\mathbf{y} \in \mathbb{R}^n$. Le problème des moindres carrés linéaires avec régularisation ℓ_1 s'écrit :

$$\underset{\mathbf{w} \in \mathbb{R}^d}{\text{minimiser}} \frac{1}{2} \|\mathbf{X}\mathbf{w} - \mathbf{y}\|^2 + \lambda \|\mathbf{w}\|_1.$$

On peut montrer que la solution de ce problème possède moins d'éléments non nuls que la solution du problème sans terme de régularisation, donnée par la pseudo-inverse.

¹Communément appelée "norme ℓ_0 " ainsi même s'il ne s'agit pas d'une norme.

4.3.3 Méthodes proximales

Nous allons maintenant décrire une classe d'algorithmes dédiée à la résolution de problèmes d'optimisation avec régularisation. Dans ces notes, on s'intéressera plus précisément à la catégorie ci-dessous.

Définition 4.3.1 (Optimisation composite) *Un problème d'optimisation composite est de la forme*

$$\underset{\mathbf{w} \in \mathbb{R}^d}{\text{minimiser}} f(\mathbf{w}) + \lambda \Omega(\mathbf{w}),$$

où $f : \mathbb{R}^d \rightarrow \mathbb{R}$ est une fonction lisse et $\mathcal{C}^{1,1}$, $\lambda > 0$, et $\Omega : \mathbb{R}^d \rightarrow \mathbb{R}$ est une fonction convexe non lisse.

Cette définition couvre nos deux exemples de la partie précédente. Pour les résoudre, nous allons adopter une approche dite proximale.

L'**approche proximale** procède selon un schéma très fréquent en optimisation : un problème donné est remplacé par une suite de sous-problèmes présumés plus faciles à résoudre². Dans le cas des méthodes proximales, on exploite la dérivabilité de f pour obtenir un problème plus simple, tandis que la structure de Ω est incorporée dans les sous-problèmes.

Algorithme 9: Méthode du gradient proximal

Initialisation : $\mathbf{w}_0 \in \mathbb{R}^d$.

Pour $k = 0, 1, \dots$

- 1. Calculer le gradient de la partie lisse du problème $\nabla f(\mathbf{w}_k)$.
- 2. Définir une taille de pas $\alpha_k > 0$.
- 3. Calculer le nouvel itéré \mathbf{w}_{k+1} tel que

$$\mathbf{w}_{k+1} \in \underset{\mathbf{w} \in \mathbb{R}^d}{\operatorname{argmin}} \left\{ f(\mathbf{w}_k) + \nabla f(\mathbf{w}_k)^T (\mathbf{w} - \mathbf{w}_k) + \frac{1}{2\alpha_k} \|\mathbf{w} - \mathbf{w}_k\|^2 + \lambda \Omega(\mathbf{w}) \right\}. \quad (4.3.2)$$

Fin

L'algorithme 9 décrit le déroulement d'une méthode proximale. Le coût de chaque itération est plus élevé que celui des méthodes que nous avons vues jusqu'à présent, car une itération comporte un calcul de gradient puis une résolution d'un sous-problème auxiliaire (4.3.2), que l'on appelle **sous-problème proximal**.

Remarque 4.3.2 Si $\Omega \equiv 0$ (Ω est la fonction nulle et il n'y a pas de régularisation), on peut montrer que la solution de (4.3.2) est unique, et donnée par

$$\mathbf{w}_{k+1} = \mathbf{w}_k - \alpha_k \nabla f(\mathbf{w}_k).$$

On reconnaît ainsi l'itération de la descente de gradient décrite par l'algorithme 3.

²Toutes les méthodes vues dans ce cours suivent implicitement cette approche.

Les méthodes de gradient proximal peuvent être munies de la plupart des outils qui peuvent être utilisés pour la descente de gradient : différentes tailles de pas, accélération, utilisation de gradients stochastiques, etc. Par ailleurs, il existe de nombreux résultats de complexité pour les méthodes proximales, principalement pour f convexe mais aussi dans le cas non convexe.

Exemple de méthode proximale : ISTA Pour terminer cette partie, nous présentons une instance de l'algorithme 9 très populaire en traitement du signal, qui permet de calculer des représentations parcimonieuses. Cette méthode résout des problèmes lisses régularisés par une norme ℓ_1 , donc de la forme

$$\underset{\mathbf{w} \in \mathbb{R}^d}{\text{minimiser}} f(\mathbf{w}) + \lambda \|\mathbf{w}\|_1$$

avec f de classe C^1 . La forme du terme de régularisation permet de caractériser directement la solution du sous-problème (4.3.2). En effet, le sous-problème proximal donné par

$$\underset{\mathbf{w} \in \mathbb{R}^d}{\text{minimiser}} \left\{ f(\mathbf{w}_k) + \nabla f(\mathbf{w}_k)^T (\mathbf{w} - \mathbf{w}_k) + \frac{1}{2\alpha_k} \|\mathbf{w} - \mathbf{w}_k\|^2 + \lambda \|\mathbf{w}\|_1 \right\},$$

possède une unique solution. Pour l'obtenir, on calcule le pas classique de l'algorithme de descente de gradient $\mathbf{w}_k - \alpha_k \nabla f(\mathbf{w}_k)$, puis on applique une fonction de *seuillage faible* (ou **soft-thresholding**) notée $s_{\alpha_k \lambda}(\bullet)$ à chacune des composantes. Cette fonction est définie par

$$\forall \mu > 0, \forall t \in \mathbb{R}, \quad s_\mu(t) = \begin{cases} t + \mu & \text{si } t < -\mu \\ t - \mu & \text{si } t > \mu \\ 0 & \text{sinon.} \end{cases}$$

Ainsi, la solution du sous-problème proximal est définie composante à composante en fonction du pas de la descente de gradient. Ce résultat est au cœur de la méthode proximale pour les problèmes avec régularisation ℓ_1 , appelée ISTA (pour *Iterative Soft-Thresholding Algorithm*) : une description de cette méthode est donnée par l'algorithme 10.

