



HAL
open science

Matrices efficaces pour le traitement du signal et l'apprentissage automatique

Luc Le Magoarou

► **To cite this version:**

Luc Le Magoarou. Matrices efficaces pour le traitement du signal et l'apprentissage automatique. Traitement du signal et de l'image [eess.SP]. INSA de Rennes, 2016. Français. NNT : 2016ISAR0008 . tel-01412558

HAL Id: tel-01412558

<https://theses.hal.science/tel-01412558>

Submitted on 8 Dec 2016

HAL is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

Thèse

UNIVERSITE
BRETAGNE
LOIRE

THESE INSA Rennes
sous le sceau de Université Bretagne Loire
pour obtenir le titre de
DOCTEUR DE L'INSA DE RENNES
Spécialité : Traitement du signal et de l'image

présentée par
Luc Le Magoarou

ECOLE DOCTORALE : MATISSE
LABORATOIRE : Inria Rennes-Bretagne Atlantique

Matrices efficaces pour le traitement du signal et l'apprentissage automatique

Thèse soutenue le 24.11.2016
devant le jury composé de :

Pierre Vandergheynst

Professeur, École Polytechnique Fédérale de Lausanne / *président*

François Malgouyres

Professeur, Université Paul Sabatier, Toulouse / *rapporteur*

Gabriel Peyré

Directeur de recherche CNRS, DMA, École Normale Supérieure / *rapporteur*

Christine Guillemot

Directeur de recherche Inria, Rennes / *examineur*

Karin Schnass

University assistant, University of Innsbruck / *examineur*

Rémi Gribonval

Directeur de recherche Inria, Rennes / *Directeur de thèse*

Matrices efficientes pour
le traitement du signal
et l'apprentissage
automatique

Luc Le Magoarou



En partenariat avec



Remerciements

Bien que j'en sois l'auteur, je n'aurais pas la prétention de m'attribuer tous les mérites des travaux effectués lors de cette thèse. Ainsi, je me dois d'être reconnaissant envers un certain nombre de personnes, cette page est l'occasion pour moi de leur exprimer ma gratitude.

En premier lieu, je remercie Rémi de m'avoir proposé un sujet si vaste et intéressant, m'avoir fait confiance et permis de réaliser cette thèse dans des conditions exceptionnelles.

En outre, je remercie François Malgouyres et Gabriel Peyré d'avoir accepté d'être rapporteurs de ce manuscrit, pour l'intérêt qu'ils ont porté à mon travail, et pour leurs remarques toujours pertinentes. De même, je remercie Christine Guillemot, Karin Schnass et Pierre Vandergheynst d'avoir accepté de faire partie du jury en tant qu'examineurs.

Je remercie par ailleurs tous mes collègues de l'équipe PANAMA d'avoir partagé un bout de chemin avec moi, et leur souhaite toute la réussite qu'ils méritent.

Je souhaiterais exprimer ma gratitude à ma famille et particulièrement à mes parents, pour avoir fait en sorte que j'aie toujours eu le choix, et m'avoir soutenu tout au long de mes études. Je mesure aujourd'hui ma chance.

Ensuite, comme il n'y a pas que le travail dans la vie, je remercie mes amis (en particulier les Rennais et tous "Les copains") et mes coéquipiers de l'ASCHB qui ont fait de ces trois années un réel plaisir.

Enfin, je ne peux pas conclure cette énumération sans remercier du plus profond de mon cœur Fanny, celle qui a témoigné le soutien le plus inconditionnel et inébranlable tout au long de cette thèse, sans qui rien n'aurait la même saveur.

Luc Le Magoarou

Table des matières

Introduction	1
I État de l'art et matrices efficientes	9
1 Optimisation	11
1.1 Problème d'optimisation	11
1.2 Convexité	14
1.3 Algorithmes	16
2 Factorisation matricielle	23
2.1 Matrices et linéarité	23
2.2 Factorisations usuelles	26
2.3 Apprentissage de dictionnaire	30
2.4 Réseaux de neurones	35
3 Matrices efficientes	41
3.1 Circuit linéaire et complexité	41
3.2 Efficience	46
3.3 Calculer et borner l'efficience	47
II Contributions	51
4 Approximation par matrices efficientes	53
4.1 Objectif général	53
4.2 Travaux liés	56
4.3 Algorithme d'approximation	61
4.4 Stratégie hiérarchique	65
5 Application à la résolution de problèmes inverses	71
5.1 Problèmes inverses	71

5.2	Application à la localisation de sources	74
6	Application à la transformée de Fourier rapide sur graphe	83
6.1	Traitement du signal sur graphe	84
6.2	FFT sur graphe par factorisation creuse multi-couche de la matrice de Fourier	88
6.3	FFT sur graphe par diagonalisation gloutonne du Laplacien	90
6.4	Comparaison des méthodes	94
7	Identifiabilité de la forme efficiente d'une matrice	103
7.1	Introduction	103
7.2	Énoncé du problème étudié	105
7.3	Résultat d'identifiabilité	107
7.4	Exemple : la FFT à partir de la matrice de Fourier	111
8	Apprentissage de matrices efficientes	115
8.1	Formulation du problème d'apprentissage	116
8.2	Application à l'apprentissage de dictionnaire	117
8.3	Vers une application aux nouveaux modèles parcimonieux	127
	Conclusion et perspectives	133
III	Annexes	139
A	Annexe du chapitre 4	141
B	Annexe du chapitre 5	145
C	Annexe du chapitre 6	147
D	Annexe du chapitre 8	151
	Bibliographie	163

Introduction

Des données de plus en plus volumineuses, sous des formes de plus en plus variées, collectées de plus en plus rapidement. C'est un véritable déluge de données numériques brutes qui s'abat sur notre monde hyperconnecté depuis une décennie. Quelques chiffres et un peu d'histoire illustreront mieux qu'un long discours ce propos. L'information numérique est représentée par une suite de symboles, pouvant prendre les valeurs "0" ou "1" dans le cas binaire, on parle alors de bit pour *binary unit*, et une suite de 8 bits forme un octet. Elle peut être stockée de différentes manières. Inventée Au XVIII^e siècle pour programmer des métiers à tisser, la carte perforée permettait, en fonction de la présence ou non de trous sur une bande de papier (nous donnons un exemple en figure 1), de stocker jusqu'à 960 bits (120 octets). Dans les années 1970, la disquette utilise une bande magnétique et permet de stocker jusqu'à 1,5 Mégaoctets (Mo), ce qui correspond à la capacité de 12500 cartes perforées. Aujourd'hui, certaines clefs USB permettent de stocker jusqu'à 256 Gigaoctets (Go), soit la capacité de 170000 disquettes ou de plus de 2 milliards de cartes perforées.

Cet accroissement exponentiel de la quantité de données qu'il est possible de stocker (et par conséquent de la quantité de données effectivement stockée) est encore plus impressionnant depuis l'apparition d'internet et sa popularisation à la fin des années 1990. En effet, on estime que la quantité de données créée tous les deux jours aujourd'hui égale celle créée au cours de toute l'histoire de l'humanité jusqu'en 2003 ([Demarthon et al., 2013](#)). La source des données produites est très variable. On distingue notamment l'utilisation quotidienne de services numériques par des particuliers (e-mails, réseaux sociaux, plateformes vidéo etc.), par exemple, l'équivalent de 20 heures de nouvelles vidéos sont publiées chaque minute sur YouTube, et presque 3 millions d'e-mails sont envoyés chaque seconde, quelques chiffres sur le sujet sont présentés en figure 2. La conduite de projets scientifiques de grande

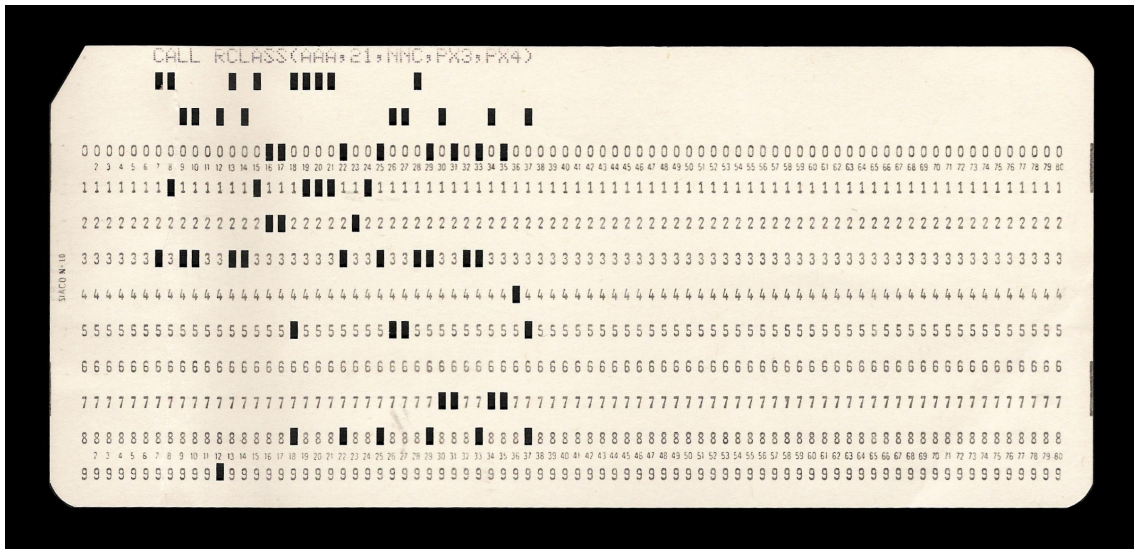


Figure 1 : Une carte perforée à 80 colonnes (source : Wikipédia). La carte perforée est un des premiers dispositifs de stockage de l’information numérique.

envergure génère elle aussi d’impressionnants volumes de données. En physique des particules, le Grand collisionneur de hadrons (LHC pour “Large Hadron Collider”), un accélérateur de particules situé à Genève, génère 25 Pétaoctets (Po, 10^{15} octets) de données brutes par an. En astronomie, le Réseau d’un kilomètre carré (SKA pour “Square Kilometer Array”), un radiotélescope actuellement en construction en Afrique du Sud et en Australie, générera 1 Exaoctet (Eo, 10^{18} octets) de données brutes par jour.

Un tel afflux massif et parfois incontrôlé de données crée des enjeux économiques, éthiques et scientifiques. *Comment exploiter cette ressource en perpétuelle croissance ?* Cette question est au centre des problématiques de l’industrie des nouvelles technologies, qui cherche à tirer profit des données massives qu’elle collecte (le terme “Big Data” est souvent utilisé dans l’industrie pour désigner ces données). *Comment préserver sa vie privée dans un monde hyperconnecté ?* Les données personnelles collectées par les entreprises fournissant des services sur internet ne sont pas forcément protégées. Le programme PRISM de l’Agence de Sécurité Nationale américaine (abrégée NSA pour “National Security Agency”) constitue un programme de surveillance électronique à grande échelle conduit avec la complicité de certains des plus grands acteurs industriels du secteur des nouvelles technologies. Il a été révélé au grand public par Edward Snowden en 2013, et est un exemple criant d’atteinte à la vie privée. *Comment traiter les données collectées afin d’en extraire de l’information utile ?* Cette dernière question, qui nous intéresse plus particulièrement dans cette thèse, est un véritable défi posé aux scientifiques du domaine des Sciences et

2016 What happens in an INTERNET MINUTE?

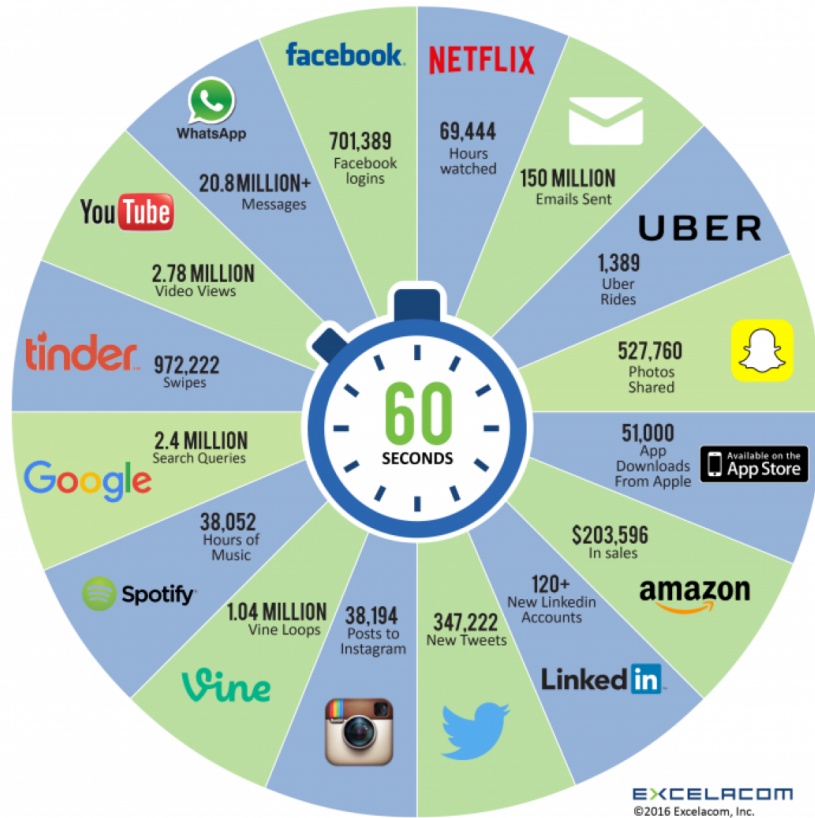


Figure 2 : Quelques chiffres sur le volume de données produites chaque minute par plusieurs services internet populaires (source : Excelacom, Inc.).

Technologies de l'Information et de la Communication (STIC). En effet, le volume sans précédent de données disponibles impose un nouveau paradigme. L'accent n'est plus forcément mis sur l'exactitude ou l'optimalité des traitements mais plutôt sur leur rapidité et leur adaptabilité.

Le traitement de données, ou science des données (qu'on peut traduire "data science") est une discipline très générale reposant sur plusieurs domaines de recherche liés aux mathématiques appliquées. L'objectif général de la science des données est d'extraire de la connaissance de données brutes. Cette thèse considère plus particulièrement deux domaines : le *traitement du signal* et l'*apprentissage automatique*, qui sont des sous-ensembles de la science des données. Ces deux disciplines se recouvrent partiellement et il est parfois difficile de les distinguer. Le traitement

du signal trouve son origine dans le domaine de l'électronique et vise à découvrir l'information contenue dans les données brutes (le signal). Les tâches les plus courantes accomplies en traitement du signal sont : le débruitage, la compression, le filtrage, le contrôle. Le traitement du signal trouve son application dans de nombreux domaines, citons par exemple les télécommunications, le traitement du son et des images, l'imagerie médicale. L'apprentissage automatique (traduit "machine learning"), quant à lui, trouve son origine dans le domaine de l'informatique et vise à analyser les données brutes (appelées données d'entraînement) afin d'apprendre un modèle, dans le but de prédire le comportement de données futures. L'apprentissage automatique s'assimile le plus souvent à de la classification de données ou de la régression (prédiction). Il s'applique notamment à la détection de fraude, la reconnaissance de la parole et de l'écriture manuscrite, l'indexation d'images et de vidéos.

Les masses de données générées à l'heure actuelle nécessitent des traitements efficaces. L'efficacité est définie en économie comme le rapport entre le résultat obtenu et les ressources utilisées pour parvenir à ce résultat. L'efficacité ne doit pas être confondue avec l'efficacités qui mesure simplement le rapport entre le résultat obtenu et le résultat espéré sans considérer les ressources utilisés. Dans le cadre du traitement de données, l'efficacité peut se définir comme le rapport entre la quantité de données traitée (mesurée en bits) et les ressources computationnelles utilisées (mesurées en nombre d'opérations arithmétiques). Un traitement de données efficace utilise un nombre réduit d'opérations arithmétiques pour manipuler une certaine quantité de données.

Dans cette thèse, on se propose de considérer un élément omniprésent dans les traitements de données numériques : les transformations linéaires. La linéarité est un concept mathématique généralisant la notion de proportionnalité à des objets plus complexes que les nombres : les vecteurs. Les transformations linéaires se composent uniquement d'additions et de multiplications par des constantes, elles sont en ce sens les plus "simples" des transformations qui s'appliquent à des vecteurs. Un nombre considérable de transformations usuelles sont linéaires, on peut citer les transformations géométriques simples (rotation, homothétie, symétrie), les transformations usuelles du traitement du signal (transformée de Fourier, en ondelettes, en cosinus discrète), la dérivée et l'intégrale discrètes. Les matrices, tableaux de nombres appelés coefficients, sont la représentation naturelle des transformations linéaires en dimension finie. En effet l'application d'une transformation linéaire d'un espace de dimension n vers un espace de dimension m est toujours équivalente à la multiplication par une matrice à m lignes et n colonnes. Les matrices jouent donc un rôle central dans la plupart des méthodes de traitement de données. La multiplication par une matrice à m lignes et n colonnes nécessite de l'ordre de $m \times n$ opérations arithmétiques. Le nombre $m \times n$ peut être très grand si les données traitées sont

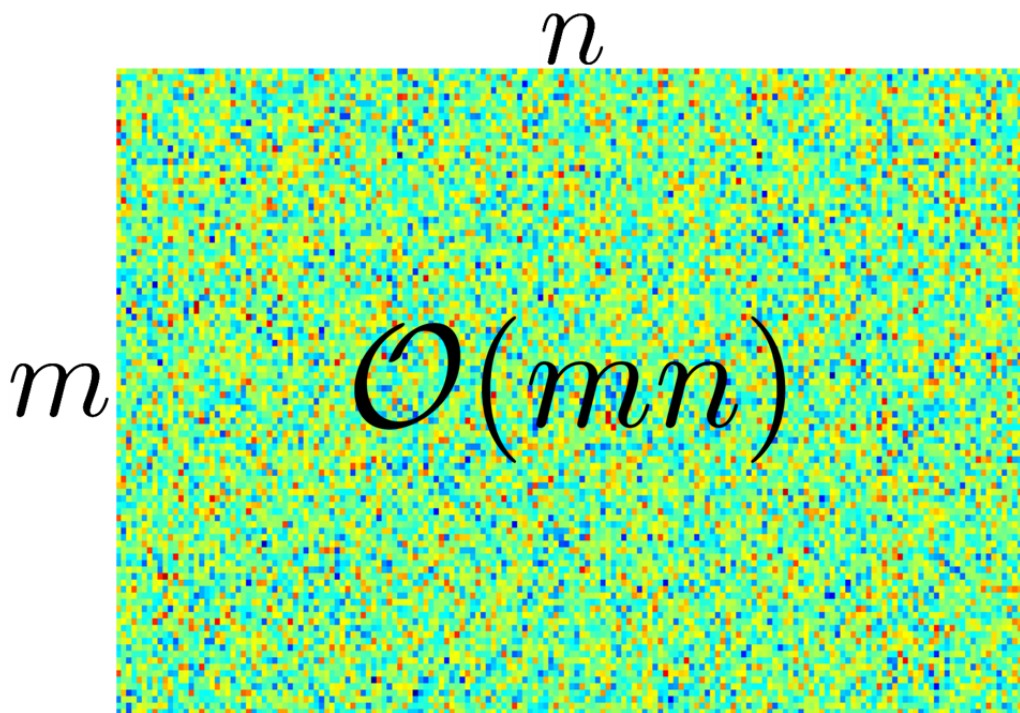


Figure 3 : Représentation d'une matrice à m lignes et n colonnes. Chaque case correspond à un coefficient, la couleur de la case correspond à la valeur du coefficient. Une telle matrice nécessite de l'ordre de mn opérations arithmétiques pour multiplier un vecteur.

volumineuses, ce qui peut limiter l'efficacité des méthodes de traitement de données. Ce constat nous mène à l'objectif général de cette thèse :

Objectif général de la thèse

Proposer une structure de matrice efficace permettant des traitements de données plus efficaces.

Organisation et contributions de la thèse. L'objectif de cette thèse tel que formulé ci-dessus est très général. Nous proposons de structurer notre réflexion, ainsi que ce document, autour des trois questions suivantes :

- *Comment définir la notion d'efficacité d'une matrice ?*

- *Peut-on approcher une matrice quelconque par une matrice efficiente ?*
- *Peut-on apprendre une matrice en la contraignant à être efficiente ?*

Afin d'apporter des réponses, ne serait-ce que partielles, à ces trois questions, cette thèse s'organise de la manière suivante :

- Première partie : État de l'art et matrices efficientes
 - Le chapitre 1 introduit les outils d'optimisation mathématique utilisés tout au long de la thèse.
 - Le chapitre 2 constitue un état de l'art des factorisations matricielles les plus communes, en mettant l'accent sur celles ayant les liens les plus forts avec ce qui est proposé dans cette thèse.
 - Le chapitre 3 est pour nous l'occasion de définir rigoureusement l'efficience d'une matrice. Pour cela, nous établissons un parallèle entre une matrice et l'algorithme permettant de faire le produit d'un vecteur par cette matrice. Ce parallèle nous permet, après avoir défini l'efficience d'un algorithme, de définir l'efficience d'une matrice. Nous discutons ensuite d'éventuels moyens de calculer ou tout du moins borner l'efficience d'une matrice.
- Deuxième partie : Contributions
 - Le chapitre 4 concerne l'approximation de matrices quelconques par des matrices efficientes. Une telle approximation implique la possibilité de remplacer une matrice par sa version efficiente afin d'accélérer les traitements effectués avec cette matrice, au prix d'une perte de précision. Nous formulons tout d'abord la tâche d'approximation comme un problème d'optimisation non-convexe. Nous proposons un algorithme basé sur les progrès récents en optimisation non-convexe afin de résoudre le problème avec des garanties de convergence vers un minimum local. Nous proposons ensuite une heuristique permettant de converger vers de "meilleurs" minima locaux.
 - Le chapitre 5 présente une application de la méthode d'approximation à un problème inverse de localisation de source dans un contexte de signaux de magnétoencéphalographie (MEG), pour lequel un compromis entre efficience et qualité d'approximation est mis en évidence.
 - Le chapitre 6 a pour objet la recherche de transformées de Fourier rapides pour le traitement du signal sur graphe. Face aux limitations de la méthode introduite au chapitre 4 dans ce contexte particulier, une autre approche spécifique à ce problème est proposée. Les deux approches sont comparées expérimentalement en fin de chapitre.

- Le chapitre 7 apporte un éclairage théorique sur l’identifiabilité de la forme efficiente d’une matrice à partir de son observation. Des conditions nécessaires et suffisantes d’identifiabilité sont données, et le cas particulier de la matrice de Fourier est plus particulièrement discuté.
- Le chapitre 8 examine le problème de l’apprentissage de matrices efficientes à partir de données d’entraînement, à des fins de reconstruction de données ou de prédiction. Nous débutons avec une modification du problème d’optimisation introduit au chapitre 4 pour prendre en compte la tâche qui nous intéresse. Ceci nous mène à un nouvel algorithme pour l’estimation de matrices efficientes. Nous appliquons cette méthode à l’apprentissage de dictionnaire, puis effectuons une expérience de débruitage d’image. Nous terminons ce chapitre par des considérations théoriques concernant les capacités de généralisation des matrices efficientes apprises, tout d’abord dans le cadre de l’apprentissage de dictionnaire, puis dans le cadre un peu plus général des modèles parcimonieux.

Première partie

État de l'art et matrices efficientes

Nous présentons dans ce chapitre un outil central et indispensable à la compréhension de cette thèse : l'optimisation mathématique. Nous commençons par définir les problèmes d'optimisation. Nous introduisons ensuite la notion de convexité et ses liens avec l'optimisation. Nous terminons par la présentation de divers algorithmes utilisés pour résoudre des problèmes d'optimisation, en insistant plus particulièrement sur les algorithmes utilisés dans cette thèse.

Sommaire

1.1	Problème d'optimisation	11
1.1.1	Formulation	12
1.1.2	Exemples	13
1.2	Convexité	14
1.2.1	Définition	15
1.2.2	Liens avec l'optimisation	15
1.2.3	Optimisation de fonctions non-convexes	16
1.3	Algorithmes	16
1.3.1	Méthodes de descente	17
1.3.2	Algorithmes proximaux	18
1.3.3	Descente par bloc de coordonnées	19
1.3.4	PALM	20
1.3.5	Descente de gradient stochastique	21

1.1. Problème d'optimisation

De nombreux problèmes rencontrés en traitement de données, dans l'industrie ou même dans la vie courante peuvent être vus comme des problèmes d'optimisation. En effet, des questions aussi diverses que *Quels paramètres du modèle expliquent le*

mieux les données observées ?, *Où placer les pylônes afin de répartir au mieux les efforts sur un pont ?* ou *Dans quel ordre effectuer un parcours à étapes afin d'en minimiser la distance totale ?* peuvent toutes trouver des réponses grâce à l'optimisation mathématique. Nous définissons mathématiquement dans la prochaine sous-section les problèmes d'optimisation, puis donnons quelques exemples dans la sous-section suivante. Nous nous appuyons fortement tout au long de ce chapitre sur le livre de Boyd et Vandenberghe (Boyd et Vandenberghe, 2004), qui constitue une excellente introduction à l'optimisation mathématique convexe.

1.1.1. Formulation

Mathématiquement parlant, résoudre un problème d'optimisation consiste à trouver la variable $x \in \mathbb{R}^n$ qui minimise une certaine fonction $f : \mathbb{R}^n \rightarrow \mathbb{R} \cup \{\infty\}$ (on peut remarquer que maximiser une fonction est équivalent à minimiser la fonction opposée). On note un tel problème de manière succincte :

$$\underset{x \in \mathbb{R}^n}{\text{minimiser}} \quad f(x), \quad (1.1)$$

où f est appelée fonction de coût, objectif ou critère, et x est appelée variable de décision ou, dans certains contextes, paramètres. On peut de manière équivalente décrire un problème d'optimisation de la façon suivante :

$$\begin{aligned} &\underset{x \in \mathbb{R}^n}{\text{minimiser}} \quad f_0(x) \\ &\text{sujet à} \quad f_i(x) \leq b_i \quad i = 1, \dots, m. \end{aligned} \quad (1.2)$$

Dans ce cas, on considère une fonction objectif à valeurs réelles $f_0 : \mathbb{R}^n \rightarrow \mathbb{R}$, et des fonctions $f_i : \mathbb{R}^n \rightarrow \mathbb{R}, i = 1, \dots, m$, appelées contraintes. On remarque que les deux formulations sont strictement équivalentes, et liées par $f(x) = f_0(x) + \sum_{i=1}^m g_i(x)$ où $g_i(x) = 0$ si $f_i(x) \leq b_i$ et $g_i(x) = \infty$ sinon (g_i est appelée fonction indicatrice de la i -ème contrainte). Poser un problème d'optimisation revient à définir le domaine $\mathcal{D} = \{x \in \mathbb{R}^n \mid f_i(x) \leq b_i, \forall i = 1, \dots, m\} = \{x \in \mathbb{R}^n \mid f(x) < \infty\}$, qui correspond à l'ensemble de variables parmi lesquelles on recherche la solution. Dans la suite de cette thèse, nous privilégierons la formulation (1.1), principalement pour sa concision.

L'objectif d'un problème d'optimisation est donc de choisir dans un domaine \mathcal{D} la ou les variables qui minimisent la fonction f . Si elles existent, ces variables sont appelées solutions, et se définissent comme les éléments de l'ensemble $\mathcal{S} = \{x^* \in \mathcal{D} \mid f(x^*) \leq f(x), \forall x \in \mathcal{D}\}$. La solution d'un problème d'optimisation est dite unique si \mathcal{S} se réduit à un élément, on note alors l'unique solution x^* , qui est telle que $\forall x \in \mathcal{D} \setminus x^*, f(x^*) < f(x)$. Unicité ou non, la valeur $f(x^*)$ est appelée valeur optimale du problème, et l'ensemble \mathcal{S} des x^* est alors appelé l'ensemble des minima globaux

du problème. On définit un minimum local comme un point $x_l^* \in \mathcal{D}$ tel qu'il existe un voisinage \mathcal{V} de x_l^* tel que $\forall x \in \mathcal{D} \cap \mathcal{V}, f(x_l^*) \leq f(x)$.

Résoudre un problème d'optimisation revient à trouver une solution, autrement dit un élément $x^* \in \mathcal{S}$. Il est parfois difficile de résoudre exactement un problème d'optimisation, il peut alors être raisonnable de chercher une solution approchée au problème, en se fixant une certaine tolérance $\tau > 0$. On définit alors l'ensemble des solutions approchées $\mathcal{S}_\tau = \{x_\tau^* \in \mathcal{D} \mid f(x_\tau^*) - f(x^*) \leq \tau, \forall x^* \in \mathcal{S}\}$, qui correspond à l'ensemble des variables dont l'image par f est supérieure à la valeur optimale d'au plus la tolérance τ . Résoudre un problème d'optimisation de manière approchée revient alors à trouver une solution approchée, autrement dit un élément $x_\tau^* \in \mathcal{S}_\tau$.

1.1.2. Exemples

Nous donnons ici trois exemples de problèmes d'optimisation, dans le but de montrer la grande diversité d'applications de l'optimisation mathématique.

Ajustement affine. Pour ce premier exemple, on considère qu'il existe une relation affine entre deux variables réelles y et z ($y = az + b$). On souhaiterait estimer les paramètres de cette relation affine (les réels a et b), à partir de $m > 2$ observations bruitées de couples $(y, z) : \{(y_1, z_1), \dots, (y_m, z_m)\}$. Pour cela, une approche possible est de choisir la variable $x = (x_1, x_2)$ qui minimise la somme des carrés des écarts entre les observations y_i et leur prédiction par le modèle $x_1 z_i + x_2$. Ceci correspond au problème d'optimisation suivant :

$$\underset{x \in \mathbb{R}^2}{\text{minimiser}} \quad f(x) = \sum_{i=1}^m (y_i - (x_1 z_i + x_2))^2.$$

Ce problème admet une solution analytique (donnée par une formule mathématique). En effet, la solution $x^* = (x_1^*, x_2^*)$ est donnée par $x_1^* = \frac{\sum_{i=1}^m (z_i - \bar{z})(y_i - \bar{y})}{\sum_{i=1}^m (z_i - \bar{z})^2}$ et $x_2^* = \bar{y} - x_1^* \bar{z}$, avec $\bar{y} = \frac{1}{m} \sum_{i=1}^m y_i$ et $\bar{z} = \frac{1}{m} \sum_{i=1}^m z_i$. Ce problème d'optimisation constitue un cas particulier de régression linéaire, qui admet une solution générale (pour $z \in \mathbb{R}^p$) s'exprimant de manière simple à l'aide de notations matricielles introduites dans les chapitres suivants. Ce type de problème est rencontré dans une multitude d'applications (en exagérant à peine, dès que des mesures d'un phénomène bruité sont utilisées, avec plus de mesures que de paramètres à estimer).

Minimisation des efforts. Pour ce deuxième exemple, on considère la répartition de n pylônes sur la longueur d'un pont supporté à ses deux extrémités. Le but est de répartir les pylônes afin de minimiser la flèche (déformation verticale) maximale au long du pont. La variable de décision $x = (x_1, \dots, x_n)$ considérée ici correspond

aux positions des pylônes (le i -ème pylône étant à la position x_i). Sachant que la flèche maximale entre le i -ème et le $i+1$ -ème pylône est proportionnelle à la quantité $(x_{i+1} - x_i)^4$, et que le pont est de longueur L , le problème de répartition se résume au problème d'optimisation suivant (avec $x_0 = 0$ et $x_{n+1} = L$ pour respecter les conditions aux extrémités) :

$$\underset{x \in \mathbb{R}^n}{\text{minimiser}} \quad f(x) = \max_{i \in \{0, \dots, n\}} (x_{i+1} - x_i).$$

Ce problème admet aussi une solution analytique, donnée par $x_i^* = \frac{iL}{n+1}$. Ceci revient à répartir uniformément les pylônes le long du pont à une distance $\frac{L}{n+1}$ de leurs deux plus proches voisins. Cette solution qui peut paraître triviale et découler du bon sens est en fait la solution d'un problème d'optimisation. Cet exemple sert à montrer que l'optimisation mathématique est un langage qu'on parle parfois sans même s'en apercevoir.

Optimisation de trajet. En guise de troisième exemple, on considère un problème d'apparence simple mais qui s'avère en réalité très complexe. Soient un point de départ et n lieux numérotés de 1 à n , on recherche le plus court chemin qui part du point de départ, visite les n lieux et revient au point de départ (on associera l'indice zéro au point de départ), en ayant connaissance des distances qui séparent les lieux (on notera $d(i, j)$ la distance qui sépare le i -ème et le j -ème lieu). Ce problème revient à optimiser un trajet à étapes, en choisissant l'ordre qui mène au plus court chemin. La variable de décision $x = (x_1, \dots, x_n)$ pour ce problème appartient à l'ensemble des images de l'ensemble des entiers 1 à n par une permutation, noté \mathcal{P}_n . Ce problème est connu sous le nom de “problème du voyageur de commerce” ou “travelling salesman problem” en anglais. Il s'énonce comme suit :

$$\underset{x \in \mathbb{R}^n}{\text{minimiser}} \quad f(x) = \sum_{i=0}^{n+1} d(x_{i+1}, x_i) + \delta_{\mathcal{P}_n}(x),$$

où $\delta_{\mathcal{P}_n}(x) = 0$ si $x \in \mathcal{P}_n$ et $\delta_{\mathcal{P}_n}(x) = +\infty$ sinon ($\delta_{\mathcal{P}_n}$ est la fonction indicatrice de l'ensemble \mathcal{P}_n). Ce problème n'a pas de solution analytique. Encore pire, on ne connaît pas de moyen de trouver une solution à ce problème en un temps raisonnable (polynomial en le nombre n de lieux à visiter (Papadimitriou, 1977)). On est alors contraint d'avoir recours à des méthodes de résolution approchée dès lors que n est grand. Malgré cela, ce problème a des applications évidentes en logistique, mais aussi en micro-électronique et en génétique (Lawler et al., 1985).

1.2. Convexité

Nous introduisons dans cette section une notion centrale de l'optimisation : la convexité. Nous commençons par la définir et discutons ensuite ses liens très forts avec l'optimisation.

1.2.1. Définition

Définition 1. Une fonction $f : \mathbb{R}^n \rightarrow \mathbb{R} \cup \infty$ est dite **convexe** si et seulement si pour tout $x, y \in \mathbb{R}^n$, et pour tout $\alpha \in [0, 1]$,

$$f(\alpha x + (1 - \alpha)y) \leq \alpha f(x) + (1 - \alpha)f(y).$$

Cette définition signifie tout simplement que la courbe d'une fonction convexe (membre de gauche de l'inégalité) est en-dessous de toutes ses cordes (membre de droite de l'inégalité).

On définit le gradient $\nabla f(x) \in \mathbb{R}^n$ comme le vecteur $(\frac{\delta f}{\delta x_1}(x), \dots, \frac{\delta f}{\delta x_n}(x))$, où $\frac{\delta f}{\delta x_i}(x)$ est la dérivée partielle de la fonction f par rapport à la i -ème composante de sa variable d'entrée évaluée en x . Le gradient n'existe que si la fonction f est différentiable. Quelques manipulations de la définition précédente ainsi que de celle du gradient mènent à la proposition suivante, particulièrement intéressante du point de vue de l'optimisation.

Proposition 1 (Condition du premier ordre). Une fonction différentiable $f : \mathbb{R}^n \rightarrow \mathbb{R} \cup \infty$ est convexe si et seulement si pour tout $x, y \in \mathbb{R}^n$,

$$f(y) \geq f(x) + \nabla f(x)^T(y - x)$$

Cette condition du premier ordre signifie que la courbe d'une fonction convexe (membre de gauche de l'inégalité) est au-dessus de toutes ses tangentes (membre de droite de l'inégalité). Ceci est illustré par la figure 4, où un exemple de graphe d'une fonction convexe et d'une fonction non-convexe est donné.

1.2.2. Liens avec l'optimisation

Regardons maintenant quelles sont les implications de la convexité en optimisation. La condition du premier ordre indique que si le gradient en un point x^* s'annule ($\nabla f(x^*) = 0$), alors l'inégalité $f(x^*) \leq f(x)$ est vérifiée pour tout $x \in \mathbb{R}^n$, et x^* est alors un minimum global de f . Cette propriété forte des fonctions convexes est très utile en optimisation, puisqu'il "suffit" alors de trouver un point d'annulation du gradient (appelé point stationnaire) pour minimiser globalement la fonction.

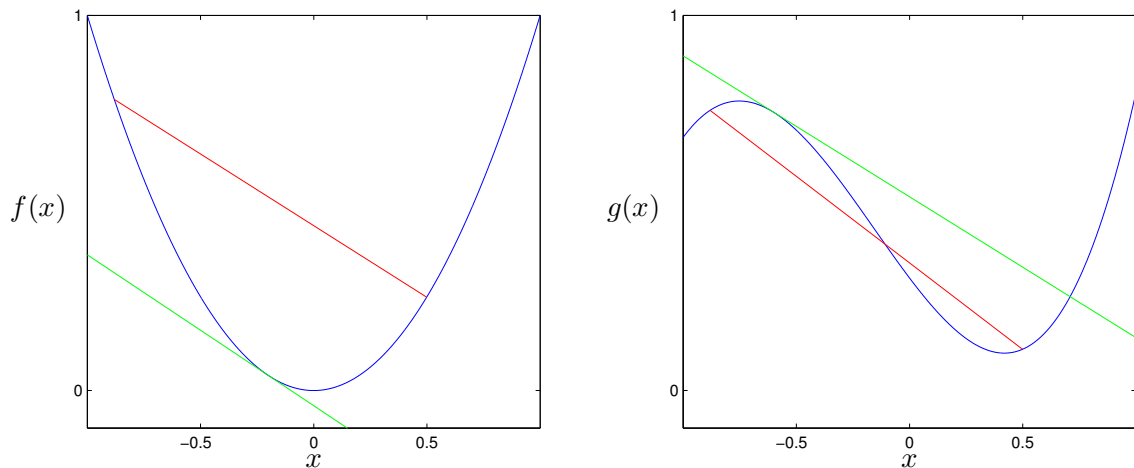


Figure 4 : Exemple d’une fonction f convexe, et d’une fonction g non-convexe. La courbe de f est en-dessous de toutes ses cordes (en rouge) et au-dessus de toutes ses tangentes (en vert), ce qui n’est pas le cas de la courbe de g .

D’autre part, un minimum local x_i^* d’une fonction différentiable vérifie $\nabla f(x_i^*) = 0$ (conséquence quasi-immédiate de la définition des minima locaux et du gradient), ce qui indique qu’un minimum local d’une fonction convexe est forcément aussi un minimum global. En d’autres termes, les fonction convexes n’ont que des minima globaux.

Comme nous le verrons dans la section suivante, et à la lumière des deux paragraphes précédents, les algorithmes d’optimisation convexes se résument bien souvent à annuler le gradient de la fonction f à minimiser.

1.2.3. Optimisation de fonctions non-convexes

La convexité est une propriété qui simplifie grandement l’optimisation. Cependant, et ce sera la plupart du temps le cas dans cette thèse, les fonctions d’intérêt à minimiser ne sont pas forcément convexes. Dans ce cas, il peut exister des minima locaux, et l’optimisation est plus difficile.

Malgré tout, il est possible de s’attaquer à des problèmes non-convexes avec des méthodes de l’optimisation convexe. Ceci se fait au prix de la perte des garanties d’optimalité globale dont jouissent la plupart des méthodes d’optimisation convexe (elles peuvent en général être remplacées par des garanties d’optimalité locale dans le cas non-convexe). On recherche alors des solutions approchées au problème considéré. De plus cela rend les méthodes de résolution sensibles à l’initialisation, et il est courant d’avoir recours à des heuristiques pour choisir une bonne initialisation.

1.3. Algorithmes

Un problème d'optimisation n'admet pas forcément de solution analytique, ou celle-ci peut être coûteuse à calculer (en réalité, la grande majorité des problèmes d'optimisation rencontrés dans les applications n'admet pas de solution analytique). Dans ce cas, on peut avoir recours à des algorithmes itératifs, qui partent d'un point initial et construisent une séquence de points (les itérés) afin de converger vers une solution exacte ou approchée du problème. Nous présentons dans cette section quelques algorithmes itératifs d'optimisation, en insistant sur ceux qu'on utilise effectivement dans cette thèse.

1.3.1. Méthodes de descente

Les méthodes de descente regroupent une grande variété d'algorithmes d'optimisation. Leur principe est simple, à partir d'un point initial $x \in \mathbb{R}^n$, elles consistent à choisir une direction de mise à jour $\Delta x \in \mathbb{R}^n$ à partir des propriétés locales de la fonction de coût f . Une fois la direction de mise à jour fixée, un pas $\lambda \in \mathbb{R}_+$ est choisi à partir du comportement de f selon cette direction, et x est mis à jour en le remplaçant par $x + \lambda\Delta x$. Cette procédure est répétée jusqu'à ce qu'un critère d'arrêt de l'algorithme soit vérifié. Les méthodes de descente tirent leur nom du fait qu'afin de minimiser f , Δx et λ sont choisis à chaque étape de l'algorithme tels que $f(x + \lambda\Delta x) < f(x)$. Un algorithme général correspondant aux méthodes de descente est donné en algorithme 1.

Algorithme 1 Algorithme de descente général

Entrée : Une fonction f à minimiser, une valeur initiale x , un critère d'arrêt c .

- 1: **while** $c = \text{False}$ **do**
- 2: Déterminer une direction de descente Δx
- 3: Choisir un pas λ
- 4: Mettre à jour $x \leftarrow x + \lambda\Delta x$
- 5: **end while**

Sortie : Le point atteint : x .

Une manière intuitive de choisir la direction de descente consiste à choisir la direction de plus forte décroissance locale de la fonction de coût f . Ceci correspond à la direction opposée au gradient ∇f . Choisir une telle direction de descente requiert bien entendu une fonction de coût différentiable. Ce choix de direction de descente donne lieu à l'algorithme très populaire de descente de gradient, donné en algorithme 2, qui sera un des piliers des algorithmes développés dans cette thèse. L'algorithme de descente de gradient converge vers un point stationnaire de la fonction

Algorithme 2 Algorithme de descente de gradient

Entrée : Une fonction f à minimiser, une valeur initiale x , un critère d'arrêt c .

- 1: **while** $c = \text{False}$ **do**
- 2: $\Delta x \leftarrow -\nabla f(x)$
- 3: Choisir un pas λ
- 4: Mettre à jour $x \leftarrow x + \lambda\Delta x$
- 5: **end while**

Sortie : Le point atteint : x .

de coût (qui s'avère être un minimum global pour des fonctions de coût convexes), si le pas λ est bien choisi (plus petit qu'un pas critique λ_c qui dépend des propriétés de f).

Il existe bien évidemment d'autres choix pour la direction de descente, que nous n'évoquons pas ici, ainsi que de multiples méthodes de détermination du pas. Des revues bien plus détaillées des méthodes de descente sont disponibles, voir (Boyd et Vandenberghe, 2004, Chapitre 9) par exemple.

1.3.2. Algorithmes proximaux

Les méthodes de descente ont l'avantage d'être relativement simples et intuitive et jouissent de garanties d'optimalité dans la plupart des cas. Cependant, elles nécessitent de disposer d'un moyen d'évaluer les directions en terme de variation locale de la fonction de coût. Ceci n'est possible que si f est différentiable.

Dans le cas contraire, une stratégie différente doit être adoptée. Par exemple, on peut vouloir faire décroître l'objectif f tout en restant proche de l'itéré courant x . On procède alors à ce qu'on appelle une itération proximale :

$$x \leftarrow \text{prox}_{\lambda f}(x) \triangleq \underset{u}{\operatorname{argmin}} \left(f(u) + \frac{1}{2\lambda} \|u - x\|_2^2 \right),$$

où le second terme de la somme, appelé terme de régularisation, incite à rester proche de l'itéré courant, et λ peut être vu comme un pas (plus λ est grand, moins l'accent est mis sur la proximité avec l'itéré courant). L'itération proximale requiert de pouvoir minimiser $f(u) + \frac{1}{2\lambda} \|u - x\|_2^2$, ce qui n'est possible que pour des fonctions f relativement simples. Appliquer de manière successive l'itération proximale donne naissance à l'algorithme du point proximal, donné en algorithme 3. Cet algorithme converge vers un minimum local de la fonction f (qui s'avère être un minimum global pour les fonctions de coût convexes). Cependant, cet algorithme trouve peu d'applications, car rares sont les cas où il est plus aisé de minimiser $f(u) + \frac{1}{2\lambda} \|u - x\|_2^2$ que de minimiser directement $f(u)$, auquel cas l'algorithme du point proximal a un intérêt (sinon un minimum global peut être atteint directement).

Algorithme 3 Algorithme du point proximal

Entrée : Une fonction f à minimiser, une valeur initiale x , un critère d'arrêt c .

- 1: **while** $c = \text{False}$ **do**
- 2: Choisir un pas λ
- 3: Mettre à jour $x \leftarrow \text{prox}_{\lambda f}(x)$
- 4: **end while**

Sortie : Le point atteint : x .

Par contre, utiliser l'itération proximale prend tout son sens par exemple lorsque l'objectif prend la forme $f = g + h$, où g est une fonction différentiable et h une fonction non-différentiable dont l'itération proximale est simple à calculer. Une stratégie possible est alors d'approcher localement f par \hat{f}_x , puis de minimiser \hat{f}_x en restant dans un voisinage de x pour que l'approximation reste bonne, et itérer. Plus concrètement, il est possible de définir $\hat{f}_x = \hat{g}_x + h$, où $\hat{g}_x(u) = g(x) + \nabla g(x)^T(u - x)$ est l'approximation linéaire de g au point x . L'itération proximale correspondant à cette approximation locale prend alors la forme $\text{prox}_{\lambda \hat{f}_x}(x) = \text{prox}_{\lambda h}(x - \lambda \nabla g(x))$, et revient à effectuer un pas de descente de gradient sur g suivi d'une itération proximale sur h . L'algorithme engendré est donné en algorithme 4, il converge vers un minimum local pour un pas λ plus petit qu'un pas critique λ_c qui dépend des propriétés de g .

Plus de détails sur les algorithmes proximaux (autres interprétations, preuves de convergence, etc.) sont donnés par exemple dans l'article de revue ([Parikh et Boyd, 2014](#)).

Algorithme 4 Algorithme de descente de gradient proximale

Entrée : Une fonction $f = g + h$ à minimiser, une valeur initiale x , un critère d'arrêt c .

- 1: **while** $c = \text{False}$ **do**
- 2: $\Delta x \leftarrow -\nabla g(x)$
- 3: Choisir un pas λ
- 4: Mettre à jour $x \leftarrow \text{prox}_{\lambda h}(x + \lambda \Delta x)$
- 5: **end while**

Sortie : Le point atteint : x .

1.3.3. Descente par bloc de coordonnées

Dans certaines situations, il peut être préférable de ne pas mettre à jour tout x d'un coup, mais de le diviser en p blocs de coordonnées x_1, \dots, x_p que l'on met à jour

alternativement. Une telle mise à jour par bloc est justifiée soit parce que cette division est naturelle (x regroupe p paramètres distincts), soit parce qu’il est plus simple ou moins coûteux de ne mettre à jour que quelques coordonnées à la fois (x est alors divisé artificiellement en p blocs). L’algorithme de descente par bloc de coordonnées donné en algorithme 5 est très similaire à l’algorithme de descente classique, sauf qu’une sélection du bloc à mettre à jour est ajoutée au début de chaque itération. Cette sélection est en général cyclique (elle prend la forme d’une boucle), on parle alors de descente par bloc de coordonnées cyclique. Cependant, la sélection du bloc peut aussi être faite de manière aléatoire, ce qui peut être plus adapté pour certaines applications (Nesterov, 2012). Dans le cas d’une direction de descente choisie opposée au gradient, cet algorithme converge vers un point stationnaire de l’objectif (Wright, 2015).

Algorithme 5 Algorithme de descente par bloc de coordonnées

Entrée : Une fonction f à minimiser, une valeur initiale divisée en p blocs de coordonnées $x = (x_1, \dots, x_p)$, un critère d’arrêt c .

- 1: **while** $c = \text{False}$ **do**
- 2: Déterminer un bloc à mettre à jour i
- 3: Déterminer une direction de descente Δx_i
- 4: Choisir un pas λ
- 5: Mettre à jour $x_i \leftarrow x_i + \lambda \Delta x_i$
- 6: **end while**

Sortie : Le point atteint : x .

1.3.4. PALM

Présentons maintenant un algorithme sur lequel s’appuient fortement les méthodes développées dans cette thèse. L’algorithme en question, baptisé PALM, pour “Proximal Alternating Linearized Minimization” (Bolte et al., 2014), s’apparente à une combinaison des idées de la descente de gradient proximal et de la descente par bloc de coordonnées. L’algorithme PALM s’applique à des fonctions de coût de la forme $f = g + h$ avec $h(x) = \sum_{i=1}^p h_i(x_i)$, où l’itération proximale associée à chaque h_i est simple à calculer, et g est différentiable. L’idée principale de PALM est d’effectuer une itération de descente de gradient proximale sur chaque bloc, en mettant à jour de manière cyclique les blocs de coordonnées. Pour cela, on définit $\nabla_i g$ comme le gradient de la fonction g restreint aux coordonnées du i -ème bloc. L’algorithme PALM est donné en algorithme 6.

Cet algorithme jouit de garanties de convergence vers un point critique de l’objectif f (généralisation du concept de point stationnaire aux fonction non-différentiables,

Algorithme 6 Algorithme PALM (descente proximale par bloc de coordonnées)

Entrée : Une fonction $f = g + h$ à minimiser, une valeur initiale divisée en p blocs de coordonnées $x = (x_1, \dots, x_p)$, un critère d'arrêt c .

- 1: **while** $c = \text{False}$ **do**
- 2: Déterminer un bloc à mettre à jour i (de manière cyclique)
- 3: $\Delta x_i \leftarrow -\nabla_i g(x_i)$
- 4: Choisir un pas λ
- 5: Mettre à jour $x_i \leftarrow \text{prox}_{\lambda h_i}(x_i + \lambda \Delta x_i)$
- 6: **end while**

Sortie : Le point atteint : x .

voir (Bolte et al., 2014) pour une définition précise). Nous donnons plus de précisions sur les conditions de convergence de cet algorithme dans la partie contributions de cette thèse.

1.3.5. Descente de gradient stochastique

Il est courant en traitement du signal et en apprentissage automatique de rencontrer des problèmes d'optimisations dont l'objectif est la moyenne de N fonctions différentiables f_1, \dots, f_N dépendant chacune d'une observation ou d'une mesure indépendante du même phénomène. On a alors une fonction de coût de la forme $f = \sum_{i=1}^N f_i$. Le nombre d'observations N pouvant être très grand pour certaines applications, il peut être difficile ne serait-ce que de calculer le gradient de la fonction globale f , pour effectuer par exemple une descente de gradient. Dans ce cas, afin de réduire le coût de calcul, il est possible d'estimer ce gradient en ne considérant qu'un petit nombre L d'observations à la fois. Insérer une telle estimation du gradient dans un algorithme de descente de gradient classique donne lieu à un algorithme appelé descente de gradient stochastique, donné en algorithme 7. Il est préférable que le choix du sous-ensemble \mathcal{B} des observations considérées change à chaque itération de l'algorithme (en effectuant un tirage aléatoire par exemple), afin de profiter de toute la diversité de l'ensemble des observations. L'algorithme de descente de gradient stochastique jouit de garanties de convergences similaires à la descente de gradient classique, sous certaines conditions (Bottou, 1998).

Algorithme 7 Algorithme de descente de gradient stochastique

Entrée : Une fonction $f = \frac{1}{N} \sum_{i=1}^N f_i$ à minimiser, une valeur initiale x , un critère d'arrêt c .

- 1: **while** $c = \text{False}$ **do**
- 2: Choisir L indices $\{i_1, \dots, i_L\} \triangleq \mathcal{B} \subset \{1, \dots, N\}$
- 3: $\Delta x \leftarrow -\frac{1}{L} \sum_{j=1}^L \nabla f_{i_j}(x)$
- 4: Choisir un pas λ
- 5: Mettre à jour $x \leftarrow x + \lambda \Delta x$
- 6: **end while**

Sortie : Le point atteint : x .

Conclusion. Nous avons profité de ce chapitre pour effectuer une introduction aux outils d'optimisation mathématique utilisés dans cette thèse. Nous avons débuté par la définition de notions très générales concernant l'optimisation pour aboutir finalement à la description d'un algorithme fondamental à cette thèse. La présentation faite ici des outils est volontairement partielle, l'accent étant mis sur les algorithmes, leurs motivations et principes généraux plutôt que sur les détails d'implémentation et preuves de convergence, qui sont disponibles dans les références (Boyd et Vandenberghe, 2004; Parikh et Boyd, 2014; Wright, 2015; Bolte et al., 2014).

Factorisation matricielle

Ce chapitre introduit les factorisations matricielles, éléments cruciaux de cette thèse. Nous débutons par une brève introduction aux matrices et à leurs liens avec la notion de linéarité. Nous présentons ensuite diverses factorisations usuelles. Nous terminons par présenter deux domaines de recherche en lien avec la factorisation matricielle : l'apprentissage de dictionnaire et les réseaux de neurones.

Sommaire

2.1 Matrices et linéarité	23
2.1.1 Matrices	24
2.1.2 Linéarité et produit matriciel	24
2.1.3 Stockage de données	25
2.2 Factorisations usuelles	26
2.2.1 Résolution de systèmes linéaires	27
2.2.2 Décomposition en valeurs singulières	29
2.3 Apprentissage de dictionnaire	30
2.3.1 Historique	30
2.3.2 Formulation en tant que factorisation matricielle	31
2.3.3 Résolution du problème	32
2.4 Réseaux de neurones	35
2.4.1 Principe	35
2.4.2 Optimisation	37

2.1. Matrices et linéarité

Cette section vise principalement à introduire les notations et conventions utilisées tout au long de la thèse pour désigner les matrices et les objets s'y rapportant. Nous rappelons aussi brièvement le lien unissant le produit matriciel et les transformations

linéaires en dimension finie. Nous terminons par évoquer une autre fonction possible des matrices : le stockage de données. Cette section s'adresse principalement au lecteur n'étant pas, ou peu, familier avec les notations standards utilisées en algèbre linéaire.

2.1.1. Matrices

Dans le cas général, une matrice \mathbf{A} est un tableau de nombres à m lignes et n colonnes. Les nombres constituant la matrice, appelés coefficients ou entrées, peuvent être complexes, on note alors $\mathbf{A} \in \mathbb{C}^{m \times n}$, ou réels, on note alors $\mathbf{A} \in \mathbb{R}^{m \times n}$. Nous considérons dans cette thèse, sauf mention contraire, des matrices à coefficients réels. La i -ème ligne de \mathbf{A} constitue un n -uplet de nombres, appelé vecteur, et est notée \mathbf{a}^i (on a $\mathbf{a}^i \in \mathbb{R}^n$). De manière similaire, la j -ème colonne de \mathbf{A} est notée \mathbf{a}_j (on a $\mathbf{a}_j \in \mathbb{R}^m$). Le coefficient de \mathbf{A} situé à la i -ème ligne et à la j -ème colonne est noté a_{ij} (on a $a_{ij} \in \mathbb{R}$). Toutes ces notations sont résumées par l'équation suivante :

$$\mathbf{A} = \begin{pmatrix} a_{11} & \dots & a_{1j} & \dots & a_{1n} \\ \vdots & & \vdots & & \vdots \\ a_{i1} & \dots & a_{ij} & \dots & a_{in} \\ \vdots & & \vdots & & \vdots \\ a_{m1} & \dots & a_{mj} & \dots & a_{mn} \end{pmatrix} \mathbf{a}^i$$

\mathbf{a}_j

Une matrice $\mathbf{B} \in \mathbb{R}^{n \times m}$ est dite matrice transposée de \mathbf{A} et notée \mathbf{A}^T si et seulement si l'on a $b_{ji} = a_{ij}$, $\forall i \in \{1, \dots, m\}, j \in \{1, \dots, n\}$. On pourra noter $\mathbf{A} = (\mathbf{a}_1 | \dots | \mathbf{a}_n)$ pour signifier que \mathbf{A} correspond à la juxtaposition horizontale de ses colonnes et $\mathbf{A} = (\mathbf{a}^1 | \dots | \mathbf{a}^m)^T$ pour signifier qu'elle correspond à la juxtaposition verticale de ses lignes.

2.1.2. Linéarité et produit matriciel

Nous rappelons ici brièvement la correspondance existant entre les transformations linéaires en dimension finie et le produit matriciel. Considérons une application $f : E \rightarrow F$, où E et F sont des espaces vectoriels de dimension finie (cela signifiant

qu'il existe une base \mathcal{B} de E (respectivement une base \mathcal{C} de F) comportant un nombre fini de vecteurs (n pour E et m pour F). L'application (ou transformation) f est dite linéaire si elle vérifie les deux propriétés suivantes :

- *Additivité* : $f(\mathbf{x} + \mathbf{y}) = f(\mathbf{x}) + f(\mathbf{y})$, $\forall \mathbf{x} \in E, \forall \mathbf{y} \in E$.
- *Homogénéité* : $f(\alpha \mathbf{x}) = \alpha f(\mathbf{x})$, $\forall \mathbf{x} \in E, \forall \alpha \in \mathbb{R}$.

la linéarité est une propriété puissante, pouvant être vue comme une généralisation de la notion de proportionnalité, et induisant divers phénomènes étudiés par l'*algèbre linéaire*. Pour une définition mathématiquement rigoureuse des concepts d'espace vectoriel et de base, ainsi qu'une introduction à l'algèbre linéaire dans son ensemble, voir, par exemple, (Strang, 2005).

Nous nous intéressons ici à une conséquence particulière de la linéarité en dimension finie : le fait de pouvoir représenter toute application linéaire $f : E \rightarrow F$, où E est de dimension n et F de dimension m par une matrice $\mathbf{A} \in \mathbb{R}^{m \times n}$. Par abus de langage, et pour simplifier les notations, nous confondons volontairement dans cette thèse un vecteur et ses coordonnées dans la base implicitement utilisée. Par exemple, pour un vecteur quelconque de E , la notation \mathbf{x} désigne le vecteur $\sum_{j=1}^n x_j \mathbf{b}_j$, où $\mathcal{B} = \{\mathbf{b}_1, \dots, \mathbf{b}_n\}$ désigne la base de E implicitement utilisée. La linéarité de f implique que $f(\mathbf{x}) = \sum_{j=1}^n x_j f(\mathbf{b}_j)$. Ainsi, il suffit de connaître les images par f des vecteurs de la base \mathcal{B} (exprimées dans une base quelconque \mathcal{C} de F) et les coordonnées de \mathbf{x} dans cette même base pour calculer $\mathbf{y} \triangleq f(\mathbf{x})$. Définissons la matrice \mathbf{A} dont les colonnes correspondent aux images par f des vecteurs de la base \mathcal{B} (exprimées dans la base \mathcal{C}) : $\mathbf{A} = (f(\mathbf{b}_1) | \dots | f(\mathbf{b}_n))$. Cette matrice spécifie entièrement f , et le produit matrice/vecteur est défini de façon à obtenir : $\mathbf{y} = \mathbf{A}\mathbf{x}$. Quelques manipulations mènent à la formule suivante : $y_i = \sum_{k=1}^n a_{ik} x_k$. De même, on définit le produit de deux matrices en juxtaposant horizontalement les images \mathbf{y}_j de vecteurs \mathbf{x}_j , ce qui donne $\mathbf{Y} = \mathbf{A}\mathbf{X}$ avec $y_{ij} = \sum_{k=1}^n a_{ik} x_{kj}$. Le produit matrice/vecteur peut donc être vu comme le moyen de calculer l'image d'un vecteur par une application linéaire. Réciproquement, toute matrice \mathbf{A} correspond à la représentation d'une application linéaire. Cette interprétation n'est pas unique, et nous présentons dans la sous-section suivante une autre manière de voir les matrices.

2.1.3. Stockage de données

De manière plus directe que dans la sous-section précédente, une matrice $\mathbf{A} \in \mathbb{R}^{m \times n}$ peut simplement être vue comme un tableau dans lequel des données (mn nombre réels) sont stockées, sans aucun lien avec une quelconque notion d'application linéaire associée. Par exemple, dans un contexte médical, \mathbf{A} peut être la juxtaposition de n vecteurs d'intérêt $\mathbf{a}_1, \dots, \mathbf{a}_n \in \mathbb{R}^m$, où chaque \mathbf{a}_i regroupe les mesures de m paramètres physiologiques (température, tension, etc.) sur un même

patient. Chaque colonne de la matrice \mathbf{A} correspond alors à un patient différent, et chaque ligne à un paramètre physiologique différent. Dans un contexte totalement différent, la matrice \mathbf{A} peut par exemple représenter une image numérique. Dans ce cas, l'entrée a_{ij} de la matrice représente l'intensité lumineuse du pixel situé à la i -ème ligne et à la j -ème colonne de l'image. Les indices de la matrice correspondent alors aux directions spatiales de l'image (horizontale et verticale).

