# Signal separation in convolutive mixtures : contributions to blind separation of sparse sources and adaptive subtraction of seismic multiples

Yves-Marie Batany

# THÈSE DE DOCTORAT

de l'Université de recherche Paris Sciences et Lettres
PSL Research University

**Préparée dans le cadre d'une cotutelle entre
MINES ParisTech et
Universidade Estadual de Campinas (UNICAMP)**

Signal separation in convolutive mixtures: contributions to blind separation of sparse sources and adaptive subtraction of seismic multiples

Séparation de signaux en mélanges convolutifs : contributions à la séparation aveugle de sources parcimonieuses et à la soustraction adaptative des réflexions multiples en sismique

### École doctorale n°398

GÉOSCIENCES, RESSOURCES NATURELLES ET ENVIRONNEMENT

**Spécialité** GÉOSCIENCES ET GÉOINGÉNIERIE

Soutenue par **Yves-Marie BATANY**
le 14 novembre 2016

Dirigée par **Hervé CHAURIS**
et **João Marcos T. ROMANO**

**COMPOSITION DU JURY :**

M. Christian JUTTEN
Université Joseph Fourier     président

M. Yannick BERTHOUMIEU
Université de Bordeaux     rapporteur

M. Laurent DUVAL
IFPEN     examinateur

M. Antonio PICA
CGG     examinateur

M. Hervé CHAURIS
MINES ParisTech     examinateur

Mme Daniela DONNO
MINES ParisTech     examinatrice

M. Leonardo T. DUARTE
UNICAMP     examinateur

M. João Marcos T. ROMANO
UNICAMP     examinateur

Yves-Marie Batany

# Separação de sinais em misturas convolutivas: contribuições para a separação cega de fontes esparsas e em subtração adaptativa de reflexões múltiplas em sísmica

# Signal separation in convolutive mixtures: contributions to blind separation of sparse sources and adaptive subtraction of seismic multiples

Campinas

2016

UNIVERSIDADE ESTADUAL DE CAMPINAS

Faculdade de Engenharia Elétrica e de Computação

Yves-Marie Batany

# Separação de sinais em misturas convolutivas: contribuições para a separação cega de fontes esparsas e em subtração adaptativa de reflexões múltiplas em sísmica

# Signal separation in convolutive mixtures: contributions to blind separation of sparse sources and adaptive subtraction of seismic multiples

Tese apresentada à Faculdade de Engenharia Elétrica e de Computação da Universidade Estadual de Campinas como parte dos requisitos exigidos para a obtenção do título de Doutor em Engenharia Elétrica, na Área de Telecomunicações e Telemática.

Thesis presented to the Faculty of Electrical and Computer Engineering of the University of Campinas in partial fulfillment of the requirements for the degree of Doctor, in the area of Telecommunications and Telematics.

Orientador: Prof. Dr. João Marcos Travassos Romano
Co-orientador Prof. Dr. Leonardo Tomazeli Duarte

Este exemplar corresponde à versão final da tese defendida pelo aluno Yves-Marie Batany, e orientada pelos Professores João Marcos Travassos Romano e Leonardo Tomazeli Duarte

Campinas

2016

Ficha catalográfica
Universidade Estadual de Campinas
Biblioteca da Área de Engenharia e Arquitetura
Luciana Pietrosanto Milla - CRB 8/8129

# COMISSÃO JULGADORA – TESE DE DOUTORADO

**Candidato**: Yves-Marie Batany RA: 151587

**Data da Defesa**: 14 de novembro de 2016

**Título da Tese**: "Signal Separation in Convolutive Mixture: Contributions to Blind Separation of Sparse Sources and Adaptive Subtraction of Seismic Multiples".

Prof. Dr. João Marcos Travassos Romano (Presidente, FEEC/UNICAMP)
Prof. Dr. Hervé Chauris (MINES ParisTech)
Prof. Dr. Christian Jutten (Université Joseph Fourier)
Prof. Dr. Yannick Berthoumieu (Université de Bordeaux)
Dr. Laurent Duval (IFPEN, France)

A ata de defesa, com as respectivas assinaturas dos membros da Comissão Julgadora, encontra-se no processo de vida acadêmica do aluno.

# Acknowledgements

I'm going to start my acknowledgements with two personnal anecdotes. Both are connected, in a certain way, to my PhD.

During the end of 2012, I was listening to a radio playing some bossa nova music. The music was inspiring and I was suddenly passionated by the portuguese language. I opened an internet browser and looked for the lyrics: I could not understand a single sentence. I felt desperate thinking how sad it was that I would certainly never be able to understand the sense of these words on my own!
Later in 2013, during my internship at the Schlumberger research center in Cambridge, I start thinking about doing a PhD. The reasons to do it were not clear yet, but mainly it was because most people that had a job I liked had done a PhD. At that time, I was completely fascinated by sparsity. It was a hot topic in both signal processing and geophysics. So I started looking for a proposal. I was more keen on finding something in France. One morning, a colleague (I should say a friend) arrived at work saying: "I found the perfect PhD for you! It's in Paris and it has the word *sparse* in the proposal!" I took a look into the PhD proposal and I realized that the subject was proposed by my Master's professors. In my mind it was clear: the topic was great and I really enjoyed the lectures I had with them. Suddenly, I also realized that the project was a double degree with a Brazilian university: that is what we call an opportunity!

Today, I have finalized my PhD, I can speak portuguese and enjoy the lyrics from brazilian music. I would like to thank the Universe for filling the gaps in such a unexpected way. The unpredictable is what keeps me moving forward. It is the foundation of my faith that tomorrow will bring something new.

But of course, the Universe is not the only pillar of this achievement. Four people have been present from the very beginning in this PhD, even before I start it, and I am deeply grateful for every piece of knowledge and the time they dedicated to me.

I would like to thank **Daniela** for her constant support. Thanks for never letting me down, even after deadlines or after one of my long periods of silence (yes, I may have a problem with that... and I'm working on it).
I would like to thank **Leonardo** for all the passionate discussions. I could not imagine when I first arrived in Limeira that a friendship would start. Thank you for making me discover so many things.
I would like to thank **João** for making this collaboration happen. Without you, FEEC would miss its spirit. I have good memories of our conversations during coffee time.
I would like to thank **Hervé** for fighting on my side to make science. There aren't that many directors who dive into the depths of the engine of codes with their student to see what's happening.
Thanks all of you for your patience.

I also would like to express my gratitude to Yannick Berthoumieu, Christian Jutten, Laurent Duval and Antonio Pica for reviewing my work and improving the final manuscript with their comments.

# Abstract [English]

The recovery of correlated signals from their linear combinations is a challenging task and has many applications in signal processing. We focus on two problems that are the blind separation of sparse sources and the adaptive subtraction of multiple events in seismic processing. A special focus is put on convolutive mixtures: for both problems, finite impulse response filters can indeed be estimated for the recovery of the desired signals.

For instantaneous and convolutive mixing models, we address the necessary and sufficient conditions for the exact extraction and separation of sparse sources by using the $\ell_0$ pseudo-norm as a contrast function. Equivalences between sparse component analysis and disjoint component analysis are investigated.

For adaptive multiple subtraction, we discuss the limits of methods based on independent component analysis and we highlight equivalence with $\ell_p$-norm-based methods. We investigate how other regularization parameters may have more influence on the estimation of the desired primaries. Finally, we propose to improve the robustness of adaptive subtraction by estimating the extracting convolutive filters directly in the curvelet domain. Computation and memory costs are limited by using the uniform discrete curvelet transform.

**Keywords:**   seismic processing • seismic multiples • adaptive filtering • blind source separation • sparse component analysis • curvelet transform.

# Résumé [Français]

La séparation de signaux corrélés à partir de leurs combinaisons linéaires est une tâche difficile et possède plusieurs applications en traitement du signal. Nous étudions deux problèmes, à savoir la séparation aveugle de sources parcimonieuses et le filtrage adaptatif des réflexions multiples en acquisition sismique. Un intérêt particulier est porté sur les mélanges convolutifs : pour ces deux problèmes, des filtres à réponses impulsionnelles finies peuvent être estimés afin de récupérer les signaux désirés.

Pour les modèles de mélange instantanés et convolutifs, nous donnons les conditions nécessaires et suffisantes pour l'extraction et la séparation exactes de sources parcimonieuses en utilisant la pseudo-norme $\ell_0$ comme une fonction de contraste. Des équivalences entre l'analyse en composantes parcimonieuses et l'analyse en composantes disjointes sont examinées.

Pour la soustraction adaptative des réflexions multiples en sismiques, nous discutons les limites des méthodes basées sur l'analyse en composantes indépendantes et nous soulignons l'équivalence avec les méthodes basées sur les normes $\ell_p$. Nous examinons de quelle manière les paramètres de régularisation peuvent être plus décisifs pour l'estimation des primaires. Enfin, nous proposons une amélioration de la robustesse de la soustraction adaptative en estimant les filtres adaptatifs directement dans le domaine des curvelets. Les coûts en calcul et en mémoire peuvent être atténués par l'utilisation de la transformée en curvelet discrète et uniforme.

**Mots Clés :**   traitement sismique ● réflexions multiples ● filtrage adaptatif ● séparation aveugle de sources ● analyse en composantes parcimonieuses ● transformée en curvelet.

# Resumo [Português]

A separação de sinais correlacionados a partir de suas combinações lineares é uma tarefa difícil e tem diversas aplicações diferentes em processamento de sinais. Nesta tese, estudamos dois problemas, que são a separação cega de fontes esparsas e a filtragem adaptativa de reflexões múltiplas em aquisição sísmica. Um interesse particular é dado às misturas convolutivas. Para esses dois problemas, filtros de resposta ao impulso finita podem ser estimados para recuperar os sinais desejados.

Para os modelos instantâneos e convolutivos, apresentamos as condições necessárias e suficientes para a extração e a separação exatas de fontes esparsas usando a pseudo-norma $\ell_0$ como uma função de contraste. Equivalência entre a análise de componentes esparsas e a análise de componentes disjuntas são examinadas.

Para a subtração adaptativa de reflexões múltiplas em sísmica, discutimos os limites de métodos baseados em análise de componentes independentes e realçamos equivalências entre os métodos baseados em normas $\ell_p$. Investigamos de qual maneira os parâmetros de regularização podem ser mais decisivos para a estimação das primárias. Finalmente, propomos uma melhora da robustez da subtração adaptativa com a estimação de filtros adaptativos diretamente no domínio de curvelets. Os custos de cálculo e de memória podem ser atenuados com o uso da transformada curvelet discreta e uniforme.

**Palavras-chave:** processamento sísmico • reflexões múltiplas • filtragem adaptativa • separação cega de fontes • análise de componentes esparsas • transformada curvelet.

# List of symbols and abbreviations

| | |
|---|---|
| $\boldsymbol{A}$ | the mixing system |
| $\mathcal{C}\boldsymbol{s}[\mu]$ | the curvelet transform of a signal $\boldsymbol{s}$ |
| $\boldsymbol{d}[x]$ | an observation vector $\in \mathbb{R}^M$ |
| $f_i(s_i)$ | the PDF of a source (signal) $\boldsymbol{s}_i$ |
| $F_i(s_i)$ | the CDF of a source (signal) $\boldsymbol{s}_i$ |
| $\mathcal{F}\boldsymbol{s}[\omega]$ | the Fourier transform of a signal $\boldsymbol{s}$ |
| $\mathcal{F}^{-1}$ | the inverse Fourier transform |
| $g_i$ | a non-linear function |
| $\boldsymbol{H}$ | the mapping system |
| $h_f$ | the differential entropy associated to the PDF $f$ |
| $I(\mathrm{s}_1, \mathrm{s}_2)$ | the mutual information between $\boldsymbol{s_1}$ and $\boldsymbol{s_2}$ |
| $\boldsymbol{m}[x]$ | the predicted multiples |
| $M$ | the number of observations |
| $N$ | the number of sources |
| $\boldsymbol{p}[x]$ | the primaries |
| $Q_s$ | the negentropy of s |
| $\mathcal{R}\boldsymbol{s}[t, q]$ | the Radon transform of a signal $\boldsymbol{s}$ |
| $\boldsymbol{s}[x]$ | a source (signal) vector $\in \mathbb{R}^N$ |
| s | a random variable |
| $\tilde{s}$ | an outcome of a random variable s |
| $\bar{\mathrm{s}}$ | a centralized random variable |
| $\breve{\mathrm{s}}$ | a normalized random variable |
| $\breve{\bar{\mathrm{s}}}$ | a standardized random variable |
| $\mathcal{S}^{\downarrow}$ | a down-sampling operator |
| $\boldsymbol{u}_i$ | an extracted vector $\boldsymbol{u} = \boldsymbol{w}^T \boldsymbol{X}$ |
| $\boldsymbol{W}$ | the separating system |
| $x$ | index position |
| $y, z$ | index position (used for convolution) |
| $\boldsymbol{\Delta}$ | a diagonal matrix |
| $\kappa_n$ | the n-th cumulent |
| $\mu_n$ | the n-th moment |
| $\boldsymbol{\Pi}$ | a permutation matrix |

| | |
|---|---|
| $\odot$ | the Hadamard product |
| $\otimes$ | the correlation product |
| $*$ | the convolution product |
| $\|\boldsymbol{s}\|_p$ | the $\ell_p$-norm of a vector $\boldsymbol{s}$ |
| $\boldsymbol{\nabla}\phi$ | the gradient of an objective function $\phi$ |

| | |
|---|---|
| BSE | bind source extraction |
| BSS | blind source separation |
| CDF | cumulative density function |
| CMP | common mid-point gather |
| COG | common offset gather |
| CSG | common shot gather |
| DCA | disjoint component analysis |
| DCT | discret curvelet transform |
| DE | differential evolution |
| DO | disjoint orthogonal |
| FIR | finite impulse response |
| ICA | independent component analysis |
| MIMO | multiple-input and multiple-output |
| MISO | multiple-input and single-output |
| NMO | normal move-out |
| PCA | principal component analysis |
| PDF | probability density function |
| SCA | sparse component analysis |
| SL0 | smooth $\ell_0$ |
| UDCT | uniform discrete curvelet transform |

# Contents

# Part I

# Opening

# Chapter 1

# Introduction

## Contents

## Résumé du chapitre [français]

Ce premier chapitre introduit le contexte, les objectifs et les contributions de la thèse. L'exploration sismique consiste à estimer les paramètres de la sub-surface de la Terre solide en analysant sa réponse à des vibrations. Les données enregistrées à la surface sont introduites an sein d'un problème inverse afin d'estimer les paramètres physiques des structures en sub-surface (tels que la vitesse de propagation des ondes, l'atténuation. . . ).

La plupart des méthodes font l'hypothèse que les données enregistrées ne contiennent que des événements primaires n'ayant été réfléchis qu'une seule fois vers la surface. Lors d'une acquisition en mer, la surface libre entre l'eau et l'air agit comme un miroir sur lequel les ondes se réfléchissent. Il en résulte dans les données des événements cohérents, appelés "réflexions multiples", qui doivent être soigneusement enlevés avant de pouvoir passer au problème inverse.

Il existe plusieurs méthodes pour atténuer les multiples. Une partie d'entre elles consiste à prédire les réflexions multiples et à les soustraire des données. Cependant, les prédictions ne sont jamais parfaites et une étape de soustraction adaptative doit être mise en place. Si cette étape n'est pas faite avec soin, on prend le risque d'atténuer les primaires, c'est-à-dire le signal pertinent. Récemment, il a été proposé d'utiliser des méthodes venant du problème de séparation aveugle de sources pour effectuer le filtrage adaptatif des multiples. Ces méthodes intéressantes viennent compléter les méthodes plus traditionnelles basées sur l'usage de normes. Elles font l'hypothèse que les primaires et les multiples sont statistiquement indépendantes.

La séparation aveugle de source est un problème d'une grande occurrence. L'analyse en composante indépendante s'est développée pour résoudre ce problème. Néanmoins, lorsque les sources ne peuvent pas être modélisées comme des variables indépendantes, on doit se tourner vers d'autres solutions. L'analyse en composante parcimonieuse fait l'hypothèse que les sources peuvent être représentées par un nombre restreint de coefficients.

Les contributions de la thèse sont les suivantes. Dans un premier temps, nous présentons et discutons une condition nécessaire et suffisante pour l'extraction de source parcimonieuse dans des mélanges instantanés et convolutifs. Dans un second temps, nous donnons une unification de plusieurs méthodes de soustraction adaptative permettant une meilleur comparaison de celles-ci. Finalement, nous proposons une méthode pour calculer un filtre à réponse impulsionnelle dans le domaine des curvelets avec un coût de calcul maîtrisé.

## Resumo do capítulo [português]

Esse primeiro capitulo introduz o contexto, os objetivos e as contribuições desta tese de doutorado. A exploração sísmica consiste em estimar os parâmetros da subsuperfícíe da Terra sólida analisando a sua resposta a vibrações. Os dados registrados da superfície são introduzidos a um problema inverso com o objetivo de estimar os parâmetros físicos das estruturas da subsuperficíe, como a velocidade da propagação de ondas, a atenuação, etc.

A maioria dos métodos são baseados na hipótese de que os dados registrados contêm unicamente eventos primários tendo sido refletidos apenas uma vez em direção à superfície. Durante uma aquisição marítima, a superfície livre entre a água e o ar funciona como um espelho sobre o qual as ondas se refletem. Esses eventos coerentes, chamados de reflexões múltiplas, devem ser cuidadosamente removidos previamente a resolução do problema inverso em questão.

A atenuação das múltiplas pode ser feita através de vários métodos. Uma parte deles consiste em prever e subtrair as reflexões múltiplas. No entanto, as previsões nunca são perfeitas e uma etapa de subtração adaptativa deve ser estabelecida. Se esta etapa não é feita com cuidado, corre-se o risco de atenuar as primárias, ou seja o sinal pertinente. Foi proposta recentemente a possibilidade da utilização dos métodos provenientes do problema de separação cega de fontes para efetuar a filtragem adaptativa das múltiplas. Esses métodos interessantes vieram completar os métodos mais tradicionais baseados no uso de normas. Eles são baseados na hipótese de que as primarias e as múltiplas são estatisticamente independentes.

A separação cega de fontes é um problema de alta ocorrência. A análise de componentes independentes foi desenvolvida para resolver esse problema. No entanto, enquanto as fontes não podem ser modelizadas como variáveis independentes, devemos nos voltar a outras soluções. A análise de componentes esparsos é baseada na hipótese de que as fontes podem ser representadas por um numero restrito de coeficientes.

As contribuições desta tese de doutorado são as seguintes. Inicialmente será apresentada e discutida uma condição necessária e suficiente para a extração de fontes esparsas nos casos de misturas instantâneas e convolutivas. Em seguida, será apresentada uma unificação de vários métodos de subtração adaptativa permitindo uma melhor comparação entre eles. Por fim, um método será proposto para calcular um filtro de resposta ao impulso no domínio das curvelets com um custo de cálculo controlado.

## 1.1   Motivations and contextualization

### 1.1.1   Seismic exploration

Seismic exploration intends to retrieve the subsurface parameters (such as velocities, densities, attenuation,...)  controlling the wave propagation inside the solid Earth [Kearey et al., 2002]. This is of particular interest for academic geosciences, oil and gas industry, geothermal activities or $C0_2$ storage. Seismic exploration is a noninvasive exploration technique based on the inverse problem theory that helps to understand rock properties via wave propagation analysis [Tarantola, 1987]. In a traditional survey, seismic waves are excited from the surface and propagate inside the subsurface. Because of contrast in acoustic impedance, part of the energy is reflected back to the surface and the wave-field response can be recorded at the surface [Aki and Richards, 2002].

Marine acquisitions are distinct from land acquisitions [Sheriff and Geldart, 1995]. The recording device is called a hydrophone in the former and a geophone in the latter. It measures, at a given location, the displacement or the acceleration as a function of time. A seismic source[1] is used to create the signal propagating inside the solid Earth. In 3D modern acquisitions, several parallel lines of receivers are used in the field. Figure 1.1 shows a 2D layered Earth model and several ray paths between the seismic source (magenta star) and a receiver (magenta triangle). Ray paths represent the main trajectories along which seismic energy propagates. For a 2D acquisition when sources and receivers are located at the surface and aligned, the seismic data are sorted in a cube comporting three fundamental coordinates: the source position, the offset which is the distance between the source and the receiver, and the time. The final seismic image is a 2D section of the solid Earth in depth. In 3D acquisition, the source position and the offset are two dimensional vectors and the final seismic image is a 3D volume of the solid Earth in depth. The processing from the data recorded in time to real depth (the seismic image) relies on the considered velocity model.

By nature, seismic exploration relies on underdetermined inverse problems, for which infinite numbers of solutions exist [Tarantola, 1987]. The main reasons of this underdeterminacy are related to the surface acquisition, the limited data and the uneven illumination [Xie et al., 2006]. The art of seismic exploration lives between understanding the forward modeling of waves propagation, engineering the recording system for the pertinent physical variables and properly inverting the seismic data. At one side, the forward model gives what the observation should be given a specific model. On the other side, the inversion methods try to reduce the number of solutions by using regularization techniques or a proper processing of the data.

One of the objectives of the present thesis aims at improving the signal to noise ratio of seismic images by attenuating one specific kind of reverberant coherent noise called *multiples*. This typical noise occurring in marine acquisition must indeed be properly removed for several reasons.

### 1.1.2   Multiple events

Seismic events in data sets are classified depending on their nature. For instance we distinguish body waves from surface waves, P-waves from S-waves, Rayleigh waves from Love waves, or refracted waves from reflection waves [Aki and Richards, 2002]. Another crucial distinction is made between what is called primary events (or primaries) and multiple events (or multiples). A multiple event has been reflected downward one or

---

[1]The terminology *seismic source* designates the particular waveform generated by the device (e.g. dynamite or water-gun).

**Figure 1.1:** *Geometry of a seismic exploration survey with a representation of a layered solid Earth. The magenta star indicates the seismic source position and the magenta triangles indicate the receiver positions. The blue line indicates a primary event ray path, the red lines indicate several surface related multiple ray paths, the green line represents an internal multiple ray path.*

more times into the Earth by a strong reflector. Red lines in figure 1.1 shows several ray paths corresponding to multiple events. Figure 1.2 shows how the multiples look like in a real common shot gather. It is of prime importance to process carefully this energy because imaging artifacts may appear during the sequence of imaging steps from the raw data to the final seismic image [Wiggins, 1988; Verschuur, 2013b]. A strong reflector exists because of a strong contrast of impedance (product of density and velocity) between two layers. This can occur between two layers made of different rocks. For instance in marine acquisition, the water free surface between the water layer and the air acts almost as a mirror and reflects back all the upcoming energy. We can write the recorded data $\boldsymbol{d}[x]$ as

$$\boldsymbol{d}[x] = \boldsymbol{p_0}[x] + \boldsymbol{m_0}[x], \tag{1.1}$$

where $\boldsymbol{p_0}[x]$ contains the energy of the primaries and $\boldsymbol{m_0}[x]$ contains the energy of the multiples. The variable vector $x$ indicates the position in the data cube.

Multiple events are themselves categorized depending on the reflector they relate to or on their signature in the data [Verschuur, 2013a]. Those categories may overlap and a single multiple event may belong to more than one category. Section 4.2 will be dedicated to a more precise overview of multiples, however it is good to have an idea of the kind of multiples that will be treated in the present work. Specifically, we will treat the multiples for which a prediction can be obtained.

Source ghosts and receiver ghosts are due to the free surface close to the seismic source or the receiver, respectively [Verschuur, 2013a]. They must be removed as they create a duplicated shifted version of any seismic event. They are not always considered as multiples.

Surface-related multiples refer to the presence of one strong reflector [Verschuur, 2013b]. All multiples that are removed if this reflector disappears are said to be related to this surface. If the strong reflector is the free surface, one speaks about free-surface-related multiples. Internal multiples are due to one or two strong reflectors

**Figure 1.2:** *a) A common shot gather from a marine seismic acquisition in which several orders of multiples can be identified. b) A prediction of the multiples.*

below the free surface [Yilmaz, 2001]. They are often referred to as short period multiples because they create close duplicated events. Intrabed multiples refer to reflections created inside a single layer. Interbed multiples refer to reflections created between two different layers. Pegleg multiples may refer to different kind of multiples. Surface related multiples can be predicted from the data (see also section 4.4).

Depending on the context, the term $\boldsymbol{m_0}[x]$ in equation 1.1 may contains only part of the entire multiple energy. For instance, we may consider only the water surface related multiples and in that case, internal multiples are contained in the term $\boldsymbol{p_0}[x]$. In a general sense, the term $\boldsymbol{p_0}[x]$ is considered as the signal and the term $\boldsymbol{m_0}[x]$ as the noise.

As we will see in the next section, most of the common imaging steps in the processing chain consider that the data are free from multiples. This is the reason why multiples has been usually considered as coherent noise in seismic imaging. It is worth mentioning that there is a change of paradigm in the community for using multiples as a signal. However, those methods still have some limits and multiples cannot be used in all steps of the process [Berkhout, 2016; Weglein, 2015] (see also section 4.6).

### 1.1.3 Basic production workflow of seismic images

A classical workflow from the raw data to the final seismic image consists of various steps [Yilmaz, 2001; Robein, 2010]. It is worth mentioning that most of them are included inside iterative loops: the decision maker quality controls and adapts each step during the process. The accuracy of each step determines the accuracy of the final image. We emphasize hereafter only a few key steps in order to highlight the effect of multiple removal in some key steps [Sheriff and Geldart, 1995]. The set of traces associated to one source is called a common shot gather (CSG). The set of

traces associated to one pair of source-receiver with a fixed offset distance is called a common offset gather (COG). The set of traces associated to the physical middle point between the source and the receiver is called a common mid-point gather (CMP).

**Pre-processing and editing.** Acquisition starts in the field and one differentiates marine acquisitions from land acquisitions. In both, the data are never acquired perfectly and a few pre-processing steps are needed, such as trace edition, basic source signature deconvolution, geometry correction, gain correction and interpolation to end up with spatially coherent and regularly sampled data.

**Stacking.** Stacking is one of the main reason why we acquire multi-offset data. It is the core of a lot of process and quality control checking point. In CMP gathers, primary events are supposed to follow a normal move-out (NMO) travel time. Stacking basically refers to the summation of several traces in CMP gathers, in order to check or enhance coherency and to allow both velocity analysis and noise reduction. Other processes can be said to be pre-stack or post-stack, depending if they are done before or after stacking.

**Migration.** Migration is the linearization of forward modeling, in which only primaries must be considered [Weglein, 2015]. It aims at determining the optimal model perturbation (i.e. the reflectivity) to match the observed reflected data. In this context, multiples create additional events that will be treated as primaries [Mulder and ten Kroode, 2002; Li and Symes, 2007].

**Tomography.** Beyond migration results, the background velocity model controls the kinematics of wave propagation. Tomography methods for retrieving a correct velocity model can be developed in data domain (travel-time tomography [Bishop et al., 1985], waveform inversion [Virieux and Operto, 2009]) or in image domain (migration velocity analysis [Sava and Biondi, 2004]). As the velocity usually increases with depth in the solid Earth, multiples generally present lower apparent velocity than primaries. A wrong interpretation of multiples leads to a lower velocity estimation of the subsurface model.

**Full waveform inversion.** Differently from migration and tomography, full waveform inversion (FWI) does not make distinction between large scales (obtained tomography) and fine scales (obtained by migration) [Virieux and Operto, 2009]. FWI is the ultimate technique for retrieving all scales at once. In principle, multiples are integrated into the FWI framework. In practice they create local minima and slow down the convergence performance, enforcing the need of a good initial velocity model [Brossier et al., 2009].

### 1.1.4 Multiple removal methods and adaptive multiple subtraction

As we explained previously, multiples are considered as a noise for most imaging steps. Various methods exist in order to attenuate or remove them, and are generally divided into two categories [Verschuur, 2013a]. The first class contains filtering techniques. Those methods are mainly based on some transform domains in which the primaries and the multiples are located in different area. Hence, multiples can be muted in their corresponding area. The second class contains multiple removal methods based on a prediction of the multiples. Generally, the prediction is done via a feedback model, basically considering that multiples are echos of the primaries [Weglein et al., 1997].

This feedback model, as used in the surface related multiple elimination (SRME) technique [Verschuur et al., 1992], is able to give a good prediction of the multiples by considering the recorded data only, hence avoiding the need of a velocity model.

The predictions generated by the feed-back model approaches are never perfect [Abma et al., 2005]. Amplitude and phase errors always remain, mainly due to acquisition limitation and complex Earth structure. Hence, the prediction methods require an adaptive subtraction step to adapt the prediction of the multiples to the data [Verschuur and Berkhout, 1997; Rickett et al., 2001; Guitton and Verschuur, 2004]. At this condition only, the multiples can be efficiently removed from the data. Adaptive subtraction has developed as an active research area and is considered as the main challenge in multiple removal. In summary, adaptive subtraction looks for a local convolutive filter that is applied to the imperfect predicted multiples to better match the true multiples.

Two main approaches exist for adaptive subtraction namely matching filter techniques [Verschuur and Berkhout, 1997; Guitton and Verschuur, 2004] and pattern recognition methods [Spitz, 2000; Guitton et al., 2001]. Several issues remain in adaptive subtraction, in particular if the signal and the noise are somehow correlated. For matching filter techniques, if the primaries and the multiples overlap, the minimum energy assumption fails and the method tends to over-estimate the multiple and destroys part of the primary events. For pattern recognition approach, multiples are removed if they share the same pattern as the primaries. Concerning matching filters, several different metics have been investigated to tackle this issue.

Adaptive subtraction methods tend to have several crucial parameters to be tuned. Those parameters can be linked to the objective function, the window size or the filter size. The estimated primaries are highly dependent to the choice of those parameters meaning that the methods are not robust enough. The choice of finding a good set of parameters is generally left to the operator processing the data.

Figure 1.3 shows a synthetic example with one primary event and one multiple event, clearly overlapping. In this example the prediction has the correct amplitude and the correct wavelet, but a wrong time shift. The adaptive subtraction consists of recovering the correct position of the multiple before its subtraction from the data. A common approach is to minimize the $\ell_2$-norm of the primaries defined as $\sqrt{\sum_x (p_x)^2}$ [Verschuur and Berkhout, 1997]. Figure 1.4 shows the objective functional, function of two parameters (amplitude and time shift). The minimum does not correspond to the exact recovery of the true primary event. A common alternative is to use the $\ell_1$-norm defined as $\sum_x |p_x|$ [Guitton and Verschuur, 2004]. Figure 1.5 shows that this approach also leads to a wrong primary estimate. Figure 1.6 shows that an approach based on an approximation of the $\ell_0$ pseudo-norm[2], further developed in the next sections, could lead to a better estimate.

In a few recent works [Kaplan and Innanen, 2008; Donno, 2011; Liu and Dragoset, 2013], it has been proposed to analyze and perform adaptive subtraction as a blind source separation (BSS) problem. The data are the sum of primary and multiple events while the prediction is a filtered version of the multiples. The forward filters between the true multiples and the prediction are unknown and the BSS framework is potentially suited.

---

[2]In this thesis, $\ell_0$ is refereed to as "pseudo-norm" because it is not a norm. This terminlogy is often use in the literature in this context for designing $\ell_0$. However, strictly speaking, $\ell_0$ is not a pseudo-norm. We refer the interested reader to appendix 8.1 for more details.

**Figure 1.3:** *Synthetic example of adaptive multiple subtraction in a local window with a) one primary event and b) one multiple event. The data c) are the sum of the signal and the noise. The prediction d) gives a correct estimate of the multiple, except for a time shift.*



**Figure 1.4:** *Adaptive multiple subtraction performed on the synthetic example of figure 1.3. Two parameters have to be optimized: time shift and amplitude. a) The $\ell_2$-norm objective function. b) Objective function along the white dashed line. c) Primary estimated by the $\ell_2$-norm optimization.*

**Figure 1.5:** *Same as figure 1.4 but with the $\ell_1$ norm.*



**Figure 1.6:** *Same as figure 1.4 but with an approximation of the $\ell_0$ pseudo-norm.*

### 1.1.5 Blind source separation and independent component analysis

Blind source separation (BSS) is a general and widely referenced problem occurring in a lot of science areas [Hyvärinen et al., 2001; Comon and Jutten, 2010]. In its simplest formulation, a BSS problem considers that the observations are combinations of the original signals. Those original signals that one wishes to recover are commonly referred to as the sources in the BSS community. This terminology should not be confused with the term seismic sources as we used and defined before. Solving a BSS problem requires the recovery of the original signal from the observations without any assumption on the mixing process. The term *blind* refers indeed to the fact that no assumption is made on the mixing process. Figure 1.7 shows an example of a simple blind source separation problem with two sources (original signals) and two mixtures. Each mixture is an unknown linear combinations of the two sources. Based on the assumption that the sources are statistically independent, one can build an objective function (also named a contrast function) for which the local minima correspond to the recovery (the separation) of the original sources [Comon and Jutten, 2010]. Table 1.1 summarizes some equivalences of terms used in the geophysics and BSS communities.

Independent component analysis (ICA) has emerged as a powerful tool for solving BSS problems [Comon, 1994; Hyvärinen et al., 2001]. The main hypothesis of ICA is that the original signals (the sources) are statistically independent. Two random variables are independent if the knowledge about one does not change the knowledge about the other. Hence, statistical independence can be seen as a generalization of correlation, which is a measure of the linear statistical dependence. ICA-based separation algorithms try to find outputs that are as much independent as possible. Measuring quantitatively the independence between variables is not trivial and several measures and algorithms have been developed such as FastICA [Hyvärinen, 1999], Infomax [Bell and Sejnowski, 1995] or JADE [Cardoso, 1999]. Some of these ICA-based methods have been used for adaptive subtraction of multiple events [Liu and Dragoset, 2013].

As we indicated before, one of the main issue in adaptive subtraction is the presence of overlapping events that may be locally correlated. In that case, the assumption for ICA-based method is not valid anymore and other methods must be investigated. Sparse component analysis (SCA) have shown capacities for dealing with this particular issue [Deville, 2014]. The second part (chapters 2 and 3) of the present thesis is dedicated to SCA.

### 1.1.6 Sparse component analysis

Sparsity is a key concept in a lot of domains such as signal processing, data mining, compression or inverse problems in general. It is a powerful concept for regularization purpose for reducing the number of solutions or for finding a more realistic solution. A signal (for instance an image or a seismic data cube) is said to be sparse in a given representation (or a dictionary) if a small number of coefficients is sufficient to explain most of it. The most common measure of sparsity is the $\ell_0$ pseudo-norm counting the number of non-zero coefficients of a given vector [Hurley and Rickard, 2009].

With modern dense acquisition, seismic events can be seen as band-limited local plane waves. Hence, seismic data are sparse in the FK domain, in which band-limited plane waves are represented by straight lines. Also, multi-scale transforms have shown to be efficient for representing signal in a sparse manner [Mallat, 1998]. In particular, for seismic purpose, the curvelet transform has shown its ability to sparsely represent the data [Candès and Demanet, 2005; Herrmann and Moghaddam, 2004].

Recently, a huge interest on sparse representation appeared after the work of Candès

**Figure 1.7:** *Example of a BSS problem. The two original sources are pages of a book. Each observation is a linear superposition of the original pages. The desired separated sources are recovered with a permutation ambiguity. No assumption on the mixing system is required for separating the sources.*

| BLIND SOURCE SEPARATION | EXPLORATION GEOPHYSICS |
|---|---|
| Source | Signal or parameters |
| Mixing system | Seismic source |
| Contrast function | Objective (or cost) function |
| Separation or extraction | Parameter (or signal) estimation |

**Table 1.1:** *Short list of particular terms used in blind source separation and geophysics, in order to avoid confusion in the rest of the thesis.*

et al. [2006b], showing that a sparse prior can be sufficient for solving under-determined problems. These works had a huge impact and a new research field, named compressive sensing, has emerged from it [Candès and Wakin, 2008]. Naturally, sparseness also arises in the field of blind source separation. Bofill and Zibulevsky [2001] have seen that the mixing system can be easily identified if the original signals are sparse, and a lot of works have followed constituting what is called sparse component analysis (SCA) [Bofill and Zibulevsky, 2001; Li et al., 2003a; Georgiev et al., 2005; Deville, 2014].

The objective of the present thesis is the removal of multiple events and more specifically the improvement of the adaptive subtraction step indispensable on prediction based methods. As we said before, adaptive subtraction can be formulated in a blind source separation framework for which independent component analysis is an accomplished direction for solutions. However, a main challenge in adaptive subtraction is the presence of correlated primaries and multiples for which the hypothesis of statistical independence fails. Sparse component analysis has emerged as a tool for solving BSS problems when the source of interest are sparse in a given dictionary. We see it as a path for improving adaptive subtraction. Figure 1.6 shows the objective function based on the $\ell_0$ pseudo-norm (smoothed) for the example presented in figure 1.4. We see that the objective function is not convex any-longer, but the global minimum corresponds to the correct recovery of the primary event.

Ultimately in this thesis, we will use the BSS framework in its convolutive form, with signal sparsity in the curvelet domain. To our knowledge, a few results exist on the conditions under which SCA is able to separate sparse signals in convolutive mixtures.

## 1.2 Contributions

The present thesis has been motivated by the adaptive multiple subtraction problem. In particular, the methods based on independent component analysis (ICA) have been investigated in more details than before. In parallel, sparse component analysis (SCA) became really attractive from a theoretical perspective. Independent contributions have been given in both theoretical and applied areas but we intend to make connections between them. The present work is motivated by the following questions :

- How far can SCA be used for separating correlated signals?
- In which manner do ICA methods really improve adaptive filtering?
- Is the curvelet transform limited to amplitude recovery?

**Theoretical contribution for sparse component analysis**

Our first contribution concerns sparse component analysis (SCA). We provide a necessary and sufficient condition to the use of the $\ell_0$-norm as a contrast function in blind source extraction and separation. In other terms, we provide the condition under which minimizing the $\ell_0$ pseudo-norm leads to a perfect signal recovery. Also,

we discuss the similarity between SCA and another method called disjoint component analysis (DCA) for the blind separation of sparse sources. To do so, we expand auto-regressive process to inter-regressive process. From that, we are able to show that for SCA also, correlated events can be an issue. However, this assertion is true for one really specific kind of strong linear dependency (namely an inter-regressive process).

### A better view of adaptive subtraction methods

Lately, some works proposed to use ICA as a new tool for adaptive multiple subtraction. In particular, it has been written that those techniques can handle the separation of correlated events. We present a coherent analysis of all the adaptive subtraction methods that clearly show that this assertion is not true. We show that all those methods tend to minimize a non-linear correlation between the signal and the noise. We emphasize the crucial role of regularization parameter such as filter size and window size in the performance of adaptive filtering.

### Improving robustness of adaptive subtraction with curvelets

Curvelet transform has been used for adaptive multiple subtraction in the amplitude recovery model only, and not in a convolutive form. We discuss a convolutional theorem for discrete curvelet transform from which we extract a way to compute FIR filtering by limiting the computational time. The method makes use of the theoretical work about the recovery condition with SCA and it gives a better overview of the convolution theorem with curvelets.

These contributions are communicated in the following list of articles:

- **Published article:**
  Y.-M. Batany, L. Tomazeli Duarte, D. Donno, J. M. T. Romano, and H. Chauris. Adaptive multiple subtraction: Unification and comparison of matching filters based on the lq-norm and statistical independence. *Geophysics*, 81(1):V43–V54, 2016.

- **Conferences:**
  Y.-M. Batany, D. Donno, L. Tomazeli Duarte, H. Chauris, and J. M. T. Romano. A necessary and sufficient condition for the blind extraction of the sparsest source in convolutive mixtures. In *European Signal Processing Conference (EUSIPCO)*, 2016.

  Y.-M. Batany, L. Tomazeli Duarte, D. Donno, J. M. T. Romano, and H. Chauris. Comparison of matching filters for adaptive multiple subtraction: Lq-norm versus statistical independence. In *78th EAGE Conference and Exhibition*, 2016.

- **Book chapter:**
  L. T. Duarte, Y.-M. Batany, and J. M. T. Romano. Blind source separation: principles of independent and sparse component analysis. In CRC Press, editor, Signals and Images: Advances and results in speech, estimation, compression, recognition, filtering and processing, chapter 1. 2015.

- **In preparation:**
  About the equivalence between SCA and DCA for blind source separation. *IEEE Transactions on Signal processing.*

## 1.3 Outline of the thesis

As explained before, the present thesis aims at improving adaptive subtraction by gathering different methods in a common framework and ultimately developing adaptive subtraction fully in the curvelet domain. Blind source separation represents a solid theoretical foundation for this purpose and sparse component analysis is a excellent starting point from a theoretical perspective.

The thesis is divided into four parts. Part I (chapter 1) contains the present introduction. Part II (chapters 2 and 3) is dedicated to blind source separation methods and specifically to methods promoting the sparsity of the signals. Chapter 2 is a general introduction on inverse problem, blind source separation and independent component analysis. Chapter 3 focuses on sparse component analysis and contains our results on the necessary and sufficient conditions for a correct recovery in convolutive mixtures.