Algorithme 10: ISTA: Iterative Soft-Thresholding Algorithm.

Initialisation : $\mathbf{w}_0 \in \mathbb{R}^d$.

Pour $k = 0, 1, \dots$

- 1. Calculer le gradient de la partie lisse du problème $\nabla f(\mathbf{w}_k)$.
- 2. Définir une taille de pas $\alpha_k > 0$.
- 3. Calculer le nouvel itéré \mathbf{w}_{k+1} composante par composante selon la formule :

$$[\mathbf{w}_{k+1}]_i = \begin{cases} [\mathbf{w}_k - \alpha_k \nabla f(\mathbf{w}_k)]_i + \alpha_k \lambda & \text{si } [\mathbf{w}_k - \alpha_k \nabla f(\mathbf{w}_k)]_i < -\alpha_k \lambda \\ [\mathbf{w}_k - \alpha_k \nabla f(\mathbf{w}_k)]_i - \alpha_k \lambda & \text{si } [\mathbf{w}_k - \alpha_k \nabla f(\mathbf{w}_k)]_i > \alpha_k \lambda \\ 0 & \text{si } [\mathbf{w}_k - \alpha_k \nabla f(\mathbf{w}_k)]_i \in [-\alpha_k \lambda, \alpha_k \lambda]. \end{cases} \quad (4.3.3)$$

Fin

La définition de la fonction de “soft-thresholding” force certaines des composantes du nouvel itéré à être nulles, ce qui produira au final une solution plus parcimonieuse que pour un problème non régularisé.

Remarque 4.3.3 *La méthode ISTA a été améliorée via l'introduction d'un terme de momentum, ce qui a conduit à un nouvel algorithme appelé FISTA (Fast ISTA) : il s'agit aujourd'hui de la variante d'ISTA la plus utilisée.*

4.4 Conclusion

Les fonctions non lisses (et plus particulièrement non dérивables) sont fréquentes en optimisation, et il peut s'avérer difficile de construire des algorithmes adaptés à leur minimisation. Dans certains cas, la structure de ces fonctions est connue, et des reformulations lisses peuvent être utilisées; dans d'autres situations, on peut avoir recours à des dérivées généralisées telles que les sous-gradients, qui permettent de généraliser en partie les stratégies basées sur les gradients en optimisation lisse.

Les aspects non lisses se retrouvent fréquemment dans les formes régularisées des problèmes d'optimisation. Le but d'une régularisation est de favoriser les modèles possédant certaines propriétés structurelles : cela est réalisé par l'ajout d'un terme dans la fonction objectif, qui est typiquement indépendant des données. Lorsque la fonction objectif originelle est lisse, un problème d'optimisation avec régularisation peut être traité au moyen d'un algorithme de gradient proximal : ce dernier procède par résolution successive de sous-problèmes impliquant la fonction de régularisation. Dans le cas important de la régularisation ℓ_1 , où l'on cherche à obtenir une solution parcimonieuse, l'algorithme de gradient proximal ISTA est la méthode la plus utilisée; pour des régularisations lisses, en revanche, il est possible d'appliquer des techniques d'optimisation lisses. C'est par exemple le cas avec la régularisation ℓ_2 , dont le but est de réduire la variance du modèle par rapport aux données.

Bibliographie

- [1] L. Bottou, F. E. Curtis, and J. Nocedal. Optimization Methods for Large-Scale Machine Learning. *SIAM Rev.*, 60:223–311, 2018.
- [2] S. Boyd and L. Vandenberghe. *Convex optimization*. Cambridge University Press, Cambridge, United Kingdom, 2004.
- [3] S. Boyd and L. Vandenberghe. *Introduction to Applied Linear Algebra - Vectors, Matrices and Least Squares*. Cambridge University Press, Cambridge, United Kingdom, 2018.
- [4] Yu. Nesterov. A method for solving convex optimization problems with convergence rate $\mathcal{O}(1/k^2)$. *Soviet Mathematics Doklady*, 27:372–376, 1983.
- [5] J. Nocedal and S. J. Wright. *Numerical Optimization*. Springer Series in Operations Research and Financial Engineering. Springer-Verlag, New York, second edition, 2006.
- [6] S. J. Wright. Optimization algorithms for data analysis. In A. C. Gilbert M. W. Mahoney, J. C. Duchi, editor, *The mathematics of data*, number 25 in IAS/Park City Mathematics Series. AMS, IAS/Park City Mathematics Institute, and Society for Industrial and Applied Mathematics, Princeton, 2018.

Annexe A

Notations et bases mathématiques

A.1 Notations

A.1.1 Conventions de notation générique

- Les scalaires seront représentés par des lettres minuscules : $a, b, c, \alpha, \beta, \gamma$.
- Les vecteurs seront représentés par des lettres minuscules **en gras** : $\mathbf{a}, \mathbf{b}, \mathbf{c}, \boldsymbol{\alpha}, \boldsymbol{\beta}, \boldsymbol{\gamma}$.
- Les lettres majuscules en gras seront utilisées pour les matrices : $\mathbf{A}, \mathbf{B}, \mathbf{C}$.
- Les lettres majuscules cursives seront utilisées pour les ensembles : $\mathcal{A}, \mathcal{B}, \mathcal{C}$.
- La définition d'une nouvelle quantité ou d'un nouvel opérateur sera indiquée par $:=$.
- On utilisera les quantificateurs \forall (pour tout), \exists (il existe), $\exists!$ (il existe un(e) unique) et \in (appartient à).
- Les opérateurs Σ et Π désigneront des sommes et des produits, respectivement. En l'absence d'ambiguïté, on pourra omettre les indices de début et de fin dans une somme ou un produit afin d'alléger les notations. On pourra de même utiliser un seul symbole de notation pour plusieurs indices et ainsi écrire de manière équivalente $\sum_{i=1}^m \sum_{j=1}^n$, $\sum_i \sum_j$ ou $\sum_{i,j}$ si le contexte le permet. Enfin, la notation $i = 1, \dots, m$ remplacera parfois les conditions $i \in \mathbb{N}$, $1 \leq i \leq m$.