L'interprétation la plus appropriée à donner à une matrice (représentation d'une application linéaire ou stockage de données) se déduit la plupart du temps trivialement du contexte. Parfois les deux interprétations sont valides, et apportent alors des éclairages différents sur la matrice que l'on manipule.

2.2. Factorisations usuelles

Factoriser une matrice $\mathbf{A} \in \mathbb{R}^{m \times n}$, c'est trouver $K \geq 2$ matrices $\mathbf{F}_1, \dots, \mathbf{F}_K$ qu'on appelle les facteurs, telles que $\mathbf{A} \approx \mathbf{F}_1 \dots \mathbf{F}_K$. On dit que la factorisation est exacte si l'égalité est stricte : $\mathbf{A} = \mathbf{F}_1 \dots \mathbf{F}_K$ (on parle sinon de factorisation approchée). La factorisation matricielle peut être vue (au moins de manière conceptuelle) comme la résolution exacte ou approchée d'un problème d'optimisation de la forme suivante :

$$\underset{\mathbf{F}_1, \dots, \mathbf{F}_K}{\text{minimiser}} \quad d(\mathbf{A}, \mathbf{F}_1 \dots \mathbf{F}_K) + \sum_{k=1}^K g_k(\mathbf{F}_k),$$

où $d(.,.)$ désigne une fonction qui mesure la distance entre \mathbf{A} et le produit des facteurs $\mathbf{F}_1 \dots \mathbf{F}_K$, et qui est nulle si la factorisation est exacte, et les g_k sont des pénalités qui imposent une certaine structure aux facteurs (ce sont le plus souvent des fonctions indicatrices d'ensembles, nulles si le facteur \mathbf{F}_k a la bonne structure, et infinies sinon). Dans la plupart des cas, si la factorisation recherchée est exacte, le problème d'optimisation est implicite, et jamais explicitement posé en ces termes, mais il est toujours sous-jacent.

Les motivations pour factoriser une matrice sont variées. La factorisation nous donne une manière alternative d'exprimer la matrice \mathbf{A} , sous la forme du produit des facteurs $\mathbf{F}_1 \dots \mathbf{F}_K$. Cette forme peut être utile soit simplement à faciliter l'interprétation de la matrice, soit à faciliter une certaine tâche associée à la matrice (résoudre un système linéaire, un problème d'approximation, extraire des caractéristiques spécifiques de la matrice, etc.). Nous présentons dans les sous-sections et sections suivantes de ce chapitre, quelques factorisations matricielles ainsi que leur contexte d'utilisation.

2.2.1. Résolution de systèmes linéaires

Soient n valeurs réelles inconnues x_1, \dots, x_n . Une équation linéaire en ces inconnues est de la forme $a_1x_1 + \dots + a_nx_n = b$ (qui peut s'écrire de manière plus succincte $\mathbf{a}^T \mathbf{x} = b$). On appelle système de m équations linéaires à n inconnues un ensemble d'équations $\{a_{i1}x_1 + \dots + a_{in}x_n = b_i, i = 1, \dots, m\}$. Un tel système s'écrit de manière équivalente en utilisant les notations matricielles $\mathbf{A}\mathbf{x} = \mathbf{b}$, avec $\mathbf{A} \in \mathbb{R}^{m \times n}$. Résoudre le système revient à trouver le (ou les) vecteur $\mathbf{x} \in \mathbb{R}^n$ satisfaisant l'équation matricielle, \mathbf{A} et \mathbf{b} étant connus. Les systèmes linéaires sont présents dans quasiment tous les domaines des sciences et de l'ingénierie, si bien que leur importance peut difficilement être surestimée. Nous présentons dans cette section des factorisations matricielles utiles notamment (mais pas exclusivement) à faciliter la résolution de tels systèmes linéaires.

Factorisation LU. Considérons pour plus de simplicité d'exposition une matrice carrée $\mathbf{A} \in \mathbb{R}^{n \times n}$ (même si la factorisation LU peut aussi s'appliquer aux matrices rectangulaires). Le but de la factorisation LU est d'obtenir une matrice $\mathbf{L} \in \mathbb{R}^{n \times n}$ triangulaire inférieure ($l_{ij} = 0$ si $j > i$) et une matrice $\mathbf{U} \in \mathbb{R}^{n \times n}$ triangulaire supérieure ($u_{ij} = 0$ si $i > j$) telles que $\mathbf{A} = \mathbf{L}\mathbf{U}$. Effectuer une factorisation LU correspond donc implicitement à la résolution exacte du problème d'optimisation suivant :

$$\underset{\mathbf{L}, \mathbf{U}}{\text{minimiser}} \quad d(\mathbf{A}, \mathbf{L}\mathbf{U}) + g_l(\mathbf{L}) + g_u(\mathbf{U}),$$

où g_l (respectivement g_u) est la fonction indicatrice de l'ensemble des matrices triangulaires inférieures (respectivement supérieures). Un exemple graphique de factorisation LU est donné en figure 5.

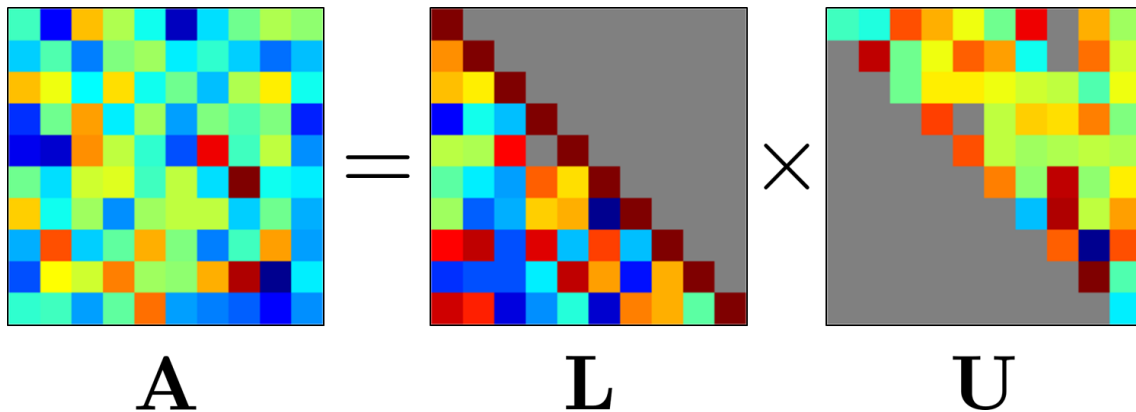


Figure 5 : Un exemple de factorisation LU, pour une matrice $\mathbf{A} \in \mathbb{R}^{10 \times 10}$. Les zones colorées correspondent à des entrées non-nulles, les grises à des entrées nulles.

Lors de la résolution d'un système de la forme $\mathbf{Ax} = \mathbf{b}$, disposer une telle factorisation de \mathbf{A} permet de ré-exprimer le système sous la forme $\mathbf{LUx} = \mathbf{b}$, qui peut se résoudre en deux étapes, tout d'abord en considérant le système $\mathbf{Ly} = \mathbf{b}$, puis le système $\mathbf{Ux} = \mathbf{y}$. Ce sont deux systèmes triangulaires simples et relativement rapides à résoudre (en $\mathcal{O}(n^2)$), par la technique dite de substitution avant/arrière, où les entrées de \mathbf{y} puis de \mathbf{x} sont déterminées une par une, en peu d'opérations.

Pour obtenir une factorisation LU, on a recours à une variante de l'algorithme d'élimination de Gauss, où l'on cherche à annuler les entrées de \mathbf{A} situées sous la diagonale par des opérations élémentaires sur ses lignes. Il existe des variantes de la factorisation LU spécifiques aux matrices symétriques (appelée factorisation LDLT), ou aux matrices symétriques définies positives (appelée factorisation de Cholesky). Ce type de factorisations coûte $\mathcal{O}(n^3)$ opérations arithmétiques ($\mathcal{O}(x)$ signifiant de l'ordre de x). Pour plus de détails sur la factorisation LU (algorithmes explicites, propriétés, complexité) et ses variantes, voir (Golub et Van Loan, 2012).

Factorisation QR. Le but de la factorisation QR est d'obtenir une matrice $\mathbf{Q} \in \mathbb{R}^{m \times m}$ orthogonale ($\mathbf{Q}^T \mathbf{Q} = \mathbf{Id}$) et une matrice $\mathbf{R} \in \mathbb{R}^{m \times n}$ triangulaire supérieure telles que $\mathbf{A} = \mathbf{QR}$. La factorisation QR correspond donc au problème d'optimisation suivant :

$$\underset{\mathbf{Q}, \mathbf{R}}{\text{minimiser}} \quad d(\mathbf{A}, \mathbf{QR}) + g_q(\mathbf{Q}) + g_r(\mathbf{R}),$$

où g_q est l'indicatrice de l'ensemble des matrices orthogonales et g_r celle de l'ensemble des matrices triangulaires supérieures. Un exemple de factorisation QR est donné en figure 6.

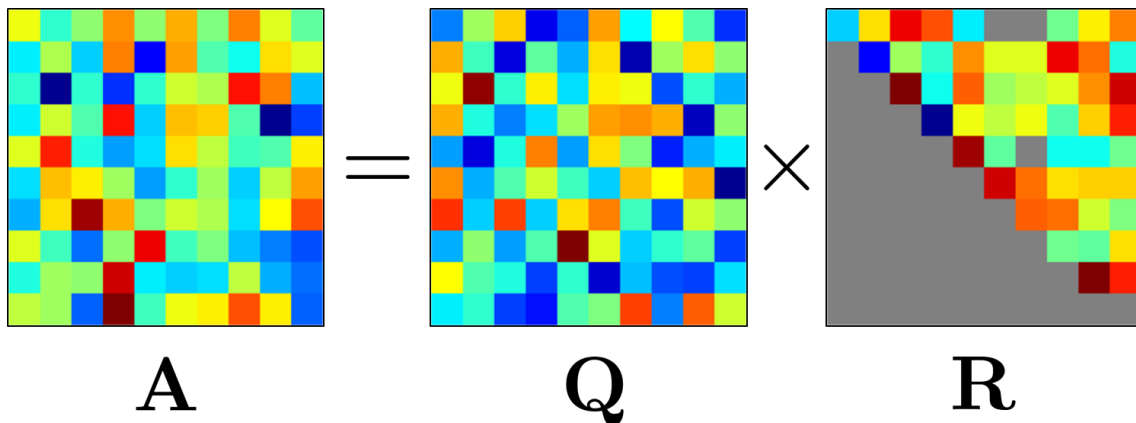


Figure 6 : Un exemple de factorisation QR, pour une matrice $\mathbf{A} \in \mathbb{R}^{10 \times 10}$. Les zones colorées correspondent à des entrées non-nulles, les grises à des entrées nulles.

Comme son homologue la factorisation LU, la factorisation QR permet de faciliter la résolution de systèmes linéaires. En effet, lorsque l'on dispose d'une factorisation

$\mathbf{A} = \mathbf{QR}$, le système $\mathbf{Ax} = \mathbf{b}$ est équivalent à $\mathbf{Rx} = \mathbf{c}$ avec $\mathbf{c} = \mathbf{Q}^T \mathbf{b}$ (\mathbf{Q} étant une matrice orthogonale), et on peut donc se ramener à la résolution d'un système triangulaire (en $\mathcal{O}(n^2)$).

On peut obtenir une factorisation QR de plusieurs manières. La première se résume à construire progressivement une base orthogonale de l'espace des colonnes de la matrice \mathbf{A} , elle est appelée procédé d'orthogonalisation de Gram-Schmidt. Deux autres méthodes existent, qui agissent par multiplications successives de la matrice \mathbf{A} par des matrices orthogonales élémentaires afin de la transformer progressivement en matrice triangulaire supérieure. Ces méthodes sont connues sous les noms de méthode de Householder (Householder, 1958) et méthode de Givens (Givens, 1958) selon le type de matrices élémentaires utilisé. Ces méthodes sont plus stables numériquement que le procédé de Gram-Schmidt. Effectuer une factorisation QR d'une matrice carrée coûte $\mathcal{O}(n^3)$ opérations arithmétiques, ce qui est du même ordre de grandeur que la factorisation LU. Cependant, utiliser la factorisation QR au lieu de la factorisation LU est avantageux pour la résolution de systèmes linéaires non-carrés, ou mal conditionnés. Une présentation plus détaillée des factorisation QR et LU, ainsi que des comparaisons précises en terme de coût de calcul et stabilité est disponible dans l'ouvrage de référence concernant le calcul matriciel (Golub et Van Loan, 2012).

2.2.2. Décomposition en valeurs singulières

Présentons maintenant une factorisation matricielle d'une extrême importance, tant sur le plan théorique qu'en terme d'applications : la décomposition en valeur singulière (souvent abrégée SVD, pour "Singular Value Decomposition"). Le but de la décomposition en valeurs singulière est d'obtenir des matrices $\mathbf{U} \in \mathbb{R}^{m \times m}$ et $\mathbf{V} \in \mathbb{R}^{n \times n}$ orthogonales et une matrice $\mathbf{\Sigma} \in \mathbb{R}^{m \times n}$ diagonale positive ($\sigma_{ij} = 0$ si $i \neq j$, $\sigma_{ii} \geq 0$) telles que $\mathbf{A} = \mathbf{U}\mathbf{\Sigma}\mathbf{V}^T$ (en général, les entrées diagonales de $\mathbf{\Sigma}$ sont ordonnées en ordre décroissant, telles que si $k \leq l$, $\sigma_{kk} \geq \sigma_{ll}$). La décomposition en valeurs singulières correspond donc au problème d'optimisation suivant :

$$\underset{\mathbf{U}, \mathbf{\Sigma}, \mathbf{V}}{\text{minimiser}} \quad d(\mathbf{A}, \mathbf{U}\mathbf{\Sigma}\mathbf{V}^T) + g_u(\mathbf{U}) + g_\sigma(\mathbf{\Sigma}) + g_v(\mathbf{V}),$$

où g_u et g_v sont indicatrices de l'ensemble des matrices orthogonales et g_σ est l'indicatrice de l'ensemble des matrices diagonales positives à coefficients diagonaux décroissants. Un exemple de décomposition en valeurs singulières est donné en figure 7.

La SVD a beaucoup d'applications, qu'il serait difficilement possible d'énumérer. Géométriquement parlant, la SVD donne accès à une base orthogonale pour l'espace engendré par les colonnes de \mathbf{A} (les colonnes de \mathbf{U}), et à une base orthogonale pour l'espace engendré par les lignes de \mathbf{A} (les colonnes de \mathbf{V}), bases dans lesquelles

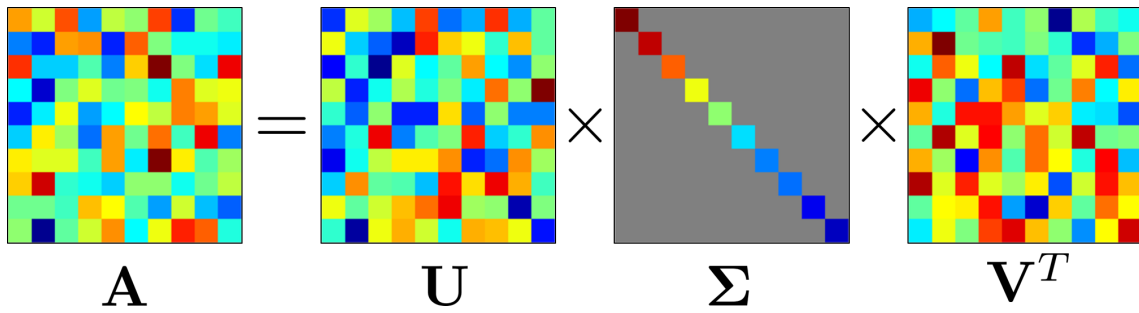


Figure 7 : Un exemple de décomposition en valeurs singulières pour une matrice $\mathbf{A} \in \mathbb{R}^{10 \times 10}$. Les zones colorées correspondent à des entrées non-nulles, les grises à des entrées nulles.

l'application linéaire associée à \mathbf{A} se résume à une simple dilatation. Elle permet donc d'interpréter la transformation linéaire associée à \mathbf{A} comme la composition d'une rotation (par l'application de \mathbf{V}^T), d'une dilatation (par l'application de $\mathbf{\Sigma}$), et d'une seconde rotation (par l'application de \mathbf{U}). La SVD permet aussi d'écrire $\mathbf{A} = \sum_{k=1}^{\min(m,n)} \sigma_{kk} \mathbf{u}_k \mathbf{v}_k^T$, la matrice \mathbf{A} étant décomposée en une somme de matrices de rang 1 (de la forme \mathbf{bc}^T) classées par ordre d'importance. Tronquer cette somme (ne garder que les r premiers termes) permet d'effectuer une approximation de rang faible de la matrice \mathbf{A} (qui se trouve être optimale, par le théorème d'Eckart-Young (Eckart et Young, 1936)). De plus, lorsque les colonnes de \mathbf{A} sont des vecteurs d'intérêt constituant un ensemble de données, les premières colonnes de \mathbf{U} peuvent être interprétées comme les directions principales (les plus importantes) de l'ensemble de données.

Calculer la SVD peut être fait de plusieurs manières, qui sont similaires à celles utilisées pour la factorisation QR, et coûtent $\mathcal{O}(n^3)$ opérations arithmétiques pour une matrice carrée. Pour plus de détails sur ces méthodes de calcul, voir (Golub et Van Loan, 2012).

2.3. Apprentissage de dictionnaire

Présentons maintenant l'apprentissage de dictionnaire, un domaine de recherche assez récent s'apparentant à un type de factorisation matricielle qui diffère des factorisations usuelles présentées lors de la section précédente.

2.3.1. Historique

Un signal discret résulte d'un nombre fini m de mesures d'une grandeur physique ou de tout phénomène quantifiable, il est représenté mathématiquement par un vec-

teur $\mathbf{a} \in \mathbb{R}^m$. Traditionnellement, le traitement du signal repose sur des changements de la base dans laquelle on exprime le signal, appelés transformées. On peut citer par exemple la transformée de Fourier (Fourier, 1822), mais aussi la transformée en ondelettes (Mallat, 1989) ou la transformée en cosinus discrète (Rao et Yip, 1990). Ces changements de base reviennent à exprimer le signal $\mathbf{a} = \mathbf{T}\boldsymbol{\gamma}$, où $\mathbf{T} \in \mathbb{R}^{n \times n}$ ($\mathbf{T} \in \mathbb{C}^{n \times n}$ pour la transformée de Fourier) est en général une matrice orthogonale dont les colonnes forment la base d'intérêt et $\boldsymbol{\gamma}$ est appelée la représentation du signal dans cette base. Afin que le changement de base ait un intérêt, la représentation $\boldsymbol{\gamma}$ du signal doit être plus "simple" que le signal lui-même, la notion de simplicité étant à définir au regard de l'application visée. Historiquement, la base \mathbf{T} était déterminée analytiquement, à partir des propriétés de la classe de signaux d'intérêt considérée. Au contraire, l'apprentissage de dictionnaire consiste à déterminer une matrice $\mathbf{D} \in \mathbb{R}^{m \times d}$ appelée le dictionnaire (avec en général $d > m$), à partir d'une collection de signaux $\mathbf{a}_1, \dots, \mathbf{a}_n$ appelée les données d'entraînement, telle que $\mathbf{a}_i \approx \mathbf{D}\boldsymbol{\gamma}_i$, $\forall i$. L'apprentissage de dictionnaire trouve ses racines dans un article fondateur (Olshausen et Field, 1996), dont la principale contribution est de montrer que le choix d'un dictionnaire à partir de petits morceaux d'images naturelles en imposant une représentation $\boldsymbol{\gamma}_i$ parcimonieuse (avec peu d'entrées non-nulles) résultait de colonnes reproduisant de manière fidèle le comportement des récepteurs du cortex visuel. Ceci constitue une justification biologique au fait de chercher des représentations dont la simplicité est mathématiquement traduite par la parcimonie. Pour une historique plus développée de l'apprentissage de dictionnaire, voir (Rubinstein et al., 2010a).

2.3.2. Formulation en tant que factorisation matricielle

Formulons le problème d'apprentissage de dictionnaire comme un problème de factorisation matricielle. Considérons une matrice $\mathbf{A} \in \mathbb{R}^{m \times n}$ constituée par la juxtaposition de n signaux d'entraînement de dimension m . Le but de l'apprentissage de dictionnaire est de trouver une représentation plus simple de ces signaux. Mathématiquement parlant, cette simplicité se traduit par le fait que les signaux doivent être bien approchés par une combinaison linéaire d'un petit nombre de vecteurs (appelés atomes) formant les colonnes d'une matrice $\mathbf{D} \in \mathbb{R}^{m \times d}$ appelée dictionnaire. Ceci s'écrit $\mathbf{a}_i \approx \sum_{j=1}^d \gamma_{ji} \mathbf{d}_j$, ou sous forme matricielle $\mathbf{a}_i \approx \mathbf{D}\boldsymbol{\gamma}_i$, où le nombre d'entrées non-nulles de $\boldsymbol{\gamma}_i$, noté $\|\boldsymbol{\gamma}_i\|_0$ et appelé norme ℓ_0 , est faible (on dit que $\boldsymbol{\gamma}_i$ est un vecteur creux, ou parcimonieux). Le nombre d'atomes d constituant le dictionnaire, ainsi que la parcimonie maximale tolérée des vecteurs $\boldsymbol{\gamma}_i$ (que l'on notera k) sont des paramètres du problème d'apprentissage de dictionnaire. En général, les dimensions suivent la relation $k \ll m < d \ll n$, de manière à ce que la représentation obtenue à l'aide du dictionnaire soit simple sans pour autant que le dictionnaire lui-même soit trop volumineux. L'apprentissage de dictionnaire

correspond donc au problème de factorisation matricielle suivant :

$$\underset{\mathbf{D}, \mathbf{\Gamma}}{\text{minimiser}} \quad \frac{1}{2} \|\mathbf{A} - \mathbf{D}\mathbf{\Gamma}\|_F^2 + g_d(\mathbf{D}) + g_\gamma(\mathbf{\Gamma}), \quad (\text{P1})$$

où g_d est une pénalité permettant d'imposer une certaine structure au dictionnaire, g_γ une pénalité imposant la parcimonie aux colonnes de la matrice des coefficients $\mathbf{\Gamma}$, et la mesure de distance choisie $\|\cdot\|_F^2$ correspond au carré de la norme de Frobenius ($\|\mathbf{E}\|_F^2 = \sum_{i,j} e_{ij}^2$). Contrairement aux factorisations usuelles présentées précédemment, l'apprentissage de dictionnaire correspond à une factorisation approchée de la matrice \mathbf{A} . Un exemple d'apprentissage de dictionnaire est donné en figure 8.

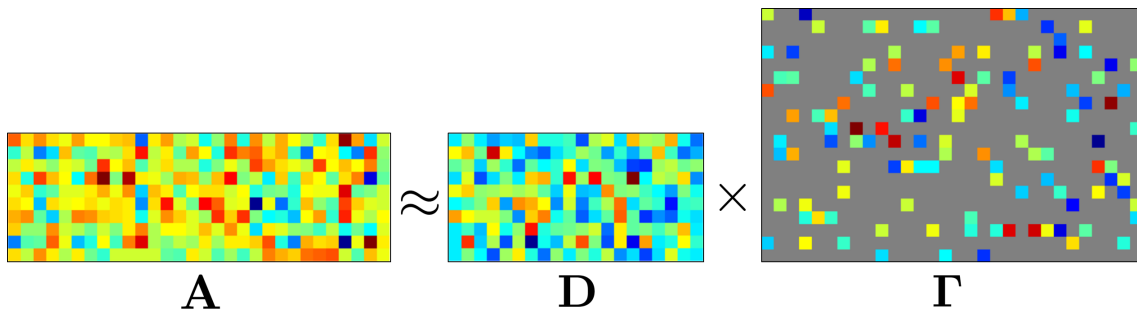


Figure 8 : Un exemple d'apprentissage de dictionnaire pour une matrice $\mathbf{A} \in \mathbb{R}^{10 \times 30}$. Le dictionnaire est constitué de $d = 20$ atomes et chaque colonne de \mathbf{A} est approchée par une combinaison linéaire de $k = 4$ atomes de ce dictionnaire. Les zones colorées correspondent à des entrées non-nulles, les grises à des entrées nulles.

2.3.3. Résolution du problème

Examinons maintenant plus en détails le problème d'apprentissage de dictionnaire (P1), et diverses façons de le résoudre de manière approchée.

Stratégie générale

La principale difficulté de ce problème vient du fait que le terme $\frac{1}{2} \|\mathbf{A} - \mathbf{D}\mathbf{\Gamma}\|_F^2$ (appelé attache aux données) est non-convexe lorsque l'on considère les variables \mathbf{D} et $\mathbf{\Gamma}$ ensemble. Cependant, il est convexe par rapport à chacune de ces variables lorsque l'autre est fixée. Les principaux algorithmes d'apprentissage de dictionnaire profitent de cette convexité partielle, et considèrent séparément le sous-problème correspondant aux coefficients :

$$\underset{\mathbf{\Gamma}}{\text{minimiser}} \quad \frac{1}{2} \|\mathbf{A} - \mathbf{D}\mathbf{\Gamma}\|_F^2 + g_\gamma(\mathbf{\Gamma}), \quad (\text{SP1})$$

et celui correspondant au dictionnaire :

$$\underset{\mathbf{D}}{\text{minimiser}} \quad \frac{1}{2} \|\mathbf{A} - \mathbf{D}\mathbf{\Gamma}\|_F^2 + g_d(\mathbf{D}). \quad (\text{SP2})$$

Ainsi, les algorithmes d'apprentissage de dictionnaire procèdent par mise à jour alternée des coefficients $\mathbf{\Gamma}$ (considérant le sous-problème (SP1)) et du dictionnaire \mathbf{D} (considérant le sous-problème (SP2)). Ce principe de mise à jour alternée est le même que celui utilisé dans la méthode de descente par bloc de coordonnées présentée au chapitre précédent (les blocs correspondant aux matrices \mathbf{D} et $\mathbf{\Gamma}$), et est donné en algorithme 8.

Algorithme 8 Algorithme d'apprentissage de dictionnaire général

Entrée : Une matrice de données \mathbf{A} , un dictionnaire initial \mathbf{D} , un critère d'arrêt c .

1: **while** $c = \text{False}$ **do**

2: Mettre à jour les coefficients en considérant (SP1) : $\mathbf{\Gamma} \leftarrow \hat{\mathbf{\Gamma}}$

3: Mettre à jour le dictionnaire en considérant (SP2) : $\mathbf{D} \leftarrow \hat{\mathbf{D}}$

4: **end while**

Sortie : Le dictionnaire \mathbf{D} et les coefficients $\mathbf{\Gamma}$.

Résolution des sous-problèmes

Concrètement, la pénalité g_γ imposant la parcimonie des coefficients est soit la fonction indicatrice des matrices dont les colonnes ont moins de k entrées non-nulles, soit une pénalité de la forme $\lambda_0 \sum_{i=1}^n \|\gamma_i\|_0$. Ces formes particulières de pénalité rendent le sous-problème (SP1) non-convexe, et même très difficile à résoudre exactement. Malgré tout, on peut le résoudre de manière approchée, ou le remplacer par un problème approché plus facile à résoudre exactement, par exemple en considérant $\hat{g}_\gamma(\mathbf{\Gamma}) = \lambda_1 \sum_{i=1}^n \|\gamma_i\|_1$ à la place de g_γ , où $\|\mathbf{x}\|_1 = \sum_i |x_i|$ est appelée la norme ℓ_1 , et est connue pour favoriser la parcimonie. Ce type de problème est connu sous le nom de codage parcimonieux (“sparse coding” en anglais), et de multiples méthodes de résolution approchées ont été développées, notamment dans les deux dernières décennies (voir (Tropp et Wright, 2010) pour un catalogue détaillé de méthodes). On peut citer par exemple les algorithmes “Matching Pursuit” (MP) (Mallat et Zhang, 1993), “Orthogonal Matching Pursuit” (OMP) (Tropp et Gilbert, 2007), “Compressive Sampling Matching Pursuit” CoSaMP (Needell et Tropp, 2009) ou “Iterative Hard Thresholding” IHT (Blumensath et Davies, 2008) pour le problème avec la norme ℓ_0 . De même, on peut citer les algorithmes “Iterative Soft Thresholding” (IST) (Daubechies et al., 2004) et “Fast Iterative Shrinkage-Thresholding Algorithm” (FISTA) (Beck et Teboulle, 2009) pour le problème avec la norme ℓ_1 .

En ce qui concerne le sous-problème (SP2), la pénalité g_d varie en fonction de la structure souhaitée pour le dictionnaire. Dans les cas les plus simples (par exemple si

g_d est l'indicatrice des matrices dont les colonnes sont normalisées), le sous-problème est convexe, mais peut être coûteux à résoudre de manière exacte, comme nous allons le voir dans les paragraphes suivants.

Algorithmes

Les algorithmes d'apprentissage de dictionnaire diffèrent dans leur façon de mettre à jour les coefficients et le dictionnaire. Nous présentons ici par ordre chronologique trois algorithmes d'apprentissage de dictionnaire représentatifs de l'évolution de ce domaine.

MOD. La méthode des directions optimales (Engan et al., 1999) (abrégée MOD pour “Method of Optimal Directions”) est un des premiers algorithmes d'apprentissage de dictionnaire connus, et un des plus simples. Elle consiste tout d'abord à mettre à jour les coefficients en résolvant de manière approchée (SP1) en utilisant un algorithme de codage parcimonieux (Tropp et Wright, 2010). Ensuite le dictionnaire est mis à jour en résolvant de manière exacte (SP2), ce qui revient à calculer une matrice appelée la pseudo-inverse de la matrice $\mathbf{\Gamma}$. Ce type de mise à jour du dictionnaire a un coût élevé, en $\mathcal{O}(nd^2)$ ($\mathcal{O}(x)$ signifiant de l'ordre de x).

KSVD. KSVD (Aharon et al., 2006) est sûrement l'algorithme d'apprentissage de dictionnaire le plus utilisé. Il utilise la même mise à jour des coefficients que la méthode des directions optimales. Le dictionnaire est lui mis à jour colonne par colonne en utilisant le fait que $\mathbf{D}\mathbf{\Gamma} = \sum_{j=1}^d \mathbf{d}_j(\boldsymbol{\gamma}^j)^T$, qui permet de ramener le sous-problème à la résolution du problème suivant pour tout i :

$$\underset{\mathbf{d}_i, \boldsymbol{\gamma}^i}{\text{minimiser}} \quad \|\mathbf{E}_i - \mathbf{d}_i(\boldsymbol{\gamma}^i)^T\|_F^2,$$

où $\mathbf{E}_i = \mathbf{A} - \sum_{j \neq i} \mathbf{d}_j(\boldsymbol{\gamma}^j)^T$. Ce problème est alors résolu via une approximation de rang 1 prenant en compte la parcimonie des coefficients $\boldsymbol{\gamma}^i$. La solution est donnée par la SVD tronquée de \mathbf{E}_i à laquelle on a enlevé les colonnes correspondant aux entrées nulles de $\boldsymbol{\gamma}^i$ (d'où l'acronyme SVD dans le nom de l'algorithme). KSVD est moins coûteux que la méthode des directions optimales en raison de la mise à jour atome par atome du dictionnaire. En effet, la complexité de ce type de mise à jour du dictionnaire varie en $\mathcal{O}(ndk)$ (où k est la parcimonie de chaque colonne de la matrice $\mathbf{\Gamma}$) (Rubinstein et al., 2008).

ODL. Récemment introduite, la méthode d'apprentissage de dictionnaire en ligne (Mairal et al., 2010) (que l'on abrège ODL pour “Online Dictionary Learning”) se distingue des deux algorithmes précédents. En effet, les signaux y sont considérés un

par un (ou par petits groupes), ce qui revient à ne prendre en compte qu'une colonne de la matrice \mathbf{A} à chaque itération. Les coefficients correspondant au signal courant sont calculés à l'aide d'un algorithme de codage parcimonieux, puis le dictionnaire est mis à jour de manière à prendre en compte l'arrivée du nouveau signal. La mise à jour du dictionnaire est effectuée de manière exacte (à l'aide d'une résolution par bloc de coordonnées, atome par atome), ce qui est rendu possible du fait de la proximité des coefficients d'une itération à l'autre (seule la colonne correspondant au signal courant est ajoutée, les autres étant fixées). ODL est un algorithme très peu coûteux, particulièrement adapté au cas où le nombre de signaux d'entraînement n est très élevé, voire infini (lorsque l'on considère que de nouveaux signaux sont pris en compte au fil des itérations). Une telle mise à jour du dictionnaire ne coûte en effet que $\mathcal{O}(md^2)$ opération arithmétique.

Les algorithmes présentés ici mettent en évidence une tendance à considérer des mises à jour du dictionnaire de moins en moins coûteuses. Cependant, les dictionnaires appris par ces méthodes sont des matrices en général denses et coûteuses à manipuler en grande dimension. D'autres méthodes plus récentes introduites en détails dans la partie contributions de cette thèse (au chapitres 4 et 8) sont spécialement conçues pour obtenir des dictionnaires moins coûteux à manipuler, ce qui est aussi l'un des objectifs de cette thèse. Pour une introduction plus conséquente à l'apprentissage de dictionnaire, avec des exemples d'applications, voir les articles de revue ([Rubinstein et al., 2010a](#); [Tosic et Frossard, 2011](#)).

2.4. Réseaux de neurones

Nous présentons dans cette section un modèle introduit dans le domaine de l'apprentissage automatique n'étant pas à proprement parler une factorisation matricielle, mais partageant de nombreuses similarités avec elles, ainsi qu'avec les modèles examinés dans cette thèse, notamment du point de vue de l'optimisation : les réseaux de neurones artificiels. Nous nous appuyons durant toute cette section sur le récent ouvrage consacré aux réseaux de neurones profonds ([Bengio et Courville, 2016](#)).

2.4.1. Principe

Le but des réseaux de neurones est d'effectuer une tâche de prédiction. Supposons qu'il existe une transformation $f : \mathbb{R}^p \rightarrow \mathbb{R}^m$ latente liant des entrées $\mathbf{x} \in \mathbb{R}^p$ à des sorties $\mathbf{a} = f(\mathbf{x}) \in \mathbb{R}^m$. L'objectif d'un réseau de neurones artificiels est de construire une transformation \hat{f} proche de f , afin de prédire les valeurs de la sortie à partir de l'entrée. Cette fonction de prédiction \hat{f} est construite d'une manière inspirée de la biologie, et de la façon dont les neurones traitent les signaux électriques.

Les neurones sont organisés en J couches successives, la première au contact de l'entrée et la dernière au contact de la sortie. Chaque neurone agit en effectuant une combinaison linéaire des activations des neurones de la couche précédente, puis applique une non-linéarité afin d'obtenir sa propre activation. La transformation liant deux couches successives s'écrit donc $\boldsymbol{\alpha}_j = \sigma_j(\mathbf{W}_j \boldsymbol{\alpha}_{j-1})$, où $\mathbf{W}_j \in \mathbb{R}^{n_j \times n_{j-1}}$ est la matrice de poids représentant la connexion entre les couches, la non-linéarité $\sigma_j(\cdot)$ s'applique coordonnée par coordonnée, n_j correspond au nombre de neurones et $\boldsymbol{\alpha}_j$ aux activations de la j -ème couche. Considérant toutes les couches, La fonction de prédiction globale prend la forme $\hat{f}(\mathbf{x}) \triangleq \sigma_J(\mathbf{W}_J \dots \sigma_2(\mathbf{W}_2 \sigma_1(\mathbf{W}_1 \mathbf{x})) \dots)$. Une fois l'architecture du réseau fixée (nombre de couches, taille des couches et forme des non-linéarités), les paramètres du réseau (les poids $\mathbf{W}_1, \dots, \mathbf{W}_J$) sont à choisir dans le but d'approcher au mieux la transformation latente f . On dispose pour cela de données pour entraîner le réseau, sous la forme de n couples entrée/sortie $(\mathbf{x}_i, \mathbf{a}_i = f(\mathbf{x}_i))$, $i = 1, \dots, n$. L'apprentissage des paramètres se fait alors par le biais de la résolution approchée du problème d'optimisation suivant :

$$\underset{\mathbf{W}_1, \dots, \mathbf{W}_J}{\text{minimiser}} \quad \frac{1}{2n} \sum_{i=1}^n \|\mathbf{a}_i - \sigma_J(\mathbf{W}_J \dots \sigma_2(\mathbf{W}_2 \sigma_1(\mathbf{W}_1 \mathbf{x}_i)) \dots)\|_2^2 + \sum_{j=1}^J g_j(\mathbf{W}_j), \quad (\text{P2})$$

où le terme d'attache au données correspond à la moyenne de l'erreur de prédiction commise sur les données d'entraînement (mesurée ici avec la norme euclidienne, mais d'autres choix sont possibles), et les pénalités g_j , $j = 1, \dots, J$ permettent d'imposer une structure aux poids. On espère que les poids ainsi appris permettent d'effectuer de bonnes prédictions non seulement sur les données d'entraînement mais aussi au-delà, sur de nouvelles données. Ce problème d'optimisation peut se reformuler en considérant les matrices $\mathbf{X} \in \mathbb{R}^{p \times n}$ formée par la concaténation des entrées et $\mathbf{A} \in \mathbb{R}^{m \times n}$ formée par la concaténation des sorties :

$$\underset{\mathbf{W}_1, \dots, \mathbf{W}_J}{\text{minimiser}} \quad \frac{1}{2n} \|\mathbf{A} - \sigma_J(\mathbf{W}_J \dots \sigma_2(\mathbf{W}_2 \sigma_1(\mathbf{W}_1 \mathbf{X})) \dots)\|_F^2 + \sum_{j=1}^J g_j(\mathbf{W}_j),$$

où l'on voit bien les similarités avec les problèmes de factorisation matricielle présentés dans les sections précédentes (la seule différence étant la présence des non-linéarités). Un exemple de réseau de neurones est donné en figure 9, mettant en évidence les similarités avec les factorisations matricielles.

2.4.2. Optimisation

Les réseaux de neurones artificiels ont rencontré de nombreux succès, spécialement ces dernières années. Ils existent pourtant depuis plusieurs décennies, mais leur

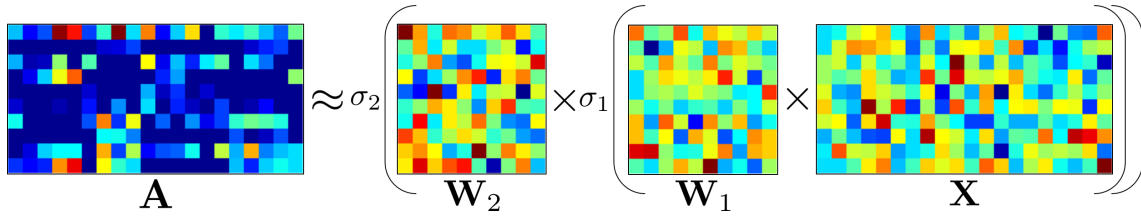


Figure 9 : Un exemple de réseau de neurones à 2 couches avec des poids $\mathbf{W}_1, \mathbf{W}_2 \in \mathbb{R}^{10 \times 10}$. On représente ici les données d'entraînement constituée de $n = 20$ exemples, rassemblés dans les matrices $\mathbf{A}, \mathbf{X} \in \mathbb{R}^{10 \times 20}$.

développement a été longtemps ralenti par la difficulté du problème d'optimisation (P2) associé. Ce problème est en effet non-convexe, et intrinsèquement coûteux car dépendant d'un grand nombre de paramètres (l'optimisation se faisant sur J matrices pleines, potentiellement de grande dimension). De nombreuses techniques et heuristiques ont été développées ces dernières années permettant de surmonter ces deux difficultés, si bien qu'aujourd'hui les réseaux de neurones constituent l'état de l'art dans de nombreuses tâches de reconnaissance automatique et de classification (LeCun et al., 2015). Nous présentons ici certaines de ces techniques, en particulier celles pouvant être adaptées au problème étudié dans cette thèse, que nous introduisons dans la partie contributions, au chapitre 4.

Descente de gradient stochastique. L'algorithme d'optimisation le plus utilisé pour l'entraînement des réseaux de neurones est très simple, il s'agit d'une descente de gradient. Cependant, pour alléger le coût, on ne considère qu'un petit nombre l d'exemples à la fois, on effectue alors ce qu'on appelle une descente de gradient stochastique (présentée au chapitre 1). Afin d'être plus clair, rassemblons tous les paramètres du réseau dans un même vecteur $\boldsymbol{\theta} \triangleq (\mathbf{W}_1, \dots, \mathbf{W}_J)$, et définissons le coût individuel correspondant à l'exemple (\mathbf{x}, \mathbf{a}) considérant ces paramètres : $L(\hat{f}(\mathbf{x}, \boldsymbol{\theta}), \mathbf{a}) = \frac{1}{2} \|\mathbf{a} - \sigma_J(\mathbf{W}_J \dots \sigma_2(\mathbf{W}_2 \sigma_1(\mathbf{W}_1 \mathbf{x})) \dots)\|_2^2 + \sum_{j=1}^J g_j(\mathbf{W}_j)$. On peut remarquer que l'objectif du problème (P2) à minimiser est une moyenne de ces coûts individuels et prend la forme $\frac{1}{n} \sum_{i=1}^n L(\hat{f}(\mathbf{x}_i, \boldsymbol{\theta}), \mathbf{a}_i)$. Ceci n'est autre qu'une estimation de l'espérance du coût individuel par rapport à la distribution de probabilités sous-jacente des entrées/sorties. Il est alors possible de calculer cette estimation en effectuant une moyenne sur moins d'exemple (l en l'occurrence), afin de réduire le coût de calcul du gradient, ce qui revient à considérer un sous-objectif de la forme $\frac{1}{l} \sum_{i \in \mathcal{B}} L(\hat{f}(\mathbf{x}_i, \boldsymbol{\theta}), \mathbf{a}_i)$, avec $|\mathcal{B}| = l$. Le choix du sous-ensemble \mathcal{B} des exemples considérés change à chaque itération de l'algorithme, afin de profiter de toute la diversité de l'ensemble d'entraînement (Bottou, 1998). L'algorithme de descente de gradient stochastique pour les réseaux de neurones est donné en algorithme 9.

Algorithme 9 Descente de gradient stochastique pour les réseaux de neurones

Entrée : Des données \mathbf{A} , \mathbf{X} , des paramètres initiaux $\boldsymbol{\theta}$, un critère d'arrêt c .

- 1: **while** $c = \text{False}$ **do**
- 2: Choisir l exemples d'entraînement $\{(\mathbf{x}_i, \mathbf{a}_i)\}_{i \in \mathcal{B}}$
- 3: Calculer le gradient : $\mathbf{g} \leftarrow \frac{1}{l} \sum_{i \in \mathcal{B}} \nabla_{\boldsymbol{\theta}} L(\hat{f}(\mathbf{x}_i, \boldsymbol{\theta}), \mathbf{a}_i)$
- 4: Choisir un pas λ
- 5: Mettre à jour les paramètres : $\boldsymbol{\theta} \leftarrow \boldsymbol{\theta} - \lambda \mathbf{g}$
- 6: **end while**

Sortie : Les paramètres $\boldsymbol{\theta}$.

Back-propagation. Malgré l'utilisation de l'algorithme de descente de gradient stochastique, entraîner un réseau de neurones reste coûteux, du fait du nombre important de paramètres à mettre à jour (les poids $\mathbf{W}_1, \dots, \mathbf{W}_J$). Il s'agit donc de calculer le gradient \mathbf{g} par rapport à l'ensemble des paramètres de manière intelligente afin de ne pas gaspiller de ressources. Un tel calcul intelligent est possible en mettant à profit la structure particulière du réseau, par une technique connue sous le nom de "back-propagation" (forme abrégée de "backward propagation", signifiant propagation arrière en anglais). En considérant un seul exemple d'entraînement (\mathbf{x}, \mathbf{a}) (le calcul doit être fait en parallèle sur plusieurs exemples pour construire \mathbf{g}), la technique de back-propagation consiste à calculer le vecteur de l'erreur commise par le réseau sur l'exemple considéré $\mathbf{e} \triangleq \sigma_J(\mathbf{W}_J \dots \sigma_2(\mathbf{W}_2 \sigma_1(\mathbf{W}_1 \mathbf{x})) \dots) - \mathbf{a}$ (cette phase est appelée "forward propagation"). Il est ensuite simple à partir de ce vecteur de calculer le gradient par rapport aux poids de la dernière couche, puis de réutiliser des calculs intermédiaires pour calculer le gradient par rapport à la couche précédente et ainsi de suite (phase de back-propagation). L'algorithme se base simplement sur l'application de la règle de dérivation des fonctions composées, permettant de ne pas recalculer inutilement plusieurs fois les mêmes quantités présentes dans les expressions des gradients par rapport aux couches successives du réseau. Une implémentation concrète de l'algorithme de back-propagation dans le cas du réseau pris en exemple dans cette section est donné en algorithme 10 (où le symbole \odot correspond au produit coordonnée par coordonnée).

Pré-entraînement. La non-convexité de la fonction de coût utilisée pour l'entraînement des réseaux de neurones implique que les algorithmes de descente de gradient peuvent converger vers des minima locaux. Ceci signifie que les performances d'un réseau de neurones dépendent de l'initialisation des poids lors de l'entraînement. Il est donc nécessaire d'initialiser les poids dans une région favorable afin d'obtenir le réseau le plus performant possible. À ce titre, l'astuce la plus célèbre qui a causé un regain d'intérêt pour les réseaux de neurones à la fin des années 2000 est le pré-entraînement couche par couche ([Hinton et Salakhutdinov](#),

Algorithme 10 Back-propagation considérant un exemple d'entraînement

Entrée : Un exemple d'entraînement (\mathbf{x}, \mathbf{a}) , un réseau de neurone avec ses paramètres $\boldsymbol{\theta} = (\mathbf{W}_1, \dots, \mathbf{W}_J)$ et ses activations $\boldsymbol{\alpha}_1, \dots, \boldsymbol{\alpha}_J$.

- 1: Calculer l'erreur commise : $\mathbf{g}_J \leftarrow \mathbf{e} \triangleq \sigma_J(\mathbf{W}_J \dots \sigma_2(\mathbf{W}_2 \sigma_1(\mathbf{W}_1 \mathbf{x})) \dots) - \mathbf{a}$
- 2: **for** $j = J : -1 : 1$ **do**
- 3: Prendre en compte l'effet de la non-linéarité : $\mathbf{g}_j \leftarrow \mathbf{g}_j \odot \sigma'_j(\boldsymbol{\alpha}_j)$
- 4: Calculer le gradient par rapport aux paramètres de la couche courante :
 $\nabla_{\mathbf{w}_j} L(\hat{f}(\mathbf{x}, \boldsymbol{\theta}), \mathbf{a}) \leftarrow \mathbf{g}_j \boldsymbol{\alpha}_{j-1}^T + \nabla_{\mathbf{w}_j} g_j(\mathbf{W}_j)$
- 5: Propager à la couche précédente : $\mathbf{g}_{j-1} \leftarrow \mathbf{W}_j^T \mathbf{g}_j$
- 6: **end for**

Sortie : Le gradient par rapport à tous les paramètres en considérant l'exemple fourni : $\nabla_{\boldsymbol{\theta}} L(\hat{f}(\mathbf{x}, \boldsymbol{\theta}), \mathbf{a}) = (\nabla_{\mathbf{w}_1} L, \dots, \nabla_{\mathbf{w}_J} L)$.

2006; Bengio et al., 2007). Il consiste à pré-entraîner les couches du réseau individuellement de la première à la J -ème, de manière à ce que la sortie de chaque couche permette de bien reconstruire son entrée. Pour cela, lors du pré-entraînement, l'entrée de la j -ème couche est prise comme la sortie (les activations) de la $j - 1$ -ème couche préalablement pré-entraînée (la première couche disposant des données d'entraînement en entrée). Le pré-entraînement impose donc une initialisation dans une région où toutes les couches "préservent l'information". Il est suivi de l'entraînement classique par descente de gradient stochastique (appelé "global fine tuning"). Il existe d'autres techniques permettant de lutter contre la non-convexité de la fonction de coût, concernant par exemple le choix du pas utilisé pour la descente de gradient, ou l'introduction d'inertie entre les itérations, que nous n'évoquons pas ici, voir (Bengio et Courville, 2016) pour une revue exhaustive.

Conclusion. Nous avons présenté dans ce chapitre de manière unifiée les factorisations matricielles, en insistant sur le problème d'optimisation résolu, explicitement ou non. Nous avons ainsi discuté un certain nombre de factorisations matricielles, illustrant leur grande diversité, tant sur le plan des objectifs que sur celui des techniques utilisées. Nous avons mis l'accent sur les factorisations ayant un lien avec le problème principal de cette thèse (introduit au chapitre 4), étant lui-même un problème de factorisation matricielle.

Matrices efficaces

Ce chapitre a pour but principal d'introduire l'objet mathématique éponyme de la thèse : les matrices efficaces. La notion d'efficacité étant étroitement liée à la complexité du produit matrice/vecteur, nous présentons tout d'abord les notions d'algorithme et de circuit linéaire. Après avoir défini l'efficacité, nous discutons la possibilité de la calculer et de la borner.

Sommaire

3.1	Circuit linéaire et complexité	41
3.1.1	Définitions	42
3.1.2	Exemples	43
3.2	Efficacité	46
3.2.1	Efficacité d'un circuit linéaire	46
3.2.2	Efficacité d'une matrice	47
3.3	Calculer et borner l'efficacité	47
3.3.1	Calcul	48
3.3.2	Bornes	48

3.1. Circuit linéaire et complexité

Ayant défini le produit matrice/vecteur dans le chapitre précédent, nous nous intéressons ici à la quantité de calcul nécessaire pour effectuer ce produit. Pour cela, nous étudions l'algorithme qui calcule effectivement ce produit, via la notion de circuit linéaire.

3.1.1. Définitions

Considérant une matrice \mathbf{A} et un vecteur d'entrée \mathbf{x} , l'algorithme qui calcule le produit $\mathbf{y} = \mathbf{A}\mathbf{x}$ est simplement la suite d'opérations arithmétiques paramétrées par les coefficients de \mathbf{A} , permettant d'obtenir les coefficients de \mathbf{y} à partir de ceux de \mathbf{x} . Nous considérons ici seulement des *algorithmes linéaires*, dont les instructions comprennent uniquement des additions et des multiplications par des constantes (paramétrées par les coefficients de \mathbf{A}). Tout algorithme linéaire est associé à un *circuit linéaire* (Lokam, 2009) qui permet de mieux visualiser l'algorithme. Les circuits linéaires sont des objets mathématiques utilisant un autre objet, un graphe, dont nous donnons la définition ici.

Définition 2. *Un graphe orienté G est un couple $(\mathcal{S}, \mathcal{A})$ où \mathcal{S} est appelé l'ensemble des sommets du graphe et $\mathcal{A} \subseteq \{(x, y) \in \mathcal{S}^2, x \neq y\}$ est appelé l'ensemble des arcs du graphe. Un chemin de longueur l dans le graphe G est une suite d'arcs $a_1 = (d_1, f_1), \dots, a_l = (d_l, f_l)$ telle que $\forall i \in \{1, \dots, l-1\}, f_i = d_{i+1}$. Le départ et l'arrivée du chemin sont, respectivement, les sommets d_1 et f_l . Un graphe orienté G est dit **acyclique** s'il n'existe pas de chemin dans G ayant pour départ et arrivée le même sommet. Un sommet s est dit **sommet d'entrée** du graphe G s'il n'existe pas de chemin dans G ayant pour arrivée s . De même, Un sommet s est dit **sommet de sortie** du graphe G s'il n'existe pas de chemin dans G ayant pour départ s . L'ensemble des arcs entrants du sommet s est l'ensemble $\{(x, y) \in \mathcal{A}, y = s\}$. De même, L'ensemble des arcs sortants du sommet s est l'ensemble $\{(x, y) \in \mathcal{A}, x = s\}$.*

Un graphe non orienté G est un couple $(\mathcal{S}, \mathcal{A})$ où $\mathcal{A} \subseteq \mathcal{P}_2(\mathcal{S})$, où $\mathcal{P}_2(\mathcal{S})$ est l'ensemble des parties à deux éléments de \mathcal{S} , est appelé l'ensemble des arêtes du graphe (ce sont des arcs non orientés).

Cette définition des graphes nous permet d'introduire maintenant la notion de circuit linéaire.

Définition 3. *Un circuit linéaire C est un graphe orienté acyclique $G_C = (\mathcal{S}, \mathcal{A})$ dans lequel chaque arc $a \in \mathcal{A}$ est associé à un nombre non-nul. Si $s \in \mathcal{S}$ est un sommet ayant ses arcs entrants associés aux nombres $\lambda_1, \dots, \lambda_k$ provenant des sommets s_1, \dots, s_k , alors s calcule la valeur $v(s) = \lambda_1 v(s_1) + \dots + \lambda_k v(s_k)$, où $v(s_i)$ est la valeur calculée au sommet s_i (en pratique, par abus de notation, on confondra souvent un sommet et sa valeur).*

La complexité (ou taille) d'un circuit linéaire correspond à son nombre d'arcs $|\mathcal{A}|$. La profondeur d'un circuit correspond à la longueur du plus long chemin d'un sommet d'entrée à un sommet de sortie du circuit.

*Un circuit linéaire C à n entrées et m sorties est dit **associé** à la matrice $\mathbf{A} \in \mathbb{R}^{m \times n}$ si et seulement si, en désignant \mathbf{x} le vecteur rassemblant les valeurs des sommets d'entrées et \mathbf{y} le vecteur rassemblant les valeurs des sommets de sortie, on a $\mathbf{y} = \mathbf{A}\mathbf{x}$,*

pour toute valeur de \mathbf{x} . On pourra alors noter $C = C_{\mathbf{A}}$. L'ensemble des circuits linéaires associés à une matrice \mathbf{A} sera noté $\mathcal{C}_{\mathbf{A}}$.

Cette définition des circuits linéaires peut sembler complexe, mais se révèle très intuitive, comme illustré dans la sous-section suivante. Il est à noter que la complexité d'un circuit linéaire est toujours du même ordre de grandeur que le nombre d'opérations arithmétiques induites par ce circuit (Lokam, 2009).

3.1.2. Exemples

Nous exposons ici divers circuits associés à des matrices particulières. Nous commençons par présenter un circuit universel qui s'associe à toute matrice quelconque. Nous présentons ensuite des circuits associés à des matrices de rang faible, puis à des matrices correspondant à des transformées rapides.

Circuit naïf

Considérant une matrice quelconque $\mathbf{A} \in \mathbb{R}^{m \times n}$, la manière la plus directe de calculer le produit de \mathbf{A} avec un vecteur quelconque $\mathbf{x} \in \mathbb{R}^n$ revient simplement à calculer explicitement $y_i = \sum_{k=1}^n a_{ik}x_k$, pour tout $i \in \{1, \dots, m\}$. En terme de circuit linéaire, cette manière de calculer le produit matrice/vecteur revient à considérer les sommets $\mathcal{S} = \{x_1, \dots, x_n, y_1, \dots, y_m\}$, et les arcs $\mathcal{A} = \{(x_j, y_i), \forall i \in \{1, \dots, m\}, j \in \{1, \dots, n\}\}$, avec l'arc (x_j, y_i) associé au nombre a_{ij} . Ce circuit est appelé *circuit naïf* dans le reste de ce document, il peut être associé à n'importe quelle matrice quelconque. Le circuit naïf associé à une matrice quelconque $\mathbf{A} \in \mathbb{R}^{3 \times 4}$ est représenté en Figure 10.

Le circuit naïf est de profondeur 1, et de complexité $|\mathcal{A}| = mn$. En terme d'opérations arithmétiques, un tel circuit induit exactement mn multiplications et $m(n - 1)$ additions.

Circuit associé aux matrices de rang faible

Considérons maintenant une matrice $\mathbf{A} \in \mathbb{R}^{m \times n}$ de rang au plus r . Cela signifie qu'on peut écrire $\mathbf{A} = \mathbf{B}\mathbf{C}^T$, avec $\mathbf{B} \in \mathbb{R}^{m \times r}$ et $\mathbf{C} \in \mathbb{R}^{n \times r}$. Dans ce cas, le calcul de $\mathbf{y} = \mathbf{A}\mathbf{x} = \mathbf{B}\mathbf{C}^T\mathbf{x}$ peut être effectué en deux étapes. En effet, il est possible de calculer en premier lieu $\mathbf{z} = \mathbf{C}^T\mathbf{x}$, puis dans un second temps $\mathbf{y} = \mathbf{B}\mathbf{z}$. Chacune des ces deux étapes peut être effectuée avec le circuit naïf. Une telle décomposition revient à considérer un circuit linéaire avec les sommets $\mathcal{S} = \{x_1, \dots, x_n, z_1, \dots, z_r, y_1, \dots, y_m\}$, et les arcs $\mathcal{A} = \{(x_j, z_i), \forall i \in \{1, \dots, r\}, j \in \{1, \dots, n\}\} \cup \{(z_k, y_l), \forall k \in \{1, \dots, r\}, l \in \{1, \dots, m\}\}$ avec l'arc (x_j, z_i) associé au nombre c_{ji} et l'arc (z_k, y_l) associé au nombre b_{lk} . Ce circuit est représenté en Figure 11 dans le cas $m = 3$, $n = 4$ et $r = 1$.

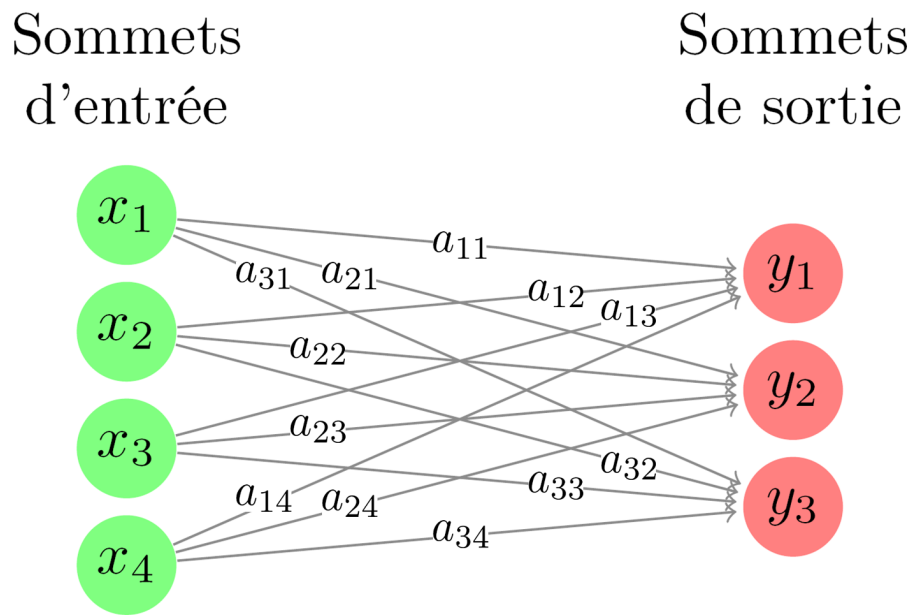


Figure 10 : Circuit naïf associé à une matrice quelconque $\mathbf{A} \in \mathbb{R}^{3 \times 4}$.

Ce circuit associé aux matrices de rang faible, correspondant à une factorisation de la matrice \mathbf{A} en deux matrices de plus petite taille \mathbf{B} et \mathbf{C} , est de profondeur 2, et de complexité $|\mathcal{A}| = r(m+n)$. Il est donc profitable, en terme d'opérations arithmétiques, d'utiliser ce circuit à la place du circuit naïf dès lors que $r < \frac{mn}{m+n}$.

Circuit associé aux transformées rapides

Les transformées sont des changements de base, utilisés dans la quasi-totalité des méthodes de traitement du signal. En dimension finie, appliquer une transformée (ou son inverse) revient à multiplier un vecteur d'entrée par une matrice carrée $\mathbf{A} \in \mathbb{K}^{n \times n}$, où $\mathbb{K} = \mathbb{R}$ ou $\mathbb{K} = \mathbb{C}$. La transformée la plus utilisée reste certainement la transformée de Fourier (Fourier, 1822), qui permet d'accéder au contenu fréquentiel du signal d'entrée. On peut citer aussi la transformée en ondelettes (Mallat, 2008), la transformée en cosinus discrète (Rao et Yip, 1990) (DCT pour *Discrete Cosine Transform*) ou la transformée de Hadamard (Ritter, 1996).