Part III (chapters 4, 5 and 6) contains our advances on adaptive subtraction methods for prediction based multiple removal. Chapter 4 is a more precise and detailed introduction on multiple elimination methods. Chapter 5 presents our analysis thus allowing to gather all methods on a common framework. Chapter 6 details our approach of convolutive adaptive subtraction fully in the curvelet domain. Finally, the last part IV (chapter 7) draws perspectives and conclusions.

# Part II

# Blind source separation and sparsity

# Chapter 2

# Prelude on blind source (signal) separation

## Contents

Some of the ideas developed in this chapter 2 have been published in an introductory chapter of a book [Duarte et al., 2015].

## Résumé du chapitre [français]

Le chapitre 2 présente une introduction à la séparation aveugle de sources (SAS). Les concepts qui y sont développés seront utilisés tout au long du manuscrit. Dans un problème de SAS, les signaux observés (les observations) sont des combinaisons inconnues de signaux à estimer (les sources). Dans ce contexte, le terme de source ne doit pas être confondu avec le terme de "source sismique" utilisé en géophysique. Les méthodes de résolution de problème de SAS se distinguent principalement par les hypothèses faites sur les sources ou sur le modèle de mélange.

La section 2.2 présente des généralités sur les problèmes inverses et sur les méthodes d'optimisation mathématique. Les méthodes de régularisation de problèmes mal-posés sont aussi évoquées, ainsi que les problèmes multi-objectifs avec l'analyse de courbe de Pareto.

La section 2.3 discute les concepts liés aux statistiques d'ordre supérieur. Les statistiques du second ordre se limitent essentiellement à l'utilisation de la moyenne et de la variance des signaux traités. Elles sont généralement insuffisantes pour résoudre des problèmes de SAS. D'autres caractéristiques dérivées des densités de probabilités doivent êtres prises en compte pour une description adéquate des signaux.

La section 2.4 présente en détails les problèmes d'extraction et de séparation aveugle de sources. Nous nous intéressons ici aux mélanges instantanés et convolutifs, deux cas particuliers de problème linéaire. Dans un cas déterminé, l'extraction et la séparation peuvent être effectuées en adaptant un système linéaire, inverse du système de mélange. L'optimisation de ce système est faite par optimisation d'une fonction dite de contraste.

La section 2.5 introduit l'analyse en composantes indépendantes comme méthode de résolution de problèmes de séparation aveugle. En effet, il a été montré que l'hypothèse d'indépendance statistique des signaux originaux est une hypothèse suffisante pour pouvoir séparer les sources.

Finalement, la section 2.6 introduit l'analyse en composantes parcimonieuses comme méthode de résolution de problèmes de séparation aveugle. Les résultats importants de la littérature y sont présentés, en particulier les conditions suffisantes pour que les sources puissent être retrouvées.

## Resumo do capítulo [português]

O capítulo 2 apresenta uma introdução à separação cega de fontes (SCF). Os conceitos que aqui são desenvolvidos serão úteis ao longo de todo o manuscrito. Em um problema de SCF, os sinais observados (as observações) são combinações desconhecidas de sinais a serem estimados (as fontes). Dentro desse contexto, o termo "fonte" não deve ser confundido com o termo "fonte sísmica" utilizado em geofísica. Os métodos de resolução de problemas de SCF se distinguem principalmente pelas hipóteses feitas sobre as fontes ou sobre o modelo de mistura.

A seção 2.2 apresenta generalidades sobre os problemas inversos e sobre os métodos de otimização matemática. Os métodos de regularização de problemas mal postos são também evocados, assim como os problemas multi-objetivos com a análise da curva de Pareto.

A seção 2.3 discute os conceitos ligados às estatísticas de ordem superior. As estatísticas de segunda ordem limitam-se essencialmente à utilização da média e da variância dos sinais tratados. Elas são geralmente insuficientes para resolver problemas SCF. Outras características derivadas das densidades das probabilidades devem ser levadas em conta para uma descrição adequada de sinal.

A seção 2.4 apresenta em detalhes os problemas de extração e de separação cega de fontes. Nós nos interessamos nessa seção pelo caso de misturas instantâneas e convolutivas, dois casos particulares de problema linear. Em um caso determinado, a extração e a separação podem ser efetuadas adaptando um sistema linear, inverso do sistema de mistura. A otimização desse sistema é feita através da otimização de uma função dita "de contraste".

A seção 2.5 introduz a análise de componentes independentes como método de resolução de problemas de separação cega. Foi mostrado que, de fato, a hipótese de independência estatística de sinais originais é uma hipótese suficiente para poder separar as fontes.

Por fim, a seção 2.6 introduz a análise de componentes esparsos como método de resolução de problema de separação cega. Os resultados importantes da literatura são aqui apresentados, particularmente as condições suficientes para que as fontes possam ser reencontradas.

## 2.1    Introduction

The present chapter aims at introducing the blind source[1] separation (BSS) problem and the key concepts around independent component analysis (ICA) and sparse component analysis (SCA). Basically, ICA and SCA are two different tools for solving BSS problems. ICA is based on the assumption that the sources are statistically independent while SCA is based on the assumption that the sources are sparse. The inverse theory in presented in a general setting in section 2.2, along with optimization and regularization techniques. Useful definitions and concepts of the statistical theory are introduced in section 2.3. Section 2.4 is dedicated to the BSS problem while sections 2.5 and 2.6 are dedicated to ICA and SCA, respectively. This chapter is the mainstay of both chapter 3, dedicated to SCA, and chapter 5, dedicated to the separation of primaries and multiples events in seismic acquisitions. Throughout the thesis, we focus on the objective functions of the different methods and the conditions assuring a perfect separation of the original sources.

## 2.2    Inverse problems and optimization

### 2.2.1    Generality on inverse problems

In a wide sense, an inverse problem consists in recovering a certain model from some physical observations (data) denoted $d$. The model is generally defined by a set of parameters $s$. The formulation of the direct problem gives an operator $\mathcal{A} : s \to d$ able to represent the mapping between the parametrized model and the observed data [Tarantola, 1987]. We write the general formulation of a inverse problem as

$$\text{find } s^* \quad \text{such that} \quad d \approx \mathcal{A}(s^*), \tag{2.1}$$

where $s^*$ represent the best solution. If the direct problem is linear, $\mathcal{A}$ can be represented by a matrix and the direct model is simply

$$\boldsymbol{d} = \boldsymbol{A}\boldsymbol{s}, \tag{2.2}$$

where $\boldsymbol{d} \in \mathbb{R}^M$ and $\boldsymbol{s} \in \mathbb{R}^N$ are two vectors and $\boldsymbol{A}$ is a $M \times N$ matrix. Equation 2.2 represents a linear system of $M$ linear equations and we are facing a linear inverse problem.

**Definition 1.** *An inverse problem is well-posed if: (i) a solution exists, (ii) the solution is unique, (iii) the solution is stable.*

If the problem is not well-posed, we say that it is ill-posed. For instance, problem 2.2 is ill-posed if the singular values of $\boldsymbol{A}$ decay gradually to zero or if the ratio between the largest and the smallest nonzero singular values is high [Hansen, 2008]. Also, if there are less equations that parameters, i.e. $M < N$, the problem is underdetermined and so ill-posed.

Figure 2.1 shows four basic examples of linear inverse problems with $N = 2$ parameters and $M$ observations. Each observation leads to one black linear constraint in the parameter space. If there are as many independent equations as parameters such that $\boldsymbol{A}$ is invertible, the problem is determined and a single solution exists at the crossing point. If there are more equations than parameters, the problem is over-determined

---

[1]We remind here again that the word "source" in blind source separation has the meaning of any original signal. For instance, sources can be temporal series or images. It must not be confused with the term "seismic source" that has a specific meaning within the geophysical community.

**Figure 2.1:** *Example of inverse problem in $\mathbb{R}^2$ with linear constraints (black lines). (a) A determined problem with two constraints. (b) A noisy overdetermined problem with four constraints. (c)-(d) An underdetermined problem with a single constraint. The $\ell_2$ and $\ell_1$ balls (in blue) gives two different solutions.*

and the solution should be at the position minimizing the distances from all linear constraints. When the number of equations is not sufficient, there exist many solutions $s^*$ such that equation 2.2 is valid. Such a problem is underdetermined and some prior information must be added to uniquely find a unique solution $s^*$. For instance, the physical informations available in a particular context can be jointed, or the sparsest solution may be preferred. As proposed in the example, we can choose the model minimizing a certain energy measurement such as the $\ell_2$ or the $\ell_1$-norm.

The problem 2.2 appears in a lot of domains and recasts into several categories, depending on the knowledge or on the size of the problem. For instance in sparse representation, $A$ represents a known overcomplete dictionary in which one wants to represent the data $d$ sparsely [Donoho and Elad, 2003]. In compressive sensing, the sensing matrix $A$ is designed to reduce the number of data in a way that the recovery of a sparse signal is possible [Baraniuk, 2007]. When both $A$ and $s$ are totally unknown, one faces a bilinear problem known as blind source separation (BSS) that will be discussed in the next sections and chapters [Comon and Jutten, 2010].

### 2.2.2 Optimization and regularization schemes

By writing $d \approx As$, we set that the modeled observations should be closed to the actual observations. The measure of the distance between the data $d$ and the modeled data $As$ is the starting point of all inverse problem solvers. A classical measure is a $\ell_q$-norm defined as

$$\|s\|_p = \left( \sum_x |s_x|^p \right)^{1/p}, \qquad p \geq 1. \tag{2.3}$$

This definition can be expanded for $p < 1$, hence defining a quasi-norm that is non-convex (see appendix 8.1). For $p = 2$ the norm is the classical Euclidean distance and for $p = 1$ the norm is the sum of the absolute values. More precisely we have

$$\|s\|_2 = \sqrt{\sum_x s_x^2} \qquad \text{and} \qquad \|s\|_1 = \sum_x |s_x|. \tag{2.4}$$

For an overdetermined linear problem, the least-square solution minimizing the $\ell_2$-norm of the residual is given by

$$s^*_{\ell_2} = (A^T A)^{-1} A^T d, \tag{2.5}$$

where $^T$ indicates the transpose operator.

When the problem is ill-posed the solution is not stable. It means a small perturbation in the observations leads to large changes in the solution. A classical approach is

|  | Data fitting term | Regularization term |
|---|---|---|
| Damped least-square | $\|\boldsymbol{d} - \boldsymbol{As}\|_2$ | $\|\boldsymbol{s}\|_2$ |
| Tikhonov | $\|\boldsymbol{d} - \boldsymbol{As}\|_2$ | $\|\boldsymbol{Bs}\|_2$ |
| LASSO | $\|\boldsymbol{d} - \boldsymbol{As}\|_2$ | $\|\boldsymbol{s}\|_1$ |
| Akaike Information Criterion | $\|\boldsymbol{d} - \boldsymbol{As}\|_2$ | $\|\boldsymbol{s}\|_0$ |
| Total variation | $\|\boldsymbol{d} - \boldsymbol{As}\|_2$ | $\|\nabla\boldsymbol{s}\|_1$ |

**Table 2.1:** *Some common data fitting and regularization terms for inverse problems.*

to add a regularization term to the data fitting term (the distance). We build an objective function of the form

$$\min \quad \phi = \|\boldsymbol{d} - \boldsymbol{As}\| + \epsilon \, \phi_r\left(\boldsymbol{s}\right), \tag{2.6}$$

where the second term penalizes solutions too far from the prior information contained in $\phi_r$. The notation $\|\cdot\|$ without specification denotes a distance in the general sense. A similar strategy is used for underdetermined problems written as

$$\min \quad \phi_r\left(\boldsymbol{s}\right) \quad \text{such that} \quad \|\boldsymbol{d} - \boldsymbol{As}\|_2 < \epsilon. \tag{2.7}$$

The parameter $\epsilon$ controls the trade-off between the fitting and the regularization terms. When $\epsilon = 0$, only the data fitting term is relevant. When $\epsilon = \infty$, only the regularization term is important. Table 2.1 presents several classical regularization schemes. For most of them, the data fitting term is the $\ell_2$-norm and only the regularization term changes. The Tikhonov regularization strategy has an analytic solution given by

$$\boldsymbol{s}^*_{\ell_2,\epsilon} = (\boldsymbol{A}^T\boldsymbol{A} + \epsilon \, \boldsymbol{B}^T\boldsymbol{B})^{-1}\boldsymbol{A}^T\boldsymbol{d}. \tag{2.8}$$

For the damped least-square solution, we have $\boldsymbol{B} = \boldsymbol{I}$. Generally, the optimization strategies do not have any analytical solution and optimization scheme must be invoked.

**Definition 2.** *A function $\phi$ is convex if*

$$\phi(\lambda\boldsymbol{s}_1 + (1 - \lambda)\boldsymbol{s}_2) \leq \lambda\phi(\boldsymbol{s}_1) + (1 - \lambda)\phi(\boldsymbol{s}_2), \tag{2.9}$$

*for any $\lambda \in [0, 1]$.*

Convexity is a useful property for an objective function, because the minimum (the solution) is located in a valley thus making the use of gradient based methods possible [Boyd and Vandenberghe, 2004]. An iterative strategy can be use, by following the opposite direction of the gradient of the objective function $\boldsymbol{\nabla}\phi$

$$\boldsymbol{s}^{(k)} \leftarrow \boldsymbol{s}^{(k-1)} - \alpha\boldsymbol{\nabla}\phi. \tag{2.10}$$

Because of convexity, the solution after convergence does not depend on the initial solution $\boldsymbol{s}^0$. If the function is non-convex, evolutionary computation can be used. For instance evolutionary algorithms based on a population of solutions can be used for solving non-convex problems [Deb, 2001].

Most of the times, the two terms in the problem 2.6 are conflicting, meaning that one decreases when the other increases. The framework of multi-objective Pareto analysis with a combination of two objectives is useful to understand the behavior of the solution and to correctly choose the tradeoff parameter $\epsilon$.

**Definition 3.** *A solution $\boldsymbol{s}_1$ dominates another $\boldsymbol{s}_2$ if $\phi_i(\boldsymbol{s}_1) < \phi_i(\boldsymbol{s}_2)$ for all objective $\phi_i$. We denote this situation $\boldsymbol{s}_1 \prec \boldsymbol{s}_2$. The Pareto frontier is the set of all non-dominated solutions.*

**Figure 2.2:** *Pareto's framework for inverse problem regularization. The parameter $\epsilon$ gives the slope of the solution line. The Pareto's frontier is not always convex.*

Figure 2.2 shows an example of an ill-posed problem with LASSO regularization scheme. In a two-objectives Pareto framework, the parameter $\epsilon$ gives the slope of the line $\|\boldsymbol{d} - \boldsymbol{As}\|_2 \propto -\epsilon \|s\|_1$, tangent to the Pareto frontier. The solution 1 gives more importance to the sparsity of the solution but does not penalize enough a large value of the residuals. The solution 3 gives more importance to the residuals, but the solution associated is not sparse enough. The solution 2 is a good compromise between minimizing the residuals and having a sparse solution.

## 2.3 Second and higher order statistics

In the previous section, we introduced the problem of parameter estimation within the inverse problem framework. In this section, we introduce a few basic statistical concepts that will be needed in the following [Cover and Thomas, 2006; T. Romano et al., 2010; Comon and Jutten, 2010].

### 2.3.1 Characterization of a single random variable

We consider a continuous random variable, denoted s, taking values in the real axis $\mathbb{R}$ [Leon-Garcia, 2008]. An outcome of s is denoted $\tilde{s}$. This random variable can be fully characterized in different manners by its cumulative density function, its probability density function, its first characteristic function or its second characteristic function. All those representations are equivalent in the sense that they uniquely represent the same random variable. In the following, the probability of a specific event $\tilde{s}$ to occur is denoted as $Pr(\tilde{s})$.

**Definition 4.** *Let* s *be a random variable. Its cumulative density function (CDF) is defined such that*

$$F(s) = Pr(\tilde{s} \leq s). \tag{2.11}$$

**Definition 5.** *Let* s *be a random variable with a CDF $F(s)$. Its probability density*

**Figure 2.3:** *(a) CDF of a generalized Gaussian variable, (b) its PDF with $a = 1$, $p = 1$ (blue) and $a = 5$, $p = 4$ (red).*

*function (PDF) is defined such that*

$$f(s) = \frac{d}{ds} F(s). \tag{2.12}$$

From the above definitions, one can show that

$$Pr(a \leq \tilde{s} \leq b) = \int_a^b f(s) ds. \tag{2.13}$$

The PDF of a random variable gives the probability of the random variable to be in a certain interval. There is a large number of useful PDFs. However, it is out of the scope of the present introduction to give a deep overview of all classes of famous PDFs. We focus on one particular class of PDFs, linked to $\ell_p$-norms defined in equation 2.3. The Gamma function is denoted $\Gamma(s)$[2].

**Definition 6.** *The generalized Gaussian distribution is defined such that*

$$f_{\mathcal{N}_p}(s) = \frac{p}{2a\Gamma(1/p)} e^{-\left(\frac{|s-\mu|}{a}\right)^p}, \tag{2.14}$$

*where $a > 0$ is a scale parameter, $p > 0$ is a shaping parameter and $\mu \in \mathbb{R}$ is a position parameter. The Gaussian distribution and the Laplacian distribution appear when $p = 2$ and $p = 1$, respectively.*

Figure 2.3 shows two examples of the generalized Gaussian distribution.

From the PDF, two useful functions can be defined, namely the first and second characteristic functions. From those two functions, two sets of coefficients describing and synthesizing in a useful way a given random variable are defined. The expectation operator with respect to the PDF $f$ is denoted $\mathbb{E}_f$.

**Definition 7.** *Let s be a random variable with a PDF $f(s)$. Its first characteristic function is defined such that*

$$\psi_1(\omega) = \mathbb{E}_f \left\{ e^{is\omega} \right\} = \int_{-\infty}^{\infty} e^{is\omega} f(s) ds. \tag{2.15}$$

The first characteristic function is the inverse Fourier transform of the PDF. Its Taylor expansion gives the definition of the *ordinary moments* $\mu_n$ of the random variable such that

$$\psi_1(\omega) = \mu_1 i\omega - \frac{\mu_2 \omega^2}{2} + \ldots \quad \Rightarrow \quad \mu_n = (-i)^n \left. \left| \frac{d^n}{d\omega^n} \psi_1(\omega) \right. \right|_{\omega=0}, \tag{2.16}$$

---

[2] The Gamma function is defined as $\Gamma(s) = \int_0^{\infty} t^{s-1} e^{-t} dt$.

where $i^2 = -1$. Equivalently, the $n$-th moment can be defined such that

$$\mu_n = \mathbb{E}_f\{s^n\} = \int_{-\infty}^{\infty} s^n f(s) ds. \tag{2.17}$$

A *central* random variable is constructed by subtracting the mean as $\bar{s} = s - \mu_1$. A *normalized* random variable is constructed by dividing by the square root of the second moment (i.e. the standard deviation) as $\breve{s} = s/\sqrt{\mu_2}$. We define the *central moments* and the *standardized moments* (or normalized central moments) as

$$\overline{\mu}_n = \mathbb{E}_f\left\{\bar{s}^n\right\} = \mathbb{E}_f\{(s - \mu_1)^n\} = \int_{-\infty}^{\infty} (s - \mu_1)^n f(s) ds, \tag{2.18}$$

$$\breve{\overline{\mu}}_n = \mathbb{E}_f\left\{\breve{s}^n\right\} = \mathbb{E}_f\left\{\left(\frac{s - \mu_1}{\sqrt{\mu_2}}\right)^n\right\} = \int_{-\infty}^{\infty} \left(\frac{s - \mu_1}{\sqrt{\mu_2}}\right)^n f(s) ds. \tag{2.19}$$

The first ordinary moment $\mu_1$ is called the *mean*. The second central moment $\overline{\mu}_2$ is called the *variance* and measure the dispersion around the mean. The square root of the variance is the *standard deviation*. The third standardized moment is called the *skewness* and it measures the symmetry of the distribution around the mean. Random variables with symmetrical PDF have zero valued skewness. The fourth standardized moment is called the *unnormalized kurtosis* and it measures the length of the tail of the distribution. The *normalized kurtosis* , simply called the kurtosis in the following, is defined as

$$\text{kurt}(s) = \breve{\overline{\mu}}_4 - 3. \tag{2.20}$$

The Gaussian distribution has a kurtosis equal to zero. If the kurtosis is superior or inferior than zero, the distribution is said to be leptokurtic or platykurtic, respectively.

**Definition 8.** *Let* s *be a random variable with a first characteristic function $\psi_1(\omega)$. Its second characteristic function is defined such that*

$$\psi_2(\omega) = \log \psi_1(\omega). \tag{2.21}$$

The Taylor expansion of the second characteristic function gives the definition of the cumulant of a random variable s. Once more, we have

$$\psi_2(\omega) = \kappa_1 i\omega - \frac{\kappa_2 \omega^2}{2} + \dots \quad \Rightarrow \quad \kappa_n = (-i)^n \left|\frac{d^n}{d\omega^n}\psi_2(\omega)\right|_{\omega=0}. \tag{2.22}$$

Second-order statistic (SOS) focusses only on the first and the second moments, in other words on the mean and the variance. This choice is not arbitrary and has deep links with the normal (or Gaussian) distribution that is fully described by its mean and its variance. As opposed to SOS, high-order statistics (HOS) explicitly make use of higher moments.

In its fundamental article, Shannon [1954] discusses ways to measure the quantity of information contained in random variables. He defines the entropy for discrete random variables and extends the concept to continuous random variables.

**Definition 9.** *Let* s *be a random variable with a PDF $f(s)$. The differential entropy of* s *is defined such that*

$$h_f = -\mathbb{E}_f\{\log[f(s)]\} = -\int_{-\infty}^{\infty} f(s) \log[f(s)] ds. \tag{2.23}$$

The differential entropy can be negative. For a fixed valued variance, the differential entropy is maximal for a Gaussian random variable.

**Definition 10.** *Let* s *be a random variable with a differential entropy* $h_f$. *The negentropy of* s *is defined such that*

$$Q_s = h_{\mathcal{N}_2} - h_f, \tag{2.24}$$

*where* $h_{\mathcal{N}_2}$ *is the differential entropy of a random vector following a Gaussian distribution with the same mean and variance as* s.

Finally, it is often practical to define a distance or divergence between two PDFs, especially for analyzing how close a random variable is from another. The Kullback-Leibler (KL) divergence is one of the possible definitions of this distance.

**Definition 11.** *Let* $f_1(s)$ *and* $f_2(s) \neq 0$ *be two PDFs. The Kullback-Leibler divergence is defined such that*

$$D_{KL}(f_1\|f_2) = \int_{-\infty}^{\infty} f_1(s) \log \frac{f_1(s)}{f_2(s)} ds. \tag{2.25}$$

In the above definition, the two PDFs are non commutative. One can show that we necessarily have $D_{KL} \geq 0$ and $D_{KL} = 0$ if and only if $f_1(s) = f_2(s)$.

### 2.3.2   Characterization of a set of random variables

So far we only considered a single random variable. From now on, we consider a set of random variables $s_1$, $s_2$, ..., $s_N$ or in other terms a random vector **s**. For each random variable, we can define its own PDF $f_1$, $f_2$, ... also called its marginal distribution. The joint cumulative density function is crucial when considering several random variables.

**Definition 12.** *Let* $\{s_i\}_{i=1}^N$ *be a set of random variables. The joint cumulative density function is defined as*

$$F(s_1, s_2, \dots) = Pr(\tilde{s}_1 < s_1, \tilde{s}_2 < s_2, \dots). \tag{2.26}$$

**Definition 13.** *Let* $\{s_i\}_{i=1}^N$ *be a set of random variables. The joint probability density function is defined as*

$$f(s_1, s_2, \dots) = \frac{\partial^N}{\partial s_1 \partial s_2 \dots} F(s_1, s_2, \dots). \tag{2.27}$$

The generalization of moments and cumulants can be obtained for a set of random variables. They are called cross-moments and cross-cumulants and constitute the *moment tensor* and the *cumulant tensor*. The first and second characteristic functions can also be generalized in the same way.

**Definition 14.** *Let* $s_1$ *and* $s_2$ *be two random variables. The covariance between those two random variables is defined as*

$$cov(s_1, s_2) = \mathbb{E}\{\bar{s}_1 \bar{s}_2\}, \tag{2.28}$$

*and their correlation coefficient is defined as*

$$corr(s_1, s_2) = \mathbb{E}\{\breve{s}_1 \breve{s}_2\}. \tag{2.29}$$

**Definition 15.** *Two random variables are said to be non-correlated if their covariance is null.*

**Definition 16.** *Let* $\{s_i\}_{i=1}^{N}$ *be a set of random variables. The covariance matrix* $\mathbf{\Sigma}$ *contains all the covariances computed for all pairs of random variables such that*

$$\mathbf{\Sigma} = \begin{bmatrix} \mathbb{E}\{\bar{s}_1\bar{s}_1\} & \mathbb{E}\{\bar{s}_1\bar{s}_2\} & \dots \\ \mathbb{E}\{\bar{s}_2\bar{s}_1\} & \mathbb{E}\{\bar{s}_2\bar{s}_2\} & \\ \vdots & & \ddots \end{bmatrix}. \tag{2.30}$$

*The correlation matrix* $\Theta$ *contains all the correlation coefficients icomputed for all pairs of random variables such that*

$$\mathbf{\Theta} = \begin{bmatrix} \mathbb{E}\{\breve{\bar{s}}_1\breve{\bar{s}}_1\} & \mathbb{E}\{\breve{\bar{s}}_1\breve{\bar{s}}_2\} & \dots \\ \mathbb{E}\{\breve{\bar{s}}_2\breve{\bar{s}}_1\} & \mathbb{E}\{\breve{\bar{s}}_2\breve{\bar{s}}_2\} & \\ \vdots & & \ddots \end{bmatrix}. \tag{2.31}$$

From the two above definitions, we see that the covariance matrix of a set of un-correlated variables is a diagonal matrix containing the variances and their correlation matrix is the identity matrix.

**Definition 17.** *Let* $\{s_i\}_{i=1}^{N}$ *be a set of random variables. The joint differential entropy is defined as*

$$h_{s_1,s_2,\dots} = -\int f(s_1, s_2, \dots) \log f(s_1, s_2, \dots) ds_1 ds_2 \dots \tag{2.32}$$

**Definition 18.** *The mutual information is defined for two random variables.*

$$I(s_1, s_2) = h_{s_1} + h_{s_2} - h_{s_1,s_2}. \tag{2.33}$$

**Definition 19.** *The random variables* $\{s_i\}_{i=1}^{N}$ *are said to be statistically independent if and only if*

$$f(s_1, \dots, s_N) = \prod_n f(s_n). \tag{2.34}$$

Statistical independence is a stronger concept than non-correlation. Statistically independent variables are necessarily un-correlated but un-correlated variables may be dependent. The mutual information between two variables is null if and only if the two random variables are statistically independent. The generalization of mutual information for several random variables is not straightforward and several definitions has emerged [McGill, 1954; Watanabe, 1960].

All the important concepts have been defined. In the next sections, we will discuss the blind source separation problem and two approaches for solving it, namely independent component analysis and sparse component analysis

## 2.4 Blind source extraction and separation

Blind source extraction (BSE) and separation (BSS) problems arose in many applications such as speech and audio processing [Asano et al., 2003; Ozerov and Fevotte, 2010], medical imaging [Vigario and Oja, 2008], geophysics [Ikelle, 2010] or astrophysics [Cardoso et al., 2002]. When the original signals of interest cannot be recorded directly but only some combinations or mixtures of them, we are facing a source separation problem. The term *separation* refers to the recovery of all sources (see subsection 2.4.2) while the term *extraction* refers to the recovery of one source (see subsection 2.4.3). If really few prior information is available on the mixing process, the problem is referred to as *blind* source separation (or extraction). We remind the

reader again that the term "source" has to be understood in the sense of "signal" in a wide sense.

The first tool that has been historically developed for solving BSS problems is independent component analysis (ICA). This class of technique makes the assumption that the original signals of interest can be modeled as statistically independent random variables [Comon, 1994; Amari and Cichocki, 1998; Hyvärinen and Oja, 2000]. More recently, others techniques have emerged based on other assumptions on the sources. One of them, namely sparse component analysis (SCA), makes the assumption that the sources are sparse in a given representation [Bofill and Zibulevsky, 2001; Gribonval and Lesage, 2006; Deville, 2014]. This technique is particularly valuable for dealing with partially correlated sources.

## 2.4.1   The mixing model

A BSS problem consists of recovering a set of original signals denoted $\boldsymbol{s}[x] \in \mathbb{R}^N$ through $M$ linear combinations of these sources denoted $\boldsymbol{d}[x] \in \mathbb{R}^M$, without any assumption on the mixing process. Here, $x$ denotes an index that can be of any dimension. For instance, for time series $x$ represents the time and for images $x$ represents the pixel location. We denote $X$ the number of indexes. For instance, if $x = \{x_1, x_2\}$, the number of indices is $X = X_1 \times X_2$. When the summation symbol $\sum_{x=1}^{X}$ is used, it implicitly means the summation over all dimensions, i.e. $\sum_{x_1=1}^{X_1} \sum_{x_2=1}^{X_2}$.

BSS is a general problem that can be adapted for denoising or for dealing with multi-channel observations [Comon and Jutten, 2010]. For an instantaneous mixture the observations at index $x$ are linear combinations of the sources taken at the same index. The instantaneous mixing equation can be written as

$$\boldsymbol{d}[x] = \boldsymbol{A}\boldsymbol{s}[x], \tag{2.35}$$

where the mixing matrix $\boldsymbol{A} \in \mathbb{R}^{M \times N}$ is unknown. If we develop the equation 2.35 we have

$$\begin{bmatrix} d_1[1] & d_1[2] & \ldots \\ d_2[1] & d_2[2] & \ldots \\ \vdots & \vdots & \ldots \\ d_M[1] & d_M[2] & \ldots \end{bmatrix} = \begin{bmatrix} a_{11} & a_{12} & \cdots & a_{1N} \\ a_{21} & a_{22} & \cdots & a_{2N} \\ \vdots & \vdots & \ddots & \vdots \\ a_{M1} & a_{M2} & \cdots & a_{MN} \end{bmatrix} \begin{bmatrix} s_1[1] & s_1[2] & \ldots \\ s_2[1] & s_2[2] & \ldots \\ \vdots & \vdots & \ldots \\ s_N[1] & s_N[2] & \ldots \end{bmatrix}.$$

Equation 2.35 is sometime referred to as the noiseless ICA model. When the number of observations is larger than the number of sources, i.e. $M > N$, the problem is said to be overdetermined. When the number of sources is equal to the number of observations, i.e. $M = N$, the BSS problem is said to be determined. When the number of observations is less than the number of sources, i.e. $M < N$, the problem is said to be underdetermined.

For a finite impulse response multiple-input multiple-output (FIR-MIMO) convolutive mixture [Pope and Bogner, 1996; Pedersen et al., 2008], the convolutive mixing equation is given by

$$\boldsymbol{d}[x] = \sum_{y=0}^{Y} \boldsymbol{A}[y]\boldsymbol{s}[x - y], \tag{2.36}$$

where the sum over $y$ is done over the memory length $Y$ of the mixing system. In a similar manner as for $X$ before, the summation is considered for all dimensions. For instance if $x = \{x_1, x_2\}$ the notation $\sum_{y=0}^{Y}$ is equivalent to $\sum_{y_1=0}^{Y_1} \sum_{y_2=0}^{Y_2}$. As

for the instantaneous BSS problem, all matrices $\boldsymbol{A}[y] \in \mathbb{R}^{M \times N}$ are unknown in the convolutive BSS problem. In the convolutive BSS problem, a filter ambiguity appears. We consider the following assumption.

**Assumption 1.** *Considering the mixing model in equation 2.36, the system $\{\boldsymbol{A}[y]\}_{y=0}^{Y}$ is invertible.*

For an instantaneous mixing system, the invertibility of the system just means that the mixing matrix $\boldsymbol{A}$ is invertible. For a convolutive FIR-MIMO system, invertibility is more difficult to address. We refer to Rajagopal and Potter [2003], Castella and Pesquet [2004], Castella and Moreau [2009] and in particular to Law et al. [2009] for a complete description of conditions ensuring the identifiability of FIR-MIMO systems.

### 2.4.2 Blind source separation

We described the mixing model for BSS problem. Under assumption 1, the mixing system of an overdetermined or determined problem can be inverted. Solving a BSS consists of finding a separating matrix $\boldsymbol{W} \in \mathbb{R}^{N \times M}$ that gives an estimate of the sources such that

$$\boldsymbol{u}[x] = \boldsymbol{W}\boldsymbol{d}[x] = \boldsymbol{W}\boldsymbol{A}\boldsymbol{s}[x] = \boldsymbol{H}\boldsymbol{s}[x], \tag{2.37}$$

where $\boldsymbol{H} = \boldsymbol{W}\boldsymbol{A}$ represents the global mapping between the true sources and the estimates. An instantaneous problem can be resolved up to two ambiguities. Firstly, the amplitude of each original signal cannot be recovered. Secondly, the order of the signals cannot be recovered. The solution set is defined as follows.

**Definition 20.** *The set of solutions $\mathcal{S}$ of an instantaneous BSS problem is defined such that*

$$\mathcal{S} = \{\hat{\boldsymbol{s}} \ : \ \hat{\boldsymbol{s}}[x] = \boldsymbol{\Pi}\boldsymbol{\Delta}\boldsymbol{s}[x]\}, \tag{2.38}$$

*where $\boldsymbol{\Pi}$ is a permutation matrix[3] and $\boldsymbol{\Delta}$ is a diagonal matrix. The BSS problem is solved when $\boldsymbol{u}[x] \in \mathcal{S}$.*

Under assumption 1, the separation of sources in a determined convolutive BSS problem can be achieved by finding a multiple-input multiples-output (MIMO) separating system of $Z + 1$ extraction matrices $\boldsymbol{W}[z] \in \mathbb{R}^{N \times M}$, with $Z \geq Y$[4], such that

$$\boldsymbol{u}[x] = \sum_{z=0}^{Z} \boldsymbol{W}^T[z]\boldsymbol{d}[x - z] \tag{2.39}$$

$$= \sum_{z=0}^{Z} \boldsymbol{W}^T[z] \sum_{y=0}^{Y} \boldsymbol{A}[y]\boldsymbol{s}[x - z - y] \tag{2.40}$$

$$= \sum_{l=0}^{L} \boldsymbol{H}^T[l]\boldsymbol{s}[x - l], \tag{2.41}$$

where we defined $L = Y + Z$ and the vectors $\boldsymbol{h}[l]$ are the mapping vectors between the original sources and the extracted signal defined such that $\boldsymbol{H}[l]^T = \sum_{z=0}^{Z} \boldsymbol{W}[z]^T \boldsymbol{A}_{z-y}$ with $0 < z - y \leq Y$.

---

[3]A permutation matrix is a square matrix containing only a single coefficient per row and per column. For instance in $\mathbb{R}^3$ the following matrix

$$\begin{bmatrix} 0 & 1 & 0 \\ 1 & 0 & 0 \\ 0 & 0 & 1 \end{bmatrix}$$

is a permutation matrix.

[4]If $x$ is multi-dimensional, $Z \geq Y$ stand for $Z_1 \geq Y_1, Z_2 \geq Y_2, \ldots$

**Definition 21** (Separability)**.** *The mixing model is said to be separable by $\phi(\boldsymbol{W})$ if each local minimum $\boldsymbol{W}^*$ of $\phi$ is such that $\boldsymbol{W}^* = \boldsymbol{\Pi}\boldsymbol{\Delta}\boldsymbol{A}$.*

**Definition 22.** *A function $\phi(\boldsymbol{W})$ is called a contrast function if its minimum $\boldsymbol{W^*}$ extracts the sources, i.e. if $\boldsymbol{u}^*[x] = \boldsymbol{W}^*\boldsymbol{d}[x]$ is in $\mathcal{S}$.*

The definition 22 is at the core of all method solving BSS problems. Under some assumptions, contrast functions are constructed and optimized in order to solve the problem. As we will see later, the next chapter 3 will be dedicated to one specific function, namely the $\ell_0$ pseudo-norm, and the conditions under which this function is a contrast function for BSS problems.

The separability is a more tight concept than the contrast function. If the mixing model is separable by $\phi$, then $\phi$ is a contrast function. However, if a function is a contrast, all its local minima are not necessarily corresponding to an extraction. If the problem is separable, then it is identifiable.

### 2.4.3   Blind source extraction

Sometimes, it is not necessary to recover all the sources, but only one of them in the set. This problem refers to as blind source extraction (BSE) [Delfosse and Loubaton, 1995]. It can be solved under assumption 1 by finding an extracting vector $\boldsymbol{w}$ such that

$$u[x] = \boldsymbol{w}^T\boldsymbol{d}[x] = \boldsymbol{w}^T\boldsymbol{A}\boldsymbol{s}[x] = \boldsymbol{h}^T\boldsymbol{s}[x], \tag{2.42}$$

where $\boldsymbol{h}$ is a global mapping vector. Similarly to definition 20, the solution set $\mathcal{S}_n$ of a BSE problem contains all the scaled versions of the signal $\boldsymbol{s}_n$ except the zero vector. For the extraction of a source in a determined convolutive BSE problem, a system of extracting vector $\{\boldsymbol{w}[z]\}$ is used instead of a system of separation matrices.

**Definition 23.** *The set of solutions $\mathcal{S}_q$ of an instantaneous BSE problem is defined such that*

$$\mathcal{S}_q = \{\boldsymbol{u} \; : \; \boldsymbol{u} = \alpha\delta_k(\boldsymbol{s}_q), \; \forall\alpha \in \mathbb{R}\backslash\{0\}, \; \forall k \in \mathbb{Z}\}, \tag{2.43}$$

*where $\delta_k$ denotes the time shift operator. We call $\mathcal{G}_q$ the set of all mapping vectors such that $\{\boldsymbol{h}_j\}_{j=0}^J \in \mathcal{G}_q \;\Leftrightarrow\; \boldsymbol{u} \in \mathcal{S}_q$. The set $\mathcal{G}_q$ denotes the solution set of global mapping corresponding to the correct extraction of the source $\boldsymbol{s}_q$. The BSE problem is equivalently solved when $\boldsymbol{u} \in \mathcal{S}_q$ or when $\{\boldsymbol{h}_j\}_{j=0}^J \in \mathcal{G}_q$.*

**Definition 24.** *A function $\phi(\boldsymbol{w})$ is called a contrast function for the extraction of $\boldsymbol{s}_q$ if its minimum $\boldsymbol{w^*}$ extracts the source $\boldsymbol{s}_q$, i.e. if $\boldsymbol{u}^*[x] = \boldsymbol{w}^*\boldsymbol{d}[x]$ is in $\mathcal{S}_q$.*

### 2.4.4   Performance measurement in BSE and BSS

Before going further into the methods for solving BSS problems, it is important to define some performance measures [Vincent et al., 2006, 2012], determining the quality of the recovered sources compared to the original sources. These distances must give a quantitative measure of how close are our estimates to the true sources. Those indexes should be insensitive to amplitude, permutation and delay ambiguities. The classic signal to noise ratio ($SNR$) is given by

$$SNR = 10\log\frac{\|\boldsymbol{s}_n\|_2}{\|\boldsymbol{s}_n - \hat{\boldsymbol{s}}_n\|_2}, \tag{2.44}$$

**Figure 2.4:** *The value of the ISI performance measure in (left) $\mathbb{R}^2$ and (right) $\mathbb{R}^3$ for a single extraction vector.*

where ambiguities must be carefully handled before the computation of such a ratio. In particular, the amplitude ambiguity must not interfere and

$$\hat{\boldsymbol{s}}_n = \alpha \boldsymbol{u}^* \quad \text{with} \quad \alpha = \left(\boldsymbol{u}^{*T}\boldsymbol{u}^*\right)^{-1} s_n \boldsymbol{u}^{*T}, \tag{2.45}$$

where $\boldsymbol{u}^*$ is the extracted signal after optimization of the contrast function and $\alpha$ is the optimal scaling factor in the least-squares sense.

When the inverse system exists and the global mappings are known, the intersymbol interference (*ISI*) is a particularly attractive performance measure [Lambert, 1999]. For an instantaneous BSE problem with global mapping $\boldsymbol{h}^T = \boldsymbol{w}^T \boldsymbol{A}$, it is defined as

$$ISI = \frac{\sum_n h_n^2 - \max h_n^2}{\max h_n^2}. \tag{2.46}$$

We always have $ISI \geq 0$ with equality if and only if the sources are recovered up to a scale ambiguity. For a convolutive BSE problem with global mapping $\boldsymbol{H}[y]$

$$ISI = \frac{\sum_n H_n[y]^2 - \max H_n[y]^2}{\max H_n[y]^2}. \tag{2.47}$$

Here, the *ISI* is also unchanged with delay. An important remark on the *ISI* measure, compared to *SNR*, is that it refers to the mixing system which is blind and not to the sources themselves. Figure 2.4 shows the value of the *ISI* for $N = 2$ and $N = 3$ sources. Its value is minimum and equal to 0 on the main axis only, when the mapping vector corresponds to the extraction of one source. Its value does not depend on the norm of $\boldsymbol{h}$, but only on the direction.