A.1.2 Notations scalaires et vectorielles

- L'ensemble des entiers naturels sera noté \mathbb{N} , l'ensemble des entiers relatifs sera noté \mathbb{Z} .
- L'ensemble des réels sera noté \mathbb{R} . L'ensemble des réels positifs sera noté \mathbb{R}_+ et l'ensemble des réels strictement positifs sera noté \mathbb{R}_{++} .
- On notera \mathbb{R}^d l'ensemble des vecteurs à d composantes réelles, et on considérera toujours que d est un entier supérieur ou égal à 1.

- Un vecteur $\mathbf{w} \in \mathbb{R}^d$ sera pensé (par convention) comme un vecteur colonne. On notera $w_i \in \mathbb{R}$ sa i -ème coordonnée dans la base canonique de \mathbb{R}^d . On aura ainsi $\mathbf{w} = \begin{bmatrix} w_1 \\ \vdots \\ w_d \end{bmatrix}$, que l'on notera plus succinctement $\mathbf{w} = [w_i]_{1 \leq i \leq d}$.
- Étant donné un vecteur (colonne) $\mathbf{w} \in \mathbb{R}^d$, le vecteur ligne correspondant sera noté \mathbf{w}^T . On aura donc $\mathbf{w}^T = [w_1 \ \cdots \ w_n]$ et $[\mathbf{w}^T]^T = \mathbf{w}$.
- Pour tout $d \geq 1$, les vecteurs $\mathbf{0}_{\mathbb{R}^d}$ et $\mathbf{1}_{\mathbb{R}^d}$ représentent les vecteurs colonnes de \mathbb{R}^d dont tous les éléments sont égaux à 0 ou 1, respectivement. En fonction du contexte, on pourra noter simplement $\mathbf{0}$ ou $\mathbf{1}$.

A.1.3 Notations matricielles

- On notera $\mathbb{R}^{n \times d}$ l'ensemble des matrices à n lignes et d colonnes à coefficients réels, où n et d seront des entiers supérieurs ou égaux à 1. On pourra considérer un vecteur de \mathbb{R}^n comme une matrice de $\mathbb{R}^{n \times 1}$, et vice versa. Une matrice $\mathbf{A} \in \mathbb{R}^{n \times n}$ est dite carrée (dans le cas général, on parlera de matrice rectangulaire).
- Étant donnée une matrice $\mathbf{A} \in \mathbb{R}^{d \times d}$, on notera \mathbf{A}_{ij} le coefficient en ligne i et colonne j de la matrice. La diagonale de \mathbf{A} est formée par l'ensembles des coefficients \mathbf{A}_{ii} pour $i = 1, \dots, \min\{n, d\}$. La notation $[\mathbf{A}_{ij}]_{\substack{1 \leq i \leq n \\ 1 \leq j \leq d}}$ sera donc équivalente à \mathbf{A} . Sans ambiguïté sur la taille de la matrice, on notera simplement $[\mathbf{A}_{ij}]$.
- Selon les besoins, on utilisera \mathbf{a}_i^T pour la i -ème ligne de \mathbf{A} ou \mathbf{a}_j pour la j -ème colonne de \mathbf{A} . Selon le cas, on aura donc $\mathbf{A} = \begin{bmatrix} \mathbf{a}_1^T \\ \vdots \\ \mathbf{a}_n^T \end{bmatrix}$ ou $\mathbf{A} = [\mathbf{a}_1 \ \cdots \ \mathbf{a}_d]$.
- Pour une matrice $\mathbf{A} = [\mathbf{A}_{ij}] \in \mathbb{R}^{n \times d}$, la matrice transposée de \mathbf{A} , notée \mathbf{A}^T , est la matrice de $\mathbb{R}^{d \times n}$ telle que

$$\forall i = 1 \dots n, \ \forall j = 1 \dots d, \quad \mathbf{A}_{ji}^T = \mathbf{A}_{ij}.$$

Cette notation généralise donc la correspondance entre vecteurs lignes et vecteurs colonnes.

- Pour tout $d \geq 1$, la matrice \mathbf{I}_d représentera la matrice identité de $\mathbb{R}^{d \times d}$ (avec des 1 sur la diagonale et des 0 partout ailleurs), tandis que les matrices $\mathbf{0}_d$ et $\mathbf{1}_d$ représenteront les matrices dont tous les éléments sont égaux à 0 ou 1, respectivement. De manière plus générale, les notations $\mathbf{0}_{n,d}$ et $\mathbf{1}_{n,d}$ désignent les matrices de $\mathbb{R}^{n \times d}$ ne contenant respectivement que des 0 et des 1. Hors ambiguïté, on utilisera la notation $\mathbf{0}$ et $\mathbf{1}$ sans préciser les dimensions.

A.2 Éléments de mathématiques

Les fondements mathématiques de l'optimisation se trouvent dans l'analyse réelle, et en particulier dans le calcul différentiel. Les structures d'algèbre linéaire dans \mathbb{R}^d jouent également un rôle

prépondérant en optimisation (et en sciences des données). Cette section regroupe les résultats de base qui seront utilisés dans le cours.