Appliquer ces transformées en utilisant le circuit naïf impliquerait de l'ordre de n^2 opérations arithmétiques (on notera $\mathcal{O}(n^2)$ dans la suite de ce document). Dans le cas où n est grand, et où la transformée doit s'appliquer à un grand nombre de vecteurs d'entrées, ce coût peut s'avérer prohibitif. Cependant, il existe des algorithmes linéaires rapides permettant d'appliquer ces transformées sans faire autant de calculs. De tels algorithmes existent pour la transformée de Fourier sous le nom de transformée de Fourier rapide (Cooley et Tukey, 1965) (FFT pour *Fast Fourier*

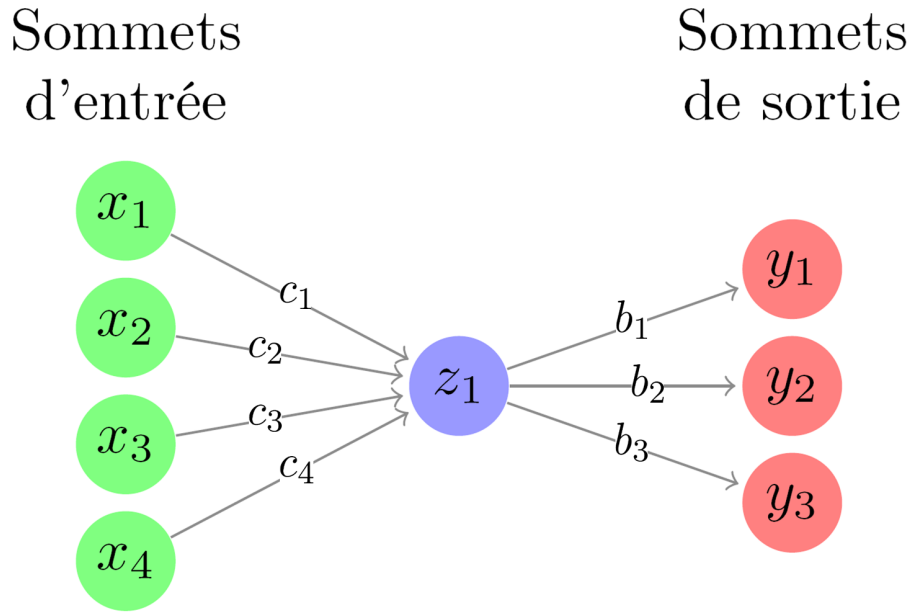


Figure 11 : Circuit associé à une matrice de rang 1, $\mathbf{A} = \mathbf{bc}^T \in \mathbb{R}^{3 \times 4}$.

Transform), pour la transformée en ondelettes (Mallat, 1989), pour la DCT (Chen et al., 1977), pour la transformée de Hadamard (Shanks, 1969) et bien d'autres. Ils ont en commun de décomposer l'application de la transformée en J étapes successives peu coûteuses, chaque étape étant une transformation linéaire simple correspondant à une multiplication par une matrice ayant peu d'entrées non-nulles (la matrice est alors dite creuse). Implicitement, une telle manière de procéder revient à considérer une factorisation de \mathbf{A} en J matrices creuses $\mathbf{S}_1, \dots, \mathbf{S}_J$: $\mathbf{A} = \mathbf{S}_J \mathbf{S}_{J-1} \dots \mathbf{S}_2 \mathbf{S}_1 = \prod_{j=1}^J \mathbf{S}_j$, et à effectuer le calcul de $\mathbf{y} = \mathbf{Ax}$ par étapes : d'abord $\mathbf{s} = \mathbf{S}_1 \mathbf{x}$, puis $\mathbf{t} = \mathbf{S}_2 \mathbf{s}$ et ainsi de suite jusqu'à \mathbf{y} . En terme de circuits linéaires, cela revient à considérer des sommets correspondant à \mathbf{x} , \mathbf{y} et tous les vecteur intermédiaires, et des arcs correspondant aux entrées non-nulles des facteurs \mathbf{S}_j , $j \in \{1, \dots, J\}$. Un exemple dans le cas de la transformée de Fourier rapide est donnée en Figure 12.

Les circuits associés à des transformées rapides ont une profondeur égale à k (le nombre d'étapes), et une complexité égale à $|\mathcal{A}| = \sum_{j=1}^J \|\mathbf{S}_j\|_0$, où $\|\mathbf{S}_j\|_0$ donne le nombre de coefficients non-nuls de \mathbf{S}_j . Par exemple dans le cas de la transformée de Fourier rapide quand n est une puissance de 2, on a $k = \log_2 n$ et $\|\mathbf{S}_j\|_0 = 2n$, $\forall j$, ce qui donne une complexité de $2n \log_2 n$.

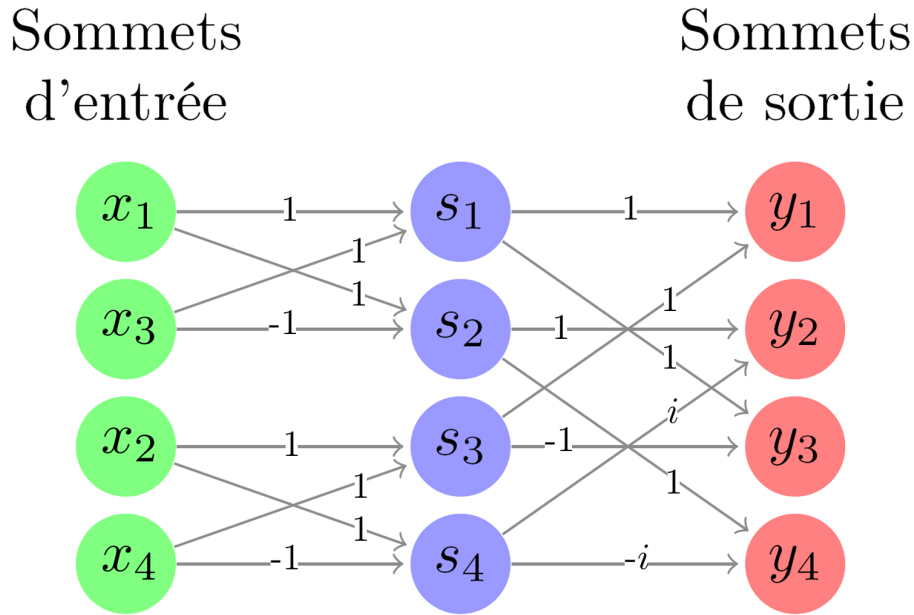


Figure 12 : Circuit associé à la transformée de Fourier rapide en dimension 4.

3.2. Efficience

Cette section nous permet d'introduire la notion centrale de cette thèse : l'efficience. Nous définissons tout d'abord l'efficience d'un circuit linéaire, puis, en nous appuyant sur elle, celle d'une matrice.

3.2.1. Efficience d'un circuit linéaire

Comme il est indiqué en introduction de ce document, l'efficience est une notion prenant en compte l'objectif visé mais aussi les ressources mises en œuvre pour l'atteindre, contrairement à l'efficacité qui ne prend en compte que l'objectif. Elle se définit comme le rapport entre le résultat obtenu et les ressources utilisées pour y parvenir. Nous transposons ici cette notion aux circuits linéaires.

Définition 4. L'efficience d'un circuit linéaire $C = (\mathcal{S}, \mathcal{A})$ ayant n entrées et m sorties, notée $\mathcal{E}(C)$, est égale au rapport entre la quantité mn et la complexité du circuit $|\mathcal{A}|$:

$$\mathcal{E}(C) \triangleq \frac{mn}{|\mathcal{A}|}.$$

L'efficience d'un circuit est donc une quantité positive. Elle correspond au rapport entre la complexité du circuit naïf pour une matrice de taille $m \times n$ et la complexité du circuit considéré. Pour un circuit quelconque C , une efficience inférieure à 1

implique qu'il existe un autre circuit (le circuit naïf) ayant exactement les mêmes relations entrées/sorties à moindre coût. Au contraire, une efficacité supérieure à 1 indique que le circuit considéré est "meilleur" en terme de quantité de calculs que le circuit naïf correspondant (une efficacité égale à k indique que le circuit C effectue k fois moins de calculs que le circuit naïf). La taille du circuit $|\mathcal{A}|$ est en général la mesure utilisée en informatique théorique pour quantifier sa complexité (Lokam, 2009). Nous introduisons dans cette thèse la notion d'efficacité du circuit afin d'avoir une mesure relative par rapport au circuit naïf, qui permet de quantifier directement le gain en complexité apporté par un circuit.

3.2.2. Efficacité d'une matrice

Définissons maintenant la notion d'efficacité relative à une matrice. Il est plus délicat de définir l'efficacité dans le cas d'une matrice que dans le cas d'un circuit. En effet, pour une matrice \mathbf{A} , il n'y a pas de notion immédiate de quantité de ressources utilisée pour effectuer le produit matrice/vecteur \mathbf{Ax} , car celle-ci dépend du circuit utilisé. C'est pourquoi nous nous appuyons sur la notion d'efficacité des circuits pour introduire l'efficacité d'une matrice.

Définition 5. L'efficacité d'une matrice \mathbf{A} , notée $\mathcal{E}(\mathbf{A})$, est égale à l'efficacité du circuit associé à \mathbf{A} le plus efficace :

$$\mathcal{E}(\mathbf{A}) \triangleq \max_{C \in \mathcal{C}_{\mathbf{A}}} \mathcal{E}(C).$$

L'efficacité d'une matrice est donc, comme celle d'un circuit, une quantité positive. Cependant $\mathcal{E}(\mathbf{A})$ ne peut pas être inférieure à 1, car il existe toujours un circuit naïf d'efficacité 1 dans $\mathcal{C}_{\mathbf{A}}$. Une matrice d'efficacité 1 est une matrice pour laquelle le circuit naïf est le meilleur en terme de quantité de calculs pour effectuer un produit matrice/vecteur. Une matrice d'efficacité k est une matrice pour laquelle le meilleur circuit linéaire associé est k fois plus efficace que le circuit naïf. Réciproquement, une matrice pour laquelle il existe un circuit associé d'efficacité k est d'efficacité supérieure à k .

3.3. Calculer et borner l'efficacité

Définir l'efficacité d'une matrice, comme nous l'avons fait dans la section précédente, amène à se poser plusieurs questions : *Peut-on calculer l'efficacité d'une matrice ?* ou plus modestement : *Peut-on borner l'efficacité d'une matrice ?* Nous nous intéressons dans cette section à ces deux questions. Nous nous appuyons pour cela sur les travaux effectués par la communauté de l'informatique théorique, et plus précisément de la théorie de la complexité. Nous considérons tout d'abord la question du calcul de l'efficacité, puis celle – plus explorée – des bornes sur l'efficacité.

3.3.1. Calcul

Calculer l'efficacité d'une matrice \mathbf{A} revient à déterminer la taille du circuit associé à \mathbf{A} le plus petit. Ce calcul est l'objet d'un article de Morgenstern ([Morgenstern, 1975](#)), dans lequel il introduit le *degré de liberté additif*, noté $D^+(\mathbf{A})$, quantité égale au nombre minimum d'additions requis pour calculer $\mathbf{A}\mathbf{x}$ pour un vecteur \mathbf{x} quelconque. Cette quantité dépend d'une décomposition additive de la matrice $\mathbf{A} = \mathbf{A}_1 + \mathbf{A}_2$, de la manière suivante : $D^+(\mathbf{A}) = \inf_{\mathbf{A}=\mathbf{A}_1+\mathbf{A}_2} f(\mathbf{A}_1, \mathbf{A}_2)$ (voir l'article original ([Morgenstern, 1975](#)) pour une expression détaillée). Le degré de liberté additif repose sur la minimisation d'une fonction définie sur toutes les matrices \mathbf{A}_1 et \mathbf{A}_2 telle que leur somme est égale à \mathbf{A} , il n'est donc a priori pas constructible, et il n'est pas possible de l'utiliser pour connaître le nombre minimum d'additions requis, et par conséquent l'efficacité de \mathbf{A} . Il n'y a par ailleurs, à ce jour et à la connaissance de l'auteur, pas d'autres travaux cherchant à calculer exactement le nombre minimum d'opérations arithmétiques requis pour un produit matrice/vecteur.

3.3.2. Bornes

Calculer exactement l'efficacité d'une matrice \mathbf{A} s'avérant une tâche difficile, il est possible d'envisager un problème plus accessible : borner l'efficacité, $B_{\text{inf}} \leq \mathcal{E}(\mathbf{A}) \leq B_{\text{sup}}$. Les bornes inférieures B_{inf} sont données par les algorithmes, on sait par exemple que l'efficacité de la matrice de Fourier en dimension n est supérieure à $\frac{n}{2 \log n}$. Les bornes supérieures B_{sup} , quant à elles, s'obtiennent en étudiant la taille minimale d'un circuit associé à \mathbf{A} . De telles bornes sont l'objet d'une littérature assez abondante émanant de la communauté de l'informatique théorique, et plus précisément de la théorie de la complexité.

Nous résumons ici les principaux résultats concernant les bornes de complexité, tout d'abord en considérant des circuits linéaires généraux, puis des circuits plus contraints.

Cas général

Borner la taille minimale d'un circuit associé à une matrice \mathbf{A} est un problème important en théorie de la complexité. Il a été démontré par Valiant ([Valiant, 1977](#)) que cette taille minimale est liée à la notion de *rigidité* de \mathbf{A} que nous définissons maintenant.

Définition 6. La rigidité d'une matrice \mathbf{A} pour un rang cible r , notée $\mathcal{R}_{\mathbf{A}}(r)$, est égale au nombre minimum d'entrées de \mathbf{A} qu'il faudrait modifier pour que \mathbf{A} soit de rang au plus r :

$$\mathcal{R}_{\mathbf{A}}(r) \triangleq \min\{\|\mathbf{C}\|_0, \text{rang}(\mathbf{A} + \mathbf{C}) \leq r\}.$$

La rigidité d'une matrice mesure donc la robustesse de son rang. Valiant prouve en particulier que pour une matrice carrée $\mathbf{A} \in \mathbb{R}^{n \times n}$, si $\mathcal{R}_{\mathbf{A}}(\epsilon n) \geq n^{1+\delta}$ avec ϵ et δ des constantes positives, alors un circuit associé à \mathbf{A} de profondeur $\log_2 n$ doit être de taille $\Omega(n \log \log n)$. Une borne inférieure sur la rigidité d'une matrice permet donc de déduire une borne inférieure sur la taille des circuits associés et donc sur l'efficacité de la matrice en question.

De nombreux auteurs ont alors étudié la rigidité de certaines matrices particulières : les matrices de Cauchy, de Vandermonde (dont fait partie la matrice de Fourier), et de Hadamard. Cependant, les meilleurs résultats obtenus à ce jour sont du type $\mathcal{R}_{\mathbf{A}}(r) = \Omega(\frac{n^2}{r} \log \frac{n}{r})$, et se réduisent à $\mathcal{R}_{\mathbf{A}}(\epsilon n) = \Omega(\frac{n}{\epsilon} \log \epsilon)$ dans le cas où $r = \epsilon n$. Ils ne permettent donc pas de formuler de bornes sur la taille minimale d'un circuit associé à \mathbf{A} . Il a été proposé d'utiliser d'autres propriétés analogues à la rigidité pour étudier ce type de problème, sans toutefois d'avancées significatives, voir (Lokam, 2009) pour une revue détaillée de ces approches.

Circuit contraint

Borner la complexité d'un circuit associé à une matrice dans le cas général n'étant pas chose aisée, il est possible de restreindre le type de circuit considéré par l'analyse, afin de simplifier le problème. Une restriction possible est de ne considérer que les circuits dont les coefficients sont bornés par une constante θ . Ce scénario est justifié par le fait que les circuits effectivement utilisés pour calculer des transformations linéaires ne permettent d'effectuer que des multiplications d'amplitude et de précision finies. Ce modèle restreint a été utilisé pour borner la taille d'un circuit associé à une matrice \mathbf{A} par une quantité dépendant du déterminant $\det(\mathbf{A})$. Les bornes obtenues sont $\Omega(n \log n)$ pour les matrices de Hadamard généralisées (dont fait partie la matrice de Fourier), dans le cas de circuits de profondeur logarithmique dans (Pudlák, 2000), et sans restriction sur la profondeur dans (Morgenstern, 1973). Cependant, ces résultats ne sont pas totalement satisfaisants, car il est facile de construire une matrice de déterminant arbitrairement élevé pouvant être associée à un circuit de taille linéaire (une matrice diagonale).

En conclusion de cette section, nous avons vu que calculer, et même borner l'efficacité d'une matrice n'est pas une chose facile. Nous proposons donc dans les chapitres suivants de fixer *a priori* une certaine efficacité, afin de l'imposer à des matrices avec lesquelles nous approchons ou estimons des matrices dignes d'intérêt.

Conclusion. Nous avons défini dans ce chapitre l'efficacité d'une matrice, concept central de cette thèse. Nous avons ensuite discuté de moyens de calculer ou de borner cette efficacité, et abouti à la conclusion que ces deux tâches sont difficiles en général.

Deuxième partie

Contributions

Approximation par matrices efficaces

Nous étudions dans ce chapitre l'approximation d'une matrice d'intérêt par une matrice efficace. Disposer d'une telle approximation permet d'effectuer de manière plus efficace diverses tâches mettant en jeu la matrice de départ, telles que la multiplication par cette matrice. Nous proposons une méthode d'approximation basée sur un algorithme d'optimisation non convexe, ainsi qu'une stratégie d'initialisation appropriée. Ce chapitre est dans sa majeure partie basé sur un article publié dans la revue IEEE Journal of Selected Topics in Signal Processing ([Le Magoarou et Gribonval, 2016b](#)).

Sommaire

4.1	Objectif général	53
4.1.1	Formulation en tant que factorisation matricielle multi-couche creuse	54
4.1.2	Avantages de la structure multi-couche creuse	55
4.2	Travaux liés	56
4.3	Algorithme d'approximation	61
4.3.1	Cas général	61
4.3.2	Cas particulier utilisé	62
4.4	Stratégie hiérarchique	65
4.4.1	Parallèle avec l'apprentissage profond	66
4.4.2	Algorithme hiérarchique proposé	66
4.4.3	Illustration : retrouver la transformée de Hadamard rapide	69

4.1. Objectif général

Notre objectif est d'obtenir une approximation efficace d'une matrice d'intérêt $\mathbf{A} \in \mathbb{R}^{m \times n}$. Un problème d'approximation matricielle s'exprime en terme d'optimisation de la façon suivante :

$$\underset{\hat{\mathbf{A}}}{\text{minimiser}} \quad d(\mathbf{A}, \hat{\mathbf{A}}) + g(\hat{\mathbf{A}}), \quad (\text{PA})$$

où d est la fonction mesurant une distance entre la matrice d'intérêt et son approximation, et g est la fonction de pénalité posant des contraintes sur la structure de l'approximation $\hat{\mathbf{A}}$. Le défi consiste à choisir la distance d à utiliser ainsi qu'à concevoir une fonction g imposant une approximation efficace, que l'on peut gérer du point de vue de l'optimisation.

Nous avons vu au chapitre précédent qu'il est difficile de trouver le circuit le plus efficace associé à une matrice donnée (et même de donner une borne inférieure sur sa taille). Nous proposons dans cette thèse de nous attaquer à un problème plus abordable, et d'approcher une certaine matrice d'intérêt \mathbf{A} par une matrice $\hat{\mathbf{A}}$ choisie dans un ensemble \mathcal{S} dont on sait par construction que tous les éléments sont d'efficacité supérieure à une certaine borne \mathcal{E}_{\min} (autrement dit on a l'implication $\mathbf{X} \in \mathcal{E} \Rightarrow \mathcal{E}(\mathbf{X}) \geq \mathcal{E}_{\min}$). Cela signifie de manière équivalente que l'on dispose d'une borne supérieure sur la taille du plus petit circuit associé à chaque élément de l'ensemble \mathcal{S} . Nous explicitons la forme choisie pour cet ensemble et formulons (PA) comme un problème de factorisation matricielle relativement simple dans la sous-section suivante.

4.1.1. Formulation en tant que factorisation matricielle multi-couche creuse

Comme évoqué au chapitre précédent, les deux principales familles de matrices pour lesquelles un circuit associé de petite taille est connu sont les matrices de rang faible et les transformées usuelles du traitement du signal. Dans les deux cas, l'algorithme rapide associé au circuit de petite taille peut être vu comme la multiplication successive par $J \geq 2$ matrices $\mathbf{S}_1, \dots, \mathbf{S}_J$ contenant peu d'entrées non-nulles ($\|\mathbf{S}_j\|_0 \leq p_j, j = 1, \dots, J$). Un exemple illustrant ceci, dans le cas de la transformée de Hadamard est donné en Figure 13. L'existence d'un circuit de petite taille et celle d'une factorisation en matrices creuses sont même équivalentes, ceci étant dû à la définition des circuits linéaires donnée au chapitre précédent.

Revenant à notre problème d'approximation (PA), ceci mène tout naturellement à choisir l'approximation $\hat{\mathbf{A}}$ dans un ensemble contenant uniquement des matrices factorisables en matrices creuses, et ainsi avoir $\hat{\mathbf{A}} = \mathbf{S}_J \dots \mathbf{S}_1$. On aboutit alors à un problème de factorisation en matrices creuses :

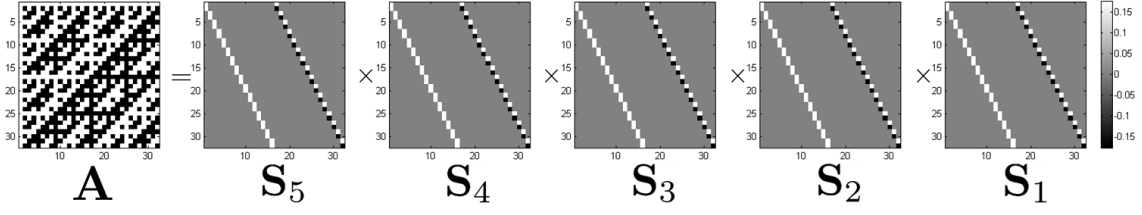


Figure 13 : La matrice de Hadamard de taille $n = 32$ et la factorisation en matrices creuses correspondant à l’algorithme rapide associé. La matrice de Hadamard est totalement dense, elle contient $n^2 = 1024$ entrées non-nulles (le circuit naïf associé est donc de taille 1024). Les facteurs creux contiennent chacun $2n = 64$ entrées non-nulles (le circuit correspondant à la multiplication successive par ces facteurs est donc de taille $2n \log_2 n = 320$).

$$\underset{\mathbf{S}_1, \dots, \mathbf{S}_J}{\text{minimiser}} \quad d(\mathbf{A}, \mathbf{S}_J \dots \mathbf{S}_1) + \sum_{j=1}^J g_j(\mathbf{S}_j), \quad (\text{PF})$$

où la pénalité g_j promeut la parcimonie du facteur \mathbf{S}_j . Par exemple, on peut tout simplement choisir

$$g_j(\mathbf{S}) = \begin{cases} 0 & \text{si } \mathbf{S} \in \mathcal{S}_j, \\ +\infty & \text{sinon,} \end{cases} \quad (4.1)$$

avec $\mathcal{S}_j = \{\mathbf{M} \in \mathbb{R}^{m_j \times m_{j-1}}, \|\mathbf{M}\|_0 \leq p_j\}$, auquel cas les pénalités sont les fonctions indicatrices d’ensembles de matrices creuses dont la parcimonie est paramétrée par les nombres p_j , $j = 1, \dots, J$, et le nombre de lignes par les nombres m_j , $j = 1, \dots, J-1$ (on a $m_J = m$ et $m_0 = n$ afin que $\hat{\mathbf{A}}$ soit de la bonne taille). On pourra alors noter $g_j(\mathbf{S}) = \delta_{\mathcal{S}_j}(\mathbf{S})$.

Une telle structure multi-couche creuse $\hat{\mathbf{A}} = \mathbf{S}_J \dots \mathbf{S}_1$ sera appelée *Flexible Multi-layer Sparse Transform*, et on utilisera l’acronyme FA μ ST pour désigner un tel objet, pour des raisons de concision. Elle garantit l’existence d’un circuit de taille $T_{\text{fac}} = \sum_{j=1}^J \|\mathbf{S}_j\|_0$ associé à l’approximation $\hat{\mathbf{A}}$ (correspondant à la multiplication successive par les facteurs creux), et donne ainsi une borne inférieure sur l’efficacité de cette approximation, que l’on appelle *gain relatif en complexité* (RCG). On a en effet :

$$\mathcal{E}(\hat{\mathbf{A}}) \geq \frac{mn}{T_{\text{fac}}} \triangleq \text{RCG}. \quad (4.2)$$

Le problème de factorisation (PF) dépend de nombreux paramètres. En effet, le nombre de facteurs J , leur taille m_j et leur parcimonie p_j sont à déterminer. Pour cela, nous nous inspirerons en pratique des algorithmes rapides associés aux transformées usuelles dans les parties expérimentales de cette thèse.

4.1.2. Avantages de la structure multi-couche creuse

Disposer d'une approximation $\hat{\mathbf{A}}$ sous la forme d'un produit de matrices creuses $\mathbf{S}_J \dots \mathbf{S}_1$ est avantageux tant que le gain relatif en complexité occasionné est suffisamment élevé. Cette quantité correspond en effet à un rapport de complexité entre la matrice d'intérêt \mathbf{A} et son approximation efficace $\hat{\mathbf{A}}$ pour plusieurs tâches. Ceci est vrai notamment pour le *coût de stockage*, la *vitesse de multiplication*, et la *pertinence statistique*.

Coût de stockage

En utilisant une structure de données en liste de coordonnées (COO (Jones et al., 2001)) pour stocker les facteurs creux $\mathbf{S}_1, \dots, \mathbf{S}_J$, une FA μ ST occupe un espace de taille $\mathcal{O}(T_{\text{fac}})$ (de l'ordre de T_{fac}). En effet chaque entrée non-nulle peut être localisée en utilisant seulement trois entiers (un pour le facteur, un pour la ligne et un pour la colonne), ce qui fait un total d'un nombre réel et trois entiers par entrée non-nulle. De plus, on a besoin de $J + 1$ entiers supplémentaires pour préciser la taille des facteurs. D'autre part, stocker la matrice d'intérêt \mathbf{A} nécessite mn nombres réels. En conclusion, le gain de stockage de l'approximation efficace est bien de l'ordre de la quantité RCG.

Vitesse de multiplication

Multiplier un vecteur $\mathbf{x} \in \mathbb{R}^n$ par une FA μ ST (en utilisant le circuit correspondant à la multiplication successive par les facteurs creux) nécessite au plus T_{fac} multiplications scalaires, alors qu'effectuer la même opération avec une matrice dense (en utilisant le circuit naïf) nécessite mn multiplications scalaires. Cette fois encore, le gain en terme de vitesse de multiplication est bien de l'ordre de la quantité RCG.

Pertinence statistique

Un autre avantage intéressant, quoique moins évident des FA μ ST intervient lorsque l'on n'a pas directement accès à l'opérateur \mathbf{A} , et que la tâche est de l'estimer, comme par exemple dans le cas de l'apprentissage de dictionnaire (on sort alors du cadre du problème (PF)). Dans ce cas, imposer à l'estimation une forme multi-couche creuse permet de réduire le nombre de paramètres à estimer (T_{fac} au lieu de mn dans le cas général). Plus spécifiquement, ceci permet de réduire le nombre d'échantillons nécessaires à une estimation fiable de l'opérateur. Le gain est une fois de plus de l'ordre de la quantité RCG, et ceci est validé expérimentalement et de manière théorique au chapitre 8.

4.2. Travaux liés

De multiples problèmes de factorisation matricielle similaires à (PF) ont été étudiés dans différents domaines. Certains sont des outils classiques de l’algèbre linéaire numérique tels que la décomposition en valeurs singulières (SVD) tronquée, alors que d’autres ont émergé plus récemment en traitement du signal ou en apprentissage automatique. Nous les présentons ici en mettant l’accent sur leurs liens avec le problème (PF) considéré dans cette thèse.

La SVD tronquée

Afin de réduire la complexité calculatoire associée à un opérateur linéaire d’intérêt \mathbf{A} , l’approche la plus classique consiste certainement à en calculer une approximation $\hat{\mathbf{A}}$ de rang k (avec $k \ll \min(m, n)$) à l’aide de la SVD tronquée. Ceci revient à ne considérer que les vecteurs singuliers associés aux k plus grandes valeurs singulières de \mathbf{A} . Une telle approche correspond à une approximation sous forme factorisée $\hat{\mathbf{A}} = \mathbf{S}_2 \mathbf{S}_1$ avec $\mathbf{S}_2 \in \mathbb{R}^{m \times k}$ et $\mathbf{S}_1 \in \mathbb{R}^{k \times n}$, où les matrices \mathbf{S}_1 et \mathbf{S}_2 sont pleines. On a donc $T_{\text{fac}} = k(m+n)$ et l’approximation est en effet plus efficace que l’opérateur d’origine si k est assez petit. Le calcul de la SVD tronquée est de complexité $\mathcal{O}(k(m+n)^2)$ en utilisant la méthode des puissances (Mises et Pollaczek-Geiringer, 1929).

Une comparaison entre la SVD tronquée (deux facteurs denses et de petite taille) et des FA μ ST dont la structure est inspirée des transformées rapides ($J > 2$ facteurs creux de taille comparable à celle de \mathbf{A}), mettant en évidence un compromis entre erreur relative d’approximation et gain relatif en complexité est présenté en Figure 14. On voit sur cette figure que les FA μ ST inspirées par les transformées rapides mènent à de bien meilleurs compromis. Des résultats similaires sont observés sur différentes matrices d’intérêt correspondant à diverses applications, ce qui motive le fait de considérer $J > 2$ facteurs creux dans le problème (PF).

Approximations locales de rang faible

Face aux limitations de l’approximation de rang faible par SVD tronquée illustrée par la figure précédente, la communauté de l’algèbre linéaire numérique a développé des méthodes d’approximation d’opérateurs par des juxtapositions de p matrices de rang faible (disons de rang k pour simplifier) de plus petite taille (appelées patches). Ce paradigme général englobe plusieurs techniques d’approximation conçues lors des trente dernières années, notamment la méthode rapide des multipôles (FMM pour “Fast Multipole Method” (Rokhlin, 1985)), les matrices hiérarchiques ou H-matrices (Hackbusch, 1999) ou d’autres (Candès et al., 2007). Ces méthodes sont utilisées pour la compression d’opérateurs utilisés pour la résolution numérique d’équations différentielles, étant la discrétisation d’un noyau continu sous-jacent. Toute la diffi-

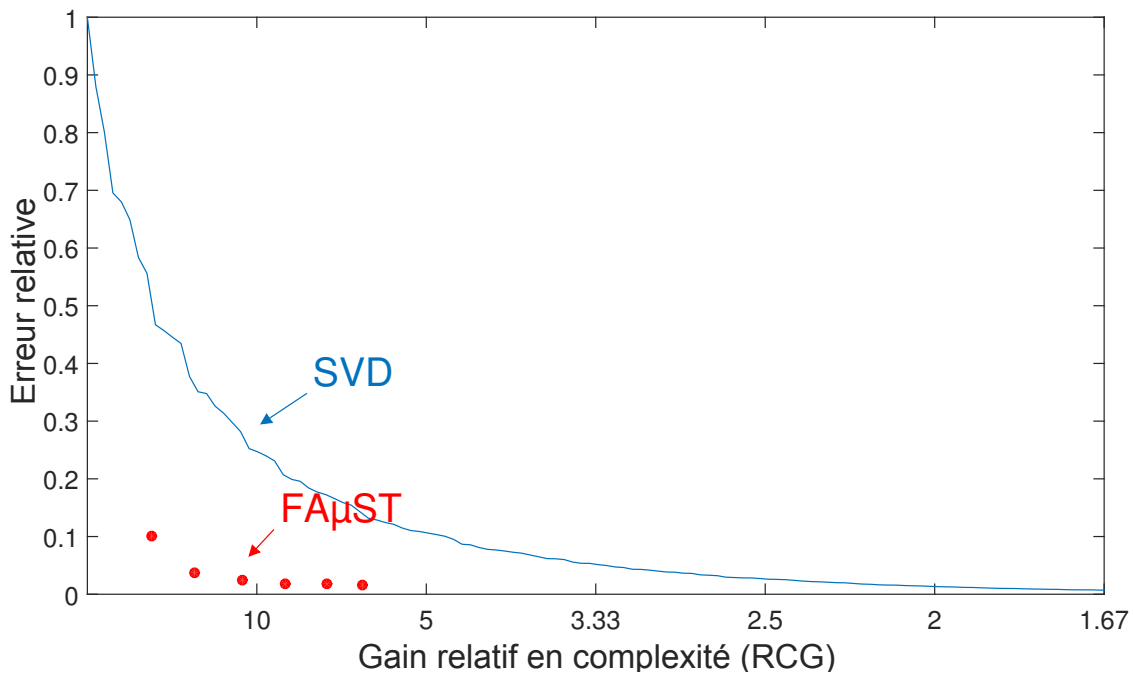


Figure 14 : Comparaison entre quatre FA μ ST correspondant à différentes configuration du problème (PF) et la SVD tronquée. L’approximation d’une matrice de taille 204×8193 correspondant à l’opérateur direct d’un problème inverse de magnéto-encéphalographie (MEG) est considérée pour cet exemple.

culté de ces techniques réside dans le choix du support des patches utilisés pour l’approximation, qui est gouverné par la régularité locale du noyau sous-jacent. Ceci peut être vu comme un problème d’approximation similaire à (PF), avec $\mathbf{A} \approx \hat{\mathbf{A}} = \mathbf{S}_2 \mathbf{S}_1$, où les facteurs $\mathbf{S}_1 \in \mathbb{R}^{m \times pk}$ et $\mathbf{S}_2 \in \mathbb{R}^{pk \times n}$ sont creux, leur supports étant fixés de manière analytique et correspondant à la localisation des patches. Les k premières colonnes de \mathbf{S}_2 ont le même support qui définit la localisation verticale ainsi que la hauteur du premier patch, alors que les k premières lignes de \mathbf{S}_1 ont le même support qui définit la localisation horizontale ainsi que la largeur du premier patch, et il en va de même pour les patches suivants. Nous nous intéressons plus particulièrement dans cette thèse à des approches où le support des facteurs est déterminé *par optimisation et non pas analytiquement*.

Compression basée sur des ondelettes

Ce paradigme de compression d’opérateur, introduit par (Beylkin et al., 1991), est basé sur l’utilisation de transformées en ondelettes orthogonales dans le domaine des colonnes (associée à une matrice $\Phi_1 \in \mathbb{R}^{m \times m}$), et dans le domaine des lignes (associée

à une matrice $\Phi_2 \in \mathbb{R}^{n \times n}$) de l'opérateur $\mathbf{A} \in \mathbb{R}^{m \times n}$ à compresser. Par orthogonalité des matrices Φ_1 et Φ_2 on a $\mathbf{A} = \Phi_1 \Phi_1^T \mathbf{A} \Phi_2 \Phi_2^T$. La technique de compression repose sur le fait que la matrice $\mathbf{B} \triangleq \Phi_1^T \mathbf{A} \Phi_2$ est compressible (ce qui n'est vrai que si les bases d'ondelettes sont bien choisies). Ceci signifie qu'on a $\mathbf{B} \approx \hat{\mathbf{B}}$ où $\hat{\mathbf{B}}$ est une matrice creuse, ce qui implique $\mathbf{A} \approx \hat{\mathbf{A}} = \Phi_1 \hat{\mathbf{B}} \Phi_2^T$. L'approximation $\hat{\mathbf{A}}$ est efficace si la matrice $\hat{\mathbf{B}}$ est assez creuse et si les transformées en ondelettes Φ_1 et Φ_2 sont associées à des algorithmes rapides. Cette méthode peut être vue comme une approximation de l'opérateur d'origine par une FA μ ST, les facteurs creux à gauche de $\hat{\mathbf{B}}$ correspondant à l'algorithme rapide associé à Φ_1 , et les facteurs creux à droite de $\hat{\mathbf{B}}$ correspondant à l'algorithme rapide associé à Φ_2 .

Apprentissage de dictionnaire

L'apprentissage de dictionnaire classique, tel que présenté au chapitre concernant les factorisations matricielles, souffre du fait que les dictionnaires appris sont en général des matrices denses, qui peuvent se révéler coûteuses à manipuler pour des signaux de grande dimension.

Afin d'apprendre des dictionnaires plus efficaces, deux tendances ont récemment émergé, donnant naissance à des approches liées au problème de factorisation multicouche creuse (PF). La méthode sparse KSVD ([Rubinstein et al., 2010b](#)) (abrégée KSVDS) consiste à apprendre un dictionnaire dont les colonnes sont des combinaisons linéaires de peu d'atomes d'un autre dictionnaire fixe appelé *dictionnaire de base* \mathbf{D}_{base} . Le dictionnaire de base est associé avec un algorithme rapide, ce qui signifie qu'il est factorisable en matrices creuses (c'est une FA μ ST). Cette méthode peut être vue comme une instance du problème (PF) où les $J - 1$ facteurs les plus à gauche correspondent à l'algorithme rapide associé à \mathbf{D}_{base} , et le facteur le plus à droite correspond à la représentation des atomes dans le dictionnaire de base, le dictionnaire appris s'exprime alors $\mathbf{D} = \mathbf{D}_{\text{base}} \mathbf{S}_1$, où \mathbf{S}_1 est creux. Cette approche est limitée par le fait d'utiliser un dictionnaire de base qui réduit l'adaptabilité du dictionnaire appris (en fixant a priori les facteurs creux correspondant au dictionnaire de base).

Un autre approche ([Chabiron et al., 2014](#)) consiste à apprendre un dictionnaire dans lequel chaque atome est la composition de plusieurs convolutions avec des noyaux creux à supports connus, tels que le dictionnaire dans son ensemble soit manipulable rapidement. Ceci peut être vu comme une instance du problème (PF) où les pénalités g_j contraignent les facteurs à être des matrices circulantes, à support fixé. Cette formulation est celle qui se rapproche le plus de ce qui est proposé dans cette thèse, mais est restreinte au cas où les matrices circulantes creuses permettent une bonne approximation, et nécessite de préciser les supports des facteurs creux.

Statistique – analyse de facteur

Un problème proche de (PF) consiste à diagonaliser approximativement une matrice de covariance \mathbf{A} à l'aide d'une matrice orthogonale associée à un algorithme rapide $\mathbf{U} = \mathbf{S}_J \dots \mathbf{S}_1$, afin de disposer d'une transformée rapide permettant de décorréler les variables d'entrée. Ceci est fait en pratique (Lee et al., 2008; Cao et al., 2011) de manière gloutonne, en utilisant pour les facteurs $\mathbf{S}_1, \dots, \mathbf{S}_J$ des rotations élémentaires telles que les rotations de Givens (Givens, 1958). Nous nous inspirons de ces techniques dans cette thèse dans le chapitre concernant la transformée de Fourier rapide sur graphe.

Apprentissage automatique

Des modèles similaires ont été explorés de différents points de vue en apprentissage automatique. Par exemple, la factorisation en matrice non-négatives creuses multifacteur (Lyu et Wang, 2013) peut être vue comme une instance du problème (PF) avec une mesure de distance spécifique au problème (la divergence de Kullback-Leibler) et des facteurs \mathbf{S}_j contraints à être non négatifs. La structure spécifique du problème (notamment la non-négativité) permet d'utiliser lors de l'optimisation un algorithme multiplicatif, alors que les approches présentées dans cette thèse sont plus générales, utilisant des itérations proximales.

Apprentissage profond

Dans le contexte des réseaux de neurones profonds, des garanties d'identifiabilité sur la structure du réseau ont été établies dans le cas d'un modèle génératif où les couches consécutives sont liées par des matrices de poids creuses et aléatoires, et les non-linéarités négligées (Neyshabur et Panigrahy, 2013; Arora et al., 2013). Le réseau peut alors être vu simplement comme une succession de multiplications par des matrices creuses, chaque facteur représentant une couche. Ces approches pourraient donc être utilisées pour obtenir des garanties théoriques concernant l'identifiabilité d'une solution au problème (PF). Cependant, les résultats de ces travaux ne s'appliquent qu'au cas où les facteurs creux sont aléatoires (chaque colonne ayant d entrées non-nulles égales à ± 1), et nécessitent une parcimonie extrêmement forte (d de l'ordre de $n^{\frac{1}{6}}$) pour parvenir à des garanties, ce qui ne correspond pas aux cas pratiques rencontrés dans cette thèse.

Traitement du signal sur graphes

Le traitement du signal sur graphe est un domaine de recherche récent, dont le but est de généraliser les outils du traitement du signal classique aux signaux définis sur les sommets d'un graphe (Shuman et al., 2013). Nous en donnons une introduction

plus conséquente dans le chapitre consacré à la transformée de Fourier rapide sur graphe. Des problèmes de factorisation matricielle similaires à (PF) ont émergé dans ce domaine, principalement dans le but de définir des transformées en ondelette pour les signaux sur graphe.

Une première approche (Rustamov et Guibas, 2013) consiste à définir des ondelettes basées sur des données d'entraînements constituées de signaux sur le graphe d'intérêt (une sorte d'apprentissage de dictionnaire). Les contraintes imposées aux ondelettes s'apparentent à une structure multi-couche creuse, dans laquelle chaque facteur est de plus contraint à être un bloc élémentaire d'ondelette "liftée". L'optimisation utilisée emprunte des techniques utilisées en apprentissage profond (il s'agit d'entraînement couche par couche d'auto-encodeurs (Bengio et al., 2007)).

Une seconde approche (Kondor et al., 2014) consiste à diagonaliser partiellement le Laplacien du graphe en question à l'aide de rotations élémentaires. Plus précisément, la base dans laquelle le Laplacien est exprimée est modifiée de manière gloutonne, à l'aide d'une matrice orthogonale creuse (rotation élémentaire), ce qui définit une transformée multi-couche creuse. La procédure engendre un Laplacien diagonal pour certaines dimensions (correspondant aux ondelettes), et dense pour les dimensions restantes (correspondant à la fonction d'échelle). Nous proposons une approche similaire pour définir une transformée de Fourier rapide sur graphe au chapitre 6.

4.3. Algorithme d'approximation

Nous présentons dans cette section l'algorithme d'optimisation utilisé pour résoudre (PF) de manière approchée. Celui-ci est basé sur l'algorithme PALM (Bolte et al., 2014) présenté dans le chapitre dédié à l'optimisation.

4.3.1. Cas général

Afin de proposer un algorithme adapté au problème (PF), commençons par l'analyser du point de vue de l'optimisation. La variable de décision du problème est ici naturellement divisée en J blocs, correspondant aux facteurs $\mathbf{S}_1, \dots, \mathbf{S}_J$. La fonction de coût peut par ailleurs être séparée en un terme d'attache aux données et un terme de pénalité imposant la parcimonie des facteurs de la manière suivante :

$$f(\mathbf{S}_1, \dots, \mathbf{S}_J) \triangleq \underbrace{d(\mathbf{A}, \mathbf{S}_J \dots \mathbf{S}_1)}_{\text{Attache aux données}} + \underbrace{\sum_{j=1}^J g_j(\mathbf{S}_j)}_{\text{Pénalités}}. \quad (4.3)$$

Le terme correspondant aux pénalités est séparable (c'est une somme de J fonctions ne dépendant chacune que d'un seul facteur). Ainsi, si tant est que l'opérateur

proximal de ces fonctions soit calculable, et que le terme d'attache aux données soit différentiable, la fonction de coût f entre dans le cadre de l'algorithme PALM (Bolte et al., 2014). Une simple adaptation de cet algorithme mène alors à l'algorithme `palM4MSA`, acronyme signifiant PALM pour l'approximation creuse multi-couche (l'acronyme venant de "PALM for multi-layer sparse approximation" en anglais). Celui-ci est présenté en Algorithme 11, où nous considérons des pénalités étant des fonctions indicatrices d'ensembles de contraintes (prenant la forme indiquée en équation(4.1)). Les facteurs sont mis à jour de manière alternée, par une itération de descente de gradient projetée ($P_{S_j}(\cdot)$ est l'opérateur de projection sur l'ensemble de contrainte S_j). Une telle mise à jour n'est autre qu'une descente de gradient proximale, lorsque les pénalités correspondent à des fonctions indicatrices d'ensembles.

Algorithme 11 PALM pour l'approximation creuse multi-couche (version générale).

Entrée : Une matrice à approcher \mathbf{A} ; un nombre de facteurs J ; les ensembles de contraintes correspondant $\{S_j\}_{j=1}^J$; les facteurs initiaux $\{\mathbf{S}_j^0\}_{j=1}^J$, λ^0 ; un critère d'arrêt c .

```

1: while  $c = \text{False}$  do
2:   for  $j = 1$  to  $J$  do
3:     Fixer un pas  $c_j$ 
4:      $\mathbf{S}_j \leftarrow P_{S_j}(\mathbf{S}_j - c_j \cdot \nabla_{\mathbf{S}_j} d(\mathbf{A}, \mathbf{S}_J \dots \mathbf{S}_1))$ 
5:   end for
6: end while

```

Sortie : La factorisation obtenue : $\{\mathbf{S}_j\}_{j=1}^J = \text{palM4MSA}(\mathbf{A}, J, \{S_j\}_{j=1}^J, \dots)$

Cet algorithme jouit des mêmes conditions de convergence vers un point critique de PALM, dès que les conditions suffisantes suivantes sont remplies :

- (i) Les pénalités g_j , $j = 1, \dots, J$ sont des fonctions propres et semi-continues inférieurement (Rockafellar, 2015).
- (ii) Le terme d'attache aux données d est lisse.
- (iii) La fonction de coût f est semi-algébrique (Bolte et al., 2014, Definition 5.1).
- (iv) Le gradient du terme d'attache aux données par rapport au j -ème facteur $\nabla_{\mathbf{S}_j} d$ est Lipschitz-continu pour tout j , et ce avec un module de Lipschitz $L_j(\mathbf{S}_1, \dots, \mathbf{S}_{j-1}, \mathbf{S}_{j+1}, \dots, \mathbf{S}_J)$ (que l'on notera L_j par abus de notation).
- (v) Le pas c_j^i choisi à la i -ème itération pour le j -ème facteur est tel que $\forall i, j$, $1/c_j^i > L_j$ (l'inégalité n'a pas besoin d'être stricte si la pénalité g_j est convexe).

Nous présentons dans la sous-section suivante le cas particulier qui nous intéressera dans cette thèse, et montrons qu'il remplit effectivement toutes ces conditions.

4.3.2. Cas particulier utilisé

Nous exposons ici plus concrètement l'algorithme principal de cette thèse, en précisant les choix de conception ainsi que les détails de la formulation retenue.

Fonction de coût et ambiguïté d'échelle

Nous considérons dans cette thèse un terme d'attache aux données utilisant la norme de Frobenius $\|\cdot\|_F$ (avec $\|\mathbf{X}\|_F^2 = \sum_{i,j} x_{ij}^2$), de la forme $d(\mathbf{A}, \mathbf{S}_J \dots \mathbf{S}_1) = \frac{1}{2} \|\mathbf{A} - \mathbf{S}_J \dots \mathbf{S}_1\|_F^2$. Ce type de formulation donne lieu à une ambiguïté d'échelle entre les facteurs. En d'autres termes, on a $d(\mathbf{A}, \mathbf{S}_1, \dots, \mathbf{S}_J) = d(\mathbf{A}, \lambda_1 \mathbf{S}_1, \dots, \lambda_J \mathbf{S}_J)$ dès que $\prod_j \lambda_j = 1$. Afin d'éviter d'aboutir à des solutions avec des facteurs dont la norme est très petite ou très grande, on peut choisir les contraintes de manière à normaliser les facteurs. Par exemple, on peut choisir des pénalités qui sont des fonctions indicatrices d'ensembles de matrices creuses et normalisées, du type $\mathcal{S}_j = \{\mathbf{M} \in \mathbb{R}^{m_j \times m_{j-1}}, \|\mathbf{M}\|_0 \leq p_j, \|\mathbf{M}\|_F = 1\}$. Cependant, un produit de matrices normalisées est forcément de norme inférieure à 1, et ne permet pas d'approcher n'importe quelle matrice \mathbf{A} . Afin de remédier à cela, on introduit un scalaire multiplicatif positif λ dans le terme d'attache aux données qui devient $d(\mathbf{A}, \lambda \mathbf{S}_J \dots \mathbf{S}_1) = \frac{1}{2} \|\mathbf{A} - \lambda \mathbf{S}_J \dots \mathbf{S}_1\|_F^2$. Ce scalaire peut alors être intégré très simplement à l'algorithme d'optimisation, en le considérant comme un bloc de variables supplémentaire dans la formulation de PALM, auquel seule la contrainte de positivité est appliquée. On aboutit alors à la fonction de coût suivante :

$$f(\lambda, \mathbf{S}_1, \dots, \mathbf{S}_J) = \frac{1}{2} \|\mathbf{A} - \lambda \mathbf{S}_J \dots \mathbf{S}_1\|_F^2 + \sum_{j=1}^J \delta_{\mathcal{S}_j}(\mathbf{S}_j) + \delta_{\mathbb{R}_+}(\lambda). \quad (4.4)$$

Il est alors assez direct de vérifier à partir des définitions des notions correspondantes que les conditions de convergence (i), (ii) et (iii) de PALM sont remplies.

Gradients et modules de Lipschitz

L'algorithme PALM nécessite le gradient du terme d'attache aux données par rapport à tous les blocs de variables considérés (les facteurs $\mathbf{S}_1, \dots, \mathbf{S}_J$ et le scalaire λ). De plus, le gradient doit être Lipschitz par rapport à ces blocs afin de garantir la convergence.

Nous donnons ici une expression générale du gradient par rapport au j -ème facteurs \mathbf{S}_j . Pour cela, nous introduisons les notations $\mathbf{L}_j \triangleq \mathbf{S}_J \dots \mathbf{S}_{j+1}$ pour désigner

le produit de tous les facteurs qui se trouvent à gauche du j -ème facteur et $\mathbf{R}_j \triangleq \mathbf{S}_{j-1} \dots \mathbf{S}_1$ pour désigner le produit de tous les facteurs qui se trouvent à droite du j -ème facteur (on a par convention $\mathbf{L}_J = \mathbf{Id}$ et $\mathbf{R}_1 = \mathbf{Id}$). En utilisant les formules classiques de dérivation d'expressions matricielles (Petersen et Pedersen, 2012), on obtient :

$$\nabla_{\mathbf{S}_j} d(\mathbf{A}, \lambda \mathbf{S}_J \dots \mathbf{S}_1) = \lambda \mathbf{L}_j^T (\lambda \mathbf{L}_j \mathbf{S}_j \mathbf{R}_j - \mathbf{A}) \mathbf{R}_j^T.$$

De plus, ce gradient est Lipschitz-continu, avec pour expression pour le module de Lipschitz $L_j \leq \lambda^2 \|\mathbf{R}_j\|_2^2 \cdot \|\mathbf{L}_j\|_2^2$. Ceci vient de l'inégalité

$$\|\lambda \mathbf{L}_j^T (\lambda \mathbf{L}_j \mathbf{S}_j^1 \mathbf{R}_j - \mathbf{A}) \mathbf{R}_j^T - \lambda \mathbf{L}_j^T (\lambda \mathbf{L}_j \mathbf{S}_j^2 \mathbf{R}_j - \mathbf{A}) \mathbf{R}_j^T\|_F \leq \lambda^2 \|\mathbf{R}_j\|_2^2 \cdot \|\mathbf{L}_j\|_2^2 \cdot \|\mathbf{S}_j^1 - \mathbf{S}_j^2\|_F,$$

dont la preuve est donnée en annexe A. La condition de convergence (iv) est donc remplie pour les blocs de variables correspondant aux facteurs. La condition (v) est aussi remplie tant que le pas choisi est inférieur à l'inverse de ce module de Lipschitz (par exemple, si l'on a $c_j < 1/L_j$). La mise à jour du j -ème facteur prend alors la forme $\mathbf{S}_j \leftarrow P_{\mathcal{S}_j}(\mathbf{S}_j - c_j \lambda \mathbf{L}_j^T (\lambda \mathbf{L}_j \mathbf{S}_j \mathbf{R}_j - \mathbf{A}) \mathbf{R}_j^T)$.

Donnons maintenant l'expression du gradient par rapport au scalaire λ . Pour cela, nous introduisons la notation $\tilde{\mathbf{A}} \triangleq \mathbf{S}_J \dots \mathbf{S}_1$. Le gradient est alors très simple à calculer, et l'on obtient

$$\nabla_{\lambda} d(\mathbf{A}, \lambda \mathbf{S}_J \dots \mathbf{S}_1) = \lambda \text{Tr}(\tilde{\mathbf{A}}^T \tilde{\mathbf{A}}) - \text{Tr}(\tilde{\mathbf{A}}^T \mathbf{A})$$

Ce gradient est Lipschitz-continu, et $\text{Tr}(\tilde{\mathbf{A}}^T \tilde{\mathbf{A}})$ constitue un module de Lipschitz (ceci vient de l'égalité $|\lambda_1 \text{Tr}(\tilde{\mathbf{A}}^T \tilde{\mathbf{A}}) - \lambda_2 \text{Tr}(\tilde{\mathbf{A}}^T \tilde{\mathbf{A}})| = \text{Tr}(\tilde{\mathbf{A}}^T \tilde{\mathbf{A}}) \cdot |\lambda_1 - \lambda_2|$). La condition de convergence (iv) est donc aussi remplie pour le bloc de variables correspondant au scalaire λ . La contrainte de positivité imposée à λ est convexe, donc la condition de convergence (v) peut être appliquée avec une inégalité non-strictes, et on peut prendre $c_{\lambda} \leq 1/\text{Tr}(\tilde{\mathbf{A}}^T \tilde{\mathbf{A}})$. La mise à jour du scalaire λ prend alors la forme $\lambda \leftarrow P_{\mathbb{R}_+}(\lambda - c_{\lambda} \lambda \text{Tr}(\tilde{\mathbf{A}}^T \tilde{\mathbf{A}}) + c_{\lambda} \text{Tr}(\tilde{\mathbf{A}}^T \mathbf{A}))$. Prendre le pas maximal $c_{\lambda} = 1/\text{Tr}(\tilde{\mathbf{A}}^T \tilde{\mathbf{A}})$ mène alors à la mise à jour $\lambda \leftarrow P_{\mathbb{R}_+}(\text{Tr}(\tilde{\mathbf{A}}^T \mathbf{A})/\text{Tr}(\tilde{\mathbf{A}}^T \tilde{\mathbf{A}}))$.

Opérateurs de projection

Examinons maintenant plus en détails les opérateurs de projection mis en jeu. L'opérateur de projection $P_{\mathcal{E}}(\cdot)$ sur un ensemble \mathcal{E} doit respecter la relation $P_{\mathcal{E}}(\mathbf{X}) \in \underset{\mathbf{M} \in \mathcal{E}}{\text{argmin}} \|\mathbf{X} - \mathbf{M}\|_F$ pour tout \mathbf{X} .

Concernant les facteurs $\mathbf{S}_1, \dots, \mathbf{S}_J$, il s'agit d'effectuer la projection sur des ensembles de la forme $\mathcal{S}_j = \{\mathbf{M} \in \mathbb{R}^{m_j \times m_{j-1}}, \|\mathbf{M}\|_0 \leq p_j, \|\mathbf{M}\|_F = 1\}$. Dans ce cas l'opérateur de projection $P_{\mathcal{S}_j}(\cdot)$ ne garde que les p_j plus grandes entrées en valeur absolue de son argument (et met les autres à zéro, on appelle cela un seuillage dur) et normalise le résultat (voir preuve en annexe A). Le cadre de PALM permet

de considérer d'autres contraintes, plus ou moins pertinentes en fonction de l'application envisagée (parcimonie structurée, support partiellement fixé, matrice de Toeplitz, etc.). Les expressions des opérateurs de projection associés et les preuves correspondantes sont données en annexe A.

Concernant le scalaire λ , il s'agit d'effectuer une projection sur l'ensemble des réels positifs \mathbb{R}_+ . L'opérateur $P_{\mathbb{R}_+}(\cdot)$ met à zéro son argument si celui-ci est négatif ou nul, et rend son argument si celui-ci est positif.

Une version de l'algorithme PALM pour l'approximation creuse multi-couche prenant en compte les expressions données dans cette sous-section est donné en algorithme 12.

Algorithme 12 PALM pour l'approximation creuse multi-couche (version particulière).

Entrée : Une matrice à approcher \mathbf{A} ; un nombre de facteurs J ; les ensembles de contraintes correspondant $\{\mathcal{S}_j\}_{j=1}^J$; les facteurs initiaux $\{\mathbf{S}_j^0\}_{j=1}^J, \lambda^0$; un nombre d'itérations N .

- 1: **for** $i = 0$ to $N - 1$ **do**
- 2: **for** $j = 1$ to J **do**
- 3: $\mathbf{L}_j \leftarrow \prod_{\ell=j+1}^J \mathbf{S}_\ell^i$
- 4: $\mathbf{R}_j \leftarrow \prod_{\ell=1}^{j-1} \mathbf{S}_\ell^{i+1}$
- 5: Choisir un pas c_j^i tel que $1/c_j^i > (\lambda^i)^2 \|\mathbf{R}_j\|_2^2 \cdot \|\mathbf{L}_j\|_2^2$
- 6: $\mathbf{S}_j^{i+1} \leftarrow P_{\mathcal{S}_j} \left(\mathbf{S}_j^i - c_j^i \lambda^i \mathbf{L}_j^T (\lambda^i \mathbf{L}_j \mathbf{S}_j^i \mathbf{R}_j - \mathbf{A}) \mathbf{R}_j^T \right)$
- 7: **end for**
- 8: $\tilde{\mathbf{A}} \leftarrow \prod_{j=1}^J \mathbf{S}_j^{i+1}$
- 9: $\lambda^{i+1} \leftarrow P_{\mathbb{R}_+} \left(\frac{\text{Tr}(\mathbf{A}^T \tilde{\mathbf{A}})}{\text{Tr}(\tilde{\mathbf{A}}^T \tilde{\mathbf{A}})} \right)$
- 10: **end for**

Sortie : La factorisation obtenue : $\lambda^N, \{\mathbf{S}_j^N\}_{j=1}^J = \text{palm4MSA}(\mathbf{A}, J, \{\mathcal{E}_j\}_{j=1}^J, \dots)$

Coût de l'algorithme

Examinons le coût de l'algorithme 12. L'essentiel de la complexité de cet algorithme vient de la mise à jour des facteurs (ligne 6). Celle-ci a un coût en $\mathcal{O}(T_{\text{fac}}n + n^2)$, où $T_{\text{fac}} = \sum_{j=1}^J \|\mathbf{S}_j\|_0$ (en supposant que m et n sont du même ordre de grandeur). On peut par ailleurs supposer que la quantité T_{fac} est au moins de l'ordre de n , le coût de la mise à jour peut alors être exprimé en $\mathcal{O}(T_{\text{fac}}n)$. Le coût d'une itération complète, qui consiste à mettre à jour les J facteurs est alors en $\mathcal{O}(JT_{\text{fac}}n)$.

4.4. Stratégie hiérarchique

L'algorithme `pa1m4MSA` proposé dans la section précédente a pour but d'approcher une matrice d'intérêt \mathbf{A} par un produit de matrices creuses. Il jouit de garanties de convergence vers un point critique de son objectif (que l'on pourra désigner comme un minimum local par abus de langage par la suite). Cependant, il ne jouit d'aucune garantie théorique sur la qualité du point critique atteint, qui dépend uniquement de l'initialisation choisie. Un bon point critique correspond à une bonne qualité d'approximation, c'est-à-dire un terme d'attache aux données faible. Une bonne qualité d'approximation est indispensable pour n'importe quelle application. Malheureusement, les stratégies d'initialisation simples telles qu'une initialisation à zéro ($\mathbf{S}_J = \dots = \mathbf{S}_2 = \mathbf{Id}$ et $\mathbf{S}_1 = \mathbf{0}$) ou une initialisation aléatoire mènent empiriquement à de mauvais minima locaux. Ceci est d'autant plus vrai que le nombre de facteurs J est élevé. Ce constat est illustré par une expérience très simple, où l'on considère la matrice $\mathbf{A} \in \mathbb{R}^{n \times n}$ correspondant à la transformée de Hadamard en dimension n . On sait qu'il existe une factorisation en matrices creuses $\mathbf{A} = \mathbf{S}_J \dots \mathbf{S}_1$ comme illustré à la Figure 13, avec $J = \log_2 n$ et où chaque facteur est carré et contient $2n$ entrées non-nulles. Appeler l'algorithme `pa1m4MSA` sur cette matrice \mathbf{A} , avec des contraintes correspondant à la factorisation connue mène à de mauvais résultats en pratique (un terme d'attache aux données élevé).