### 2.4.5 Principal component analysis and its limit

Principal component analysis (PCA or Karhunen-Loève transform) is a convenient method for data analysis and matrix factorization based on second-order statistics. One could think of using PCA for solving BSS problems. It is well known that, actually, PCA cannot solve BSS problems but it is a common tool for whitening the observations [Comon and Jutten, 2010]. PCA algorithms look for a new orthogonal frame in which it is more efficient to observe the data in term of variances. The first component is given by

$$\boldsymbol{w}_1 = \arg\min \left\|\boldsymbol{w}^T \boldsymbol{d}[x]\right\|_2 \quad \text{such that} \quad \|\boldsymbol{w}\|_2 = 1, \tag{2.48}$$

and the second component as

$$\boldsymbol{w}_2 = \arg\min \left\| \boldsymbol{w}^T \boldsymbol{d}[x] \right\|_2 \quad \text{such that} \quad \|\boldsymbol{w}\|_2 = 1 \text{ and } \boldsymbol{w}^T \boldsymbol{w}_1 = 0. \tag{2.49}$$

A PCA algorithm leads to uncorrelated components. Is this assumption enough for separating sources? Actually, no, and this is easy to see. We are looking for estimated sources such that $\boldsymbol{\Sigma} = \hat{\boldsymbol{S}}^T \boldsymbol{S} = \boldsymbol{I}$. If we consider a rotation matrix $\boldsymbol{R}_{\boldsymbol{\theta}}$ that we apply on the estimated sources, we have

$$\boldsymbol{\Sigma}_{\boldsymbol{\theta}} = (\boldsymbol{R}_{\boldsymbol{\theta}} \hat{\boldsymbol{S}})^T \boldsymbol{R}_{\boldsymbol{\theta}} \boldsymbol{S} = \hat{\boldsymbol{S}}^T \boldsymbol{R}_{\boldsymbol{\theta}}^T \boldsymbol{R}_{\boldsymbol{\theta}} \boldsymbol{S} = \boldsymbol{I} = \boldsymbol{\Sigma}. \tag{2.50}$$

The correlation matrix is unchanged for any chosen rotation matrix.

Hence, PCA algorithms cannot solve BSS problems under the simple un-correlation assumption because of the rotation ambiguity. In other words, second-order statistics are not sufficient and, as we will see in the next section, we must use higher order statistics. Nevertheless, PCA and related are commonly used on the observations for whitening. It makes ICA algorithms faster and more robust.

## 2.5   Independent component analysis

As explained in the previous subsection, the rotation ambiguity shows that the prior information of decorrelation is not strong enough for a correct separation in BSS problems. In other words, sources that are only uncorrelated cannot be separated. The main result of ICA methods for solving BSS problems is that statistically independent sources with any kind of distribution, except the Gaussian distribution can be separated [Comon, 1994]. The main restriction, indeed, is that at most one source can have a Gaussian distribution.

**Assumption 2.** *The sources (signals) are statistically independent and at most one of them is a Gaussian random variable.*

Based on this assumption, a number of algorithms can be developed. We distinguish two classes of algorithms. The first is based on maximizing the likelihood while the second is based on maximizing the distance from Gaussian distribution.

### 2.5.1   Infomax and maximum likelihood

As proposed by Bell and Sejnowski [1995], maximizing the mutual information $I$ between the inputs and the outputs of the neural network in figure 2.5 leads to the recovery of estimated sources that are statistically independent (see also Linsker [1989]). It is important to note that, in the context of BSS, the outputs $\boldsymbol{z}_i = G_0(\boldsymbol{u}_i)$ of the neural network are auxiliary variables used to optimize the statistical independence between the estimated sources $\boldsymbol{u}_i$. Often, logistic functions are used such as the sigmoid function

$$G_0(s) = \frac{1}{1 + \mathrm{e}^{-\lambda s}}, \quad \text{with} \quad g_0(s) = G_0'(s) = \frac{\lambda \mathrm{e}^{-\lambda s}}{(1 + \mathrm{e}^{-\lambda s})^2}, \tag{2.51}$$

where $\lambda$ is a shaping parameter.

It can be shown [Bell and Sejnowski, 1995] that maximizing the mutual information $I$ is actually equivalent to minimizing the following objective function

$$\phi_{IM} = -\mathbb{E}\{\log(|\boldsymbol{J}|)\}, \tag{2.52}$$

**Figure 2.5:** *The Infomax neural network of a BSS linear mixing model and the separation.*

where $\boldsymbol{J}$ is the Jacobian of the neural network defined as

$$\boldsymbol{J} = \det \begin{bmatrix} \frac{\partial z_1}{\partial d_1} & \frac{\partial z_1}{\partial d_2} & \cdots \\ \frac{\partial z_2}{\partial d_1} & \frac{\partial z_2}{\partial d_2} & \\ \vdots & & \ddots \end{bmatrix}. \tag{2.53}$$

Cardoso [1997] demonstrates that the Infomax method can be interpreted as the maximum likelihood estimation of the sources with the estimated CDF $G_0(s)$ [Xu, 2005].

### 2.5.2 Maximization of non-Gaussianity

Independent sources can also be recovered by maximizing their distance to the Gaussian distribution. This can be explained by the central limit theorem (CLT) [Leon-Garcia, 2008]. Due to the CLT, the mixture are more Gaussian than the signals. This also explains why Gaussian sources cannot be recovered (except one).

As explained in Comon and Jutten [2010] or Hyvärinen et al. [2001], recovering estimates $\hat{\boldsymbol{S}}$ as far as possible from a Gaussian distribution is a valid method for recovering independent sources, after whitening of the observations. This measure of Gaussianity can be achieved by computing the negentropy $Q_{\hat{\boldsymbol{S}}}$ of the estimated sources. As proposed by Hyvärinen and Oja [2000], we can approximate the maximization of negentropy by minimizing the following objective function

$$\phi_Q = -\left( \mathbb{E}\{g_i(\hat{\boldsymbol{S}}_0)\} - \mathbb{E}\{g_i(\boldsymbol{N}_0)\} \right)^2, \tag{2.54}$$

where $\hat{\boldsymbol{S}}_0$ has zero mean vector and unit co-variance matrix and $\boldsymbol{N}_0$ is a random standardized Gaussian vector. The function $g_i(\cdot)$ can be chosen within, for instance,

the following set of non-quadratic functions [Li and Lu, 2013]

$$g_1(s) \quad = \quad -\exp\left(-\frac{s^2}{2}\right), \tag{2.55}$$

$$g_2(s) \quad = \quad \log(\cosh(s)), \tag{2.56}$$

$$g_3(s) \quad = \quad \sqrt{1+s^2} - 1. \tag{2.57}$$

The kurtosis can also be used to measure the distance to a Gaussian (see equation 2.20).

### 2.5.3   ICA with convolutive mixture

The major part of works about ICA and BSS considers the instantaneous mixing model. However, a non-negligible part is dedicated to handle the convolutive case (see [Comon and Jutten, 2010], chapter 8). In most works, the convolutive problem is transformed in the frequency domain by slicing the observations in the time domain and using the short-time Fourier transform [Pedersen et al., 2008; Sawada et al., 2010]. We can write

$$\mathcal{F}\boldsymbol{d}[\omega, \tau] = \boldsymbol{A}[\omega, \tau] \cdot \mathcal{F}\boldsymbol{s}[\omega, \tau], \tag{2.58}$$

where $\mathcal{F}$ denotes the short-time Fourier transform, $\omega$ the frequency and $\tau$ the time index of the considered bin. Hence, the convolutive model is reduced to several instantaneous models, one per frequency and time bin. The main drawback of this strategy is that the permutation ambiguity, present at each frequency bin, must be carefully addressed for an accurate reconstruction of the sources [Duong et al., 2009]. Benichoux et al. [2012] discuss the filter ambiguity.

### 2.5.4   Recovery conditions with ICA

We focus here on two definitions, namely *separability* and *uniqueness* [Comon, 1994; Cao and Liu, 1996; Eriksson and Koivunen, 2004]. The model in equation 2.35 is said to be identifiable if we can retrieve the mixing matrix up to scaling and permutation [Deville, 2014]. The model in equation 2.35 is said to be separable if we can retrieve the sources up to scaling and permutation. One can notice that separability implies identifiability. They are equivalent in the case of determined problems. The model in equation 2.35 is said to be unique, if the decomposition is unique up to scaling and permutation. The following results can be found in the literature.

**Theorem 1** ([Comon, 1994], [Eriksson and Koivunen, 2004])**.** *The model in equation 2.35 is separable if the mixing matrix $\boldsymbol{A}$ is of full column rank and at most one source variable is normal.*

**Theorem 2** ([Eriksson and Koivunen, 2004])**.** *If the model is separable, then the ICA model is unique.*

Eriksson and Koivunen [2006] extend these results for complex-valued sources. The study of uniqueness and separability of model is an active area of research [Murillo-Fuentes and Boloix-Tortosa, 2010; Wang et al., 2015].

Considering now the Infomax and maximum likelihood contrast functions, one generally assumes that the true PDFs are known. In practice, however, the PDF of the sources are generally unknown or difficult to model. The next theorem named the one-bit matching ICA theorem gives a useful result.

**Theorem 3** ([Xu, 2005])**.** *If the kurtoses of the modeling PDFs $f_0(u_i)$ have the same sign as the true PDFs $f(s_i)$, then the sources can be separated.*

This result is of important value. It means that we do not have to perfectly model the true PDF of the sources we want to recover. Roughly knowing the PDF is sufficient for an exact recovery, at least in the noise-less case.

## 2.6 Sparse component analysis

When the assumption of statistical independent variables is no-longer valid, ICA may fail to recover the original sources. Other prior should be used such as the sparsity of the sources in a domain or in a dictionary $\mathcal{D}$. Sparse component analysis (SCA) has been developed with the purpose of decomposing the observations such that $\boldsymbol{d}[x] = \hat{\boldsymbol{A}}\hat{\boldsymbol{s}}[x]$, where the estimated sources $\hat{\boldsymbol{s}}$ are sparse. SCA can be used for solving BSS problems [Bofill and Zibulevsky, 2001; Gribonval and Lesage, 2006; Mourad and Reilly, 2010; Nadalin, 2011; Bobin et al., 2015].

### 2.6.1 Sparsity of signals

The $\ell_0$ pseudo-norm is often used as a measure of sparseness [Hurley and Rickard, 2009]. Strictly speaking, it is neither a norm nor a pseudo-norm. However, the terminology of pseudo-norm is often seen in the literature. We maintain this abuse of terminology (see also appendix 8.1 for more details).

**Definition 25.** *The $\ell_0$ pseudo-norm of a vector $\boldsymbol{s}$ is defined such that*

$$\|\boldsymbol{s}\|_0 = \lim_{p \to 0} \|\boldsymbol{s}\|_p = \# \left\{ s_x \mid s_x \neq 0 \right\}, \tag{2.59}$$

*which is the number of non-zero coefficients in the vector $\boldsymbol{s}$.*

As for ICA, the conditions for a correct estimation of the sources can be investigated in SCA. Because SCA can be used for solving under-determined problems, system recovery and source recovery are often treated separately. First, the recovery of the mixing model $\boldsymbol{A}$ from the data studies the conditions underwhich blind identifiability is possible. If this mixing model $\hat{\boldsymbol{A}}$ can be inverted, the sources can be recovered directly at this stage. However, if the problem is underdetermined, we need to estimate the sources $\hat{\boldsymbol{s}}[x]$ from the data $\boldsymbol{d}$ and the estimated mixing model $\hat{\boldsymbol{A}}$.

Figure 2.6 shows an determined example with two sources, for which the ICA assumption is not adequate for a correct extraction or separation. The two sources (figure 2.6a) have a strong correlation in the background (the points) and some sparse lines. The minimum of the kurtosis contrast function corresponds to an incorrect signal (figure 2.6b) while the minimum of the $\ell_0$ pseudo-norm corresponds to the correct source extraction (figure 2.6c). The kurtosis objective function separates the sources based on a statistical independence assumption, which is clearly not valid in this example. The objective function has two local minima, one for each source, associated to an incorrect source extraction. The $\ell_0$ pseudo-norm has clearly three strong local minima. Two of them correspond to a correct source extraction. The extracting source $\hat{\boldsymbol{u}_0}$ associated to the parasitic minima is canceling the correlated background made of points.

### 2.6.2 Matrix recovery (identifiability)

The following assumption is very important for a lot of methods and results. It sets the statistical dependence between the sources such that only one source can be active, with a coefficient different from zero, at a time.

a)



b)



c)



**Figure 2.6:** *Synthetic example. a) The two initial sparse images $s_1$ and $s_2$. b) The kurtosis-based contrast function and the extracted image. c) The $\ell_0$ pseudo-norm contrast function and the three images corresponding to the three lower minima. The objective functions are shown in the global mapping space.*

**Assumption 3** (Disjointness of the sources.)**.** *At most one signal is active at a time.*

Such a set of signals are also said to be disjoint orthogonal (DO). DO variables are decorrelated but statistically dependent variables. This aspect will be discussed in chapter 3.

Bofill and Zibulevsky [2001] is one of the first work showing that disjoint orthogonal sources can be recovered from fewer observations in instantaneous mixtures. He shows that after the mixing process, the data samples are located along lines corresponding to the columns of $\boldsymbol{A}$. This idea has been expanded for simple delayed mixtures in Bofill [2003]. Around the same time, Jourjine et al. [2000] propose that any number of sources can be separated from only two mixtures, if the sources satisfy assumption 3 in the frequency domain (see also [Rickard and Yilmaz, 2002]). They developed the popular DUET algorithm based on clustering [Scott, 2007]. These results can be summarized by the following theorem.

**Theorem 4.** *If the sources are disjoint orthogonal, then the mixing system is identifiable.*

In order to relax the DO assumption and possibly letting several sources to be active at a time, Georgiev et al. [2007] propose sufficient conditions, on the source only, for the recovery of both $\boldsymbol{S}$ and $\boldsymbol{A}$ in the underdetermined case. They consider column sparsity and a subspace clustering algorithm but their conditions are difficult to verify at first sight (see also Georgiev et al. [2005, 2007]). Mishali and Eldar [2009] develop algorithms around the same idea (see also [Lindenbaum et al., 2015]).

Sun and Xin [2011] propose a necessary and sufficient for the underdetermined BSS problem of $m + 1$ sources, under the non-negative assumption of both $\boldsymbol{A}$ and $\boldsymbol{S}$. The two conditions forming the NSC are: i) on $\boldsymbol{S}$, single source zones exist for each source; ii) at most $m - 1$ sources are active at the same time; iii) on $\boldsymbol{A}$, the mixing matrix is degenerated. They extend their result for $n$ sources.

Duarte et al. [2011] propose a sufficient condition for the extraction and separation of sparse source in instantaneous mixtures. Our work in chapter 3 will expand these results.

**Theorem 5** ([Duarte et al., 2011])**.** *If $\|\boldsymbol{s_1}\|_{\ell_0} < \frac{1}{2}\|\boldsymbol{s_2}\|_{\ell_0}$, then the $\ell_0$ pseudo-norm is a contrast function for the extraction of $\boldsymbol{s_1}$.*

### 2.6.3 Source recovery (separability)

Once the mixing model $\hat{\boldsymbol{A}}$ has been estimated, we must recover the original sources. If the mixing model is overdetermined, the separation matrix can be estimated by $\boldsymbol{W} = \hat{\boldsymbol{A}}^{-1}$ and the sources are simply estimated by

$$\hat{\boldsymbol{s}}[x] = \boldsymbol{W}\boldsymbol{d}[x] = \hat{\boldsymbol{A}}^{-1}\boldsymbol{d}[x]. \tag{2.60}$$

However, if the mixing model is underdetermined, we must investigate a signal recovery problem. This problem has received a lot of attention, particularly the $\ell_0$ pseudo-norm optimization problem $P_0$:

$$\min \|\boldsymbol{s}[x]\|_0 \quad \text{such that} \quad \boldsymbol{d}[x] \approx \boldsymbol{A}\boldsymbol{s}[x]. \tag{2.61}$$

This problem is quite general in signal processing. A complete research has emerged from the recent results in this area, namely compressive sensing.

The sparse recovery problem has received a lot of attention, starting from the work of Donoho and Elad [2003] or Gribonval and Nielsen [2003]. In particular, a central

question has been to know under which conditions the problem 2.61 can be solved by convex optimization. In particular, it has been proved that under some conditions on the matrix $\boldsymbol{A}$ and if $\boldsymbol{s}$ is sparse enough, the problem 2.61 can be solved by a $\ell_1$-norm minimization problem $P_1$ [Li et al., 2003b]:

$$\min \|\boldsymbol{s}[x]\|_1 \quad \text{such that} \quad \boldsymbol{d}[x] \approx \boldsymbol{A}\boldsymbol{s}[x]. \tag{2.62}$$

**Theorem 6** (Sufficient condition for the equivalence of $P_0$ and $P_1$ [Donoho and Elad, 2003]). *Less than 50% concentration implies equivalence.*

Li and Amari [2010] tackle the problem in a probabilistic framework and found interesting results and estimates (see also [Li et al., 2006]).

**Theorem 7** (Necessary and sufficient condition for the equivalence of $P_0$ and $P_1$ [Li and Amari, 2010]). *$\hat{\boldsymbol{s}}_{\ell_1} = \hat{\boldsymbol{s}}_{\ell_0}$ if and only if*

$$\max \sum_{i=1}^{N_s} [\text{sign}(s_i)\delta_i]_+ < \frac{1}{2}, \tag{2.63}$$

*where $\boldsymbol{\delta}$ is such that $\boldsymbol{A}\boldsymbol{\delta} = \boldsymbol{0}$ and $\|\boldsymbol{\delta}\|_1 = 1$.*

However, all the work on compressive sensing is dedicated to design the matrix $\boldsymbol{A}$ with specific properties. In BSS, such an approach is not valid as one cannot assume a priori that $\boldsymbol{A}$ follows such properties. Saab et al. [2007b] point out that the solution of a $\ell_p$-norm minimisation program (basis-pursuit) necessarily gives an estimate of $\boldsymbol{S}$ that is $k$-column-sparse, $k \leq N$. This result is true only for a real-valued mixing matrix and fails for complex-valued mixing matrices.

## 2.7   Conclusion

In this chapter 2, we provide an introduction to the BSS problem which can be seen as a bilinear inverse problem with an unknown direct problem. In particular, we focused on the contrast functions (objective functions) used for tackling this specific problem. Higher-order statistics form the theoretical background for methods based on ICA that are able to separate sources that are statistically independent. When this assumption fails, other prior information must be used such as sparsity for SCA methods. We have reviewed the literature about SCA conditions. However, what is yet not clear are the exact necessary and sufficient conditions under which the SCA method is valid. The next chapter 3 is dedicated to this aspect.

# Chapter 3

# Blind extraction and separation of sparse sources

## Contents

Part of the results in this chapter has been presented at the European Signal Precessing Conference (EUSIPCO) in Budapest [Batany et al., 2016a]. More contents have been added here, particularly concerning disjoint component analysis (DCA).

## Résumé du chapitre [français]

Le chapitre 3 traite de la séparation aveugle de sources parcimonieuses pour des mélanges déterminés instantanés et convolutifs. Lorsque l'hypothèse d'indépendance statistique des sources n'est pas valable, en d'autres termes lorsque les signaux à estimer sont statistiquement dépendants , on peut (ou l'on doit) s'appuyer sur d'autres caractéristiques des signaux. L'hypothèse de parcimonie des sources est souvent utilisée et mène à l'analyse en composantes parcimonieuses. Un signal est dit parcimonieux sous une représentation donnée lorsqu'un nombre faible de coefficients permet de le représenter.

L'hypothèse forte d'orthogonalité disjointe est d'abord traitée dans la section 3.2. Elle correspond au cas où au plus une seule source peut être active à la fois. Il est montré que, dans le cas déterminé, la méthode Infomax peut être utilisée pour séparer des sources orthogonalement disjointes. Il est ainsi mis en exergue le fait que l'indépendance statistique n'est pas une condition nécessaire pour la séparation de sources.

La section 3.3 introduit le concept de *processus inter-regressif* qui peut être vu comme une généralisation de la notion de processus auto-regressif pour plusieurs signaux.

Dans les sections 3.4 et 3.5, nous analysons une condition nécessaire et suffisante à l'extraction et à la séparation de sources parcimonieuses, en utilisant la pseudo-norme $\ell_0$ comme fonction de contraste. Ces résultats utilisent directement la définition des processus inter-regressifs.

Finalement, la section 3.6 propose un algorithme évolutionniste de type "évolution différentielle" pour résoudre des problèmes d'extraction et de séparation de source basé sur une version lissée de la pseudo-norme $\ell_0$. Plusieurs exemples sont traités et analysés. La présence de bruit est examinée et une méthode d'analyse de Pareto est proposée pour déterminer le niveau de bruit à partir des données uniquement.

## Resumo do capítulo [português]

O capítulo 3 trata da separação cega de fontes esparsas para mistura determinada instantânea e convolutiva. Quando a hipótese da independência estatística de fontes não é verificada, devemos nos basear em outras características dos sinais. A hipótese da esparsidade das fontes é freqüentemente utilizada e leva à análise de componentes esparsos. Um sinal é chamado de esparso sob uma dada representação na medida em que um número pequeno de coeficientes permite representá-lo.

A forte hipótese da ortogonalidade disjunta entre as fontes é discutida inicialmente na seção 3.2. Ela corresponde ao caso onde no máximo uma única fonte pode estar ativa em um dado instante. Demonstra-se que, no caso determinado, o método Infomax pode ser utilizado para separar fontes ortogonalmente disjuntas. É igualmente salientado o fato de que a independência estatística não é uma condição necessária para a separação de fontes.

A seção 3.3 introduz o conceito de *processo inter-regressivo*, que pode ser visto como uma generalização da noção de um processo auto-regressivo para vários sinais.

Na seção 3.4 e 3.5 analisamos uma condição necessária e suficiente para a extração e separação de fontes esparsas utilizando a pseudo norma $\ell_0$ como função de contraste. Esses resultados utilizam diretamente a definição de processos inter-regressivos.

Por fim, a seção 3.6 sugere um algoritmo evolutivo do tipo "evolução diferencial" para resolver problemas de extração e de separação de fontes baseados na versão suave da pseudo norma $\ell_0$. Vários exemplos são discutidos e analisados. A presença de ruído é examinada e um método de análise de Pareto é sugerido para determinar o nível de ruído exclusivamente a partir dos dados.

## 3.1  Introduction

In various situations, the observations we make of physical processes are mixtures of several sources (signals) and one would like to recover the original sources. In a general setting, we can write this problem as $\boldsymbol{d} = \mathcal{A}(\boldsymbol{s})$ where $\boldsymbol{d}$ is the vector of observations, $\boldsymbol{s}$ is the source vector (the original signal) and $\mathcal{A}$ represents the mixing process. When few information is available about the mixing process and all the sources must be recovered, the problem is referred to as blind source separation (BSS). If only one source must be recovered, the problem is called blind source extraction (BSE) [Comon and Jutten, 2010].

In order to recover the original sources with little prior information on the mixing process, some assumptions must be added on the original sources to be recovered. Independent component analysis (ICA) makes the assumption that the sources are statistically independent and one can prove that this assumption is sufficient to ensure the recoverability of the sources for linear problems [Comon, 1994].

Based on another assumption, namely the sparsity of the sources, sparse component analysis (SCA) has shown to be able to solve BSS problems when the sources are active over a restricted support. However, to our knowledge, the exact necessary and sufficient conditions on the sources for their recovery have not been described yet in details. The present chapter 3 attempts to fill this lack. In particular, we focus on convolutive mixtures. Besides this theoretical work, the present chapter also provides some numerical analysis in order to illustrate the conditions derived herein.

Concerning the chapter organization, section 3.2 is dedicated to present the strong sparsity assumption of disjoint orthogonal sources. Section 3.3 defines a new concept, namely *inter-regressive process*, that will be used in the next sections. Section 3.4 presents the necessary and sufficient conditions for the correct extraction of a sparse source. Section 3.5 presents the necessary and sufficient conditions for the correct separation of sparse sources. Finally, section 3.6 describes a differential evolution (DE) algorithm for solving the extraction and the separation of sources. A method for dealing with noisy observations is presented. Some examples are provided and discussed.

## 3.2  Blind separation of disjoint orthogonal variables

ICA methods assume that the sources of interest are statistically independent and one can show that this is a *sufficient* condition for the separability of sources [Comon, 1994]. In this section, we focus on a particular kind of statistically dependent class of sources, namely disjoint orthogonal (DO) random variables. This class has been of huge interest in SCA [Bofill and Zibulevsky, 2001; Scott, 2007]. Hereafter, we show that ICA-based contrast functions are also contrast functions for the separation of this particular class of sources. This result emphasizes that the assumption of statistical independence is definitely not a necessary condition [Caiafa, 2012]. In this section we consider disjoint orthogonal (DO) variables. The case where both $s_1$ and $s_2$ are equal to zero is not considered here because, in such a case, all the observation coefficients $\boldsymbol{d} = \boldsymbol{As}$ are equal to zero and can be easily discarded.

### 3.2.1 Mutual information of two DO variables

The joint PDF $f_{1,2}(s_1, s_2)$ of two disjoint orthogonal random variables $s_1 \sim f_1(s_1)$ and $s_2 \sim f_2(s_2)$ can be written as (see figure 3.1a)

$$f_{1,2}(s_1, s_2) = \psi_1(s_1)\delta(s_2) + \psi_2(s_2)\delta(s_1), \tag{3.1}$$

where $\delta(s_i)$ are Dirac distributions and $\psi_i(s_i)$ are positive continuous functions defined such that

$$\psi_i(0) = 0, \tag{3.2}$$

$$\int_{-\infty}^{\infty} ds_i \ \psi_i(s_i) = P_i, \tag{3.3}$$

with $P_1 + P_2 = 1$. The value $P_i$ represents the probability of $s_i$ to be active (i.e. non-zero). Equation 3.1 means that if $s_1 = 0$ then $s_2 \neq 0$ and if $s_1 \neq 0$ then $s_2 = 0$. The marginal probability distribution of each random variable is given by

$$f_i(s_i) = \psi_i(s_i) + (1 - P_i)\delta(s_i). \tag{3.4}$$

The mutual information of two random variables is defined such that

$$I(s_1, s_2) = \iint_{\mathbb{R}^2} ds_1 ds_2 \ f_{1,2}(s_1, s_2) \log \frac{f_{1,2}(s_1, s_2)}{f_1(s_1)f_2(s_2)}. \tag{3.5}$$

In the case of two disjoint orthogonal random variables, the domain of integration can be restricted to the two axes without the origin. We define the two sets $\mathcal{A}_1 = \{s_1, s_2 | s_1 \neq 0, s_2 = 0\}$ and $\mathcal{A}_2 = \{s_1, s_2 | s_1 = 0, s_2 \neq 0\}$. The value of the joint PDF is null outside $\mathcal{A}_1 \cup \mathcal{A}_2$. Hence we have

$$I(s_1, s_2) = \int_{\mathcal{A}_1} ds_1 \ \psi_1(s_1) \log \frac{\psi_1(s_1)}{\psi_1(s_1)(1 - P_2)} + \int_{\mathcal{A}_2} ds_2 \ \psi_2(s_2) \log \frac{\psi_2(s_2)}{\psi_2(s_2)(1 - P_1)} \tag{3.6}$$

$$= -P_1 \log(P_1) - P_2 \log(P_2). \tag{3.7}$$

We see that $0 \leq I(s_1, s_2) \leq \log 2$, where the lower bound is reached for the trivial case when one of the $P_i$ is null (meaning that one of the random variable is null) and the higher bound is reached when $P_1 = P_2 = 1/2$ (figure 3.1b). As expected, disjoint orthogonal random variables are uncorrelated but statistically dependent. Also, their mutual information does not depend on the function $\psi_i$ but only on their respective sparsity coefficients.

### 3.2.2 Blind separation of DO variables with $\ell_p$-norm optimization

In this subsection, we show that DO variables can be separated with an Infomax framework. It emphasizes that dependent variables can be separated by ICA-based methods.

**Theorem 8.** *Disjoint orthogonal sources can be separated with an Infomax network by using any generalized super-Gaussian (or leptokurtic, $p < 2$) non-linear function.*

*Proof.* We consider that the observation matrix $\boldsymbol{D} = \boldsymbol{AS}$ has been whitened. For two DO random variables, this means that a rotational ambiguity remains. The separation can be restrained to orthogonal matrices $\boldsymbol{W}$ for which $\det \boldsymbol{W} = \pm 1$. We can write

$$|\det \boldsymbol{J}| = \left| \det \boldsymbol{W} \prod_{n=1}^{N} G'(z_n) \right| = \prod_{n=1}^{N} G'(z_n). \tag{3.8}$$

**Figure 3.1:** *a) Joint PDF (blue color) of two disjoint orthogonal random variables and their marginal PDF (red color). b) The mutual information of two disjoint orthogonal variables, function of their sparsity repartition.*

The Infomax objective function (equation 2.52) can be written such that

$$\phi_{IM} \quad = \quad -\sum_{n=1}^{N} \mathbb{E}\left\{\log G'(z_n)\right\}. \tag{3.9}$$

Our study can also be restrained to orthogonal global mapping $\boldsymbol{H}$.

We now consider that $\{\boldsymbol{s}_n\}_{n=1}^{N}$ are DO random variables with sparsity coefficients $\alpha_n$ (probability of having a non-zero coefficient) such that $\sum_{n=1}^{N} \alpha_n = 1$. Because of this property, we can split the arguments inside the function $G'$ such as

$$\mathbb{E}\left\{\log G'(z_n)\right\} \quad = \quad \sum_{n'=1}^{N} \alpha_{n'}\mathbb{E}_{n'}\left\{\log G'(h_{nn'}s_{n'})\right\}, \tag{3.10}$$

where $\mathbb{E}_{n'}$ indicates that the considered expectation is computed for the active support of $s_{n'}$ only. Let us consider that the non-linear function $G$ is the CDF of a symmetric centered generalized normal distribution (equation 2.14) with a shape parameter $p$ and a scale parameter $\sigma$ such that

$$G'(z) = \underbrace{\frac{p}{2a\Gamma(1/p)}}_{=C_1} \exp\left[-\left(\frac{|z|}{a}\right)^p\right]. \tag{3.11}$$

This gives

$$\mathbb{E}\left\{\log G'(z_n)\right\} = \sum_{n'=1}^{N} \alpha_{n'}\mathbb{E}_{n'}\left\{\log C_1 - \frac{|h_{nn'}s_{n'}|^p}{a}\right\} \tag{3.12}$$

$$= \underbrace{\sum_{n'=1}^{N} \alpha_{n'}\log C_1}_{=C_2} - \frac{1}{\sigma}\sum_{n'=1}^{Q} \alpha_{n'}|h_{nn'}|\mathbb{E}_{n'}\left\{|s_{n'}|^p\right\}. \tag{3.13}$$

For a discrete vector $\boldsymbol{s}_{n'} \in \mathbb{R}_x^N$, we can write

$$\mathbb{E}_{n'}\{|s_{n'}|^p\} = \frac{1}{\alpha_{n'}N_x}\|\boldsymbol{s}_{n'}\|_p^p \tag{3.14}$$

**Figure 3.2:** *The blue circle indicates the $\ell_2$ unit-ball. The black squares indicates two $\ell_1$ balls, bounding the $\ell_2$ unit-ball.*

as the inactive support does not change the value of the norm. Then, we can express the Infomax objective function as

$$\phi_{IM} = -NC_2 + \frac{1}{aN_x} \sum_{n'=1}^{Q} N \|\boldsymbol{s}_{n'}\|_p^p \sum_{n=1}^{N} |h_{nn'}|^p \tag{3.15}$$

$$= -NC_2 + \frac{1}{aN_x} \sum_{n'=1}^{N} \|\boldsymbol{s}_n\|_p^p \|\boldsymbol{h}_n\|_p^p. \tag{3.16}$$

We see that $\phi_{IM}$ is a function of $\boldsymbol{H}$ only and is minimized when the $\|\boldsymbol{h}[n]\|_q$ are minimized. Because $\boldsymbol{H}$ is an orthogonal matrix, all the vectors $\boldsymbol{h}[q]$ have a unit $\ell_2$-norm and are located on the unit sphere. As shown in figure 3.2 for $\mathbb{R}^2$, the minima of the $\ell_p$-norm of the vectors $\boldsymbol{h}[q]$, with $p < 2$, are located on the axes. Then, the minimum of $\phi_{IM}$ is reached when $\boldsymbol{H}$ is a permutation matrix. ∎

Equivalently, we could also maximize $\phi_{IM}$ with a generalized Gaussian CDF having $p > 2$ in order to reach a maximum when $\boldsymbol{H} = \boldsymbol{\Pi}$. In practice, a sigmoid function $G_0$ is used in place of a generalized normal CDF. A sigmoid is a typical super-Gaussian function making a smooth transition between an $\ell_1$ and an $\ell_2$-norm solution [Batany et al., 2016b].

## 3.3 Inter-regressive processes

In several cases, the assumption of DO sources is not valid. Our main question is how much can this assumption be relaxed while assuming the sources as sparse signals. In other terms, how much can the sources overlap? We introduce in this section a new concept, named *inter-regressive processes*, that will be useful in the next sections. It can be seen as an extension of auto-regressive processes.

The properties describing a set of source signals $\boldsymbol{s}_1, \boldsymbol{s}_2, \ldots \boldsymbol{s}_N$ can be divided in two categories. A first category contains the properties of each single signal $\boldsymbol{s}_n$ taken independently from the others. For instance, the kurtosis is defined for a single source,

independently from the others. A second category contains the properties linked to the relations between several signals, i.e. by considering $\boldsymbol{s}[x]$. For instance, the covariance structure or the measures related to statistical dependence are properties defining the way all the sources interact between them.

The $\ell_p$ norms, defined as a distance, can belong to both categories. The norm of a single source $\|\boldsymbol{s}_n\|_p$ can be considered, as well as the norm $\|\boldsymbol{s}[x]\|_p$ of all of them at a single index $x$. In the next sections, we will mainly considered the norm of each single source, and especially the $\ell_0$ pseudo-norm $\|\boldsymbol{s}_n\|_0$.

Auto-regressive processes refer to the first category of signal properties. They have been extensively used in the signal processing literature and can be defined as follow [Mitra and Kaiser, 1993]. We emphasize that our definition is valid for a noiseless process.

**Definition 26.** *A signal $\boldsymbol{s}_i$ is said to be an auto-regressive process of order $D$ if it exists a set of $D + 1$ parameters $c_d$ such that*

$$\sum_{d=0}^{D} c_d s_i[x - d] = 0, \tag{3.17}$$

*where at least two parameters $c_d$ are non-null.*

For auto-regressive processes, the present $s_i[x]$ is fully determined by the past $s_i[x-1]$, $s_i[x-2]$, ... In other words, the convolution of the signal $\boldsymbol{s_i}$ by the filter $\boldsymbol{c}$ is equal to zero. From a geometric point of view, any set of $D + 1$ consecutive coefficients extracted from an auto-regressive process is located in a hyperplane in $\mathbb{R}^{D+1}$ defined by its normal vector $\boldsymbol{c} = \{c_d\}_{d=0}^{D}$.

We propose to extend the concept of *auto-regressive* processes to the second category of properties by introducing the concept of *inter-regressive* processes.

**Definition 27.** *A set of $N$ signals is said to be an inter-regressive process of order $D$ if it exists a set of $N \times (D + 1)$ parameters $c_{id}$ such that*

$$\sum_{i=1}^{N} \sum_{d=0}^{D} c_{id} s_i[x - d] = 0, \tag{3.18}$$

*where at least two parameters $c_{id}$ are non-null.*

From a geometric point of view, any set of $D + 1$ consecutive source vectors extracted from an inter-regressive process and forming a set of $N(D+1)$ coefficients is located in a hyperplane in $\mathbb{R}^{N(D+1)}$. For an inter-regressive process of order $D = 0$, equation 3.18 can be written as $\boldsymbol{c}^T \boldsymbol{s}[x] = 0$. When there is only one source, $N = 1$, an inter-regressive process becomes an auto-regressive process.

Both of the above definitions generally consider that equations 3.17 and 3.18 must be true for all indices $x$ or for a closed support of consecutive indices. For our purpose, we propose to expand these definitions. A signal is said to yield an auto-regressive process of order $D$ and length $E \in \mathbb{N}$ if equation 3.17 is true for $E$ indices $x$, possibly not consecutive. Equivalently, signals are said to yield an inter-regressive process of order $D$ and length $E \in \mathbb{N}$ if equation 3.18 is true for $E$ indices $x$, possibly not consecutive.

Figure 3.3 shows the construction of both auto-regressive and inter-regressive processes of hypothetical length $E = 2$. For the auto-regressive process (figure 3.3a), each red square is fully determined by the value of all antecedent grey dots, with a unique set of parameters $\{c_d\}$. When a signal holding an auto-regressive process is convolved with the filter $\boldsymbol{c}$, it results zero-valued coefficients (blue squares). For the inter-regressive

a) Auto-regressive process



b) Inter-regressive process



**Figure 3.3:** *Examples of construction of a) an auto-regressive process of order $D = 7$ and length $E = 2$ and b) an inter-regressive process of order $D = 5$ and length $E = 2$ for $N = 3$ signals. For both a) and b), the result of the convolution between the source or the source vector with the set of coefficients defining the inter-regressive process is shown. Blue squares indicate zero-value coefficient.*

process (figure 3.3b), each red square is fully determined by the value of all antecedent and concomitant grey dots, with a unique set of parameters $\{c_{id}\}$. When a set of sources holding an auto-regressive process are convolved with the filter $\boldsymbol{c}$, it results zero-valued coefficients (blue squares).

**Figure 3.4:** *Example of three sources having an inter-regressive process of order* $D = 0$ *and length* $E = 6$. *In the right hand figure, the grey dots circled in black in the source space belong to this inter-regressive process and are inside the same hyperplane.*

Figure 3.4 shows an example of an inter-regressive process of order $D = 0$ and length $E = 6$ for three sources. The position of each sample is shown in the source space. The samples belonging to the inter-regressive process are located inside the same hyperplane.

## 3.4   Blind source extraction of the sparsest source

In this section, we discuss necessary and sufficient conditions, on the sources only, to use the $\ell_0$ pseudo-norm as a contrast function in linear BSE problems, for both instantaneous and convolutive mixtures. In other words, we discuss the conditions under which the solution of the $\ell_0$ pseudo-norm minimization problem

$$\{\boldsymbol{w}_l\}^* = \arg\min_{\{\boldsymbol{w}_l\}} \left\| \boldsymbol{u}[x] = \sum_{l=0}^{L} \boldsymbol{w}_l^T \boldsymbol{x}[x - l] \right\|_0, \tag{3.19}$$

extracts a signal $\boldsymbol{y}^* \in \mathcal{S}_1$ corresponding to the recovery of the sparsest source.

**Assumption 4.** *Without loss of generality, we consider that the sources are sorted in order of decreasing sparsity such that*

$$\|\boldsymbol{s}_1\|_0 < \|\boldsymbol{s}_2\|_0 \leq \cdots \leq \|\boldsymbol{s}_N\|_0. \tag{3.20}$$

*We emphasize that* $\|\boldsymbol{s}_1\|_0 < \|\boldsymbol{s}_i\|_0$, $\forall i \neq 1$, *to avoid any competition between the extraction of the sparsest source* $\boldsymbol{s}_1$ *and another source.*

Because only the $\ell_0$ pseudo-norm is considered in the following developments, there is no way to distinguish between two sources that have the same $\ell_0$ pseudo-norm. This is the reason why we need to consider $\|\boldsymbol{s}_1\|_0 < \|\boldsymbol{s}_i\|_0$, $\forall i \neq 1$. If this condition is not valid, then the following theorems fails. For the sake of clarity, the instantaneous case is treated first and then the generalization to the convolutive case is presented. Both proofs are similar.