Pour approfondir ces notions, on pourra consulter les liens suivants :

- Pour l'algèbre linéaire :
 - <https://www.ceremade.dauphine.fr/~carlier/polyalgebre.pdf> (en français);
 - <http://vmls-book.stanford.edu/vmls.pdf> (Chapitres 1 à 3, en anglais).
- Pour le calcul différentiel :
 - https://www.ceremade.dauphine.fr/~bouin/ens1819/Cours_Bolley.pdf (en français);
 - https://sebastianraschka.com/pdf/books/dlb/appendix_d_calculus.pdf (en anglais).

A.2.1 Algèbre linéaire vectorielle

On considérera toujours l'espace des vecteurs \mathbb{R}^d muni de sa structure d'espace vectoriel normé de dimension d . On définit donc les opérations suivantes :

- Pour tous $\mathbf{x}, \mathbf{y} \in \mathbb{R}^d$, la somme des vecteurs \mathbf{x} et \mathbf{y} est notée $\mathbf{x} + \mathbf{y} = [x_i + y_i]_{1 \leq i \leq d}$;
- Pour tout $\lambda \in \mathbb{R}$, on définit $\lambda \mathbf{x} \stackrel{n}{=} \lambda \cdot \mathbf{x} = [\lambda x_i]_{1 \leq i \leq d}$. Dans ce contexte, un tel réel λ sera appelé un *scalaire*.

A l'aide de ces opérations, nous pouvons donc construire des **combinaisons linéaires** de vecteurs de \mathbb{R}^d dont le résultat sera un vecteur de \mathbb{R}^d , à savoir des vecteurs de la forme $\sum_{i=1}^p \lambda_i \mathbf{x}_i$, où $\mathbf{x}_i \in \mathbb{R}^d$ et $\lambda_i \in \mathbb{R}$ pour tout $i = 1, \dots, p$.

Pour ce qui est de l'espace des matrices $\mathbb{R}^{n \times d}$, on peut également le munir d'une structure d'espace vectoriel normé de dimension nd . On définit donc les opérations suivantes :

- Pour tous $\mathbf{A}, \mathbf{B} \in \mathbb{R}^{n \times d}$, la somme des matrices \mathbf{A} et \mathbf{B} est notée $\mathbf{A} + \mathbf{B} = [\mathbf{A}_{ij} + \mathbf{B}_{ij}]_{\substack{1 \leq i \leq n \\ 1 \leq j \leq d}}$;
- Pour tout scalaire $\lambda \in \mathbb{R}$, on définit $\lambda \mathbf{A} \stackrel{n}{=} \lambda \cdot \mathbf{A} = [\lambda \mathbf{A}_{ij}]_{\substack{1 \leq i \leq n \\ 1 \leq j \leq d}}$.

Définition A.2.1 Un ensemble $\mathcal{S} \subseteq \mathbb{R}^d$ vérifiant les conditions

1. $\mathbf{0}_d \in \mathcal{S}$;
2. $\forall (\mathbf{x}, \mathbf{y}) \in \mathcal{S}, \mathbf{x} + \mathbf{y} \in \mathcal{S}$;
3. $\forall \mathbf{x} \in \mathcal{S}, \forall \lambda \in \mathbb{R}, \lambda \mathbf{x} \in \mathcal{S}$.

s'appelle un sous-espace vectoriel de \mathbb{R}^d .

Définition A.2.2 Soient $\mathbf{x}_1, \dots, \mathbf{x}_p$ p vecteurs de \mathbb{R}^d . Le **sous-espace engendré** par les vecteurs $\mathbf{x}_1, \dots, \mathbf{x}_p$, noté $\text{vect}(\mathbf{x}_1, \dots, \mathbf{x}_p)$, est le sous-espace vectoriel

$$\text{vect}(\mathbf{x}_1, \dots, \mathbf{x}_p) := \left\{ \mathbf{x} = \sum_{i=1}^p \alpha_i \mathbf{x}_i \middle| \alpha_i \in \mathbb{R} \ \forall i \right\}.$$

On rappelle ensuite les différentes propriétés notables des familles de vecteurs.

Définition A.2.3 • Une famille libre $\{\mathbf{x}_i\}_{i=1}^k$ de vecteurs de \mathbb{R}^n est telle que les vecteurs sont linéairement indépendants : pour tous scalaires $\lambda_1, \dots, \lambda_k$ tels que $\sum_{i=1}^k \lambda_i \mathbf{x}_i = 0$, alors $\lambda_1 = \dots = \lambda_k = 0$. On a nécessairement $k \leq n$.

- Une famille liée est une famille de vecteurs qui n'est pas libre.
- Une famille génératrice de vecteurs de \mathbb{R}^n est un ensemble de vecteurs $\{\mathbf{x}_i\}$ tel que le sous-espace engendré par ces vecteurs soit égal à \mathbb{R}^n .
- Une base de \mathbb{R}^n est une famille de vecteurs $\{\mathbf{x}_i\}_{i=1}^n$ qui est à la fois libre et génératrice. Tout vecteur de \mathbb{R}^n s'écrit alors de manière unique comme combinaison linéaire des \mathbf{x}_i . Toute base de \mathbb{R}^n comporte exactement n vecteurs.

Comme la taille d'une base de \mathbb{R}^n est n , on dit que cet espace vectoriel est de dimension n . En conséquence, la dimension d'un sous-espace vectoriel est au plus n .

Exemple A.2.1 Tout vecteur \mathbf{x} de \mathbb{R}^n s'écrit $\mathbf{x} = \sum_{i=1}^n x_i e_i$, où $e_i = [0 \cdots 0 1 0 \cdots 0]^T$ est le i -ème vecteur de la base canonique (le coefficient 1 se trouvant en i -ème position).

Norme et produit scalaire L'utilisation d'une norme, et du produit scalaire associé, permet de mesurer les distances entre vecteurs, ce qui sera particulièrement utile pour montrer la convergence d'une suite de points générés par un algorithme vers une solution d'un problème d'optimisation.