Nous proposons dans cette section une stratégie d'initialisation permettant de remédier empiriquement à ce problème, et de converger en pratique vers de meilleurs minima locaux.

4.4.1. Parallèle avec l'apprentissage profond

Des problèmes de convergence similaires ont été rencontrés dans le passé par la communauté de l'apprentissage automatique. En effet, pendant des années, l'optimisation des réseaux de neurones comprenant un nombre élevé de couche (appelés réseaux de neurones profonds, voir le chapitre concernant les factorisation matricielles pour une introduction) n'aboutissait pas à des résultats satisfaisants. En conséquence, les réseaux de neurones profonds étaient bien souvent délaissés au profit de réseaux ne comprenant qu'un très petit nombre de couches. Cependant, il y a une dizaine d'années, un papier fondateur ([Hinton et Salakhutdinov, 2006](#)) a proposé une méthode d'initialisation ingénieuse permettant d'améliorer significativement les performances des réseaux de neurones profonds. Celle-ci consiste à considérer les couches du réseaux une par une, en les optimisant par descente de gradient stochastique de manière à ce qu'elles soient capables de bien reconstruire leurs données d'entrées (chacune prenant comme données d'entrée la sortie de la couche précédente, la première couche prenant les données d'entraînement), cette phase est

appelée le pré-entraînement. Après le pré-entraînement qui constitue une initialisation, le réseaux est entraîné globalement (toutes les couches en même temps) de manière classique par descente de gradient stochastique. Cette stratégie d'initialisation heuristique a montré de bons résultats dans de nombreuses applications (Bengio et al., 2007). Elle est d'ailleurs un des facteurs à l'origine du retour en grâce des réseaux de neurones profonds au cours de ces dernières années.

4.4.2. Algorithme hiérarchique proposé

Inspirés par ce type de stratégies d'initialisation, et en remarquant que `palm4MSA` permet d'atteindre de meilleurs minima locaux lorsqu'il est appelé avec un nombre de facteurs J faible, nous proposons ici une stratégie ne mettant en jeu que des factorisations avec $J = 2$ facteurs. En effet, imaginons que l'on dispose d'une matrice $\mathbf{A} = \mathbf{S}_N \dots \mathbf{S}_1$ étant le produit de N facteurs creux, et que le but est de retrouver ces facteurs creux $\mathbf{S}_1, \dots, \mathbf{S}_N$. On peut alors remarquer que la matrice d'intérêt est aussi le produit $\mathbf{A} = \mathbf{R}_1 \mathbf{S}_1$ de seulement deux facteurs avec $\mathbf{R}_1 = \mathbf{S}_N \dots \mathbf{S}_2$ (on appelle alors \mathbf{R}_1 le premier résidu). On peut alors chercher à retrouver approximativement cette factorisation (plus simple) en appelant `palm4MSA` avec $J = 2$ facteurs sur un terme d'attache aux données du type $\frac{1}{2} \|\mathbf{A} - \mathbf{R}_1 \mathbf{S}_1\|_F^2$. Ceci donne une approximation du premier facteur creux \mathbf{S}_1 . On peut ensuite effectuer le même type de factorisation sur le résidu \mathbf{R}_1 avec un terme d'attache aux données du type $\frac{1}{2} \|\mathbf{R}_1 - \mathbf{R}_2 \mathbf{S}_2\|_F^2$ pour retrouver le second facteur creux \mathbf{S}_2 ainsi qu'un second résidu $\mathbf{R}_2 \approx \mathbf{S}_N \dots \mathbf{S}_3$, et ainsi de suite jusqu'à ce que tous les facteurs creux soient exhibés. Ainsi, la factorisation en J facteurs creux est remplacée par $J - 1$ factorisation en seulement deux facteurs. De plus, on peut insérer entre chaque factorisation en deux facteurs une étape d'optimisation globale concernant tous les facteurs creux introduits jusque là ainsi que le résidu courant afin de rester proche de factoriser effectivement la matrice d'intérêt \mathbf{A} . Ceci peut être fait en appelant après la ℓ -ème factorisation en deux `palm4MSA` avec un terme d'attache aux données du type $\frac{1}{2} \|\mathbf{A} - \mathbf{R}_\ell \mathbf{S}_\ell \dots \mathbf{S}_1\|_F^2$, en initialisant les facteurs à leur valeur courante.

Choix des contraintes de parcimonie

Comment fixer la parcimonie des facteurs dans toutes les factorisations utilisant `palm4MSA` mises en jeu dans la procédure introduite au paragraphe précédent ? Les transformées rapides utilisées en traitement du signal correspondent toutes au même type de factorisation en matrices creuses. Par exemple, la transformée de Hadamard $\mathbf{A} \in \mathbb{R}^{n \times n}$ est factorisable en $\mathbf{A} = \mathbf{S}_J \dots \mathbf{S}_1$ avec $J = \log_2 n$ facteurs, chacun ayant $2n$ entrées non-nulles (2 par ligne et par colonne). De plus, la matrice $\mathbf{S}_J \dots \mathbf{S}_{\ell+1}$ (qui correspond au ℓ -ème résidu si l'on considère une factorisation hiérarchique comme présentée au paragraphe précédent) contient $n^2 \times (\frac{1}{2})^\ell$ entrées

non-nulles.

Inspirés par ces transformées rapides, nous proposons dans cette thèse de garder constante la parcimonie des facteurs creux (on a donc $\|\mathbf{S}_\ell\|_0 = p$) et de décroître exponentiellement la parcimonie du ℓ -ème résidu avec ℓ (on a donc $\|\mathbf{R}_\ell\|_0 = p_0 \times \theta^\ell$).

Détails d'implémentation

La stratégie hiérarchique proposée est résumée en Algorithme 13, où les ensembles de contraintes imposés aux facteurs doivent être spécifiés à chaque étape. \mathcal{R}_ℓ désigne l'ensemble de contrainte imposé au résidu à la ℓ -ème étape \mathbf{R}_ℓ , et \mathcal{S}_ℓ désigne l'ensemble de contrainte imposé au facteur creux à la ℓ -ème étape.

Pour la ℓ -ème factorisation en deux facteurs (ligne 3 de l'algorithme), les facteurs sont initialisés à zéro, pour des raisons de simplicité (l'initialisation aléatoire n'aboutissant pas à de meilleurs résultats). Ceci est précisé par l'argument `init=zéro` lors de l'appel à `palm4MSA`. S'agissant de l'optimisation globale (ligne 5 de l'algorithme), les facteurs sont initialisés à leur valeur courante (issue de la précédente factorisation en deux facteurs ou de la précédente optimisation globale). Ceci est précisé par l'argument `init=courant` lors de l'appel à `palm4MSA`.

Si l'on compare la stratégie hiérarchique présentée ici et l'entraînement des réseaux de neurones profonds, la factorisation en deux facteurs (ligne 3) correspond à l'étape de pré-entraînement couche par couche et l'optimisation globale (ligne 5) correspond à l'étape d'entraînement global.

Concernant le choix du pas lors des appels à `palm4MSA`, il est fixé pour toutes les expériences effectuées dans cette thèse à $0,99/L_j$, où L_j est le module de Lipschitz du gradient à l'itération courante (pour rappel, la convergence est garantie pour des pas inférieur à $1/L_j$).

Il est à noter que la stratégie hiérarchique peut tout aussi bien être appliquée dans l'autre sens (en commençant par les facteurs de gauche), en transposant la matrice d'intérêt. Nous ne présentons ici que la version commençant par les facteurs de droite parce que les notations induites sont plus simples. On peut aussi noter que d'autres critères d'arrêt qu'un nombre total de facteurs peuvent être utilisés. Par exemple, on pourrait imaginer de continuer à effectuer des factorisations en deux facteurs tant que l'erreur d'approximation à l'étape d'optimisation globale n'excède pas un certain seuil pré-défini. De plus, on pourrait imaginer d'autres types de factorisations hiérarchiques, où les deux facteurs issus de la première factorisation en deux constituent des résidus et sont factorisés à leur tour de manière itérative.

Coût de la stratégie hiérarchique

Il est possible d'analyser le coût de la stratégie hiérarchique (algorithme 13) à partir du coût de l'algorithme `palm4MSA` (en $\mathcal{O}(JT_{\text{fac}}n)$ pour J facteurs et une par-

Algorithme 13 Factorisation hiérarchique.

Entrée : Une matrice à approcher \mathbf{A} ; un nombre de facteurs J ; les ensembles de contraintes correspondant \mathcal{R}_ℓ and \mathcal{S}_ℓ , $\ell \in \{1 \dots J - 1\}$.

- 1: $\mathbf{R}_0 \leftarrow \mathbf{A}$
- 2: **for** $\ell = 1$ to $J - 1$ **do**
- 3: Factorisation du résidu $\mathbf{R}_{\ell-1}$ en 2 facteurs :
 $\lambda', \{\mathbf{F}_2, \mathbf{F}_1\} = \text{palm4MSA}(\mathbf{R}_{\ell-1}, 2, \{\mathcal{R}_\ell, \mathcal{S}_\ell\}, \text{init=zéro})$
- 4: $\mathbf{R}_\ell \leftarrow \lambda' \mathbf{F}_2$ and $\mathbf{S}_\ell \leftarrow \mathbf{F}_1$
- 5: Optimisation globale :
 $\lambda, \{\mathbf{R}_\ell, \{\mathbf{S}_j\}_{j=1}^\ell\} = \text{palm4MSA}(\mathbf{A}, \ell + 1, \{\mathcal{R}_\ell, \{\mathcal{S}_j\}_{j=1}^\ell\}, \text{init=courant})$
- 6: **end for**
- 7: $\mathbf{S}_J \leftarrow \mathbf{R}_{J-1}$

Sortie : La factorisation estimée : $\lambda, \{\mathbf{S}_j\}_{j=1}^J$.

cimonie totale de T_{fac}). Premièrement, il est clair que l'étape d'optimisation globale (ligne 5) constitue l'étape la plus coûteuse de l'algorithme. Celle-ci correspond, à la ℓ -ième étape, à un appel à `palm4MSA` avec $\ell + 1$ facteurs et une parcimonie totale donnée par les ensembles de contraintes \mathcal{R}_ℓ et $\{\mathcal{S}_j\}_{j=1}^\ell$, notée T_ℓ . La ℓ -ième étape de la stratégie hiérarchique à donc un coût en $\mathcal{O}((\ell + 1)T_\ell n)$. Lors des premières étapes, le résidu est une matrice assez dense, on peut donc supposer que T_ℓ est de l'ordre de n^2 lorsque ℓ est faible. Ceci signifie que les premières étapes de la factorisation hiérarchique ont un coût en $\mathcal{O}(n^3)$, qui domine le coût total de cette approche.

4.4.3. Illustration : retrouver la transformée de Hadamard rapide

Afin d'illustrer la méthode de factorisation hiérarchique présentée ici, et de démontrer son intérêt par rapport à la stratégie naïve qui consiste à factoriser directement une matrice d'intérêt en J facteurs creux, nous proposons de factoriser la matrice $\mathbf{A} \in \mathbb{R}^{n \times n}$ correspondant à la transformée de Hadamard en dimension $n = 2^N$. La méthode est appelée avec $J = N$ facteurs, $\mathcal{R}_\ell = \{\mathbf{R} \in \mathbb{R}^{n \times n}, \|\mathbf{R}\|_0 \leq \frac{n^2}{2^\ell}, \|\mathbf{R}\|_F = 1\}$, et $\mathcal{S}_\ell = \{\mathbf{S} \in \mathbb{R}^{n \times n}, \|\mathbf{S}\|_0 \leq 2n, \|\mathbf{S}\|_F = 1\}$.

Contrairement au cas où l'on applique directement `palm4MSA` avec $J = N$, une factorisation exacte est atteinte. En effet la première factorisation en deux ($\ell = 1$) aboutit à une factorisation exacte $\mathbf{A} = \mathbf{R}_1 \mathbf{S}_1$. Les étapes suivantes ($\ell > 1$) aboutissent elles aussi à des factorisations exactes.

Ceci est illustré en figure 15 en dimension $n = 32$. La factorisation est exacte et *aussi efficiente que la référence* montrée en figure 13. La factorisation prend moins d'une seconde pour s'effectuer pour cette dimension. Des factorisations de la matrice de Hadamard de dimension croissante jusqu'à $n = 1024$ ont montré le même comportement, prenant un temps allant jusqu'à dix minutes, et semblant

évoluer en $\mathcal{O}(n^2)$ (alors que la borne de complexité théorique est en $\mathcal{O}(n^3)$). Ceci peut s'expliquer par le fait que le comportement en $\mathcal{O}(n^3)$ ne peut être observé qu'en plus grande dimension, ou par la spécificité de la matrice de Hadamard.

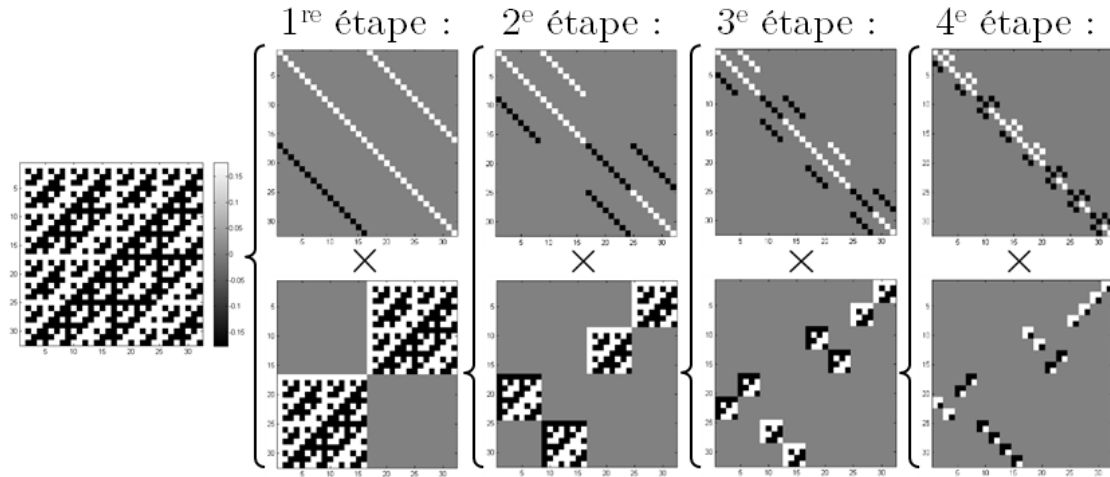


Figure 15 : Factorisation hiérarchique de la matrice de Hadamard de taille 32×32 . La matrice est factorisée de manière itérative en deux facteurs, jusqu'à avoir $J = 5$ facteurs, chacun ayant $s = 64$ entrées non-nulles.

Conclusion. Nous avons présenté dans ce chapitre le problème principal de cette thèse : l'approximation par matrices efficientes. Nous avons tout d'abord formulé le problème puis l'avons exprimé comme un problème de factorisation matricielle. Nous avons ensuite effectué un inventaire des travaux liés au problème considéré. Nous avons enfin proposé un algorithme d'optimisation afin d'attaquer le problème en question, puis une stratégie d'initialisation permettant d'atteindre de meilleures solutions.

Application à la résolution de problèmes inverses

Nous présentons dans ce chapitre une première application de l'approximation par matrice efficiente : la résolution rapide de problèmes inverses. Nous commençons par une brève introduction aux problèmes inverses linéaires, en soulignant le rôle que pourraient jouer les approximations efficaces dans l'accélération de leur résolution. Nous illustrons ensuite ceci par des expériences sur un problème de localisation de sources en magnéto-encéphalographie. Ce chapitre est dans sa majeure partie basé sur un article publié dans la revue IEEE Journal of Selected Topics in Signal Processing ([Le Magoarou et Gribonval, 2016b](#)), ainsi que sur un article présenté à la conférence European Signal Processing Conference 2015 ([Le Magoarou et al., 2015](#)).

Sommaire

5.1 Problèmes inverses	71
5.1.1 Formulation	72
5.1.2 Algorithmes de résolution	72
5.1.3 Accélération à l'aide de matrices efficaces	73
5.2 Application à la localisation de sources	74
5.2.1 Contexte	74
5.2.2 Approximations efficaces de la matrice de mesure	75
5.2.3 Expérience de localisation de sources	79

5.1. Problèmes inverses

Nous débutons ce chapitre par une brève introduction aux problèmes inverses, qui nous permet de mettre en évidence le potentiel des approximations efficaces dans ce domaine.

5.1.1. Formulation

Considérons un système pour lequel on dispose d'un modèle, c'est-à-dire que les paramètres intrinsèques du système et les grandeurs observables (par la mesure) de ce système sont liés par une relation mathématique connue. Dans le cadre d'un tel système, le but d'un problème inverse est de retrouver les causes (paramètres) d'un phénomène à partir de l'observation de leurs conséquences (observations mesurées). Nous nous intéressons dans cette thèse à une classe spécifique de problèmes inverses dits linéaires, pour lesquels paramètres $\boldsymbol{\gamma} \in \mathbb{R}^n$ et observations $\mathbf{y} \in \mathbb{R}^m$ sont liés par une application linéaire connue $\mathbf{M} \in \mathbb{R}^{m \times n}$. On a alors :

$$\mathbf{y} \approx \mathbf{M}\boldsymbol{\gamma},$$

où l'égalité n'est pas exacte car les observations sont en général bruitées, et le modèle peut également être une approximation.

Les problèmes inverses linéaires ont attiré beaucoup d'attention au cours des cinquante dernières années. Récemment, on constate un regain d'intérêt pour les problèmes inverses linéaires sous-déterminés, avec plus de paramètres que d'observations ($m < n$). Dans cette configuration, une infinité de paramètres différents peuvent mener à l'observation. Il est alors courant de privilégier les paramètres les plus simples (avec le moins d'entrées non-nulles). Ceci revient à tenter de résoudre le problème d'optimisation suivant :

$$\begin{aligned} & \underset{\boldsymbol{\gamma}}{\text{minimiser}} && \|\boldsymbol{\gamma}\|_0, \\ & \text{sujet à} && \frac{1}{2} \|\mathbf{y} - \mathbf{M}\boldsymbol{\gamma}\|_2^2 \leq \epsilon. \end{aligned} \tag{P0}$$

Ce type de problème est à la base des techniques d'échantillonnage compressé (Candès et al., 2006; Donoho, 2006). Malheureusement, il est très difficile de le résoudre de manière exacte (le problème est dit NP-difficile, et nécessite de tester tous les supports possibles pour être résolu (Tropp et Wright, 2010)), il faut alors recourir à des approximations. Nous présentons les principaux algorithmes destinés à résoudre de manière approchée ce type de problèmes inverses dans la sous-section suivante.

5.1.2. Algorithmes de résolution

De nombreux algorithmes ont été proposés au cours des deux dernières décennies afin de résoudre (P0) de manière approchée. Ce sont des algorithmes itératifs au cours desquels une estimation initiale des paramètres $\boldsymbol{\gamma}_0$ est mise à jour de manière à obtenir une séquence que l'on espère convergente vers une solution acceptable (on dénotera $\boldsymbol{\gamma}_k$ l'estimation des paramètres après la k -ième itération).

Malgré leur grande diversité, ces algorithmes impliquent tous le calcul de la quantité $\mathbf{M}^T(\mathbf{y} - \mathbf{M}\boldsymbol{\gamma}_k)$, vue comme :

- la corrélation du résidu à la k -ième itération $\mathbf{r}_k \triangleq \mathbf{y} - \mathbf{M}\boldsymbol{\gamma}_k$ avec les colonnes de la matrice de mesure. Cette interprétation est à la base des algorithmes gloutons tels que Matching Pursuit (MP) (Mallat et Zhang, 1993), Orthogonal Matching Pursuit (OMP) (Tropp et Gilbert, 2007), ou leurs nombreuses variantes (Donoho et al., 2012; Needell et Tropp, 2009). Dans toutes ces méthodes, l'estimation est initialisée à zéro et son support est mis à jour en lui ajoutant le ou les indices correspondant aux colonnes de \mathbf{M} les plus corrélées avec le résidu. Par exemple, à chaque itération d'OMP, l'indice de la colonne la plus corrélée est ajouté au support et les valeurs des coefficients dans le support sont déterminées par la résolution d'un problème des moindres carrés.
- le gradient du terme d'attache aux données $\frac{1}{2} \|\mathbf{y} - \mathbf{M}\boldsymbol{\gamma}_k\|_2^2$ par rapport aux paramètres courants $\boldsymbol{\gamma}_k$. Cette interprétation est à la base des algorithmes de seuillage itératifs. Ces algorithmes sont destinés à résoudre le problème d'optimisation lié à (P0) suivant :

$$\underset{\boldsymbol{\gamma}}{\text{minimiser}} \quad \frac{1}{2} \|\mathbf{y} - \mathbf{M}\boldsymbol{\gamma}\|_2^2 + \mu \cdot g(\boldsymbol{\gamma}).$$

Ils opèrent par descente de gradient proximale. Le choix le plus évident au regard de (P0) pour la pénalité g est la pseudo-norme ℓ_0 ($g(\boldsymbol{\gamma}) = \|\boldsymbol{\gamma}\|_0$), l'algorithme engendré est appelé Iterative Hard Thresholding (IHT) (Blumensath et Davies, 2010). Cependant, il est aussi possible de remplacer la pseudo-norme ℓ_0 par la norme ℓ_1 ($g(\boldsymbol{\gamma}) = \|\boldsymbol{\gamma}\|_1$). Ceci permet d'aboutir à un problème d'optimisation convexe proche de (P0) (on parle de relaxation convexe). Dans ce cas, l'algorithme engendré (qui correspond à l'algorithme de descente de gradient proximale présenté au chapitre 1) est appelé Iterative Soft Thresholding (IST) (Daubechies et al., 2004).

Il est à noter que le calcul de la quantité $\mathbf{M}^T(\mathbf{y} - \mathbf{M}\boldsymbol{\gamma}_k)$ constitue invariablement la procédure la plus coûteuse engendrée par tous ces algorithmes (Escande et Weiss, 2015). Pour plus de détails concernant ces algorithmes (énoncés explicites, garanties théoriques), voir les articles de revue (Bruckstein et al., 2009; Tropp et Wright, 2010).

5.1.3. Accélération à l'aide de matrices efficaces

On peut remarquer que le calcul de la quantité $\mathbf{M}^T(\mathbf{y} - \mathbf{M}\boldsymbol{\gamma}_k)$ repose principalement sur des produits matrice/vecteur mettant en jeu la matrice de mesure \mathbf{M} et sa transposée. Ainsi, disposer d'une approximation $\widehat{\mathbf{M}}$ pour laquelle on connaît

un circuit efficient associé permettrait d'accélérer le calcul de cette quantité, et, par conséquent, la résolution des problèmes inverses.

Nous proposons donc d'approcher la matrice de mesure \mathbf{M} par une FA μ ST $\widehat{\mathbf{M}} = \mathbf{S}_J \dots \mathbf{S}_1$. Le gain attendu en terme de coût des algorithmes de résolution est de l'ordre du gain relatif en complexité RCG (pour rappel, $\text{RCG} = \frac{mn}{T_{\text{fac}}}$, où $T_{\text{fac}} = \sum_{j=1}^J \|\mathbf{S}_j\|_0$).

Il est à noter qu'au cours de ces dernières années, plusieurs travaux ont eu pour but d'accélérer la résolution de problèmes inverses similaires à (P0). Ces méthodes (Gregor et LeCun, 2010; Makhzani et Frey, 2013; Sprechmann et al., 2015) sont basées sur des données d'entraînement et consistent à apprendre la fonction qui relie les observations et les paramètres via un réseau de neurones correspondant à un nombre fixe et peu élevé d'itérations d'un algorithme de seuillage itératif. Ce que nous proposons ici permet de réduire le coût des itérations des algorithmes, alors que ces approches permettent de réduire le nombre d'itérations nécessaires à calculer les paramètres. Les deux approches agissent sur des composants différents des algorithmes et peuvent donc parfaitement, en théorie, être combinées.

5.2. Application à la localisation de sources

Nous nous intéressons dans cette section à un problème inverse linéaire particulier, afin d'illustrer l'apport potentiel des approximations efficientes dans ce domaine.

5.2.1. Contexte

Les expériences conduites dans cette section se placent dans le contexte de l'imagerie cérébrale fonctionnelle utilisant des signaux captés par magnéto-encéphalographie (MEG) et électro-encéphalographie (EEG). Ce type d'imagerie permet d'appréhender le fonctionnement interne du cerveau de manière non-invasive, via la résolution d'un problème inverse particulier appelé problème inverse bio-électromagnétique. Le but de ce problème est de localiser des sources actives dans le cerveau à partir de mesures de l'activité électromagnétique à la surface du crâne. Plus concrètement, ceci revient à considérer que :

- Les paramètres $\boldsymbol{\gamma} \in \mathbb{R}^n$ correspondent aux sources d'activité électromagnétique dans le cerveau au moment de la mesure. Nous nous intéresserons ici à la localisation des sources dans le cerveau, ce qui correspond au support du vecteur $\boldsymbol{\gamma}$ (il y a n localisations possibles).
- Les observations $\mathbf{y} \in \mathbb{R}^m$ correspondent aux mesures du champ électromagnétique effectuées par les capteurs à la surface du crâne, après propagation depuis les sources (il y a m capteurs).

- La matrice de mesure $\mathbf{M} \in \mathbb{R}^{m \times n}$ correspond à la physique de la propagation des ondes électromagnétiques dans le cerveau et à travers le crâne depuis les sources vers les capteurs, on a $\mathbf{y} \approx \mathbf{M}\boldsymbol{\gamma}$. Cette matrice est calculée à partir des équations de Maxwell, en utilisant une méthode des éléments finis de frontière (BEM) implémentée par le logiciel MNE ([Gramfort et al., 2014](#)).

Il y a en général beaucoup plus de positions possibles pour les sources que de capteurs à la surface du crâne ($n \gg m$), le problème inverse est donc sous-déterminé. De plus, il est raisonnable de supposer qu'uniquement un petit nombre de sources sont actives en même temps, il est donc commun de résoudre ce problème inverse en utilisant le type de techniques de présentées dans la section précédente ([Haufe et al., 2008](#); [Gramfort et al., 2012](#)).

Le plan de cette section est simple. Nous considérons une configuration avec $n = 8193$ localisations possibles pour les sources et $m = 204$ capteurs. Nous débutons par effectuer un grand nombre de factorisations de la matrice de mesure $\mathbf{M} \in \mathbb{R}^{204 \times 8193}$ correspondante à l'aide de l'algorithme de factorisation hiérarchique présenté au chapitre précédent, afin d'évaluer le compromis erreur/complexité qu'il est possible d'atteindre. Nous utilisons ensuite ces approximations efficaces pour la résolution du problème inverse de localisation de source, ceci afin d'évaluer le potentiel apport des approximations efficaces dans le domaine des problèmes inverses.

Remarque. Dans le cas considéré ici, les sources et les capteurs ne sont pas disposés sur un échantillonnage régulier de l'espace. Il est alors difficile de définir des bases d'ondelettes orthogonales, à la fois dans le domaine des sources et dans le domaine des capteurs. Dans ce contexte, il paraît peu opportun d'utiliser les méthodes classiques de compression d'opérateur basées sur des ondelettes ([Beylkin et al., 1991](#)). Nous présentons tout de même le résultat d'une telle compression (comparée à la méthode proposée ici) en annexe B.

5.2.2. Approximations efficaces de la matrice de mesure

Dans cette sous-section, notre objectif est d'effectuer des approximations efficaces de la matrice de mesure \mathbf{M} de la forme $\widehat{\mathbf{M}} = \mathbf{S}_J \dots \mathbf{S}_1$. Pour cela, l'algorithme de factorisation hiérarchique présenté au chapitre précédent est utilisé. Il consiste à effectuer $J - 1$ factorisations en deux facteurs pour aboutir à la factorisation complète en J facteurs. En pratique, cela signifie qu'il faut choisir les paramètres correspondant au nombre total de facteurs J ainsi qu'aux ensembles de contraintes $\mathcal{R}_\ell, \mathcal{S}_\ell, \ell = 1, \dots, J - 1$ du résidu et du facteur creux à chaque factorisation en deux. Pour cela, nous allons tester un grand nombre de paramètres menant à des factorisations différentes, afin de voir émerger un compromis entre gain en complexité relative et qualité d'approximation.

Configurations

Nous avons évoqué au chapitre précédent une manière de fixer les paramètres de la factorisation hiérarchique, inspirée des transformées rapides usuelles. Cependant, cette paramétrisation n'est pas directement applicable ici car nous avons affaire à une matrice rectangulaire. Nous proposons donc la configuration suivante, correspondant à une légère adaptation des configurations précédemment évoquées.

Nous imposons que le facteur creux le plus à droite \mathbf{S}_1 , qui émerge lors de la première factorisation en deux, soit de la même taille que la matrice de mesure \mathbf{M} , mais avec des colonnes k -parcimonieuses, ce qui correspond à l'ensemble de contrainte $\mathcal{S}_1 = \{\mathbf{S} \in \mathbb{R}^{204 \times 8193}, \|\mathbf{s}_i\|_0 \leq k, \|\mathbf{S}\|_F = 1\}$. Nous imposons par ailleurs que les autres facteurs creux émergeant lors des factorisations en deux suivantes soient carrés, et ne contiennent pas plus de s entrées non-nulles, ce qui correspond à des ensembles de contrainte de forme $\mathcal{S}_\ell = \{\mathbf{S} \in \mathbb{R}^{204 \times 204}, \|\mathbf{S}\|_0 \leq s, \|\mathbf{S}\|_F = 1\}$ pour $\ell = 2, \dots, J - 1$.

Concernant le résidu à chaque factorisation en deux \mathbf{R}_ℓ , $\ell \in \{1, \dots, J - 1\}$, nous imposons qu'il soit carré, et que le nombre d'entrées non-nulles qu'il contient décroisse géométriquement avec ℓ , comme évoqué au chapitre précédent. Ceci correspond à des ensembles de contrainte de forme $\mathcal{R}_\ell = \{\mathbf{R} \in \mathbb{R}^{204 \times 204}, \|\mathbf{R}\|_0 \leq P\rho^{\ell-1}, \|\mathbf{R}\|_F = 1\}$ pour $\ell = 1, \dots, J - 1$.

En résumé, les paramètres de la factorisation sont le nombre de facteurs J , la parcimonie par colonne du premier facteur creux k , la parcimonie des autres facteurs creux s , la parcimonie du premier résidu P et le taux de décroissance de la parcimonie du résidu ρ . Nous faisons varier ces paramètres de la façon suivante :

- Nombre de facteurs $J \in \{2, \dots, 10\}$.
- Parcimonie par colonne du premier facteur creux $k \in \{5, 10, 15, 20, 25, 30\}$.
- Parcimonie des autres facteurs creux $s \in \{2m, 4m, 8m\}$.

Les paramètres contrôlant le comportement du résidu ont empiriquement peu d'influence, nous les fixons à $P = 1.4 \times m^2$ et $\rho = 0.8$. La configuration de factorisation est résumée en figure 16, où la parcimonie de chaque facteur de la factorisation finale $\mathbf{S}_J \dots \mathbf{S}_1$ est indiquée.

Résultats de factorisation

Le résultat des 150 configurations de factorisations différentes impliquées par la paramétrisation utilisée est donné en figure 17. Sur cette figure, l'axe des abscisses représente le gain relatif en complexité RCG, et l'axe des ordonnées l'erreur relative d'approximation, calculée via l'expression $\|\mathbf{M} - \widehat{\mathbf{M}}\|_2 / \|\mathbf{M}\|_2$. Cette mesure correspond à l'erreur relative en norme d'opérateur, elle est générale et très standard.

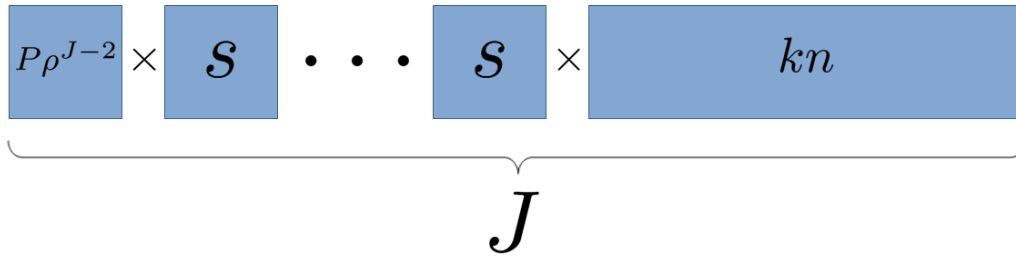


Figure 16 : Configuration de factorisation : chaque facteur (de \mathbf{S}_J à l'extrême gauche à \mathbf{S}_1 à l'extrême droite) est représenté avec sa parcimonie totale.

On pourrait imaginer mesurer l'erreur d'approximation par une fonction $d(\mathbf{M}, \widehat{\mathbf{M}})$ spécialement conçue pour une application aux problèmes inverses.

En analysant cette figure, on observe que :

- Le gain en complexité relative RCG est essentiellement contrôlé par le paramètre k qui fixe la parcimonie du facteur creux \mathbf{S}_1 . Ceci semble parfaitement naturel, car ce facteur est bien plus grand que tous les autres (il est de taille 204×8193 alors que les autres sont de taille 204×204).
- Le compromis entre gain en complexité relative et qualité d'approximation, pour une valeur de k donnée, est principalement contrôlé par le nombre de facteurs J . Un grand nombre de facteurs mène à un gain en complexité relative plus élevé, mais un trop grand nombre de facteurs mène à une plus grande erreur relative. Il est intéressant de remarquer que prendre seulement $J = 2$ facteurs ne mène jamais au meilleur compromis, ce qui montre l'intérêt de réelles approximation *multi-couches*.
- Encore pour une valeur de k fixée, on peut distinguer des courbes proches les unes des autres correspondant à différentes valeurs du paramètre s qui contrôle la parcimonie des autres facteurs. Ce paramètre contrôle l'espacement horizontal entre deux points consécutifs sur la même courbe. Prendre un s grand permet d'effectuer des factorisations avec un nombre de facteurs J élevé sans faire augmenter l'erreur d'approximation, mais en retour mène à un gain en complexité relative plus faible pour le même nombre de facteurs.

En résumé, on peut distinguer un compromis entre gain en complexité relative et qualité d'approximation auquel on s'attendait. La configuration menant à la meilleure qualité d'approximation pour chaque valeur de k est indiqué en figure 17, et ceci donne $\widehat{\mathbf{M}}_{25}$, $\widehat{\mathbf{M}}_{16}$, $\widehat{\mathbf{M}}_{11}$, $\widehat{\mathbf{M}}_8$, $\widehat{\mathbf{M}}_7$, $\widehat{\mathbf{M}}_6$, où l'indice indique le gain en complexité relative atteint par chaque configuration (arrondi à l'entier le plus proche). Par exemple, l'approximation efficiente $\widehat{\mathbf{M}}_6$ permet de multiplier des vecteurs (sous forme directe ou transposée) avec 6 fois moins d'opération arithmétiques que \mathbf{M}

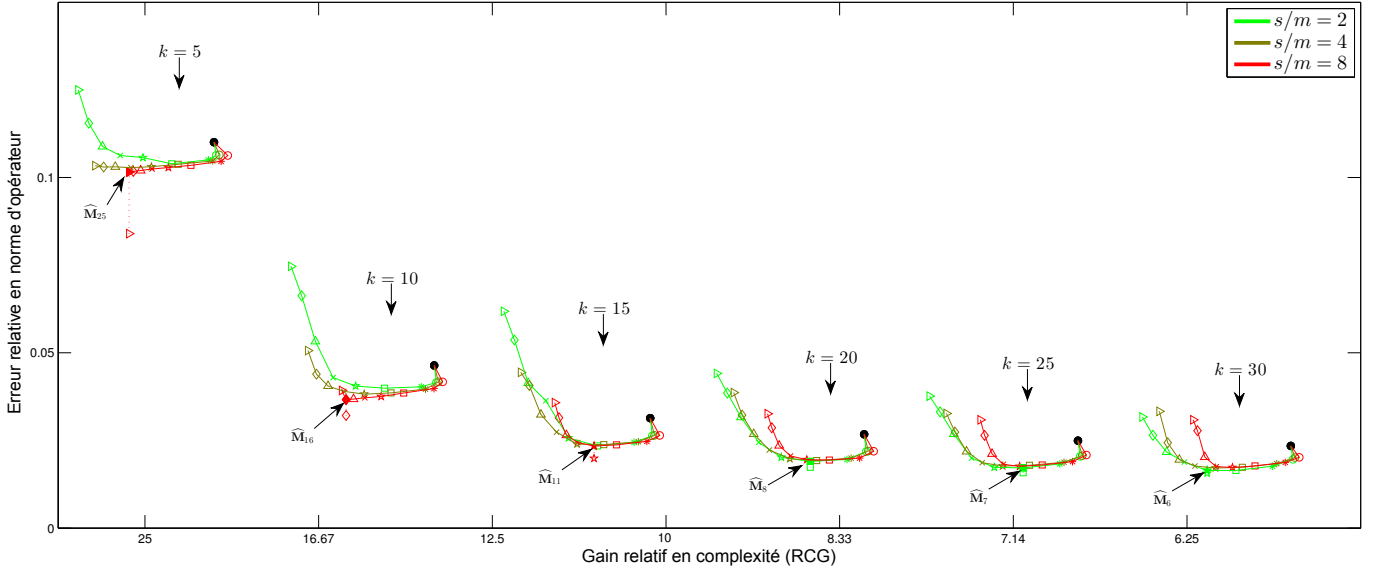


Figure 17 : Résultats de multiples factorisations d’une matrice de mesure de taille $m \times n = 204 \times 8193$ correspondant à un problème inverse bio-électromagnétique. Chaque point correspond à une configuration de factorisation différente. La forme des symboles indique le nombre total de facteurs J (\bullet : $J = 2$, \circ : $J = 3$, $*$: $J = 4$, \square : $J = 5$, \star : $J = 6$, \times : $J = 7$, \triangle : $J = 8$, \diamond : $J = 9$, \triangleright : $J = 10$), et leur couleur la valeur du paramètre s contrôlant leur parcimonie.

(économisant 84% de calculs), et $\widehat{\mathbf{M}}_{25}$ permet de multiplier des vecteurs avec 25 fois moins d’opérations arithmétiques que \mathbf{M} (économisant 96% de calculs). Ces six matrices sont celles que nous allons utiliser pour la résolution rapide d’un problème inverse de localisation de source dans la sous-section suivante.

Remarques.

- Une qualité d’approximation très légèrement supérieure peut être atteinte en imposant une contrainte de parcimonie globale au facteur creux le plus à droite \mathbf{S}_1 , de type $\mathcal{S}_1 = \{\mathbf{S} \in \mathbb{R}^{204 \times 8193}, \|\mathbf{S}\|_0 \leq kn, \|\mathbf{S}\|_F = 1\}$. Ceci est représenté en figure 17 par les points reliés par une ligne pointillée aux six matrices évoquées au paragraphe précédent. Cependant une telle contrainte autorise l’apparition de colonnes nulles dans les approximations efficaces, ce qui n’est pas désirable pour l’application de localisation de source considérée à la sous-section suivante (cela donnerait lieu à des localisations impossibles à retrouver).
- Clairement, il n’est pas envisageable de procéder à un grand nombre de factorisation de la matrice de mesure à chaque fois qu’on a besoin de résoudre un problème inverse. Malheureusement, Le grand nombre de paramètres nécessaires à la factorisation creuse multi-couche rend difficile leur réglage. Afin de limiter

ce problème, nous proposons une paramétrisation simplifiée et une méthode alternative au chapitre 6, et évoquons des pistes d'amélioration en conclusion de cette thèse.

5.2.3. Expérience de localisation de sources

Nous nous intéressons maintenant à l'impact du remplacement de la matrice de mesure $\mathbf{M} \in \mathbb{R}^{204 \times 8193}$ par une approximation efficace dans le cadre d'une expérience de localisation de sources. Pour cette expérience synthétique, deux sources localisées uniformément aléatoirement dans le cerveau sont activées avec des poids aléatoires gaussiens, ce qui donne un vecteur $\boldsymbol{\gamma} \in \mathbb{R}^{8193}$ 2-parcimonieux, dont le support indique la localisation des sources. Dans le domaine de l'imagerie magnéto-encéphalographique, il est courant de considérer un petit nombre de sources (typiquement deux ou trois). À l'observation de $y \triangleq \mathbf{M}\boldsymbol{\gamma}$, l'objectif est de retrouver (le support de) $\boldsymbol{\gamma}$. L'expérience que nous menons ici revient à résoudre le problème inverse pour obtenir $\hat{\boldsymbol{\gamma}} \in \mathbb{R}^{8193}$ à partir des observations $\mathbf{y} \in \mathbb{R}^{204}$, en utilisant la vraie matrice de mesure \mathbf{M} ou une approximation efficace lors de la résolution, afin de comparer les deux cas.

Trois algorithmes de résolution ont été testés : Orthogonal Matching Pursuit (OMP) (Tropp et Gilbert, 2007) (en choisissant 2 atomes), ℓ_1 -regularized least squares (ℓ_1 ls) (Kim et al., 2007) et Iterative Hard Thresholding (IHT) (Blumensath et Davies, 2008)). Les trois méthodes ont mené à des résultats qualitativement similaires, nous présenterons alors ici seulement les résultats obtenus avec OMP.

Les matrices utilisées lors de la résolution sont la vraie matrice de mesure \mathbf{M} et ses approximations efficaces $\widehat{\mathbf{M}}_{25}$, $\widehat{\mathbf{M}}_{16}$, $\widehat{\mathbf{M}}_{11}$, $\widehat{\mathbf{M}}_8$, $\widehat{\mathbf{M}}_7$ et $\widehat{\mathbf{M}}_6$. Le gain attendu en terme de complexité de résolution est de l'ordre du gain en complexité relative RCG. En effet le coût d'OMP est dominé par des multiplications avec la transposée de la matrice de mesure \mathbf{M}^T .

Trois configurations ont été distinguées lors de la génération de l'emplacement des sources, en terme de distance d (en centimètres) entre les deux sources. Pour chacune de ces configurations, 500 vecteurs $\mathbf{y} = \mathbf{M}\boldsymbol{\gamma}$ ont été générés, et OMP a été appliqué avec chacune des matrices de résolution considérées. La distance entre chaque source réelle et la source retrouvée la plus proche a ensuite été mesurée.

Résultats de localisation

La figure 18 donne les statistiques de cette distance. L'analyse de cette figure donne lieu aux commentaires suivants :

- Comme attendu, la localisation des sources est mieux estimée lorsque les sources réelles sont séparées par une plus grande distance d , et ceci indépendamment du choix de la matrice de résolution.

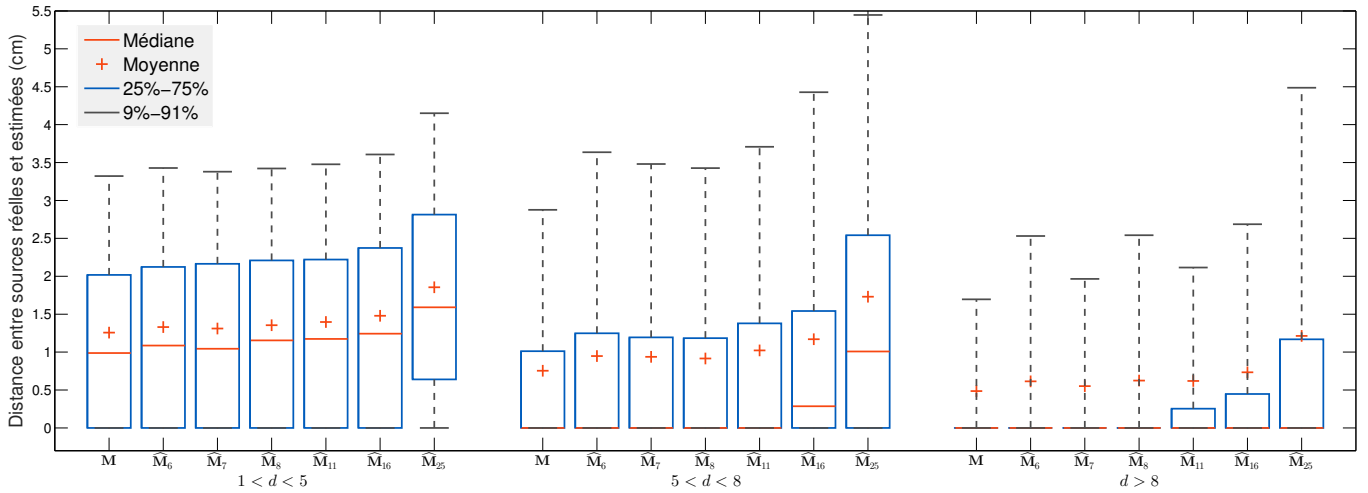


Figure 18 : Performance de localisation (distance entre sources réelles et retrouvées), obtenue avec différentes matrices de résolution, et pour différentes distance entre les sources réelles.

- Plus significativement, l'estimation est quasiment aussi performante en utilisant les approximations efficaces $\widehat{\mathbf{M}}_6$, $\widehat{\mathbf{M}}_7$, $\widehat{\mathbf{M}}_8$ and $\widehat{\mathbf{M}}_{11}$ qu'en utilisant la vraie matrice de mesure \mathbf{M} , alors que les approximations efficaces sont bien plus avantageuses en terme de complexité calculatoire. Par exemple, dans le cas de sources réelles séparées d'une distance comprise entre cinq et huit centimètres ($5 < d < 8$), les FA μ ST permettent de retrouver exactement la localisation des sources dans plus de 50% des cas, ce qui est quasiment aussi bien qu'en utilisant la matrice \mathbf{M} .
- La performance des deux autres approximation efficaces $\widehat{\mathbf{M}}_{16}$ et $\widehat{\mathbf{M}}_{25}$ est un peu moins bonne, mais leur utilisation est encore plus avantageuse en terme de complexité calculatoire (16 et 25 fois moins de calculs). Cependant les résultats obtenus en utilisant ces approximations sont tout de même acceptables, puisque par exemple, dans le cas de source bien séparées ($d > 8$), elles permettent de retrouver exactement la localisation des sources dans plus de 50% des cas.

Ces observations confirment qu'il est possible de remplacer la matrice de mesure d'un problème inverse par une approximation efficace, ce qui permet d'obtenir un gain substantiel en terme de complexité calculatoire, et ceci sans occasionner une perte de précision trop importante.

Gains effectifs en temps de résolution

Afin de mesurer le gain en complexité obtenu lors de l'utilisation d'une approximation efficiente $\widehat{\mathbf{M}}$ en lieu et place de la véritable matrice de mesure \mathbf{M} , nous avons jusque là considéré uniquement le gain relatif en complexité RCG. Celui-ci correspond au nombre de multiplications scalaires impliquées dans un produit par \mathbf{M} ou \mathbf{M}^T divisé par le nombre de multiplications scalaires impliquées dans un produit par $\widehat{\mathbf{M}}$ ou $\widehat{\mathbf{M}}^T$. On espère que ce gain en terme de nombre de multiplications se répercute sur le temps nécessaire à effectuer une multiplication, et, par voie de conséquence, sur le temps de résolution du problème inverse. Afin de vérifier cela, nous comparons le gain effectif en temps obtenu sur la résolution du problème de localisation de source (temps médian de la résolution avec \mathbf{M} divisé par le temps médian de la résolution avec $\widehat{\mathbf{M}}$) avec le gain relatif en complexité RCG (pour les FA μ ST $\widehat{\mathbf{M}}_6, \widehat{\mathbf{M}}_7, \widehat{\mathbf{M}}_8, \widehat{\mathbf{M}}_{11}, \widehat{\mathbf{M}}_{16}$, et $\widehat{\mathbf{M}}_{25}$) en figure 19. Pour ces mesures de temps, les calculs avec les FA μ ST sont implémentés par une bibliothèque C++ disponible à l'adresse <http://faust.gforge.inria.fr>. À l'analyse de cette figure, on voit que le gain effectif en temps est inférieur à RCG. Par exemple, le gain effectif en temps de la FA μ ST $\widehat{\mathbf{M}}_{16}$ est de 4.8, alors que son RCG est de 16. Cette différence est certainement due à l'implémentation de la multiplication par les FA μ ST, moins optimisée que celle de la multiplication par des matrices denses telles que \mathbf{M} , qui utilise en général des bibliothèques spécialisées (Eigen en l'occurrence dans le cas présenté ici (Guennebaud et al., 2010)). Le gain en temps de résolution maximal obtenu est proche de 6, pour la matrice efficiente $\widehat{\mathbf{M}}_{25}$ (on résout le problème inverse six fois plus vite avec la FA μ ST $\widehat{\mathbf{M}}_{25}$ qu'avec la véritable matrice de mesure \mathbf{M}). Ce gain effectif maximal est environ quatre fois inférieur au gain théorique RCG maximal. Un des défis futurs qui se pose afin de rendre les FA μ ST plus attrayantes est de réduire cette différence entre le gain théorique et le gain effectif.

Conclusion. Nous avons présenté dans ce chapitre une première application de l'approximation par matrice efficientes : les problèmes inverses. Nous avons tout d'abord introduit les problèmes inverses, puis indiqué dans quelle mesure l'approximation par matrice efficiente pouvait en accélérer la résolution. Nous avons ensuite réalisé une étude de cas sur un problème de localisation de source cérébrale, confirmant que les approximations efficientes sont en effet potentiellement intéressantes dans ce domaine.

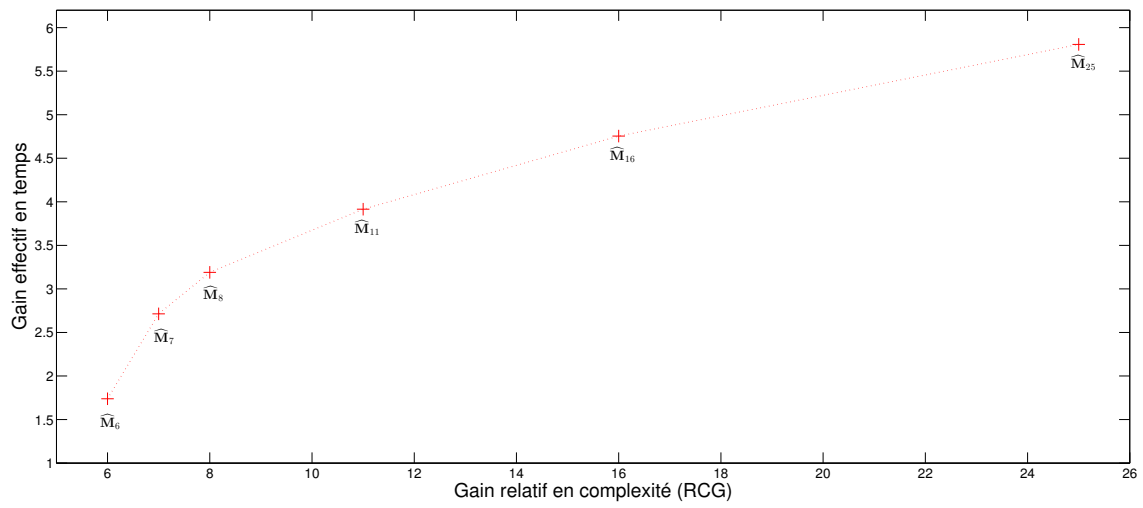


Figure 19 : Comparaison entre le gain théorique RCG et le gain en temps observé en pratique sur la résolution du problème inverse.

Application à la transformée de Fourier rapide sur graphe

Nous appliquons dans ce chapitre l'approximation par matrice efficace à la transformée de Fourier sur graphe. Nous commençons par motiver cette application en introduisant brièvement le traitement du signal sur graphe, en mettant en évidence le besoin d'une transformée de Fourier rapide dans ce contexte. Nous proposons ensuite deux approches distinctes permettant d'aboutir à une transformée de Fourier rapide approchée sur graphe, et les comparons de manière expérimentale. Ce chapitre est en partie basé sur un article présenté à la conférence IEEE International Conference on Acoustics, Speech and Signal Processing 2016 ([Le Magoarou et Gribonval, 2016a](#)).

Sommaire

6.1	Traitement du signal sur graphe	84
6.1.1	Graphes et matrices associées	84
6.1.2	Signaux et transformée de Fourier sur graphe	85
6.1.3	Transformée de Fourier rapide	88
6.2	FFT sur graphe par factorisation creuse multi-couche de la matrice de Fourier	88
6.2.1	Principe général	88
6.2.2	Ensembles de contraintes	89
6.3	FFT sur graphe par diagonalisation gloutonne du Laplacien	90
6.3.1	Principe général	91
6.3.2	Travaux liés	92
6.3.3	Algorithme proposé	92
6.4	Comparaison des méthodes	94
6.4.1	Évaluation des transformées de Fourier approchées	95
6.4.2	Exemple d'application : filtrage	98

6.1. Traitement du signal sur graphe

Le traitement du signal sur graphe est un domaine de recherche récent (Shuman et al., 2013), dont le but est de généraliser les outils du traitement du signal classique (filtrage, analyse fréquentielle, ondelettes, etc.) à des signaux définis sur les nœuds d'un graphe. Nous en donnons ici un bref aperçu, en nous concentrant sur l'opérateur qui nous intéresse plus particulièrement : la transformée de Fourier sur graphe.

6.1.1. Graphes et matrices associées

Nous commençons par donner ici quelques définitions utiles à la compréhension du chapitre. Pour une référence plus complète sur les graphes, voir (Bondy et Murty, 1976).

Définition 7. *Un graphe non-orienté \mathcal{G} est constitué d'un ensemble de nœuds $\mathcal{V} = \{v_1, \dots, v_n\}$ et d'un ensemble \mathcal{E} d'arêtes de la forme $\{v_i, v_j\}$, où $v_i, v_j \in \mathcal{V}$ sont deux nœuds distincts du graphe. On utilisera la notation $\mathcal{G} = (\mathcal{V}, \mathcal{E})$. Le graphe est dit pondéré s'il existe une fonction $w : \mathcal{E} \rightarrow \mathbb{R}_+^*$ qui à chaque arête $\{v_i, v_j\} \in \mathcal{E}$ associe un poids strictement positif $w(\{v_i, v_j\})$, que l'on notera w_{ij} . L'ensemble des poids définit la matrice d'adjacence $\mathbf{W} \in \mathbb{R}^{n \times n}$ du graphe ($w_{ij} = 0$ si $\{v_i, v_j\} \notin \mathcal{E}$).*

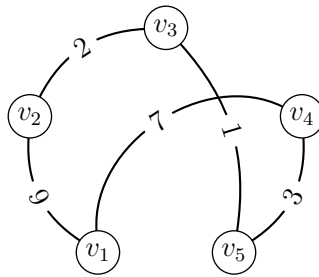
Dans cette définition, le poids w_{ij} peut être interprété comme une mesure de proximité entre les nœuds v_i et v_j . Nous considérons dans cette thèse, sauf mention contraire, des graphes non-orientés pondérés. Selon les définitions données ci-dessus, un graphe est entièrement décrit par sa matrice d'adjacence. Par exemple, considérons le graphe quelconque donné par :

- l'ensemble de nœuds $\mathcal{V} = \{v_1, v_2, v_3, v_4, v_5\}$,
- l'ensemble d'arêtes $\mathcal{E} = \left\{ \{v_1, v_2\}, \{v_1, v_4\}, \{v_2, v_3\}, \{v_3, v_5\}, \{v_4, v_5\} \right\}$,
- les poids $w(\{v_1, v_2\}) = 9$, $w(\{v_1, v_4\}) = 7$, $w(\{v_2, v_3\}) = 2$, $w(\{v_3, v_5\}) = 1$, $w(\{v_4, v_5\}) = 3$.

Il correspond au schéma suivant

qui est équivalent à la matrice d'adjacence

$$\mathbf{W} = \begin{pmatrix} 0 & 9 & 0 & 7 & 0 \\ 9 & 0 & 2 & 0 & 0 \\ 0 & 2 & 0 & 0 & 1 \\ 7 & 0 & 0 & 0 & 3 \\ 0 & 0 & 1 & 3 & 0 \end{pmatrix}.$$



À partir de la matrice d'adjacence, on définit deux autres matrices importantes dans l'analyse des graphes.

Définition 8. La matrice des degrés $\mathbf{D} \in \mathbb{R}^{n \times n}$ d'un graphe pondéré \mathcal{G} est une matrice diagonale dont l'entrée d_{ii} correspond à la somme des poids des arêtes contenant le nœud v_i , ceci s'exprime $d_{ii} \triangleq \sum_{k=1}^n w_{ik}$. Le Laplacien combinatoire (aussi appelée la matrice laplacienne combinatoire) du graphe correspond à la différence de la matrice des degrés et de la matrice d'adjacence du graphe : $\mathbf{L} \triangleq \mathbf{D} - \mathbf{W}$.

La diagonale de la matrice des degrés peut s'interpréter comme une fonction qui associe à chaque nœud du graphe son importance dans le graphe, mesurée par la force de connexions liant chaque nœud au reste du graphe. Il existe par ailleurs d'autres types de matrices laplaciennes (Laplacien normalisé, Laplacien de marche aléatoire), mais nous ne considérons dans ce chapitre que le Laplacien combinatoire, que nous appellerons par la suite simplement Laplacien, pour des raisons de simplicité. Par exemple, en considérant le graphe évoqué ci-dessus on a

$$\mathbf{D} = \begin{pmatrix} 16 & 0 & 0 & 0 & 0 \\ 0 & 11 & 0 & 0 & 0 \\ 0 & 0 & 3 & 0 & 0 \\ 0 & 0 & 0 & 10 & 0 \\ 0 & 0 & 0 & 0 & 4 \end{pmatrix}, \mathbf{L} = \begin{pmatrix} 16 & -9 & 0 & -7 & 0 \\ -9 & 11 & -2 & 0 & 0 \\ 0 & -2 & 3 & 0 & -1 \\ -7 & 0 & 0 & 10 & -3 \\ 0 & 0 & -1 & -3 & 4 \end{pmatrix}.$$

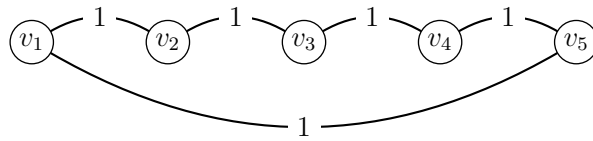
6.1.2. Signaux et transformée de Fourier sur graphe

Définissons maintenant l'objet d'intérêt principal de ce chapitre : les signaux sur graphe.

Définition 9. Un signal sur un graphe $\mathcal{G} = (\mathcal{V}, \mathcal{E})$ est une fonction $f : \mathcal{V} \rightarrow \mathbb{R}$ qui à chaque nœud du graphe associe un nombre réel qui correspond à la valeur prise par le signal sur ce nœud. Il se représente par un vecteur $\mathbf{f} \in \mathbb{R}^n$ de forme $\mathbf{f} = (f(v_1), \dots, f(v_n))^T$.

Un signal sur graphe peut donc s'interpréter par un signal défini sur un domaine discret dont la topologie est donnée par le graphe.

Traitement du signal classique. Dans le cas du traitement du signal discret classique, on suppose implicitement que le domaine de définition du signal correspond à une grille régulière, issue d'un échantillonnage uniforme de la grandeur définissant le domaine, le plus souvent le temps ou l'espace. En 1D, ceci correspond à un graphe où chaque nœud est lié par deux arêtes de même poids aux deux nœuds correspondant aux deux points directement adjacents dans le domaine de définition du signal. Par exemple si l'on considère $n = 5$ points dans le domaine de définition, le traitement du signal classique revient à considérer le graphe donné par le schéma et les matrices suivants (en considérant une extension périodique du signal) :



$$\mathbf{W}_c = \begin{pmatrix} 0 & 1 & 0 & 0 & 1 \\ 1 & 0 & 1 & 0 & 0 \\ 0 & 1 & 0 & 1 & 0 \\ 0 & 0 & 1 & 0 & 1 \\ 1 & 0 & 0 & 1 & 0 \end{pmatrix}, \mathbf{L}_c = \begin{pmatrix} 2 & -1 & 0 & 0 & -1 \\ -1 & 2 & -1 & 0 & 0 \\ 0 & -1 & 2 & -1 & 0 \\ 0 & 0 & -1 & 2 & -1 \\ -1 & 0 & 0 & -1 & 2 \end{pmatrix}.$$

Il est aisé de remarquer que dans ce cas, la multiplication d'un signal quelconque \mathbf{f} par le Laplacien classique \mathbf{L}_c correspond à une analyse de la variation locale du signal. En effet, en notant $\mathbf{g} \triangleq \mathbf{L}_c \mathbf{f}$, on a

$$g_i = (f_i - f_{i-1}) + (f_i - f_{i+1}),$$

en considérant que $f_0 = f_n$ et $f_{n+1} = f_1$. Le Laplacien étant une matrice symétrique définie positive, il peut donc se décomposer en vecteurs/valeurs propres de la manière suivante :

$$\mathbf{L}_c = \mathbf{U}_c \mathbf{\Lambda}_c \mathbf{U}_c^T,$$

où $\mathbf{U}_c \in \mathbb{R}^{n \times n}$ est une matrice orthogonale de vecteurs propres et $\mathbf{\Lambda}_c \in \mathbb{R}^{n \times n}$ est une matrice diagonale de valeurs propres dont tous les éléments diagonaux sont supérieurs à 0. Dans cette décomposition, les vecteurs propres (colonnes de \mathbf{U}_c) associés aux valeurs propres les plus petites correspondent aux directions engendrant peu de variations locales (les basses fréquences) et les vecteurs propres associés aux valeurs propres les plus grandes correspondent aux directions engendrant beaucoup de variations locales (les hautes fréquences). Les colonnes de la matrice orthogonale \mathbf{U}_c constituent alors une base de Fourier pour la grille régulière considéré en traitement du signal classique.