### 3.4.1 Instantaneous mixtures

**Theorem 9.** *The $\ell_0$ pseudo-norm is a contrast function for the extraction of the sparsest source $\boldsymbol{s}_1$ if and only if the sources do not have any inter-regressive process of order $0$ with a length higher than or equal to the size of the inactive support of $\boldsymbol{s}_1$.*

*Proof.* From definitions 23 and 24, the $\ell_0$ pseudo-norm is a contrast function for the extraction of $\boldsymbol{s}_1$ if and only if

$$\|\boldsymbol{s}_1\|_0 < \|\boldsymbol{u}\|_0 \qquad \forall \boldsymbol{u} \notin \mathcal{S}_1,$$

i.e. if and only if

$$\|\boldsymbol{s}_1\|_0 < \#\{u[x] = \boldsymbol{h}^T \boldsymbol{s}[x] : u[x] \neq 0\} \qquad \forall \boldsymbol{h} \notin \mathcal{G}_1,$$
$$\|\boldsymbol{s}_1\|_0 < N - \#\{u[x] = \boldsymbol{h}^T \boldsymbol{s}[x] : u[x] = 0\} \qquad \forall \boldsymbol{h} \notin \mathcal{G}_1,$$
$$N - \|\boldsymbol{s}_1\|_0 > \#\{u[n] = \boldsymbol{h}^T \boldsymbol{s}[x] : u[x] = 0\} \qquad \forall \boldsymbol{h} \notin \mathcal{G}_1,$$
$$N - \|\boldsymbol{s}_1\|_0 > E^*,$$

where we defined

$$E^* = \max \left[ \#\{u[x] = \boldsymbol{h}^T \boldsymbol{s}[x] : u[x] = 0\}, \ \boldsymbol{h} \notin \mathcal{G}_1 \right].$$

$N - \|\boldsymbol{s}_1\|_0$ is the number of null values of $\boldsymbol{s}_1$. $E^*$ is the maximum length of an inter-regressive process of order 0 among the sources. ∎

From theorem 9, one can re-derive the sufficiency of the condition proposed by Duarte et al. [2011] (see appendix 8.3 for details).

### 3.4.2 Convolutive mixtures

**Theorem 10.** *The $\ell_0$ pseudo-norm is a contrast function for the extraction of the sparsest source $\boldsymbol{s}_1$ if and only if the sources do not have any inter-regressive process of order $J = K + L$ with a length higher than or equal to the size of the inactive support of $\boldsymbol{s}_1$.*

*Proof.* From definitions 23 and 24, the $\ell_0$ pseudo-norm is a contrast function for the extraction of $\boldsymbol{s}_1$ if and only if

$$\|\boldsymbol{s}_1\|_0 < \|\boldsymbol{u}\|_0 \qquad \forall \boldsymbol{u} \notin \mathcal{S}_1,$$

i.e. if and only if

$$\|\boldsymbol{s}_1\|_0 < \#\{u[x] = \sum_{j=0}^{J} \boldsymbol{h}_j^T \boldsymbol{s}[x-j] : u[x] \neq 0\} \qquad \forall \boldsymbol{h} \notin \mathcal{G}_1,$$

$$\|\boldsymbol{s}_1\|_0 < N - \#\{u[x] = \sum_{j=0}^{J} \boldsymbol{h}_j^T \boldsymbol{s}[x-j] : u[x] = 0\} \qquad \forall \boldsymbol{h} \notin \mathcal{G}_1,$$

$$N - \|\boldsymbol{s}_1\|_0 > \#\{u[x] = \sum_{j=0}^{J} \boldsymbol{h}_j^T \boldsymbol{s}[x-j] : u[x] = 0\} \qquad \forall \boldsymbol{h} \notin \mathcal{G}_1,$$

$$N - \|\boldsymbol{s}_1\|_0 > E^*,$$

where we defined

$$E^* = \max \left[ \#\{ u[x] = \sum_{j=0}^{J} \boldsymbol{h}_j^T \boldsymbol{s}[x-j] : u[x] = 0\}, \ \boldsymbol{h} \notin \mathcal{G}_1 \right].$$

$N - \|\boldsymbol{s}_1\|_0$ is the number of null values of $\boldsymbol{s}_1$. $E^*$ is the maximum length of an inter-regressive process of order $J = K + L$ among the sources.                    ■

In both theorems 9 and 10, the assumption of W-disjoint orthogonality of the sources is not necessary and the sources can overlap. This will be shown in the next section. SCA is often presented as a method able to separate signals violating the independence assumption. We emphasis here that this assertion is true below the limit defined by theorems 9 and 10: the limit for SCA-based BSE is one kind of strong linear dependency among the sources, named an inter-regressive process.

## 3.5  Blind source separation of sparse sources

In the previous sections, we focused on the blind extraction of the sparsest source in both instantaneous and convolutive mixtures. Our analysis was based on the minimization of the $\ell_0$ pseudo-norm and we discussed the conditions under which this pseudo-norm can be used as a contrast function for a correct blind extraction. Delfosse and Loubaton [1995] show that a BSS problem can be solved as a sequence of BSE problems in which each source is extracted after the other [Papadias, 2000]. However, this method can be sensible to the noise level because each step supposes a perfect extraction.

It is valuable to specify an objective function that is a contrast function for the separation. A first idea is to build a function such that

$$\min_{\boldsymbol{W}} \quad \sum_i \|\boldsymbol{u}_i\|_0 + \epsilon \ \phi_r(\boldsymbol{u}_1, \boldsymbol{u}_2, \dots), \tag{3.21}$$

where the term $\phi_r$ guaranties that the same source cannot be extracted twice. Without such a term, the sparsest source is extracted several time. This term could be a robust measure of the rank of the matrix $\boldsymbol{W}$. It could be also a penalized function of the correlation of the separated sources making sure that two sources are not too correlated. In practice, this term may be difficult to specify, especially the trade-off parameter $\epsilon$.

The Hadamard product given by $\boldsymbol{u} = \boldsymbol{u_1} \odot \boldsymbol{u_2}$ is defined such that $u[x] = u_1[x] \times u_2[x]$. We propose a different objective function based on the Hadamard products of the extracted sources. We define

$$\boldsymbol{\Upsilon} = \begin{bmatrix} \|\bar{\boldsymbol{u}}_1 \odot \bar{\boldsymbol{u}}_1\|_0 & \|\bar{\boldsymbol{u}}_1 \odot \bar{\boldsymbol{u}}_2\|_0 & \cdots \\ \|\bar{\boldsymbol{u}}_2 \odot \bar{\boldsymbol{u}}_1\|_0 & \|\bar{\boldsymbol{u}}_2 \odot \bar{\boldsymbol{u}}_2\|_0 & \cdots \\ \vdots & & \ddots \end{bmatrix}. \tag{3.22}$$

The diagonal of the matrix $\boldsymbol{\Upsilon}$ contains the $\ell_0$ pseudo-norm as we have

$$\Upsilon_{ii} = \|\bar{\boldsymbol{u}}_i \odot \bar{\boldsymbol{u}}_i\|_0 = \|\bar{\boldsymbol{u}}_i\|_0, \tag{3.23}$$

and the off-diagonal terms are the cross-Hadamard-products between any pair of extracted vectors. Equation 3.21 mainly considers the trace of $\boldsymbol{\Upsilon}$. We propose to build an objective function based on the off-diagonal terms such as

$$\min \quad \phi_{\Upsilon} = \sum_i^N \sum_{j=i+1}^N \Upsilon_{ij}. \tag{3.24}$$

This method is also known as disjoint component analysis (DCA) [Anemüller, 2007; Mei and Mertins, 2008; Nose-Filho, 2015].

The function $\phi_\Upsilon$ is the sum of the off-diagonal terms. The second sum is made such that only the upper or lower part of the matrix is considered because $\boldsymbol{\Upsilon}$ is symmetric.

We start with the trivial case of disjoint orthogonal sources and propose a sufficient condition. The general case is discussed in hereafter.

**Proposition 1.** *We consider the set of global mapping vectors. The following equivalence holds*

$$\boldsymbol{h}_i \odot \boldsymbol{h}_j = \boldsymbol{0} \quad \forall i, j \neq i \quad \Leftrightarrow \quad \{\boldsymbol{h}_i\} \in \mathcal{G}. \tag{3.25}$$

**Theorem 11.** *The function $\phi_\Upsilon$ (defined in equation 3.24) is a contrast function for the blind separation of disjoint orthogonal sources.*

*Proof.* For disjoint orthogonal sources, we have $\phi_\Upsilon^* = 0$ as all the off-diagonal terms of $\boldsymbol{\Upsilon}$ vanish. Any other set of extracted vectors is a linear combination of at least two sources. We have

$$\phi_\Upsilon = \sum_i \sum_j \left\| \bar{\boldsymbol{u}}_i \odot \bar{\boldsymbol{u}}_j \right\|_0 = \sum_i \sum_j \left\| \boldsymbol{h}_i^T \boldsymbol{S} \odot \boldsymbol{h}_j^T \boldsymbol{S} \right\|_0. \tag{3.26}$$

Because the sources are disjoint orthogonal, all the cross-products vanish and only the auto-products remain. We can write

$$\phi_\Upsilon = \sum_i \sum_j \left\| \left[ \boldsymbol{h}_i \odot \boldsymbol{h}_j \right]^T \left[ \boldsymbol{S} \odot \boldsymbol{S} \right] \right\|_0. \tag{3.27}$$

By using proposition 1, we see that at least one term must be non-zero and we have $\phi_\Upsilon > \phi_\Upsilon^* = 0$. The sufficiency of the disjointness condition is proved. ∎

**Conjecture 1.** *The necessary and sufficient conditions for the separation of sources with DCA are the same as the necessary and sufficient conditions for the extraction of all sources in SCA.*

In appendix 8.4, we prove this conjecture for the case with $N = 2$ sources and $M = 2$ observations. However, we have not been able to find a proper proof for the general case. We have performed tests on synthetic data, that seem to show that the presence of an inter-regressive process is a theoretical limit. In the next section, we will present an evolutionary algorithm for solving equations 3.19 and 3.24.

## 3.6 Differential evolution algorithm

The two contrast functions presented in equations 3.19 and 3.24 (respectively for BSE and BSS) are both using the $\ell_0$ pseudo-norm. Therefore, they are non robust to the presence of noise in the observation or the extraction (separation) of almost sparse signal. Solving these problems is combinatorial and intractable in practice with large data. However, different approximations of the $\ell_0$ pseudo-norm exist. In this work, we do not use the equivalence between the $\ell_0$ pseudo-norm and the $\ell_1$-norm, because the desired signal may not be sufficiently sparse to validate the sparsity assumption needed for an equivalence (see, for instance, the example in the introductory chapter 1). Instead, we use the smooth version proposed by Naini et al. [2007] and named SL0 (smooth $\ell_0$). It is defined for a vector $\boldsymbol{u} \in \mathbb{R}^N$ with Gaussian kernel such that

$$\|\boldsymbol{u}\|_{0,\sigma} = N - \sum_x e^{-\bar{u}[x]^2/2\sigma^2}, \tag{3.28}$$

**Figure 3.5:** *Shape of several norms in $\mathbb{R}^2$. The first two lines shows the $\ell_p$-norms for different values of p. The third line shows the smooth $\ell_0$ pseudo-norm for different values of the shaping parameter $\sigma$.*

| Parameter | Notation | Value |
|---|---|---|
| Mutation | $F$ | $[0.4, 1]$ |
| Crossover | $Cr$ | $[0, 1]$ |
| Number of individuals | $N_{pop}$ | $10\times N \times M \times Z$ |

**Table 3.1:** *Key parameters of the DE algorithm.*

where $\sigma$ is a shaping parameter and the vector $\boldsymbol{u}$ is normalized to have unit $\ell_2$ norm. We have

$$\lim_{\sigma \to 0} \|\boldsymbol{u}\|_{0,\sigma} = \|\boldsymbol{u}\|_0 \,, \tag{3.29}$$

showing that we can play with the parameter $\sigma$ in order to adapt the contrast function to be more robust to the noise. Nevertheless, the approximation of the functions described in equations 3.19 and 3.24 with SL0 are still highly non-convex and a lot of local minima exist. In this section, we propose a differential evolution (DE) algorithm [Das and Suganthan, 2001] for solving the approximations of equations 3.19 and 3.24. We use this algorithm to show the veracity of our previous theorems. Compared to other evolutionary algorithm [Zhou et al., 2011], DE has really few number of parameters (see table 3.1). Figure 3.5 shows the shape of several $\ell_p$-norms for both $p \geq 2$ and $p < 2$, as well as the shape of the SL0 approximation.

### 3.6.1 Generalities about DE

Differential evolution (DE) is a stochastic evolutionary optimization algorithm initially proposed by Storn and Price [1997]. As any other stochastic algorithm, DE is particularly efficient for optimizing non-convex and multi-modal functions for which traditional gradient-based methods can only provide local optima.

A lot of evolutionary meta-heuristics have been developed in the last decades [Deb, 2001]. Basically, they are all based on a population of individuals. Each individual represents a candidate solution. At each generation, individuals are perturbed, mixed and selected until a convergence criterion is reached. The basic DE algorithm (as for most evolutionary scheme) can be summarized as

> Initialization;
> **while** *DE not converged* **do**
> | Mutation;
> | Crossover;
> | Selection
> **end**

The number of individual is denoted $N_{pop}$. In the initialization step, the full parameter space should be uniformly covered. The following notation is used for designating DE algorithms: $DE/a/b/c$, where $a$ refers to the mode of selection of individual to be perturbated (*best* or *random*), $b$ is an integer that refers to the number of differences used (see equation 3.31 hereafter) and $c$ refers to the cross-over method (*exp*, *bin* or *dir*) [Storn and Price, 1997; Yang et al., 2015].

#### Diversity in evolutionnary meta-heuritstics

The diversity of the population is a key property in all evolutionary algorithms. A population with high diversity means that the individual are far away from each other, in average. If all individuals are concentrated in a local area of the parameter space, we say that the population has lost its diversity. If this happen, the mutation and crossover processes are not effective any more to create new original individuals. This means that there is low chance to be able to cover the entire parameter space. When the population has lost its diversity, we observe a high "consanguinity" between the parents and the children. In other terms, the population is stuck in a local minima. Several measures exist for monitoring the diversity of the population and a lot of procedures exist to re-inject diversity in the population [Yang et al., 2015].

### 3.6.2 DE for blind source extraction

In the DE algorithm for instantaneous BSE, each individual is an extracting vector $\boldsymbol{w}$ acting on the observations and giving an extracted vector $u[x] = \boldsymbol{w}^T \boldsymbol{d}[x]$. There is a fixed number of $N_{pop}$ individuals at each generation. The upper script indicates the arbitrary index of the individual as

$$\text{population:} \quad \begin{bmatrix} \boldsymbol{w}^1 \\ \boldsymbol{w}^2 \\ \vdots \\ \boldsymbol{w}^k \\ \vdots \\ \boldsymbol{w}^{N_{pop}} \end{bmatrix}. \tag{3.30}$$

For each individual $\boldsymbol{w}^k$, a mutant vector $\boldsymbol{r}^k$ is created by combining three randomly selected other individuals in the population such that

$$\boldsymbol{r}^k = \boldsymbol{w}^{i_1^k} + F \times (\boldsymbol{w}^{i_3^k} - \boldsymbol{w}^{i_2^k}), \qquad i_1^k, i_2^k, i_3^k \neq k. \tag{3.31}$$

For each individual $\boldsymbol{w}^k$, the mutant vector $\boldsymbol{r}^k$ is used to create a trial vector $\boldsymbol{u}^k$ such that

$$\boldsymbol{u}^k[x] = \begin{cases} \boldsymbol{r}^k[x] & \text{with probability } Cr, \\ \boldsymbol{w}^k[x] & \text{with probability } 1 - Cr. \end{cases} \tag{3.32}$$

For convolutive BSE, each individual represents an extracting system, i.e. a set of vectors

$$\text{population:} \quad \begin{bmatrix} \{\boldsymbol{w}_l^1\}_{l=0}^L \\ \{\boldsymbol{w}_l^2\}_{l=0}^L \\ \vdots \\ \{\boldsymbol{w}_l^k\}_{l=0}^L \\ \vdots \\ \{\boldsymbol{w}_l^{N_{pop}}\}_{l=0}^L \end{bmatrix}. \tag{3.33}$$

**Constraining the search space**

When based on sparsity, instantaneous BSE problems have a scale ambiguity and convolutive BSE problems also have a shift ambiguity. Those ambiguities are taken into account in the formulation of the solution spaces $\mathcal{S}_n$ corresponding to a correct extraction. We can also use these ambiguities to reduce the search space of the DE algorithm and improve performance. In particular, this avoids the population to identify twice the same solution.

Vectors $\boldsymbol{w}^k$ are systematically projected back on half of the unit hypersphere. For instance in $\mathbb{R}^2$, the lower arc can be discarded, and in $\mathbb{R}^3$ the half lower semi-sphere can also be discarded. We see that the size of the search space is divided by 2. The spherical coordinates in $\mathbb{R}^N$ gives a parametrization

$$\begin{aligned} w_1 &= \cos\theta_1 \\ w_2 &= \sin\theta_1 \cos\theta_2 \\ w_3 &= \sin\theta_1 \sin\theta_2 \cos\theta_3 \\ &\vdots \\ w_{N-1} &= \sin\theta_1 \sin\theta_2 \cdots \cos\theta_{N-1} \\ w_N &= \sin\theta_1 \sin\theta_2 \cdots \sin\theta_{N-1}, \end{aligned} \tag{3.34}$$

where $\theta_{N-1} \in [0, 2\pi]$ and $\theta_1 \cdots \theta_{N-2} \in [0, \pi]$ are given by

$$\begin{aligned} \theta_1 &= \arccos \frac{w_1}{\|[w_1 \ldots w_N]\|_2} \\ \theta_2 &= \arccos \frac{w_2}{\|[w_2 \ldots w_N]\|_2} \\ &\vdots \\ \theta_{N-2} &= \arccos \frac{w_{N-2}}{\|[w_{N-2} \ldots w_N]\|_2} \\ \theta_{N-1} &= \begin{cases} \arccos \frac{w_{N-1}}{\|[w_{N-1} \ w_N]\|_2} & w_n \geq 0, \\ 2\pi - \arccos \frac{w_{N-1}}{\|[w_{N-1} \ w_N]\|_2} & w_n < 0. \end{cases} \end{aligned} \tag{3.35}$$

We can simply constrain the search by limiting $\theta_{N-1}$ in the interval $[0, \pi]$. Once the vector $\boldsymbol{w}$ has been normalized, we have to find its correct sign inside the limited search space.

The shift ambiguity is more tricky to handle, especially for keeping similar solutions close to each other in the search space. For instance the two following extracting systems

$$\boldsymbol{w}[z] = \begin{bmatrix} 1 & 0 & 0 & 0 \\ 1 & 0 & 0 & 0 \\ 1 & 0 & 0 & 0 \end{bmatrix} \quad \text{and} \quad \boldsymbol{w}[z] = \begin{bmatrix} 0 & 1 & 0 & 0 \\ 0 & 1 & 0 & 0 \\ 0 & 1 & 0 & 0 \end{bmatrix}$$

are equivalent as they extract the same signal, but shifted by one index. Our parametrization does not take into account the shift ambiguity. This problem that we identify as a topological problem is out of our scope but we want to give here an idea of the problem. It is important to notice that as we do not tackle the shifting ambiguity, the search space may be redundant and we loose the unicity of the solution. For instance, let us take the shifting ambiguity in $\mathbb{R}^3$. First the three points should be connected. The vector $[0, 1, 1]$ should be close to $[1, 1, 0]$. In other word, two exact same solutions belonging to $\mathcal{S}$ can be far away from each other.

**Stoping convergence criterion for DE**

A simple convergence criterion could be a maximum number of iteration, i.e. the DE algorithm stops after $N_{ite}$ generations. However, this approach can lead to unnecessary computations, as the population could have reached the minimum before. Finding an adequate criterion for stoping the DE algorithm is not trivial.

A first idea can be to say that the individuals of a DE algorithm converge with high probability to the same location after a certain time. This aspect can be useful for proposing a convergence criterion. For instance, we could measure the average distance between all individuals and stop the algorithm when this value is small. However, this approach is computationally costly because it needs the computation of half a square distance matrix at each generation, containing all the distances $\left\| \boldsymbol{w}^{k_1} - \boldsymbol{w}^{k_2} \right\|_2$ for all couple $k_1, k_2 \in [1, N_{pop}]$.

Also, the population could have reached several minimum location and the individuals are not necessarily close to each other. We prefer to use a simple criterion, easy to compute, which is the stationary of the best element in the population. If the value of the fitness of the best element is constant over $N_\Delta$ generations, we say that the DE has converged.

**A simple gradient step**

The DE algorithm is efficient for computing efficiently a good approximation of the solution. To increase performance, we add a gradient descent step to the best individual. The gradient of the objective function with respect to a parameter is given by

$$\frac{\partial \left\| \boldsymbol{u} \right\|_{0,\sigma}}{\partial w_z[m]} = \frac{1}{\sigma^2} \sum_{x=1}^{X} d_m[x - z] u[x] \mathrm{e}^{-u[x]^2 / 2\sigma^2}. \tag{3.36}$$

The best individual is updated in opposite direction of the gradient until convergence.

**Comments on initialization**

For a blind problem, no assumption is made on the mixing model and by consequence no assumption can be made on the extracting vector (or the separation matrix). The

**Figure 3.6:** *a) A uniform initialization of the population over the unit cube. b) The projection of this initialization on the $\ell_2$ unit-sphere.*

initial population should uniformly cover the parameter space with all the individuals. A naive idea is to cover the parameter space by taking each parameter from a uniform distribution, for instance in the interval $[-1, +1]$. However this naive method leads to an incorrect initialization.

For instance, the extraction is done by a vector $\boldsymbol{w} \in \mathbb{R}^3$ that can be parametrized in spherical coordinates by 2 parameters $\theta_1$ and $\theta_2$ in order to get rid of the scale ambiguity. Then, the uniform initialization must be done on the unit-sphere and not on the unit cube. Figure 3.6a shows an uniform initialization of the population in a $\mathbb{R}^3$ parameter space. Figure 3.6b displays the projection of this initialization after projection on the $\ell_2$ unit-sphere. Clearly, we see areas in which the density of individuals is higher than elsewhere. This initialization is not correct.

The correct way to initialize uniformly the vector over the unit sphere has been proposed by Muller [1959]. Let all $\mathrm{w}_n$ be independent Gaussian random variables, then $\mathbf{w}/\|\mathbf{w}\|_2$ is uniform over the unit sphere.

### 3.6.3 DE for blind source separation

In the DE algorithm for instantaneous BSS, each individual is an extracting matrix $\boldsymbol{W}$ acting on the observations and giving a set of extracted vector $\boldsymbol{u}[x] = \boldsymbol{W}^T \boldsymbol{d}[x]$. There is a fixed number of $N_{pop}$ individuals at each generation. The upper script indicates the arbitrary index of the individual as

$$\text{population:} \quad \begin{bmatrix} \boldsymbol{W}^1 \\ \boldsymbol{W}^2 \\ \vdots \\ \boldsymbol{W}^k \\ \vdots \\ \boldsymbol{W}^{N_{pop}} \end{bmatrix}. \tag{3.37}$$

For convolutive BSS, each individual represents a separating system, i.e. a set of

vectors

$$
\text{population:} \quad
\begin{bmatrix}
\{\boldsymbol{W}_l^1\}_{l=0}^L \\
\{\boldsymbol{W}_l^2\}_{l=0}^L \\
\vdots \\
\{\boldsymbol{W}_l^k\}_{l=0}^L \\
\vdots \\
\{\boldsymbol{W}_l^{N_{pop}}\}_{l=0}^L
\end{bmatrix} .
\tag{3.38}
$$

### 3.6.4 A strategy for dealing with noisy data

In the previous sections, we mainly considered the noiseless mixture model $\boldsymbol{d}[x] = \boldsymbol{A}\boldsymbol{s}[x]$ of equation 2.35, where the sparse sources were containing a lot of exact zero coefficients. However, in practice, observations may be contaminated with independent noise and the sources may not be perfectly sparse.

We model the noisy data by the following equation

$$
\boldsymbol{d}[x] = \boldsymbol{A}\boldsymbol{s}[x] + \boldsymbol{n}[x],
\tag{3.39}
$$

where $\boldsymbol{n}[x]$ is the noise vector. This additive noise can be, for instance, a random Gaussian vector with zero mean vector $\boldsymbol{\mu}_1 = \boldsymbol{0}$ and a diagonal variance matrix $\boldsymbol{\Sigma} = \boldsymbol{\sigma}^T \boldsymbol{I}$ representing uncorrelated components. Almost sparse sources contain a lot of small values and we consider the following model

$$
\boldsymbol{s}[x] = \boldsymbol{s}_0[x] + \boldsymbol{\delta}[x],
\tag{3.40}
$$

where $\boldsymbol{s}_0[x]$ is a perfectly sparse source vector containing a lot of exact zero-values and $\boldsymbol{\delta}[x]$ contains all the small perturbations.

The separation of sources is performed by a separating matrix $\boldsymbol{W}$ as in equation 2.37. It gives with the two precedent models

$$
\boldsymbol{u}[x] = \boldsymbol{W}^T \boldsymbol{d}[x] = \boldsymbol{W}^T \left( \boldsymbol{A} \left( \boldsymbol{s}_0[x] + \boldsymbol{\delta}[x] \right) + \boldsymbol{n}[x] \right)
\tag{3.41}
$$

$$
= \boldsymbol{W}^T \boldsymbol{A} \boldsymbol{s}_0[x] + \underbrace{\boldsymbol{W}^T \boldsymbol{A} \boldsymbol{\delta}[x] + \boldsymbol{W}^T \boldsymbol{n}[x]}_{\boldsymbol{e}_1[x] + \boldsymbol{e}_2[x]},
\tag{3.42}
$$

where the two terms $\boldsymbol{e}_1[x] = \boldsymbol{W}^T \boldsymbol{A} \boldsymbol{\delta}[x]$ and $\boldsymbol{e}_2[x] = \boldsymbol{W}^T \boldsymbol{n}[x]$ are two residual terms due to the almost sparsity of the sources and the noise, respectively.

Both terms lead to a deviation from the perfect sparse model. The hyperplanes due to the presence of zeros are not crossing at a unique point anymore, but are approximating each other around a small area, as we saw in figure 2.1 with the noisy inverse problem. The smooth $\ell_0$ pseudo-norm is particularly useful to handle noisy observation via the shaping parameter $\sigma$. If the level of noise is known, Naini et al. [2007] propose to use a fixed value of $\sigma$. Unfortunately, the level of noise is not always known a priori. We propose a strategy for identifying the presence of noise in the data or in the model by the analysis of a Pareto curve [Kim and Weck, 2005; Campigotto et al., 2014].

Ideally, we want to take a small value for $\sigma$ to be as close as possible from the exact $\ell_0$ pseudo-norm estimate. However, the presence of noise or almost sparse sources requires to relax and increase this value. If $\sigma$ is taken too high, the solution is close to the $\ell_2$ norm solution and is not accurate enough. We propose to evaluate the best $\sigma$ via the analysis of the function

$$
\phi_0^*(\sigma),
\tag{3.43}
$$

where $\phi_0^*$ is the value of the global minimum estimated by the DE algorithm for a fixed $\sigma$.

With a perfectly sparse and noiseless model, the function $\phi_0^*(\sigma)$ presents a continuous sigmoid shape. For small values of $\sigma$, $\phi_0^*$ tends to be close to $\|s_1\|_0$ because only the coefficient closed to zero are contributing to the computation. For high values of $\sigma$, all values contribute and the $\phi_0^*$ is closed to zero.

When noise is present in the observation, the Pareto curve of $\phi_0^*$ changes its shape into a specific form. Instead of presenting a single stair, it shows two stairs. For low $\sigma$, $\phi_0^*$ is close to the number of indexes, i.e. to the length of the signals. The two stair shapes are typical of noisy data. The best $\sigma$ can be evaluated by

$$\sigma^* = \min_{\sigma} \quad \phi_0^*(\sigma) + \epsilon \ \log \sigma. \tag{3.44}$$

This second optimization problem must be carefully addressed because the Pareto curve is not convex. In particular, the search must be limited to an interval of small $\sigma$ values. This method is expensive as several runs of the DE algorithm must be achieved.

## 3.7 Synthetic examples for BSE and BSS

In this section, we develop and analyze several synthetic examples. Basic examples help to understand the structure of the cost functions defined in equation 3.19 and 3.24, especially the location of local and global minima. They are also of useful help for understanding the behavior of the DE algorithm. Subsection 3.7.1 gives the first example with $N = 2$ sources and $M = 2$ observations in instantaneous mixture. This simple example is particularly appreciable because both the BSE and the BSS objective functions can be displayed. Subsection 3.7.2 considers $N = 3$ sources and $M = 3$ observations in instantaneous mixture. Only the BSE objective function can be displayed. Subsection 3.7.3 present a problem with $N = 2$ sources and $M = 3$ in a convolutive mixture. These examples will confirm our theorems.

### 3.7.1 $2 \times 2$ instantaneous mixture

**Sources and mixtures.** The first example is a simple instantaneous BSS problem with $N = 2$ sources and $M = 2$ sources. The two sources are 1D signals with $\|s_1\|_0 = 56\%$ and $\|s_2\|_0 = 75\%$. Sparsity is expressed in percentage, as the ratio of the number of active coefficients over the size of the signal. In particular, they are not disjoint orthogonal and share a common support of $30\%$ in which a $12.5\%$ inter-regressive process is added. An IR process is added by randomly choosing the coefficients $[c_1, c_2]$ of the process and constraining $c_1 s_1[x] + c_2 s_2[x] = 0$. Figure 3.7a displays the two signals and figure 3.7b shows their cross-plot. In the cross-plot, the samples are aligned along lines corresponding to the silent zone of $s_1$ (vertical) and $s_2$ (horizontal). The third line corresponds to the inter-regressive process, where $s_2[x] \propto s_1[x]$.

The mixing system is set such that

$$\boldsymbol{A} = \begin{bmatrix} 2 & -1 \\ 1 & 1 \end{bmatrix}. \tag{3.45}$$

**Extracting and separating systems.** The extraction is obtained by optimizing the following extracting vector

$$\boldsymbol{w}^T = \begin{bmatrix} w_1 & w_2 \end{bmatrix} = \begin{bmatrix} \cos\theta_1 & \sin\theta_1 \end{bmatrix}, \tag{3.46}$$

**Figure 3.7:** *First synthetic example with $N = 2$ sources. a) $s_1$ and $s_2$ are sparse signals but not disjoint. b) Cross-plot of $s_1$ and $s_2$.*



**Figure 3.8:** *The SL0 objective function defined in equation 3.19 for the BSS of sparse sources with different shaping parameter $\sigma$.*

where the amplitude ambiguity is used to reduce the number of parameter from 2 to 1. The parameter space can be limited to $[0 \ \pi[$. The separation is obtained by optimizing the following separating matrix

$$\boldsymbol{W} = \begin{bmatrix} \boldsymbol{w}_1^T \\ \boldsymbol{w}_2^T \end{bmatrix} = \begin{bmatrix} \cos\theta_1 & \sin\theta_1 \\ \cos\theta_2 & \sin\theta_2 \end{bmatrix}, \tag{3.47}$$

where each extracting vector can be parametrized by a single parameter $\theta_i$. The parameter space can be limited to $[0 \ \pi[\times[0 \ \pi[$.

**Results and discussion.** Figure 3.8 shows the contrast function 3.19 for the extraction of the sparsest sources, for different values of the shaping parameter $\sigma$. Three local minima exist: two are associated to the extraction of each source and one is due to the presence of an inter-regressive process. When $\sigma$ is really small, the objective function becomes closer to the exact $\ell_0$ contrast function and a lot of small parasitic minima appear. All those parasitic local minima are due to the cancellation of at least one coefficient in the extracted vector. Strong local minima correspond to the cancellation of several coefficients: each corresponds to the extraction of a sparse source. Also, when $\sigma$ is too large, the exact minimum corresponding to the extraction of $s_2$ is merged with the parasitic minima.

The challenge of BSS compared to BSE is the direct recovery of all sources at the same time. As we explained before, we could try to identify all the local minima of the BSE contrast function. It will require a multimodal optimization scheme for

**Figure 3.9:**  *a) The naive BSS objective function based on the $\ell_0$ pseudo norm (equation 3.21). b) The correct BSS objective function for the separation of sparse sources as defined in equation 3.24.*



**Figure 3.10:**  *The SL0 objective function defined in equation 3.24 for the BSS of sparse sources with different shaping parameters $\sigma$.*

identifying the two local minima. The use of a contrast function for the separation is a valuable approach as it simplify the optimization. The naive BSS objective function (equation 3.21) is the sum of the $\ell_0$ pseudo-norm of the two extracted vectors. The result is presented in figure 3.9-a). Clearly, this function is symmetric with respect to the main diagonal. One can observes several vertical and horizontal lines, each corresponding to the cancellation of one coefficient in 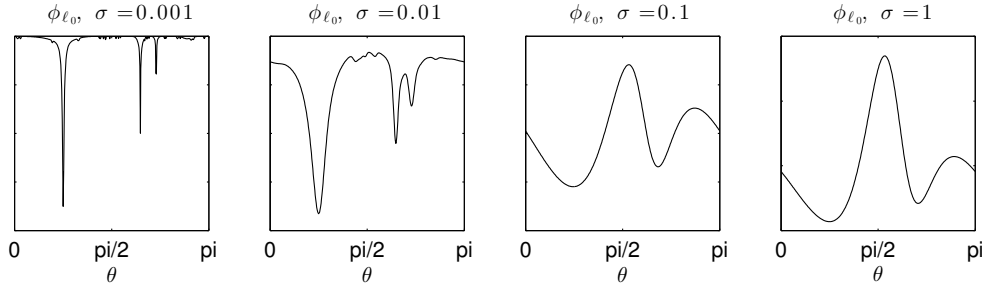the extracted vector. Two strong lines correspond to the extraction of $s_1$ and $s_2$. Without any regularization term as proposed in equation 3.24, we see that the global minimum correspond to the extraction of twice the sparsest source (in that case twice $s_1$).

In figure 3.9-b), the contrast function $\phi_\Upsilon$ is shown for the same BSS problem. Clearly, the global minimum has been removed from the main diagonal and corresponds to the separation of the two desired sources. Figure 3.10 shows $\phi_\Upsilon$ for different values of the $\sigma$ parameter. For small values of $\sigma$, the twice-the-same parasitic minima are still local minima, but not global minima. As $\sigma$ increases, the function get closer to the $\ell_2$ norm contrast function, and the twice-the-same parasitic minima become local maxima. However, in the presence of strong correlation in the common support (i.e. a strong inter-regressive process), the $\ell_2$ norm fails to recover the correct sources, and small values of $\sigma$ must be used.

### 3.7.2 $3 \times 3$ **instantaneous mixture**

**Sources and mixtures.** The case with 3 sources and 3 observations is also convenient for displaying the BSE objective function. As before, the mixing matrix $\boldsymbol{A}$ is constrained to be well-conditioned. We consider at first the following arbitrary mixing matrix

$$\boldsymbol{A} = \begin{bmatrix} 2 & 1 & 1 \\ -2 & 2 & -2 \\ -1 & 2 & 1 \end{bmatrix} \tag{3.48}$$

and three one-dimension sources ($N = 128$) with the following sparsity:

$$\|\boldsymbol{s}_1\|_0 = 25\%, \quad \|\boldsymbol{s}_2\|_0 = 50\%, \quad \|\boldsymbol{s}_3\|_0 = 50\%. \tag{3.49}$$

**Results and discussion.** The BSE contrast function is shown in figure 3.11 for different values of the shaping parameter $\sigma$. For small values of $\sigma$, the function is almost equal to the size of the signals $N$ everywhere, except on some lines. Three lines with lower values appear. They cross on a single point between each other and create three strong local minima. Each of these local minima corresponds to the extraction of a sparse source. In this case, the contrast function for blind separation cannot be displayed.

The effect of having an inter-regressive process among the sources is shown in figure 3.12. Coefficients verifying equation 3.18 were added among the sources. We see that a parasitic minimum appears. The minima corresponding to the extraction of the sources are at the crossing point of three main "lines". The minimum related to the inter-regressive process is at the crossing point of several lines.

We add some Gaussian noise to the observations for the BSE problem and we observe the behavior of the objective function with different shaping parameters (figure 3.13). The variance of the noise for each observation is fixed to be a percentage of the variance of the observations. The contrast function is shown in figure 3.13 for different values of shaping parameters and different variances of noise. We see that as the noise level increases, the lines are not crossing on a single point anymore, but are passing nearby a point. The shaping parameter $\sigma$ is then useful to open the valley of each line and for determining a unique crossing point.

Our strategy for dealing with noisy observations is presented in figure 3.14. The same mixing system of equation 3.48 is used. The same amount of sparsity is used for the three sources. In particular, the number of active coefficients of $\boldsymbol{s}_1$ is given by $\|\boldsymbol{s}_1\|_0 = 0.25 \times 512 = 128$. We perturb the observation with a level of noise $R_n$ and we compute the curve $\phi_0(\sigma)$ by using the DE algorithm for different value of $\sigma$.

When there is no noise (figure 3.14a), the curve $\phi_0$ is mainly flat with a constant value of $\|\boldsymbol{s}_1\|_0 = 128$ until $\sigma \approx .1$. This is the range where the approximation by $SL0$ is accurate enough. After this value, the curve decreases to 0. When some noise is added (figure 3.14b), we see perturbations on the left hand side part of the curve. This behavior is due to instabilities of the DE algorithm for dealing with really small values of $\sigma$. For higher values of noise (figures 3.14c to e), we see that the curve $\phi_0(\sigma)$ presents a typical shape of two sigmoids. The starting point of the second sigmoid (indicated by red arrows) indicates the correct value for $\sigma$. This value correspond to a correct trade-off.

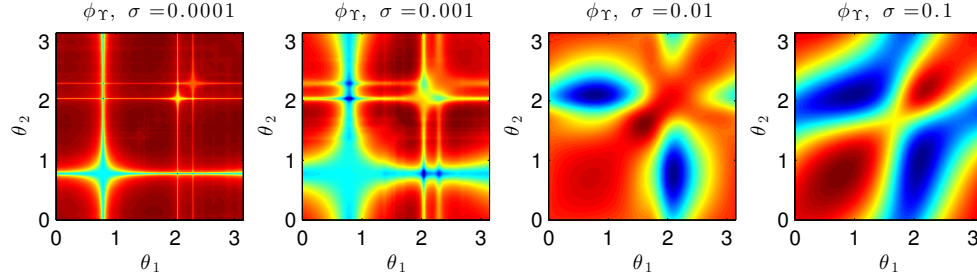**Figure 3.11:** *The SL0 objective function defined in equation 3.19 for the BSS of sparse sources with different shaping parameter $\sigma$.*



**Figure 3.12:** *Effect of the presence of an increasing inter-regressive process on the SL0 objective function defined in equation 3.19 for the BSS of sparse sources. The parasitic local minimum associated to the inter-regressive process is shown by a white arrow.*

**Figure 3.13:** *Magnification of the SL0 objective function defined in equation 3.19 for the BSS of sparse sources with different level of noise $R_n$ and shaping parameter $\sigma_0$.*

### 3.7.3 $3 \times 2$ convolutive mixture

**Sources and mixtures.** For a convolutive problem, the noiseless convolutive mixture model (equation 2.36) can be written such that

$$\boldsymbol{d}[x] = \sum_{y=0}^{Y} \boldsymbol{A}[y]\boldsymbol{s}[x - y].$$

where the mixing system is a set of matrices $\{\boldsymbol{A}[y]\}_{y=0}^{Y}$. Our first example considers $M = 3$ sources and $N = 2$ observations. The mixing system is of length $Y = 4$ and is explicitly given by

$$\boldsymbol{A}[0] = \begin{bmatrix} 0 & 4 \\ 1 & -4 \\ 0 & 5 \end{bmatrix},$$

$$\boldsymbol{A}[1] = \begin{bmatrix} 8 & -2 \\ 4 & -1 \\ 9 & -4 \end{bmatrix}, \quad \boldsymbol{A}[2] = \begin{bmatrix} -4 & -4 \\ 2 & 2 \\ 2 & 2 \end{bmatrix}, \quad \boldsymbol{A}[3] = \begin{bmatrix} 2 & -2 \\ 1 & 1 \\ -1 & -1 \end{bmatrix}, \quad \boldsymbol{A}[4] = \begin{bmatrix} -1 & 0 \\ 0 & 0 \\ 1 & 1 \end{bmatrix}.$$

**Figure 3.14:** *Strategy for identifying the level of noise in a BSS problem based on sparsity. The outsider points are due to the fact that the DE algorithm did not converged. Red arrows indicates the typical curvature due to the presence of noise in the observations.*

For this system, an inverse system exists (see section 2.4.1 for the invertibility of FIR-MIMO systems). The blind extraction is achieved by

$$\boldsymbol{u}[x] = \sum_{z=0}^{Z} \boldsymbol{w}^T[z]\boldsymbol{d}[x-z].$$

Invertibility ensures that an inverse exists, but it does not ensure that the solution is unique: two solution may exist and so two local minima associated. It is worth mentioning that the sources *and* the observations are zero-padded. Otherwise, one may loose coefficients at the edges of the observations.