Définition A.2.4 La **norme euclidienne** $\|\cdot\|$ sur \mathbb{R}^n est définie pour tout vecteur $\mathbf{x} \in \mathbb{R}^n$ par

$$\|\mathbf{x}\| := \sqrt{\sum_{i=1}^n x_i^2}.$$

Remarque A.2.1 Il s'agit bien d'une norme, car elle vérifie les quatres axiomes d'une norme :

1. $\forall \mathbf{x}, \mathbf{y} \in \mathbb{R}^n, \|\mathbf{x} + \mathbf{y}\| \leq \|\mathbf{x}\| + \|\mathbf{y}\|;$
2. $\|\mathbf{x}\| = 0 \Leftrightarrow \mathbf{x} = \mathbf{0}_{\mathbb{R}^n};$
3. $\forall \mathbf{x}, \|\mathbf{x}\| \geq 0;$
4. $\forall \mathbf{x} \in \mathbb{R}^n, \forall \lambda \in \mathbb{R}, \|\lambda \mathbf{x}\| = |\lambda| \|\mathbf{x}\|.$

On dira qu'un vecteur $\mathbf{x} \in \mathbb{R}^n$ est unitaire si $\|\mathbf{x}\| = 1$.

Définition A.2.5 Pour tous vecteurs $\mathbf{x}, \mathbf{y} \in \mathbb{R}^n$, le **produit scalaire** dérivé de la norme euclidienne est une fonction de \mathbf{x} et \mathbf{y} , notée $\mathbf{x}^T \mathbf{y}$, et définie par

$$\mathbf{x}^T \mathbf{y} := \sum_{i=1}^n x_i y_i.$$

Deux vecteurs \mathbf{x} et \mathbf{y} tels que $\mathbf{x}^T \mathbf{y} = 0$ seront dits orthogonaux.

On notera en particulier que $\mathbf{y}^T \mathbf{x} = \mathbf{x}^T \mathbf{y}$. Ce produit scalaire définit donc un “produit” entre un vecteur ligne et un vecteur colonne.

Proposition A.2.1 Soient \mathbf{x} et \mathbf{y} deux vecteurs de \mathbb{R}^n . Alors, on a les propriétés suivantes :

- i) $\|\mathbf{x} + \mathbf{y}\|^2 = \|\mathbf{x}\|^2 + 2\mathbf{x}^T \mathbf{y} + \|\mathbf{y}\|^2$;
- ii) $\|\mathbf{x} - \mathbf{y}\|^2 = \|\mathbf{x}\|^2 - 2\mathbf{x}^T \mathbf{y} + \|\mathbf{y}\|^2$;
- iii) $\|\mathbf{x}\|^2 + \|\mathbf{y}\|^2 = \frac{1}{4} (\|\mathbf{x} + \mathbf{y}\|^2 + \|\mathbf{x} - \mathbf{y}\|^2)$;
- iv) **Inégalité de Cauchy-Schwarz** :

$$\forall \mathbf{x}, \mathbf{y} \in \mathbb{R}^n, \quad \mathbf{x}^T \mathbf{y} \leq \|\mathbf{x}\| \|\mathbf{y}\|.$$

Remarque A.2.2 La dernière inégalité est un résultat très important, qui s'utilise aussi bien en algèbre linéaire qu'en analyse fonctionnelle. Comme on le verra, elle apparaît également dans la preuve des inégalités de Taylor, qui sont fondamentales en optimisation.

A.2.2 Algèbre linéaire matricielle

Sur les espaces matriciels, on définit également un produit matriciel entre matrices de dimensions compatibles. Pour toutes matrices $\mathbf{A} \in \mathbb{R}^{m \times n}$ et $\mathbf{B} \in \mathbb{R}^{n \times p}$, la matrice produit \mathbf{AB} est définie comme la matrice $\mathbf{C} \in \mathbb{R}^{m \times p}$ telle que

$$\forall i = 1, \dots, m, \forall j = 1, \dots, p, \quad C_{ij} = \sum_{k=1}^n \mathbf{A}_{ik} \mathbf{B}_{kj}.$$

Par analogie, le produit d'une matrice $\mathbf{A} \in \mathbb{R}^{m \times n}$ avec un vecteur $\mathbf{x} \in \mathbb{R}^n$ sera le vecteur $\mathbf{y} \in \mathbb{R}^m$ donné par

$$\forall i = 1, \dots, m, \quad y_i = \sum_{j=1}^n \mathbf{A}_{ij} x_j.$$

Remarque A.2.3 On notera que la notation du produit scalaire sur \mathbb{R}^n correspond au produit de deux matrices de taille $1 \times n$ et $n \times 1$, dont le résultat est une matrice 1×1 , c'est-à-dire un scalaire.

Lorsque l'on travaille avec des matrices, on s'intéresse généralement aux sous-espaces définis ci-dessous.

Définition A.2.6 (Sous-espaces matriciels) Soit une matrice $\mathbf{A} \in \mathbb{R}^{m \times n}$, on définit les deux sous-espaces suivants :

- Le noyau de \mathbf{A} est le sous-espace vectoriel

$$\ker(\mathbf{A}) := \{\mathbf{x} \in \mathbb{R}^n \mid \mathbf{Ax} = \mathbf{0}_m\}$$

- L'image de \mathbf{A} est le sous-espace vectoriel

$$\text{Im}(\mathbf{A}) := \{\mathbf{y} \in \mathbb{R}^m \mid \exists \mathbf{x} \in \mathbb{R}^n, \mathbf{y} = \mathbf{Ax}\}$$

La dimension de ce sous-espace vectoriel s'appelle le **rang** de \mathbf{A} . On la note $\text{rang}(\mathbf{A})$ et on a $\text{rang}(\mathbf{A}) \leq \min\{m, n\}$.