Il est à noter que la matrice de Fourier utilisée habituellement en traitement du signal est une matrice complexe. Cependant, le Laplacien étant réel et symétrique, il existe une base orthogonale réelle de vecteurs propres, qui bien que ne jouissant pas de toutes les propriétés mathématiques intéressantes de la base complexe, lui est équivalente en terme de traitement du signal (décomposition fréquentielle, filtrage).

Généralisation aux graphes quelconques. Dans le cas d'un graphe quelconque, l'interprétation de la multiplication par le Laplacien \mathbf{L} comme une analyse de la variation locale du signal par rapport à la topologie définie par le graphe tient toujours. En effet, si $\mathbf{g} = \mathbf{L}\mathbf{f}$ on a

$$g_i = \sum_{k=1}^n w_{ik} \cdot (f_i - f_k),$$

la i -ème coordonnée du vecteur \mathbf{g} correspond à une somme des différences entre la valeur que prend le signal \mathbf{f} au i -ème nœud du graphe et celles qu'il prend aux nœuds adjacents, pondérées par le poids des arêtes reliant le i -ème nœud à chacun de ses voisins. La transformée de Fourier sur graphe peut alors être définie (Shuman et al., 2013) par analogie avec le cas classique : *une base de Fourier est une base orthogonale de vecteurs propres du Laplacien du graphe*. On a en effet comme dans le cas classique

$$\mathbf{L} = \mathbf{U}\mathbf{\Lambda}\mathbf{U}^T,$$

avec \mathbf{U} orthogonale et $\mathbf{\Lambda}$ diagonale positive. Pour un signal sur graphe quelconque \mathbf{x} , sa transformée de Fourier sur graphe prend alors par définition la forme

$$\mathbf{y} = \mathbf{U}^T \mathbf{x},$$

et l'orthogonalité de la matrice \mathbf{U} permet de retrouver le signal à partir de sa transformée de Fourier de la manière suivante :

$$\mathbf{x} = \mathbf{U}\mathbf{y}.$$

La transformée de Fourier est à la base des techniques de traitement du signal sur graphe. Par exemple on peut définir des ondelettes sur graphe dans le domaine de Fourier du graphe, en terme de filtrages passe-bande successifs (Hammond et al., 2011).

Il est à noter qu'il est possible de définir la transformée de Fourier autrement, en utilisant une autre analogie avec le cas du traitement du signal classique. Par exemple, une approche récente (Sandryhaila et Moura, 2013) consiste à définir la transformée de Fourier sur graphe comme la base de vecteurs propres généralisés de la matrice d'adjacence du graphe transposée, qui s'interprète comme l'opérateur de retard pour les signaux sur graphe. Nous ne considérons dans ce chapitre que la

transformée de Fourier définie par (Shuman et al., 2013), comme la base des vecteurs propres du Laplacien. Pour une comparaison de ces deux analogies, et plus de détails concernant le traitement du signal sur graphe dans son ensemble, voir (Tremblay, 2014).

6.1.3. Transformée de Fourier rapide

La matrice de Fourier \mathbf{U} étant dense en général, le calcul de la transformée de Fourier et de son inverse nécessite a priori $\mathcal{O}(n^2)$ opérations arithmétiques. Cependant, dans le cas de la transformée de Fourier classique, un algorithme rapide permettant d'effectuer la multiplication par \mathbf{U}_c en $\mathcal{O}(n \log n)$ existe : la transformée de Fourier rapide (souvent abrégée FFT pour Fast Fourier Transform (Cooley et Tukey, 1965)). Il correspond à une factorisation de la matrice de Fourier classique de la forme

$$\mathbf{U}_c = \mathbf{F}_I \dots \mathbf{F}_1,$$

avec $I = \log_2 n$ et chacun des facteurs \mathbf{F}_i ayant $2n$ entrées non-nulles (lorsque n est une puissance de deux). Peut-on généraliser ce type de factorisation creuse multi-couche (de manière approchée) à des matrices de Fourier \mathbf{U} correspondant à des graphes quelconque ? C'est à cette question que nous allons nous intéresser dans ce chapitre. Nous proposons deux méthodes pour parvenir à cette généralisation. La première consiste à directement approcher la matrice de Fourier \mathbf{U} par une matrice efficace. La seconde consiste à diagonaliser de manière approchée le Laplacien \mathbf{L} à l'aide d'une matrice efficace. Nous allons tout d'abord présenter les deux méthodes, puis les comparer expérimentalement et discuter leur potentiels respectifs en terme d'applications.

6.2. FFT sur graphe par factorisation creuse multi-couche de la matrice de Fourier

Nous décrivons dans cette section une méthode d'approximation directe de la matrice de Fourier sur graphe \mathbf{U} . Celle-ci consiste en une factorisation de la matrice \mathbf{U} en J facteurs creux $\mathbf{S}_1, \dots, \mathbf{S}_J$.

6.2.1. Principe général

Cette méthode n'est autre qu'une application directe de la méthode d'approximation par matrice efficace introduite au chapitre 4. Elle repose sur la résolution approchée du problème d'optimisation

$$\underset{\mathbf{S}_1, \dots, \mathbf{S}_J}{\text{minimiser}} \quad \frac{1}{2} \|\mathbf{U} - \mathbf{S}_J \dots \mathbf{S}_1\|_F^2 + \sum_{j=1}^J \delta_{\mathcal{S}_j}(\mathbf{S}_j), \quad (\text{PF})$$

à l'aide de la stratégie de factorisation hiérarchique présentée au chapitre 4, afin d'obtenir une approximation efficiente $\hat{\mathbf{U}} \triangleq \mathbf{S}_J \dots \mathbf{S}_1$. Il est alors nécessaire de spécifier les ensembles de contraintes correspondant aux facteurs creux \mathcal{S}_ℓ et au résidu \mathcal{R}_ℓ pour chacune des factorisations en deux facteurs ($\ell = 1, \dots, J - 1$) qui composent la stratégie hiérarchique.

6.2.2. Ensembles de contraintes

Comme précisé dans la section précédente, la matrice de Fourier \mathbf{U}_c utilisée en traitement du signal classique peut se factoriser de manière exacte en $I = \log_2 n$ facteurs ayant chacun $2n$ entrées non-nulles (lorsque n est une puissance de deux) : $\mathbf{U}_c = \mathbf{F}_I \dots \mathbf{F}_1$. Dans le contexte de la stratégie hiérarchique, une telle factorisation correspondant à la FFT classique pourrait être retrouvée en effectuant $I - 1$ factorisations en deux facteurs, le facteur de droite (le facteur creux) ayant toujours $2n$ entrées non-nulles et le facteur de gauche (le résidu) ayant $\frac{n^2}{2^\ell}$ entrées non-nulles à la ℓ -ième factorisation en deux. En effet, on a $\|\mathbf{F}_\ell\|_0 = 2n$ et $\|\mathbf{F}_I \dots \mathbf{F}_{\ell+1}\|_0 = \frac{n^2}{2^\ell}$ pour tout $\ell \in \{1, \dots, I - 1\}$. Adopter cette stratégie hiérarchique revient à diviser le nombre d'entrées non-nulles du résidu par deux à chaque étape, on dit alors qu'on adopte un taux de décroissance de deux. Le résultat en terme de parcimonie des facteurs obtenu avec une telle stratégie est décrit par la figure 20.

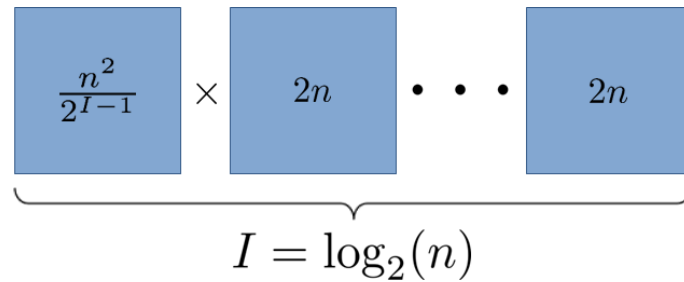


Figure 20 : Configuration de factorisation idéale correspondant au cas de la grille régulière 1D. Chaque facteur est représenté par un carré, le nombre d'entrées non-nulle correspondant écrit à l'intérieur.

Ici, notre objectif est de généraliser la factorisation creuse multi-couche de la matrice de Fourier classique \mathbf{U}_c à des matrices de Fourier \mathbf{U} correspondant à des graphes quelconques, de manière à avoir une FFT *approchée* pour les signaux sur graphe. Dans la mesure où dans le cas de graphes quelconques, on ne sait pas si une

telle factorisation exacte existe, on propose de relâcher les contraintes de parcimonie appliquées lors de la factorisation des trois façons suivantes :

- Le nombre de facteurs J est réduit d'une constante C_1 ($J = \log_2 n - C_1$). Ceci mène potentiellement à des factorisations de meilleure qualité (en terme d'approximation), car le nombre total d'entrées non-nulle est réduit à chaque étape de la factorisation hiérarchique (plus de facteurs signifie moins d'entrées non-nulles). De plus, ceci permet de réduire le temps de factorisation, celui-ci étant proportionnel au nombre de facteurs en utilisant la stratégie hiérarchique.
- Le taux de décroissance de la parcimonie du résidu est réduit à C_2 avec $1 < C_2 \leq 2$. Ceci revient à diviser le degré de parcimonie du résidu par un nombre plus petit que deux à chaque étape de la factorisation hiérarchique, ce qui mène potentiellement à une factorisation de meilleure qualité (il y a plus d'entrées non-nulles dans le résidu à chaque étape).
- La parcimonie totale de chaque facteur est de plus multiplié par une constante C_3 avec $C_3 \geq 1$. Ceci permet aussi, en autorisant plus d'entrées non-nulles dans les facteurs, d'obtenir des factorisations potentiellement de meilleure qualité.

La configuration de factorisation en terme de parcimonie est donc contrôlée uniquement par les trois paramètres C_1 , C_2 et C_3 . le cas idéal de la FFT classique correspond aux paramètres $C_1 = 0$, $C_2 = 2$ et $C_3 = 1$, et d'autres valeurs pour ces paramètres mènent à des factorisation relâchées, avec plus d'entrées non-nulles, mais plus précises. Le résultat en terme de parcimonie des facteurs obtenus avec cette stratégie relâchée est décrit par la figure 21. Une FA μ ST dotée d'une telle configuration de parcimonie contient au total $C_3 \cdot (2n(\log_2 n - C_1) + n^2 \cdot C_2^{1+C_1-\log_2 n})$ entrées non-nulle, ce qui correspond à une complexité en $\mathcal{O}(n^{2-\log_2 C_2})$ pour $1 < C_2 < 2$ et en $\mathcal{O}(n \log_2 n)$ pour $C_2 = 2$.

$$\underbrace{\frac{n^2 \cdot C_3}{C_2^{J-1}} \times 2n \cdot C_3 \cdot \dots \cdot 2n \cdot C_3}_{J = \log_2(n) - C_1}$$

Figure 21 : Configuration de factorisation relâchée appliquée aux matrices de Fourier de graphes quelconques. Chaque facteur est représenté par un carré, le nombre d'entrées non-nulle correspondant est écrit à l'intérieur.

6.3. FFT sur graphe par diagonalisation gloutonne du Laplacien

Nous avons présenté dans la section précédente une méthode permettant, à partir de la matrice de Fourier \mathbf{U} associée à un graphe quelconque, d'obtenir une transformée de Fourier rapide approchée pour les signaux définis sur ce graphe. Cependant, cette méthode souffre de deux limitations majeures :

- (L1) Elle nécessite au préalable une diagonalisation complète du Laplacien. En effet l'algorithme proposé prend en entrée la matrice de Fourier \mathbf{U} , qui n'est autre que la matrice de vecteurs propres du Laplacien $\mathbf{L} \in \mathbb{R}^{n \times n}$. Diagonaliser le Laplacien peut être très coûteux lorsque n est grand (la diagonalisation coûte $\mathcal{O}(n^3)$ opérations arithmétique en utilisant les méthodes classiques (Watkins, 2007)), ce qui rend la méthode peu pratique.
- (L2) Elle fournit des approximations efficaces $\hat{\mathbf{U}} = \mathbf{S}_J \dots \mathbf{S}_1$ qui ne sont pas orthogonales. En effet elle repose sur des projections sur des ensembles de contraintes, et il est difficile de projeter sur des ensembles de matrices à la fois creuses *et* orthogonales. Ceci mène à des transformées de Fourier approchées qui ne sont pas facilement inversibles, ou au moins pas de manière efficace. Ceci peut être problématique dans des applications où une reconstruction du signal à partir de ses coefficients de Fourier est requise.

Afin de remédier à ces deux inconvénients, nous proposons ici une autre méthode permettant d'obtenir des transformées de Fourier rapides approchées sur graphe.

6.3.1. Principe général

Il est possible de s'affranchir de la limitation (L1) en diagonalisant directement le Laplacien \mathbf{L} de manière approximative par une FA μ ST, sans effectuer au préalable une diagonalisation exacte. Ceci correspond à la résolution approchée d'un problème d'optimisation de la forme suivante :

$$\underset{\mathbf{D}, \mathbf{S}_1, \dots, \mathbf{S}_J}{\text{minimiser}} \quad \|\mathbf{L} - \mathbf{S}_1 \dots \mathbf{S}_J \mathbf{D} \mathbf{S}_J^T \dots \mathbf{S}_1^T\|_F^2 + \sum_{j=1}^J \delta_{\mathcal{S}_j}(\mathbf{S}_j) + \delta_{\mathcal{D}}(\mathbf{D}), \quad (\text{PD})$$

où les ensembles \mathcal{S}_j sont des ensembles de matrices creuses et \mathcal{D} est un ensemble de matrices diagonales.

Concernant la limitation (L2), une solution est de considérer des ensembles de contraintes \mathcal{S}_j constitués de matrices à la fois creuses et orthogonales (souvent appelées rotations élémentaires). Cependant, comme évoqué plus tôt, il est difficile

d'effectuer des projections sur ce type d'ensembles, ce qui rend difficile d'intégrer ce type de contraintes dans l'algorithme de factorisation hiérarchique proposé dans cette thèse. Afin de surmonter cette difficulté, nous proposons dans cette section une stratégie gloutonne ne nécessitant pas de projections, où la meilleure rotation élémentaire (celle qui permet de décroître le plus la fonction de coût) est choisie à chaque étape, et fixée une fois pour toutes.

6.3.2. Travaux liés

L'idée d'utiliser des rotations élémentaires afin d'obtenir une transformée orthogonale efficace n'est pas nouvelle. En effet, ces dernières années, plusieurs travaux sont allés dans ce sens. Le premier, dans le domaine des statistiques (Lee et al., 2008), introduit une méthode d'estimation de matrice de covariance en faisant l'hypothèse que la matrice des vecteurs propres de la covariance est le produit d'un nombre raisonnable de rotations élémentaires appelées rotations de Givens (Givens, 1958). Le travail suivant (Cao et al., 2011) poursuit le même objectif, mais avec une formulation légèrement différente basée sur l'estimation du maximum de vraisemblance de la matrice de covariance. Une autre méthode (Kondor et al., 2014), proche de ce qui est proposé ici, a pour but de définir des ondelettes efficaces sur graphe, via une technique appelé Factorisation Matricielle Multi-résolution (MMF), qui consiste en une factorisation gloutonne du Laplacien à l'aide de rotations élémentaires, avec des contraintes spécifiques liée à la multi-résolution. Dans un article récent (Rusu et al., 2016), il est proposé d'apprendre des dictionnaires orthogonaux efficaces, produits d'un nombre raisonnable de réflexions de Householder (Householder, 1958) (un type de matrices orthogonales élémentaires), à l'aide d'un algorithme d'optimisation alternée. Enfin, l'algorithme que nous proposons ici peut être vu comme une version tronquée de l'algorithme de valeurs propres de Jacobi (Jacobi, 1846; Golub et Van der Vorst, 2000).

6.3.3. Algorithme proposé

Ensembles de contrainte. Nous imposons que chaque facteur creux \mathbf{S}_j appartienne à l'ensemble \mathcal{R}_G des rotations de Givens (Givens, 1958) (ce qui correspond aux ensembles de contraintes $\mathcal{S}_1 = \dots = \mathcal{S}_J = \mathcal{R}_G$). En dimension n , une rotation de Givens est une transformation linéaire qui n'agit pas sur $n - 2$ coordonnées, et effectue une rotation d'un angle $\theta \in [0; 2\pi[$ par rapport aux deux coordonnées restantes (appelons les p et q). Les rotations de Givens correspondent donc à des matrices de la forme

et ainsi de suite jusqu'à la J -ième et dernière étape. La stratégie gloutonne revient donc à résoudre une séquence de J sous-problèmes très similaires, tous de la forme :

$$\underset{\mathbf{S} \in \mathcal{R}_G}{\text{minimiser}} \quad \|\mathbf{S}^T \mathbf{L}_j \mathbf{S}\|_{\text{offdiag}}^2, \quad (\text{SP})$$

avec $\mathbf{L}_j \triangleq \mathbf{S}_{j-1}^T \mathbf{L}_{j-1} \mathbf{S}_{j-1}$. Un résumé de cette stratégie gloutonne est donné en algorithme 14.

Algorithme 14 Diagonalisation gloutonne approchée du Laplacien

Entrée : Le Laplacien \mathbf{L} , un nombre de facteurs J .

- 1: $\mathbf{S}_j \leftarrow \mathbf{I} \forall j, \mathbf{L}_1 \leftarrow \mathbf{L}$
- 2: **for** $j = 1$ to J **do**
- 3: $\mathbf{S}_j \leftarrow \underset{\mathbf{S} \in \mathcal{R}_G}{\text{argmin}} \|\mathbf{S}^T \mathbf{L}_j \mathbf{S}\|_{\text{offdiag}}^2$
- 4: $\mathbf{L}_{j+1} \leftarrow \mathbf{S}_j^T \mathbf{L}_j \mathbf{S}_j$
- 5: **end for**
- 6: $\mathbf{D} \leftarrow \text{diag}(\mathbf{L}_{J+1})$

Sortie : La factorisation estimée : $\{\mathbf{S}_j\}_{j=1}^J, \mathbf{D}$.

Résolution du sous-problème. La stratégie gloutonne nécessite de résoudre J fois le problème (SP) (à la ligne 3 de l'algorithme 14). La solution de ce problème est donnée par la rotation de Givens $\mathbf{G}_{p,q,\theta}$, où les indices p et q correspondent à l'entrée de \mathbf{L}_j la plus grande en valeur absolue (notée l_{pq}^j), et $\theta = \frac{1}{2} \arctan\left(\frac{l_{qq}^j - l_{pp}^j}{2l_{pq}^j}\right) + (2k+1)\frac{\pi}{4}$, $k \in \mathbb{Z}$. On a alors $\|\mathbf{L}_{j+1}\|_{\text{offdiag}}^2 = \|\mathbf{L}_j\|_{\text{offdiag}}^2 - 2(l_{pq}^j)^2$ (la preuve est donnée en annexe C). La procédure correspondant à la résolution du sous-problème est donnée en algorithme 15. Le coût de cette procédure est dominé par la ligne 1 de cet algorithme, de complexité $\mathcal{O}(n^2)$. Cependant, il est possible de réduire ce coût. En effet, les matrices \mathbf{L}_j et \mathbf{L}_{j+1} ne diffèrent qu'aux lignes et colonnes p et q . On peut alors réutiliser les calculs effectués à l'itération précédente à l'itération courante, par une technique inspirée de (Cao et al., 2011) (un algorithme complet prenant en compte cette technique est donné en annexe C). En effectuant ceci, la complexité au pire cas reste $\mathcal{O}(n^2)$ mais est en pratique $\mathcal{O}(n)$ pour la plupart des itérations (dès que les coordonnées p et q sélectionnées à l'itération courante sont toutes deux différentes de celles sélectionnées à l'itération précédente).

6.4. Comparaison des méthodes

Les méthodes présentées lors des deux sections précédentes permettent toutes les deux d'obtenir une approximation efficace de la transformée de Fourier sur

Algorithme 15 Résolution du sous-problème (SP)**Entrée :** La matrice \mathbf{L}_j .

- 1: $(p, q) \leftarrow \operatorname{argmax}_{(r,s) \in [n]^2} (l_{rs}^j)^2$
- 2: $\theta \leftarrow \frac{1}{2} \arctan\left(\frac{l_{qq}^j - l_{pp}^j}{2l_{pq}^j}\right) + \frac{\pi}{4}$
- 3: $\mathbf{S}_j \leftarrow \mathbf{G}_{p,q,\theta}$

Sortie : La matrice $\mathbf{S}_j = \operatorname{argmin}_{\mathbf{S} \in \mathcal{R}_G} \|\mathbf{S}^T \mathbf{L}_j \mathbf{S}\|_{\text{offdiag}}^2$.

graphe. Cependant, elles diffèrent grandement, aussi bien du point de vue de leurs caractéristiques intrinsèques que de leurs propriétés empiriques. Nous mettons ceci en évidence dans cette section.

6.4.1. Évaluation des transformées de Fourier approchées

Nous profitons de cette sous-section pour évaluer les transformées de Fourier rapides approchées obtenues à l'aide des deux méthodes, par le biais de deux mesures de performances, en considérant divers graphes de diverses tailles.

Graphes considérés. Nous considérons pour cette expérience plusieurs familles de graphes aléatoires parmi les plus répandues. Nous générons tous les graphes à l'aide de la librairie “Graph Signal Processing” (GSPBOX) (Perraudin et al., 2014). Les graphes utilisés sont :

- **Erdős-Rényi** : un graphe totalement aléatoire où chaque paire de nœuds est connectée avec probabilité $p = 0.1$.
- **Community** : un graphe constitué de $\sqrt{n}/2$ communautés de tailles aléatoires. Chaque communauté est elle-même un graphe dont les nœuds correspondent à des points placés de manière aléatoire sur le disque unité et connectés avec un poids inversement proportionnel à leur distance, si celle-ci est inférieure à un certain seuil. De plus, des arêtes aléatoires entre n'importe quelle paire de nœuds apparaissent avec une probabilité de $1/n$.
- **Sensor** : un graphe dont les nœuds correspondent à des points placés aléatoirement dans un carré de côté 1, et connectés avec un poids inversement proportionnel à leur distance, si celle-ci est inférieure à un certain seuil.
- **Path** : le graphe chemin régulier, où chaque nœud est connecté aux deux nœuds adjacents par des arêtes de poids égaux.

Nous prenons pour toutes ces familles des graphes de taille n variable, avec $n \in \{128, 256, 512, 1024\}$ nœuds. Des exemples de visualisations de graphes considérés pour cette expérience sont montrés en figure 22.

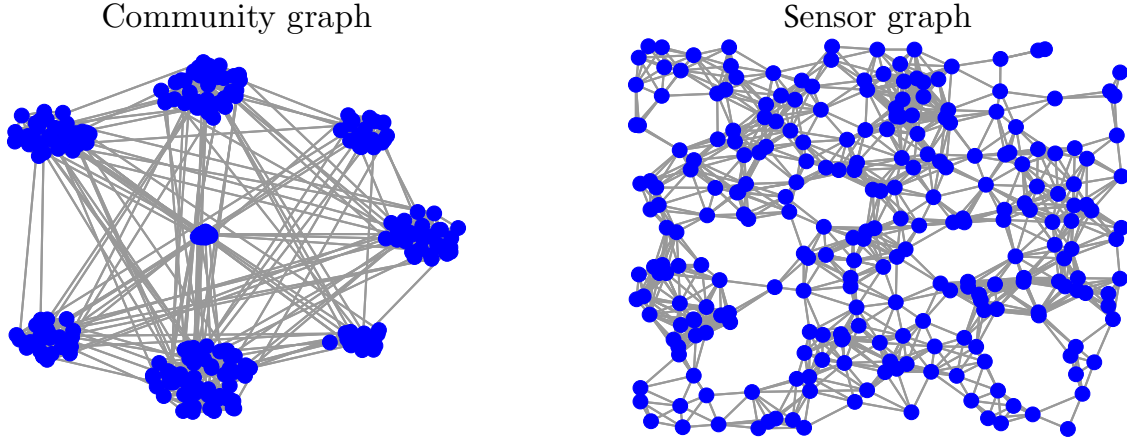


Figure 22 : Exemples d’un graphe de communautés (“Community”) et d’un graphe de réseau de capteurs (“Sensor”), de taille $n = 256$.

Calcul des transformées de Fourier approchées. On considère pour chaque configuration deux transformées de Fourier rapides approchées. La première est une approximation efficace $\hat{\mathbf{U}}_{\text{PALM}}$ calculée à l’aide de la méthode de factorisation de la matrice de Fourier, avec les paramètres de factorisations donnés par $C_1 = 3$, $C_2 = 1.66$ et $C_3 = 1.4$. La seconde est une approximation efficace $\hat{\mathbf{U}}_{\text{Givens}}$ calculée à l’aide de la méthode de diagonalisation gloutonne du Laplacien, pour laquelle le nombre de rotations de Givens J est pris tel que le gain relatif en complexité RCG soit le même pour les deux méthodes (dans le cas de la méthode de diagonalisation gloutonne, on a $\text{RCG} = \frac{\|\mathbf{U}\|_0}{4J}$).

Mesures de performance. Comment mesurer la pertinence d’une transformée de Fourier rapide approchée sur graphe $\hat{\mathbf{U}}$? On utilise pour cela deux mesures. Premièrement, on définit l’erreur relative sur le calcul de la transformée de Fourier :

$$\text{err}_c(\hat{\mathbf{U}}) \triangleq \frac{\|\mathbf{U} - \hat{\mathbf{U}}\|_F}{\|\mathbf{U}\|_F},$$

(en effectuant une permutation des colonnes de $\hat{\mathbf{U}}$ au préalable dans le cas de la méthode de diagonalisation gloutonne, de manière ordonner les modes de Fourier approchés). Deuxièmement, on définit l’erreur relative de diagonalisation du Laplacien :

$$\text{err}_d(\hat{\mathbf{U}}) \triangleq \frac{\|\hat{\mathbf{U}}^T \mathbf{L} \hat{\mathbf{U}}\|_{\text{offdiag}}}{\|\mathbf{L}\|_F}$$

Résultats. Les résultats de cette expérience (en moyenne sur 10 réalisations) sont donnés par le tableau 1. Ils suscitent plusieurs commentaires :

- Premièrement, il est très clair que la méthode de factorisation de la matrice de Fourier fournit des FFT approchées meilleures en terme d’erreur relative de calcul de la transformée de Fourier, par rapport à la méthode de diagonalisation gloutonne du Laplacien. Inversement, la méthode de diagonalisation gloutonne du Laplacien fournit des FFT approchées meilleures en terme d’erreur relative de diagonalisation. Ceci est vrai pour toutes les configurations testées. Cela paraît d’ailleurs assez logique car les mesures de performance utilisées sont très proches des fonctions de coût considérées par les deux méthodes.
- Deuxièmement, concernant les différents graphes considérés, on constate que les deux méthodes montrent de meilleures performances pour les graphes “Sensor” et “Community”, et des performances un peu dégradées pour les graphes “Erdős-Rényi” et “Path”. On s’attendait effectivement à ce qu’il soit difficile d’obtenir une transformée de Fourier rapide approchée pour le graphe “Erdős-Rényi”, celui-ci étant totalement aléatoire et non-structuré, mais cela paraît plus étonnant pour le graphe “Path”, dont on sait que la transformée de Fourier admet une factorisation creuse multi-couche.
- Troisièmement, on remarque que pour la méthode de diagonalisation gloutonne du Laplacien, l’erreur relative n’augmente pas lorsque la taille des graphes augmente, alors que le gain relatif en complexité augmente lui avec la taille des graphes (on considère ici des transformées de Fourier rapides approchées avec une complexité en $\mathcal{O}(n^{1.27})$). Cela signifie que plus la taille des graphes augmente, plus il est profitable d’utiliser les transformées de Fourier rapides approchées, à erreur constante. C’est typiquement le comportement attendu pour une transformée rapide. Le même comportement est observable pour la méthode de factorisation de la matrice de Fourier, pour des transformées rapides approchées de complexité plus élevée.

En résumé, on peut retenir de cette expérience que les transformées de Fourier rapides approchées calculées à l’aide de la méthode de factorisation de la véritable matrice de Fourier sont plus adaptées aux cas où la connaissance des coefficients de la transformée de Fourier d’un signal \mathbf{f} est importante ($\mathbf{U}^T \mathbf{f}$ et $\hat{\mathbf{U}}^T \mathbf{f}$ sont proches). Alternativement, les transformées de Fourier rapides approchées calculées à l’aide de la méthode de diagonalisation gloutonne du Laplacien sont plus adaptées au cas où l’interprétation physique de la transformée de Fourier d’un signal \mathbf{f} est importante ($\hat{\mathbf{U}}^T \mathbf{f}$ est proche d’une projection de \mathbf{f} sur les espaces propres du Laplacien \mathbf{L}). L’une ou l’autre des méthodes peut donc être la plus adaptée en fonction de l’application considérée. De plus, les méthodes présentées ici semblent montrer de

meilleures performances pour les graphes les plus réalistes (“Sensor” et “Community”), susceptibles de ressembler aux graphes rencontrés dans des applications réelles.

		Erdős-Rényi	Community	Sensor	Path
$n = 128, \text{RCG} = 2.7$	$\text{err}_c(\hat{\mathbf{U}})$	0.31 1.19	0.05 0.79	0.12 0.76	0.37 0.99
	$\text{err}_d(\hat{\mathbf{U}})$	0.37 0.09	0.07 0.03	0.16 0.04	0.39 0.06
$n = 256, \text{RCG} = 4.4$	$\text{err}_c(\hat{\mathbf{U}})$	0.47 1.29	0.10 0.93	0.18 0.95	0.56 1.11
	$\text{err}_d(\hat{\mathbf{U}})$	0.46 0.08	0.12 0.03	0.23 0.04	0.47 0.05
$n = 512, \text{RCG} = 7.3$	$\text{err}_c(\hat{\mathbf{U}})$	0.60 1.34	0.19 1.07	0.22 1.07	0.68 1.19
	$\text{err}_d(\hat{\mathbf{U}})$	0.45 0.07	0.21 0.03	0.26 0.04	0.42 0.05
$n = 1024, \text{RCG} = 12.2$	$\text{err}_c(\hat{\mathbf{U}})$	0.70 1.37	0.30 1.19	0.24 1.18	0.76 1.25
	$\text{err}_d(\hat{\mathbf{U}})$	0.40 0.05	0.28 0.04	0.27 0.04	0.38 0.04

Tableau 1.: Résultats de l’évaluation des transformées de Fourier rapides approchées. Les résultats sont donnés pour plusieurs familles de graphes aléatoires classiques (la moyenne sur 10 réalisations est considérée), plusieurs tailles de graphes, et les deux mesures de performances introduites ci-dessus. Pour chaque configuration, le résultat de la FFT approchée calculée à l’aide de la méthode de factorisation de la matrice de Fourier est donné à gauche de la barre verticale, et le résultat de la FFT approchée calculée à l’aide de la méthode de diagonalisation gloutonne du Laplacien est donné à droite de la barre verticale (la gain en complexité RCG étant le même pour les deux méthodes), le meilleur résultat étant donné en caractères gras.

6.4.2. Exemple d’application : filtrage

Nous nous intéressons ici au filtrage de signaux sur graphe, application assez immédiate pour laquelle les transformées de Fourier rapides approchées que nous proposons sont particulièrement adaptées. En effet pour un signal sur graphe \mathbf{f} , on obtient sa version filtrée \mathbf{g} par le calcul suivant :

$$\mathbf{g} = \mathbf{U}\mathbf{H}\mathbf{U}^T\mathbf{f},$$

où \mathbf{H} est une matrice diagonale représentant la réponse en fréquence du filtre. On voit alors que l’opération de filtrage nécessite l’application de la transformée de Fourier et de son inverse. Une approximation efficace de celle-ci permettrait donc d’obtenir un gain en complexité sur l’opération de filtrage de l’ordre du gain en complexité relative RCG de la FA μ ST correspondante.

On effectue ici une expérience de filtrage sur graphe à l'aide de transformées de Fourier rapides approchées, afin d'évaluer le potentiel des méthodes présentées dans les deux sections précédentes dans ce contexte. Pour cela, on considère un graphe classique représentant le réseau routier du Minnesota (appelé Minnesota road graph), constitué de $n = 2642$ nœuds. On a alors un Laplacien $\mathbf{L} \in \mathbb{R}^{2642 \times 2642}$, avec $\mathbf{L} = \mathbf{U}\mathbf{\Lambda}\mathbf{U}^T$, où \mathbf{U} représente la transformée de Fourier sur le graphe et la matrice diagonale $\mathbf{\Lambda}$ contient les valeurs propres du Laplacien (rangées par ordre croissant).

Modèle de signal. Pour cette expérience, un signal basse fréquence \mathbf{f} est généré aléatoirement dans le domaine de Fourier du graphe : les composantes de son spectre \mathbf{s} sont indépendantes et suivent une loi normale, $s_i \sim \mathcal{N}(0, \theta_i)$ avec $\theta_i = \exp(-\lambda_i)$ où λ_i est la i -ème valeur propre du Laplacien. On obtient alors un signal dans le domaine des nœuds du graphe $\mathbf{f} = \mathbf{U}\mathbf{s}$, qui est ensuite corrompu par un bruit blanc Gaussien \mathbf{n} avec $n_i \sim \mathcal{N}(0, \sigma)$, pour obtenir le signal bruité $\tilde{\mathbf{f}} = \mathbf{f} + \mathbf{n}$.

Filtrage. Le signal bruité est filtré passe-bas afin de le débruiter. Le filtre utilisé a une réponse fréquentielle donnée par $h_i = 1/(1 + \gamma\lambda_i)$ avec $\gamma = 3$, où h_i est la i -ème entrée sur la diagonale de la matrice \mathbf{H} de filtrage. Le filtrage est effectué à l'aide de la véritable matrice de Fourier \mathbf{U} (on a alors $\mathbf{f} = \mathbf{U}\mathbf{H}\mathbf{U}^T\tilde{\mathbf{f}}$) ou d'une approximation efficace de celle-ci. On considère une approximation efficace $\hat{\mathbf{U}}_{\text{PALM}}$ calculée à l'aide de la méthode de factorisation de la matrice de Fourier, avec les paramètres de factorisations donnés par $C_1 = 3$, $C_2 = 1.4$ et $C_3 = 1.5$ (on a alors $\hat{\mathbf{f}} = \hat{\mathbf{U}}_{\text{PALM}}\mathbf{H}\hat{\mathbf{U}}_{\text{PALM}}^T\tilde{\mathbf{f}}$). Ceci résulte en une approximation efficace pour laquelle le gain en complexité relative RCG est proche de 8. On considère aussi trois approximations efficaces calculées à l'aide de la méthode de diagonalisation gloutonne du Laplacien avec $J \in \{50000, 100000, 250000\}$, on les note $\hat{\mathbf{U}}_{50}$, $\hat{\mathbf{U}}_{100}$ et $\hat{\mathbf{U}}_{250}$. Ces approximations ont des gains en complexité relative égaux respectivement à 35, 17 et 7.

Résultats. Les résultats en terme de rapport signal à bruit (SNR) en décibels :

$$\text{SNR} = 10 \log_{10} \left(\frac{\|\hat{\mathbf{f}}\|_2^2}{\|\mathbf{f} - \hat{\mathbf{f}}\|_2^2} \right),$$

sont donnés dans le tableau 2 pour différents niveaux de bruit σ , et en moyenne pour 100 réalisations. On peut y voir que les FA μ ST obtenues à l'aide de la stratégie de diagonalisation gloutonne ($\hat{\mathbf{U}}_{50}$, $\hat{\mathbf{U}}_{100}$ et $\hat{\mathbf{U}}_{250}$) mènent à de meilleures performances que celles obtenues à l'aide de la stratégie de factorisation ($\hat{\mathbf{U}}_{\text{PALM}}$), et ceci même pour un RCG bien plus grand. Par exemple $\hat{\mathbf{U}}_{50}$ a un RCG de 35, elle est donc plus de quatre fois plus efficace que $\hat{\mathbf{U}}_{\text{PALM}}$ (RCG de 8), et pourtant, elle permet un meilleur débruitage (un exemple est donné en figure 23). Ceci s'explique certainement par le fait que les FA μ ST obtenues par la stratégie gloutonne sont des matrices orthogonales, et permettent donc une parfaite reconstruction des signaux

à partir de leur transformée de Fourier approchée, ce qui n'est pas le cas des FA μ ST obtenues avec la stratégie de factorisation. Cette propriété semble cruciale pour les applications où des signaux doivent être reconstruits, telles que le filtrage. De plus, les résultats de débruitage obtenus avec les FA μ ST sont proches de ceux obtenus en utilisant la véritable matrice de Fourier, ce qui montre le potentiel des FA μ ST dans ce contexte de filtrage sur graphe.

	RCG	$\sigma = 0.2$	$\sigma = 0.25$	$\sigma = 0.3$	$\sigma = 0.4$	$\sigma = 0.5$	$\sigma = 0.6$
Signal bruité	\mathbf{X}	5.29	3.43	1.82	-0.68	-2.65	-4.25
Filtré avec \mathbf{U}	1	5.58	5.40	5.17	4.55	3.95	3.24
Filtré avec $\hat{\mathbf{U}}_{250}$	7	5.55	5.36	5.15	4.53	3.94	3.24
Filtré avec $\hat{\mathbf{U}}_{100}$	17	5.47	5.29	5.09	4.49	3.92	3.24
Filtré avec $\hat{\mathbf{U}}_{50}$	35	5.33	5.16	4.97	4.41	3.87	3.22
Filtré avec $\hat{\mathbf{U}}_{\text{PALM}}$	8	5.12	4.98	4.70	4.23	3.59	2.98

Tableau 2 .: Résultats de filtrage, le SNR en décibels moyen sur 100 signaux générés aléatoirement est donné pour différents niveaux de bruit. Le SNR moyen des signaux bruités est donné à la première ligne. Le SNR moyen des signaux filtrés est donné aux lignes suivantes (chaque ligne correspond à une transformée de Fourier différente).

Discussion. Comme précisé au début de cette sous-section, le filtrage constitue une des applications les plus immédiates, et nous permet de fournir une preuve de concept intéressante. Cependant, ce n'est pas une application pour laquelle les transformées de Fourier rapides approchées apportent une solution à un problème non-consideré auparavant. En effet, il existe déjà un autre type de méthodes permettant d'effectuer du filtrage rapide sur graphe, sans nécessiter de transformée de Fourier rapide. Celles-ci consistent à approcher la matrice de filtrage $\mathbf{U}\mathbf{H}\mathbf{U}^T$ par un polynôme de degré r du Laplacien : $\mathbf{U}\mathbf{H}\mathbf{U}^T \approx \sum_{i=1}^r c_i \mathbf{L}^i$ (Hammond et al., 2011). L'opération de filtrage est alors de complexité $\mathcal{O}(r \|\mathbf{L}\|_0)$. Ces méthodes sont plus rapides à mettre en œuvre que celles présentées dans cette thèse elles ne nécessitent pas de factorisation ni de diagonalisation, et le calcul des coefficients c_i du polynôme est rapide. Cependant, elles sont moins flexibles car un seul filtre est approché par jeu de coefficients calculé, alors qu'une approximation de la transformée de Fourier permet d'approcher rapidement n'importe quel filtre. De plus, elles ne s'appliquent bien qu'aux filtres dont la réponse fréquentielle est bien approchée par un polynôme de degré faible, alors que les méthodes développées ici sont indépendantes du type de filtre considéré. Une comparaison quantitative entre les deux types de méthodes serait très intéressante à effectuer.

Conclusion. Nous nous sommes intéressés dans ce chapitre au problème de la transformée de Fourier rapide sur graphe. Après avoir introduit brièvement le traitement

du signal sur graphe, nous avons proposé deux méthodes permettant d'obtenir des transformées de Fourier rapides approchées sur graphe. Nous avons ensuite comparé ces méthodes de manière expérimentale, puis les avons appliquées au filtrage de signaux sur graphes.

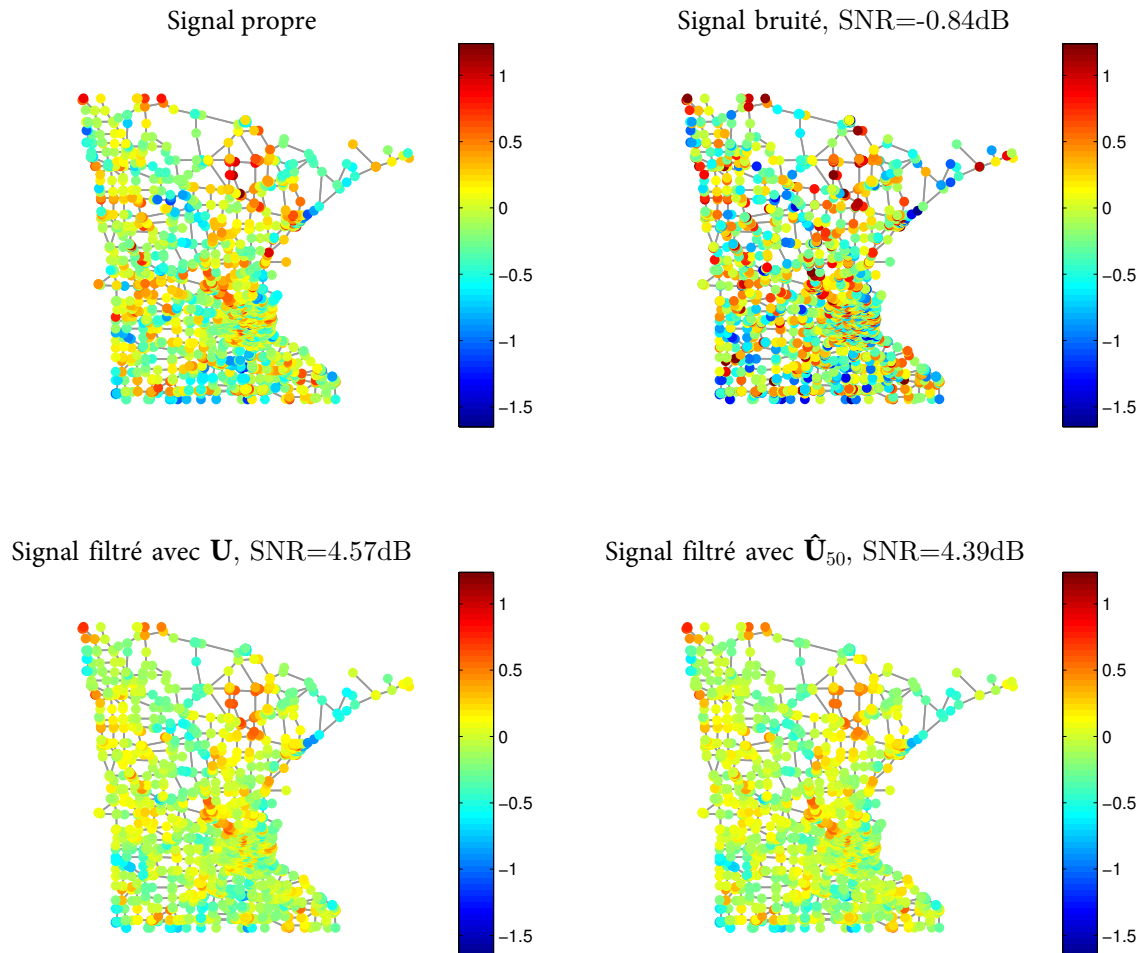


Figure 23 : Exemple de filtrage sur le graphe Minnesota, la transformée de Fourier rapide approchée $\hat{\mathbf{U}}_{50}$ est approximativement 35 fois plus efficace que la véritable matrice de Fourier \mathbf{U} .

Identifiabilité de la forme efficiente d'une matrice

Nous nous intéressons dans ce chapitre à l'identifiabilité de facteurs creux à partir de l'observation de leur produit. Nous commençons par introduire le problème considéré et le relier à l'état de l'art. Nous donnons ensuite des conditions nécessaire et suffisante pour que les facteurs soient identifiables. Nous terminons par discuter l'identifiabilité des factorisations correspondant aux algorithmes rapides associés aux transformées usuelles du traitement du signal via la stratégie hiérarchique de factorisation proposée dans cette thèse.

Sommaire

7.1	Introduction	103
7.2	Énoncé du problème étudié	105
7.3	Résultat d'identifiabilité	107
7.3.1	Identifiabilité en complétion de matrice de rang 1	107
7.3.2	Condition suffisante d'identifiabilité	109
7.3.3	Condition nécessaire d'identifiabilité	109
7.4	Exemple : la FFT à partir de la matrice de Fourier	111

7.1. Introduction

Est-il possible de retrouver des facteurs creux $\mathbf{S}_1, \dots, \mathbf{S}_J$ à partir de l'observation de leur produit $\mathbf{A} = \mathbf{S}_J \dots \mathbf{S}_1 \triangleq f(\mathbf{S}_1, \dots, \mathbf{S}_J)$? Si tel est le cas, on dit que la factorisation est *identifiable*. Dans un contexte de factorisation creuse multi-couche, l'identifiabilité correspond donc à l'existence d'une solution unique à un système de mn équations polynomiales (une par entrée observée de la matrice \mathbf{A}) de la forme

$$a_{ij} = \sum_{k_1, \dots, k_{J-1}} (\mathbf{S}_J)_{ik_{J-1}} (\mathbf{S}_{J-1})_{k_{J-1}k_{J-2}} \dots (\mathbf{S}_{J-1})_{k_1j}, \quad (7.1)$$

où $(\mathbf{S}_j)_{kl}$ désigne l'entrée de la matrice \mathbf{S}_j située à la k -ième ligne et à la l -ième colonne (nous utilisons cette notation au lieu de s_{jkl} qui correspond à la notation standard utilisée dans le reste du document pour désigner les entrées d'une matrice mais est moins lisible). L'étude des systèmes d'équations polynomiales s'avère fastidieuse dans le cas général (Sturmfels, 2002). Cependant il est possible d'exploiter la forme spécifique des polynômes impliqués dans le système pour obtenir des garanties d'unicité de la solution (à certaines ambiguïtés inhérentes au problème près, que nous précisons à la section 7.2). Nous présentons dans les paragraphes suivants deux approches allant dans ce sens, tout d'abord une approche classique reposant sur un procédé d'élévation (lifting), puis l'approche proposée dans cette thèse basée sur des techniques de complétion de matrice.

Approche d'élévation. On peut remarquer ici que les équations (7.1) ont une forme particulière, chaque monôme étant le produit d'exactly une entrée de chaque facteur creux. Ceci est dû au fait que la fonction f liant les facteurs et l'observation est multilinéaire (linéaire par rapport à chaque facteur, les autres étant fixés). Le problème d'identifier les facteurs $\mathbf{S}_1, \dots, \mathbf{S}_J$ à partir de l'observation de leur produit \mathbf{A} peut donc être vu comme un *problème inverse multilinéaire*. Il est possible de reformuler ce type de problème comme la récupération d'un tenseur de rang 1 $\mathbf{T} \triangleq \text{vec}(\mathbf{S}_1) \otimes \dots \otimes \text{vec}(\mathbf{S}_J)$ à partir d'observations linéaires de ce tenseur $\mathbf{A} = \mathcal{L}(\mathbf{W})$ où \mathcal{L} est un opérateur linéaire fixe et connu. On appelle cette reformulation une *élévation* ("lifting" en anglais) du problème. Un tel procédé a été utilisé afin de donner des conditions d'identifiabilité pour le problème qui nous intéresse lorsque le support des facteurs creux est connu (Malgouyres et Landsberg, 2016). Des techniques similaires ont aussi été utilisées pour donner des conditions d'identifiabilité dans le contexte des *problèmes inverses bilinéaires* (ce qui correspond à $J = 2$ facteurs), dans des domaines tels que la déconvolution aveugle (Choudhary et Mitra, 2014; Li et al., 2015) ou la reconstruction de phase (Candès et al., 2013). Ces techniques se concentrent sur les propriétés de l'opérateur linéaire d'observation \mathcal{L} (plus précisément de son noyau) pour obtenir des garanties d'identifiabilité. Les garanties obtenues sont alors abstraites et difficiles à vérifier concrètement.

Approche proposée. Nous étudions ici l'identifiabilité de facteurs creux à partir de l'observation de leur produit via la méthode de factorisation hiérarchique proposée dans cette thèse. Ceci revient à étudier théoriquement uniquement des factorisations à 2 facteurs. Les facteurs sont alors identifiés un par un, et ceci peut être vu comme la transformation d'un problème inverse multilinéaire en une séquence de problèmes inverse bilinéaires. Notre étude repose sur la généralisation de techniques utilisées en complétion de matrice (Király et Tomioka, 2012; Cosse et Demanet, 2015). L'intérêt principal de l'approche proposée ici par rapport à l'approche d'élévation (Malgouyres et Landsberg, 2016) est qu'elle mène à des conditions concrètement vérifiables, comme illustré dans le cas de la transformée de Fourier de dimension $n = 2^k$ en fin de chapitre.

7.2. Énoncé du problème étudié

La stratégie de factorisation hiérarchique introduite dans cette thèse a pour but de factoriser une matrice d'intérêt $\mathbf{A} \in \mathbb{C}^{m \times n}$ en J facteurs creux $\mathbf{S}_1, \dots, \mathbf{S}_J$ en effectuant $J-1$ factorisations en seulement deux facteurs, où le niveau de parcimonie des facteurs est donné en paramètre. Nous considérons dans ce chapitre théorique des matrices complexes pour plus de généralité. Dans un régime de fonctionnement idéal, et en supposant qu'il existe effectivement une factorisation exacte de la matrice d'intérêt en matrices creuses $\mathbf{A} = \mathbf{S}_J \dots \mathbf{S}_1$, la factorisation hiérarchique permettrait à la première factorisation en deux de retrouver le premier facteur creux \mathbf{S}_1 , via la factorisation $\mathbf{A} = \mathbf{R}_1 \mathbf{S}_1$ avec $\mathbf{R}_1 = \mathbf{S}_J \dots \mathbf{S}_2$. De manière plus générale la ℓ -ième factorisation en deux permettrait de retrouver le ℓ -ième facteur creux via la factorisation $\mathbf{R}_{\ell-1} = \mathbf{R}_\ell \mathbf{S}_\ell$ avec $\mathbf{R}_{\ell-1} = \mathbf{S}_J \dots \mathbf{S}_\ell$ et $\mathbf{R}_\ell = \mathbf{S}_J \dots \mathbf{S}_{\ell+1}$.

Il est alors clair que l'identifiabilité de la factorisation globale en J facteur découle de l'identifiabilité des factorisations en deux facteurs. Nous étudions donc ici le problème général d'identifiabilité de deux matrices creuses $\mathbf{X} \in \mathbb{C}^{m \times r}$ et $\mathbf{Y} \in \mathbb{C}^{r \times n}$, à partir de l'observation de leur produit $\mathbf{Z} \triangleq \mathbf{X}\mathbf{Y} \in \mathbb{C}^{m \times n}$.

Pour cela, l'angle d'attaque que nous choisissons est assez différent de celui employé par (Malgouyres et Landsberg, 2016). En effet, nous commençons par exprimer l'observation de la façon suivante : $\mathbf{Z} = \sum_{i=1}^r \mathbf{x}_i \cdot (\mathbf{y}^i)^T$, où \mathbf{x}_i est la i -ième colonne de \mathbf{X} et \mathbf{y}^i est la i -ième ligne de \mathbf{Y} . En introduisant la notation $\mathbf{C}_i \triangleq \mathbf{x}_i \cdot (\mathbf{y}^i)^T$, pour désigner la matrice de rang 1 constitué par le produit externe de la i -ième colonne de \mathbf{X} avec la i -ième ligne de \mathbf{Y} , on a ainsi

$$\mathbf{Z} = \sum_{i=1}^r \mathbf{C}_i,$$

où l'observation est explicitement décrite comme une somme de contributions de rang 1. Le lemme suivant indique que l'identifiabilité de ces contributions \mathbf{C}_i , $i \in \{1, \dots, r\}$ est équivalente à l'identifiabilité des facteurs \mathbf{X} et \mathbf{Y} , à quelques ambiguïtés près.

Ambiguïtés. La numérotation des contributions de rang 1 est arbitraire, on a donc ce qu'on appelle une ambiguïté de permutation (on peut permuer deux colonnes de \mathbf{X} et effectuer la même permutation sur les lignes de \mathbf{Y} sans changer le produit $\mathbf{X}\mathbf{Y}$). On peut aussi multiplier n'importe quelle colonne de \mathbf{X} par un coefficient $\alpha \in \mathbb{C}^*$, et multiplier la ligne de \mathbf{Y} correspondante par $1/\alpha$ sans changer le produit $\mathbf{X}\mathbf{Y}$, on appelle cela l'ambiguïté de phase et d'échelle (tout nombre complexe non-nul α étant le produit d'une phase $e^{i\theta}$ et d'une échelle strictement positive $|\alpha|$). On ne peut pas éviter l'ambiguïté de permutation. Par contre, on peut supprimer les ambiguïtés de phase et d'échelle simplement en imposant que la première entrée non-nulle de chaque colonne de \mathbf{X} soit égale à 1.

Lemme 1. *Deux vecteurs \mathbf{a} and \mathbf{b} non nuls sont uniquement déterminés par leur produit externe $\mathbf{R} \triangleq \mathbf{a}\mathbf{b}^T$ (aux ambiguïtés classiques près).*

Démonstration. Supposons que \mathbf{R} puisse être factorisé de deux manière différentes :

$$\mathbf{R} = \mathbf{a}\mathbf{b}^T = \mathbf{e}\mathbf{f}^T.$$

Premièrement, il est aisé de remarquer que les supports de \mathbf{a} et \mathbf{e} sont identiques, de même que ceux de \mathbf{b} et \mathbf{f} . En effet, le support de la matrice \mathbf{R} n'est autre que le produit cartésien des supports des deux vecteurs dont elle est le produit externe. Désignons par p l'indice correspondant à la première entrée non-nulle de \mathbf{a} et \mathbf{e} . Si l'on impose

$$a_p = e_p = 1 \text{ (pour supprimer les ambiguïtés de phase et d'échelle),}$$

alors on a nécessairement

$$\mathbf{a}(\mathbf{b})^T = \mathbf{e}(\mathbf{f})^T \Rightarrow \mathbf{b} = \mathbf{f} \Rightarrow a_i = e_i, \forall i,$$

ce qui prouve le lemme. □

Nous nous concentrons donc sur l'identifiabilité des contributions de rang 1 \mathbf{C}_i , $i \in \{1, \dots, r\}$ dans la suite de ce chapitre.

Afin de simplifier le problème, nous considérons que les supports des contributions de rang 1 (donc les supports des facteurs \mathbf{X} et \mathbf{Y}) sont connus. Plus précisément, on désigne les supports par $\mathcal{C}_i \triangleq \text{supp}(\mathbf{C}_i) = \text{supp}(\mathbf{x}_i) \times \text{supp}(\mathbf{y}_i)$, et on considère que l'on connaît \mathcal{C}_i , pour tout $i \in \{1, \dots, r\}$ (cette hypothèse est aussi faite dans (Malgouyres et Landsberg, 2016)). Faire cette hypothèse de connaissance des supports revient à décomposer le problème d'identifiabilité des facteurs creux en deux étapes : premièrement l'identifiabilité des supports et ensuite l'identifiabilité des entrée non-nulles. Nous considérons dans la dernière section de ce chapitre l'identifiabilité de la transformée de Fourier rapide à partir de l'observation de la matrice de Fourier, et montrons que dans ce cas, les supports des contributions de rang 1 sont identifiables à partir de la connaissance de leur degré de parcimonie, ce qui justifie cette approche.

Contributions. Dans ce chapitre, nous donnons une condition suffisante pour l'identifiabilité de deux facteurs creux à partir de l'observation de leur produit, en supposant que leurs supports sont connus. Cette condition suffisante a l'avantage d'être concrète et simple à vérifier. Nous donnons ensuite une condition nécessaire d'identifiabilité, toujours à supports connus. Nous considérons enfin un cas précis, celui de la transformée de Fourier discrète, où les supports des facteurs creux correspondant à la transformée de Fourier rapide sont identifiables si l'on a connaissance du degré de parcimonie des facteurs, et ces supports remplissent la condition suffisante d'identifiabilité, menant à une factorisation globale théoriquement identifiable via la stratégie hiérarchique présentée dans cette thèse.

7.3. Résultat d'identifiabilité

Nous donnons dans cette section une condition suffisante sur les supports \mathcal{C}_i , $i \in \{1, \dots, r\}$ pour que les contributions de rang 1 associées soient identifiables. Nous nous appuyons pour cela sur un résultat d'identifiabilité en complétion de matrice de rang 1 que nous présentons dans la sous-section suivante.

7.3.1. Identifiabilité en complétion de matrice de rang 1

Considérons une matrice $\mathbf{R} \in \mathbb{C}^{m \times n}$ de rang 1 (on a donc $\mathbf{R} = \mathbf{a}\mathbf{b}^T$) dont on n'observe que les entrées indexées par l'ensemble \mathcal{R} (c'est-à-dire que $(i, j) \in \mathcal{R}$ si et seulement si l'entrée r_{ij} est observée). Le but de la complétion de matrice est de retrouver la matrice \mathbf{R} toute entière à partir des entrées observées uniquement. Ceci est équivalent, d'après le lemme 1 à retrouver les vecteurs \mathbf{a} et \mathbf{b} (en fixant $a_p = 1$ pour éviter les ambiguïtés classiques).

Par exemple si l'on a

$$\mathbf{R} = \begin{pmatrix} 3 & 1 & 2 \\ 9 & 3 & 6 \\ 6 & 2 & 4 \end{pmatrix} = \begin{pmatrix} 1 \\ 3 \\ 2 \end{pmatrix} \cdot (3 \ 1 \ 2)$$

et on observe

$$\begin{pmatrix} 3 & * & 2 \\ * & 3 & * \\ 6 & * & 4 \end{pmatrix} \quad (\text{S1})$$

(où le symbole $*$ désigne une entrée non-observée), la complétion revient à résoudre le système d'équations polynomiales suivant (sachant que $a_1 = 1$)

$$\begin{aligned} a_1 b_1 &= 3 \\ a_1 b_3 &= 2 \\ a_2 b_2 &= 3, \\ a_3 b_1 &= 6 \\ a_3 b_3 &= 4 \end{aligned}$$

et elle est identifiable si la solution de ce système est unique. Ici, il est aisé de voir que ce n'est pas le cas, car n'importe quelles valeurs de a_2 et b_2 sont possibles tant que $a_2 b_2 = 3$. Par contre si l'observation prend la forme

$$\begin{pmatrix} 3 & 1 & 2 \\ * & 3 & * \\ 6 & * & * \end{pmatrix} \quad (\text{S2})$$

la complétion correspond au système suivant

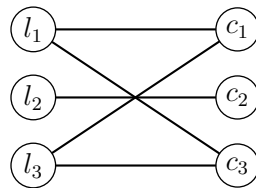
$$\begin{aligned} a_1 b_1 &= 3 \\ a_1 b_2 &= 1 \\ a_1 b_3 &= 2, \\ a_2 b_2 &= 3 \\ a_3 b_1 &= 6 \end{aligned}$$

qui quant à lui mène à une solution unique immédiatement.

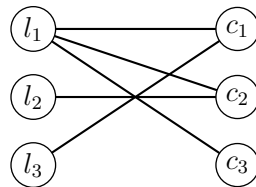
Quelle est la différence entre les deux cas ? Assez simplement, dans un cas il est possible de procéder par élimination pour résoudre le système et pas dans l'autre. Le succès ou non de la procédure d'élimination peut être déterminé a priori à partir de l'ensemble \mathcal{R} des entrées observées. Pour cela, nous introduisons le graphe bipartite $\mathcal{G}(\{l_1, \dots, l_m\}, \{c_1, \dots, c_n\}, \mathcal{R})$ associé au problème de complétion, dont les nœuds correspondent aux lignes et aux colonnes de \mathbf{R} (donc aux inconnues \mathbf{a} et \mathbf{b}), et les arêtes correspondent aux entrées observées. On a alors la proposition suivante (voir (Király et Tomioka, 2012) pour plus de détails).