**Results and discussion.** Figure 3.15 displays the result of our algorithm for extracting a sparse source, with an extracting system of length $L = 7$. This example shows a perfect recovery of the sparse source, for the case of sources verifying the assumption of theorems 9 and 10. We emphasize, here again, that the conditions described by the proposed theorems are necessary and sufficient. If the condition is not valid, then the minimum of the objective function does not correspond to a correct sparse source recovery. Instead, a parasitic vector having a lot of zeros but not corresponding to any of the original source is extracted. If the condition is valid then the minimum of the objective function correspond to the correct source extraction.

However, even if the condition is valid, the DE algorithm may fail to recovery the correct source. This is due to the random nature of evolutionary algorithms. A necessary trade-off exists for choosing the correct value of $\sigma$ for an accurate approximation of the $\ell_0$ pseudo-norm by SL0. If this value is too small, then the "lines" objective function (see figure 3.11) are really tiny valleys that are difficult to reach for an individual. The DE algorithm will need a lot of generations to converge. On the other size,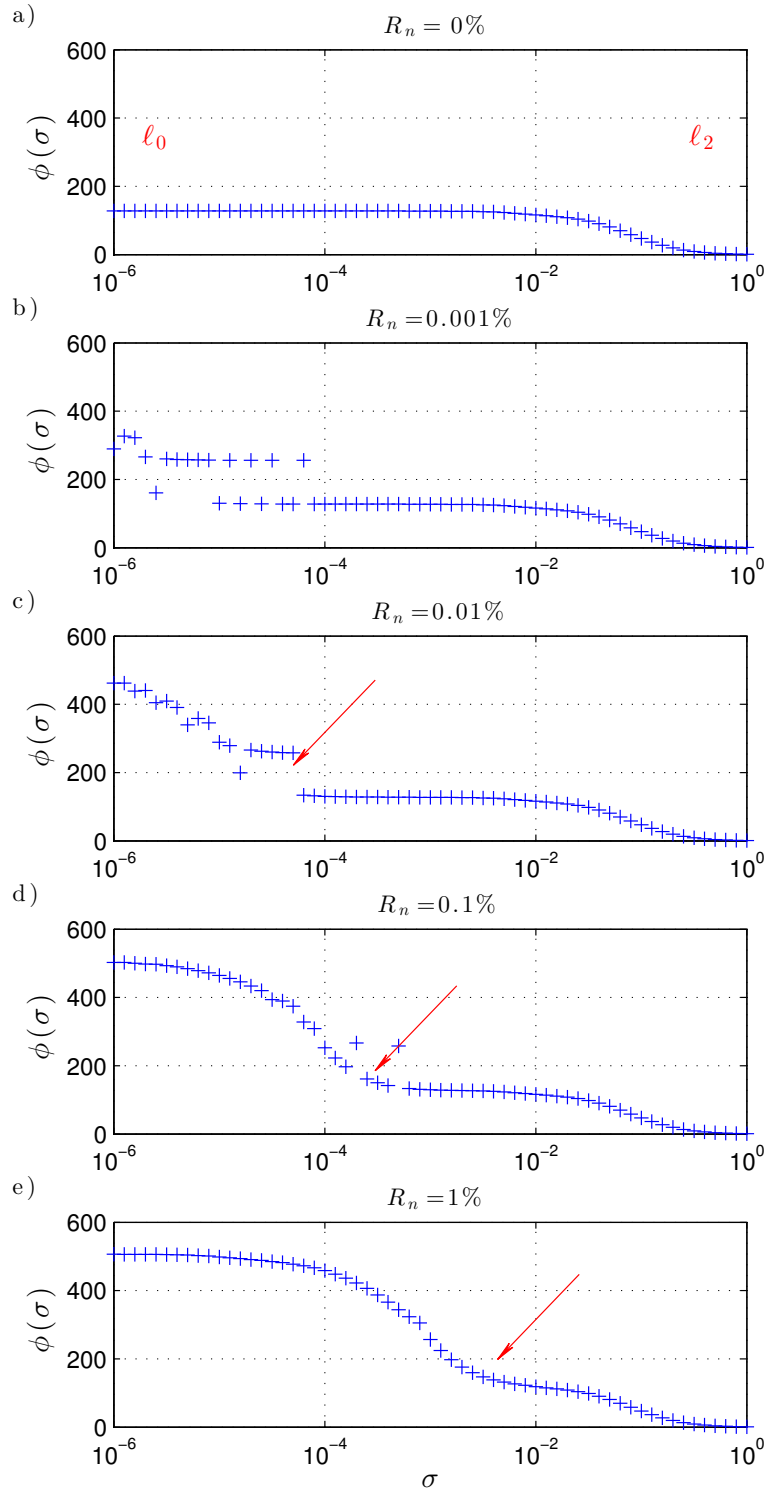 if the value of $\sigma$ is too high, then the minimum is slightly shifted from its original position and does not correspond to an accurate sparse source recovery.

One can notice that the conditions proposed by theorems 9 and 10 are actually quite weak. In practice, for random signals, the conditions are easy to satisfy and the probability of presence of an inter-regressive process is small.

When the signal has some coherence, the probability of having an inter-regressive process increases. This is quite interesting because coherency is a structured relationship inside the signal itself. It is actually equivalent to some kind of dependency. The limit of SCA and DCA is, indeed, the presence of a strong dependency among the source, named an inter-regressive process. In chapter 6, this aspect will be discussed in the particular case of primaries and multiples.

In figure 3.16, we add a really strong inter-regressive process among the source. It is a simple strict correlation between several coefficients. The condition is not valid anymore and the minimum of the objective function corresponds to a vector containing a lot of zeros, but it is not the desired sparse source.

## 3.8 Conclusion

In this chapter 3, we have analyzed the necessary and sufficient conditions for the extraction and the separation of sparse sources in determined linear mixtures. We have introduced the definition of inter-regressive processes, in a similar fashion as auto-regressive processes are already defined in the literature. In summary, an inter-regressive process denotes the presence of a correlation cluster among the sources. We have used a deterministic framework for counting the number of samples present in each inter-regressive process. Sparse sources can be extracted and separated by

**Figure 3.15:** *Example of a successful blind extraction of a sparse source from convolutive mixture with $Q = 2$ sources, $R = 3$ observations and $N = 128$ samples. The observations are mixed by filters of size $K = 4$. The extraction of the signal $\mathbf{y}^*$ is achieved by finding 3 filters of size $L = 7$. We observe perfect recovery in this case.*

**Figure 3.16:** *Same as figure 3.15, except that the recovery is not perfect.*

SCA and DCA if no strong inter-regressive process exists between the sources. If an inter-regressive process does exist, it will create a strong parasitic minimum in the SCA and DCA objective functions. When a strong correlation exists among the sources, such as an inter-regressive process, classical convex sparse measures (for instance the $\ell_1$-norm) are perturbed by the presence of the inter-regressive process. Here comes the need for finer objective functions such as SL0 to distinguish between the true minimum and the parasitic minimum.   We presented a DE algorithm for optimizing the set of extracting and separating coefficients. We have focus in particular on convolutive mixtures.

# Part III

# Adaptive subtraction of multiple events

# Chapter 4

# Multiple events in seismic acquisition

## Contents

## Résumé du chapitre [français]

Le chapitre 4 donne une vue d'ensemble des méthodes d'élimination des réflexions multiples pour les données sismiques. En effet, la plupart des méthodes d'imagerie sismique font l'hypothèse que seuls les événements primaires sont présents dans les données et donc que les événements multiples sont indésirables.

Après une classification succincte des multiples dans la section 4.2, les techniques d'éliminations sont divisées en deux groupes. La section 4.3 présente les méthodes ne faisant pas de prédiction des multiples. Il s'agit principalement de techniques de filtrage *a priori* dans divers domaines mathématiques, tels que le filtrage FK ou le filtrage Radon. En effet, après transformation dans ces domaines, l'énergie des primaires et l'énergie des multiples se concentrent dans des régions différentes. Il est alors plus facile d'y atténuer les réflexions multiples puis de réaliser la transformation inverse vers le domaine d'origine.

La section 4.4 présente les méthodes qui modélisent et calculent une prédiction des multiples. Essentiellement, les données peuvent être auto-corrélées afin de générer des échos proches cinétiquement des multiples enregistrés. Ce modèle des multiples sera ensuite soustrait aux données pour avoir une estimation des primaires. Les méthodes SRME et EPSI sont fondées sur cette idée. Elles peuvent être étendues à la prédiction des multiples internes.

Comme nous l'avons évoqué, les méthodes du second groupe ne prédisent jamais parfaitement les multiples, et l'on doit ajouter une étape d'adaptation pour que le modèle corresponde aux vrais multiples. La section 4.5 est dédiée à cette étape appelée filtrage adaptatif. Plusieurs méthodes y sont décrites succinctement, telles que les méthodes basées sur les normes $\ell_p$ ou les méthodes basées sur l'indépendance statistique.

## Resumo do capítulo [português]

O capitulo 4 apresenta um conjunto de vários métodos de eliminação de reflexões múltiplas para dados sísmicos. A maior parte dos métodos de imagens sísmicas, de fato, baseiam-se na hipótese que apenas os eventos primárias são presentes nos dados e que, portanto, os eventos múltiplos são indesejáveis.

Depois de uma classificação sucinta das múltiplas na seção 4.2, as técnicas de eliminação são divididas em dois grupos. A seção 4.3 apresenta os métodos que não fazem previsões de múltiplas. Tratam-se principalmente de técnicas de filtragem a priori dentro de diversos domínios matemáticos, como a filtragem FK ou a filtragem Radon. Após transformações dentro desses domínios, a energia das primárias e a energia das múltiplas concentram-se em regiões diferentes. Portanto, torna-se mais fácil atenuar as reflexões múltiplas depois de realizar a transformação inversa em direção ao domínio de origem.

A seção 4.4 apresenta os métodos que modelam e calculam uma previsão de múltiplas. Os dados podem ser basicamente auto correlacionados a fím de gerar ecos próximos cineticamentes de múltiplas registradas. Esse modelo de múltiplas será, em seguida, subtraído dos dados para poder-se estimar as primárias. Os métodos SRME e EPSI são baseados nessa idéia. Eles podem ser estendidos às múltiplas internas.

Como foi mencionado anteriormente, os métodos do segundo grupo nunca prevêem perfeitamente as múltiplas, e nós devemos adicionar uma etapa de adaptação para que o modelo corresponda às verdadeiras múltiplas. A seção 4.5 é dedicada a essa etapa chamada "filtragem adaptativa". Vários métodos são aqui descritos de maneira sucinta, como os métodos baseados nas normas $\ell_p$ e os métodos baseados na independência estatística.

## 4.1   Introduction

As introduced in chapter 1, the data recorded in seismic acquisitions can be considered as the superposition of primary and multiple events such as

$$\boldsymbol{d}[x] = \boldsymbol{p_0}[x] + \boldsymbol{m_0}[x], \tag{4.1}$$

where $\boldsymbol{d}$, $\boldsymbol{p}_0$ and $\boldsymbol{m}_0$ are the data, the true primaries and the true multiples, and the index $x$ indicates the location in the data cube. However, most of the conventional imaging and processing techniques developed by the exploration community consider that the data are multiple free. This is the case, for instance, for reverse time migration (RTM) or velocity analysis. Therefore, it is crucial to adequately remove the multiple energy which is considered, here, as a noise. An effective way to have an overview of the importance of multiples in seismic acquisition is to refer to the four special section editions of the international journal *The Leading Edge*, in *1999*, *2005*, *2011* and *2015*, dedicated to multiples.

In this chapter 4, section 4.2 is dedicated to a rapid classification of multiples. In sections 4.3 and 4.4 we develop the different methods existing to eliminate multiples. A first class of methods (section 4.3) considers some mathematical transformations of the data able to make primaries and multiples almost disjoint orthogonal. This kind of strategy allows for filtering the multiples in the transformed domain. A second class of methods (section 4.4) based on the feedback model, allows estimating the kinematic of the multiples. These methods generally need an adaptive subtraction step to adequately remove the noise (section 4.5). Adaptive subtraction of multiples will be the core subject of chapter 5. To finish, section 4.6 briefly presents some recent considerations about using multiples as a signal and not as a noise. As for the previous chapters, we will focus on the objective functions and the optimization scheme of the presented methods.

## 4.2   Classification of multiples

Several classifications exist for multiples. These classifications are motivated by the different methods able to tackle the problem of multiple elimination. Each method is indeed able to tackle only part of the multiples. Hence, there is no single definition of multiples. We start with the definition proposed by Weglein et al. [1997]: "a multiple is a seismic event that has experienced two or more upward reflections". More details can be found in Yilmaz [2001] or Verschuur [2013b]. Moreover, in table 4.1 inspired by Wong [2012], we present the various types of multiples and the most suitable methods to tackle their removal.

**Source ghost and receiver ghost**

Source ghosts and receiver ghosts are due to the location of seismic sources and receivers, respectively. Those devices are indeed not exactly located at the free surface, but deeper for practical reasons. On the source side, the energy will propagate in all direction and reflect rapidly on the free surface. On the receiver side, the receiver records the upward energy and the downward energy reflected by the free surface. The ray paths of one source ghost and one receiver ghost are presented in figure 4.1-a) and b).

Depending on the definition of multiples, ghosts may be not classified within the multiple category (see for instance the definition by Weglein et al. [1997]). Also, dedicated methods are used to eliminate ghosts and their signature can be easily

| Multiple characteristic | Method | Pros | Cons |
| --- | --- | --- | --- |
| short-period | Predictive deconvolution | Fast | Periodicity not always perfect (e.g. if dip) |
| different move-out | FK filtering | Fast | Not optimal for near offset |
| | Radon filtering | Fast | Not optimal for near offset |
| surface related (free surface) | SRME | Accurate prediction | Require dense acquisition |
| | EPSI | Accurate prediction | Computationally expensive require effective identification of first primaries |
| surface related (identified reflector) | SRME | Accurate prediction | Downward continuation Need dense grid and interpolation |
| internal multiples | ISS | | |
| | Jakubowicz's approach | | |

**Table 4.1:** *Principal methods for multiple attenuation (Inspired from Wong [2012]).*

implemented in forward modeling (see e.g. Verschuur [2013a]). As multiples, ghosts create small repetitions in the data, but without the periodicity specific to multiples.

### Surface related multiples

Surface-related multiples refer to the presence of one strong reflector [Yilmaz, 2001; Verschuur, 2013b]. All multiples that are removed if this reflector disappears are said to be related to this surface. Generally, the strong reflector is the free-surface, and one speak about *free-surface-related multiples*. The *order* indicates the number of time the energy has been reflecting on the surface. One speaks about first order multiples, second order multiples, etc. The ray paths of a first order and a second order surface related multiple are presented in figure 4.1-d) and e).

### Internal multiples

Internal multiples are due to one or two strong reflectors below the free surface [Yilmaz, 2001]. They are often referred to as short period multiples because they create close duplicated events. *Intrabed multiples* refer to reflections created inside a single layer. *Interbed multiples* refer to reflections created between two different layers.

*Pegleg multiples* may refer to different kind of multiples. They are due to two different reflectors acting at two different depths. The ray paths of internal multiples are presented in figure 4.1-f) to h).

## 4.3   Removal methods based on filtering (no prediction)

When they reach a receiver, multiples have generally travelled shallower than the primaries arriving at the same time, as the wave propagation velocities generally increase with depth. The apparent velocities of multiples are then lower than the ones for primaries and they present different dips (more horizontal) in common midpoint gathers (CMP) [Yilmaz, 2001]. In other words, primaries have less move-out than multiples. Move-out correction can enhance this effect. Stacking CMPs after normal move-out (NMO) correction can reduce the presence of multiples in the stacked trace, but more advanced methods exist to filter the multiples before stacking and increasing the signal to noise ratio for imaging process.

**Assumption 5.** *Multiples have a smaller move-out compared to primaries.*

Assumption 5 is the foundation for many methods, mainly based on some kind of filtering. These techniques are often non-adaptive and use prior informations on the mapping of multiples and primaries in a transformed domain.

### 4.3.1   Stacking and fancy stacking

Multi-offset surveys allow a lot of different processing for enhancing the signal to noise ratio. Perhaps the most simple one is the sum of several traces after NMO correction, called stacking. If the noise from trace to trace is not correlated, stacking has shown to be effective for removing noise. However, if the noise is correlated in space, i.e. from trace to trace, the assumption fails. After NMO correction we can write the

**Figure 4.1:** *(a) source ghost ; (b) receiver ghost ; (c) Primary ray path ; (d) first order multiple ray path ; (e) second order multiple ray path ; (f) intrabed multiple ray path ; (g) peg-leg multiple ray path ; (h) deeper surface related multiple. The red star indicates the seismic source position.*

stacking operation as:

$$p[t] = \sum_x \boldsymbol{d}_{NMO}[x,t] = \sum_x \boldsymbol{p}_{NMO}[x,t] + \underbrace{\sum_x \boldsymbol{n}_{NMO}[x,t]}_{\approx 0}. \qquad (4.2)$$

Under assumption 5, part of the multiples should be removed.

### 4.3.2 Pre-stack FK filtering

The 2D Fourier domain is commonly called FK domain in seismic (where F denotes the temporal frequency $\omega$ and K denotes the spatial frequency or wavenumber $k$). The

continuous 2D Fourier transform is defined as

$$\mathcal{F}\boldsymbol{d}(k,\omega) = \iint \boldsymbol{d}(x,t)\mathrm{e}^{i(\omega t - kx)}dtdx, \tag{4.3}$$

where the minus sign is used for conveniency. The discrete Fourier transform (see figure 4.2b) is defined as

$$\mathcal{F}\boldsymbol{d}[k,\omega] = \frac{1}{N_x N_t} \sum \sum \boldsymbol{d}[x,t]\mathrm{e}^{i\left(\frac{\omega t}{N_t} - \frac{kx}{N_x}\cdot\right)}. \tag{4.4}$$

After NMO correction, the primaries should be flat (i.e. not depend on the receiver position) and the multiples should be over-corrected and go upward. Hence primaries and multiples map into different areas of the FK domain. Muting can be done in the FK domain such as

$$\mathcal{F}\boldsymbol{p} = \mathcal{F}\boldsymbol{w} \odot \mathcal{F}\boldsymbol{d_{NMO}}. \tag{4.5}$$

This distinction is not valid at near offset, as both primaries and multiples are usually horizontal in this region, whatever the velocity used for NMO correction.

### 4.3.3   Pre-stack Radon filtering

The Radon transform (also called slant-stack transform) is a popular transform for seismic processing. It aims at identifying curves inside an image. In the continuous case, it is defined as [Durrani and Bisset, 1984; Verschuur, 2013b] (see figure 4.2c)

$$\mathcal{R}\boldsymbol{d}(q,\tau) = \int \boldsymbol{d}(x, t = \psi(q,\tau))dx \quad \text{with} \begin{cases} \psi = \tau + qx & \text{linear}, \\ \psi = \tau + qx^2 & \text{parabolic}, \\ \psi = \sqrt{\tau^2 + \frac{x^2}{q^2}} & \text{hyperbolic}, \end{cases} \tag{4.6}$$

where the curve $\psi$ gives the name of the transform. The parameter $q$ is linked to velocity (to the shape of the curve) while $\tau$ gives the travel-time at a specific spatial position. Under assumption 5, primaries and multiples map into different areas and a muting approach is well-suited. Parabolic Radon filtering is often applied after NMO correction while hyperbolic Radon filtering can be applied before or after NMO correction. Because of limited offsets, high-resolution versions are used but the inverse of those transforms can lead to artifacts. Inversion and regularization schemes are therefore needed for a more precise inverse transform [Sacchi and Ulrych, 1995; Maeland, 2003; Abbad et al., 2011]. The main drawback of these techniques is the difficulty to separate of primaries and multiples at near offsets, as they are both horizontal in this area.

### 4.3.4   Post-stack filtering

The previous methods for filtering multiple events are applied before stacking, and generally after move-out correction. However, it is also possible to perform multiple attenuation on stacked sections. The multiples have a tendency to create events flatter than the primary events [Verschuur, 2013a].

**Assumption 6.** *In a stacked section, multiples may have a different dip compared to primaries.*

Assumption 6 can be used in the image domain. For instance, strategies based on FK or curvelet filtering may be used to get rid of this parasitic energy.

**Figure 4.2:** *a) a common shot gather, b) its 2D Fourier transform (only the amplitude), c) its linear Radon transform. Red color indicates positive values and blue color indicates negative values.*

## 4.4 Removal methods with prediction

In this section, we discuss several methods that have been developed for removing multiples, either surface-related multiples or internal multiples. Those methods are referred to as *prediction methods* because they create a model (a prediction) of the multiple events [Verschuur, 2013a]. As we will see, due to acquisition limitations, those methods cannot perfectly predict the multiple events and an adaptive subtraction step is needed to adapt the prediction to the data.

### 4.4.1 Predictive deconvolution

Multiples are echos of the original primaries and they appear with a certain regularity, or periodicity. By exploiting this specific property, one can remove part of the multiples.

**Assumption 7.** *The multiples are periodic.*

The mixing model based on assumption 7 can be written as

$$\boldsymbol{d}[t] = (\boldsymbol{a} * \boldsymbol{p})[t], \tag{4.7}$$

where $\boldsymbol{a}$ can be modeled as a sum of Dirac delta functions, with lags between the Dirac functions defining the periodicity and with amplitude coefficients defining the attenuation (see figure 4.3). Deconvolution of the data can be designed as

$$\boldsymbol{p}[t] = (\boldsymbol{w} * \boldsymbol{d})[t], \tag{4.8}$$

where constraints are required on the filter $\boldsymbol{w}$, based on the model. The deconvolution parameters can be obtained by least-square inversion [Peacock and Treitel, 1969] or $\ell_1$-norm minimization [Li et al., 2016]. Multichannel predictive deconvolution can improved the results by exploiting the spatial coherence of adjacent traces [Porsani and Ursin, 2007].

### 4.4.2 Surface related multiple elimination (SRME)

Two main methods extensively used today for predicting multiples have emerged in the 90's, both around the same time. At Delft University, Berkhout and Verschuur

**Figure 4.3:** *Mixing model of multiples generation for classical predictive deconvolution technique.*

have developed the popular surface related multiple elimination (SRME) method ([Verschuur et al., 1992; Berkhout and Verschuur, 1997]). Weglein developed the inverse-scattering series (ISS) method ([Weglein et al., 1997]). Levin [2008] discusses how they differ in some practical aspect but shows that they are theoretically really close.

The introduction of the surface related multiple elimination (SRME) method has been a breakthrough in the processing of multiples generated by a specific surface. This method is fully data-driven, meaning that no velocity model is needed for creating an accurate model of the multiples. The core of SRME lives in the observation that the 3D data cube $\boldsymbol{d}[x]$ allows us to generate multiples by auto-convolution of the seismic data, where the index $x = [h, t, x_s]$ contains the offset $h$, the time $t$ and the source position $x_s$. As we will see, those auto-convolutions lead to a filter indetermination that will need to be addressed by adaptive filtering. In practice, this filter indetermination is not the only reason why one needs the adaptive subtraction, as we shall discuss late. We present here the SRME fundamental equation.
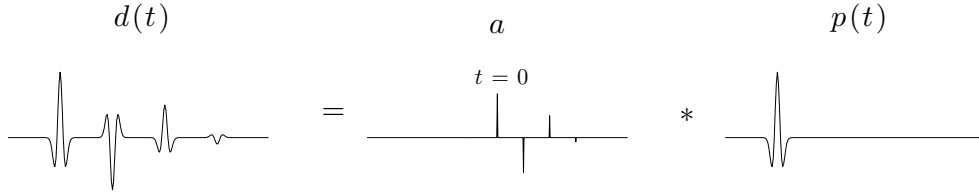
As we defined before, a surface related multiple is due to the presence of a strong reflector that literally reflects downward part of the energy. If we consider the point at which the seismic ray has been reflected, we see that multiples are the kinematic concatenation of several hypothetical primaries (figure 4.4). By hypothetical primaries, we mean here a primary that could have been generated by placing the seismic source at the reflection point. The strategy of seismic surveys is to move the source and the receivers along the desired profile, such that at the end, those primaries are not hypothetical but actually recorded at another moment. If we are able to identify where this "multiple recorded as a primary' is located in the data, then we have our prediction of the multiples.

For a zero-offset acquisition, the prediction equation is even easier to obtain: the multiples have been recorded as primaries in the same trace. For a linear system, the recorded primaries are $\boldsymbol{p}[t] = (\boldsymbol{a} * \boldsymbol{s_p})[t]$ where $\boldsymbol{a}$ is the seismic source and $\boldsymbol{s_p}$ the primary Green's function[1]. If the strong reflector is at the position of the source and receiver, any energy recorded at some instant will be re-injected downward in the same system (having the same Green's function). We can write the trace such that

$$\boldsymbol{d}[t] = \underbrace{\boldsymbol{p}[t]}_{primaries} \quad \underbrace{-(\boldsymbol{p} * \boldsymbol{s_p})[t]}_{\substack{1st\ order \\ multiple}} \quad \underbrace{+(\boldsymbol{p} * \boldsymbol{s_p} * \boldsymbol{s_p})[t]}_{\substack{2nd\ order \\ multiple}} \quad - \dots, \tag{4.9}$$

where a factor $-1$ is included because of the reflection from the perfect water free

---

[1]Our notation changes from conventional notations in SRME works. This choice is motivated by the precedent part on blind source separation with which we want to make the link. In particular, the seismic source is generally unknown and can be seen as the mixing system, while the Green's function is the signal we want to recover. The reflectivity is generally considered as sparse.

surface. We can verify that the primaries can be obtained by [Weglein et al., 1997]

$$p[t] = \underbrace{d[t]}_{data} - \Bigg( \underbrace{(w * d * d)[t]}_{\substack{1st\ order \\ multiple}} \quad \underbrace{-(w * w * d * d * d)[t]}_{\substack{2nd\ order \\ multiple}} \quad + \dots \Bigg), \qquad (4.10)$$

where we defined $w$ such that $(w * a)[t] = -\delta[t]$. We see that each order of multiple can be predicted by auto-convolution of the data. A filter indetermination remains with the surface operator.

In 2D acquisition, the structure of the 3D data cube must be used to carefully generate multiples. In particular, the prediction is generally done in the frequency domain which also makes the SRME equation more readable. Similarly to equation 4.10, we can write in the frequency domain

$$P_\omega = D_\omega - W(\omega)D_\omega^2 + W^2(\omega)D_\omega^3 - \dots, \qquad (4.11)$$

where the lines of the matrices contain monochromatic receiver gathers and the columns of the matrices contain monochromatic shot gathers. Clearly, an adaptive step is needed to adapt the prediction to the data, if the source is unknown. Also, because of acquisition geometry, the prediction is not perfect and the perturbations must be tackled by the adaptive subtraction step. SRME can be generalized to 3D modern acquisitions [Pica et al., 2005; Reshef et al., 2006; Dragoset et al., 2010].



**Figure 4.4:** *Principle of surface related multiple prediction. a) For free surface multiple, the ray path can be decomposed into the superposition of two primary ray paths. b) For an internal multiple, the ray path can be decomposed with three ray path, one of them (magenta) having a negative contribution.*

### 4.4.3 Estimation of primaries by sparse inversion (EPSI)

The estimation of primaries by sparse inversion (EPSI) method has been proposed by van Groenestijn [2010] and applied in diverse contexts [van Groenestijn and Verschuur, 2009, 2010]. EPSI is a continuation of the Delft University method. EPSI is an iterative full waveform inversion based on the exact same prediction model as SRME. Hence, EPSI tries to avoid any adaptive subtraction step by using what they called a sparse inversion. The EPSI formula can be express in the frequency domain with the matrix notation such as

$$D_\omega = S_{p,\omega} \left[ \mathcal{F}a_\omega - D_\omega \right]. \qquad (4.12)$$

where $D_\omega$ and $S_{p,\omega}$ are monochromatic matrices and $\mathcal{F}a_\omega$ is the Fourier transform of the seismic source at frequency $\omega$ (a scalar in that case). The minus sign is used because a perfect reflection is considered at the surface. The optimization problem for EPSI is set as

$$\text{find } S_{p,\omega} \text{ and } a(\omega) \qquad \text{such that} \qquad D_\omega \approx \hat{S}_{p,\omega} \left[ \hat{a}(\omega) - D_\omega \right]. \qquad (4.13)$$

Because this problem is highly underdetermined, some prior must be added. EPSI makes the assumption that the reflectivity is sparse. By looking at equation 4.12, we see that EPSI is a bilinear inverse problem for which the primary Green's function $S_p$ and the seismic source $a$ must be iteratively estimated until convergence. In the original work by van Groenestijn and Verschuur [2009], the operator imposing the reflectivity sparse is a simple threshold. Verschuur [2013a] includes ghosts in the EPSI framework. Savels et al. [2011] shows application on real data. It is worth mentioning that EPSI does not need a near-offset interpolation pre-processing, compared to SRME.

Lin and Herrmann [2013] improved the robustness and the stability of the original algorithm by formulating EPSI in a biconvex optimization scheme. As in the original algorithm, the Green's function and the source are updated sequentially. The seismic source is the result of a $\ell_2$ optimization scheme while the Green's function is the result of a LASSO optimization scheme written as

$$\min \quad \left\| \boldsymbol{D} - \hat{\boldsymbol{D}} \right\|_2 + \epsilon \left\| \hat{\boldsymbol{S}}_p \right\|_1 . \tag{4.14}$$

Feng et al. [2013] include a similar strategy in the curvelet domain. Recently, EPSI has received a lot of attention for dimension reduction [Jumah and Herrmann, 2014; Tu and Herrmann, 2015].

### 4.4.4   Wave field extrapolation

The beauty and usefulness of SRME and related techniques is the fact that they are fully data driven. In other words, no velocity model is needed to make the prediction of the multiples. However, if we know some information about the subsurface such as the position of a strong reflector, the multiples can be predicted by wave field extrapolation. Historically, this method has proceeded SRME. The data recorded at the surface are extrapolated to the position of identified strong reflector. This step constructs the data at the position of virtual sources. Then, the extrapolated wavefield can be used to construct a prediction of the multiples [Wiggins, 1999].

### 4.4.5   Internal multiples

Jakubowicz [1998] has described a method for predicting interbed multiples. His approach is a generalization of SRME prediction (proposed by Verschuur et al. [1992] and Verschuur and Berkhout [1997]). We saw in section 4.4.2 that water-surface related multiples can be predicted by decomposing the ray path of each multiple into several primaries. Each primary is merged in a positive manner to contribute to the full ray path. In a similar way, an internal multiple can also be decomposed into several primaries, but some of them must be merged in a negative manner to remove some paths (figure 4.4). In the first formulation of the method, the surface acting as a strong reflector must be identified. Hung and Wang [2012] use ideas from ISS to improve the Jakubowicz's approach such that the strong reflector does not need to be identified.

## 4.5   Adaptive multiple subtraction [quick overview]

Multiple attenuation is crucial for improving the quality of seismic images, especially in marine acquisitions. As shown in section 4.4, several techniques exist to provide a prediction of these multiples. Unfortunately, none of these prediction-based methods can provide a perfect prediction of the multiples because of phase, wavelet or space-shift errors [Abma et al., 2005]. Therefore, a second step, usually referred to as

adaptive multiple subtraction, is required to accommodate the prediction to the actual multiples before the subtraction. The most common solutions are based on matching-filter approaches [Verschuur and Berkhout, 1997; Rickett et al., 2001; Guitton and Verschuur, 2004] and on prediction-error filters either in the frequency [Spitz, 1999] or in the time domain [Guitton, 2005]. Table 4.2 aims at giving a list of these methods with their main properties.

### 4.5.1 Generalities about adaptive subtraction

Most of adaptive multiple subtraction schemes rely on a linear convolutive model to reshape the predicted multiples. However they may differ in the following aspects:

- the objective function to be optimized,
- the domain to perform the optimization,
- the strategy to overcome the non-stationarity of the filter,
- and the strategy to exploit the space-time coherence of the seismic signal.

Often, non-stationarity and space-time coherence are handled with a common strategy. However it is important to keep in mind that non-stationarity is a difficulty to overcome whereas the space-time coherence is an asset to capitalize on.

Because of its computational efficiency, the $\ell_2$-norm is the most commonly employed objective function in adaptive multiple subtraction. The resulting filter, which is known as least-squares or Wiener filter, works under the assumption that primaries and multiples are orthogonal in the considered domain [Verschuur, 2013b]. However in practice, and in particular in the time-offset domain, this assumption fails and it may lead to an over-attenuation of the estimated primaries. For this reason, some works consider some $\ell_1$-norm based filters that seem to overcome the problem by promoting a sparser solution of the estimated primaries [Guitton and Verschuur, 2004]. Interestingly, $\ell_q$-norms have been considered as a regularization term [Costagliola et al., 2011]. Moreover, a Bayesian framework has also been investigated by Saab et al. [2007a].

More recently, other works proposed to use independent component analysis (ICA) [Comon and Jutten, 2010] to separate primaries and multiples. This approach has led to the use of new objective functions associated with methods such as geometric-based ICA [Lu, 2006], FastICA [Kaplan and Innanen, 2008], kurtosis-based methods [Donno, 2011], Infomax [Liu and Dragoset, 2013], negentropy maximization [Li and Lu, 2013]. The first works [Lu, 2006; Kaplan and Innanen, 2008; Donno, 2011] on ICA-based adaptive multiple subtraction operate in a two-step fashion. They comprise an estimation of the shape of the filter using a classic $\ell_2$-norm matching filter or a histogram method to correct for time delay, followed by a more precise adjustment of its amplitude using ICA. More recent works [Liu and Dragoset, 2013; Li and Lu, 2013] have proposed to directly rely on a convolutive modeling with objective functions based on statistical independence.

The domain in which the matching filter is performed is decisive in adaptive multiple subtraction and a lot of effort has been done to search for domains where primaries and multiples do not overlap. Usually, the adaptive multiple subtraction procedure is carried out in the time-offset domain where the orthogonal assumption fails. Other domains have been proposed such as dip-domain [Donno, 2011], wavelet-domain [Ahmed, 2007; Ventosa et al., 2012], curvelet-domain [Herrmann et al., 2007, 2008; Donno et al., 2010], Radon domain [Li and Lu, 2014], frequency domain [Spitz, 1999], adjoint domains (the first derivative along with the Hilbert transform and its first derivative) [Wang, 2003]. In this paper we only consider the space-time domain but our conclusions hold for other domains.

After defining a proper objective function and a suitable domain to perform the adaptive multiple subtraction procedure, the last issue to overcome is the non-stationarity of the primaries and the multiples [Guitton, 2005; Fomel, 2007b, 2009]. This means that the statistical features of the data are not steady with respect to the time or the offset and so neither the filter we aim to recover [Velis, 2003]. However, the spatial and temporal coherence of the seismic signal prevents from drastic changes, and smooth variations can be assumed. Hence, most of the time, the signal is considered as stationary in a small data window in which a unique filter can be obtained. This operation is then repeated on several over-lapping windows to complete the full data length. Finally we would like to mention that one, two or three dimensional data windows – and so filters – can be considered [Wang, 2003; Donno, 2011], since the seismic signal is locally coherent in the full data cube.

We consider the following linear model for adaptive multiple subtraction

$$\begin{cases} \boldsymbol{d} = \boldsymbol{p} + \boldsymbol{m}, \\ \hat{\boldsymbol{m}} = \boldsymbol{w} * \boldsymbol{m}, \end{cases} \tag{4.15}$$

where a hat $\hat{\cdot}$ indicates a final estimation and $*$ denotes the convolution product that can be either 1D, 2D or 3D, according to the dimension of $\boldsymbol{w}$ and $\boldsymbol{m}$. The estimate of the primaries is then given by

$$\hat{\boldsymbol{p}} = \boldsymbol{d} - \hat{\boldsymbol{m}} = \boldsymbol{d} - \boldsymbol{w} * \boldsymbol{m}. \tag{4.16}$$

The model described by equations 4.15 and 4.16 will be used in chapters 5 and 6.

## 4.5.2 Toy synthetics for adaptive multiple subtraction

### Crossing events

Figure 4.5 shows a synthetic gather in which the primary and the multiple events are crossing. The two events have different dips. This situation can occur in gathers at far offsets. Multiples present a different move-out because their apparent velocity is lower (they have travelled shallower). This kind of situation can also occur in post stack.

In this example, we use the exact multiple (figure 4.5b) as the prediction. This strategy helps us to underline the limit of the least-squares filtering. A filter of 16 ms is computed in the least-squares sense and applied on the prediction for each trace. As result, we see that the primary estimated has been damaged where the multiple is crossing. This is what is called *over-attenuation*. Because the $\ell_2$-norm approach tends to minimize the energy of the primary events, a solution full of zeros is a perfect solution. The matching filter does not adapt the prediction to the multiple event, but to the primary event.

### Parallel events

Figure 4.6 shows a synthetic gather in which the primary and the multiple are parallel. This situation can occur at near offset. Parallel events are also visible at far offsets when considering a small window.

In this example, we use the actual multiple (figure 4.6b) as the prediction. A filter of 68 ms is computed in the least-squares sense and applied on the prediction for each trace. As result, we see that the primary estimated has been damaged where the primary is close enough to the multiple, within the filter length. Here again, over-attenuation occurs. The filter is not adapting the prediction to the multiple event in the data, but to the primary event.

| Authors [year] | Model | Domain | Optimization | Regularization |
|---|---|---|---|---|
| Verschuur et al. [1992] | conv. | time | $\ell_2$-norm | 1D |
| Guitton and Verschuur [2004] | conv. | time | $\ell_1$-norm | 1D/2D |
| | | id.$\ell_{1/2}$-norm | $\phi = \sum g_\epsilon(p)$ | id. |
| Lu [2006] | amp. | time | ICA | 2 steps / 2D |
| Saab et al. [2007a] | Full W. | curvelet | norms | 2D |
| Kaplan and Innanen [2008] | amp. | time | ICA | 2D |
| Herrmann et al. [2008] | Full W. | curvelet | norms | 2D |
| Fomel [2009] | conv. | time | $\ell_2$-norm | 2D |
| Neelamani et al. [2010] | $\mathbb{C}$-amp. | curvelet $\mathbb{C}$ | $\ell_1$-norm | local |
| Donno [2011] | amp. | dip | kurtosis | 2D |
| Costagliola et al. [2011] | conv. | time | $\ell_q$-norm | 1D |
| Li and Lu [2013] | conv. | time | Negentropy | 1D |
| Liu and Dragoset [2013] | conv. | time | Infomax | 1D |
| Ventosa et al. [2012] | amp. | wavelet | $\ell_2$-norm | local |
| Pham et al. [2014] | conv. | wavelet | norms $\ell_1$, $\ell_2$ | 2D |
| Liu and Kostov [2015] | conv. | time | AIC | 2D |
| Wu and Hung [2015] | amp. | curvelet | norms $\ell_2$, $\ell_1$ | local |

**Table 4.2:** *Tentative of exhaustive list of literature about matching filter based adaptive subtraction of multiples.*

**Figure 4.5:** *Synthetic example of crossing events. a) Primary event. b) Multiple event. c) Data. d) Result of the least-squares estimation with a filter of 16 ms for each trace.*



**Figure 4.6:** *Synthetic example of parallel events. a) Primary event. b) Multiple event. c) Data. d) Result of the least-squares estimation with a filter of 68 ms for each trace.*

## 4.6   Imaging with multiples

In the previous sections, we introduced methods for removing multiples from the data. In the last years, a lot of efforts have been done for using multiples as a signal. The main reason is that multiples could better illuminate some parts of the sub-surface and enforce resolution [Berkhout and Verschuur, 2006]. For instance, the missing short offsets can be recovered [Curry and Shan, 2010], or the cross line resolution can be improved. It is out of our scope to describe all existing methods, but we want to give to the reader some perspectives on using multiples in seismic imaging. A good starting point is the *2015* special edition of *The Leading Edge* entitled "Multiples from attenuation to imaging".

### 4.6.1   Full waveform inversion

Commonly in seismic imaging there exists a distinction between large scale imaging (tomography) and fine scale imaging (migration). Differently from that, full waveform inversion (FWI) consists, as its name suggests, of the inversion of the whole seismic data set [Tarantola, 1987; Brossier et al., 2009; Virieux and Operto, 2009]. Theoretically, the multiples are embraced in the forward modeling operator $L$ and the objective function minimizes the distance between the data $\boldsymbol{d}$ and the modeled data $L\boldsymbol{\theta}$, where $\theta$ are the parameters of the model. However it is well known that this objective

function contains a lot of local minima due among other to cycle skipping or the presence of multiples. The initial model is crucial for a good result [Bunks et al., 1995]. Compared to classical tomography approaches, FWI does not need any pointing of primaries to get a macro-model but this method is really sensitive to the presence of low frequencies in the data [Pratt et al., 1996]. Today, a lot of efforts are done in the forward modeling (e.g. for including visco-elastic equations by Brossier [2011]) and in the use of alternative objective functions [Métivier et al., 2016].