Théorème A.2.1 (Théorème du rang) Pour toute matrice $A \in \mathbb{R}^{m \times n}$, on a

$$\dim(\ker(A)) + \text{rang}(A) = n.$$

Définition A.2.7 (Normes matricielles) On définit sur $\mathbb{R}^{m \times n}$ la norme d'opérateur $\|\cdot\|$ et la norme de Frobenius $\|\cdot\|_F$ par

$$\forall A \in \mathbb{R}^{m \times n}, \begin{cases} \|A\| &:= \max_{\substack{x \in \mathbb{R}^n \\ x \neq 0_n}} \frac{\|Ax\|}{\|x\|} = \max_{\substack{x \in \mathbb{R}^n \\ \|x\|=1}} \|Ax\| \\ \|A\|_F &:= \sqrt{\sum_{\substack{1 \leq i \leq m \\ 1 \leq j \leq n}} A_{ij}^2}. \end{cases}$$

Définition A.2.8 (Matrice symétrique) Une matrice carrée $A \in \mathbb{R}^{n \times n}$ est dite symétrique si elle vérifie $A^T = A$.

L'ensemble des matrices symétriques de taille $n \times n$ sera noté \mathcal{S}^n .

Définition A.2.9 (Matrice inversible) Une matrice carrée $A \in \mathbb{R}^{n \times n}$ est dite inversible s'il existe $B \in \mathbb{R}^{n \times n}$ telle que $BA = AB = I_n$ (où l'on rappelle que I_n désigne la matrice identité de $\mathbb{R}^{n \times n}$).

Si elle existe, une telle matrice B est unique : elle est appelée **l'inverse de A** et on la note A^{-1} .

Définition A.2.10 (Matrice (semi-)définie positive) Une matrice carrée $A \in \mathbb{R}^{n \times n}$ symétrique est dite **semi-définie positive** si

$$\forall x \in \mathbb{R}^n, \quad x^T Ax \geq 0,$$

ce que l'on notera $A \succeq 0$.

Une telle matrice est dite **définie positive** lorsque $x^T Ax > 0$ pour tout vecteur x non nul, ce que l'on notera $A \succ 0$.

Définition A.2.11 (Matrice orthogonale) Une matrice carrée $P \in \mathbb{R}^{n \times n}$ est dite **orthogonale** si $P^T = P^{-1}$.

Par extension, on dira que $Q \in \mathbb{R}^{m \times n}$ avec $m \leq n$ est orthogonale si $QQ^T = I_m$ (les colonnes de Q sont donc orthonormées dans \mathbb{R}^m).

On notera que lorsque $Q \in \mathbb{R}^{n \times n}$ est une matrice orthogonale, alors sa transposée Q^T est également orthogonale (ce résultat n'est pas vrai pour une matrice rectangulaire). On utilisera fréquemment la propriété des matrices orthogonales énoncée ci-dessous.

Lemme A.2.1 Soit une matrice $A \in \mathbb{R}^{m \times n}$ et $U \in \mathbb{R}^{m \times m}$, $V \in \mathbb{R}^{n \times n}$ des matrices orthogonales, respectivement de $\mathbb{R}^{m \times m}$ et $\mathbb{R}^{n \times n}$. On a

$$\|A\| = \|UA\| = \|AV\| \quad \text{et} \quad \|A\|_F = \|UA\|_F = \|AV\|_F,$$

c'est-à-dire que la multiplication par une matrice orthogonale ne modifie pas la norme d'opérateur.

Par corollaire immédiat du lemme précédent, on note qu'une matrice $Q \in \mathbb{R}^{m \times n}$ orthogonale avec $m \leq n$ vérifie nécessairement $\|Q\| = 1$ et $\|Q\|_F = \sqrt{m}$.

Définition A.2.12 (Valeur propre) Soit une matrice $A \in \mathbb{R}^{n \times n}$. On dit que $\lambda \in \mathbb{R}$ est une **valeur propre de A** si

$$\exists v \in \mathbb{R}^n, v \neq 0_n, \quad Av = \lambda v.$$

Le vecteur v est appelé un **vecteur propre associé à la valeur propre λ** . L'ensemble des valeurs propres de A s'appelle le **spectre de A** .

Le sous-espace engendré par les vecteurs propres associés à la même valeur propre d'une matrice s'appelle un **sous-espace propre**. Sa dimension correspond à l'ordre de multiplicité de la valeur propre relativement à la matrice.

Proposition A.2.2 Pour toute matrice $A \in \mathbb{R}^{n \times n}$, on a les propriétés suivantes :

- La matrice A possède n valeurs propres complexes mais pas nécessairement réelles.
- Si la matrice A est semi-définie positive (respectivement définie positive), alors ses valeurs propres sont réelles positives (respectivement strictement positives).
- Le noyau de A est engendré par les vecteurs propres associés à la valeur propre 0.

Théorème A.2.2 (Théorème spectral) Toute matrice carrée $A \in \mathbb{R}^{n \times n}$ symétrique admet une décomposition dite **spectrale** de la forme :

$$A = P \Lambda P^{-1},$$

où $P \in \mathbb{R}^{n \times n}$ est une matrice orthogonale, dont les colonnes p_1, \dots, p_n forment une base orthonormée de vecteurs propres, et $\Lambda \in \mathbb{R}^{n \times n}$ est une matrice diagonale qui contient les n valeurs propres de A $\lambda_1, \dots, \lambda_n$ sur la diagonale.

Il est à noter que la décomposition spectrale n'est pas unique. En revanche, l'ensemble des valeurs propres est unique, que l'on prenne en compte les ordres de multiplicité ou non.

Géométriquement parlant, on voit ainsi que, pour tout vecteur $x \in \mathbb{R}^n$ décomposé dans la base des p_i que l'on multiplie par A , les composantes de ce vecteur associées aux plus grandes valeurs propres¹ seront augmentées, tandis que celles associées aux valeurs propres de petite magnitude seront réduites (voire annihilées dans le cas d'une valeur propre nulle).