Proposition 2. *La complétion d'une matrice \mathbf{R} de rang 1 avec l'ensemble d'entrées observées \mathcal{R} est identifiable si et seulement si le graphe \mathcal{G} associé à \mathcal{R} est connecté.*

Cette proposition indique s'il est possible de procéder par élimination afin de résoudre de manière unique le système d'équations polynomiales associé au problème de complétion. En effet, ceci est possible uniquement lorsque toutes les inconnues (les nœuds du graphes) sont connectées par les équations (les arêtes du graphes, mettant en jeu deux variables) correspondant aux entrées observées. Par exemple, la situation (S1) correspond au graphe suivant qui n'est pas connecté :



et la situation (S2) correspond au graphe suivant qui lui est connecté :



7.3.2. Condition suffisante d'identifiabilité

Utilisons maintenant le résultat de la proposition 2 dans le cas qui nous intéresse, c'est-à-dire lorsque le but est d'identifier r contributions de rang 1 à partir de l'observation de leur somme $\mathbf{Z} = \sum_{i=1}^r \mathbf{C}_i$, leur supports \mathcal{C}_i , $i = 1, \dots, r$ étant connus. Pour cela, introduisons le *support observé* correspondant à chaque contribution $\underline{\mathcal{C}}_i \triangleq \mathcal{C}_i \cap (\overline{\cup_{j \neq i} \mathcal{C}_j})$, où $\overline{\mathcal{C}}$ désigne le complément de l'ensemble d'indices \mathcal{C} . Le support observé $\underline{\mathcal{C}}_i$ rassemble donc les indices des entrées où l'observation \mathbf{Z} coïncide avec la i -ème contribution de rang 1. On associe aussi à chaque contribution le graphe bipartite $\mathcal{G}_i(\text{supp}(\mathbf{x}_i), \text{supp}(\mathbf{y}_i), \underline{\mathcal{C}}_i)$. Ceci nous permet d'énoncer la condition suffisante d'identifiabilité suivante.

Lemme 2. *Les contributions de rang 1 \mathbf{C}_i , $i \in \{1, \dots, r\}$ sont toutes identifiables si les graphes \mathcal{G}_i associés à leurs supports observés $\underline{\mathcal{C}}_i$ sont tous connectés.*

Démonstration. Ce lemme découle directement de l'application de la proposition 2 à chaque contribution de rang 1 de manière indépendante. \square

Un cas particulier d'application de ce lemme, qui nous intéressera dans la prochaine section, est énoncé par le corollaire suivant.

Corollaire 1. *Les contributions de rang 1 sont identifiables si leurs supports sont disjoints : $\mathcal{C}_i \cap \mathcal{C}_j = \emptyset$, $\forall (i, j) \in \{1, \dots, r\}^2$, $i \neq j$.*

Ce corollaire correspond au cas où chaque contribution de rang 1 est observée entièrement. Ceci signifie que les graphes \mathcal{G}_i associés sont tous complets, et par conséquent tous connectés.

7.3.3. Condition nécessaire d'identifiabilité

Le lemme 2 nous donne une condition suffisante d'identifiabilité concrète et relativement simple à vérifier, en supposant les supports des facteurs connus. Cependant, cette condition n'est pas nécessaire et il existe des cas où les contributions de rang 1 sont identifiables alors que la condition du lemme 2 n'est pas remplie. Ceci est dû au fait que, contrairement au cas de la complétion de matrice de rang 1, dans lequel une entrée est soit observée soit non-observée, nous avons ici de l'information à propos des entrées des contributions de rang 1 qui ne sont pas directement observées (via leur somme avec d'autres contributions). Il est par ailleurs assez aisé de construire des exemples simples pour lesquels la condition du lemme 2 n'est pas remplie mais l'identifiabilité est quand même possible, par un algorithme glouton par exemple. Nous donnons en figure 24 des exemples de situations où deux contributions de rang 1 sont identifiables, la condition du lemme 2 étant remplie ou non.

Une condition nécessaire assez immédiate pour que l'identifiabilité tienne est donnée par le lemme suivant.

$$\left(\begin{array}{|cc|cc|} \hline \times & \times & \times & \times \\ \times & \times & \times & \times \\ \times & \times & \times & \times \\ \times & \times & \times & \times \\ \hline \end{array} \right) \left(\begin{array}{|ccc|c|} \hline \times & \times & \times & 0 \\ \times & \times & \times & \times \\ \times & \times & \times & \times \\ 0 & \times & \times & \times \\ \hline \end{array} \right) \left(\begin{array}{|cccc|} \hline \times & \times & \times & \times \\ \times & \times & \times & \times \\ \times & \times & \times & \times \\ \times & \times & \times & \times \\ \hline \end{array} \right)$$

Figure 24 : Trois exemples de situations où les contributions de rang 1 sont identifiables. Pour ces trois exemples, une observation $\mathbf{Z} \in \mathbb{R}^{4 \times 4}$ est représentée, le symbole \times désignant n'importe quel nombre réel. On suppose ici que \mathbf{Z} est la somme de deux contributions de rang 1 dont les supports connus sont représentés par les rectangles rouge et vert. La matrice la plus à gauche représente un cas où la condition du corollaire 1 est remplie. La matrice du milieu représente un cas où la condition du lemme 2 est remplie. Enfin, la matrice la plus à droite représente un cas où la condition suffisante du lemme 2 n'est pas remplie, mais les contributions sont tout de même identifiables (ceci parce que la contribution rouge est identifiable par complétion de matrice de rang 1, et la soustraire à \mathbf{Z} mène à une contribution verte totalement observée).

Lemme 3. *Si les contributions de rang 1 sont identifiables, alors on a*

$$\left| \bigcup_{i=1}^r \text{supp}(\mathbf{C}_i) \right| \geq \|\mathbf{X}\|_0 + \|\mathbf{Y}\|_0 - r$$

Démonstration. Afin que les contributions de rang 1 soient identifiables, il faut que le système d'équations polynomiales correspondant ne soit pas sous-déterminé. Pour cela, il faut absolument qu'il y ait plus d'équations que d'inconnues (Sturmfels, 2002). Or, pour la i -ème contribution, il y a exactement $\|\mathbf{x}_i\|_0 + \|\mathbf{y}^i\|_0 - 1$ inconnues (pour rappel on fixe à 1 la première entrée non-nulle de chaque colonne \mathbf{x}_i afin de se débarrasser des ambiguïtés de phase et d'échelle). Le nombre total d'inconnues est donc égal à $\sum_{i=1}^r (\|\mathbf{x}_i\|_0 + \|\mathbf{y}^i\|_0 - 1)$. Quant aux équations, il y en a une pour chaque entrée de la matrice observée \mathbf{Z} dont la position appartient au support d'au moins une contribution de rang 1. Ceci donne un nombre d'équations total égal au cardinal de l'union des support des contributions de rang 1 : $|\bigcup_{i=1}^r \text{supp}(\mathbf{C}_i)|$ (cette quantité est supérieure ou égale à $\|\mathbf{Z}\|_0$ car des contributions peuvent s'annuler). Ainsi une condition nécessaire d'identifiabilité s'exprime comme

$$\left| \bigcup_{i=1}^r \text{supp}(\mathbf{C}_i) \right| \geq \sum_{i=1}^r (\|\mathbf{x}_i\|_0 + \|\mathbf{y}^i\|_0 - 1).$$

Il suffit alors de distribuer la somme dans le membre de droite pour aboutir au lemme. \square

Ce lemme indique que les supports des contributions de rang 1 doivent être relativement disjoints, de manière à maximiser le membre de gauche de l'inégalité. Sinon l'information donnée par l'observation des entrées de la matrice \mathbf{Z} n'est pas suffisante pour reconstruire ces contributions. Afin d'aboutir à une condition nécessaire plus intéressante (se rapprochant au maximum de la condition suffisante), il faudrait généraliser le type d'argument utilisé dans (Király et Tomioka, 2012). Ceci n'est pas fait ici.

7.4. Exemple : la FFT à partir de la matrice de Fourier

On étudie ici l'identifiabilité de la transformée de Fourier rapide (FFT) classique de type papillon ("butterfly radix 2 FFT"), à partir de l'observation de la matrice de Fourier $\mathbf{F} \in \mathbb{C}^{n \times n}$ avec $f_{ij} = W^{(i-1)(j-1)}$ et $W = e^{\frac{2\pi i}{n}}$ (où $i = \sqrt{-1}$), en supposant que n est une puissance de deux, pour des raisons de simplicité.

Factorisation de référence. La FFT classique peut être vue comme une procédure à étapes correspondant à une succession de multiplications par des matrices creuses \mathbf{F}_j , $j \in \{1, \dots, J = \log_2 n\}$, de telle manière qu'on a $\mathbf{F} = \mathbf{F}_J \dots \mathbf{F}_1$. Nous appelons ces facteurs \mathbf{F}_j les facteurs de références. Ils ont les propriétés suivantes :

- Chaque facteur a exactement deux entrées non-nulles par ligne et par colonne, et est à permutation près bloc-diagonal, les blocs de taille 2×2 correspondant aux papillons.
- Les produits partiels $\mathbf{F}_J \dots \mathbf{F}_\ell$, $\ell = 2, \dots, J - 1$ ont exactement $2^{J-\ell+1}$ entrées non-nulles par lignes et par colonnes.

On s'intéresse dans cette section à l'identifiabilité de la factorisation de référence correspondant à la FFT via la stratégie hiérarchique proposée dans cette thèse, dans laquelle une succession de factorisation en deux facteurs est effectuée, le degré de parcimonie des facteurs étant donné en paramètre.

Hypothèses. Nous considérons ici que les degrés de parcimonie sont donnés par colonne pour la facteur de gauche et par ligne pour le facteur de droite, ce qui revient à fixer la taille des contributions de rang 1 à chaque factorisation. Plus précisément, à chaque factorisation en deux du type $\mathbf{R}_{j-1} \approx \mathbf{X}\mathbf{Y}$, \mathbf{X} est de taille $n \times n$ et chacune de ses colonnes a au plus k entrées non-nulles ($\|\mathbf{x}_i\|_0 \leq k$, $i = 1, \dots, n$) et \mathbf{Y} est aussi de taille $n \times n$ et chacune de ses lignes a au plus l entrées non-nulles ($\|\mathbf{y}^i\|_0 \leq l$, $i = 1, \dots, n$). Nous posons la question suivante : avec des degrés de parcimonie k et l appropriés et connus à chaque étape, les facteurs creux de référence \mathbf{F}_j , $j \in \{1, \dots, J = \log_2 n\}$ correspondant à la FFT sont-ils identifiables via la stratégie hiérarchique ?

Afin de répondre à cette question, nous commençons par énoncer un lemme qui nous sera utile par la suite.

Lemme 4. *Supposons qu'une matrice \mathbf{C} de taille $n \times n$ puisse être factorisée en $\mathbf{C} = \mathbf{A}\mathbf{B}$, \mathbf{A} et \mathbf{B} étant carrées. Si $\|\mathbf{C}\|_0 = \sum_{i=1}^n \|\mathbf{a}_i\|_0 \cdot \|\mathbf{b}^i\|_0$, alors les supports des contributions de rang 1 $\mathbf{C}_i \triangleq \mathbf{a}_i \cdot (\mathbf{b}^i)^T$ sont nécessairement disjoints.*

Démonstration. On a $\mathbf{C} = \sum_{i=1}^n \mathbf{a}_i (\mathbf{b}^i)^T$, et le support de chaque contribution de rang 1 $\mathbf{a}_i (\mathbf{b}^i)^T$ est de cardinal (d'aire) $\|\mathbf{a}_i\|_0 \cdot \|\mathbf{b}^i\|_0$. Le support de \mathbf{C} étant de cardinal $\|\mathbf{C}\|_0$, et étant inclus dans l'union des supports des contributions de rang 1 ($\text{supp}(\mathbf{C}) \subset \bigcup_{i=1}^n \text{supp}(\mathbf{a}_i (\mathbf{b}^i)^T)$), on a forcément $\|\mathbf{C}\|_0 \leq \sum_{i=1}^n \|\mathbf{a}_i\|_0 \cdot \|\mathbf{b}^i\|_0$. Par le principe d'inclusion-exclusion, l'égalité n'est possible que si les supports des contributions de rang 1 sont disjoints. \square

De ce lemme assez général on peut déduire le corollaire suivant plus adapté à notre étude de la matrice de Fourier.

Corollaire 2. *Supposons qu'une matrice \mathbf{C} de taille $n \times n$ dont toutes les entrées sont non-nulles puisse être factorisée en $\mathbf{C} = \mathbf{A}\mathbf{B}$, \mathbf{A} et \mathbf{B} étant carrées. Si $\|\mathbf{a}_i\|_0 \leq n/s$ et $\|\mathbf{b}^i\|_0 \leq s$, $\forall i \in \{1, \dots, n\}$, alors \mathbf{C} peut nécessairement être décomposée en n contributions de rang 1 $\mathbf{C}_i \triangleq \mathbf{a}_i \cdot (\mathbf{b}^i)^T$ à supports disjoints de taille $n/s \times s$.*

Démonstration. Ce corollaire découle simplement du lemme 4, dans le cas où $\|\mathbf{a}_i\|_0 \leq n/s$ et $\|\mathbf{b}^i\|_0 \leq s$, $\forall i \in \{1, \dots, n\}$. Pour que la matrice \mathbf{C} soit totalement pleine, il faut que $\sum_{i=1}^n \|\mathbf{a}_i\|_0 \cdot \|\mathbf{b}^i\|_0 = n^2$. Ceci n'est possible que si $\|\mathbf{a}_i\|_0 = n/s$ et $\|\mathbf{b}^i\|_0 = s$, $\forall i \in \{1, \dots, n\}$, et si les contributions de rang 1 ont des supports disjoints. \square

Ce corollaire implique en particulier que si une factorisation de la matrice de Fourier de la forme $\mathbf{F} = \mathbf{X}\mathbf{Y}$ existe, avec \mathbf{X} ayant $n/2$ entrées non-nulles par colonne et \mathbf{Y} ayant 2 entrées non-nulles par ligne, alors les supports de \mathbf{X} et \mathbf{Y} sont tels que les contributions de rang 1 correspondantes ont des supports disjoints. Ceci signifie que le corollaire 1 est applicable si on connaît ces supports.

Le lemme suivant nous indique qu'une telle décomposition de \mathbf{F} en n blocs de rang 1 de taille $n/2 \times 2$ est en fait unique, et les supports sont bien identifiables.

Lemme 5. *Considérons la matrice de Fourier $\mathbf{F} \in \mathbb{R}^{n \times n}$, avec $n = 2^J$. Les assertions suivantes sont vérifiées :*

1. *Il existe une factorisation du type $\mathbf{F} = \mathbf{A}\mathbf{B}$, \mathbf{A} et \mathbf{B} étant carrées, avec $\|\mathbf{a}_i\|_0 \leq n/2$ et $\|\mathbf{b}^i\|_0 \leq 2$, $\forall i \in \{1, \dots, n\}$.*
2. *\mathbf{F} peut s'écrire $\mathbf{F} = \sum_{i=1}^n \mathbf{C}_i$, où les \mathbf{C}_i sont des contributions de rang 1 à supports disjoints de taille $n/2 \times 2$.*

3. Cette décomposition est unique : les contributions \mathbf{C}_i sont identifiables (à permutation près).

Démonstration. La preuve du premier point est immédiate, en considérant la factorisation de référence, avec $\mathbf{A} = \mathbf{F}_J \dots \mathbf{F}_2$ et $\mathbf{B} = \mathbf{F}_1$, on a alors $\|\mathbf{a}_i\|_0 = n/2$ et $\|\mathbf{b}^i\|_0 = 2$, $\forall i \in \{1, \dots, n\}$. Le second point découle du premier en appliquant le corollaire 2, on a alors $\text{supp}(\mathbf{C}_i) = \text{supp}(\mathbf{a}_i) \times \text{supp}(\mathbf{b}_i)$, et les supports des \mathbf{C}_i sont disjoints et de taille $n/2 \times 2$. Afin de prouver leur identifiabilité (le troisième point du lemme), on peut remarquer que les entrées de la matrice de Fourier s'expriment $f_{ij} = W^{(i-1)(j-1)}$ avec $W = e^{\frac{2\pi i}{n}}$. Décomposer cette matrice en blocs de rang 1 de taille $n/2 \times 2$ revient à trouver pour chaque colonne $i \in \{1, \dots, n\}$ une autre colonne $k \in \{1, \dots, n\} \setminus i$ et $n/2$ lignes distinctes $j_1, \dots, j_{n/2}$ telles que $(f_{j_1 i}, \dots, f_{j_{n/2} i}) = C(f_{j_1 k}, \dots, f_{j_{n/2} k})$, avec C une constante complexe. Ceci se réécrit de la manière suivante :

$$\begin{aligned}
& (W^{(i-1)(j_1-1)}, \dots, W^{(i-1)(j_{n/2}-1)}) = C(W^{(k-1)(j_1-1)}, \dots, W^{(k-1)(j_{n/2}-1)}) \\
\Leftrightarrow & \quad W^{i(j_1-1)-k(j_1-1)} = \dots = W^{i(j_{n/2}-1)-k(j_{n/2}-1)} = C \\
\Leftrightarrow & \quad W^{j_1(i-k)} = \dots = W^{j_{n/2}(i-k)} = C.W^{(i-k)} \\
\Leftrightarrow & \quad W^{j_1(i-k)} = \dots = W^{j_{n/2}(i-k)} \\
\Leftrightarrow & \quad j_1(i-k) \equiv \dots \equiv j_{n/2}(i-k) \pmod{n}. \\
\Leftrightarrow & \quad (j_l - j_m)(i-k) \equiv 0 \pmod{n}, \forall l, m \in \{1, \dots, n/2\}, l \neq m.
\end{aligned} \tag{7.2}$$

Or $|i-k| \in \{1, \dots, n-1\}$ et les indices des lignes j_l sont distincts et appartiennent à $\{1, \dots, n\}$. Si $|i-k| < n/2$, l'égalité de congruence ne peut être satisfaite que si $|j_l - j_m| > 2$, $\forall (l, m) \in \{1, \dots, n/2\}^2, l \neq m$, et il est impossible de trouver autant de tels indices distincts sans qu'au moins l'un d'eux soit plus grand que n (par le principe des tiroirs). La même chose peut être conclue pour $|i-k| > n/2$ en remarquant que ceci est équivalent à $n - |i-k| < n/2$. On a alors forcément $|i-k| = n/2$ et $|j_l - j_m| = 2p$, avec $p \in \mathbb{N}^*$ (les indices des lignes choisies sont de même parité). Ceci implique que la seule manière de décomposer la matrice de Fourier en n blocs de rang 1 de taille $n/2 \times 2$ est de construire les blocs en mettant ensemble la i -ème colonne avec la $i \pm n/2$ -ème, et toutes les lignes d'indices de même parité (c'est exactement comme cela que la FFT "butterfly" classique est construite). \square

Ce lemme combiné avec le lemme 4 implique que si une factorisation $\mathbf{F} = \mathbf{X}\mathbf{Y}$ existe, avec \mathbf{X} ayant $n/2$ entrées non-nulles par colonne et \mathbf{Y} ayant 2 entrées non-nulles par ligne, alors \mathbf{X} est de même support que $\mathbf{F}_J \dots \mathbf{F}_2$ et \mathbf{Y} est de même support que \mathbf{F}_1 . Ceci indique que les supports de \mathbf{F}_1 et de $\mathbf{F}_J \dots \mathbf{F}_2$ sont identifiables à partir de leur degré de parcimonie (par ligne et colonne), à permutation près.

Ensuite, on peut remarquer que si ces supports sont connus, la condition suffisante d'identifiabilité du lemme 2 s'applique (on est même ici dans le cas du corollaire 1), on a alors forcément $\mathbf{X} = \mathbf{F}_J \dots \mathbf{F}_2$ et $\mathbf{Y} = \mathbf{F}_1$ aux ambiguïtés classiques près. De

plus, la matrice $\mathbf{F}_J \dots \mathbf{F}_2$ a l'intéressante propriété d'être à permutation près bloc-diagonale, les blocs étant la matrice de Fourier de taille $n/2$. Le même argument peut alors être utilisé pour prouver l'identifiabilité de \mathbf{F}_2 et $\mathbf{F}_J \dots \mathbf{F}_3$ par factorisation de $\mathbf{F}_J \dots \mathbf{F}_2$, et ainsi de suite, menant à une factorisation complète identifiable, via la stratégie hiérarchique proposée dans cette thèse. Cependant, identifiable ne signifie pas que l'algorithme empirique implémentant la stratégie hiérarchique permet de retrouver les facteurs creux de référence correspondant à la FFT, cela signifie simplement que le minimum global de la fonction de coût utilisée par la factorisation hiérarchique à chaque étape lorsqu'elle est appliquée avec les contraintes adéquates correspond effectivement à la factorisation de référence. Ceci est résumé par le théorème suivant.

Théorème 3. *Soit $\mathbf{F} = \mathbf{F}_J \dots \mathbf{F}_1$ la matrice de Fourier en dimension $n = 2^J$, les facteurs de référence $\mathbf{F}_1, \dots, \mathbf{F}_J$ correspondant à l'algorithme de FFT "butterfly" classique. L'ensemble des minima globaux du problème d'optimisation*

$$\underset{\mathbf{X}, \mathbf{Y}}{\text{minimiser}} \quad \frac{1}{2} \|\mathbf{S}_J \dots \mathbf{S}_\ell - \mathbf{X}\mathbf{Y}\|_F^2 + \delta_{\mathcal{X}}(\mathbf{X}) + \delta_{\mathcal{Y}}(\mathbf{Y}),$$

avec $\mathcal{X} = \{\mathbf{A} \in \mathbb{C}^{n \times n}, \|\mathbf{a}_i\|_0 \leq \frac{n}{2^\ell}, \forall i\}$ et $\mathcal{Y} = \{\mathbf{A} \in \mathbb{C}^{n \times n}, \|\mathbf{a}^i\|_0 \leq 2, \forall i\}$ correspond pour tout $\ell \in \{1, \dots, J-1\}$ aux ambiguïtés classiques près à $\mathbf{S}_J \dots \mathbf{S}_{\ell+1}, \mathbf{S}_\ell$.

Démonstration. Ce lemme est prouvé en combinant les lemmes 4, 5 et le corollaire 1 comme décrit au paragraphe précédent. \square

Conclusion. Nous avons étudié dans ce chapitre de manière théorique l'identifiabilité de facteurs creux à partir de l'observation de leur produit, via la stratégie de factorisation hiérarchique. Pour cela, nous avons étudié l'identifiabilité de deux facteurs creux à partir de l'observation de leur produit, en supposant les supports des facteurs connus. Nous avons alors donné des conditions suffisantes et d'autres nécessaires sur ces supports pour que l'identifiabilité tienne. Nous avons ensuite analysé l'exemple concret de la transformée de Fourier, pour lequel nous avons montré que les supports des facteurs creux correspondant à la transformée de Fourier rapide sont également identifiables et remplissent la condition suffisante pour que les facteurs soient identifiables à supports connus, ce qui signifie que la factorisation correspondant à la FFT est unique sous ces hypothèses. En terme de perspectives, une analyse théorique des fonctions de coût utilisées dans cette thèse pour effectuer les factorisations, et notamment de leurs points critiques serait intéressante. En effet, ceci permettrait dans l'idéal de déterminer dans quelles conditions les algorithmes proposés convergent effectivement vers le minimum global, ou un minimum local satisfaisant en un sens à définir.

Apprentissage de matrices efficaces

Nous étudions dans ce chapitre l'apprentissage de matrices efficaces à partir de données d'entraînement. Imposer l'efficacité lors d'une tâche d'apprentissage permet non seulement de faciliter la manipulation de la matrice apprise, mais aussi son estimation. Nous formulons tout d'abord le problème de manière générale, puis donnons un algorithme explicite dans le cadre de l'apprentissage de dictionnaire. Nous terminons par une mise en évidence théorique des bienfaits de l'efficacité pour l'estimation en elle-même. Ce chapitre est dans sa majeure partie basé sur un article publié dans la revue IEEE Journal of Selected Topics in Signal Processing ([Le Ma-goarou et Gribonval, 2016b](#)), ainsi que sur un article présenté à la conférence IEEE International Conference on Acoustics, Speech and Signal Processing 2015 ([Le Ma-goarou et Gribonval, 2015a](#)).

Sommaire

8.1	Formulation du problème d'apprentissage	116
8.2	Application à l'apprentissage de dictionnaire	117
8.2.1	Apprentissage de dictionnaires efficaces	117
8.2.2	Expérience de débruitage d'image	119
8.2.3	Bornes de généralisation pour l'apprentissage de dictionnaire	123
8.2.4	Dimension de couverture des matrices efficaces	125
8.3	Vers une application aux nouveaux modèles parcimonieux	127
8.3.1	Algorithme d'apprentissage général	128
8.3.2	Bornes de généralisation pour l'apprentissage de matrices efficaces	129

8.1. Formulation du problème d'apprentissage

En traitement de données, disposer d'un modèle mathématique permettant d'analyser ou de décrire les données est souvent un objectif intermédiaire facilitant l'interprétation, la classification ou la visualisation de ces données (Bengio et al., 2013). Les paramètres du modèle peuvent être déterminés en se basant sur une connaissance approfondie de la nature des données d'intérêt (images, sons, données médicales, etc.) ou inférés à partir d'un ensemble de données (les données d'entraînement), on parle alors d'apprentissage automatique (Bishop, 2006).

Dans ce chapitre, on suppose que l'on souhaite modéliser des données $\mathbf{y} \in \mathbb{R}^l$ qui suivent une distribution de probabilité inconnue μ (on notera $\mathbf{y} \sim \mu$). Pour cela, on utilise un modèle paramétré par une matrice $\mathbf{A} \in \mathbb{R}^{m \times n}$, et dont l'erreur est donnée par la fonction $f_{\mathbf{A}} : \mathbb{R}^l \rightarrow \mathbb{R}^+$ ($f_{\mathbf{A}}(\mathbf{y})$ mesure l'erreur commise par le modèle sur l'échantillon \mathbf{y}). Idéalement, on souhaiterait avoir un modèle le plus précis possible, dont l'erreur moyenne est minimale pour les données considérées. Ce but est atteint si l'on parvient à résoudre le problème d'optimisation suivant :

$$\underset{\mathbf{A}}{\text{minimiser}} \quad \mathbb{E}_{\mathbf{y} \sim \mu} [f_{\mathbf{A}}(\mathbf{y})] + g(\mathbf{A}),$$

où la fonction g est une pénalité permettant de contraindre les paramètres estimés. Cependant, il est impossible en pratique de résoudre ce problème d'optimisation, car le calcul de l'espérance nécessiterait un nombre infini d'échantillons. On utilise alors en pratique un ensemble de données d'entraînement constitué de N échantillons \mathbf{y}_i , $i = 1, \dots, N$, et on remplace simplement l'espérance par une moyenne empirique pour aboutir au problème d'optimisation suivant :

$$\underset{\mathbf{A}}{\text{minimiser}} \quad \frac{1}{N} \sum_{i=1}^N f_{\mathbf{A}}(\mathbf{y}_i) + g(\mathbf{A}). \quad (\text{PE})$$

La résolution d'un tel problème est appelée *minimisation du risque empirique* (Vapnik, 1991), et nous utiliserons la notation $\hat{\mathbf{A}}$ pour désigner la matrice apprise par celle-ci. La qualité de cette matrice apprise sur l'ensemble des données d'entraînement peut être évaluée par la quantité $\frac{1}{N} \sum_{i=1}^N f_{\hat{\mathbf{A}}}(\mathbf{y}_i)$. Si l'on désirait connaître la qualité de la matrice apprise sur n'importe quelle donnée suivant la distribution μ , il faudrait calculer la quantité $\mathbb{E}_{\mathbf{y} \sim \mu} [f_{\hat{\mathbf{A}}}(\mathbf{y})]$. Cette espérance étant impossible à calculer en pratique, on utilise un ensemble de T données \mathbf{y}_j , $j = 1, \dots, T$, non considérées pendant l'apprentissage (appelées données de test) pour l'approcher par la moyenne empirique $\frac{1}{T} \sum_{j=1}^T f_{\hat{\mathbf{A}}}(\mathbf{y}_j)$.

La pertinence statistique de l'approche de minimisation du risque empirique (par

rapport à la minimisation de l'espérance) peut être évaluée par la quantité suivante :

$$\sup_{\mathbf{A} \in \mathcal{A}} \left| \frac{1}{N} \sum_{i=1}^N f_{\mathbf{A}}(\mathbf{y}_i) - \mathbb{E}_{\mathbf{y} \sim \mu} [f_{\mathbf{A}}(\mathbf{y})] \right|,$$

où l'ensemble \mathcal{A} est défini comme l'ensemble des matrices \mathbf{A} pour lesquelles la pénalité $g(\mathbf{A})$ est inférieure à l'infini ($\mathcal{A} = \{\mathbf{A} \in \mathbb{R}^{m \times n}, g(\mathbf{A}) < +\infty\}$), on a donc nécessairement $\hat{\mathbf{A}} \in \mathcal{A}$. Cette quantité représente la différence maximale pour la classe de paramètres considérée entre l'objectif idéal (avec l'espérance) et l'objectif réel utilisé en pratique (avec la moyenne empirique). Elle dépend notamment de la taille de la classe de paramètres \mathcal{A} considérée, et permet de contrôler l'erreur de généralisation commise en ayant recours à la minimisation du risque empirique (Bousquet et Bottou, 2008).

Nous nous intéressons dans ce chapitre au cas où l'ensemble \mathcal{A} contient uniquement des matrices efficientes. Ceci revient à imposer que la matrice apprise $\hat{\mathbf{A}}$ soit factorisable en un produit de matrices creuse : $\hat{\mathbf{A}} = \mathbf{S}_J \dots \mathbf{S}_1$. Une telle contrainte permet non seulement de manipuler la matrice apprise de manière plus aisée (comme précisé au chapitre 4), mais dans un cadre d'apprentissage, elle permet aussi de profiter de conditions plus favorables pour l'apprentissage, en terme d'erreur de généralisation, comme discuté au cours de ce chapitre.

8.2. Application à l'apprentissage de dictionnaire

Nous considérons dans cette section un premier problème d'apprentissage précis : l'apprentissage de dictionnaire. Ce problème est introduit dans le chapitre 2. Il est particulièrement adapté à une stratégie de factorisation hiérarchique, comme nous allons le démontrer dans la sous-section suivante.

8.2.1. Apprentissage de dictionnaires efficientes

Classiquement, on peut distinguer deux manières de choisir un dictionnaire pour la représentation parcimonieuse de signaux d'intérêts (Rubinstein et al., 2010a).

Dictionnaires analytiques. Historiquement, le seul moyen de construire un dictionnaire était d'analyser mathématiquement les données d'intérêt afin de déduire une formule mathématique "simple" permettant de former le dictionnaire. Les dictionnaires construits de cette façon sont appelés *dictionnaires analytiques* (Rubinstein et al., 2010a), et rassemblent notamment la transformée de Fourier, de Hadamard en cosinus discrète ou en ondelettes. Grâce à la structure inhérente à leur mode de construction, les dictionnaires analytiques sont efficientes, car associés à des algorithmes rapides tels que la transformée de Fourier rapide (Cooley et Tukey, 1965) ou la transformée en ondelettes discrète (Mallat, 1989).

Dictionnaires appris. D'autre part, le développement des moyens de calcul modernes a autorisé l'émergence de méthodes d'apprentissage automatique qui permettent d'apprendre un dictionnaire à partir de données d'entraînement (Engan et al., 1999; Aharon et al., 2006; Mairal et al., 2010). Concrètement, à partir de données d'entraînement $\mathbf{y}_i \in \mathbb{R}^m$, $i = 1, \dots, N$ constituant les colonnes d'une matrice $\mathbf{Y} \in \mathbb{R}^{m \times N}$, le principe de l'apprentissage de dictionnaire est d'approcher \mathbf{Y} par le produit d'un dictionnaire $\mathbf{D} \in \mathbb{R}^{m \times n}$ et d'une matrice de coefficients $\mathbf{\Gamma} \in \mathbb{R}^{n \times N}$ dont les colonnes sont parcimonieuses (voir chapitre 2) :

$$\mathbf{Y} \approx \mathbf{D}\mathbf{\Gamma}.$$

De tels dictionnaires appris sont en général bien adaptés aux données considérées. Cependant, ce sont en général des matrices denses sans structure apparente, qui ne sont pas associées à des algorithmes rapides.

Le meilleur des deux mondes. Peut-on concevoir des dictionnaires aussi bien adaptés aux données que les dictionnaires appris tout en étant aussi efficaces que les dictionnaires analytiques ? Cette question a commencé à intéresser les chercheurs récemment (Rubinstein et al., 2010b; Chabiron et al., 2014), et se résume à la possibilité ou non d'apprendre des dictionnaires efficaces. Pour rappel, le problème d'optimisation associé à l'apprentissage de dictionnaire est le suivant :

$$\underset{\mathbf{D}, \mathbf{\Gamma}}{\text{minimiser}} \quad \frac{1}{2} \|\mathbf{Y} - \mathbf{D}\mathbf{\Gamma}\|_F^2 + g_d(\mathbf{D}) + g_\gamma(\mathbf{\Gamma}),$$

où la pénalité g_γ impose des colonnes parcimonieuses à $\mathbf{\Gamma}$ et la pénalité g_d permet d'imposer des contraintes au dictionnaire \mathbf{D} . Nous proposons ici de prendre comme pénalité g_d la fonction indicatrice des matrices factorisables en matrices creuses, afin d'aboutir à un dictionnaire efficace de forme

$$\mathbf{D} = \mathbf{S}_J \dots \mathbf{S}_1.$$

Il est de coutume de séparer le problème d'apprentissage de dictionnaire en deux sous-problèmes (voir chapitre 2) traités de manière alternée, l'un correspondant aux coefficients et l'autre au dictionnaire. Le problème correspondant aux coefficients est le suivant :

$$\underset{\mathbf{\Gamma}}{\text{minimiser}} \quad \frac{1}{2} \|\mathbf{Y} - \mathbf{D}\mathbf{\Gamma}\|_F^2 + g_\gamma(\mathbf{\Gamma}),$$

et celui correspondant au dictionnaire prend la forme suivante :

$$\underset{\mathbf{D}}{\text{minimiser}} \quad \frac{1}{N} \sum_{i=1}^N f_{\mathbf{D}}(\mathbf{y}_i) + g_d(\mathbf{D}),$$

avec $f_{\mathbf{D}}(\mathbf{y}_i) = \frac{1}{2} \|\mathbf{y}_i - \mathbf{D}\boldsymbol{\gamma}_i\|_2^2$. On reconnaît ici un problème très similaire au problème (PE).

L'algorithme que nous proposons pour l'apprentissage de dictionnaire efficient est simplement une adaptation de la stratégie de factorisation hiérarchique proposée au chapitre 4, dans laquelle une mise à jour des coefficients est insérée. Le sous-problème correspondant au dictionnaire est géré comme dans l'algorithme 13 (factorisation en deux puis optimisation globale), et le sous-problème correspondant aux coefficients est géré de manière classique (après l'étape d'optimisation globale). Ceci est résumé par l'algorithme 16 dans lequel :

- l'étape de factorisation (ligne 3) est effectuée exactement comme dans l'algorithme 13.
- l'étape d'optimisation globale (ligne 4) est spécifique au problème d'apprentissage de dictionnaire (elle prend en compte la matrice de coefficients $\mathbf{\Gamma}$). Le problème rencontré correspond alors à une factorisation matricielle multicouche de forme :

$$\underset{\mathbf{R}_\ell, \mathbf{S}_\ell, \dots, \mathbf{S}_1}{\text{minimiser}} \quad \|\mathbf{Y} - \mathbf{R}_\ell \mathbf{S}_\ell \dots \mathbf{S}_1 \mathbf{\Gamma}\|_F^2 + \delta_{\mathcal{R}_\ell}(\mathbf{R}_\ell) + \sum_{j=1}^{\ell} \delta_{\mathcal{S}_j}(\mathbf{S}_j),$$

et peut être traité par l'algorithme `palm4MSA` présenté au chapitre 4. La notation `palm4MSA(\mathbf{Y} , $\ell + 2$, $\{\mathcal{R}_\ell, \{\mathcal{S}_j\}_{j=1}^{\ell}, \{\mathbf{\Gamma}_{\ell-1}\}\}$, init=courant)` signifie qu'on effectue une factorisation en $\ell + 2$ facteurs de la matrice de données \mathbf{Y} , le facteur le plus à droite correspondant aux coefficients étant gardé fixe, les autres appartenant aux ensembles de contraintes \mathcal{R}_ℓ (pour le résidu) et $\{\mathcal{S}_j\}_{j=1}^{\ell}$ (pour les facteurs creux).

- l'étape de mise à jour des coefficients (ligne 5) est effectuée de manière classique. La notation générique `sparseCoding(\mathbf{Y} , $\mathbf{R}_\ell \mathbf{S}_\ell \dots \mathbf{S}_1$)` signifie qu'il est possible d'utiliser n'importe quel algorithme de décomposition parcimonieuse (OMP, IHT, ISTA, etc.) afin de calculer les coefficients parcimonieux correspondant aux données \mathbf{Y} avec le dictionnaire $\mathbf{R}_\ell \mathbf{S}_\ell \dots \mathbf{S}_1$.

8.2.2. Expérience de débruitage d'image

Afin d'illustrer les avantages des dictionnaires efficientes (appelés dictionnaires FA μ ST) par rapport aux dictionnaires appris classiques, nous effectuons ici une expérience de débruitage d'image. Le protocole expérimental que nous suivons ici correspond à une version simplifiée du débruitage d'image à base d'apprentissage de dictionnaire classique (Elad et Aharon, 2006). Premièrement, $N = 10000$ morceaux ("patches") \mathbf{y}_i de taille 8×8 (8 pixels de haut et 8 de large, on a alors $m = 64$, $\mathbf{y}_i \in \mathbb{R}^{64}$) sont pris au hasard dans une image bruitée de taille 512×512 , ils

Algorithme 16 Factorisation hiérarchique pour l'apprentissage de dictionnaire.

Entrée : Un ensemble de données d'entraînement $\{\mathbf{y}_i\}_{i=1}^N$; une estimation initiale du dictionnaire \mathbf{D}_0 et des coefficients $\mathbf{\Gamma}_0$; un nombre de facteurs J ; les ensembles de contraintes correspondant \mathcal{R}_ℓ and \mathcal{S}_ℓ , $\ell \in \{1 \dots J - 1\}$.

- 1: $\mathbf{R}_0 \leftarrow \mathbf{D}_0$
- 2: **for** $\ell = 1$ to $J - 1$ **do**
- 3: Factorisation du résidu $\mathbf{R}_{\ell-1}$ en 2 facteurs :
 $\lambda, \{\mathbf{F}_2, \mathbf{F}_1\} = \text{palm4MSA}(\mathbf{R}_{\ell-1}, 2, \{\mathcal{R}_\ell, \mathcal{S}_\ell\}, \text{init=zéro})$
 $\mathbf{R}_\ell \leftarrow \lambda \mathbf{F}_2$ and $\mathbf{S}_\ell \leftarrow \mathbf{F}_1$
- 4: Optimisation globale :
 $\lambda, \{\mathbf{R}_\ell, \{\mathbf{S}_j\}_{j=1}^\ell, \mathbf{\Gamma}_{\ell-1}\} = \text{palm4MSA}(\mathbf{Y}, \ell + 2, \{\mathcal{R}_\ell, \{\mathcal{S}_j\}_{j=1}^\ell, \{\mathbf{\Gamma}_{\ell-1}\}\}, \text{init=courant})$
- 5: Mise à jour des coefficients :
 $\mathbf{\Gamma}_\ell = \text{sparseCoding}(\mathbf{Y}, \mathbf{R}_\ell \mathbf{S}_\ell \dots \mathbf{S}_1)$
- 6: **end for**
- 7: $\mathbf{S}_J \leftarrow \mathbf{R}_{J-1}$

Sortie : La factorisation estimée : $\lambda, \{\mathbf{S}_j\}_{j=1}^J$.

constituent les colonnes d'une matrice de données $\mathbf{Y} \in \mathbb{R}^{64 \times 10000}$. Le bruit considéré est blanc Gaussien, de variance $\sigma \in \{10, 15, 20, 30, 50\}$. Un dictionnaire est ensuite appris sur la matrice de donnée, puis utilisé pour débruiter l'image complète en effectuant une représentation parcimonieuse avec l'algorithme OMP (Mallat et Zhang, 1993) de tous les morceaux de taille 8×8 contenus dans l'image (pas seulement ceux contenus dans la matrice \mathbf{Y}), en autorisant chaque morceau à être la combinaison linéaire de 5 atomes du dictionnaire appris. L'image est ensuite reconstruite en effectuant la moyenne des morceaux se recouvrant.

Configuration des dictionnaires FA μ ST. Nous avons testé diverses configurations de l'algorithme 16 sur les morceaux d'image bruités. Le nombre n d'atomes du dictionnaire (de colonnes) a été choisi dans l'ensemble $\{128, 256, 512\}$. Inspirés par les transformées rapides, le nombre de facteurs J est pris proche du logarithme de la dimension des signaux rencontrés : nous prenons $J = 4$. Le facteur creux le plus à droite est pris de la même taille que le dictionnaire : $\mathbf{S}_1 \in \mathbb{R}^{m \times n}$, et tous les autres facteurs creux sont pris carrés : $\mathbf{S}_J, \dots, \mathbf{S}_2 \in \mathbb{R}^{m \times m}$. L'algorithme 16 est utilisé, avec pour dictionnaire initial \mathbf{D}_0 un dictionnaire appris par l'algorithme K-SVD (Aharon et al., 2006). Nous effectuons la mise à jour des coefficient (ligne 5 de l'algorithme) à l'aide de l'algorithme OMP (Mallat et Zhang, 1993), en autorisant chaque morceau d'image à utiliser 5 atomes du dictionnaire. Nous prenons des ensembles de contrainte de la forme :

$$- \mathcal{S}_1 = \{\mathbf{S} \in \mathbb{R}^{64 \times n}, \|\mathbf{s}_i\|_0 \leq k, \|\mathbf{S}\|_F = 1\}.$$

- $\mathcal{S}_\ell = \{\mathbf{S} \in \mathbb{R}^{64 \times 64}, \|\mathbf{S}\|_0 \leq s, \|\mathbf{S}\|_F = 1\}$ pour $\ell = 2, \dots, J - 1$
- $\mathcal{R}_\ell = \{\mathbf{R} \in \mathbb{R}^{64 \times 64}, \|\mathbf{R}\|_0 \leq P\rho^{\ell-1}, \|\mathbf{R}\|_F = 1\}$ pour $\ell = 1, \dots, J - 1$

avec $k = s/m \in \{2, 3, 6, 12\}$, $\rho \in \{0.4, 0.5, 0.7, 0.9\}$ et $P = 64^2$. Ceci correspond à un total de 16 configurations de factorisation par taille de dictionnaire n , chacune menant à un gain relatif en complexité différent. Le critère d'arrêt pour l'algorithme palm4MSA est un nombre d'itérations, fixé à 50. On peut noter qu'il est un peu abusif d'utiliser l'algorithme OMP durant l'apprentissage de dictionnaire et à l'étape de débruitage car les dictionnaire FA μ ST qu'on obtient n'a pas ses colonnes normalisées. Cependant, nous utilisons quand même cet algorithme, ce qui revient à utiliser une sorte d'OMP avec des poids, où certains atomes ont plus de chance d'être utilisés que d'autres.

Comparaison à des méthodes de l'état de l'art. La méthode d'apprentissage de dictionnaires efficients est comparé à l'*apprentissage de dictionnaire dense* (DDL pour dense dictionary learning). Nous utilisons K-SVD pour effectuer l'apprentissage de dictionnaire dense. D'autres méthodes ont été testées (par exemple l'apprentissage de dictionnaire en ligne (Mairal et al., 2010)), menant à des résultats qualitativement similaires. L'apprentissage de dictionnaire dense est effectué en suivant exactement le même protocole expérimental que pour les dictionnaire FA μ ST (dictionnaire appris sur 10000 morceaux d'image, et débruitage en autorisant 5 atomes par morceau). Nous utilisons l'implémentation de K-SVD décrite dans (Rubinstein et al., 2008), avec 50 itérations (suffisant pour converger en pratique).

Il est intéressant de noter que de meilleures performance de débruitage peuvent être obtenues si l'on insère l'apprentissage de dictionnaire dans un protocole de débruitage plus sophistiqué, par exemple celui utilisé dans (Elad et Aharon, 2006). Les algorithmes de débruitage les plus performants reposent en effet sur une technique de moyennage intelligente appelée agrégation. Le but de l'expérience menée ici est d'illustrer le potentiel des dictionnaires FA μ ST pour le débruitage. Le protocole expérimental est simplifié au maximum afin de mettre en évidence les bonnes capacités de généralisation de tels dictionnaires par rapport aux dictionnaires denses.

Nous comparons aussi les dictionnaires FA μ ST à un dictionnaire analytique : la transformée en cosinus discrète (DCT). Nous utilisons une DCT dite surcomplète avec le même nombre d'atomes que les dictionnaires FA μ ST ($n \in \{128, 256, 512\}$).

Résultats. L'expérience est effectuée sur une base de données d'images standards utilisées en traitement d'image (prise sur le site www.imageprocessingplace.com). La base est constituée de 12 images de taille 512×512 . La figure 26 montre les résultats de débruitage obtenus pour trois images particulières : celle pour laquelle les dictionnaires FA μ ST aboutissent aux plus mauvaises performances (l'image "Mandrill"), celle pour laquelle ils aboutissent aux meilleures performances (l'image "WomanDarkHair") et enfin une pour laquelle ils aboutissent aux performances les plus

représentatives du comportement typique sur toute la base (l'image "Pirate"). Ces trois images sont représentées sur la figure 25.

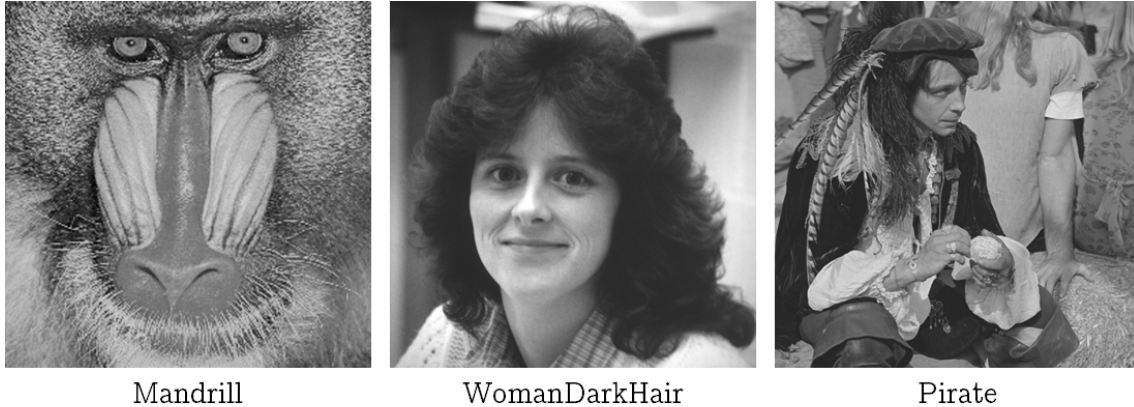


Figure 25 : Les trois images pour lesquelles les résultats de débruitage sont analysés.

Sur la figure 26, chaque croix correspond au débruitage effectué avec un dictionnaire $FA\mu ST$ différent. La couleur de la croix correspond au niveau de bruit considéré. La position de la croix sur l'axe des abscisses correspond au nombre de paramètres du dictionnaire $FA\mu ST$, égal à la quantité $\sum_{j=1}^J \|\mathbf{S}_j\|_0$. La position de la croix sur l'axe des ordonnées correspond à la différence en terme de PSNR (Peak Signal to Noise Ratio, mesure très utilisée en traitement d'image) par rapport au débruitage effectué avec un dictionnaire dense appris à l'aide de l'algorithme K-SVD sur les mêmes morceaux d'image (on note cette quantité $\Delta PSNR$). Enfin, les lignes pointillées correspondent au débruitage effectué avec un dictionnaire analytique (la DCT), leur hauteur correspondant à la différence en terme de PSNR avec le débruitage effectué avec un dictionnaire dense appris par K-SVD. Cette figure nous permet de tirer plusieurs conclusions :

- Premièrement, il est clair que pour le protocole de débruitage simplifié utilisé ici, les dictionnaires $FA\mu ST$ obtiennent de meilleures performances de débruitage que les dictionnaires dense lorsque le niveau de bruit est élevé, c'est-à-dire lorsqu'on a $\sigma = 30$ et $\sigma = 50$. Ceci peut s'expliquer par le fait que lorsque les morceaux d'image utilisés lors de l'apprentissage du dictionnaire sont très bruités, la forte adaptabilité des dictionnaires denses leur est nuisible. En effet, ils ont tendance à apprendre le bruit, ce phénomène étant bien connu sous le nom de sur-apprentissage. Les dictionnaires $FA\mu ST$ quant à eux, semblent moins affectés par ce problème grâce à leur plus grande structure. D'autre part, lorsque le niveau de bruit est faible, le relatif manque d'adaptabilité des dictionnaires $FA\mu ST$ par rapport aux dictionnaires denses est préjudiciable. En effet, ceci est particulièrement vrai pour les images très

texturées (Mandrill typiquement), pour lesquelles le dictionnaire doit être très flexible pour s'adapter aux détails fins de l'image. Les dictionnaires FA μ ST permettent aussi un meilleur débruitage que la DCT lorsque le niveau de bruit est élevé.

- Deuxièmement, il semble que les dictionnaires FA μ ST les plus efficaces (avec le moins de paramètres) mènent à de meilleures performances que les dictionnaires FA μ ST moins efficaces. Ceci peut sembler assez contre-intuitif car ils sont moins adaptables aux données d'entraînement. Cependant, ceci peut s'expliquer par le fait que leur petit nombre de paramètres leur évite le sur-apprentissage. Cependant, ceci n'est pas vrai lorsque le niveau de bruit est faible ou lorsque le nombre de paramètres est vraiment trop petit, car dans ce cas le manque d'adaptabilité du dictionnaire est vraiment trop important.

8.2.3. Bornes de généralisation pour l'apprentissage de dictionnaire

Nous avons vu à la fin de la section précédente que des dictionnaires FA μ ST très efficaces (avec peu de paramètres) peuvent mener à de meilleures performances que des dictionnaires denses. Ceci peut sembler contre-intuitif au premier abord, mais nous donnons dans cette section une explication théorique possible à ce phénomène.

Il est possible d'expliquer théoriquement les bonnes performances de dictionnaire FA μ ST par leur bonne capacité de généralisation. L'erreur de généralisation mesure la différence entre le risque empirique effectivement considéré par les algorithmes d'apprentissage de dictionnaire et l'espérance du risque correspondant (Vainsencher et al., 2011; Maurer et Pontil, 2010). Récemment (Gribonval et al., 2015), une borne générale sur l'erreur de généralisation commise en apprentissage de dictionnaire a été donnée, sous la forme du théorème donné ci-dessous (voir (Gribonval et al., 2015) pour un résultat plus détaillé). Le risque considéré par ce théorème est défini comme

$$f_{\mathbf{D}}(\mathbf{y}) \triangleq \inf_{\boldsymbol{\gamma}} \frac{1}{2} \|\mathbf{y} - \mathbf{D}\boldsymbol{\gamma}\|_2^2 + g_s(\boldsymbol{\gamma}),$$

il revient à considérer la meilleure représentation $\boldsymbol{\gamma}$ possible pour chaque échantillon, par rapport à la pénalité g_s choisie pour imposer la parcimonie aux coefficients.

Théorème 4. (Gribonval et al., 2015). *Dans le cas de l'apprentissage de dictionnaire, on a avec forte probabilité sur le tirage des échantillons \mathbf{y}_i , $i = 1, \dots, N$,*

$$\sup_{\mathbf{D} \in \mathcal{D}} \left| \frac{1}{N} \sum_{i=1}^N f_{\mathbf{D}}(\mathbf{y}_i) - \mathbb{E}_{\mathbf{y} \sim \mu} [f_{\mathbf{D}}(\mathbf{y})] \right| \leq \eta(N, \mathcal{D}),$$

où \mathcal{D} est la classe de dictionnaires considérés, avec $\eta = \mathcal{O} \left(\sqrt{d(\mathcal{D}) \frac{\log N}{N}} \right)$.

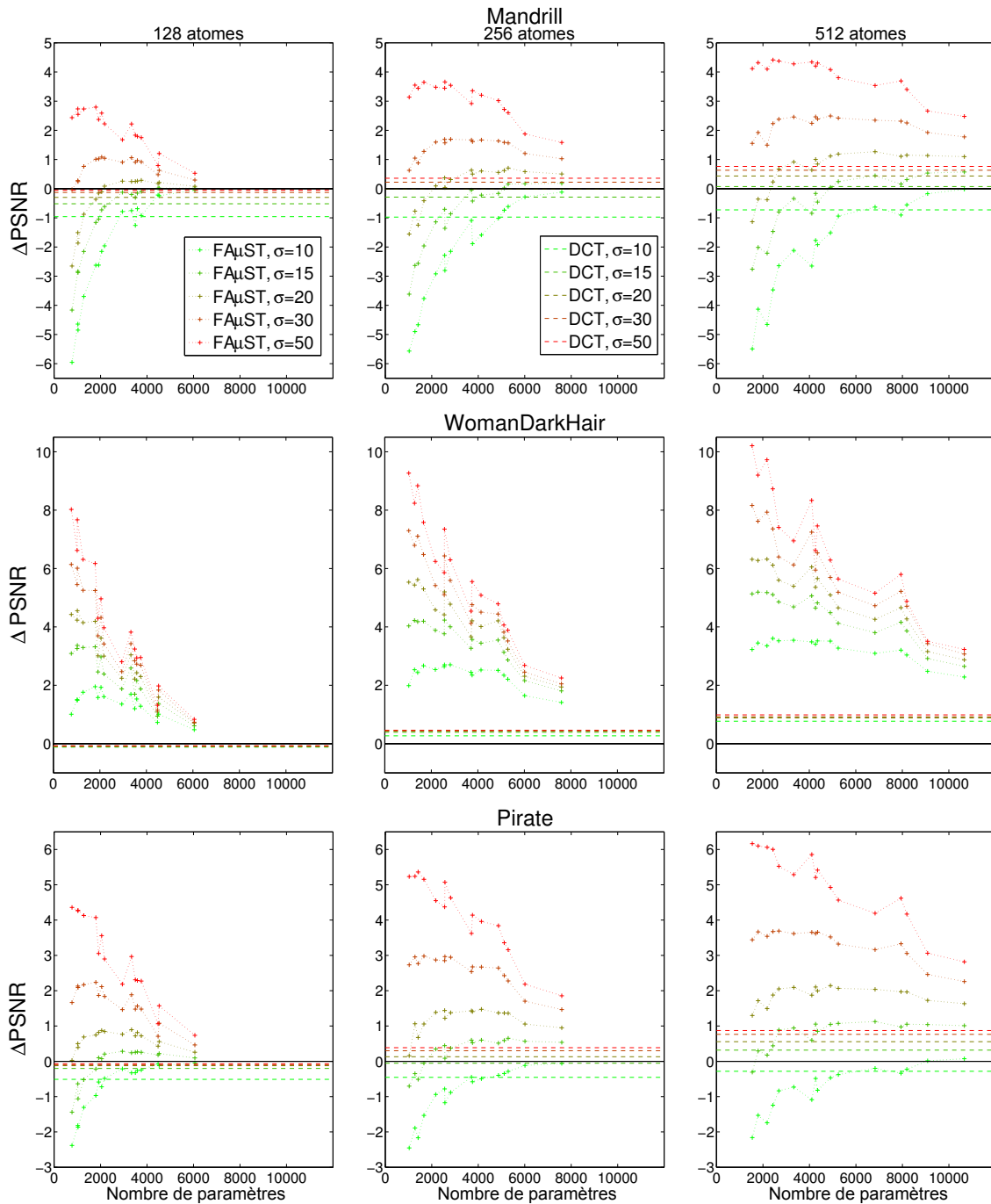


Figure 26 : Résultats de débruitage. La performance relative des dictionnaires $\text{FA}\mu\text{ST}$, comparée à celle des dictionnaires dense et celle de la DCT est représentée. Sur ces figures, l'axe des abscisses correspond à la parcimonie des dictionnaires $\text{FA}\mu\text{ST}$ et l'axe des ordonnées à la différence en terme de PSNR par rapport au débruitage effectué avec un dictionnaire dense appris.

Ce théorème indique en premier lieu que l'erreur de généralisation est réduite lorsque le nombre d'échantillons N qui constituent les données d'entraînement augmente. Ceci paraît parfaitement logique, la moyenne empirique étant calculée sur plus d'échantillon, elle est plus précise. Le théorème indique aussi que la borne sur l'erreur de généralisation dépend de la classe de dictionnaires considérée, et plus précisément de sa dimension de couverture $d(\mathcal{D})$, que nous définissons maintenant.

Définition 10. *Le nombre de couverture $\mathcal{N}(\mathcal{A}, \epsilon)$ d'un ensemble \mathcal{A} est le nombre minimal de boules de rayon ϵ requises pour couvrir l'ensemble \mathcal{A} . Si l'on désigne par $B_\epsilon(\mathbf{c})$ une boule de rayon ϵ centrée en \mathbf{c} , la définition prend la forme*

$$\mathcal{N}(\mathcal{A}, \epsilon) \triangleq \min\{|\mathcal{B}| : \mathcal{A} \subset \bigcup_{\mathbf{c} \in \mathcal{B}} B_\epsilon(\mathbf{c})\}.$$

Le nombre de couverture dépend de la métrique choisie pour définir les boules de rayon ϵ .

Définition 11. *La dimension de couverture $d(\mathcal{A})$ d'un ensemble \mathcal{A} est définie à partir du nombre de couverture de l'ensemble \mathcal{A} comme la quantité*

$$d(\mathcal{A}) \triangleq \limsup_{\epsilon \rightarrow 0} \frac{\log \mathcal{N}(\mathcal{A}, \epsilon)}{\log 1/\epsilon}.$$

Ces définitions peuvent paraître compliquées au premier abord, mais la dimension de couverture d'un ensemble correspond intuitivement au nombre de degrés de libertés de l'ensemble considéré. Par exemple, pour la classe de dictionnaires denses la plus souvent considérée en apprentissage de dictionnaire classique $\mathcal{D}_c = \{\mathbf{D} \in \mathbb{R}^{m \times n}, \|\mathbf{d}_i\|_2 = 1, \forall i\}$, on a $d(\mathcal{D}_c) \leq mn$ (Gribonval et al., 2015). Ceci mène à une borne de généralisation $\eta = \mathcal{O}\left(\sqrt{mn \frac{\log N}{N}}\right)$ pour l'apprentissage de dictionnaire dense. Nous étudions dans la prochaine sous-section ce qu'il en est pour l'apprentissage de dictionnaire efficient.

8.2.4. Dimension de couverture des matrices efficientes

Notre but est ici de donner une borne sur la dimension de couverture de l'ensemble des matrices efficientes, qui nous permettra ensuite d'avoir une borne de généralisation pour l'apprentissage de dictionnaire efficient. Pour cela, nous commençons par définir l'ensemble des matrices efficientes

$$\mathcal{A}_{\text{FA}\mu\text{ST}} \triangleq \{\mathbf{A} \in \mathbb{R}^{m \times n}, \mathbf{A} = \mathbf{S}_J \dots \mathbf{S}_1, \mathbf{S}_j \in \mathbb{R}^{a_j \times a_{j-1}}, \|\mathbf{S}_j\|_0 \leq s_j, \|\mathbf{S}_j\|_F = 1, \forall j\},$$

où l'on a $a_J = m$ et $a_0 = n$. On peut noter que l'on s'intéresse potentiellement à un sous-ensemble de $\mathcal{A}_{\text{FA}\mu\text{ST}}$ lors de l'apprentissage de dictionnaire efficient (ou toute autre tâche d'apprentissage de matrice efficiente). On établit le théorème suivant concernant la dimension de couverture de cet ensemble.