### 4.6.2 Migration with multiples

Classical migration is a linearization of forward modeling and considers that only the primaries are present in the data. When multiples are added, the problem becomes non-linear and multiples have to be handled with a lot of good care [Wong et al., 2015]. Currently, efforts are done for including multiples in the migration velocity analysis (MVA) framework for retrieving large scales [Cocher et al., 2015] that can lead to a better initialization of the macro-model for FWI [Diaz and Sava, 2013].

### 4.6.3 Marchenko approach

From the reciprocity theorem, seismic interferometry estimates the Green's functions between receivers by cross-correlation of their respective traces. In particular, this process can be performed with passive noise, without any active source. If one cross-correlated one trace at location $x$ with other traces, a virtual source and its associated virtual wave-field is created coming from this location $x$ [Snieder et al., 2006]

In presence of random noise, interferometry can create a virtual source at any receiver by cross-correlation with other receivers. As a step forward, Marchenko imaging can create any virtual source in the subsurface [Wapenaar et al., 2014; Meles et al., 2015]. Marchenko imaging uses internal multiples, but surface related multiples must be removed from the data [van der Neut et al., 2015].

## 4.7 Conclusion

In this chapter 4, we have discussed the definition of multiples as well as the two main classes of multiple attenuation techniques. Methods from the first class are based on non-adaptive filtering techniques and are refereed to as methods without prediction. Methods from the second class are based on a prediction of the multiple events. Variations on the feed-back model are used to produce the prediction that can be fully data-driven or based on a physical model. For methods belonging to the second class, a second step, namely the adaptive subtraction step, is needed to better accommodate the noise to the data. This step is of crucial importance, as illustrated in figures 4.5 and 4.6. If the subtraction is not performed in a correct manner, even a perfect prediction could damage primaries and part of the signal could be lost. Adaptive filtering is the subject of the next two chapters 5 and 6, with the incorporation of ideas from chapters 2 and 3.

# Chapter 5

# Inside the wildlife of multiple subtraction methods

## Contents

Most of the results of this chapter have been published in Batany et al. [2016b] and presented at the 78th EAGE Conference and Exhibition in Vienna [Batany et al., 2016c]. The content of these articles is renewed and a few ideas are added in order to make links with the other chapters.

## Résumé du chapitre [français]

Le chapitre 5 présente les c ontributions concernant la soustraction adaptative des multiples. Plusieurs méthodes sont développées dans le but de pouvoir les unifier. Grâce à cette étape, leur comparaison sera par la suite plus aisée.

La section 5.2 présente les méthodes basées sur la minimisation de normes $\ell_q$. Ces méthodes sont certainement les plus populaires. La méthode de minimisation des moindres carrés (norme $\ell_2$) a l'avantage de posséder une solution analytique. La minimisation de la norme $\ell_1$ propose quant à elle une estimation des primaires plus parcimonieuse et plus proche de la distribution statistique des réflexions primaires.

La section 5.3 concerne les méthodes venant de l'analyse en composante indépendante (méthodes ICA) et de la séparation aveugle de sources. Dans le cadre de la soustraction adaptative, l'hypothèse principale est que les réflexions primaires et multiples peuvent être modélisées comme des variables aléatoires statistiquement indépendantes. Les méthodes basées sur la néguentropie, Infomax et le kurtosis sont décrites.

Les méthodes de reconnaissance de forme ne sont pas incluses dans l'unification proposée. Il semble important de les évoquer pour avoir une meilleure perspective de leurs différences avec les méthodes de filtrage évoquées précédemment. La section 5.4 est dédiée à ces méthodes.

La contribution principale est développée dans la section 5.5. L'ensemble des méthodes basées sur les normes et sur l'indépendance statistique sont analysées, en considérant les fonctions objectives minimisées dans chaque cas. Cette façon de faire permet de séparer d'un côté la fonction optimisée et de l'autre la stratégie de fenêtrage (évoquée ensuite). Il est montré que l'ensemble des méthodes cherchent à minimiser une corrélation non-linéaire entre les primaires estimés et les multiples prédits. La non-linéarité est portée par un opérateur agissant comme un compresseur sur les primaires. On montre ainsi que certaines méthodes présentent de fortes similarités.

La dimension du filtre et sa zone d'application sont des paramètres essentiels, regroupés sous le terme de stratégies de fenêtrages. Elles sont présentées dans la section 5.6. Finalement, la section 5.7 présente plusieurs expériences faites sur des données réelles afin de valider l'analyse des méthodes et confirmer les similarités théoriques démontrées.

## Resumo do capítulo [português]

O capitulo 5 apresenta as contribuições que dizem respeito à subtração adaptativa de múltiplas. Vários métodos são desenvolvidos tendo como objetivo unificá-los. A comparação deles será facilitada graças a essa etapa.

A seção 5.2 apresenta métodos baseados na minimização de normas $\ell_q$. Esses métodos são definitivamente os mais populares. O método de minimização dos mínimos quadrados (norma $\ell_2$) tem a vantagem de possuir uma solução analítica. A minimização da norma $\ell_1$ sugere uma estimação das primarias mais esparsas e mais próximas da distribuição estatística de reflexões primárias.

A seção 5.3 diz respeito aos métodos provenientes da análise de componentes independentes (método ICA) e da separação cega de fontes. No contexto da subtração adaptativa, a hipótese principal é de que as reflexões primárias e múltiplas podem ser modelizadas como variáveis aleatórias estatisticamente independentes. Os métodos baseados na negentropia, Infomax e kurtosis são descritos.

Os métodos de reconhecimento de forma não são inclusos na unificação sugerida. É importante mencionar estes métodos para se ter uma melhor perspectiva de suas diferenças em comparação aos métodos de filtragem mencionados anteriormente. A seção 5.4 é dedicada a estes métodos.

A contribuição principal é desenvolvida na seção 5.5. O conjunto dos métodos baseados nas normas e na independência estatística são analisados considerando as funções objetivo minimizadas em cada um dos casos. Essa ação permite separar de um lado a função otimizada e do outro a estratégia de janelamento (mencionada em seguida). Será demonstrado que o conjunto dos métodos procura minimizar uma correlação não-linear entre as primarias estimadas e as múltiplas previstas. A não-linearidade é realizada por um operador agindo como um compressor sobre as primarias. Será demonstrado igualmente que certos métodos apresentam fortes similaridades.

A dimensão do filtro e de sua zona de aplicação são os parâmetros essenciais reagrupados sob o termo de estratégias de janelamento. Elas são apresentadas na seção 5.6. Por fim, a seção 5.7 apresenta varios experimentos feitos utilizando dados reais com a intenção de validar a análise teórica dos métodos e confirmar as similaridades teóricas demonstradas.

## 5.1 Introduction

It is of prime importance to accurately remove multiples in seismic data set for most of the imaging methods (see the previous chapter 4). As explained before, some methods for multiple attenuation are based on a prediction of the multiples (e.g. SRME) and require an *adaptive subtraction* step. Those methods are often refereed to as *two-step methods* (the two steps being the prediction and the adaptive subtraction). The subtraction step not only removes the filter ambiguity due to the method itself (e.g. auto-convolutions of the data for SRME) but also tackles errors (amplitude, phase, kinematic) coming from an imperfect acquisition geometry. We quote Spitz [1999]:

> [...] elimination of the multiples should be seen as a two-step process. These two steps are independent and equally important.

The terminology of *adaptive subtraction* comes from the idea that the multiples must be slightly modified before being removed from the data. Most of the methods accurately predict the kinematic of multiple events, but fail to recover the correct amplitude and phase. This slight modification of the prediction can be done by a filter. Within a linear model, the coefficients of a linear filter are optimized. Some other methods directly optimize the full waveform of the data and are called *separation method*. In that case, the coefficients of the primaries and the multiples are optimized. Adaptive subtraction methods are cheaper in term of computation, as the filters are generally small.

There is a wildlife of methods in the literature for adaptive subtraction (see table 4.2). The first objective of this chapter is to make a distinction between them with a unifying framework easier to catch, with an emphasis on methods based on statistical independence. In particular, some authors claimed that ICA based methods could overcome the overlapping problem in adaptive subtraction. Verschuur [2013a] even writes about ICA based methods

> [...] In this method multiples and primaries are not separated via matching filters, but via independent component analysis, which drops the orthogonality requirement of primaries and multiples in standard least-squares subtraction.

In this chapter, we will see the limit of this assertion and we will be able to answer the second question proposed in chapter 1: in which manner ICA methods do really improve adaptive filtering? Are they different? We will focus on the optimization problem for each method.

## 5.2 Methods based on $\ell_q$-norm matching filters

Matching filter methods are parameter estimation problems solved via optimization approaches. They consider that the prediction $\boldsymbol{m}$ of the multiples must be adapted to better match the data $\boldsymbol{d}$ and so better remove the noise. We consider the following linear model

$$\begin{cases} \boldsymbol{d} = \boldsymbol{p_0} + \boldsymbol{m_0}, \\ \hat{\boldsymbol{m}} = \boldsymbol{w} * \boldsymbol{m}, \end{cases} \tag{5.1}$$

where $\boldsymbol{m}$ are the predicted multiples, $\boldsymbol{p_0}$ and $\boldsymbol{m_0}$ are the true primaries and multiples respectively. A hat $\hat{\cdot}$ indicates an estimation and the symbol $*$ denotes a consistent

convolution product that can be either 1D, 2D or 3D, according to the dimension of $\boldsymbol{w}$ and $\boldsymbol{m}$. The estimate of the primaries is then given by

$$\hat{\boldsymbol{p}} = \boldsymbol{d} - \hat{\boldsymbol{m}} = \boldsymbol{d} - \boldsymbol{w} * \boldsymbol{m}. \tag{5.2}$$

To find a filter, we need to formulate an optimization problem of the form

$$\text{find} \quad \boldsymbol{w} \quad \text{such that} \quad \phi(\hat{\boldsymbol{p}}(\boldsymbol{w})) \quad \text{is minimum}, \tag{5.3}$$

where $\phi(\boldsymbol{w})$ is an objective function to be defined. The most common objective functions are based on the $\ell_p$-norm as we defined previously in equation 2.3 which is defined for a vector $\boldsymbol{s}$ as

$$\|\boldsymbol{s}\|_p = \left( \sum_x |s_x|^p \right)^{1/p}. \tag{5.4}$$

The $\ell_2$-norm refers to the classical Euclidean distance, whereas the $\ell_1$-norm is simply the sum of the absolute values of the vector components (see equations 2.4).

### 5.2.1 Least-squares filtering

In multiple subtraction, the most commonly adopted approach is based on the $\ell_2$-norm such as

$$\phi_{\ell_2} = \|\boldsymbol{d} - \boldsymbol{w} * \boldsymbol{m}\|_2. \tag{5.5}$$

This norm is quite convenient from a mathematical point of view because it admits an analytical solution when a linear model is considered [Haykin, 2013]. The objective function in equation 5.5 can be written in matrix form

$$\phi_{\ell_2} = \|\boldsymbol{d} - \boldsymbol{M}\boldsymbol{w}\|_2, \tag{5.6}$$

where $\boldsymbol{M}$ and $\boldsymbol{w}$ are respectively a matrix and a vector constructed such that $\boldsymbol{M}\boldsymbol{w} = \boldsymbol{w} * \boldsymbol{m}$. This construction can be done for the 1D, 2D or 3D convolutional product according to the dimension of $\boldsymbol{w}$ and $\boldsymbol{m}$. For instance, if a 1D filter with three coefficients is computed for a trace, the Toeplitz matrix is constructed as

$$\hat{\boldsymbol{m}} = \begin{bmatrix} \hat{m}_1 \\ \hat{m}_2 \\ \vdots \\ \hat{m}_{N_t-1} \\ \hat{m}_{N_t} \end{bmatrix} = \begin{bmatrix} 0 & m_1 & m_2 \\ m_1 & m_2 & m_3 \\ \vdots & \vdots & \vdots \\ m_{N_t-2} & m_{N_t-1} & m_K \\ m_{N_t-1} & m_{N_t} & 0 \end{bmatrix} \begin{bmatrix} w_1 \\ w_2 \\ w_3 \end{bmatrix} = \boldsymbol{M}\boldsymbol{w}. \tag{5.7}$$

The damped least-squares solution gives

$$\boldsymbol{w}_{\ell_2} = (\boldsymbol{M}^T\boldsymbol{M} + \zeta\boldsymbol{I})^{-1}\boldsymbol{M}^T\boldsymbol{d}, \tag{5.8}$$

where the term $\zeta\boldsymbol{I}$ regularizes the inversion of $\boldsymbol{M}^T\boldsymbol{M}$ if necessary (see subsection 2.2.2).

### 5.2.2 $\ell_p$-norm filtering

It is well known that the $\ell_2$-norm filter, also known as Wiener or least-squares filter, may lead to over-attenuation issues when primaries and multiples overlap, as shown in subsection 4.5.2 [Abma et al., 2005]. Guitton and Verschuur [2004] analyzed the use of the $\ell_1$-norm objective function

$$\phi_{\ell_1} = \|\boldsymbol{d} - \boldsymbol{w} * \boldsymbol{m}\|_1, \tag{5.9}$$

and they have shown that it may lead to a sparser estimate of the primaries. Unfortunately, a direct analytic solution does not exist for the $\ell_1$-norm. Guitton and Verschuur [2004] proposed to use the Iterative Reweighted Least Squares (IRLS) algorithm to approximate the $\ell_1$-norm solution by using the objective function

$$\phi_{\ell_{1/2}} = \|\boldsymbol{F}(\boldsymbol{d} - \boldsymbol{w} * \boldsymbol{m})\|_2 \,, \tag{5.10}$$

where $\boldsymbol{F}$ is a diagonal matrix depending on the estimated primaries $\hat{\boldsymbol{p}}$ and iteratively updated with the least-squares solution given by equation 5.8. By using a specific $\boldsymbol{F}$, they have shown that their method is equivalent to consider the following objective function

$$\phi_{\ell_{1/2}} = \mathbb{E}\left\{\sqrt{1 + (\hat{p}/\epsilon)^2} - 1\right\} \,, \tag{5.11}$$

where $\mathbb{E}\{\cdot\}$ is the expectation operator and $\epsilon$ a positive constant. Their analysis suggests to use a constant $\epsilon = \max|\boldsymbol{d}|/100$.

More generally, it is also possible to consider formulations based on the minimization of a $\ell_q$-norm objective function

$$\phi_{\ell_q} = \|\boldsymbol{d} - \boldsymbol{w} * \boldsymbol{m}\|_q \,, \tag{5.12}$$

with $q \geq 1$. This formulation was adopted, for instance, by Costagliola et al. [2011] for regularization purposes. Pang et al. [2009] studied a constrained approach for $\ell_1$ minimization. Pham et al. [2014] also proposed a more general framework able to introduce different kinds of norms for both the estimated primaries and the filter coefficients.

## 5.3   Methods based on statistical independence

More recently, some authors considered that primaries and multiples can be modeled as statistical independent variables [Kaplan and Innanen, 2008; Lu and Liu, 2009; Donno, 2011; Li and Lu, 2013]. On the same fashion as the problem in equation 5.3, we can write the optimization problem as

$$\text{find } \boldsymbol{w} \quad \text{such that} \quad \hat{\boldsymbol{p}} \text{ and } \boldsymbol{m} \text{ are independent.} \tag{5.13}$$

As we described in chapter 2, a lot of metrics exist for measuring the statistical independence of variables, and have been developed with independent component analysis (ICA). As for $\ell_p$-norms (except for $p = 2$) no analytic solution exists for this problem and lot of algorithms can be found [Hyvärinen et al., 2001; Comon and Jutten, 2010] (see also chapter 2). While all these methods have in common to try to solve the problem in equation 5.13, they differ by their objective function. For instance, Donno [2011] uses a kurtosis based function, Liu and Dragoset [2013] use an information maximization (Infomax) objective function and Li and Lu [2013] use a negentropy maximization objective function. Also, some of those method are used after a classical least-squares filtering, and used only for amplitude recovery.

### 5.3.1   Infomax

Liu and Dragoset [2013] propose to use an Infomax framework for adaptive subtraction. Let consider the neural network in figure 5.1b, which represents the matching filter approach as proposed by the linear convolutive model presented in equations 4.15 and 4.16. The difference between this network and the classical formulation of a BSS
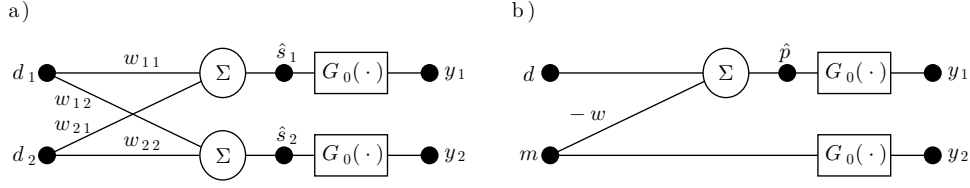
**Figure 5.1:** *(a) Neural network of Infomax algorithm with two mixtures and two sources, corresponding to the formulation of a BSS problem. (b) Adaptive multiple subtraction (equation 4.16) described as a neural network. In the image, a line indicates a convolution with the specified filter.*

problem (figure 5.1a) is that only one filter $w$ is required to perform the separation for adaptive multiple subtraction. This network could be considered as a special case of the BSS network. Therefore, the objective function of equation 2.52 also holds for the network of figure 5.1b. In this case, the statistical independence is required between the estimated primaries $\hat{p}$ and the predicted multiples $m$.

For this specific network, the Jacobian given by equation 2.53 simplifies and becomes

$$J_{IM} = \frac{\partial y_1}{\partial d}\frac{\partial y_2}{\partial m} = \frac{\partial y_1}{\partial \hat{p}}\frac{\partial \hat{p}}{\partial d}\frac{\partial y_2}{\partial m} = G_0'(\hat{p})G_0'(m). \tag{5.14}$$

Because the prediction $m$ does not change in the network, after substituting equation 5.14 in equation 2.52, the term $\mathbb{E}\{\log|G_0'(m)|\}$ is constant and the objective function to be minimized becomes

$$\phi_{IM} \quad = \quad -\mathbb{E}\{\log|G_0'(\hat{p})|\}. \tag{5.15}$$

As discussed in chapter 2, the function $G_0(\cdot)$ can be seen as an estimate of the CDF of the desired signals and the function $G_0'(\cdot)$ represents an estimate of their PDF. We can assume that the primaries follow a generalized Gaussian distribution that can either be super- or sub-Gaussian, so that its PDF is given by

$$G_0'(\hat{p}) \propto \exp(-|\hat{p}|^p). \tag{5.16}$$

The objective function to be minimized in this case becomes

$$\phi_{IM} \propto +\mathbb{E}\{|\hat{p}|^p\}, \tag{5.17}$$

that is equivalent to the minimization of the $\ell_q$-norm of the primaries as in equation 5.12.

In particular, if we assume that the primaries follow a Laplacian distribution, we can choose $G_0'(\hat{p}) \propto \exp(-|\hat{p}|)$ and so $\phi_{IM} \propto +\mathbb{E}\{|\hat{p}|\}$ which is equivalent to the minimization of the $\ell_1$-norm of the primaries as in equation 5.9. In the same way, if we assume that the primaries follow a Gaussian distribution, we can choose $G_0'(\hat{p}) \propto \exp(-\hat{p}^2)$ and so $\phi_{IM} \propto +\mathbb{E}\{\hat{p}^2\}$ which is equivalent to the minimization of the $\ell_2$-norm of the primaries as in equation 5.5.

### 5.3.2 Negentropy maximization

Li and Lu [2013] propose to use the negentropy to perform adaptive subtraction. In the case of adaptive multiple subtraction, they have shown that the objective function in equation 2.54 can be written as

$$\phi_Q = +\mathbb{E}\left\{g_i\left(\hat{p}/\sigma_{\hat{p}}\right)\right\}, \tag{5.18}$$

where $\sigma_{\hat{p}}^2$ is the variance of the estimated primaries. As they already pointed out, the use of the function $g_3(\cdot)$ in adaptive multiple subtraction leads to a formulation identical to the IRLS algorithm described in equation 5.11. Therefore, in the following we will focus on the two first non-quadratic functions $g_1$ and $g_2$ by keeping in mind that for a zero-mean signal, the normalization by $\sigma_{\hat{p}}$ is equivalent to a normalization by the $\ell_2$-norm of the estimated primaries.

In particular, if the non-quadratic function $g_2$ is used, we can write

$$\phi_Q = -\mathbb{E}\left\{\log\frac{1}{\cosh(\hat{p}/\sigma_{\hat{p}})}\right\}, \tag{5.19}$$

that is equivalent to the objective function of the Infomax matching filter in equation 5.15 with $G_0'(s) = 1/\cosh(s)$. Also, if the non-quadratic function $g_1$ is used, we have

$$\phi_Q = -\mathbb{E}\left\{\exp\left(-\frac{\hat{p}^2}{2\sigma_{\hat{p}}^2}\right)\right\}, \tag{5.20}$$

that is equivalent to the objective function of the Infomax matching filter with $G_0'(s) = \exp\left[\exp\left(-s^2/2\right)\right]$. Interestingly, it can also be seen as using the Infomax objective function after removing the log operator, such that $\phi = -\mathbb{E}\{G_0'(\hat{p})\}$ with a Gaussian prior distribution $G_0'$.

### 5.3.3 Kurtosis

The method proposed by Donno [2011] is a two steps adaptive filtering method. First, a least-square adaptive filter is applied to the prediction such that $\boldsymbol{m} \leftarrow \boldsymbol{w}_{\ell_2} * \boldsymbol{m}$. Then an ICA step is added to improve the amplitude recovery of the primaries. This step is done by a separating matrix $\boldsymbol{W}$ on the whitened observations, such that it depends only on a single parameter $\theta \in [0, \pi]$

$$\begin{bmatrix} \hat{\boldsymbol{p}}^T \\ \hat{\boldsymbol{m}}^T \end{bmatrix} = \boldsymbol{W} \begin{bmatrix} \breve{\boldsymbol{d}}^T \\ \breve{\boldsymbol{m}}^T \end{bmatrix} = \begin{bmatrix} \cos\theta & \sin\theta \\ -\sin\theta & \cos\theta \end{bmatrix} \begin{bmatrix} \breve{\boldsymbol{d}}^T \\ \breve{\boldsymbol{m}}^T \end{bmatrix}. \tag{5.21}$$

A kurtosis based function is proposed as optimization scheme

$$\min_{\theta} \quad \phi_k = -\left|\mathrm{kurt}(\hat{\boldsymbol{p}})\right|, \tag{5.22}$$

and the correct amplitude for the separated primaries is found by correlation with the data.

## 5.4 Methods based on pattern recognition

Pattern recognition methods are based on the assumption that primaries and multiples can be represented as auto-regressive processes in space. This property is often referred to as *predictability* of events in the literature [Spitz, 1999]. Indeed, the amplitude and phase (time shift) of seismic events evolve smoothly in space from one trace to another, especially if we consider a local window. Instead of predictability, one could also speak about the spatial coherency of seismic events. Comparing with the previous methods, one can see pattern recognition as an extension of matching filters with different constraints on the primaries or the multiples. These constraints are given by prediction-error filters (PEFs).

**Definition 28.** *Let $s = s[x]$ be a discrete signal and $c_s = c_s[y]$ a non zero FIR filter. The prediction-error vector is given by*

$$r = r[x] = c_s * s = \sum_y c_s[y]s[x - y]. \tag{5.23}$$

*We call $c_s$ the prediction-error filter (PEF) of $s$ with respect to an objective function $\phi$ if and only if $\phi(r)$ is minimized at $c_s$. We say that a signal $s$ is predictable by the PEF $c_s$ if $r = c_s * s \approx 0$.*

From the definition, we see that the prediction-error vector is just the residual vector associated to assuming a signal as an auto-regressive process. The PEF contains the coefficients of the AR process and is computed by minimizing an objective function. In order to avoid the trivial solution $c = 0$, one can force the vector $c$ to have a unit norm. Also, we can force the first coefficient of the PEF to be equal to one. The PEF is computed by solving

$$c_s = \arg\min_c \|c * m\|, \quad \text{such that } c[0] = 1. \tag{5.24}$$

In the literature, the term pattern refers to the sequence of amplitude and time shift function of space that can be extracted from the PEF

$$c_s = \begin{bmatrix} 1 & A_2 e^{i\phi_2} & A_3 e^{i\phi_3} & \dots \end{bmatrix}, \tag{5.25}$$

Associated with an appropriate initial time function $v[x = 0, t]$, the signal $s[x]$ can be restored over a small window. For instance, the PEF of an horizontal constant event (a perfect up-going plane wave with no angle) is given by $[+1, -1]$. A plane wave with an angle is given by $[+1, A_2 e^{i\phi_2}]$. From the PEF, the pattern can be reconstructed and the correct wavelet is found by $\ell_2$ minimization. Figure 5.2 shows basic seismic events that are perfect auto-regressive processes.

In adaptive subtraction, it is assumed that the pattern extracted from the prediction is the same as the pattern of the true multiples in the data. In other terms $c_m = c_{m_0}$, and if $\|c_m * m\|$ is minimized, then $\|c_m * m_0\|$ is minimized too. In particular, the pattern should be invariant to the change of waveform caused by an incorrect multiple prediction.

**Assumption 8.** *The PEF of the prediction is equal to the PEF of the true multiples.*

In the original form proposed by Spitz [1999], the PEF of the prediction $c_m$ and the PEF of the data $c_d$ are used to obtain the PEF of the primaries. If primaries and multiples are uncorrelated, the PEF of the primaries can be expressed as the data's PEF deconvolved by the multiple's PEF as

$$c_p = c_m^{-1} * c_d. \tag{5.26}$$

As we should have $c_p * p \approx 0$, the choice of the objective function is less important compared to a matching filter approach. The primaries can be obtained by solving [Spitz, 2000]

$$\min_w \|c_p * (d - w * m)\|_2, \tag{5.27}$$

where the first convolution is in space (for the PEF) while the second convolution is in time (for a 1D matching filter). The least-squares solution gives [Guo, 2003]

$$\hat{w} = \left( M^T C_p{}^T C_p M \right)^{-1} M^T C_p^T C_p d. \tag{5.28}$$

**Figure 5.2:** *Basic examples of predictable event (a) horizontal plane-wave, (b) plane-wave, (c) pane-wave with attenuation.*

Based on the original idea proposed by Spitz, different refinements have emerged to better tackle the signal recovery. Guitton and Cambois [1999] extended the idea in 3D to build PEFs in offset and mid-point gathers (see also [Guitton, 2003] and [Guitton, 2006]). Guitton et al. [2001] change the optimization scheme by considering the following objective function

$$\min \|d - \hat{d}\|_{\ell_2}. \tag{5.29}$$

The same approach has been used by Luo et al. [2003]. Guo [2003] slightly expanded the original idea of Spitz by adding what he calls the projection signal filter. From this early paper, not much work and improvement have been done on pattern-recognition methods. Liu and Lu [2016] recently proposed a new method based on pattern coding.

## 5.5  Unification by primary operators

It is not easy to see the difference between all the previous methods. In this section, we will unifying them, driven by the will of understanding what kind of solution we end-up with. We will make a clear distinction between the theory behind a method, the objective function and the algorithm.

All the described objective functions are convex. The first stationary condition states that the gradient should be null at the solution such that:

$$\nabla \phi(\hat{\boldsymbol{p}}^\star) = 0. \tag{5.30}$$

For a matching filter, the parameters are the coefficients of the filter $\{w_y\}$ and we write

$$\frac{\partial \phi(\boldsymbol{w})}{\partial w_y} = \frac{\partial \phi}{\partial \boldsymbol{p}} \frac{\partial \boldsymbol{p}}{\partial w_y}, \tag{5.31}$$

and

$$\boldsymbol{p}[x] = \boldsymbol{d}[x] - \sum_y w_y \boldsymbol{m}[x - y], \tag{5.32}$$

so

$$\frac{\partial \boldsymbol{p}}{\partial w_y} = -\boldsymbol{m}[x - y]. \tag{5.33}$$

**Least-squares.**   We consider the objective function proposed by Verschuur et al. [1992] and based on the $\ell_2$-norm (see equation 5.5). We have

$$\frac{\partial \phi_{\ell_2}}{\partial \boldsymbol{p}} = \frac{\boldsymbol{p}}{\|\boldsymbol{p}\|_2}. \tag{5.34}$$

**$\ell_1$-norm.**   We consider the objective function proposed by Guitton and Verschuur [2004] and based on the $\ell_1$-norm (see equation 5.9). If we consider that the derivative exists, we have

$$\frac{\partial \phi_{\ell_1}}{\partial \boldsymbol{p}} = -\mathrm{sign}\{\boldsymbol{p}\}. \tag{5.35}$$

**$\ell_{1/2}$-norm.**   We consider the objective function proposed by Guitton and Verschuur [2004] and based on the hybrid $\ell_{1/2}$-norm (see equation 5.11). We have

$$\frac{\partial \phi_{\ell_{1/2}}}{\partial \boldsymbol{p}} = \boldsymbol{p} \odot \frac{1}{\epsilon^2 \sqrt{1 + \frac{\boldsymbol{p}^2}{\epsilon^2}}}. \tag{5.36}$$

**$\ell_p$-norm.**   We consider the objective function proposed by Costagliola et al. [2011] or Pham et al. [2014] and based on the $\ell_p$-norm (see equation 5.12). We have

$$\frac{\partial \phi_{\ell_p}}{\partial \boldsymbol{p}} = \frac{\boldsymbol{p} \odot |\boldsymbol{p}|^{p-2}}{\|\boldsymbol{p}\|_p^{p-1}}. \tag{5.37}$$

**Infomax.**   We consider the objective function proposed by Liu and Dragoset [2013] and based on Information maximization (see equation 5.15). We have

$$\frac{\partial \phi_{IM}}{\partial \boldsymbol{p}} = G_0''(\boldsymbol{p}) \odot \frac{1}{G_0'(\boldsymbol{p})}. \tag{5.38}$$

**Negentropy.**   We consider the objective function proposed by Li and Lu [2013] and based on negentropy maximization (see equation 5.18). We have

$$\frac{\partial \phi_Q}{\partial \boldsymbol{p}} = G_q'(\boldsymbol{p}). \tag{5.39}$$

**Kurtosis.**   We consider the objective function originally proposed by Donno [2011] for amplitude recovery and based on Kurtosis maximization (see equation 5.22). We have for a convolutive filter

$$\frac{\partial \phi_k}{\partial \boldsymbol{p}} = \frac{4\boldsymbol{p} \odot \left(\boldsymbol{p}^2 \mu_2 - \mu_4\right)}{\mu_2^3}. \tag{5.40}$$

**Primary enhancer operator for matching filter approaches**

The cross-correlation product between two signals $\boldsymbol{s}_1[x]$ and $\boldsymbol{s}_2[x]$ is denoted $(\boldsymbol{s}_1 \otimes \boldsymbol{s}_2)[x]$. It is defined such that

$$(\boldsymbol{s}_1 \otimes \boldsymbol{s}_2)[y] = \sum_y s_1[x] s_2[y+x] \tag{5.41}$$

where the sum over $y$ can be chosen to be finite over a limited support.

We propose to unify the previous conditions by introducing the operator $\overline{G}(\cdot)$, that we call primary enhancer, which allows writing the first derivative condition as

$$\boxed{\overline{G}(\hat{\boldsymbol{p}}^{\star}) \otimes \boldsymbol{m} = \boldsymbol{0}} \, , \tag{5.42}$$

where $\otimes$ denotes the cross-correlation product, defined to be consistent with the filter size.

Figure 5.3 shows the analyzed primary enhancer operators and their application on a small seismic data window. In this context, the result of the $\ell_1$-norm matching filter can be seen as the limit of the hybrid $\ell_1/\ell_2$-norm matching filter when $\epsilon \to 0$ or also as the limit of the Infomax matching filter when $\lambda \to \infty$. Between those extreme values, the Infomax and the hybrid $\ell_1/\ell_2$-norm primary enhancer operators share strong similarities as they provide a smooth transition from the $\ell_1$-norm to the $\ell_2$-norm solution. From our observation, they are the most similar when a relation $\lambda = 1/\epsilon$ is kept. Figure 5.4 also shows how the Infomax and the hybrid $\ell_{1/2}$-norm objective functions make a smooth transition between the $\ell_2$-norm and the $\ell_1$-norm objective function depending on the value of their shaping parameter.

It is already known that the least-squares filter aims at canceling the cross-correlation between the estimated primaries and the predicted multiples in a vicinity defined by the dimensions of the filter. Equation 5.42 unifies the methods analyzed in this paper in a same fashion. They can all be seen as canceling the cross-correlation between the enhanced primaries and the predicted multiples. As we indicated in the two previous subsections, it exist equivalences between certain methods if specific shaping parameters or non-linear functions are used. In those cases, the enhanced operators are equal. Otherwise, the methods are similar and their practical differences will be discussed in the next section.

**Primary annihilator operator for pattern recognition approaches**

Let us first consider the effect of a PEF on its associated signal:

$$\frac{\partial}{\partial c_j} \sum_x g(r_x) = \sum_x \frac{\partial g(r_x)}{\partial r_x} \frac{\partial}{\partial c_j} \left( \sum_y c_y s_{x-y} \right) . = g'(r_x) s_{x-j} \tag{5.43}$$

A PEF is canceling the correlation between the signal and the residual, eventually enhanced by an operator. We can write

$$\boxed{\boldsymbol{c_p} * \hat{\boldsymbol{p}}^{\star} \otimes \boldsymbol{c_p} * \boldsymbol{m} = \boldsymbol{0}} \, , \tag{5.44}$$

where $\otimes$ denotes the cross-correlation product. Pattern recognition methods are not considered in the following analysis and tests. We wanted to make their associated objective functions more clear, compared to matching filters.

## 5.6  Comments about windowing

### 5.6.1  Windowing strategies

The regression model described in equation 5.1 is actually not stationary and the coefficients of the filter depend on the position in the data. It is more accurate to write

$$\begin{cases} \boldsymbol{d} = \boldsymbol{p_0} + \boldsymbol{m_0}, \\ \hat{\boldsymbol{m}}[x] = \sum_y \boldsymbol{w}[x, y] \boldsymbol{m}[x - y], \end{cases} \tag{5.45}$$

**Figure 5.3:** *The primary enhancer operators $\overline{G}(\cdot)$, analyzed in equations 5.34 to 5.40, are applied on the same small window of seismic data, supposedly the estimated primaries. The mean and the variance of the data have been respectively normalized to zero and one. A scaling factor have been applied for the operator of the hybrid $\ell_{1/2}$ method and the Infomax method in order to bound their value range between -1 and 1.*
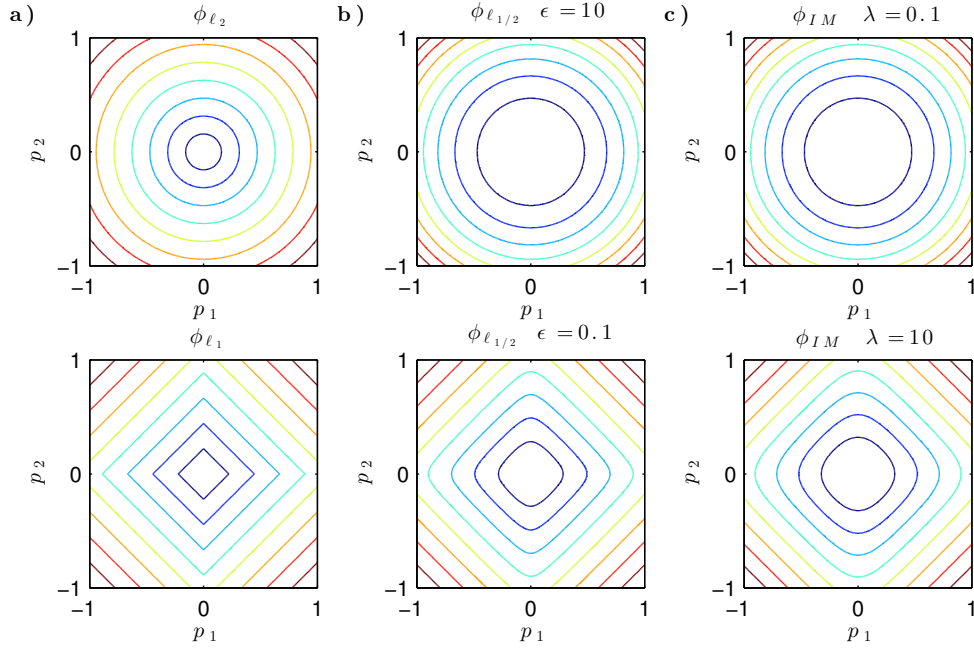
**Figure 5.4:** *Contour plots in $\mathbb{R}^2$ of: a) the $\ell_2$-norm and the $\ell_1$-norm objective functions, b) the hybrid $\ell_{1/2}$-norm objective function with $\epsilon = 10$ and $\epsilon = 0.1$, c) the Infomax based objective function with $\lambda = 0.1$ and $\lambda = 10$.*

where the coefficients depend on $x$. Fomel [2007a, 2009] describe quite well the problem of designing non-stationary regression in geophysical problem. In particular, he presents a shaping regularization method and compares it with the common Tikhonov regularization.

Because the seismic signal is not stationary neither in space and time, windowing strategies are usually used. A common approach is to divide the data into several windows in which the stationary assumption is used. In each window, eventually overlapping between each other, a filter is computed. A window can be a piece of trace, a 2D local window, or a 3D local window. The space-time coherence of the seismic signal can be used to smooth the variations of the filter and avoid drastic changes in both space and time.

Until now, in the formulation of the methods presented in this chapter, we mainly considered the 1D–1D strategy (1D filter, 1D data window) for which one wants to recover a single 1D matching filter of length $K_t$ for a segment of seismic trace of length $F_t$. However, this 1D–1D strategy does not avoid drastic change in space, but only in time. That is the reason why any matching filter based on $\ell_q$-norm or independence applied trace by trace with the 1D–1D strategy may lead to over-attenuation of the primaries if they do overlap with multiple events.

To overcome the over-attenuation problem, the 1D-2D strategy (1D filter, 2D data window) uses adjacent traces to find a 1D filter. The result of using adjacent traces in term of statistical diversity is shown in figure 5.5 in the case of crossing events. In this toy example, the prediction of the multiples is equal to the true multiples. Figure 5.5b and 5.5d respectively show the scatter plot of the primary versus the multiple (figure 5.5a) and the data versus the predicted multiple (figure 5.5c) at a single offset (in black) and in a small window (in white). We see that if a single trace containing the crossing event (in black) is considered, the primary and the multiple are highly correlated and so highly statistically dependent. Hence, any strategy

**Figure 5.5:** *Synthetic toy example of two crossing events. a) The synthetic data set containing one primary and one multiple that are overlapping at traces 20 to 30. b) Scatter plot of the primary and the multiple at trace 25 only (in black) and at traces 20 to 30 (in white). c) The prediction of the multiple that is equal, in this example, to the true multiple. d) Scatter plot of the data and the prediction of the multiple at trace 25 only (in black) and at traces 20 to 30 (in white).*

trying to make them uncorrelated (least-squares) or independent (*e.g.* Infomax) will systematically fail. In other words, over-attenuation will systematically happen with the 1D–1D strategy. However, when adjacent traces are used (in white), the primaries and the multiples became statistically independent events and a strategy forcing the independence may work. We emphasize here again that considering primaries and multiples as independent events does not help if they overlap. It is the use of adjacent traces that help overcoming the over-attenuation problem in adaptive multiple subtraction.

The 2D–2D strategy (2D filter, 2D data window) explicitly uses the coherence of the seismic signals in both space and time. It is expressed as finding a 2D matching filter of size $K_t \times K_h$, over a data window of size $F_t \times F_h$. In this case, the convolutional product is defined in two dimensions. Most of the time, we have $K_t < F_t$ and $K_h < F_h$ in order to solve a well-posed problem. A 3D strategy can also be considered using the shot number as the third dimension; the convolutional product will be defined this time on three dimensions.

**Figure 5.6:** *Curve of objective function for a training set and a validation set with increasing complexity. The validation set is used to avoid over-fitting. The complexity has to be understood as the number of parameters, here the size of the filter.*
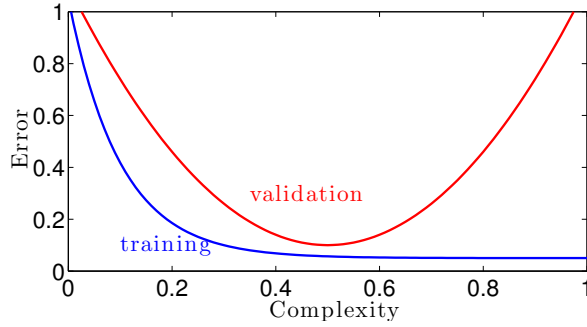
### 5.6.2 Training and validation sets

A necessary trade-off between a good fit and an over fit of the data exists in most of inverse problems [Arlot, 2010]. The same is true for adaptive subtraction. It exists a necessary trade-off between removing the predicted noise and preserving the underlying signal. In machine learning, it is common to distinguish a training set and a validation set of data. The training set is used to learn a model, while the validation set is used to monitor over-fitting. Most of the time, the same objective function is used for both sets (see figure 5.6).