Lien avec la décomposition en valeurs singulières Soit une matrice rectangulaire $A \in \mathbb{R}^{m \times n}$: dans le cas général, les dimensions de la matrice diffèrent, et on ne peut donc pas parler de valeurs propres de la matrice A . On peut en revanche considérer les deux matrices

$$A^T A \in \mathbb{R}^{n \times n} \quad \text{et} \quad A A^T \in \mathbb{R}^{m \times m}.$$

Ces matrices sont symétriques réelles, et par conséquent diagonalisables. On peut se baser sur cette propriété pour obtenir une décomposition en valeurs singulières de A (ou *SVD* d'après l'acronyme anglais).

¹On parle ici de plus grandes valeurs propres en valeur absolue, ou magnitude.

A.2.3 Calcul différentiel

Définition A.2.13 (Continuité) Une fonction $f : \mathbb{R}^n \rightarrow \mathbb{R}^m$ est dite **continue en $\mathbf{x} \in \mathbb{R}^n$** si

$$\forall \epsilon > 0, \exists \delta > 0, \forall \mathbf{y} \in \mathbb{R}^n, \|\mathbf{y} - \mathbf{x}\| < \delta \Rightarrow \|f(\mathbf{y}) - f(\mathbf{x})\| < \epsilon.$$

f est **continue** sur un ensemble $\mathcal{A} \subseteq \mathbb{R}^n$ si elle l'est en tout point de \mathcal{A} . On dira en particulier que f est continue si elle est continue sur \mathbb{R}^n .

On parlera de fonction continue même s'il s'agit en réalité de fonction uniformément continue.

Une caractérisation de la continuité au moyen des suites (dite caractérisation séquentielle) est donnée ci-dessous. Elle sera particulièrement importante dans le cadre des convergences d'algorithmes.

Définition A.2.14 (Continuité (définition séquentielle)) Une fonction $f : \mathbb{R}^n \rightarrow \mathbb{R}^m$ est continue en $\mathbf{x} \in \mathbb{R}^n$ si

$$\forall \{\mathbf{x}_n\} \in (\mathbb{R}^n)^{\mathbb{N}}, \{\mathbf{x}_n\} \rightarrow \mathbf{x}, \lim_{n \rightarrow \infty} f(\mathbf{x}_n) = f(\mathbf{x}).$$

Exemple A.2.2 Une application linéaire $f : \mathbb{R}^n \rightarrow \mathbb{R}^m$, où $f(\mathbf{x}) = \mathbf{A}\mathbf{x}$ pour tout $\mathbf{x} \in \mathbb{R}^n$ et $\mathbf{A} \in \mathbb{R}^{m \times n}$, est une fonction continue sur \mathbb{R}^n .

Définition A.2.15 (Matrice jacobienne) $f : \mathbb{R}^n \rightarrow \mathbb{R}^m$ est dérivable en $\mathbf{x} \in \mathbb{R}^n$ si il existe une matrice $\mathbf{J}_f(\mathbf{x}) \in \mathbb{R}^{m \times n}$ telle que

$$\lim_{\substack{\mathbf{z} \rightarrow \mathbf{x} \\ \mathbf{z} \neq \mathbf{x}}} \frac{\|f(\mathbf{z}) - f(\mathbf{x}) - \mathbf{J}_f(\mathbf{x})(\mathbf{z} - \mathbf{x})\|}{\|\mathbf{z} - \mathbf{x}\|} = 0.$$

- $\mathbf{J}_f(\mathbf{x})$ s'appelle la (matrice) **jacobienne** de f en \mathbf{x} , et est définie de manière unique;
- Si $f(\cdot) = (f_1(\cdot), \dots, f_m(\cdot))^T$, alors

$$\forall 1 \leq i \leq m, \forall 1 \leq j \leq n, [\mathbf{J}_f(\mathbf{x})]_{ij} = \frac{\partial f_i}{\partial x_j}(\mathbf{x}).$$

Remarque A.2.4 Deux cas particuliers sont à connaître :

- Lorsque $m = 1$, on définit le vecteur $\nabla f(\mathbf{x}) \equiv \mathbf{J}_f(\mathbf{x})^T$, que l'on appelle le **vecteur gradient**;
- Lorsqu'on a $n = m = 1$, la matrice Jacobienne et le vecteur gradient sont équivalents à un scalaire $f'(\mathbf{x}) \equiv \nabla f(\mathbf{x}) \equiv \mathbf{J}_f(\mathbf{x})^T$, que l'on appelle la dérivée de f en \mathbf{x} .

Corollaire A.2.1 Pour toute fonction $f : \mathbb{R}^n \rightarrow \mathbb{R}^m$ définie par $f(\mathbf{x}) = (f_1(\mathbf{x}), \dots, f_m(\mathbf{x}))^T$ et dérivable en $\mathbf{x} \in \mathbb{R}^n$, la matrice Jacobienne est donnée par les dérivées partielles des f_i , c'est-à-dire :

$$\forall i = 1, \dots, m, \forall j = 1, \dots, n, [J_f(\mathbf{x})]_{ij} = \frac{\partial f_i}{\partial x_j}(\mathbf{x}).$$

Dans le cas $m = 1$, on obtient ainsi que le gradient est le vecteur des dérivées partielles de f :

$$\forall i = 1, \dots, n, \nabla f(\mathbf{x}) = \left[\frac{\partial f}{\partial x_i}(\mathbf{x}) \right]_{1 \leq i \leq n}.$$

On supposera dans ces notes que les dérivées de fonctions usuelles de \mathbb{R} dans \mathbb{R} sont connues. Afin de combiner ces formules pour calculer des dérivées de fonctions plus complexes, on se basera sur la règle ci-dessous.