Théorème 5. *La dimension de couverture de l'ensemble des matrices efficientes est bornée par la parcimonie totale de la factorisation sous-jacente :*

$$d(\mathcal{A}_{\text{FA}\mu\text{ST}}) \leq \sum_{j=1}^J \|\mathbf{S}_j\|_0.$$

Démonstration. La preuve de ce théorème se base sur le nombre de couverture des ensembles élémentaires $\mathcal{E}_j = \{\mathbf{A} \in \mathbb{R}^{a_j \times a_{j+1}} : \|\mathbf{A}\|_0 \leq s_j, \|\mathbf{A}\|_F = 1\}$, qui correspondent aux facteurs creux dont le produit appartient à l'ensemble $\mathcal{A}_{\text{FA}\mu\text{ST}}$. Ces ensembles peuvent être vus comme des ensemble de vecteurs de taille $a_j \times a_{j+1}$ de parcimonie bornée par s_j . En combinant les résultats de (Gribonval et al., 2015, lemmes 15 et 16), ceci permet de borner leur nombre de couverture de cette façon (en utilisant la norme de Frobenius comme métrique) :

$$\mathcal{N}(\mathcal{E}_j, \epsilon) \leq \binom{a_j a_{j+1}}{s_j} \left(1 + \frac{2}{\epsilon}\right)^{s_j}.$$

Définissons maintenant le produit cartésien de ces ensembles élémentaires $\mathcal{M} \triangleq \mathcal{E}_1 \times \dots \times \mathcal{E}_J$. Le lemme (Gribonval et al., 2015, lemme 16) indique que le nombre de couverture du produit cartésien d'ensembles quelconques est inférieur au produit des nombres de couverture de ces ensembles (avec la métrique correspondant au maximum des distances entre éléments de chaque ensemble). Ceci nous permet d'avoir $\mathcal{N}(\mathcal{M}, \epsilon) \leq \prod_{j=1}^J \binom{a_j a_{j+1}}{s_j} \left(1 + \frac{2}{\epsilon}\right)^{s_j}$, par rapport à la métrique définie comme $d_1((\mathbf{S}_1, \dots, \mathbf{S}_J), (\mathbf{T}_1, \dots, \mathbf{T}_J)) \triangleq \max_j \|\mathbf{S}_j - \mathbf{T}_j\|_F$. Définissons maintenant une autre métrique sur cet ensemble, de la forme $d_2((\mathbf{S}_1, \dots, \mathbf{S}_J), (\mathbf{T}_1, \dots, \mathbf{T}_J)) \triangleq \sum_{j=1}^J \|\mathbf{S}_j - \mathbf{T}_j\|_F$. On a

$$d_2((\mathbf{S}_1, \dots, \mathbf{S}_J), (\mathbf{T}_1, \dots, \mathbf{T}_J)) \leq J \cdot d_1((\mathbf{S}_1, \dots, \mathbf{S}_J), (\mathbf{T}_1, \dots, \mathbf{T}_J)),$$

d'où la boule de rayon ϵ par rapport à la métrique d_1 est incluse dans la boule de rayon $J\epsilon$ par rapport à la métrique d_2 . Ceci nous permet d'avoir $\mathcal{N}(\mathcal{M}, \epsilon) \leq \prod_{j=1}^J \binom{a_j a_{j+1}}{s_j} \left(1 + \frac{2J}{\epsilon}\right)^{s_j}$ en utilisant la métrique d_2 .

On définit maintenant la fonction liant les facteurs creux et la matrice efficiente correspondante :

$$\begin{aligned} \Phi : \quad \mathcal{M} &\rightarrow \mathcal{A}_{\text{FA}\mu\text{ST}} \\ (\mathbf{S}_1, \dots, \mathbf{S}_J) &\mapsto \mathbf{S}_J \dots \mathbf{S}_1 \end{aligned}$$

Cette fonction Φ est contractante si l'on considère la métrique d_2 dans l'ensemble de départ et la distance induite par la norme de Frobenius ($d_F(\mathbf{A}, \mathbf{B}) \triangleq \|\mathbf{A} - \mathbf{B}\|_F$) dans l'ensemble d'arrivée (ceci se démontre aisément par récurrence, les facteurs étant normalisés). On en conclut que :

$$\mathcal{N}(\mathcal{A}_{\text{FA}\mu\text{ST}}, \epsilon) \leq \prod_{j=1}^J \binom{a_j a_{j+1}}{s_j} \left(1 + \frac{2J}{\epsilon}\right)^{s_j},$$

en utilisant la métrique induite par la norme de Frobenius. Afin de simplifier cette expression, nous utilisons les identités $\binom{n}{p} \leq \frac{n^p}{p!}$ et $n! \geq \sqrt{2\pi n} \left(\frac{n}{e}\right)^n$, ce qui mène à $\mathcal{N}(\mathcal{A}_{\text{FA}\mu\text{ST}}, \epsilon) \leq \prod_{j=1}^J \frac{1}{\sqrt{2\pi s_j}} \left(\frac{e}{s_j} a_j a_{j+1} \left(1 + \frac{2J}{\epsilon}\right)\right)^{s_j}$. En définissant $C_j = \frac{e \cdot a_j \cdot a_{j+1} \cdot (2J+1)}{s_j \cdot 2s_j \sqrt{2\pi s_j}}$, on a $\mathcal{N}(\mathcal{A}_{\text{FA}\mu\text{ST}}, \epsilon) \leq \left(\frac{C}{\epsilon}\right)^h$, avec $h = \sum_{j=1}^J s_j$ et $C = \max_j C_j$. Le calcul de la limite dans l'expression de la dimension de couverture mène alors à $d(\mathcal{A}_{\text{FA}\mu\text{ST}}) \leq h = \sum_{j=1}^J s_j = \sum_{j=1}^J \|\mathbf{S}_j\|_0$, ce qui prouve le théorème. \square

Cette borne sur la dimension de couverture de l'ensemble des matrices efficientes mène à une borne de généralisation $\eta = \mathcal{O}\left(\sqrt{\sum_{j=1}^J \|\mathbf{S}_j\|_0 \frac{\log N}{N}}\right)$ pour l'apprentissage de dictionnaire efficient. Pour rappel, on a $\eta = \mathcal{O}\left(\sqrt{mn \frac{\log N}{N}}\right)$ pour l'apprentissage de dictionnaire dense. Ceci implique que pour un nombre d'échantillons N fixé, la borne sur l'erreur de généralisation commise par l'apprentissage de dictionnaire est meilleure si l'on considère une classe de dictionnaires efficientes, le gain étant de l'ordre de la quantité $\sqrt{\frac{mn}{\sum_{j=1}^J \|\mathbf{S}_j\|_0}}$, qui n'est autre que la racine carrée du gain en complexité relative RCG de la classe de dictionnaires efficientes considérée.

Ceci explique probablement le succès empirique des dictionnaires FA μ ST observé lors des expériences menées à la section précédente.

8.3. Vers une application aux nouveaux modèles parcimonieux

L'apprentissage de dictionnaire fait partie d'une catégorie de modèles plus large : les modèles parcimonieux. Dans un modèle parcimonieux, une représentation parcimonieuse des données est recherchée, et ceci dans le but d'effectuer une tâche quelconque (classification, prédiction, reconstruction, etc.). D'autre part, l'apprentissage de dictionnaire classique est limité en pratique, au moins pour les deux raisons suivantes :

- **Flexibilité** : il mène à des représentations parcimonieuses orientées vers les tâches de reconstruction uniquement (telles que le débruitage, la complétion, etc.). En effet la fonction de coût utilisée pour l'apprentissage de dictionnaire vise uniquement à minimiser l'erreur de reconstruction des données d'entraînement.
- **Coût d'encodage** : calculer la représentation parcimonieuse à l'aide du dictionnaire appris (nous appelons cette procédure l'*encodage* par la suite) requiert de résoudre un problème d'optimisation non-trivial. Ceci est habituellement fait à l'aide d'algorithmes itératifs dont le coût est élevé (surtout pour

des dictionnaires denses). De plus, l'apprentissage en lui-même du dictionnaire nécessite d'effectuer un grand nombre d'encodages, son coût est par conséquent très élevé.

Modèles parcimonieux flexibles à encodage rapide. Récemment, la plupart des travaux dans le domaine des modèles parcimonieux ont pris en compte ces deux limitations. Concernant la flexibilité, il a été proposé (Mairal et al., 2008, 2009; Bagnell et Bradley, 2008) d'apprendre un dictionnaire adapté spécifiquement pour la tâche que l'on souhaite effectuer, et non pas pour reconstruire les données. Ceci est fait en modifiant le problème d'optimisation rencontré en apprentissage de dictionnaire afin de considérer explicitement la tâche à accomplir. Ces méthodes, appelées *apprentissage de dictionnaire supervisé*, ont mené à de nombreux résultats encourageants dans divers domaines (Mairal et al., 2012). Concernant le coût d'encodage, des modèles à coût d'encodage fixé ont été proposés (Sprechmann et al., 2015; Fawzi et al., 2014). L'encodeur proposé correspond alors à un petit nombre fixé d'itérations d'un algorithme de reconstruction parcimonieuse habituellement utilisé en apprentissage de dictionnaire. Ce type de modèles parcimonieux est dit *orienté processus* (process-centric).

Nous étudions dans cette section l'applicabilité des matrices efficaces à ces nouveaux modèles parcimonieux, tout d'abord en proposant un algorithme général pouvant être appliqué à ce type de modèles, puis en étudiant la généralisation de ces nouveaux type de modèles parcimonieux, afin de déterminer si les matrices efficaces sont aussi intéressantes (théoriquement) pour ces modèles que pour l'apprentissage de dictionnaire classique.

8.3.1. Algorithme d'apprentissage général

Nous proposons ici un algorithme général adapté aux situations où le but est d'apprendre une matrice efficace $\hat{\mathbf{A}} = \mathbf{S}_J \dots \mathbf{S}_1$ à partir d'un ensemble de données d'entraînement $\{\mathbf{y}_i\}_{i=1}^N$ et d'une fonction de risque paramétrée par une matrice, que nous notons $f_{\mathbf{A}}$, à minimiser sur cet ensemble. Cette situation revient à considérer le problème (PE) dans sa forme générale.

Principe général. L'algorithme que nous proposons est très proche de l'algorithme de factorisation hiérarchique proposé au chapitre 4. Il consiste simplement à partir d'une estimation initiale $\hat{\mathbf{A}}_0$, la factoriser en deux facteurs à l'aide de l'algorithme `palm4MSA` pour obtenir $\hat{\mathbf{A}}_0 \approx \mathbf{R}_1 \mathbf{S}_1$, puis mettre à jour le facteur creux \mathbf{S}_1 et le résidu \mathbf{R}_1 obtenu en prenant en compte le problème d'apprentissage spécifique auquel on s'intéresse (les données d'entraînement $\{\mathbf{y}_i\}_{i=1}^N$ et la fonction de risque $f_{\mathbf{A}}$). La même procédure est ensuite effectuée sur le résidu \mathbf{R}_1 , puis de manière récursive sur les résidus suivants. L'algorithme résultant de cette récursion est résumé par l'algorithme 17 donné ci-dessous.

Algorithme 17 Factorisation hiérarchique pour l'apprentissage.

Entrée : Un ensemble de données d'entraînement $\{\mathbf{y}_i\}_{i=1}^N$; une estimation initiale $\hat{\mathbf{A}}_0$; une fonction d'erreur $f_{\mathbf{A}}$; un nombre de facteurs J ; les ensembles de contraintes correspondant \mathcal{R}_ℓ and \mathcal{S}_ℓ , $\ell \in \{1 \dots J - 1\}$.

- 1: $\mathbf{R}_0 \leftarrow \hat{\mathbf{A}}_0$
- 2: **for** $\ell = 1$ to $J - 1$ **do**
- 3: Factorisation du résidu $\mathbf{R}_{\ell-1}$ en 2 facteurs :
 $\lambda', \{\mathbf{F}_2, \mathbf{F}_1\} = \text{palm4MSA}(\mathbf{R}_{\ell-1}, 2, \{\mathcal{R}_\ell, \mathcal{S}_\ell\}, \text{init=zéro})$
- 4: $\mathbf{R}_\ell \leftarrow \lambda' \mathbf{F}_2$ and $\mathbf{S}_\ell \leftarrow \mathbf{F}_1$
- 5: Optimisation globale adapté au problème d'apprentissage :
 $\lambda, \{\mathbf{R}_\ell, \{\mathbf{S}_j\}_{j=1}^\ell\} = \text{optGlob}(\{\mathbf{y}_i\}_{i=1}^N, f_{\mathbf{A}}, \ell + 1, \{\mathcal{R}_\ell, \{\mathcal{S}_j\}_{j=1}^\ell\}, \text{init=courant})$
- 6: **end for**
- 7: $\mathbf{S}_J \leftarrow \mathbf{R}_{J-1}$

Sortie : La factorisation estimée : $\lambda, \{\mathbf{S}_j\}_{j=1}^J$.

Optimisation globale. La ligne 5 de l'algorithme 17 correspond à une optimisation globale effectuée sur l'ensemble des facteurs creux introduit jusqu'à l'étape courante et le résidu courant. C'est là la principale différence avec l'algorithme de factorisation hiérarchique présenté au chapitre 4. La notation $\text{optGlob}(\{\mathbf{y}_i\}_{i=1}^N, f_{\mathbf{A}}, \ell + 1, \{\mathcal{R}_\ell, \{\mathcal{S}_j\}_{j=1}^\ell\}, \text{init=courant})$ désigne un algorithme générique dont le but est de résoudre de manière approchée le problème d'optimisation suivant :

$$\underset{\mathbf{R}_\ell, \mathbf{S}_\ell, \dots, \mathbf{S}_1}{\text{minimiser}} \quad \frac{1}{N} \sum_{i=1}^N f_{\mathbf{R}_\ell \mathbf{S}_\ell \dots \mathbf{S}_1}(\mathbf{y}_i) + \delta_{\mathcal{R}_\ell}(\mathbf{R}_\ell) + \sum_{j=1}^{\ell} \delta_{\mathcal{S}_j}(\mathbf{S}_j).$$

Si la fonction de risque $f_{\mathbf{R}_\ell \mathbf{S}_\ell \dots \mathbf{S}_1}$ est différentiable par rapport au produit $\mathbf{R}_\ell \mathbf{S}_\ell \dots \mathbf{S}_1$, alors PALM (Bolte et al., 2014) peut être utilisé pour ce problème, sinon une autre stratégie devra être employée. De plus, la fonction de coût comportant un terme égal à une moyenne empirique, il est possible d'utiliser des techniques de descente de gradient stochastique (Bousquet et Bottou, 2008) pour résoudre ce problème. Ceci revient à ne considérer qu'un sous-ensemble des données d'entraînement (de taille b) choisi aléatoirement à chaque mise à jour effectuée dans l'algorithme `optGlob`, ce qui est équivalent à estimer une moyenne empirique sur moins d'échantillons. On y gagne de la complexité calculatoire (un facteur $\frac{N}{b}$), au prix d'une estimation moins fiable du gradient.

8.3.2. Bornes de généralisation pour l'apprentissage de matrices efficaces

Les avantages des matrices efficaces en terme d'erreur de généralisation existent-ils aussi pour des tâches d'apprentissage correspondant à des modèles parcimonieux autres que l'apprentissage de dictionnaire? Cette question est l'objet de cette sous-section.

Modèle

Nous considérons ici un modèle parcimonieux à la fois orienté processus et supervisé. Le but du modèle est de calculer une représentation parcimonieuse des données d'entrée afin de prédire linéairement une certaine cible (ou label). On considère des données du type $\mathbf{z} \triangleq (\mathbf{x}, \mathbf{y})$, où $\mathbf{x} \in \mathcal{X} \subset \mathbb{R}^m$ correspond à la donnée d'entrée et $\mathbf{y} \in \mathcal{Y} \subset \mathbb{R}^l$ correspond à la cible. On a $\mathbf{z} \sim \mu$ (μ est la distribution conjointe des entrées et des cibles). La représentation parcimonieuse des données d'entrée est calculée à l'aide d'un encodeur explicite $e_{\mathbf{D}} : \mathcal{X} \rightarrow \mathbb{R}^d$ paramétré par un opérateur linéaire $\mathbf{D} \in \mathcal{D} \subset \mathbb{R}^{m \times d}$ appelé dictionnaire. Une prédiction de la cible $\hat{\mathbf{y}} \triangleq \mathbf{W}e_{\mathbf{D}}(\mathbf{x})$ est ensuite calculée à partir de la représentation $e_{\mathbf{D}}(\mathbf{x})$ via un opérateur linéaire $\mathbf{W} \in \mathcal{W} \subset \mathbb{R}^{l \times d}$ (appelé prédicteur). Les paramètres du modèle (le dictionnaire et le prédicteur) constituent l'hypothèse, on utilise la notation $\mathbf{H} \triangleq (\mathbf{D}, \mathbf{W})$. L'erreur de prédiction commise est calculée par la fonction de perte $\ell : \mathcal{Y} \times \mathcal{Y} \rightarrow \mathbb{R}_+$, comme $\ell(\mathbf{y}, \mathbf{W}e_{\mathbf{D}}(\mathbf{x})) \triangleq f_{\mathbf{H}}(\mathbf{z})$.

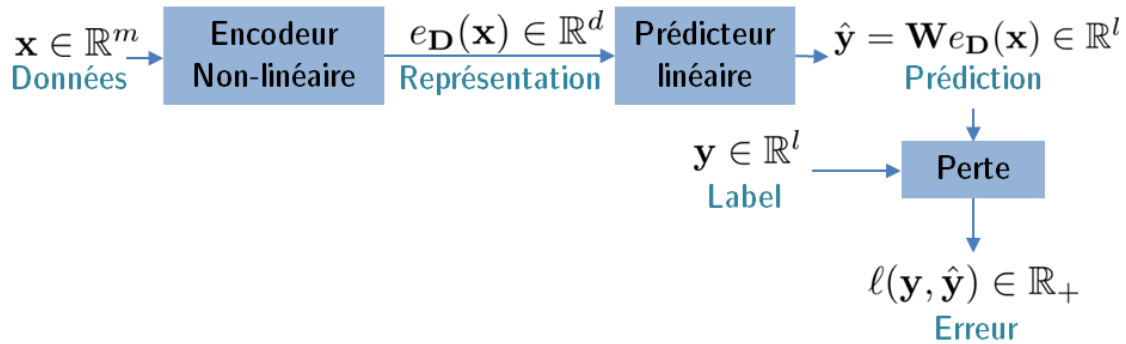


Figure 27 : Modèle parcimonieux supervisé et orienté processus.

Ce type de modèle parcimonieux orienté processus (où l'encodeur est considéré comme une fonction explicite) a été introduit pour la première fois tout récemment (Sprechmann et al., 2015). Le modèle considéré est schématisé en figure 29, et englobe en réalité plusieurs modèles parcimonieux précédemment introduits. Par exemple, l'apprentissage de dictionnaire classique correspond au cas où l'encodeur

est donné par $e_{\mathbf{D}}(\mathbf{x}) \in \arg \min_{\alpha} \|\mathbf{x} - \mathbf{D}\alpha\|_2^2 + g(\alpha)$ (il n'est pas ici explicite car c'est la solution d'un problème d'optimisation non trivial), la cible \mathbf{y} est égale à la donnée d'entrée \mathbf{x} (le but est de reconstruire les données d'entrée), le dictionnaire \mathbf{D} est égal au prédicteur \mathbf{W} (la prédiction est effectuée à l'aide de l'opérateur linéaire utilisé pour le calcul de la représentation parcimonieuse) et la fonction de perte ℓ est simplement le carré de la norme euclidienne de la différence entre l'entrée et la prédiction. En guise de second exemple, l'apprentissage de dictionnaire supervisé pour la classification binaire à l'aide d'un encodeur basé sur le seuillage doux introduit par (Fawzi et al., 2014) correspond au cas où la cible $\mathbf{y} \in \{-1, 1\}$ est un label binaire, l'encodeur $e_{\mathbf{D}}(\mathbf{x}) = h_{\alpha}(\mathbf{D}^T \mathbf{x})$ correspond à une itération de seuillage doux (h_{α} est l'opérateur de seuillage doux qui opère coordonnée par coordonnée) et la perte ℓ est par exemple la perte logistique (ou toute autre fonction de perte adaptée à la classification binaire).

Résultat

Notre but est de fournir une borne de généralisation similaire à celle du théorème 4 dans le cadre de ce nouveau type de modèles parcimonieux. Ceci permettrait de quantifier l'erreur de généralisation commise lorsque les paramètres du modèle sont appris par minimisation du risque empirique $\frac{1}{N} \sum_{i=1}^N f_{\mathbf{H}}(\mathbf{z}_i)$ à partir de données d'entraînement $\mathbf{z}_i = (\mathbf{x}_i, \mathbf{y}_i)$, $i = 1, \dots, N$. De plus, cela permettrait d'avoir un argument théorique pour l'utilisation des matrices efficaces dans un contexte d'apprentissage qui dépasse le cadre de l'apprentissage de dictionnaire classique. Nous fournissons le résultat énoncé par le théorème suivant, dont la preuve est en annexe D.

Théorème 6. *Dans le cas de l'apprentissage de modèle parcimonieux supervisé orienté processus, avec une fonction de perte $\ell(\cdot, \cdot)$ bornée et Lipschitz par rapport à son second argument, une classe d'hypothèse \mathcal{H} compacte et avec un encodeur Lipschitz par rapport au dictionnaire \mathbf{D} de module borné par L_e , on a*

$$\sup_{\mathbf{H} \in \mathcal{H}} \left| \frac{1}{N} \sum_{i=1}^N f_{\mathbf{H}}(\mathbf{z}_i) - \mathbb{E}_{\mathbf{z} \sim \mu} f_{\mathbf{H}}(\mathbf{z}) \right| \leq \eta(N, \mathcal{H}),$$

$$\text{avec } \eta = \mathcal{O} \left(\sqrt{d(\mathcal{H}) \log(L_e) \frac{\log N}{N}} \right).$$

Ce résultat est très similaire dans son énoncé au théorème 4, il indique que l'erreur de généralisation commise dépend aussi dans ce cadre de la dimension de couverture des classes choisies pour le dictionnaire et le prédicteur. Ceci est un argument théorique pour l'utilisation de dictionnaires efficaces ou de prédicteurs efficaces dans ce contexte. La borne dépend par ailleurs des propriétés de stabilité de l'encodeur par

rapport au dictionnaire. Il se trouve que les encodeurs constitués d'un nombre fixé d'itérations d'algorithme proximal (utilisés par (Sprechmann et al., 2015) et (Fawzi et al., 2014)) et les encodeurs de type “elastic-net” (utilisé par (Mairal et al., 2012)) sont stables par rapport au dictionnaire. Pour une démonstration de ce théorème, ainsi que plus de détails et un résultat plus quantitatif, voir l'annexe D.

Conclusion. Nous avons présenté dans ce chapitre l'apprentissage de matrices efficaces. Nous avons tout d'abord présenté une application à l'apprentissage de dictionnaire, ainsi que des expériences de débruitage d'image. Nous avons ensuite montré théoriquement l'avantage d'utiliser des dictionnaires efficaces dans ce cadre, en terme de généralisation. Nous avons enfin discuté de possibles applications à des modèles parcimonieux plus généraux, pour lesquels les avantages théoriques en terme de généralisation sont aussi valables.

Conclusion et perspectives

Ce chapitre de conclusion nous permet de résumer les contributions de cette thèse. Ce faisant, nous isolons quelques problématiques concrètes posées à ce stade du travail, et proposons quelques pistes d'amélioration. Enfin, nous évoquons des perspectives à long terme pour la poursuite de ce travail de thèse.

Résumé des contributions

Nous débutons la conclusion de cette thèse par un rapide résumé des contributions de celle-ci, en listant les différents apports du travail effectué.

Formulation concrète des matrices efficaces. L'objectif principal de cette thèse était de proposer une structure permettant d'effectuer des transformations linéaires de manière efficace. Concrètement parlant, cela revient à effectuer une multiplication par une matrice $\mathbf{A} \in \mathbb{R}^{m \times n}$ pour un coût de calcul réduit (moins de mn multiplications scalaires). Nous avons tout d'abord remarqué que cela n'était possible que si une factorisation existe du type

$$\mathbf{A} = \mathbf{S}_J \dots \mathbf{S}_1,$$

avec $\sum_{j=1}^J \|\mathbf{S}_j\|_0 \leq mn$. Ceci permet de conclure que la structure adaptée aux transformations linéaires efficaces correspond donc à un produit de matrices creuses.

Algorithme d'approximation par matrice efficace. Nous avons proposé un algorithme permettant d'approcher une matrice quelconque par un tel produit de matrices creuses. Le problème d'optimisation correspondant à cette tâche d'approximation est non-convexe. Nous nous basons sur l'algorithme PALM (Bolte et al., 2014) pour aboutir à un algorithme ayant une garantie de convergence vers un point stationnaire de la fonction de coût utilisée. De plus, nous avons proposé une stratégie

de factorisation hiérarchique inspirée des techniques d'optimisation utilisée pour les réseaux de neurones artificiels afin de converger vers de meilleurs minima locaux.

Application à la résolution efficiente de problèmes inverses. Nous avons appliqué cet algorithme à la résolution de problèmes inverses, en factorisant la matrice de mesures associée. Ceci est motivé par le fait que les algorithmes de résolution sont composés principalement de produits avec cette matrice et sa transposée. Nous avons effectué des expériences sur un problème de localisation de sources cérébrales, et obtenu des résultats encourageants mettant en évidence un compromis entre complexité de l'approximation efficiente et précision de résolution.

Application à la transformée de Fourier rapide sur graphe. Nous avons étudié le problème de la transformée de Fourier rapide sur graphe. Le traitement du signal sur graphe est en effet dépourvu à ce jour d'algorithme rapide pour calculer la transformée de Fourier pour des graphes généraux. Nous avons proposé d'appliquer directement la stratégie de factorisation hiérarchique à la matrice de Fourier sur graphe. Cette méthode souffrant de plusieurs limitations dans ce contexte, nous avons proposé une stratégie différente de diagonalisation gloutonne du Laplacien. Nous avons comparé de manière expérimentale ces deux méthodes permettant d'obtenir des transformées de Fourier efficaces approchées sur graphe, mettant en évidence les différences qui les caractérisent.

Étude théorique de l'identifiabilité de la factorisation. Nous avons étudié de manière théorique l'identifiabilité de facteurs creux à partir de l'observation de leur produit, via la stratégie de factorisation hiérarchique. Pour cela, nous avons étudié l'identifiabilité de deux facteurs creux à partir de l'observation de leur produit, en supposant les supports des facteurs connus. Nous avons alors donné des conditions suffisantes et d'autres nécessaires sur ces supports pour que l'identifiabilité tienne. Nous avons ensuite analysé l'exemple concret de la transformée de Fourier, pour lequel nous avons montré que les supports des facteurs creux correspondant à la transformée de Fourier rapide sont également identifiables et remplissent la condition suffisante pour que les facteurs soient identifiables, ce qui signifie que la factorisation correspondant à la FFT est unique sous ces hypothèses.

Apprentissage de matrice efficiente. Nous avons enfin étudié des problèmes d'apprentissage d'opérateur linéaire, sous la contrainte d'efficience. Nous avons proposé dans ce contexte un algorithme d'apprentissage de dictionnaire en tant que produit de matrices creuses, basé sur la stratégie hiérarchique. Nous avons appliqué celui-ci au débruitage d'image, mettant en évidence la bonne capacité de généralisation des dictionnaires efficaces. Nous avons ensuite confirmé de manière théorique cet avantage des dictionnaires efficaces en terme de généralisation. Nous avons ensuite étendu l'analyse théorique précédente à d'autres modèles parcimonieux plus généraux, ouvrant la voie à de futures applications en apprentissage automatique pour les opérateurs linéaires efficaces.

Problématiques et perspectives à court terme

Les travaux effectués lors de cette thèse constituent un premier pas intéressant vers des transformations linéaires efficaces ayant un impact réel sur des applications concrètes. Cependant, quelques améliorations sont requises pour atteindre cet objectif. En effet, les méthodes proposées ici soulèvent les problématiques suivantes :

- (P1) Le *coût de factorisation* est élevé. Ceci rend les méthodes développées dans cette thèse adaptées uniquement aux situations dans lesquelles un temps assez important peut être alloué à la factorisation. Par exemple si une seule et même matrice doit être utilisée pour effectuer des produits matrice/vecteur un grand nombre de fois, alors le coût de sa factorisation peut être amorti par les gains en complexité obtenus à chaque multiplication.
- (P2) Les *gains observés* en pratique sont inférieurs au gain en complexité relative RCG (ceci est mis en évidence au chapitre 5). Celui-ci représente le gain théorique en complexité, étant un rapport entre le nombre de multiplications scalaires requises pour effectuer un produit matrice/vecteur avec une matrice quelconque et avec une FA μ ST. Le temps pris pour effectuer la multiplication par une FA μ ST ne dépend pas uniquement du nombre de multiplications scalaires impliquées. En effet, le nombre de facteurs creux, leurs supports et l'implémentation de la multiplication matrice/vecteur utilisée entrent aussi en compte.
- (P3) Le problème d'optimisation considéré est non-convexe. On se heurte alors à des problèmes de *minima locaux*. L'algorithme que nous proposons jouit de garanties de convergence vers un point stationnaire, mais d'aucune garantie sur la qualité de ce point stationnaire. Celui-ci peut en effet être très mauvais en terme d'attache aux données. Nous avons proposé pour limiter cela la stratégie hiérarchique qui permet de converger vers de meilleurs minima locaux, au prix d'une augmentation du coût de factorisation.
- (P4) Les *paramètres* de l'algorithme de factorisation hiérarchique que nous proposons sont nombreux (taille des facteurs, parcimonie, etc.), et peuvent être difficiles à régler. Nous proposons au cours de la thèse différentes paramétrisations inspirées des factorisations correspondant aux transformées rapides usuelles, mais nous n'avons pas de garanties que cela soit le meilleur choix. De plus, une exploration exhaustive de tout l'espace des paramètres serait extrêmement coûteuse.

Plusieurs améliorations pourraient être apportées aux méthodes développées dans la thèse, afin de répondre à ces différentes problématiques.

- Premièrement les stratégies gloutonnes, comme celle que nous proposons dans le cadre de la transformée de Fourier rapide sur graphe, peuvent permettre de répondre aux problématiques **(P1)** et **(P4)**. En effet, ce type de stratégie permet de ne considérer qu'un seul facteur creux à la fois, et de le fixer une fois pour toute, ce qui réduit grandement le coût de factorisation. De plus, la paramétrisation de ce type de méthode est différente, et plus simple que celle de la stratégie hiérarchique.
- Deuxièmement, afin de réduire le coût de factorisation **(P1)**, il est possible d'imaginer approcher une matrice $\mathbf{A} \in \mathbb{R}^{m \times n}$ sans même stocker cette matrice. En effet, si l'on dispose de N observations de la forme $(\mathbf{x}_i, \mathbf{y}_i \triangleq \mathbf{A}\mathbf{x}_i)$, on peut résoudre de manière approchée le problème d'optimisation suivant

$$\underset{\mathbf{S}_1, \dots, \mathbf{S}_J}{\text{minimiser}} \quad \frac{1}{N} \sum_{i=1}^N \|\mathbf{y}_i - \mathbf{S}_J \dots \mathbf{S}_1 \mathbf{x}_i\|_2^2 + \sum_{j=1}^J \delta_{\mathcal{S}_j}(\mathbf{S}_j),$$

où l'objectif est d'obtenir une matrice efficiente $\mathbf{S}_J \dots \mathbf{S}_1$ ayant le même effet que \mathbf{A} sur les données d'entraînement. Ce problème d'optimisation rentre dans le cadre de l'apprentissage de matrice efficiente décrit au chapitre 8. Il partage par ailleurs de nombreuses similitudes avec l'apprentissage profond (les non-linéarités en moins et les contraintes de parcimonie en plus). À ce titre, les techniques d'optimisation de ce domaine pourraient être utilisées dans ce cas pour réduire encore le coût d'apprentissage (gradient stochastique, pas fixe, back-propagation, etc.).

- Troisièmement, une étude théorique plus poussée de l'algorithme PALM appliqué à la factorisation à deux facteurs creux pourrait mener à une meilleure méthode de factorisation hiérarchique. En effet une compréhension accrue des points stationnaires de la fonction de coût pourrait permettre de choisir une paramétrisation favorable en terme de parcimonie des facteurs (plus ou moins d'entrées non-nulles dans le résidu et/ou le facteur creux), afin de réduire le nombre de minima locaux, et ainsi répondre, ne serait-ce que partiellement, à la problématique **(P3)**.
- Enfin, on pourrait imaginer profiter du cadre très général de l'algorithme PALM pour imposer plus de structure aux facteurs creux lors de la stratégie hiérarchique. Ceci permettrait de répondre à la problématique **(P2)**. En effet, on pourrait contraindre les supports des facteurs creux afin de favoriser une multiplication matrice/vecteur rapide (par exemple en imposant une parcimonie par bloc afin de favoriser la contiguïté mémoire).

Perspectives à long terme

Les transformations linéaires sont omniprésentes en traitement du signal et en apprentissage automatique, et la dimension des données que ces deux domaines considèrent est en constante augmentation. À ce titre, les champs d'application possibles des matrices efficaces sont très nombreux. Nous évoquons ici deux applications pour lesquelles les matrices efficaces semblent particulièrement adaptées.

Échantillonnage compressé. Supposons que l'on dispose de données $\mathbf{x} \in \mathbb{R}^n$ parcimonieuses dans une certaine base orthogonale donnée par les colonnes de la matrice $\mathbf{B} \in \mathbb{R}^{n \times n}$ ($\|\mathbf{B}^T \mathbf{x}\|_0 \leq S$). Le principe de l'échantillonnage compressé (Candès et al., 2006; Donoho, 2006) est d'acquérir des mesures linéaires de ces données $\mathbf{y} = \mathbf{A}\mathbf{x}$, avec $m < n$, où \mathbf{A} est appelée la matrice d'acquisition. Il est alors possible de reconstruire les données pourvu que la base \mathbf{B} et les lignes de \mathbf{A} soient mutuellement incohérentes (peu corrélées). Il a été montré que si les entrées de la matrice \mathbf{A} sont tirées aléatoirement (suivant une loi normale par exemple), alors celle-ci est incohérente avec toutes les autres bases orthogonales avec forte probabilité, et on peut prendre $m = \mathcal{O}(S \log n)$ mesures. La matrice \mathbf{A} étant dans ce cas dense, l'acquisition coûte alors $\mathcal{O}(mn)$ opérations arithmétiques. Est-il possible de faire baisser ce coût d'acquisition ? On pourrait imaginer choisir une matrice d'acquisition efficace $\mathbf{A} = \mathbf{S}_J \dots \mathbf{S}_1$, où les facteurs creux $\mathbf{S}_1, \dots, \mathbf{S}_J$ seraient générés suivant un modèle de matrice creuse aléatoire. On aurait alors un coût d'acquisition en $\mathcal{O}(\sum_{j=1}^J \|\mathbf{S}_j\|_0)$, reste à savoir si le nombre de mesures m à prendre serait affecté, et dans quelle mesure. A priori, on s'attend à voir émerger un compromis entre coût d'acquisition et nombre de mesures. Il est à noter que pour ce type d'application, les problématiques (P1), (P3) et (P4) ne se posent pas, car on génère une matrice efficace, il n'y a donc pas de factorisation impliquée.

Apprentissage profond. Les réseaux de neurones profonds sont utilisés pour des tâches de plus en plus variées (LeCun et al., 2015). Dans certains cas, par exemple en reconnaissance automatique de la parole (Hinton et al., 2012), une grande partie (voire la majorité) des paramètres sont localisés dans la dernière couche du réseau. Il a été proposé d'effectuer une approximation de rang faible de la matrice des poids \mathbf{A} de cette dernière couche afin d'en réduire la complexité lors de l'utilisation du réseau. Cette approximation de rang faible est effectuée soit en post-traitement après l'apprentissage des paramètres (Xue et al., 2013), soit dès l'apprentissage, ce qui revient à modifier l'architecture du réseau (Sainath et al., 2013). Cependant, les approximations de rang faible ne sont pas optimales en terme de compromis complexité/approximation, comme nous l'avons évoqué dans cette thèse. On peut alors imaginer que cette dernière couche soit efficace, et prenne la forme $\mathbf{A} = \mathbf{S}_J \dots \mathbf{S}_1$, et non de rang faible, soit en modifiant directement son architecture avant l'apprentissage, soit en effectuant une factorisation après l'apprentissage. Il est à noter que dans le cas de l'apprentissage profond, un temps très important est en général

alloué à l'apprentissage des paramètres du réseau. On est alors dans le cas favorable évoqué dans (**P1**), l'ajout de la factorisation de la dernière couche après l'apprentissage n'augmentant pas significativement le temps requis avant l'utilisation du réseau (le temps de factorisation étant négligeable par rapport au temps d'entraînement).

Troisième partie

Annexes

Sommaire

A.1 Opérateurs de projection	141
A.1.1 Contrainte de parcimonie générale	141
A.1.2 Contraintes de parcimonie et constante par morceaux	142
A.2 Modules de Lipschitz	143

A.1. Opérateurs de projection

Dans cette annexe, nous donnons les opérateurs de projection correspondant à plusieurs ensemble de contraintes potentiellement intéressant pour des application pratiques.

A.1.1. Contrainte de parcimonie générale

Nous considérons premièrement la contrainte de parcimonie générale correspondant à l'ensemble suivant

$$\mathcal{E} := \{\mathbf{S} \in \mathbb{R}^{p \times q} : \|\mathbf{S}_{\mathcal{H}_i}\|_0 \leq s_i \forall i \in \{1, \dots, K\}, \|\mathbf{S}\|_F = 1\},$$

où $\{\mathcal{H}_1, \dots, \mathcal{H}_K\}$ constitue une partition de l'ensemble des indices, $s_i \in \mathbb{N}$, $\forall i \in \{1, \dots, K\}$, et $\mathbf{S}_{\mathcal{T}}$ est la matrice dont les entrées sont égales à celles de \mathbf{S} sur \mathcal{T} et sont mises à zéro ailleurs. Étant donné une matrice $\mathbf{U} \in \mathbb{R}^{p \times q}$, on souhaite calculer sa projection sur l'ensemble \mathcal{E} : $P_{\mathcal{E}}(\mathbf{U}) \in \arg \min_{\mathbf{S}} \{\|\mathbf{S} - \mathbf{U}\|_F^2 : \mathbf{S} \in \mathcal{E}\}$.

Proposition 7. *Formule de l'opérateur de projection.*

$$P_{\mathcal{E}}(\mathbf{U}) = \frac{\mathbf{U}_{\mathcal{I}}}{\|\mathbf{U}_{\mathcal{I}}\|_F}$$

où \mathcal{I} est l'ensemble des indices correspondant à l'union des s_i entrées de $\mathbf{U}_{\mathcal{H}_i}$ les plus grandes en valeur absolue, $\forall i \in \{1, \dots, K\}$.

Démonstration. Soit \mathbf{S} un élément de \mathcal{E} et \mathcal{J} son support. On a $\|\mathbf{S} - \mathbf{U}\|_F^2 = 1 + \|\mathbf{U}\|_F^2 - 2\langle \text{vec}(\mathbf{U}_{\mathcal{J}}), \text{vec}(\mathbf{S}) \rangle$. Pour un support donné, la matrice $\mathbf{S} \in \mathcal{E}$ qui maximise $\langle \text{vec}(\mathbf{U}_{\mathcal{J}}), \text{vec}(\mathbf{S}) \rangle$ est $\mathbf{S} = \mathbf{U}_{\mathcal{J}} / \|\mathbf{U}_{\mathcal{J}}\|_F$. Pour cette matrice, $\langle \text{vec}(\mathbf{U}_{\mathcal{J}}), \text{vec}(\mathbf{S}) \rangle = \|\mathbf{U}_{\mathcal{J}}\|_F = \sqrt{\sum_{i=1}^K \|\mathbf{U}_{\mathcal{J} \cap \mathcal{H}_i}\|_F^2}$ qui est maximisé si $\mathcal{J} \cap \mathcal{H}_i$ correspond aux s_i entrées les plus grandes en valeur absolue de \mathbf{U} dans la zone déterminée par \mathcal{H}_i , $\forall i \in \{1, \dots, K\}$. \square

A.1.2. Contraintes de parcimonie et constante par morceaux

Soient K ensembles deux à deux disjoints \mathcal{C}_i constituant des indices d'entrées de matrices, considérons l'ensemble de contrainte correspondant aux matrice normalisées qui sont constantes sur chacun des ensembles d'indices \mathcal{C}_i , nulles en dehors de ces ensembles d'indices, et avec pas plus de s zones non-nulles. En d'autres termes : $\mathcal{E}_c := \{\mathbf{S} \in \mathbb{R}^{p \times q} : \exists \tilde{\mathbf{a}} = (\tilde{a}_i)_{i=1}^K, \|\tilde{\mathbf{a}}\|_0 \leq s, \mathbf{S}_{\mathcal{C}_i} = \tilde{a}_i \forall i \in \{1, \dots, K\}, \mathbf{S}_{\overline{\cup_i \mathcal{C}_i}} = 0, \text{ and } \|\mathbf{S}\|_F = 1\}$.

Définissons $\tilde{\mathbf{u}} := (\tilde{u}_i)_{i=1}^K$ avec $\tilde{u}_i := \sum_{(m,n) \in \mathcal{C}_i} u_{mn}$, et désignons $\tilde{\mathcal{J}} \subset \{1, \dots, K\}$ le support de $\tilde{\mathbf{a}}$.

Proposition 8. *La projection de \mathbf{U} sur \mathcal{E}_c est obtenue avec $\tilde{\mathcal{J}}$ la collection des s indices i donnant le plus grand $|\tilde{u}_i|/\sqrt{|\mathcal{C}_i|}$, $\tilde{a}_i := \tilde{u}_i/\sqrt{\sum_{i \in \tilde{\mathcal{J}}} |\mathcal{C}_i| \tilde{u}_i^2}$ if $i \in \tilde{\mathcal{J}}$, $\tilde{a}_i := 0$ sinon.*

Démonstration. Soit \mathbf{S} un élément de \mathcal{E}_c , et $\tilde{\mathcal{J}} \subset \{1, \dots, K\}$ le support du $\tilde{\mathbf{a}}$ associé. En utilisant le même procédé que pour la proposition précédente on remarque que $\langle \text{vec}(\mathbf{U}), \text{vec}(\mathbf{S}) \rangle = \sum_{i \in \tilde{\mathcal{J}}} \langle \text{vec}(\mathbf{U}_{\mathcal{C}_i}), \text{vec}(\mathbf{S}) \rangle = \sum_{i \in \tilde{\mathcal{J}}} \tilde{u}_i \tilde{a}_i = \langle \tilde{\mathbf{u}}_{\tilde{\mathcal{J}}}, \tilde{\mathbf{a}} \rangle$. En effectuant le changement de variable $\tilde{b}_i = \sqrt{|\mathcal{C}_i|} \cdot \tilde{a}_i$ et $\tilde{v}_i = \tilde{u}_i / \sqrt{|\mathcal{C}_i|}$ on obtient $\langle \tilde{\mathbf{u}}_{\tilde{\mathcal{J}}}, \tilde{\mathbf{a}} \rangle = \langle \tilde{\mathbf{v}}_{\tilde{\mathcal{J}}}, \tilde{\mathbf{b}} \rangle$ avec $\tilde{\mathbf{b}} := (\tilde{b}_i)_{i=1}^K$. Étant donné $\tilde{\mathcal{J}}$, maximiser ce produit scalaire sous la contrainte $1 = \|\mathbf{S}\|_F = \|\tilde{\mathbf{b}}\|_2$ donne $\tilde{\mathbf{b}}^* := \tilde{\mathbf{v}}_{\tilde{\mathcal{J}}} / \|\tilde{\mathbf{v}}_{\tilde{\mathcal{J}}}\|_2$, et $\langle \tilde{\mathbf{v}}_{\tilde{\mathcal{J}}}, \tilde{\mathbf{b}}^* \rangle = \|\tilde{\mathbf{v}}_{\tilde{\mathcal{J}}}\|_2$. Maximiser sur $\tilde{\mathcal{J}}$ est obtenu en sélectionnant les s entrées de $\tilde{\mathbf{v}} := (\tilde{v}_i)_{i=1}^K$ les plus grandes en valeur absolue (Proposition 7). Revenir aux variables d'origine donne le résultat. \square

A.2. Modules de Lipschitz

Afin d'estimer un module de Lipschitz du gradient de la partie différentiable de l'objectif, on écrit :

$$\begin{aligned}
 & \left\| \nabla_{\mathbf{S}_j^i} H(\mathbf{L}, \mathbf{S}_1, \mathbf{R}, \lambda^i) - \nabla_{\mathbf{S}_j^i} H(\mathbf{L}, \mathbf{S}_2, \mathbf{R}, \lambda^i) \right\|_F \\
 &= (\lambda^i)^2 \left\| \mathbf{L}^T \mathbf{L} (\mathbf{S}_1 - \mathbf{S}_2) \mathbf{R} \mathbf{R}^T \right\|_F \\
 &\leq (\lambda^i)^2 \|\mathbf{R}\|_2^2 \cdot \|\mathbf{L}\|_2^2 \|\mathbf{S}_1 - \mathbf{S}_2\|_F.
 \end{aligned}$$

Annexe du chapitre 5

Dans cette annexe, nous appliquons une méthode de compression d'opérateurs basée sur des ondelettes (Beylkin et al., 1991) à la matrice de gain étudié au chapitre 5. Dans ce cas précis, il est difficile de définir des ondelettes orthogonales, car les capteurs et les sources ne sont pas disposés sur une grille régulière. De plus, les capteurs ont des orientations différentes, ce qui rend d'autant plus difficile la définition d'ondelettes dans le domaine des capteurs. Nous proposons donc ici d'utiliser des ondelettes sur graphes (Hammond et al., 2011) dans le domaine des sources, qui forment une transformée en ondelettes surcomplète Φ_2 au lieu d'une base orthogonale. Nous n'utilisons pas d'ondelettes dans le domaine des capteurs, ce qui revient à considérer $\Phi_1 = \mathbf{Id}$. On définit les ondelettes dans le domaine des sources avec 5 échelles (ce qui implique $\Phi_2 \in \mathbb{R}^{8193 \times 49158}$), en utilisant le filtre de type "Meyer". On utilise plusieurs niveaux de compression pour $\hat{\mathbf{B}}$, tels que $\|\hat{\mathbf{B}}\|_0 \in \|\mathbf{B}\|_0 \times \{\frac{1}{2}, \frac{1}{5}, \frac{1}{10}, \frac{1}{20}\}$, et plusieurs valeurs m pour l'approximation polynomiale des filtres, tels que $m \in \{10, 25, 50, 100\}$. Les compromis précision/complexité obtenus avec cette méthode de compression basée sur des ondelettes de la matrice de gain \mathbf{M} sont représentés sur la figure 28, où ils sont comparés avec ceux obtenus à l'aide de factorisations creuses multi-couche (les six approximations utilisées au chapitre 5. Il semble clair que dans cette configuration, les approximations FA μ ST mènent à de bien meilleurs compromis, et ceci principalement parce que les ondelettes sur graphe ne sont pas assez efficaces.

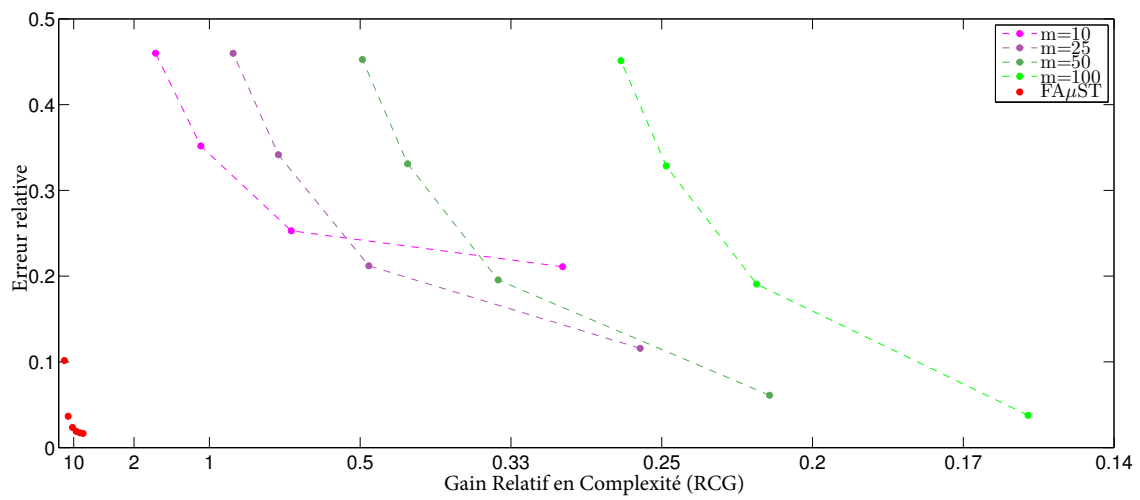


Figure 28 : Comparaison des compromis précision/complexité obtenus avec les approximations FA μ ST et la compression basée sur des ondelettes de la matrice de gain introduite au chapitre 5.

Annexe du chapitre 6

Sommaire

C.1	Solution du sous-problème	147
C.2	Diagonalisation gloutonne efficace du Laplacien	148

C.1. Solution du sous-problème

Le sous-problème rencontré au chapitre 6 est en fait assez simple à résoudre, car une fois le support (p, q) de la rotation de Givens fixé, il dépend uniquement de l'angle de rotation θ et des trois entrées l_{pp}^j , l_{pq}^j and l_{qq}^j . Pour voir cela, désignons $\mathbf{R}_j \triangleq \mathbf{S}^T \mathbf{L}_j \mathbf{S}$. Les entrées de cette matrice sont données par la formule générale $r_{uv}^j = \sum_k \sum_l s_{ku} l_{kl}^j s_{lv}$, qui mène à :

- $r_{uv}^j = l_{uv}^j$ dès que $u \neq p, u \neq q$ and $v \neq p, v \neq q$,
- $r_{pv}^j = cl_{pv}^j + sl_{qv}^j$ dès que $v \neq p, v \neq q$,
 $r_{qv}^j = -sl_{pv}^j + cl_{qv}^j$ dès que $v \neq p, v \neq q$,
tel que $(r_{pv}^j)^2 + (r_{qv}^j)^2 = (l_{pv}^j)^2 + (l_{qv}^j)^2 \forall v \neq p, v \neq q$. Et par symétrie
 $r_{up}^j = cl_{up}^j + sl_{uq}^j$ dès que $u \neq p, u \neq q$,
 $r_{uq}^j = -sl_{up}^j + cl_{uq}^j$ dès que $u \neq p, u \neq q$,
tel que $(r_{vp}^j)^2 + (r_{vq}^j)^2 = (l_{vp}^j)^2 + (l_{vq}^j)^2 \forall u \neq p, u \neq q$.
- $r_{pp}^j = c^2 l_{pp}^j + s^2 l_{qq}^j + 2csl_{pq}^j$,
- $r_{qq}^j = s^2 l_{pp}^j + c^2 l_{qq}^j - 2csl_{pq}^j$,

$$- r_{pq}^j = r_{qp}^j = cs(l_{qq}^j - l_{pp}^j) + (c^2 - s^2)l_{pq}^j.$$

On a donc

$$\mathbf{R}_j = \begin{pmatrix} \ddots & & & & & \\ & c^2 l_{pp}^j + s^2 l_{qq}^j + 2cs l_{pq}^j & & & cs(l_{qq}^j - l_{pp}^j) + (c^2 - s^2)l_{pq}^j & \\ & & \ddots & & & \\ & cs(l_{qq}^j - l_{pp}^j) + (c^2 - s^2)l_{pq}^j & & & s^2 l_{pp}^j + c^2 l_{qq}^j - 2cs l_{pq}^j & \\ & & & & & \ddots \end{pmatrix},$$

en ne montrant que les entrées ayant une influence sur la norme que l'on cherche à minimiser. On a en effet

$$\begin{aligned} \|\mathbf{R}_j\|_{\text{offdiag}}^2 &= \|\mathbf{L}_j\|_{\text{offdiag}}^2 - 2(l_{pq}^j)^2 + 2(r_{pq}^j)^2 \\ &= \|\mathbf{L}_j\|_{\text{offdiag}}^2 - 2(l_{pq}^j)^2 + 2\left[\frac{1}{2}(l_{qq}^j - l_{pp}^j) \sin(2\theta) + l_{pq}^j \cos(2\theta)\right]^2 \\ &= \|\mathbf{L}_j\|_{\text{offdiag}}^2 - 2(l_{pq}^j)^2 \\ &\quad + 2 \left[\sqrt{\frac{1}{4}(l_{qq}^j - l_{pp}^j)^2 + (l_{pq}^j)^2} \cos\left(2\theta - \arctan\left(\frac{l_{qq}^j - l_{pp}^j}{2l_{pq}^j}\right)\right) \right]^2 \end{aligned}$$

Ceci est minimisé pour $\theta = \frac{1}{2} \arctan\left(\frac{l_{qq}^j - l_{pp}^j}{2l_{pq}^j}\right) + (2k + 1)\frac{\pi}{4}$, $k \in \mathbb{Z}$, où la valeur $\|\mathbf{L}_j\|_{\text{offdiag}}^2 - 2(l_{pq}^j)^2$ est atteinte. Ceci suggère de choisir comme support (p, q) le couple qui minimise cette quantité, ce qui correspond à rechercher l'entrée la plus grande en valeur absolue dans la matrice \mathbf{L}_j .

C.2. Diagonalisation gloutonne efficiente du Laplacien

Nous donnons ici une version efficiente de l'algorithme de diagonalisation gloutonne du Laplacien donné au chapitre 6, qui ne nécessite pas de reclasser toutes les entrées de la matrice à chaque itération :

Algorithme 18 Diagonalisation gloutonne approchée du Laplacien (version efficace)

Entrée : Laplacien \mathbf{L} , nombre de facteurs désiré J .

```

1:  $\mathbf{S}_j \leftarrow \mathbf{I} \forall j, \mathbf{L}_1 \leftarrow \mathbf{L}$ 
2: for  $r = 1$  to  $n$  do
3:   for  $s = r + 1$  to  $n$  do
4:      $c_{rs} \leftarrow -2(l_{rs}^j)^2$ 
5:   end for
6:    $d(r) \leftarrow \min \mathbf{C}(r, :)$ 
7:    $e(r) \leftarrow \operatorname{argmin} \mathbf{C}(r, :)$ 
8: end for
9: for  $j = 1$  to  $J$  do
10:   $p \leftarrow \operatorname{argmin} d(r)$ 
11:   $q \leftarrow e(p)$ 
12:   $\theta \leftarrow \frac{1}{2} \arctan\left(\frac{l_{pq}^j - l_{qp}^j}{2l_{pq}^j}\right) + \frac{\pi}{4}$ 
13:   $\mathbf{S}_j \leftarrow \mathbf{G}_{p,q,\theta}$ 
14:   $\mathbf{L}_{j+1} \leftarrow \mathbf{S}_j^T \mathbf{L}_j \mathbf{S}_j$ 
15:  for  $r = p, q$  do
16:    for  $s = r + 1$  to  $n$  do
17:       $c_{rs} \leftarrow -2(l_{rs}^j)^2$ 
18:    end for
19:     $d(r) \leftarrow \min \mathbf{C}(r, :)$ 
20:     $e(r) \leftarrow \operatorname{argmin} \mathbf{C}(r, :)$ 
21:  end for
22:  for  $s = p, q$  do
23:    for  $r = 1$  to  $s - 1$  do
24:       $c_{rs} \leftarrow -2(l_{rs}^j)^2$ 
25:      if  $c_{rs} < d(r)$  then
26:         $d(r) \leftarrow c_{rs}$ 
27:         $e(r) \leftarrow s$ 
28:      else
29:        if  $e(r) = s$  then
30:           $d(r) \leftarrow \min \mathbf{C}(r, :)$ 
31:           $e(r) \leftarrow \operatorname{argmin} \mathbf{C}(r, :)$ 
32:        end if
33:      end if
34:    end for
35:  end for
36: end for
37:  $\mathbf{D} \leftarrow \operatorname{diag}(\mathbf{L}_{J+1})$ 

```

Sortie : La diagonalisation estimée : $\{\mathbf{S}_j\}_{j=1}^J, \mathbf{D}$.



Annexe du chapitre 8

Cette annexe est la reproduction d'un rapport de recherche (rédigé en anglais) contenant les preuves de la borne de généralisation donnée au chapitre 8.

Sommaire

D.1 Introduction	151
D.2 Considered setting	154
D.2.1 Model	154
D.2.2 Generalization	154
D.2.3 Assumptions	155
D.3 Generalization bound using Covering/Lipschitz	156
D.3.1 Union bound and concentration inequality	156
D.3.2 Lipschitz properties of the cost function	157
D.3.3 Final Covering/Lipschitz bound	157
D.4 Generalization bound using McDiarmid/Slepian	158
D.4.1 Application of McDiarmid's inequality	158
D.4.2 Symmetrization	159
D.4.3 Application of Slepian's lemma	159
D.4.4 Final McDiarmid/Slepian bound	160
D.5 Comparison of the two techniques	161
D.6 Stable fast encoders	162

D.1. Introduction

Sparse modeling is a general and convenient mathematical way of giving simple and meaningful descriptions of data, as sparse representations. Parameters of the model are called the dictionary, and can be learned on some training data. Sparse representations using learned dictionaries were introduced two decades ago in a

seminal paper by Olshausen and Field [Olshausen et Field \(1997\)](#), and have since then led to state of the art performance in many areas of signal processing and machine learning [Rubinstein et al. \(2010a\)](#); [Tosic et Frossard \(2011\)](#). Dictionary learning consists in finding a frame, the dictionary, in which the training data can be represented as sparsely as possible. It is usually done by solving a non-convex optimization problem [Aharon et al. \(2006\)](#); [Mairal et al. \(2010\)](#), aiming at minimizing a cost function including a term measuring the capability of the dictionary to reconstruct some training data. Despite its empirical success, this way of doing dictionary learning (that we call classical dictionary learning in the sequel) has been found to be quite limited at least in the two following aspects :

- **Flexibility** : it leads to representations that can be used with success in reconstructive tasks only (such as denoising, completion, etc.). Indeed, sparse representations computed with classically learned dictionaries perform less well in predictive tasks such as regression or classification. This is due to the fact that the optimization problem considered in classical dictionary learning explicitly targets only a good reconstruction of the data.
- **Computational cost** : computing the sparse representation using the learned dictionary (we call this procedure *encoding* in the sequel) implies solving a non-trivial optimization problem. This is usually done by resorting to iterative algorithms which are of high complexity for classical dictionaries, that are in general dense matrices.

Task-driven and computationally efficient sparse models. In recent years, most of the work in the domain of dictionary learning, and more generally sparse models, attempts to address the two aforementioned issues. Regarding flexibility, it was proposed [Mairal et al. \(2008, 2009\)](#); [Bagnell et Bradley \(2008\)](#) to learn a dictionary adapted to a specific supervised learning task, and not only to reconstruct well the training data as is done in classical dictionary learning. This is done by modifying the classical dictionary learning optimization problem in order to explicitly consider the targeted task. This new technique called task-driven dictionary learning (or supervised dictionary learning, or predictive sparse coding) can be seen as a generalization of classical dictionary learning allowing for more flexibility, and led to numerous empirical successes in several domains [Mairal et al. \(2012\)](#). Regarding computational complexity, we can identify two trends to provide more efficient sparse models. First, it was proposed to learn dictionaries with fast transform structures [Rubinstein et al. \(2010b\)](#); [Le Magoarou et Gribonval \(2015b\)](#), allowing to solve the representation problem more efficiently. It consists in looking for dictionaries being products of sparse matrices, that lead to fast dictionary/vector multiplications, that are the building blocks of all classical encoding algorithms [Mallat et Zhang \(1993\)](#); [Tropp et Gilbert \(2007\)](#); [Daubechies et al. \(2004\)](#); [Blumensath et Davies \(2008\)](#).

These techniques can thus be seen as a way to reduce the cost of each iteration of the encoding algorithm. Second, it was also proposed [Sprechmann et al. \(2015\)](#); [Fawzi et al. \(2014\)](#) to use fixed complexity encoder architectures corresponding to a fixed and low number of iterations of a proximal optimization algorithm. This so called *process-centric* view of sparse models can thus be seen as a way to reduce the number of iterations of the encoding algorithm. The two latter techniques act on different components of the encoding algorithm to reduce its overall complexity, and can perfectly, in theory, be combined. Methods allowing a fast encoding step led to encouraging results in various application domains [Sprechmann et al. \(2015\)](#); [Fawzi et al. \(2014\)](#); [Kavukcuoglu et al. \(2008\)](#); [Gregor et LeCun \(2010\)](#); [Rubinstein et al. \(2010b\)](#), and are expected to be even more needed in the future, because of the increasing dimension of signals encountered in modern data processing.