In adaptive subtraction, the training set is the local patch of data $d[x_0]$ for which we want to remove the noise. The validation set is the data $d[x_1]$ around this local patch. For a classical least-squares adaptive subtraction, the objective function is simply

$$\|\boldsymbol{p}[x_0]\|_2 + \kappa \|\boldsymbol{p}[x_1]\|_2 \,, \tag{5.46}$$

where the trade-off parameter $\kappa$ weighs the importance of the validation set. If $\kappa = 0$, the noise is locally adapted to minimise the energy of the primaries on the training set only. This may lead to over-attenuation of the primaries. If $\kappa \to \infty$, the weight is put on the validation set and the noise is not locally adapted anymore. This may lead to a poor noise removal.

As we discussed in the introduction, over-attenuation is the main issue in adaptive multiple subtraction. To overcome this issue, several strategies exist that can be seen as different definitions of the validation set. For instance in one dimension, we can use adjacent traces of the current trace as the validation set. In two dimensions, we may use a small window around the local patch.

## 5.7 Experiments and comparison

### Shape of objective functions

Both figures 5.7a and 5.7b show the same toy example. They differ only in the temporal size of the filter that we aim to recover: 0.125 s in figure a and 0.35 s in figure b. A primary event $p^\circ$ is surrounded by two multiple events $m^\circ$ and a perfect prediction $\breve{m} = m^\circ$ is used to find the $\ell_2$-norm solution $\boldsymbol{w}_{\ell_2}$ (first row, last column). The theoretical filter $\boldsymbol{w}_{th}$ is known and is a perfect Dirac at $t = 0$. In the second row, we plot the value of each objective function ($\ell_2$, $\ell_1$ and Infomax) on the line passing through these two solutions of the "filter space".

For both figure a) and b), the $\ell_2$-norm objective function has its minimum at $\boldsymbol{w}_{\ell_2}$, which is normal by construction. The main difference between the two examples is the location of the minimum of the $\ell_1$-norm objective function.

In the case of a small filter $k_t = 0.125$ s, the $\ell_1$-norm has its minimum at the theoretic solution. In that case, the $\ell_1$-norm is a better strategy. Depending on the value of $\lambda$, Infomax can retrieve the $\ell_2$ or the $\ell_1$ solution. The primaries estimated by Infomax are shown in the third row of each figure.

However, when we increase the size of the filter to $k_t = 0.35$ s, the $\ell_1$-norm has its minimum really close to the $\ell_2$-norm solution: the $\ell_1$-norm is no longer a better strategy. Once again, depending on the value of $\lambda$, Infomax makes a smooth transition between $\ell_2$ and $\ell_1$.

### 5.7.1 Results on real data set

We compare the results of matching filter methods on a 2D real marine data set. The common shot gather is presented in figures 5.8a and 5.9a and the 2D SRME prediction in figures 5.8b and 5.9b. The spatial sampling is 25 m and the time sampling is 2 ms. Minimum and maximum offsets are 225 and 4700 m, respectively. The 2D SRME prediction is realized with 600 m aperture around source and receivers to reduce aliasing. Some primary events surrounded by multiple events are clearly identifiable. A global time shift correction of 40 ms is pre-applied on the prediction but no spatial correction is necessary. Hence, the matching filter we are seeking for should mainly compensate for the surface operator due to the auto-convolutions of the data during the prediction process. In a first test we use the $\ell_2$-norm matching filter with different windowing strategies. In a second test we use the same windowing strategy with different objective functions.
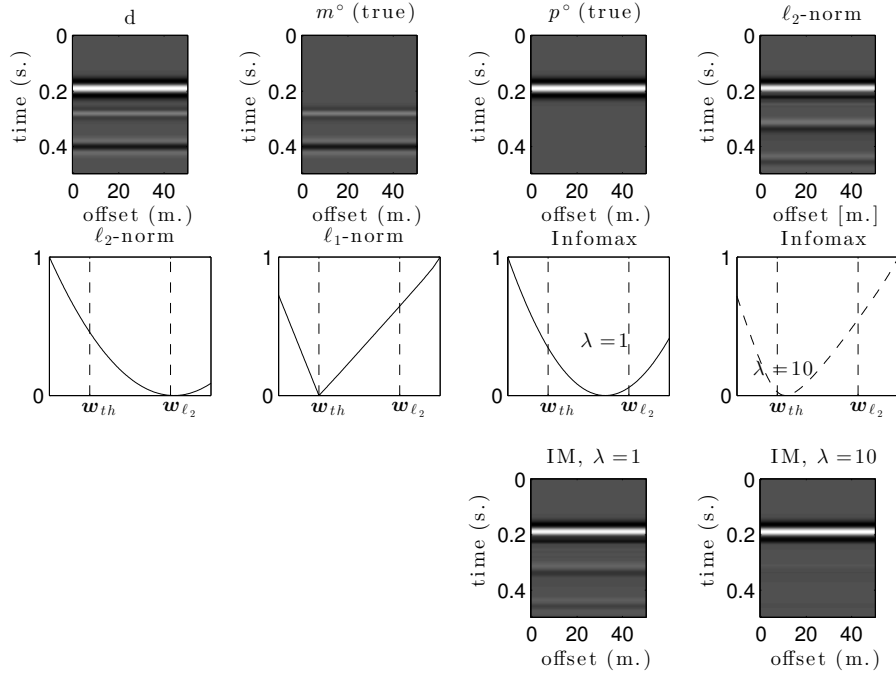
The negentropy maximization matching filter is performed with the non-linear function $g_1$ (see equation 2.55) and the IRLS algorithm proposed by Li and Lu [2013]. The hybrid $\ell_{1/2}$-norm (IRLS) has one parameter $\epsilon$ which is estimated for each window by the relation proposed by Guitton and Verschuur [2004]. Finally, our formulation of the Infomax matching filter needs the estimation of the parameter $\lambda$ that is related to the prior CDF of the primaries we tend to estimate. We propose here to use a fixed value of $\lambda$, but an adaptive scheme could also be used to take better into account the non-stationarity of the signal. As one assumes that the multiples should be removed from the signal, the primaries should have a more spiky PDF compared to the data. First, an optimum parameter $\lambda_d$ is determined to fit the data and then the value for the primaries is over-evaluated by $\lambda \approx 5\lambda_d$.

The hybrid $\ell_{1/2}$-norm and negentropy methods are implemented by using the IRLS algorithm. If an identity matrix is chosen for the initialization of the matrix of weights $\boldsymbol{F}$ in equation 5.10), they have the advantage to give the $\ell_2$-norm solution at the first iteration [Guitton and Verschuur, 2004]. A gradient method is used for the Infomax method [Liu and Dragoset, 2013] and has the advantage to actually compute the non-linear correlation between the estimated primaries and the predicted multiples for the gradient update rule. However, Infomax is generally more time consuming compared to the IRLS methods for a small matrix $\boldsymbol{M}$.

The results of the $\ell_2$-norm objective function with four windowing strategies are shown in figure 5.10 and figure 5.11. A 50% window overlapping strategy is used for all tests. In all these tests, five adjacent traces are used to compute a one dimensional filter. As expected, the increase of the temporal length of the filter leads to more attenuation of the multiples, with an eventual over-attenuation of the primaries.

The results of the objective functions described in this chapter with the same windowing

a) filter size:125ms
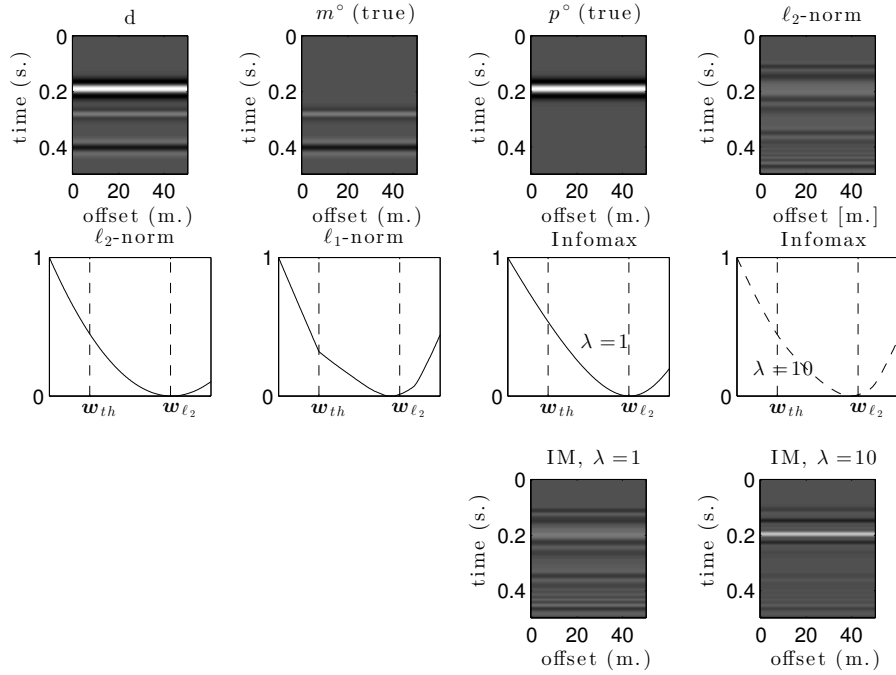


b) filter size:350ms



**Figure 5.7:** *Adaptive subtraction performed on synthetic example with two filter sizes.*

strategy are shown in figure 5.12 and figure 5.13. In all these tests, five adjacent traces are used to compute a one dimensional filter. In this example, fewer differences are visible, which is consistent with the previous theoretical analysis showing the similarities between the analyzed objective functions. The observation of fewer dissimilarities between the primaries estimated by different objective functions is also valid for other windowing strategies. This suggests that the windowing strategy has more impact on the results than the choice of an objective function. In this context, Liu and Kostov [2015] recently focused on a criterion to find a proper filter size for a given data set.

### 5.7.2 Discussion

Fundamentally, adaptive multiple subtraction is an underdetermined problem that consists of the joint recovery of a filter $\boldsymbol{w} \in \mathcal{W}$ and the primary signal $\hat{\boldsymbol{p}} \in \mathcal{P}$ from the data $\boldsymbol{d}$. Hence, for the problem to become determined in this form, it always misses a number of equations equal to the size of the filter, at least. By setting $\hat{\boldsymbol{p}} \approx \boldsymbol{0}$, the problem can become virtually overdetermined and $\boldsymbol{w}$ can be estimated with an outnumber of linear constraints. Because the primaries (what we could see as the "noise" in Wiener filtering) are not zeros (it is indeed the signal), an objective function must weight the contribution of each constraint to be able to specify which solution is the best and unique estimate $\hat{\boldsymbol{p}}$.

Most of adaptive multiple subtraction schemes consider $\ell_q$-norm matching filters for which it is assumed that the estimated primaries have minimum energy in the $\ell_q$-norm sense (equation 5.4). However, the desired geophysical solution may not coincide with the optimized solution by $\ell_q$-norms. In particular, $\ell_q$-norm objective functions have their minimum at $\hat{\boldsymbol{p}} = \boldsymbol{0}$, leading to over-attenuation problems if the outnumbered constraints in $\mathcal{P}$ are actually verifying this solution. To overcome this inherent problem, some authors recently proposed to use objective functions based on the statistical independence of primaries and multiples. However, as we have shown in this chapter, there is an equivalence between them and $\ell_q$-norm objective functions, if the right non-linear function (or parameter) is chosen to approximate independence (via Infomax or negentropy maximization). Hence, independence based objective functions share the same issue as $\ell_q$-norm objective functions because their minimum is obtained for $\hat{\boldsymbol{p}} = 0$.

From a statistical point of view, the least-squares solution ($\ell_2$-norm) assumes that primaries and multiple are uncorrelated and we must remind that correlation is a measure of linear statistical dependence. When primaries and multiples overlap, they are actually correlated and so dependent. Hence, considering primaries and multiples as independent events is not a better strategy if they do overlap. In fact, as demonstrated in this chapter, it is the use of adjacent traces that increases the statistical diversity of primaries and multiples in a given window, thus allowing to overcome the over-attenuation problem as we pointed out in figure 5.5. Moreover, we have shown that forcing the independence between the predicted multiples and the estimated primaries can be seen as a non-linear de-correlation between the same predicted multiples and the estimated primaries enhanced by a chosen operator. This operator has to be chosen to respect an a priori information about the PDF of the desired primaries.

All the methods analyzed in this chapter can be seen as adding a prior information about the statistical distribution of the primaries, so that the underdetermined adaptive multiple subtraction problem can be virtually overdetermined. If a sigmoid function is used, the Infomax method becomes really similar to the hybrid $\ell_{1/2}$-norm method as they both make a smooth transition between the $\ell_2$ and the $\ell_1$-norm solution that respectively assume a Gaussian and a Laplacian distribution. Other non-linear

**Figure 5.8:** *a) Input common shot gather from real marine data set and b) SRME predicted multiples.*



**Figure 5.9:** *Zoom of figure 5.8. The legend is the same. A primary event is indicated by a white arrow.*

**Figure 5.10:** *Primaries and multiples estimated by the $\ell_2$-norm with a 1D-2D strategy (5 adjacent traces) of: a) $F_t = 200$ ms and $K_t = 40$ ms, b) $F_t = 200$ ms and $K_t = 80$ ms, c) $F_t = 200$ ms and $K_t = 160$ ms, d) $F_t = 400$ ms and $K_t = 160$ ms.*

**Figure 5.11:** *Zoom of figure 5.10. The legend is the same. The main differences are indicated by an ellipse.*

**Figure 5.12:** *Primaries and multiples estimated with a 1D-2D strategy (5 adjacent traces) of $F_t = 200$ ms and $K_t = 80$ ms by: a) the $\ell_2$-norm, b) the $\ell_{1/2}$-norm (IRLS), c) Infomax ($\lambda = 400$), d) Negentropy maximization ($g_1$).*

**Figure 5.13:** *Zoom of figure 5.12. The legend is the same.*

functions could be used in the Infomax network, for instance an asymmetric distribution. Unfortunately, the true distribution is not known and its estimation is a difficult task. Hence, parametric methods, such as Infomax and the hybrid $\ell_{1/2}$-norm, may be challenging in practice at choosing the appropriate parameter. On the other hand, non parametric methods, such as $\ell_q$-norm or negentropy methods, are easier to use and to interpret but less flexible.

If the $\ell_1$ and the $\ell_2$ solutions are close in the parameter space $\mathcal{W}$, all the methods are expected to give similar results. It is well known that a subtle balance exists between the use o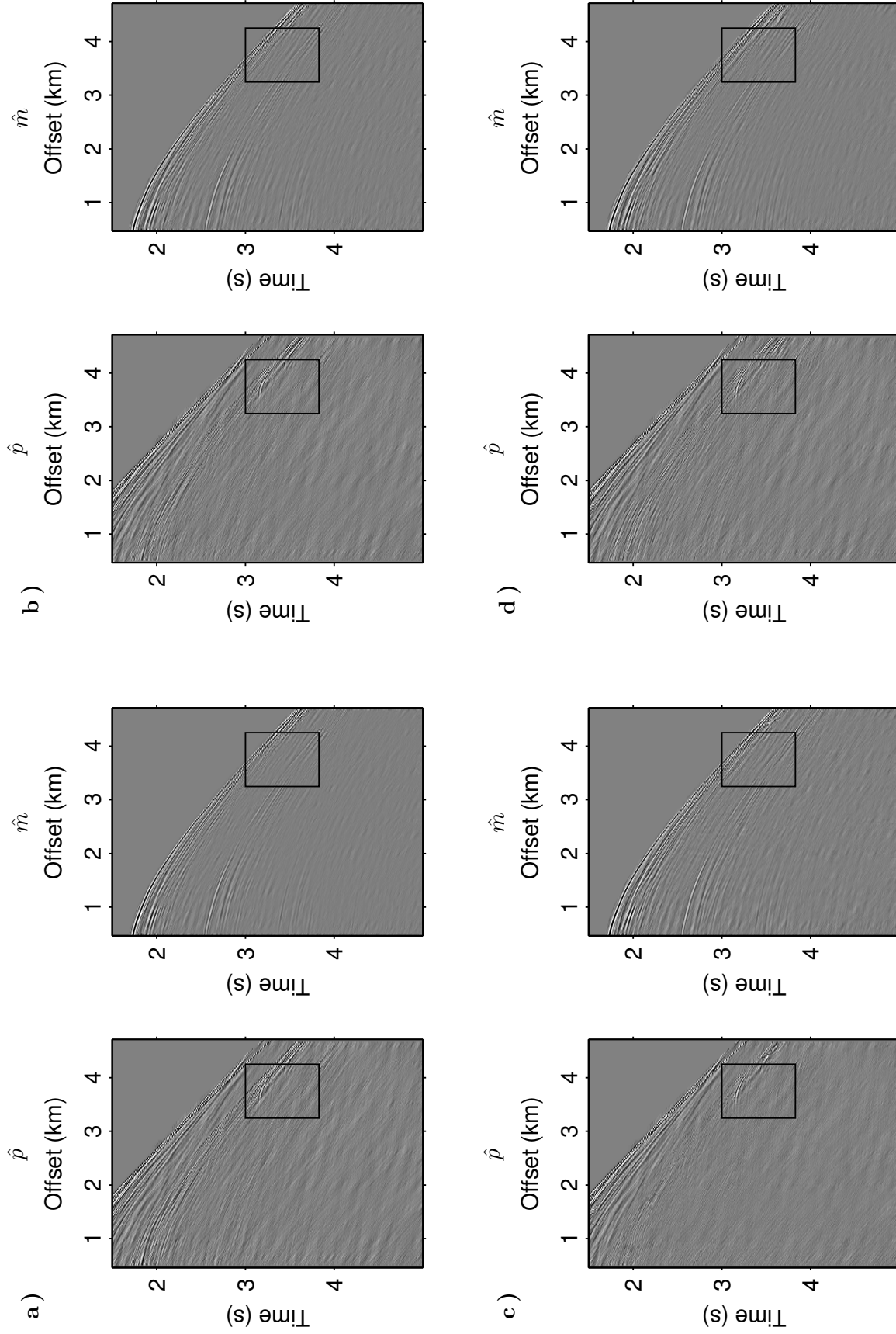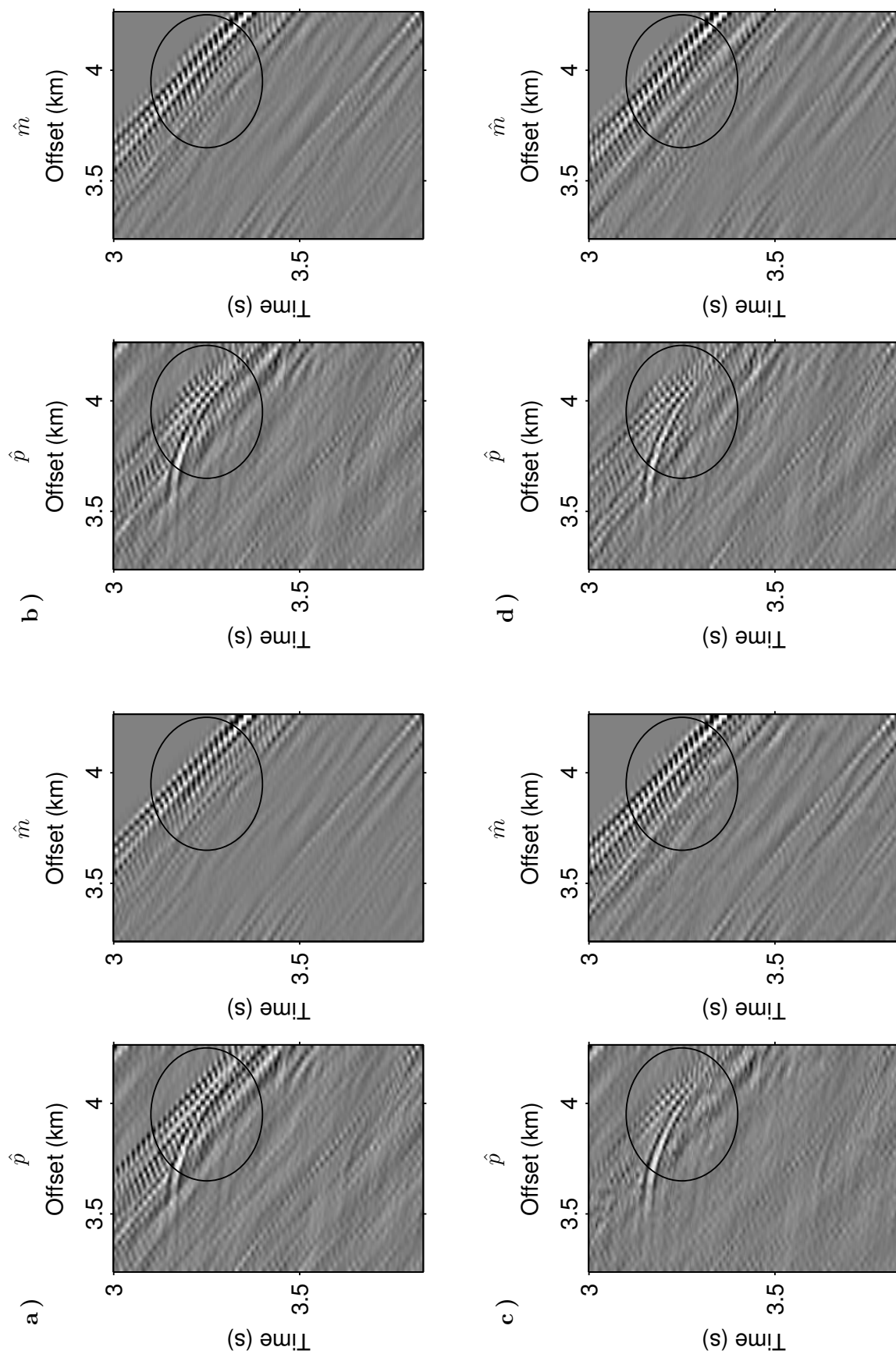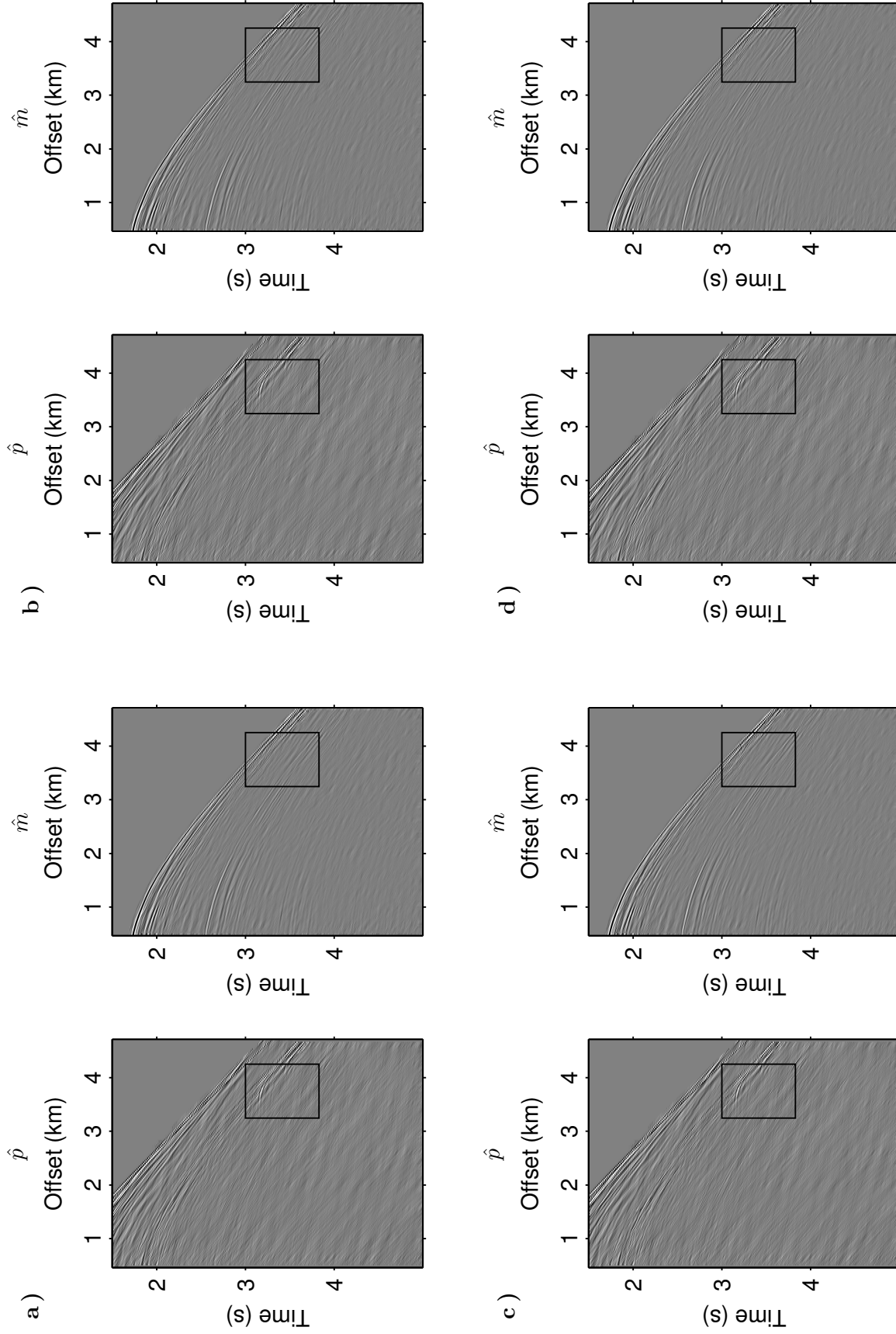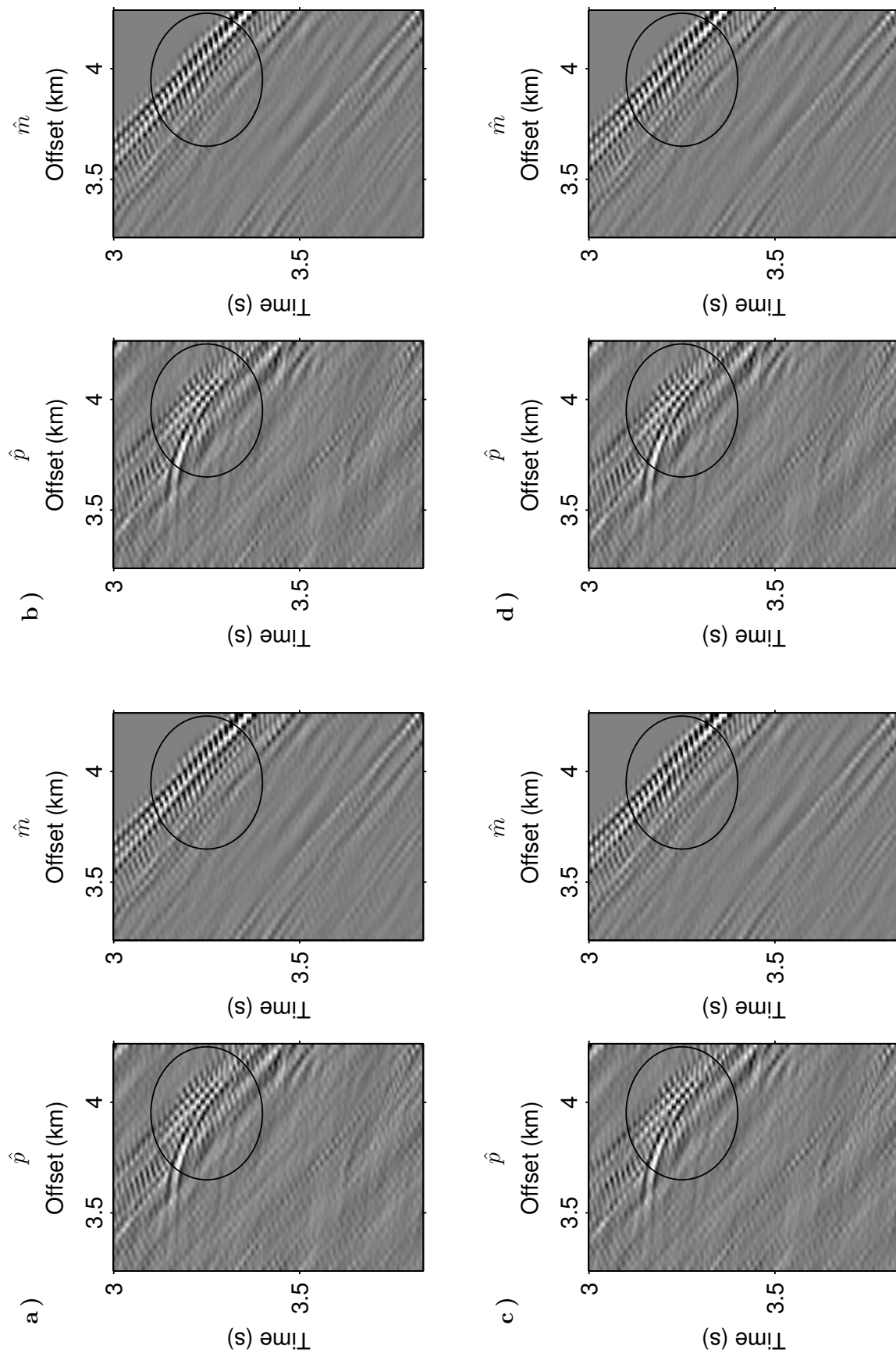f a short filter underestimating the noise and the use of a long filter overestimating it. The use of shorter filters may lead to more significant differences between two methods such as $\ell_1$ and $\ell_2$-norm, and we noticed the same behavior with the analyzed methods on synthetic examples. However, on our real data set, the statistical diversity of the windows and the need of longer filters to well attenuate the noise seem to bring closer $\ell_1$ and $\ell_2$-norm solutions, leading to fewer significant differences between the analyzed methods.

## 5.8 Conclusion

In this chapter 5, we have shown that ICA-based methods for adaptive subtraction are some kind of matching filter, with different windowing strategies. In particular, Infomax, negentropy maximization and hybrid $\ell_1/\ell_2$-norm based matching filters share strong similarities. All these techniques aim at minimizing the cross-correlation between the predicted multiples and the estimated primaries enhanced by a chosen operator. It is this operator that links all the analyzed filtering techniques. As correlation is a particular case of linear statistical dependence, the primary and multiple of a crossing event are statistically dependent. Then, forcing their statistical independence does not lead to a better solution. However, the windowing strategy, increasing the statistical diversity around the crossing event by the use of adjacent traces, is decisive as it actually allows to model primaries and multiples as independent events. In the next chapter 6, we will discuss the recovery of the FIR filter directly in the curvelet domain.

# Chapter 6

# Convolutive adaptive filtering in the curvelet domain

## Contents

## Résumé du chapitre [français]

Le chapitre 6 présente une méthodologie pour calculer des filtres à réponses impulsionnelles finies dans le domaine des curvelets. En effet, les méthodes de soustraction adaptative dans ce domaine se sont principalement concentrées sur des filtres unitaires (un coefficient réel ou complexe).

La section 6.2 présente la transformée en curvelets. Il s'agit d'une transformation redondante basée sur un découpage du domaine fréquentiel du signal. Elle a l'avantage de décrire adéquatement les événements sismiques présents dans les données et de pouvoir représenter ces données de façon parcimonieuse. La section 6.3 se consacre à la littérature autour du filtrage adaptatif dans le domaine des curvelets.

Il existe plusieurs définitions et méthodes pour calculer une transformée en curvelets. La section 6.4 présente la transformée en curvelet discrète et uniforme. En utilisant cette définition, une méthodologie est présentée à la section 6.5 pour réduire le coût de calcul des éléments nécessaire à l'obtention d'un filtre à réponse impulsionnelle finie. Des tests simples sont proposés à la section 6.6.

## Resumo do capítulo [português]

O capítulo 6 apresenta uma metodologia para calcular os filtros feitos em resposta aos impulsos finitos no domínio das curvelets. Os métodos de subtração adaptativa dentro desse domínio são principalmente concentrados nos filtros unitários (um coeficiente real ou complexo).

A seção 6.2 apresenta a transformada em curvelets. Trata-se de uma transformada redundante baseada em uma segmentação do domínio das freqüências do sinal. Essa transformação tem a vantagem de descrever adequadamente os eventos sísmicos presentes nos dados e de poder representar esses dados de maneira esparsa. A seção 6.3 é dedicada à literatura envolvendo filtragem adaptativa no domínio das curvelets.

Existem várias definições e métodos para calcular uma transformada em curvelets. A seção 6.4 apresenta a transformada em curvelets discreta e uniforme. Utilizando essa definição, uma metodologia é apresentada na seção 6.5 para reduzir o custo do cálculo dos elementos necessários para à obtenção de um filtro em resposta ao impulso finito. Testes simples são propostos na seção 6.6.

# 6.1 Introduction

As we discussed in the previous two chapters, adaptive subtraction is of prime importance for prediction-based multiple attenuation methods. Compared to separation methods, matching filter techniques have the advantages to be able to tackle larger shifts of the prediction and to be cheaper, computationally speaking. We also emphasized that some matching filter techniques require to be split into hierarchical sub-steps, from larger kinematic lag to exact amplitude recovery.

In order to perform an accurate adaptive subtraction, several choices can be made in the wide range of proposed methods. They differ essentially on: i) the objective function to be optimized, ii) the way to tackle the non-stationarity of the estimated parameters and iii) the transformed domain to choose for estimating the parameters. The previous chapter 5 has been dedicated to the objective function. We have discussed why the choice of the objective function can be less crucial than the other degrees of freedom, such as the filter size or the window size. The present chapter 6 is dedicated on the transformed domain.

The curvelet transform [Candès et al., 2006a] has shown to be efficient for a sparse representation of seismic data [Herrmann and Moghaddam, 2004; Candès and Demanet, 2005]. It has been used for adaptive subtraction in different manners, but mainly for retrieving the correct amplitude or small shifts between primaries and multiples [Herrmann et al., 2007; Neelamani et al., 2010; Wu and Hung, 2015]. In this kind of scheme, the adaptive subtraction process itself is divided into two steps: first a global filter is found and then the correct amplitudes are recovered. We believe that there is a gap to be fulfilled: FIR matching filters can be obtained directly in the curvelet domain, by promoting sparsity. Also, the curvelet transform allows for much more flexibility. For instance, gathers can be separated in direction or in scale and several filters can be computed. Donno [2011] had a similar approach by splitting gathers in different dips. In this context, Pham et al. [2015] propose a method for retrieving 2D FIR filters with dual-tree wavelets and they focus on optimization schemes (see also Pham [2015]).

In this chapter, we discuss a method to perform the search of a convolutive FIR filter directly in the curvelet domain, while limiting the computational time. Our objective is to obtain a method that is more robust in term of filter size and window size. This method is also more flexible, as several filters can be obtained by splitting scales or directions. In section 6.2 we give a description of the curvelet transform. The section 6.3 is dedicated to a review of the existing methods using curvelets for adaptive subtraction. It is important to note that several definitions and algorithms exist for the curvelet transform. The uniform discrete curvelet transform (UDCT) is presented in section 6.4. Our contribution for adaptive subtraction is explained in section 6.5. Examples are provided and discussed in section 6.6.

# 6.2 The curvelet transforms

## 6.2.1 A few words about the WNKS theorem

The Whittaker-Nyquist-Kotelnikov-Shannon (WNKS) theorem is fundamental in signal processing, and at the core of the limited redundancy of the curvelet transform. It says that any band-limited signal can be reconstructed from a finite number of samples [Mitra and Kaiser, 1993].

**Theorem 12.** *Let $s(t)$ be a band-limited signal with $\mathcal{F}\boldsymbol{s}[\omega] = 0$, $\forall \omega \notin [-1/T, +1/T]$.*

*Then, this signal can be perfectly reconstructed by*

$$s(t) = \sum_{n=-\infty}^{\infty} s[nT] \cdot \text{sinc}\left(\frac{t-nT}{T}\right). \tag{6.1}$$

*The quantity $f_N = 2/T$ is called the Nyquist rate.*

This theorem means that we can reduce the number of samples needed to fully describe any signal if we know its support in the frequency domain. In other words, it is sufficient to sample the signal at the Nyquist rate. If the signal is sampled bellow the Nyquist rate, aliasing effects may appear. It means that the frequencies higher than the Nyquist frequency (the sampling frequency divided by 2) are folded back into low frequencies. The WNKS theorem can be generalized in several dimensions [Prosser, 1966].

As we will see hereafter, the curvelet transform splits the frequency domain into several band-limited signals. Without any subsampling or decimation of each band, the redundancy of the curvelet transform should be equal to the number of band-limited filters. The WNKS theorem sets that we can actually use less coefficients if we know the limit of the band, and still be able to perfectly recover the signal in each band. This ensures the feasibility of an inverse transform and reduces the redundancy of the curvelet transform.

### 6.2.2 General considerations about curvelets

The curvelet transform is a convenient and theoretically supported way to represent seismic data, as Herrmann and Moghaddam [2004] point out. Candès and Demanet [2005] show that it is able to optimally represent wave propagators, at least for short propagation times. With a more practitioner-oriented approach, the curvelet transform locally decomposes the seismic signal (or actually any image) into band-limited local plane waves named *curvelets*. Each of these curvelets has a direction, a limited frequency content and a central position. Somehow, the curvelet transform performs what a human interpreter intuitively do when seeing seismic data: he localizes the seismic events in space and time and orders them depending on their direction and frequency content. It also corresponds to the common plane-wave approximation in small local windows taken from seismic data (see figures 6.1 and 6.2 and explanations hereafter).

The curvelet transform of a 2D signal (an image) has three parameters: the scale $j \in [0, 1, \ldots, J]$, the direction $l \in [1, 2, \ldots L_j]$ and the translation $k = [k_1, k_2]$. It is common to index those three parameters by a single one $\mu = \{j, l, k\}$ to make notations simpler. The curvelet transform of a signal $\boldsymbol{s}[x]$ is denoted $\mathcal{C}\boldsymbol{s}[\mu]$, or just $\mathcal{C}\boldsymbol{s}$ if no ambiguity exists. A single curvelet coefficient is denoted $\mathcal{C}s_\mu$ and is given by the following inner product

$$\mathcal{C}s_\mu = \langle \boldsymbol{s}[x], \boldsymbol{\varphi}_\mu[x] \rangle, \tag{6.2}$$

where $\boldsymbol{\varphi}_\mu[x]$ is a curvelet function. From the set of curvelet coefficients, the original signal can be reconstructed by summing all the curvelet functions weighted by the associated curvelet coefficient such that

$$\boldsymbol{s}[x] = \sum_\mu \mathcal{C}s_\mu \boldsymbol{\varphi}_\mu[x]. \tag{6.3}$$

As Candès et al. [2006a] point out, the curvelet functions $\boldsymbol{\varphi}_\mu$ are never explicitly computed in the image domain. Only their Fourier transform are specifically defined and used to perform the transform. In the Fourier domain, the curvelet transform

corresponds to a decomposition by a set of curvelet filters denoted $\psi_{j,l}[\omega]$. To be clear we have

$$\psi_{j,l}[\omega] = \mathcal{F}\boldsymbol{\varphi}_\mu[x], \tag{6.4}$$

so that the curvelets are defined in the frequency domain. However, in the image domain, any curvelet follows a parabolic scaling rule such that its length is almost equal to the square of its width (the width being the direction of oscillation). The figure 6.1 shows a tilling of the frequency domain for a rectangular image.

The beauty of the curvelet transform comes from a downsampling operator $\mathscr{S}^\downarrow$ which is applied after filtering by a curvelet. This downsampling operator is allowed according to the WNKS sampling theorem such that no information is lost and reconstruction is possible without aliasing effect. A curvelet transform can be roughly summarized in the two following equivalent ways

$$\mathcal{C}s_{j,l}[k] = \mathscr{S}^\downarrow_{j,l}\left(\mathscr{F}^{-1}\left(\psi_{j,l}\cdot\mathcal{F}\boldsymbol{s}\right)\right), \tag{6.5}$$

$$\text{or}\quad \mathcal{C}s_{j,l}[k] = \mathscr{S}^\downarrow_{j,l}\left(\boldsymbol{\varphi}_{j,l}*\boldsymbol{s}\right), \tag{6.6}$$

where $\mathcal{C}\boldsymbol{s}_{j,l}[k]$ contains all the curvelet coefficients at scale $j$ and direction $l$, $\mathscr{S}^\downarrow_{j,l}$ denotes the downsampling operator at scale $j$ and direction $l$ and $\mathscr{F}^{-1}$ is the inverse Fourier transform. Equations 6.5 and 6.6 are built upon the fact that convolution can be obtained as a multiplication in the frequency domain.