Théorème A.2.3 (Dérivation de fonctions composées) *Si $f : \mathbb{R}^n \mapsto \mathbb{R}^m$ et $g : \mathbb{R}^m \mapsto \mathbb{R}^p$ sont dérивables sur leurs espaces respectifs, alors $h : \mathbb{R}^n \mapsto \mathbb{R}^p$ est dérivable sur \mathbb{R}^n et*

$$\forall \mathbf{x} \in \mathbb{R}^n, \quad \mathbf{J}_h(\mathbf{x}) = \mathbf{J}_g(f(\mathbf{x}))\mathbf{J}_f(\mathbf{x}).$$

Remarque A.2.5 *La règle ci-dessus donne les cas particuliers suivants :*

- $m = p = 1 : \nabla h(\mathbf{x}) = g'(f(\mathbf{x}))\nabla f(\mathbf{x});$
- $n = m = p = 1 : h'(x) = g'(f(x))f'(x).$

Théorème A.2.4 (Théorème des accroissements finis en dimension 1) *Soit $f : [a, b] \rightarrow \mathbb{R}$. Si f est continue sur $[a, b]$ et dérivable sur (a, b) , il existe $c \in (a, b)$ tel que*

$$\frac{f(b) - f(a)}{b - a} = f'(c).$$

Définition A.2.16 (Développements de Taylor) *Si $f : [a, b] \mapsto \mathbb{R}$ est \mathcal{C}^1 sur $[a, b]$, alors*

$$\begin{aligned} f(b) &= f(a) + f'(c)(b - a) \quad \text{où } c \in [a, b] \\ f(b) &= f(a) + \int_0^1 f'(a + t(b - a))(b - a) dt. \end{aligned}$$

Théorème A.2.5 (Accroissements finis en dimension d) *Soit $f : \mathbb{R}^d \rightarrow \mathbb{R}$ avec $f \in \mathcal{C}^1(\mathbb{R}^d)$. Pour tous $\mathbf{x}, \mathbf{y} \in \mathbb{R}^d$, $\mathbf{x} \neq \mathbf{y}$, il existe $t \in (0, 1)$ tel que*

$$f(\mathbf{y}) = f(\mathbf{x}) + \nabla f(\mathbf{x} + t(\mathbf{y} - \mathbf{x}))^\top (\mathbf{y} - \mathbf{x}).$$

Définition A.2.17 (Fonction lipschitzienne) *Une fonction $f : \mathbb{R}^n \rightarrow \mathbb{R}^m$ est L -lipschitzienne sur $\mathcal{A} \subset \mathbb{R}^n$ si*

$$\forall \mathbf{x}, \mathbf{y} \in \mathcal{A}^2, \quad \|f(\mathbf{x}) - f(\mathbf{y})\| \leq L\|\mathbf{x} - \mathbf{y}\|.$$

Proposition A.2.3 *Toute fonction lipschitzienne sur un ensemble est continue sur ce même ensemble.*

Définition A.2.18 (Classes de fonctions) *Soit une fonction $f : \mathbb{R}^d \rightarrow \mathbb{R}$.*

- *On dira que la fonction f est de classe $\mathcal{C}^p(\mathbb{R}^d)$ (ou simplement \mathcal{C}^p) si elle est p -fois dérivable et que sa dérivée p -ième est une fonction continue. On parlera de classe \mathcal{C}^∞ si la fonction est dérivable une infinité de fois.*
- *On dira que la fonction f est de classe $\mathcal{C}_L^{p,p}(\mathbb{R}^d)$ (ou simplement $\mathcal{C}_L^{p,p}$) si elle est p -fois dérivable et que sa dérivée p -ième est une fonction L -Lipschitzienne.*

Théorème A.2.6 (Développement de Taylor à l'ordre 1) Soit $f \in \mathcal{C}^1(\mathbb{R}^d)$ pour tous vecteurs \mathbf{x} et \mathbf{y} de \mathbb{R}^d , on a :

$$f(\mathbf{y}) = f(\mathbf{x}) + \int_0^1 \nabla f(\mathbf{x} + t(\mathbf{y} - \mathbf{x}))^\text{T} (\mathbf{y} - \mathbf{x}) dt.$$

Si de plus $f \in \mathcal{C}_L^{1,1}(\mathbb{R}^d)$, alors on a :

$$f(\mathbf{y}) \leq f(\mathbf{x}) + \nabla f(\mathbf{x})^\text{T} (\mathbf{y} - \mathbf{x}) + \frac{L}{2} \|\mathbf{y} - \mathbf{x}\|^2. \quad (\text{A.2.1})$$

Théorème A.2.7 (Développement de Taylor à l'ordre 2) Soit $f \in \mathcal{C}^2(\mathbb{R}^d)$; pour tous vecteurs \mathbf{x} et \mathbf{y} de \mathbb{R}^d , on a :

$$f(\mathbf{y}) = f(\mathbf{x}) + \nabla f(\mathbf{x})^\text{T} (\mathbf{y} - \mathbf{x}) + \frac{1}{2} \int_0^1 (\mathbf{y} - \mathbf{x})^\text{T} \nabla^2 f(\mathbf{x} + t(\mathbf{y} - \mathbf{x})) (\mathbf{y} - \mathbf{x}) dt.$$

Si de plus $f \in \mathcal{C}_L^{2,2}(\mathbb{R}^d)$, alors on a :

$$f(\mathbf{y}) \leq f(\mathbf{x}) + \nabla f(\mathbf{x})^\text{T} (\mathbf{y} - \mathbf{x}) + \frac{1}{2} (\mathbf{y} - \mathbf{x})^\text{T} \nabla^2 f(\mathbf{x}) (\mathbf{y} - \mathbf{x}) + \frac{L}{2} \|\mathbf{y} - \mathbf{x}\|^3. \quad (\text{A.2.2})$$