Theory. Despite the ever growing number of sparse models, their theoretical side still remains to be explored in certain areas. The two main theoretical questions regarding dictionary learning are the identifiability : *Is it possible to retrieve by learning a dictionary that generated the data ?*, and the generalization : *How does a learned dictionary behave outside the training corpus ?*. Considering the classical dictionary learning model, these two questions have drawn much attention in the last five years, see [Gribonval et Schnass \(2010\)](#); [Geng et al. \(2011\)](#); [Agarwal et al. \(2013\)](#); [Gribonval et al. \(2014\)](#) for identifiability questions and [Maurer et Pontil \(2010\)](#); [Vainsencher et al. \(2011\)](#); [Gribonval et al. \(2015\)](#) for generalization properties. Considering task-driven dictionary learning, theoretical questions are way less explored. Indeed, to the best of the author’s knowledge, there are no work investigating identifiability of task-driven dictionaries, and only one investigating their generalization properties [Mehta et Gray \(2013\)](#). Moreover, the setting considered in [Mehta et Gray \(2013\)](#) allows the analysis only in the case of a LASSO encoder [Tibshirani \(1996\)](#) with a classical dense dictionary, and relies on pretty strong assumptions, that are not guaranteed to be met in practice.

Contributions and related work. We give the first generalization bounds for task-driven sparse models with fast encoders. Indeed, we propose in this paper a general framework to address the generalization properties of task-driven sparse models, for various fast encoding methods. We give in this vein a unifying view of sparse models allowing to analyze them using the same strategy. We put an emphasis on the recently introduced computationally efficient sparse models [Rubinstein et al. \(2010b\)](#); [Le Magoarou et Gribonval \(2015b\)](#); [Sprechmann et al. \(2015\)](#); [Fawzi et al. \(2014\)](#), in order to assess the influence of the techniques used to achieve computational efficiency on the generalization properties of the model.

We provide generalization bounds for these models both using covering techniques like in [Gribonval et al. \(2015\)](#) and Rademacher averages as done in [Maurer et Pontil \(2010\)](#), for more general cost functions. These two techniques give bounds appro-

priate for different settings.

In comparison to the work of [Mehta et Gray \(2013\)](#) which considers only the LASSO encoder, generalization bounds are provided here considering general fast encoding strategies, and with far less assumptions on the dictionary.

The remaining of the paper is organized as follows. The problem of interest is formulated in section D.2. Generalization bounds are provided using two different techniques in sections D.3 and D.4, and then compared in section D.5. Stability properties of the fast encoders of interest are discussed in section D.6.

D.2. Considered setting

D.2.1. Model

In this paper, we consider random data $\mathbf{z} := (\mathbf{x}, \mathbf{y})$, with $\mathbf{x} \in \mathcal{X} \subset \mathbb{R}^m$ and $\mathbf{y} \in \mathcal{Y} \subset \mathbb{R}^l$, following the probability distribution μ whose support is $\mathcal{Z} := \mathcal{X} \times \mathcal{Y}$. The considered objective is to linearly predict \mathbf{y} using a code computed from the observation of \mathbf{x} via a non-linear encoder $e_{\mathbf{D}} : \mathcal{X} \rightarrow \mathbb{R}^d$ parameterized by a dictionary $\mathbf{D} \in \mathcal{D} \subset \mathbb{R}^{m \times d}$, as $\hat{\mathbf{y}} := \mathbf{W}e_{\mathbf{D}}(\mathbf{x})$, with $\mathbf{W} \in \mathcal{W} \subset \mathbb{R}^{l \times d}$. The parameters of the model are gathered in the hypothesis $\mathbf{H} := (\mathbf{D}, \mathbf{W}) \in \mathcal{H} := \mathcal{D} \times \mathcal{W}$. We consider the loss on a sample to be measured by the function $\ell : \mathcal{Y} \times \mathcal{Y} \rightarrow \mathbb{R}_+$, as $\ell(\mathbf{y}, \mathbf{W}e_{\mathbf{D}}(\mathbf{x})) := f_{\mathbf{H}}(\mathbf{z})$. The loss class is then $\mathcal{F} = \{f_{\mathbf{H}} : \mathcal{Z} \rightarrow \mathbb{R}_+, \mathbf{H} \in \mathcal{H}\}$.

Such a *process-centric* view of sparse models (where the encoder is considered as an explicit function) has been first introduced in [Sprechmann et al. \(2015\)](#). The considered model is schematized on Figure 29, and actually encompasses many known models. For example, classical dictionary learning corresponds to the case where $\mathbf{y} = \mathbf{x}$, $\mathbf{W} = \mathbf{D}$, the encoder is the solution of an optimization problem $e_{\mathbf{D}}(\mathbf{x}) \in \arg \min_{\alpha} \|\mathbf{x} - \mathbf{D}\alpha\|_2^2 + g(\alpha)$ and the loss is the square loss. On the other hand, task-driven dictionary learning for binary classification using a fast encoder based on soft thresholding, as is introduced in [Fawzi et al. \(2014\)](#), corresponds to the case where $\mathbf{y} \in \{-1, 1\}$ is a binary label, the encoder $e_{\mathbf{D}}(\mathbf{x}) = h_{\alpha}(\mathbf{D}^T \mathbf{x})$ corresponds to one iteration of soft-thresholding and the loss is for example the logistic loss.

D.2.2. Generalization

We consider hypotheses estimated from a training set made of n samples : $\mathbf{Z} := \{\mathbf{z}_i = (\mathbf{x}_i, \mathbf{y}_i), i \in \{1, \dots, n\}\}$, via Empirical Risk Minimization (ERM), as $\mathbf{H}_n := \operatorname{argmin}_{\mathbf{H} \in \mathcal{H}} \frac{1}{n} \sum_{i=1}^n f_{\mathbf{H}}(\mathbf{z}_i) + g(\mathbf{H})$, where $g(\cdot)$ is a regularizer. To assess the relevance of this method, we would like to bound uniformly the deviation of the empirical risk

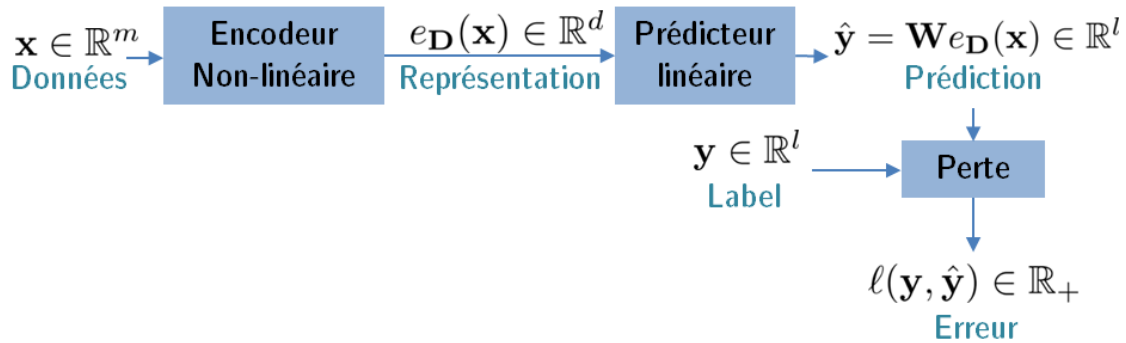


Figure 29 : Schematic architecture of the process-centric view of sparse models.

from the expected risk, namely the quantity :

$$\sup_{\mathbf{H} \in \mathcal{H}} \left| \frac{1}{n} \sum_{i=1}^n f_{\mathbf{H}}(\mathbf{z}_i) - \mathbb{E}_{\mathbf{z} \sim \mu} f_{\mathbf{H}}(\mathbf{z}) \right| := \Psi(\mathbf{Z}) \quad (\text{D.1})$$

In this paper, we propose to bound this quantity using two techniques :

- The first one, that we call *McDiarmid/Slepian*, is similar to what is done in [Maurer et Pontil \(2010\)](#) for matrix factorization, based on McDiarmid’s inequality and Slepian’s lemma.
- The second one, that we call *Covering/Lipschitz*, is based on some covering of the hypothesis class and Lipschitz continuity of the loss function, as is done for example in [Gribonval et al. \(2015\)](#) for matrix factorization.

D.2.3. Assumptions

In order to provide generalization bounds for the model we just introduced, we make the following assumptions :

- **Stable encoder** : We assume the encoder has some stability property with respect to the dictionary \mathbf{D} . The property can take various forms depending which technique is used. When using the covering/Lipschitz technique, we require the encoder $e_{\mathbf{D}}(\mathbf{x})$ to be Lipschitz with respect to \mathbf{D} with modulus $L_e(\mathbf{x})$. When using the McDiarmid/Slepian technique, we require a property of the form $\|e_{\mathbf{D}_1}(\mathbf{x}_i) - e_{\mathbf{D}_2}(\mathbf{x}_i)\|^2 \leq C_1 \cdot \|\mathbf{D}_1^T \mathbf{D}_1 - \mathbf{D}_2^T \mathbf{D}_2\|^2 + C_2 \cdot \|(\mathbf{D}_1 - \mathbf{D}_2)^T \mathbf{x}_i\|^2$. Moreover, we require the encoder to be bounded : $\|e_{\mathbf{D}}(\mathbf{x})\|_2 \leq e_{\max}(\mathbf{x})$. This assumption may seem strong, but we will see in Section D.6 that it actually holds for a wide variety of encoders.

- **Stable and bounded loss** : We assume that the loss $\ell(\cdot, \cdot)$ is Lipschitz with respect to its second argument with modulus L_ℓ and bounded by a positive constant b .
- **Compact hypothesis class** : We assume the hypothesis class \mathcal{H} is compact. In particular, it implies that the dictionary and predictor classes are compact, implying in turn $\|\mathbf{D}\|_F \leq D_{\max}$ and $\|\mathbf{W}\|_F \leq W_{\max}$.

D.3. Generalization bound using Covering/Lipschitz

In this section we present the general strategy we use to bound $\Psi(\mathbf{Z})$ using an union bound in conjunction with a discrete covering of the hypothesis class and Lipschitz properties of $f_{\mathbf{H}}$ with respect to the hypothesis.

D.3.1. Union bound and concentration inequality

This technique relies on the fact that we consider a hypothesis class \mathcal{H} with finite ϵ -covering number $\mathcal{N}(\mathcal{H}, \epsilon)$ with respect to a metric d . This means that for each $\mathbf{H} \in \mathcal{H}$, there exists $\mathbf{H}_j \in \mathcal{H}_j$, with $|\mathcal{H}_j| \leq \mathcal{N}(\mathcal{H}, \epsilon)$ such that $d(\mathbf{H}, \mathbf{H}_j) \leq \epsilon$. Knowing that, let us first rewrite the quantity we wish to bound :

$$\begin{aligned}
\Psi(\mathbf{Z}) &= \sup_{\mathbf{H} \in \mathcal{H}} \left| \frac{1}{n} \sum_{i=1}^n f_{\mathbf{H}}(\mathbf{z}_i) - \mathbb{E}_{\mathbf{z} \sim \mu} f_{\mathbf{H}}(\mathbf{z}) \right. \\
&\quad \left. - \frac{1}{n} \sum_{i=1}^n f_{\mathbf{H}_j}(\mathbf{z}_i) + \frac{1}{n} \sum_{i=1}^n f_{\mathbf{H}_j}(\mathbf{z}_i) \right. \\
&\quad \left. - \mathbb{E}_{\mathbf{z} \sim \mu} f_{\mathbf{H}_j}(\mathbf{z}) + \mathbb{E}_{\mathbf{z} \sim \mu} f_{\mathbf{H}_j}(\mathbf{z}) \right| \\
&\leq \underbrace{\sup_{1 \leq j \leq \mathcal{N}} \left| \frac{1}{n} \sum_{i=1}^n f_{\mathbf{H}_j}(\mathbf{z}_i) - \mathbb{E}_{\mathbf{z} \sim \mu} f_{\mathbf{H}_j}(\mathbf{z}) \right|}_{\text{(T1)}} \\
&\quad + \underbrace{\sup_{\mathbf{H} \in \mathcal{H}} \left| \frac{1}{n} \sum_{i=1}^n f_{\mathbf{H}}(\mathbf{z}_i) - \frac{1}{n} \sum_{i=1}^n f_{\mathbf{H}_j}(\mathbf{z}_i) \right|}_{\text{(T2)}} \\
&\quad + \underbrace{\sup_{\mathbf{H} \in \mathcal{H}} \left| \mathbb{E}_{\mathbf{z} \sim \mu} f_{\mathbf{H}}(\mathbf{z}) - \mathbb{E}_{\mathbf{z} \sim \mu} f_{\mathbf{H}_j}(\mathbf{z}) \right|}_{\text{(T3)}}.
\end{aligned} \tag{D.2}$$

Now, By a union bound and Hoeffding's inequality, (T1) is bounded by t except with probability at most $2\mathcal{N}(\mathcal{H}, \epsilon) \cdot \exp\left(\frac{-2nt^2}{b^2}\right)$. Said differently, (T1) is bounded by $b\sqrt{\frac{\ln(\mathcal{N}(\mathcal{H}, \epsilon)) + \ln(\frac{2}{\delta})}{2n}}$ except with probability at most δ . It remains to bound (T2) and (T3), knowing that $d(\mathbf{H}, \mathbf{H}_j) \leq \epsilon$. This is done in the next subsection, by studying the Lipschitz properties of the cost function.

D.3.2. Lipschitz properties of the cost function

We investigate in this subsection the Lipschitz properties of the cost function $f_{\mathbf{H}}(\mathbf{z})$ with respect to the hypothesis, using the metric d , assuming a stable encoder. We get the following result.

Théorème 9. *The cost function $f_{\mathbf{H}}(\mathbf{z})$ is Lipschitz with respect to \mathbf{H} , namely we have : $|f_{\mathbf{H}_1}(\mathbf{z}) - f_{\mathbf{H}_2}(\mathbf{z})| \leq L_{\mathbf{z}} \cdot d(\mathbf{H}_1, \mathbf{H}_2)$ with $L_{\mathbf{z}} = L_{\ell}(e_{\max}(\mathbf{x}) + W_{\max} \cdot L_e(\mathbf{x}))$ and $d(\mathbf{H}_1, \mathbf{H}_2) \triangleq \max(\|\mathbf{D}_1 - \mathbf{D}_2\|_F, \|\mathbf{W}_1 - \mathbf{W}_2\|_F)$.*

Using this theorem, and assuming we have some upper bound $L_c \geq L_{\mathbf{z}}, \forall \mathbf{z} \in \mathcal{Z}$, then we can bound (T2),

$$\underbrace{\sup_{\mathbf{H} \in \mathcal{H}} \left| \frac{1}{n} \sum_{i=1}^n f_{\mathbf{H}}(\mathbf{z}_i) - \frac{1}{n} \sum_{i=1}^n f_{\mathbf{H}_j}(\mathbf{z}_i) \right|}_{(T2)} \leq \frac{1}{n} \sum_{i=1}^n L_{\mathbf{z}_i} \cdot \epsilon \leq L_c \cdot \epsilon.$$

A very similar reasoning can be applied for (T3) :

$$\underbrace{\sup_{\mathbf{H} \in \mathcal{H}} \left| \mathbb{E}_{\mathbf{z} \sim \mu} f_{\mathbf{H}}(\mathbf{z}) - \mathbb{E}_{\mathbf{z} \sim \mu} f_{\mathbf{H}_j}(\mathbf{z}) \right|}_{(T3)} = \sup_{\mathbf{H} \in \mathcal{H}} \left| \lim_{n \rightarrow \infty} \frac{1}{n} \sum_{i=1}^n f_{\mathbf{H}}(\mathbf{z}_i) - f_{\mathbf{H}_j}(\mathbf{z}_i) \right| \leq L_c \cdot \epsilon.$$

We now have bounded all three terms in play, we thus can state a generalization bound, in the next subsection.

D.3.3. Final Covering/Lipschitz bound

First, let us remind that we decomposed the quantity $\Psi(\mathbf{Z})$ we wish to bound as

$$\Psi(\mathbf{Z}) \leq (T1) + (T2) + (T3).$$

So, combining the results of the two previous subsections we get, with probability at least $1 - \delta$,

$$\Psi(\mathbf{Z}) \leq b\sqrt{\frac{\ln(\mathcal{N}(\mathcal{H}, \epsilon)) + \ln(\frac{2}{\delta})}{2n}} + 2L_c \epsilon.$$

Assuming a particular form for the covering number, namely $\mathcal{N}(\mathcal{H}, \epsilon) = \left(\frac{C}{\epsilon}\right)^h$, where h is called the covering dimension of the hypothesis class, we get with probability at least $1 - \delta$,

$$\Psi(\mathbf{Z}) \leq b \sqrt{\frac{h \ln(C) + h \ln\left(\frac{1}{\epsilon}\right) + \ln\left(\frac{2}{\delta}\right)}{2n}} + 2L_c \epsilon.$$

Now, taking $\epsilon = \frac{1}{\sqrt{n}}$ gives with probability at least $1 - \delta$,

$$\Psi(\mathbf{Z}) \leq \frac{b}{2} \sqrt{\frac{h \log(n)}{n}} + \left(b \sqrt{\frac{h \ln(C) + \ln\left(\frac{2}{\delta}\right)}{2}} + 2L_c \right) \frac{1}{\sqrt{n}}. \quad (\text{B1})$$

Another possibility is to take $\epsilon = \frac{b \sqrt{\beta \log(n)}}{L_c \sqrt{n}}$ with $\beta = h \cdot \max(\log(\frac{2CL_c}{b}), 1)$, as is done in [Gribonval et al. \(2015\)](#), gives with probability at least $1 - \delta$,

$$\Psi(\mathbf{Z}) \leq 2b \sqrt{\frac{\beta \log(n)}{n}} + b \sqrt{\frac{\beta + \log\left(\frac{2}{\delta}\right)}{2n}}. \quad (\text{B2})$$

D.4. Generalization bound using McDiarmid/Slepian

In this section we present the general strategy we use to bound $\Psi(\mathbf{Z})$ using McDiarmid's inequality in conjunction with Slepian's lemma.

D.4.1. Application of McDiarmid's inequality

McDiarmid's inequality is a crucial ingredient of our bound. It requires to bound the following quantity :

$$\begin{aligned} & \left| \Psi(\mathbf{z}_1, \dots, \mathbf{z}_j, \dots, \mathbf{z}_n) - \Psi(\mathbf{z}_1, \dots, \mathbf{z}'_j, \dots, \mathbf{z}_n) \right| \\ &= \left| \sup_{\mathbf{H} \in \mathcal{H}} \left| \frac{1}{n} \sum_{i=1}^n f_{\mathbf{H}}(\mathbf{z}_i) - \mathbb{E}_{\mathbf{z} \sim \mu} f_{\mathbf{H}}(\mathbf{z}) \right| \right. \\ & \quad \left. - \sup_{\mathbf{H} \in \mathcal{H}} \left| \frac{1}{n} \sum_{i=1}^n f_{\mathbf{H}}(\mathbf{z}_i) - \mathbb{E}_{\mathbf{z} \sim \mu} f_{\mathbf{H}}(\mathbf{z}) + \frac{1}{n} (f_{\mathbf{H}}(\mathbf{z}'_j) - f_{\mathbf{H}}(\mathbf{z}_j)) \right| \right| \\ &\leq \sup_{\mathbf{H} \in \mathcal{H}} \frac{1}{n} |f_{\mathbf{H}}(\mathbf{z}'_j) - f_{\mathbf{H}}(\mathbf{z}_j)| \\ &= \sup_{\mathbf{H} \in \mathcal{H}} \frac{1}{n} \left| \ell(\mathbf{y}'_j, \mathbf{W}e_{\mathbf{D}}(\mathbf{x}'_j)) - \ell(\mathbf{y}_j, \mathbf{W}e_{\mathbf{D}}(\mathbf{x}_j)) \right| \end{aligned}$$

The loss function being bounded, $\ell : \mathcal{Y} \times \mathcal{Y} \rightarrow [0, b]$, then

$$\sup_{\mathbf{H} \in \mathcal{H}^n} \frac{1}{n} \left| \ell(\mathbf{y}'_j, \mathbf{W}_{e_{\mathbf{D}}}(\mathbf{x}'_j)) - \ell(\mathbf{y}_j, \mathbf{W}_{e_{\mathbf{D}}}(\mathbf{x}_j)) \right| \leq \frac{b}{n}$$

Applying McDiarmid's inequality then gives us, with probability at least $1 - \delta$,

$$\Psi(\mathbf{Z}) \leq \mathbb{E}\Psi(\mathbf{Z}) + b\sqrt{\frac{\ln(1/\delta)}{2n}}.$$

It now remains to bound $\mathbb{E}\Psi(\mathbf{Z})$, and this is the object of the next two subsections.

D.4.2. Symmetrization

Symmetrization is a pretty standard technique allowing to upperbound $\mathbb{E}\Psi(\mathbf{Z})$. We do not detail it here, it goes as follows :

$$\mathbb{E}\Psi(\mathbf{Z}) = \mathbb{E} \left[\sup_{\mathbf{H} \in \mathcal{H}} \left| \frac{1}{n} \sum_{i=1}^n f_{\mathbf{H}}(\mathbf{z}_i) - \mathbb{E}_{\mathbf{z} \sim \mu} f_{\mathbf{H}}(\mathbf{z}) \right| \right] \leq \mathcal{R}_n(\mathcal{F}, \mu) \leq \sqrt{\pi/2} \Gamma_n(\mathcal{F}, \mu),$$

where

$$\mathcal{R}_n(\mathcal{F}, \mu) \triangleq \frac{2}{n} \mathbb{E} \sup_{f \in \mathcal{F}} \sum_{i=1}^n \sigma_i f(\mathbf{z}_i)$$

is called Rademacher average (the σ_i s are i.i.d. Rademacher random variables), and

$$\Gamma_n(\mathcal{F}, \mu) \triangleq \frac{2}{n} \mathbb{E} \sup_{f \in \mathcal{F}} \sum_{i=1}^n \gamma_i f(\mathbf{z}_i)$$

is called Gaussian average (the γ_i s are i.i.d. standard Gaussian random variables) of the loss class \mathcal{F} .

The problem is now to bound the expectation of the supremum over the hypothesis class \mathcal{H} of a Gaussian process $\Omega_{\mathbf{H}} := \sum_{i=1}^n \gamma_i f_{\mathbf{H}}(\mathbf{z}_i)$. This is done using Slepian's lemma in the next subsection.

D.4.3. Application of Slepian's lemma

Let us now introduce the second main ingredient of our bound, namely Slepian's lemma. We now apply this result to our problem, considering the Gaussian process

$\Omega_{\mathbf{H}} := \sum_{i=1}^n \gamma_i f_{\mathbf{H}}(\mathbf{z}_i)$, so that $\mathbb{E}\Psi(\mathbf{Z}) \leq \frac{\sqrt{2\pi}}{n} \mathbb{E}_{\gamma} \sup_{\mathbf{H} \in \mathcal{H}} \Omega_{\mathbf{H}}$. We then have :

$$\begin{aligned}
& \mathbb{E}_{\gamma} (\Omega_{\mathbf{H}_1} - \Omega_{\mathbf{H}_2})^2 \\
&= \sum_{i=1}^n (f_{\mathbf{H}_1}(\mathbf{z}_i) - f_{\mathbf{H}_2}(\mathbf{z}_i))^2 \\
&= \sum_{i=1}^n \left(\ell(\mathbf{y}_i, \mathbf{W}_1 e_{\mathbf{D}_1}(\mathbf{x}_i)) - \ell(\mathbf{y}_i, \mathbf{W}_2 e_{\mathbf{D}_2}(\mathbf{x}_i)) \right)^2 \\
&\leq \sum_{i=1}^n L_{\ell_2}^2 \cdot \|\mathbf{W}_1 e_{\mathbf{D}_1}(\mathbf{x}_i) - \mathbf{W}_2 e_{\mathbf{D}_2}(\mathbf{x}_i)\|^2 \\
&\leq \sum_{i=1}^n L_{\ell_2}^2 \cdot \|(\mathbf{W}_1 - \mathbf{W}_2) e_{\mathbf{D}_1}(\mathbf{x}_i)\|^2 + L_{\ell_2}^2 \cdot \|\mathbf{W}_2 (e_{\mathbf{D}_1}(\mathbf{x}_i) - e_{\mathbf{D}_2}(\mathbf{x}_i))\|^2 \\
&\leq \sum_{i=1}^n L_{\ell_2}^2 \cdot e_{\max}(\mathbf{x})^2 \cdot \|\mathbf{W}_1 - \mathbf{W}_2\|_F^2 + L_{\ell_2}^2 \cdot W_{\max}^2 \cdot \|e_{\mathbf{D}_1}(\mathbf{x}_i) - e_{\mathbf{D}_2}(\mathbf{x}_i)\|^2 \\
&= \mathbb{E}_{\gamma} (\Xi_{\mathbf{H}_1} - \Xi_{\mathbf{H}_2})^2,
\end{aligned}$$

with (using the stable encoder assumption)

$$\begin{aligned}
\Xi_{\mathbf{H}} := & L_{\ell_2} \cdot e_{\max}(\mathbf{x}) \sum_{ikl} \gamma_{ikl} \cdot w_{kl} \\
& + L_{\ell_2} \cdot C_1 \cdot W_{\max} \sum_{ikl} \gamma_{ikl} \cdot \langle \mathbf{d}_k, \mathbf{d}_l \rangle \\
& + L_{\ell_2} \cdot C_2 \cdot W_{\max} \sum_{ik} \gamma_{ik} \cdot \langle \mathbf{d}_k, \mathbf{x}_i \rangle.
\end{aligned}$$

Then by Slepian's lemma we have $\mathbb{E}_{\gamma} \sup_{\mathbf{H} \in \mathcal{H}} \Omega_{\mathbf{H}} \leq \mathbb{E}_{\gamma} \sup_{\mathbf{H} \in \mathcal{H}} \Xi_{\mathbf{H}}$. Let us now bound the supremum of this Gaussian process over the hypothesis class,

$$\begin{aligned}
& \mathbb{E}_{\gamma} \sup_{\mathbf{H} \in \mathcal{H}} \Xi_{\mathbf{H}} \\
&\leq \mathbb{E}_{\gamma} \sup_{\mathbf{H} \in \mathcal{H}} L_{\ell_2} \cdot e_{\max}(\mathbf{x}) \sum_{ikl} \gamma_{ikl} \cdot w_{kl} \\
&\quad + \mathbb{E}_{\gamma} \sup_{\mathbf{H} \in \mathcal{H}} L_{\ell_2} \cdot C_1 \cdot W_{\max} \sum_{ikl} \gamma_{ikl} \cdot \langle \mathbf{d}_k, \mathbf{d}_l \rangle \\
&\quad + \mathbb{E}_{\gamma} \sup_{\mathbf{H} \in \mathcal{H}} L_{\ell_2} \cdot C_2 \cdot W_{\max} \sum_{ik} \gamma_{ik} \cdot \langle \mathbf{d}_k, \mathbf{x}_i \rangle \\
&\leq \mathbb{E}_{\gamma} L_{\ell_2} \cdot e_{\max}(\mathbf{x}) \sum_{kl} \sup_{\mathbf{W} \in \mathcal{W}} w_{kl} \sum_i \gamma_{ikl} \\
&\quad + \mathbb{E}_{\gamma} L_{\ell_2} \cdot C_1 \cdot W_{\max} \sum_{kl} \sup_{\mathbf{D} \in \mathcal{D}} \langle \mathbf{d}_k, \mathbf{d}_l \rangle \sum_i \gamma_{ikl} \\
&\quad + \mathbb{E}_{\gamma} L_{\ell_2} \cdot C_2 \cdot W_{\max} \sum_k \sup_{\mathbf{H} \in \mathcal{H}} \langle \mathbf{d}_k, \sum_i \gamma_{ik} \cdot \mathbf{x}_i \rangle \\
&= C \cdot ld \sqrt{n} \\
&\quad + C' \cdot d^2 \sqrt{n} \\
&\quad + C'' \cdot d \sqrt{n},
\end{aligned}$$

with $C = L_{\ell_2} \cdot e_{\max}(\mathbf{x}) \cdot W_{\max}$, $C' = L_{\ell_2} \cdot C_1 \cdot W_{\max} \cdot D_{\max}^2$ and $C'' = L_{\ell_2} \cdot C_2 \cdot W_{\max} \cdot D_{\max} \cdot x_{\max}$.
We thus have

$$\begin{aligned} E\Psi(\mathbf{Z}) &\leq \sqrt{2\pi} \cdot C \cdot ld \frac{1}{\sqrt{n}} + \sqrt{2\pi} \cdot C' \cdot d^2 \frac{1}{\sqrt{n}} + \sqrt{2\pi} \cdot C'' \cdot d \frac{1}{\sqrt{n}} \\ &\leq \sqrt{2\pi} \cdot C \cdot ld \frac{1}{\sqrt{n}} + \sqrt{2\pi} \cdot C''' \cdot d^2 \frac{1}{\sqrt{n}}, \end{aligned}$$

with $C''' = C' + C''$.

D.4.4. Final McDiarmid/Slepian bound

Finally, the bound we get using the McDiarmid/Slepian technique is the following.
With probability at least $1 - \delta$,

$$\Psi(\mathbf{Z}) \leq \sqrt{2\pi} \cdot C \cdot ld \frac{1}{\sqrt{n}} + \sqrt{2\pi} \cdot C''' \cdot d^2 \frac{1}{\sqrt{n}} + b \sqrt{\frac{\ln(1/\delta)}{2n}}. \quad (\text{B3})$$

D.5. Comparison of the two techniques

In this section, we compare the bounds obtained with the Covering/Lipschitz (namely (B1) and (B2)) and with the McDiarmid/Slepian (B3) techniques. For this, we introduce constants C_1 to C_6 in order to get cleaner bounds. First, using the *Covering/Lipschitz* technique, we get :

With probability at least $1 - \delta$,

$$\Psi(\mathbf{Z}) \leq \left[C_4 \sqrt{h \cdot \ln(C_5 \cdot L_c)} \right] \cdot \left(\frac{1 + \sqrt{\ln(n)}}{\sqrt{n}} \right) + C_6 \sqrt{\frac{\ln(\frac{2}{\delta})}{n}}$$

This bound :

- can depend on the ambient dimension m ,
- has suboptimal order $\mathcal{O}(\sqrt{\frac{\ln(n)}{n}})$ in the number of samples,
- does directly take into account the structure of the hypothesis class \mathcal{H} via its covering dimension h .

Now let us look at the second technique, namely the *McDiarmid/Slepian* technique, for which we got :

With probability at least $1 - \delta$,

$$\Psi(\mathbf{Z}) \leq (C_1 \cdot ld + C_2 \cdot d^2) \cdot \frac{1}{\sqrt{n}} + C_3 \sqrt{\frac{\ln(\frac{1}{\delta})}{n}}$$

This bound :

- does not depend on the ambient dimension m ,
- has optimal order $\mathcal{O}(\frac{1}{\sqrt{n}})$ in the number of samples,
- does not directly take into account the structure of the hypothesis class \mathcal{H} .

In conclusion, the two techniques give bounds quite different in nature, allowing to take into account, or not, some characteristics of the sparse model at hand.

D.6. Stable fast encoders

We relied on a stability property of the encoder with respect to the dictionary to state our main result. This property is fulfilled by the two main encoding strategies used for process-centric sparse modeling. The first one is an encoding strategy based on several iterations of a proximal optimization algorithm, as recently introduced in [Sprechmann et al. \(2015\)](#) in the process-centric view of sparse modeling. The encoder of [Fawzi et al. \(2014\)](#) is a special case of this encoding strategy. The second strategy amounts to perform encoding by solving an optimization problem, namely the elastic-net, as is used in [Mairal et al. \(2012\)](#).

Bibliographie

- ([Agarwal et al., 2013](#)) Agarwal, A., Anandkumar, A., Jain, P., Netrapalli, P., et Tandon, R. (2013). Learning sparsely used overcomplete dictionaries via alternating minimization. *CoRR*, abs/1310.7991.
- ([Aharon et al., 2006](#)) Aharon, M., Elad, M., et Bruckstein, A. (2006). K-SVD : An algorithm for designing overcomplete dictionaries for sparse representation. *Signal Processing, IEEE Transactions on*, 54(11) :4311–4322.
- ([Arora et al., 2013](#)) Arora, S., Bhaskara, A., Ge, R., et Ma, T. (2013). Provable bounds for learning some deep representations. *CoRR*, abs/1310.6343.
- ([Bagnell et Bradley, 2008](#)) Bagnell, J. A. et Bradley, D. M. (2008). Differentiable sparse coding. In Koller, D., Schuurmans, D., Bengio, Y., et Bottou, L., editors, *Advances in Neural Information Processing Systems 21*, pages 113–120.
- ([Beck et Teboulle, 2009](#)) Beck, A. et Teboulle, M. (2009). A fast iterative shrinkage-thresholding algorithm for linear inverse problems. *SIAM J. Img. Sci.*, 2(1) :183–202.
- ([Bengio et Courville, 2016](#)) Bengio, I. G. Y. et Courville, A. (2016). Deep learning. Book in preparation for MIT Press.
- ([Bengio et al., 2013](#)) Bengio, Y., Courville, A., et Vincent, P. (2013). Representation learning : A review and new perspectives. *IEEE transactions on pattern analysis and machine intelligence*, 35(8) :1798–1828.
- ([Bengio et al., 2007](#)) Bengio, Y., Lamblin, P., Popovici, D., et Larochelle, H. (2007). Greedy layer-wise training of deep networks. *Advances in neural information processing systems*, 19 :153.

- (Beylkin et al., 1991) Beylkin, G., Coifman, R., et Rokhlin, V. (1991). Fast wavelet transforms and numerical algorithms i. *Communications on pure and applied mathematics*, 44(2) :141–183.
- (Bishop, 2006) Bishop, C. M. (2006). *Pattern Recognition and Machine Learning (Information Science and Statistics)*. Springer-Verlag New York, Inc., Secaucus, NJ, USA.
- (Blumensath et Davies, 2010) Blumensath, T. et Davies, M. (2010). Normalized iterative hard thresholding : Guaranteed stability and performance. *Selected Topics in Signal Processing, IEEE Journal of*, 4(2) :298–309.
- (Blumensath et Davies, 2008) Blumensath, T. et Davies, M. E. (2008). Iterative thresholding for sparse approximations. *Journal of Fourier Analysis and Applications*, 14(5-6) :629–654.
- (Bolte et al., 2014) Bolte, J., Sabach, S., et Teboulle, M. (2014). Proximal alternating linearized minimization for nonconvex and nonsmooth problems. *Mathematical Programming*, 146(1-2) :459–494.
- (Bondy et Murty, 1976) Bondy, J. A. et Murty, U. S. R. (1976). *Graph theory with applications*, volume 290. Elsevier Science Ltd/North-Holland.
- (Bottou, 1998) Bottou, L. (1998). Online learning and stochastic approximations. *On-line learning in neural networks*, 17(9) :142.
- (Bousquet et Bottou, 2008) Bousquet, O. et Bottou, L. (2008). The tradeoffs of large scale learning. In *Advances in neural information processing systems*, pages 161–168.
- (Boyd et Vandenberghe, 2004) Boyd, S. et Vandenberghe, L. (2004). *Convex Optimization*. Cambridge University Press.
- (Bruckstein et al., 2009) Bruckstein, A. M., Donoho, D. L., et Elad, M. (2009). From Sparse Solutions of Systems of Equations to Sparse Modeling of Signals and Images. *SIAM Review*, 51(1) :34–81.
- (Candès et al., 2007) Candès, E., Demanet, L., et Ying, L. (2007). Fast computation of fourier integral operators. *SIAM Journal on Scientific Computing*, 29(6) :2464–2493.
- (Candès et al., 2006) Candès, E. J., Romberg, J., et Tao, T. (2006). Robust uncertainty principles : Exact signal reconstruction from highly incomplete frequency information. *IEEE Transactions on information theory*, 52(2) :489–509.

- (Candès et al., 2013) Candès, E. J., Strohmer, T., et Voroninski, V. (2013). Phaselift : Exact and stable signal recovery from magnitude measurements via convex programming. *Communications on Pure and Applied Mathematics*, 66(8) :1241–1274.
- (Cao et al., 2011) Cao, G., Bachega, L., et Bouman, C. (2011). The sparse matrix transform for covariance estimation and analysis of high dimensional signals. *Image Processing, IEEE Transactions on*, 20(3) :625–640.
- (Chabiron et al., 2014) Chabiron, O., Malgouyres, F., Tourneret, J.-Y., et Dobigeon, N. (2014). Toward fast transform learning. *International Journal of Computer Vision*, pages 1–22.
- (Chen et al., 1977) Chen, W.-H., Smith, C., et Fralick, S. (1977). A fast computational algorithm for the discrete cosine transform. *Communications, IEEE Transactions on*, 25(9) :1004–1009.
- (Choudhary et Mitra, 2014) Choudhary, S. et Mitra, U. (2014). Identifiability scaling laws in bilinear inverse problems. *CoRR*, abs/1402.2637.
- (Cooley et Tukey, 1965) Cooley, J. et Tukey, J. (1965). An algorithm for the machine calculation of complex Fourier series. *Mathematics of Computation*, 19(90) :297–301.
- (Cosse et Demanet, 2015) Cosse, A. et Demanet, L. (2015). Rank-one matrix completion is solved by the sum-of-squares relaxation of order two. In *Computational Advances in Multi-Sensor Adaptive Processing (CAMSAP), 2015 IEEE 6th International Workshop on*, pages 9–12. IEEE.
- (Daubechies et al., 2004) Daubechies, I., Defrise, M., et De Mol, C. (2004). An iterative thresholding algorithm for linear inverse problems with a sparsity constraint. *Communications on Pure and Applied Mathematics*, 57(11) :1413–1457.
- (Demarthon et al., 2013) Demarthon, F., Delbecq, D., et Fléchet, G. (2013). The big data revolution. *Le journal, CNRS*, 28(28) :20–27.
- (Donoho, 2006) Donoho, D. L. (2006). Compressed sensing. *IEEE Transactions on information theory*, 52(4) :1289–1306.
- (Donoho et al., 2012) Donoho, D. L., Tsaig, Y., Drori, I., et Starck, J.-L. (2012). Sparse solution of underdetermined systems of linear equations by stagewise orthogonal matching pursuit. *IEEE Transactions on Information Theory*, 58(2) :1094–1121.
- (Eckart et Young, 1936) Eckart, C. et Young, G. (1936). The approximation of one matrix by another of lower rank. *Psychometrika*, 1(3) :211–218.

- (Elad et Aharon, 2006) Elad, M. et Aharon, M. (2006). Image denoising via sparse and redundant representations over learned dictionaries. *Image Processing, IEEE Transactions on*, 15(12) :3736–3745.
- (Engan et al., 1999) Engan, K., Aase, S., et Hakon Husoy, J. (1999). Method of optimal directions for frame design. In *Acoustics, Speech, and Signal Processing, 1999. Proceedings., 1999 IEEE International Conference on*, volume 5, pages 2443–2446 vol.5.
- (Escande et Weiss, 2015) Escande, P. et Weiss, P. (2015). Real-time ℓ^1 - ℓ^2 deblurring using wavelet expansions of operators. *arXiv preprint arXiv :1512.08401*.
- (Fawzi et al., 2014) Fawzi, A., Davies, M., et Frossard, P. (2014). Dictionary learning for fast classification based on soft-thresholding. *International Journal of Computer Vision*, pages 1–16.
- (Fourier, 1822) Fourier, J. (1822). *Théorie analytique de la chaleur*.
- (Geng et al., 2011) Geng, Q., Wang, H., et Wright, J. (2011). On the local correctness of ℓ^1 minimization for dictionary learning. *CoRR*, abs/1101.5672.
- (Givens, 1958) Givens, W. (1958). Computation of plane unitary rotations transforming a general matrix to triangular form. *Journal of the Society for Industrial and Applied Mathematics*, 6(1) :pp. 26–50.
- (Golub et Van der Vorst, 2000) Golub, G. H. et Van der Vorst, H. A. (2000). Eigenvalue computation in the 20th century. *Journal of Computational and Applied Mathematics*, 123(1) :35–65.
- (Golub et Van Loan, 2012) Golub, G. H. et Van Loan, C. F. (2012). *Matrix computations*, volume 3. JHU Press.
- (Gramfort et al., 2012) Gramfort, A., Kowalski, M., et Hämmäläinen, M. S. (2012). Mixed-norm estimates for the M/EEG inverse problem using accelerated gradient methods. *Phys. Med. Biol.*, 57(7) :1937–1961.
- (Gramfort et al., 2014) Gramfort, A., Luessi, M., Larson, E., Engemann, D. A., Strohmeier, D., Brodbeck, C., Parkkonen, L., et Hämmäläinen, M. S. (2014). MNE software for processing MEG and EEG data . *NeuroImage*, 86(0) :446 – 460.
- (Gregor et LeCun, 2010) Gregor, K. et LeCun, Y. (2010). Learning fast approximations of sparse coding. In *Proceedings of the 27th Annual International Conference on Machine Learning, ICML '10*, pages 399–406.

- (Gribonval et al., 2015) Gribonval, R., Jenatton, R., Bach, F., Kleinstenber, M., et Seibert, M. (2015). Sample complexity of dictionary learning and other matrix factorizations. *Information Theory, IEEE Transactions on*, 61(6) :3469–3486.
- (Gribonval et al., 2014) Gribonval, R., Jenatton, R., et Bach, F. R. (2014). Sparse and spurious : dictionary learning with noise and outliers. *CoRR*, abs/1407.5155.
- (Gribonval et Schnass, 2010) Gribonval, R. et Schnass, K. (2010). Dictionary identification - sparse matrix-factorisation via ℓ_1 -minimisation. *IEEE Transactions on Information Theory*, 56(7) :3523–3539.
- (Guennebaud et al., 2010) Guennebaud, G., Jacob, B., et al. (2010). Eigen v3. <http://eigen.tuxfamily.org>.
- (Hackbusch, 1999) Hackbusch, W. (1999). A Sparse Matrix Arithmetic Based on H-matrices. Part I : Introduction to H-matrices. *Computing*, 62(2) :89–108.
- (Hammond et al., 2011) Hammond, D. K., Vandergheynst, P., et Gribonval, R. (2011). Wavelets on graphs via spectral graph theory. *Applied and Computational Harmonic Analysis*, 30(2) :129 – 150.
- (Haufe et al., 2008) Haufe, S., Nikulin, V. V., Ziehe, A., Müller, K.-R., et Nolte, G. (2008). Combining sparsity and rotational invariance in EEG/MEG source reconstruction. *NeuroImage*, 42(2) :726–738.
- (Hinton et al., 2012) Hinton, G., Deng, L., Yu, D., Dahl, G. E., Mohamed, A.-r., Jaitly, N., Senior, A., Vanhoucke, V., Nguyen, P., Sainath, T. N., et al. (2012). Deep neural networks for acoustic modeling in speech recognition : The shared views of four research groups. *IEEE Signal Processing Magazine*, 29(6) :82–97.
- (Hinton et Salakhutdinov, 2006) Hinton, G. E. et Salakhutdinov, R. R. (2006). Reducing the dimensionality of data with neural networks. *Science*, 313(5786) :504–507.
- (Householder, 1958) Householder, A. S. (1958). Unitary triangularization of a nonsymmetric matrix. *J. ACM*, 5(4) :339–342.
- (Jacobi, 1846) Jacobi, C. G. J. (1846). Über ein leichtes verfahren, die in der theorie der säkularstörungen vorkommenden gleichungen numerisch aufzulösen. *J. reine angew. Math.*, 30 :51–94.
- (Jones et al., 2001) Jones, E., Oliphant, T., Peterson, P., et al. (2001). SciPy : Open source scientific tools for Python.

- ([Kavukcuoglu et al., 2008](#)) Kavukcuoglu, K., Ranzato, M., et LeCun, Y. (2008). Fast inference in sparse coding algorithms with applications to object recognition. Technical Report CBLL-TR-2008-12-01, Computational and Biological Learning Lab, Courant Institute, NYU.
- ([Kim et al., 2007](#)) Kim, S.-J., Koh, K., Lustig, M., Boyd, S., et Gorinevsky, D. (2007). An interior-point method for large-scale l1-regularized least squares. *Selected Topics in Signal Processing, IEEE Journal of*, 1(4) :606–617.
- ([Király et Tomioka, 2012](#)) Király, F. et Tomioka, R. (2012). A combinatorial algebraic approach for the identifiability of low-rank matrix completion. *arXiv preprint arXiv :1206.6470*.
- ([Kondor et al., 2014](#)) Kondor, R., Teneva, N., et Garg, V. (2014). Multiresolution matrix factorization. In *Proceedings of the 31st International Conference on Machine Learning (ICML-14)*, pages 1620–1628.
- ([Lawler et al., 1985](#)) Lawler, E. L., Lenstra, J. K., Kan, A. H. G. R., et Shmoys, D. B. (1985). The traveling salesman problem ; a guided tour of combinatorial optimization. *Queueing Systems - Theory and Applications*.
- ([Le Magoarou et Gribonval, 2015a](#)) Le Magoarou, L. et Gribonval, R. (2015a). Chasing butterflies : In search of efficient dictionaries. In *2015 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pages 3287–3291.
- ([Le Magoarou et Gribonval, 2015b](#)) Le Magoarou, L. et Gribonval, R. (2015b). Flexible multi-layer sparse approximations of matrices and applications. *CoRR*, abs/1506.07300.
- ([Le Magoarou et Gribonval, 2016a](#)) Le Magoarou, L. et Gribonval, R. (2016a). Are there approximate fast fourier transforms on graphs? In *2016 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pages 4811–4815.
- ([Le Magoarou et Gribonval, 2016b](#)) Le Magoarou, L. et Gribonval, R. (2016b). Flexible multilayer sparse approximations of matrices and applications. *IEEE Journal of Selected Topics in Signal Processing*, 10(4) :688–700.
- ([Le Magoarou et al., 2015](#)) Le Magoarou, L., Gribonval, R., et Gramfort, A. (2015). FA μ ST : speeding up linear transforms for tractable inverse problems. In *23rd European Signal Processing Conference (EUSIPCO), 2015*, pages 2516–2520, Nice, France.

- ([LeCun et al., 2015](#)) LeCun, Y., Bengio, Y., et Hinton, G. (2015). Deep learning. *Nature*, 521(7553) :436–444.
- ([Lecun et al., 1998](#)) Lecun, Y., Bottou, L., Bengio, Y., et Haffner, P. (1998). Gradient-based learning applied to document recognition. *Proceedings of the IEEE*, 86(11) :2278–2324.
- ([Lee et al., 2008](#)) Lee, A. B., Nadler, B., et Wasserman, L. (2008). Treelets - an adaptive multi-scale basis for sparse unordered data. *The Annals of Applied Statistics*, 2(2) :435–471.
- ([Li et al., 2015](#)) Li, Y., Lee, K., et Bresler, Y. (2015). Identifiability in blind deconvolution with subspace or sparsity constraints.
- ([Lokam, 2009](#)) Lokam, S. V. (2009). Complexity lower bounds using linear algebra. *Foundations and Trends in Theoretical Computer Science*, 4(1-2) :1–155.
- ([Lyu et Wang, 2013](#)) Lyu, S. et Wang, X. (2013). On algorithms for sparse multi-factor NMF. In *Advances in Neural Information Processing Systems 26*, pages 602–610.
- ([Mairal et al., 2012](#)) Mairal, J., Bach, F., et Ponce, J. (2012). Task-driven dictionary learning. *Pattern Analysis and Machine Intelligence, IEEE Transactions on*, 34(4) :791–804.
- ([Mairal et al., 2010](#)) Mairal, J., Bach, F., Ponce, J., et Sapiro, G. (2010). Online learning for matrix factorization and sparse coding. *Journal of Machine Learning Research*, 11(1) :19–60.
- ([Mairal et al., 2008](#)) Mairal, J., Bach, F., Ponce, J., Sapiro, G., et Zisserman, A. (2008). Discriminative learned dictionaries for local image analysis. In *Computer Vision and Pattern Recognition, 2008. CVPR 2008. IEEE Conference on*, pages 1–8.
- ([Mairal et al., 2009](#)) Mairal, J., Ponce, J., Sapiro, G., Zisserman, A., et Bach, F. R. (2009). Supervised dictionary learning. In Koller, D., Schuurmans, D., Bengio, Y., et Bottou, L., editors, *Advances in Neural Information Processing Systems 21*, pages 1033–1040. Curran Associates, Inc.
- ([Makhzani et Frey, 2013](#)) Makhzani, A. et Frey, B. (2013). k-sparse autoencoders. *CoRR*, abs/1312.5663.
- ([Malgouyres et Landsberg, 2016](#)) Malgouyres, F. et Landsberg, J. (2016). On the identifiability and stable recovery of deep/multi-layer structured matrix factorization.

- (Mallat, 1989) Mallat, S. (1989). A theory for multiresolution signal decomposition : the wavelet representation. *Pattern Analysis and Machine Intelligence, IEEE Transactions on*, 11(7) :674–693.
- (Mallat, 2008) Mallat, S. (2008). *A Wavelet Tour of Signal Processing, Third Edition : The Sparse Way*. Academic Press, 3rd edition.
- (Mallat et Zhang, 1993) Mallat, S. et Zhang, Z. (1993). Matching pursuits with time-frequency dictionaries. *Signal Processing, IEEE Transactions on*, 41(12) :3397–3415.
- (Maurer et Pontil, 2010) Maurer, A. et Pontil, M. (2010). K-Dimensional Coding Schemes in Hilbert Spaces. *Information Theory, IEEE Transactions on*, 56(11) :5839–5846.
- (Mehta et Gray, 2013) Mehta, N. et Gray, A. G. (2013). Sparsity-based generalization bounds for predictive sparse coding. In Dasgupta, S. et Mcallester, D., editors, *Proceedings of the 30th International Conference on Machine Learning (ICML-13)*, volume 28, pages 36–44. JMLR Workshop and Conference Proceedings.
- (Mises et Pollaczek-Geiringer, 1929) Mises, R. et Pollaczek-Geiringer, H. (1929). Praktische verfahren der gleichungsauflösung. *ZAMM-Journal of Applied Mathematics and Mechanics/Zeitschrift für Angewandte Mathematik und Mechanik*, 9(2) :152–164.
- (Morgenstern, 1973) Morgenstern, J. (1973). Note on a lower bound on the linear complexity of the fast fourier transform. *J. ACM*, 20(2) :305–306.
- (Morgenstern, 1975) Morgenstern, J. (1975). The linear complexity of computation. *J. ACM*, 22(2) :184–194.
- (Needell et Tropp, 2009) Needell, D. et Tropp, J. (2009). Cosamp : Iterative signal recovery from incomplete and inaccurate samples. *Applied and Computational Harmonic Analysis*, 26(3) :301 – 321.
- (Nesterov, 2012) Nesterov, Y. (2012). Efficiency of coordinate descent methods on huge-scale optimization problems. *SIAM Journal on Optimization*, 22(2) :341–362.
- (Neyshabur et Panigrahy, 2013) Neyshabur, B. et Panigrahy, R. (2013). Sparse matrix factorization. *CoRR*, abs/1311.3315.
- (Olshausen et Field, 1996) Olshausen, B. A. et Field, D. J. (1996). Emergence of simple-cell receptive field properties by learning a sparse code for natural images. *Nature*, 381(6583) :607–609.

- (Olshausen et Field, 1997) Olshausen, B. A. et Field, D. J. (1997). Sparse coding with an overcomplete basis set : A strategy employed by v1? *Vision Research*, 37(23) :3311 – 3325.
- (Papadimitriou, 1977) Papadimitriou, C. H. (1977). The euclidean travelling salesman problem is np-complete. *Theoretical Computer Science*, 4(3) :237–244.
- (Parikh et Boyd, 2014) Parikh, N. et Boyd, S. P. (2014). Proximal algorithms. *Foundations and Trends in optimization*, 1(3) :127–239.
- (Perraudin et al., 2014) Perraudin, N., Paratte, J., Shuman, D., Kalofolias, V., Vandergheynst, P., et Hammond, D. K. (2014). GSPBOX : A toolbox for signal processing on graphs. *ArXiv e-prints*.
- (Petersen et Pedersen, 2012) Petersen, K. B. et Pedersen, M. S. (2012). The matrix cookbook. Version 20121115.
- (Pudlák, 2000) Pudlák, P. (2000). A note on the use of determinant for proving lower bounds on the size of linear circuits. *Inf. Process. Lett.*, 74(5-6) :197–201.
- (Rao et Yip, 1990) Rao, K. R. et Yip, P. (1990). *Discrete Cosine Transform : Algorithms, Advantages, Applications*. Academic Press Professional, Inc., San Diego, CA, USA.
- (Ritter, 1996) Ritter, T. (1996). Walsh-hadamard transforms : A literature survey. <http://www.ciphersbyritter.com/RES/WALHAD.HTM>.
- (Rockafellar, 2015) Rockafellar, R. T. (2015). *Convex analysis*. Princeton university press.
- (Rokhlin, 1985) Rokhlin, V. (1985). Rapid solution of integral equations of classical potential theory. *Journal of Computational Physics*, 60(2) :187 – 207.
- (Rubinstein et al., 2010a) Rubinstein, R., Bruckstein, A., et Elad, M. (2010a). Dictionaries for Sparse Representation Modeling. *Proceedings of the IEEE*, 98(6) :1045 –1057.
- (Rubinstein et al., 2008) Rubinstein, R., Zibulevsky, M., et Elad, M. (2008). Efficient Implementation of the K-SVD Algorithm using Batch Orthogonal Matching Pursuit. Technical report.
- (Rubinstein et al., 2010b) Rubinstein, R., Zibulevsky, M., et Elad, M. (2010b). Double sparsity : Learning sparse dictionaries for sparse signal approximation. *Signal Processing, IEEE Transactions on*, 58(3) :1553–1564.

- ([Rustamov et Guibas, 2013](#)) Rustamov, R. et Guibas, L. (2013). Wavelets on graphs via deep learning. In *Advances in Neural Information Processing Systems 26*, pages 998–1006. Curran Associates, Inc.
- ([Rusu et al., 2016](#)) Rusu, C., Gonzalez-Prelcic, N., et Heath, R. W. J. (2016). Fast orthonormal sparsifying transforms based on householder reflectors. *Submitted to IEEE Transactions on Signal Processing*.
- ([Sainath et al., 2013](#)) Sainath, T., Kingsbury, B., Sindhvani, V., Arisoy, E., et Ramabhadran, B. (2013). Low-rank matrix factorization for deep neural network training with high-dimensional output targets. In *Acoustics, Speech and Signal Processing (ICASSP), 2013 IEEE International Conference on*, pages 6655–6659.
- ([Sandryhaila et Moura, 2013](#)) Sandryhaila, A. et Moura, J. (2013). Discrete signal processing on graphs. *Signal Processing, IEEE Transactions on*, 61(7) :1644–656.
- ([Shanks, 1969](#)) Shanks, J. (1969). Computation of the fast walsh-fourier transform. *Computers, IEEE Transactions on*, C-18(5) :457–459.
- ([Shuman et al., 2013](#)) Shuman, D. I., Narang, S. K., Frossard, P., Ortega, A., et Vandergheynst, P. (2013). The emerging field of signal processing on graphs : Extending high-dimensional data analysis to networks and other irregular domains. *Signal Processing Magazine, IEEE*, 30(3) :83–98.
- ([Sprechmann et al., 2015](#)) Sprechmann, P., Bronstein, A., et Sapiro, G. (2015). Learning efficient sparse and low rank models. *Pattern Analysis and Machine Intelligence, IEEE Transactions on*, 37(9) :1821–1833.
- ([Strang, 2005](#)) Strang, G. (2005). *Linear Algebra and Its Applications, 4th Edition*. Brooks Cole.
- ([Sturmfels, 2002](#)) Sturmfels, B. (2002). *Solving systems of polynomial equations*. Number 97. American Mathematical Soc.
- ([Tibshirani, 1996](#)) Tibshirani, R. (1996). Regression shrinkage and selection via the lasso. *Journal of the Royal Statistical Society (Series B)*, 58 :267–288.
- ([Tosic et Frossard, 2011](#)) Tosic, I. et Frossard, P. (2011). Dictionary learning. *Signal Processing Magazine, IEEE*, 28(2) :27–38.
- ([Tremblay, 2014](#)) Tremblay, N. (2014). *Networks and signal : signal processing tools for network analysis*. Theses, Ecole normale supérieure de lyon - ENS LYON.

- ([Tropp et Gilbert, 2007](#)) Tropp, J. et Gilbert, A. (2007). Signal recovery from random measurements via orthogonal matching pursuit. *Information Theory, IEEE Transactions on*, 53(12) :4655–4666.
- ([Tropp et Wright, 2010](#)) Tropp, J. et Wright, S. (2010). Computational methods for sparse solution of linear inverse problems. *Proceedings of the IEEE*, 98(6) :948–958.
- ([Vainsencher et al., 2011](#)) Vainsencher, D., Mannor, S., et Bruckstein, A. M. (2011). The sample complexity of dictionary learning. *The Journal of Machine Learning Research*, 12 :3259–3281.
- ([Valiant, 1977](#)) Valiant, L. G. (1977). Graph-theoretic arguments in low-level complexity. In Gruska, J., editor, *Mathematical Foundations of Computer Science 1977*, volume 53 of *Lecture Notes in Computer Science*, pages 162–176. Springer Berlin Heidelberg.
- ([Vapnik, 1991](#)) Vapnik, V. (1991). Principles of risk minimization for learning theory. In *NIPS*, pages 831–838.
- ([Watkins, 2007](#)) Watkins, D. S. (2007). *The matrix eigenvalue problem : GR and Krylov subspace methods*. Siam.
- ([Wright, 2015](#)) Wright, S. J. (2015). Coordinate descent algorithms. *Mathematical Programming*, 151(1) :3–34.
- ([Xue et al., 2013](#)) Xue, J., Li, J., et Gong, Y. (2013). Restructuring of deep neural network acoustic models with singular value decomposition. In *INTERSPEECH*, pages 2365–2369.

AVIS DU JURY SUR LA REPRODUCTION DE LA THESE SOUTENUE

Titre de la thèse:

Matrices efficaces pour le traitement du signal et l'apprentissage automatique

Nom Prénom de l'auteur : LE MAGOAROU LUC

Membres du jury :

- Madame GUILLEMOT Christine
- Monsieur GRIBONVAL Rémi
- Monsieur MALGOUYRES François
- Monsieur PEYRE Gabriel
- Madame SCHNASS Karin
- Monsieur VANDERGHEYNST Pierre

Président du jury : *Prof. Pierre VANDERGHEYNST*

Date de la soutenance : 24 Novembre 2016

Reproduction de la these soutenue

Thèse pouvant être reproduite en l'état

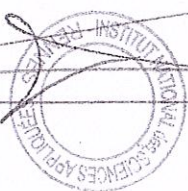
~~Thèse pouvant être reproduite après corrections suggérées~~

Fait à Rennes, le 24 Novembre 2016

Signature du président de jury

Le Directeur,

M'hamed DRISSI



A handwritten signature in black ink, enclosed in a large, loopy oval shape.

Résumé

Les matrices, en tant que représentations des applications linéaires en dimension finie, jouent un rôle central en traitement du signal et des images et en apprentissage automatique. L'application d'une matrice de rang plein à un vecteur implique a priori un nombre d'opérations arithmétiques de l'ordre du nombre d'entrées non-nulles que contient la matrice. Cependant, il existe des matrices pouvant être appliquées bien plus rapidement, cette propriété étant d'ailleurs un des fondements du succès de certaines transformations linéaires, telles que la transformée de Fourier ou la transformée en ondelettes.

Quelle est cette propriété ? Est-elle vérifiable aisément ? Peut-on approcher des matrices quelconques par des matrices ayant cette propriété ? Peut-on estimer des matrices ayant cette propriété ?

La thèse s'attaque à ces questions en explorant des applications telles que l'apprentissage de dictionnaires efficients, l'accélération des itérations d'algorithmes de résolution de problèmes inverses, ou l'analyse de Fourier rapide sur graphe.

Abstract

Matrices, as natural representation of linear mappings in finite dimension, play a crucial role in signal processing and machine learning. Multiplying a vector by a full rank matrix a priori costs of the order of the number of non-zero entries in the matrix, in terms of arithmetic operations. However, matrices exist that can be applied much faster, this property being crucial to the success of certain linear transformations, such as the Fourier transform or the wavelet transform.

What is the property that allows these matrices to be applied rapidly ? Is it easy to verify ? Can we approximate matrices with ones having this property ? Can we estimate matrices having this property ?

This thesis investigates these questions, exploring applications such as learning dictionaries with efficient implementations, accelerating the resolution of inverse problems or Fast Fourier Transform on graphs.