The curvelet functions can be defined either real-valued or complex-valued, even if the original signal is real-valued. This is simply done by splitting the two sides of the curvelet filter function $\psi_{j,l}$. It results into two complex-valued curvelet functions in the original domain, instead of a single real-valued one [Neelamani et al., 2010].

Unlike the discrete Fourier transform, they are several discrete curvelet transforms (DCTs). For instance, Candès and Demanet [2005] proposed two algorithms namely *DCT via wrapping* and *DCT via UFFT*. Nguyen and Chauris [2010] proposed the *uniform DCT* (UDCT) by using filter banks. They essentially differ on two aspects: the exact definition of the filtering functions $\psi_{j,l}$ used to tile the frequency domain (i.e. the curvelets themselves) and the downsampling operators $\mathscr{S}^\downarrow_{j,l}$ reducing the redundancy of the transform.

It is out of our scope to give a full description of the Candès's transform. However, we give a rapid idea of its transform in the following paragraph in order to better clarify the differences with the used UDCT. The description of the UDCT is done in section 6.4.

To our point of view, the main drawback of the FDCT for practical use is the downsampling operator. This operator changes at each scale $j$ and direction $l$. Also, as Nguyen and Chauris [2010] point out, the curvelets are not necessarily located on the exact same grid, as the original image. For practical applications, as for instance convolution (see section 6.5), this aspect leads to an unnecessary complication of computation. Even if the down-sampling operator is defined to optimize the redundancy of the transform, we will prefer to increase redundancy and work on a predefined grid by the use of the UDCT. We believe that this reason explains why the community has been using the curvelet transform for amplitude recovery only, and not for a complete convolutive filter recovery in adaptive subtraction.

## 6.3 Review of adaptive subtraction with curvelets

Historically, Herrmann and Moghaddam [2004] have imported the curvelet transform into the geophysics community, as well as most of the ideas from sparse recovery and

**Figure 6.1:** *Tilling of the Fourier domain performed by the discrete curvelet transforms (DCT) for a rectangular image having twice as columns as lines. The scale $j$ and the direction $l$ can be identified in the frequency domain. The lower scale $j = 0$ corresponds to the low frequency content of the data. The number of directions $L_j$ depends on the scale $j$ and the orientation. The location index $k$ emerges after the inverse Fourier transform and the down-sampling operation.*



**Figure 6.2:** *A curvelet function in the image domain from the tilling in figure 6.1 with $j = 1$ and $l = 1$.*

compressive sensing. They have shown the usefulness of this particular transform for several applications, and especially for adaptive subtraction. Herrmann et al. [2007], Saab et al. [2007a] and Herrmann et al. [2008] expand the method, firstly based on thresholding of some curvelet coefficients, into a non-linear optimization promoting sparsity. Wang et al. [2007] apply the method on real data. As Neelamani et al. [2010] point out, those methods are not only computationally expensive but they also cannot correct for high shift error.

Neelamani et al. [2010] propose to use the curvelet transform in its complex form, by splitting the real and the imaginary part. With a similar approach to Fourier transform, complex-valued filter coefficients are defined by an amplitude and a phase. The phase is able to slightly shift the curvelet atom perpendicular to its main direction. The proposed algorithm for adaptive subtraction holds in two steps. First, a classical $\ell_2$-norm approach globally adapt the prediction. Then, the correct complex-valued coefficients are found by locally minimizing the $\ell_1$-norm of the residuals. Donno [2011] proposed an approach closed to the concept of the curvelet transform, by splitting the data into dips in the $f - k$ domain. Set of filters are computed separately for each dip and the results are merged to give the final estimated primaries.

## 6.4 The uniform discrete curvelet transform (UDCT)

### 6.4.1 A few words about filter banks

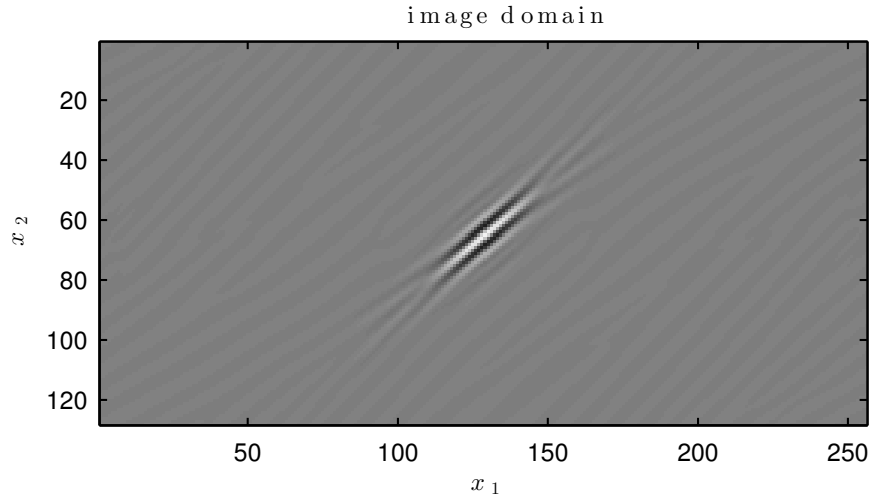A filter bank is a system of several band-pass filters followed by down-sampling operators. In other words, a filter bank performs a splitting of a signal into several frequency sub-bands. This process is often referred to as the analysis of the signal, while the reconstruction is referred to as the synthesis of the signal [Mitra and Kaiser, 1993].

With the previous definition, it is easy to see that the curvelet transform can be seen as a filter bank. The curvelet decomposition being the analysis (equation 6.5) and the reconstruction being the synthesis (equation 6.3). Nguyen and Chauris [2010] define the curvelet transform as a filter bank.

### 6.4.2 Using filter bank for the curvelet transform

The uniform discrete curvelet transform (UDCT) proposed by Nguyen and Chauris [2010] implements the curvelet transform by cascading $J$ multi-band filter banks. Each filter bank consists of one smooth low-pass filter around half the Nyquist frequency and $L_j^{vert.} + L_j^{horiz.}$ directional smooth high-pass filters. The resulting filtered signals are decimated by factors of 2 following the WNKS theorem. Hence, sequentially, a sub-level filter bank is applied on the low-frequency signal of the precedent filter bank. The number of directional high-pass filters $L_j$ at each scale is determined in order to follow the parabolic rule. The synthesis filters implementing the inverse transform are the same as the analysis filters.

Compared to the two FDCT proposed by Candès et al. [2006a], the UDCT has the great advantage to use a simpler down-sampling operator, that is a decimation by a power of 2. Hence, the translation coefficients $k \in \mathbb{Z}$ at all scales and directions are subsets of the original signal grid (see also figure 6.3). As we will see in hereafter, this property is helpful for a fast and easy computation of convolutions with the UDCT coefficients.

In the frequency domain, the curvelet filter at scale $J = 0$ is defined as an inner

coarse rectangle. The curvelet filters at higher scales are defined as the product of a rectangular band pass and a directional band-pass filter(see figure 6.1). Each function is smoothed at the edges to avoid oscillations in the image domain. This tilling is a squared version of the continuous curvelet transform [Candès et al., 2006a]. The non sub-sampled DCT is given by

$$\mathcal{C}\boldsymbol{s}_{j,l}[x] = (\boldsymbol{s} * \boldsymbol{\varphi}_{j,l})[x]. \tag{6.7}$$

In the UDCT, a difference is made between the mostly horizontal curvelets and the mostly vertical curvelets. At the end of the sequence of filter banks performing the entire curvelet transform, the decimation ratios are simply powers of 2. They depend on the scale $j$ and the main dip (*vert.* or *horiz.*) of the curvelet associated to the direction $l$. We can concatenate the decimation ratios in a matrix $\Delta_{j,l}$ such as

$$\Delta_j^{horiz.} = 2^{J-j} \begin{bmatrix} 2 & 0 \\ 0 & \frac{1}{3}2L_j^{horiz.} \end{bmatrix}, \tag{6.8}$$

and

$$\Delta_j^{vert.} = 2^{J-j} \begin{bmatrix} \frac{1}{3}2L_j^{vert.} & 0 \\ 0 & 2 \end{bmatrix}. \tag{6.9}$$

Finally, the UDCT can be constructed following

$$\boxed{\mathcal{C}\boldsymbol{s}_{j,l}[k] = (\boldsymbol{s} * \boldsymbol{\varphi}_{j,l})[\Delta_j x]}, \tag{6.10}$$

where the convolution is actually performed in the frequency domain.

## 6.5    Convolution in the curvelet domain

### 6.5.1    Convolution theorem with curvelets

The convolution theorem associated to the Fourier transform is well-known and massively used in common applications. It sets that a convolution in the original domain is equivalent to a multiplication in the Fourier domain. We consider two signals $\boldsymbol{a}(x)$ and $\boldsymbol{s}(x)$. Their Fourier transforms are denoted $\mathcal{F}\boldsymbol{a}(\omega)$ and $\mathcal{F}\boldsymbol{s}(\omega)$. From the convolution theorem, we can write

$$\mathcal{F}(\boldsymbol{a} * \boldsymbol{s})(\omega) = \mathcal{F}\boldsymbol{a}(\omega) \cdot \mathcal{F}\boldsymbol{s}(\omega). \tag{6.11}$$

For the curvelet transform, such a theorem may not be obvious at first sight. Rajendran and Rajakumar [2014] propose an analysis and a convolution theorem of the curvelet transform[1].

**Theorem 13.** *Convolution theorem for curvelet transform (Rajendran and Rajakumar [2014]: theorem 18 page 274). We consider two signals $\boldsymbol{a}(x)$ and $\boldsymbol{s}(x)$, $x \in \mathbb{R}^2$. We denote $\mathcal{C}\boldsymbol{a}(\mu)$ the continuous curvelet transform of $\boldsymbol{a}(x)$. Two convolutional products are defined such that*

$$(\boldsymbol{a} * \boldsymbol{s})(x) = \int_{\mathbb{R}^2} \boldsymbol{a}(x-y)\boldsymbol{s}(y)dy, \tag{6.12}$$

$$(\mathcal{C}\boldsymbol{a} \diamond \boldsymbol{s})(\mu) = \int_{\mathbb{R}^2} \mathcal{C}\boldsymbol{s}(j,l,k-y)\boldsymbol{s}(y)dy. \tag{6.13}$$

*We can write:*

$$\mathcal{C}(\boldsymbol{a} * \boldsymbol{s})(\mu) = (\mathcal{C}\boldsymbol{a} \diamond \boldsymbol{s})(\mu). \tag{6.14}$$

---

[1]The mathematical details of this reference are not within the competence of the author of the present thesis. This article is cited because it explicitly suggests a convolution theorem related to the curvelet transform.
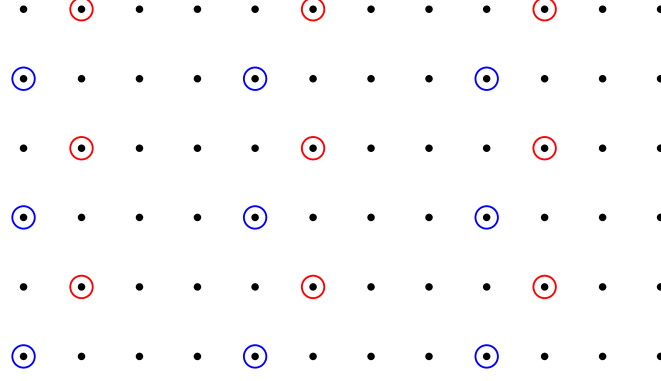
**Figure 6.3:** *The black dots indicate the original grid, as well as the position of the non subsampled curvelet coefficients. The blue circles indicate the location of the samples taken after decimation. The red circles indicates which samples must be taken for evaluating the curvelet transform of the shifted signal, for the coefficient $w_{-1,-1}$.*

In words, the curvelet transform of a convolution remains a convolution. This result actually justifies our analysis of blind source separation of sparse sources in convolutive mixtures, presented in chapter 3. Convolutional mixtures cannot be tackled in the frequency domain with the instantaneous model if the sparsity of the signal occurs in the curvelet domain.

### 6.5.2 Fast construction of the convolutional matrix

We consider the adaptive subtraction problem in the curvelet domain for which the estimated primaries are given by

$$\mathcal{C}(\hat{\boldsymbol{p}}[x]) = \mathcal{C}(\boldsymbol{d}[x]) - \sum_y w[y]\mathcal{C}(\boldsymbol{m}[x-y]). \tag{6.15}$$

In this formulation, the key point is the computation of $\mathcal{C}(\boldsymbol{m}[x-y])$, i.e. the curvelet transform of each shifted version of the prediction.

A naive way to construct all shifted version of the prediction is to actually shift the original image and compute the curvelet transform. By repeating this operation for each shift, we can build the convolutional matrix $\boldsymbol{M}$. This method is computationally expensive, because it requires $Y$ times the cost of the curvelet transform.

A second approach makes use of the interpolation formula. Once the curvelet transform of the image has been computed, the interpolation filter $\boldsymbol{h}$ can be used to compute the shifted version of the curvelet transform for the size of the filter. However, the main drawback of this method is the infinite sum in the interpolation formula. Truncated version could be used, for instance with Lanczos kernel, but the accuracy of the method will be injured. This approach has been used for instance by Chauris and Nguyen [2008].

We propose to capitalize on the useful property of the UDCT that we emphasized in the previous section: the curvelet coefficients $\mathcal{C}_{j,l}\boldsymbol{s}[k]$ are localized on the original image grid denoted $[x]$. Hence, from the non sub-sampled curvelet transform, we can construct the curvelet transform of any shifted image. This approach has the advantage to be exact. The figure 6.3 shows the non sub-sampled grid and the sub-sampling approach corresponding to the construction of the convolutive matrix $\boldsymbol{M}$.

**Figure 6.4:** *Synthetic example of two crossing events (on the left). The prediction (on the right) is perfect. The red rectangle indicates the local window where the filter is computed.*

### 6.5.3  Implementation details of the convolution matrix

Our method computes a convolutive matrix $\boldsymbol{M}$ containing the curvelet coefficients of each shifted version of the prediction such that the final multiple estimate is given by $\mathcal{C}\hat{\boldsymbol{m}} = \hat{\boldsymbol{w}}^T \boldsymbol{M}$. The number of lines is given by the size of the filter. The number of columns is given by the size of the considered local window and the redundancy of the curvelet transform. The advantage of the UDCT is that the number of coefficients to be actually stored can be reduced, and so the needed memory .

Our algorithm can be summarized step-by-step as:

- consider a 2D gather;
- compute the non sub-sampled DCT (equation 6.7);
- consider a local window;
- shift the window and apply sub-sampling for each filter coefficient (equation 6.10);
- store the vectorized results in the convolutive matrix $\boldsymbol{M}$.

The described pseudo-algorithm can be applied for each $\mu = \{j, l, k\}$, independently from the others. Hence, the needed memory is limited, as the full non sub-sampled DCT is never stored completely.

## 6.6  Experiments

The main differences between the matching filter methods are i) the objective function, ii) the parameters handling non-stationnarity and iii) the transformed domain. In this section, we use the Infomax objective function and compare the results between the image domain (here, a CSG) and the curvelet domain. This function can, indeed, be use to make a transition between the $\ell_2$-norm and the $\ell_1$-norm solutions (see chapter 5). The objective function is explicitly given by

$$\phi_{IM,\mathcal{C}} = -\mathbb{E}\left\{\log |G'_\lambda(\mathcal{C}\boldsymbol{p})|\right\}, \tag{6.16}$$

where $G'_\lambda$ is the derivative of the sigmoid function with a shaping parameter $\lambda$. The objective function is minimized with a gradient-based approach.

### 6.6.1 Adaptive search of $\lambda$.

As we explained in chapter 5, the parameter $\lambda$ used in the Infomax framework can be adapted to better match the estimate of the CDF (equivalently PDF) of the desired signal (here the primaries). Previously, we used a rough estimate of this parameter. In the results shown in this chapter, we adapt this parameter locally to the data by maximum likelihood estimation (MLE). The log-likelihood is

$$L(\lambda) = \sum_x \log g_0(d_x) \tag{6.17}$$

$$= \sum_x \log \lambda - \lambda d_x - 2 \log \left(1 + \mathrm{e}^{-\lambda d_x}\right), \tag{6.18}$$

for which the gradient is

$$\frac{\partial L(\lambda)}{\partial \lambda} = \sum_x \frac{1}{\lambda} - d_x + \frac{2d_x \mathrm{e}^{-\lambda d_x}}{1 + \mathrm{e}^{-\lambda d_x}}. \tag{6.19}$$

For each local window, the best $\lambda^*$ fitting the data is found by line search (see appendix 8.2) such that

$$\lambda^* = \arg\max L(\lambda). \tag{6.20}$$

The obtained value $\lambda^*$ is the best parameter for explaining the data with a sigmoid CDF. However, the primaries are not supposed to have the same distribution. The primaries are, indeed, more likely to have a more spiky distribution compared to the data, because the multiples must be removed from the data. In order to have an estimate of the shaping parameter of the primaries, the obtained value $\lambda^*$ is over shifted by a factor $O_F \approx 2$ to 5 such that $L = OF \times \lambda^*$. This value is used in the Infomax network.

### 6.6.2 Example on synthetics

Figure 6.4 shows a synthetic example in which the primary and the multiple events are crossing. The red rectangle indicates the considered local window centered around the crossing point. The filter is computed over this small window. In this example, the prediction is correct, except for a time shift. The shaping parameter $\lambda$ is computed and over-fitted with $O_F = 2$.

Figure 6.5 shows the results of matching filter approaches computed in the original image domain and in the curvelet domain. The Infomax objective function is optimized in both cases. The curvelet method has a better amplitude recovery, compared to the image approach and seems to better preserve the primary event. Figure 6.6 shows the errors between the true primaries and the estimated primaries for different filter sizes. The error in the image domain increases rapidly with respect to the filter size. In the curvelet domain, the method is more robust.

### 6.6.3 Discussion: IR processes between primaries and multiples

From the work presented in this thesis, an interesting question arises: does inter-regressive processes have high probability to occur between primaries and multiples?

The proposed definition of inter-regressive processes (see section 3.3) involves a particular kind of strict clustering, *i.e.* the presence of source coefficients exactly in the hyperplane defined by the IR process. The strict definition imposing a perfect cancellation of the coefficients (see equation 3.18) has been motivated by the use of the $\ell_0$ pseudo-norm in our theoretical analysis of SCA and DCA. However, in practice, the sparsity of signals is measured by some smooth functions. Hence, our definition of inter-regressive processes can be relaxed to adapt to the sparsity measure. An IR process for the $\ell_0$ pseudo-norm is also an IR process for other sparsity measure.

Another question arises when considering the curvelet transform. If the primaries and multiples yield an inter-regressive process in the original domain (e.g. CSG or COG), does the curvelet transform make it disappear? The correlation cannot disappear completely, but the curvelet transform may be sufficient to reduce the number of points in the cluster and emphasize at the same time the disjointness of the signals. Also, the curvelet transform is defined as a tight frame, meaning that $\mathcal{C}^{-1}\mathcal{C} = \mathcal{I}$, where $\mathcal{I}$ is the identity operator. However, we generally have $\mathcal{C}\mathcal{C}^{-1} \neq \mathcal{I}$ and several curvelet representations may exist for the same signal. Among them, some representations can be sparser than others. The existence of several representations for the desired signal can be investigated. It is possible that for some of them, the condition of theorem 10 is valid.

## 6.7 Conclusion

In this chapter 6, we show that matching filter approaches for multiple attenuation can be fully performed in the curvelet domain. With this approach, there is no need for splitting the adaptive multiple subtraction step into several sub-steps, such as proposed in the literature. Instead of using the classical discrete curvelet transform, we use the uniform discrete curvelet transform that allows for a better repartition of the memory and a convenient location of curvelet coefficients. In particular, FIR convolution can be well defined. Yet, the method has been used only on synthetics in order to demonstrate the feasibility. More effort must be done on real data.

**Figure 6.5:** *a) Magnification of the red window in figure 6.4. b) The estimated primary and the estimated multiple computed in the image domain. c) The estimated primary and the estimated multiple computed in the curvelet domain. d) The central trace for the true primary (black), the estimated primary in the image domain (blue) and the estimated primary in the curvelet domain (magenta).*

**Figure 6.6:** *Difference between the true and the estimated primaries for different filter sizes and for both methods (image and curvelet domain).*

# Part IV

# Closure

# Chapter 7

# Perspectives and conclusions

## Contents

In this chapter, section 7.1 draws the final conclusions. Sections 7.2 and 7.3 propose some perspectives and ideas for research direction in sparse component analysis (following our developments in part II) and adaptive subtraction of seismic multiples (following our developments in part III), respectively.

# 7.1   General conclusions

The present thesis has been motivated by the problem of adaptive subtraction of multiple events in seismic acquisition. This specific problem can be seen as a particular case of blind source separation (BSS) for which the sources are the primaries and the multiples, and the observations are the recorded data and the multiple prediction. In order to better understand the BSS framework, we have been concerned with the extraction and the separation (i.e. the recovery) of sources (i.e. signals) in convolutive mixtures, with a focus on sparse signals. In this problem, the separation of sources can be performed by optimizing some unknown finite impulse response (FIR) filter. With seismic data, sparsity can be enhanced in the curvelet domain, in which the original convolution is maintained as a convolution after curvelet transform.

Independent contributions have been obtained in sparse component analysis (SCA) and adaptive subtraction. The framework used for BSS and SCA has been chosen because it is general enough to give insights for adaptive subtraction.

## 7.1.1   Blind extraction and separation of sparse sources

When a specific contrast function (i.e. objective function) is used, such as a sparsity-based function, an underlying model is assumed. In BSS problems, this underlying model may refer to some prior information about the sources. It is of fundamental matter to question the limit of such a model in order to know exactly the conditions for which any approach is valid.

We have analyzed this topic with the $\ell_0$ pseudo-norm for determined BSS problems. In particular, we have investigated the necessary and sufficient conditions on the sources for which the $\ell_0$ pseudo-norm is a contrast function. By discussing the link between clustering and sparsity, we have shown that the $\ell_0$ pseudo-norm acts as an indicator of the presence of clusters in the observations.

**(1)** We have shown that the $\ell_0$ pseudo-norm can be used for the extraction of sparse sources, as long as no inter-regressive process bigger than the inactive parts of the sources exists. In summary, an inter-regressive process is a correlation cluster present in the original signals and can be seen as an expansion of auto-regressive processes. Combined with a multi-population evolutionary algorithm, each minimum of the contrast function can be identified. They correspond to the extraction of each source. The presence of an inter-regressive process creates a parasitic minimum in the contrast function.

**(2)** We have shown that the separation of sources can be obtained by disjoint component analysis (DCA) under the same conditions, at least for the determined case with $N = 2$ sources. We have conjectured that this result can be generalized for $N > 2$. DCA is a valid method for the separation of sparse sources and can be seen, at least with the $\ell_0$ pseudo-norm, as the correct extension of SCA for the blind separation.

**(3)** We have proposed a Differential Evolution algorithm, able to perform the non-convex optimization of the smooth $\ell_0$ pseudo-norm. This algorithm can be used for both SCA and DCA. Also, an adaptive search for the shaping parameter is able to deal with an unknown presence of additive noise.

## 7.1.2   Adaptive subtraction based on matching filter techniques

Adaptive subtraction is crucial for properly removing the multiple events with prediction based method. If not, the signal can be damaged because of an over-attenuation. Matching filter techniques are commonly used for retrieving a missing FIR filter,

underlined by the mixing model. The problem of adaptive subtraction of multiple events can be seen as a particular case of blind source separation problems. This is the main reason why we investigated sparse component analysis in the part II (chapters 2 and 3). We have identified the limit of sparse component analysis, which can help to also understand to limit and the differences with other approaches, such as independent component analysis.

The choice of the objective function is a natural question for adaptive subtraction. Recently, it has been proposed to use different metrics for adaptive subtraction. In particular, new methods based on ICA have been released. As underlined before, knowing the limit of the model underlying an objective function is of prime importance, and knowing the property of the solution, independently from the algorithm, helps.

**(1)** We have shown the limit of ICA-based method for adaptive subtraction. In particular, the methods lead to a minimization of a non-linear cross-correlation pattern between the primaries and the multiples. This non-linearity depends on the method and acts as a compressor on the estimated primaries. The assertion that ICA-based method could separate correlated signals is false.

**(2)** In order to become closer to the disjoint assumption, one may use the curvelet transform. In the literature, the curvelet transform has been used mainly for amplitude recovery. This method has the disadvantage to need a pre-determined FIR filter recovering step in order to estimate a big shift of the prediction. We believe that this limit is a consequence of the very large propagation of the discrete curvelet transform proposed by Candès et al. [2006a]. This method is indeed optimal for reducing the redundancy of the transform but not for other applications, such as the convolution. We have shown how the uniform discrete curvelet transform (UDCT) proposed by Nguyen and Chauris [2010] is a good candidate. Hence, the search of a FIR filter can be done directly in the curvelet transform, with a limited number of calculation. This is an advantage from an operator point of view. Compared to classical approaches, less tuned parameters have to be determined. Also the method shows more robustness.

## 7.2  Perspectives for sparse component analysis

### 7.2.1  Complex-valued data and under-determined problems

We only considered real-valued sources in chapter 3 dedicated to SCA and DCA. To go further, the seek of necessary and sufficient conditions can be explored for complex-valued sources (original signals). This naturally arises when signals are transformed into the frequency domain but it may occur in other fields such as in radar or magnetic resonance [Adali et al., 2011; Yang et al., 2014]. However, the $\ell_0$ pseudo-norm is only sensitive to the amplitude of a complex-valued coefficient because it is defined as the number of non-zero coefficients [Mohimani et al., 2008]. This might also be the case for other measures of sparsity (see others definitions proposed by Hurley and Rickard [2009]).

Nevertheless, for some problems, retrieving the phase is not of prime importance and the focus can be restricted to the modulus only [Fiori, 2001]. In that case, a phase ambiguity can be added to the previous known ambiguities and the definition of the solution set $\mathcal{S}$ can be easily modified. The global mapping $\boldsymbol{H}$ at the solution should be the product of a permutation matrix and a diagonal (complex) matrix. Two concepts must be investigated with a specific care about *property* and *circularity* [Adali et al., 2011].

Table 7.1 summarizes those aspects and the direction where to go. Instantaneous

|                | $\mathbb{R}$-valued | $\mathbb{C}$-valued |
|----------------|---------------------|---------------------|
| instantaneous  | chapter 3           | ?                   |
| convolutive    | chapter 3           | ?                   |

**Table 7.1:** *Table of achieved theoretical results and perspective of future works for sparse component analysis, depending on the mixing model and the value of coefficients. An interrogation mark indicates possible direction.*
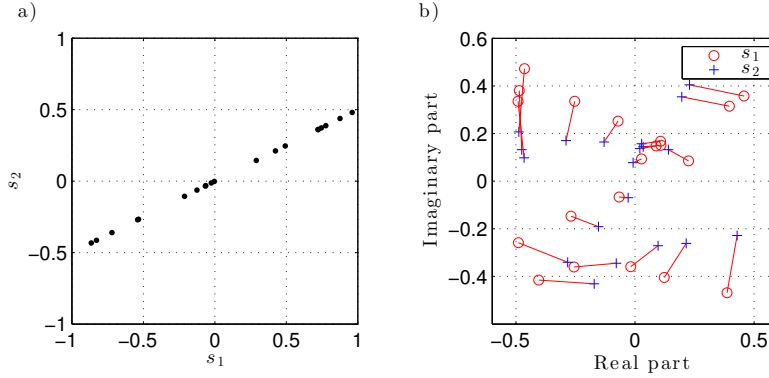


**Figure 7.1:** *a) An inter-regressive process between two sources for a real-valued problem. The parameter is $c = 0.5$. b) An inter-regressive process between two sources for a complex-valued problem. The parameter is $c = 0.8 \exp(i\pi/8)$.*

mixtures can be investigated first and then expanded to convolutive mixtures. In order to give a first idea of this direction, we consider a simple $2 \times 2$ BSS problem in $\mathbb{C}$ and we try to expand the concept of inter-regressive process for seeking necessary and sufficient conditions. We consider the mixing model

$$\begin{bmatrix} \boldsymbol{d_1}[x] \\ \boldsymbol{d_2}[x] \end{bmatrix} = \begin{bmatrix} \boldsymbol{a_{11}} & \boldsymbol{a_{12}} \\ \boldsymbol{a_{21}} & \boldsymbol{a_{22}} \end{bmatrix} \begin{bmatrix} \boldsymbol{s_1}[x] \\ \boldsymbol{s_2}[x] \end{bmatrix}, \tag{7.1}$$

for which all coefficients are complex-valued. In this simple case, an inter-regressive process can be defined as the number of coefficients such that $s_1[x] = cs_2[x]$ where $c \in \mathbb{C}$. Figure 7.1 compares the definition for real- and complex-valued signals in this simple $2 \times 2$ BSS problem. An amplitude and a rotation phase are relating the two sources for an inter-regressive process. If such a structure exists among the sources, a parasitic minimum appears in the objective function.

## 7.2.2 Inter-regressive process in different applications

In chapter 3 we defined inter-regressive process as an expansion of auto-regressive process with several signals. However, we have not discussed much their physical meaning in different contexts (see for 6.6.3 primaries and multiples). We believe additional efforts should be dedicated to the understanding of such a clustering for different contexts, and particularly in remote sensing for which SCA is of particular use. For instance, a probabilistic framework could be suited in order to determine the following probability:

$$Pr(\boldsymbol{c}^T \boldsymbol{S} = 0 \mid \boldsymbol{c}) \tag{7.2}$$

for exact inter-regressive process recovery, or:

$$Pr(\boldsymbol{c}^T \boldsymbol{S} < \epsilon \mid \boldsymbol{c}) \tag{7.3}$$

for a more robust analysis of the presence of data near a defined cluster.

|  | $\mathbb{R}$-valued | $\mathbb{C}$-valued |
|---|---|---|
| instantaneous | Herrmann et al. [2008] | Neelamani et al. [2010] |
| convolutive | chapter 6 | ? |

**Table 7.2:** *Application of the theoretical results exposed in table 7.1 for adaptive subtraction. An interrogation mark indicates not known work.*

## 7.3 Perspectives for adaptive multiple subtraction

### 7.3.1 Dealing with several orders of multiples

The multiples present in a gather are made of several orders (see figure 1.2) and we can write

$$\boldsymbol{m_0}[x] = \boldsymbol{m_0}^{(1)}[x] + \boldsymbol{m_0}^{(2)}[x] + \dots \tag{7.4}$$

In the present work, we considered the problem of adaptive subtraction for which orders of multiples are not overlapping. In order term, the subsets of indices in which each order is active are disjoints. Hence, in each small window that has been considered, the mixing model is given by

$$\begin{cases} \boldsymbol{d}[x] &= \boldsymbol{p_0}[x] + \boldsymbol{m_0}[x] \\ \boldsymbol{m}[x] &= (\boldsymbol{a} * \boldsymbol{m_0})[x]. \end{cases} \tag{7.5}$$

However, if several orders of multiples are overlapping, this assumption is not valid anymore. We must consider the following model

$$\begin{cases} \boldsymbol{d}[x] &= \boldsymbol{p_0}[x] + \boldsymbol{m_0}[x] \\ \boldsymbol{m}[x] &= (\boldsymbol{a_1} * \boldsymbol{m_0}^{(1)})[x] + (\boldsymbol{a_2} * \boldsymbol{m_0}^{(2)})[x] + \dots \end{cases} \tag{7.6}$$

in which each order of multiple has a different mixing filter. In a more BSS-like formulation, and with abuse of notations we have

$$\begin{bmatrix} \boldsymbol{d} \\ \boldsymbol{m} \end{bmatrix} = \begin{bmatrix} 1 & 1 & 1 & \dots \\ 0 & \boldsymbol{a_1}* & \boldsymbol{a_2}* & \dots \end{bmatrix} \begin{bmatrix} \boldsymbol{p} \\ \boldsymbol{m_0}^{(1)} \\ \boldsymbol{m_0}^{(2)} \\ \vdots \end{bmatrix}. \tag{7.7}$$

This model has been considered by Lu [2006] but only for amplitude recovery. Using this model for FIR filter recovery is a direction for investigations. Also, the recovery conditions for this problem must be investigated. Figures 7.2 and 7.3 show two examples in which the number of orders can be found by clustering in the observation space along axes with different angles. In the first, only one order is present and so two pics are visible in the histogram of the clustered values. In the second, three pics indicate the presence of two orders of multiples in the prediction.

From a BSS point of view, the mixing system described by equation 7.7 is an underdetermined BSS problem. In such a case, the inverse mixing system cannot be directly estimated as for determined and overdetermined problems because such an inverse does not exist. Instead, the mixing model must be identified from the data only. Then the sources are recovered by, for instance, sparse signal recovery. It is well known that if the sources are sparse enough in the considered domain, the mixing system can be identified. The observations will be indeed aligned along axes in the observation space.

However, most of the work about under-determined SCA has mainly been focussed on the instantaneous mixing model. Here again, the convolutive model has been

considered in the frequency domain, in which convolutions become multiplications. As we already explained, we expect the seismic data to be sparse in the curvelet domain. To our knowledge, no work has been done for considering sparsity in the domain where the convolution needs to be considered. This aspect is actually the continuity of the previous section 7.2. As for the present thesis, a complete work could be done with the theoretical aspects considering a fully blind problem, then it could be applied to adaptive subtraction.

### 7.3.2    Over-attenuation detection using SCA

Two assumptions are important when considering an adaptive subtraction problem: un-correlation and disjointness of primaries and multiples. Considering two random variables, the first implies $\mathbb{E}\{s_1 s_2\} = 0$ while the second implies $s_1 s_2 = 0$, which is a particular case (or a subgroup) of un-correlation.

Disjointness of the primaries and multiples (the sources) leads to clustering of the observations along two main directions. Identifying those two directions is a valid method for recovering the sources. This is in essence what SCA based on the $\ell_0$ pseudo-norm performs. Previously, we discussed how this clustering can be used for under-determined problems in order to retrieve, blindly, the mixing system. The recovery of the sources needs an efficient second step (e.g. sparse signal recovery).

We believe that clustering can be used, associated with a quantitative marker, for detecting the areas in a gather where over-attenuation may occur. In other term, SCA could be used for detecting problematic small windows. If the number of observations clustered along specific axes is larger than two, the test indicates a potential problematic area.

Based upon such a detection, the window size and the filter size could be adjusted until the data and the prediction present some kind of clustering, meaning that enough diversity is present in the observations. This method could be an alternative to the work initiated by Liu and Kostov [2015] for determining the optimal filter size.

### 7.3.3    Filters estimation from Pareto sets of filters

As discussed in chapter 5, the filter size and the window size are crucial parameters for adaptive subtraction, and their values influence the results of the estimated primaries. Liu and Kostov [2015] propose to use the Akaike information criterion (AIC) defined such that

$$AIC = \log \|\boldsymbol{p}\|_2 + \kappa \|\boldsymbol{w}\|_0 \tag{7.8}$$

for determining the correct filter size $\|\boldsymbol{w}\|_0$. This method is a classical regularization method, but it defines a discrete Pareto set instead of a continuous one. In this method, the parameter $\kappa$ has a large influence and it must be pre-determined by the operator with a series of trial values.

With this method, the operator does not have to give the filter size, but still has to find the correct regularization parameter $\kappa$. We believe that the method can be improved by examining the Pareto set of filter solutions and some clustering techniques. Figure 7.4 shows a synthetic example in which two different classes of filters can be identified in the AIC/Pareto curve: in each side of the drop, there is one group of filters clustered around the same filter. This information is the relevant one for this example. The barycenter of each cluster could be used as an optimal filter and proposed to the decision maker.

**Figure 7.2:** *Synthetic example for a single-order estimation in adaptive subtraction. a) The data contain several events, with one order of multiples. b) The prediction is good, except for the amplitude that must be recovered. c) The cross-plot between the observations and the prediction. d) The clustering measure, function of one angle parameter θ, shows the number of orders. In this case, the inverse can be directly estimated.*



**Figure 7.3:** *Synthetic example for a two-order estimation in adaptive subtraction. a) The data contain several events, with two orders of multiples. b) The prediction is good, except for two different amplitudes that must be recovered. c) The cross-plot between the observations and the prediction. d) The clustering measure, function of one angle parameter θ, shows the number of orders. In this case, the inverse cannot be directly estimated, and sparse recovery should be investigated.*

**Figure 7.4:** *Synthetic example with two parallel events for the estimation of group of filters from the AIC curve. For each filter size a least-square problem is performed in order to build the AIC curve. Each filter solution is plotted. We can see two clusters in this family of filters.*

### 7.3.4 Pattern-recognition and matching filter combined

As discussed in chapter 5, two main methods exist for adaptive subtraction: matching filter and pattern recognition techniques. Both have their pros and cons and one could try to use both in a classical form such as

$$\min_{w} \|\mathcal{C}(\boldsymbol{d} - \boldsymbol{w} * \boldsymbol{m})\|_{\ell_1} \quad \text{s. t.} \quad \boldsymbol{c_p} * \boldsymbol{p} < \epsilon. \tag{7.9}$$
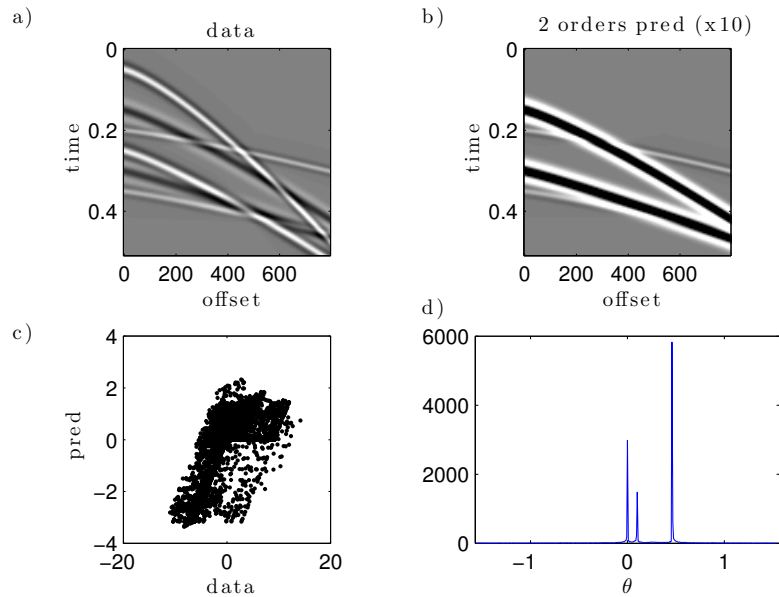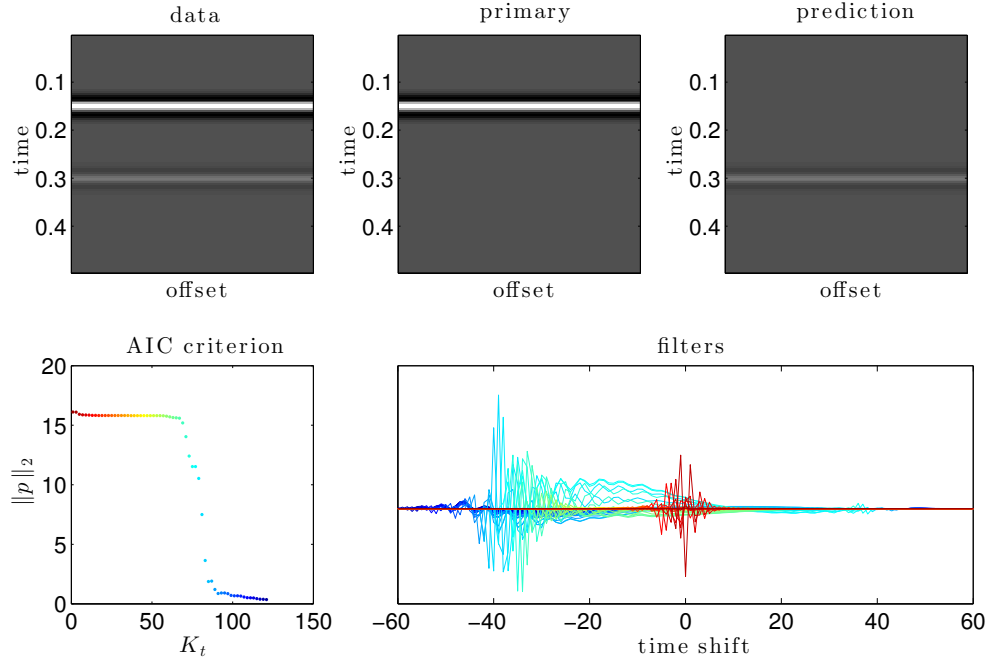
Because pattern recognition often fails at near offset, one could use higher values of $\epsilon$ at near offset and smaller values at far offset. This way, the over-attenuation issue of the main $\ell_1$-norm term could be compensated by the constraints of the auto-regressive process obtained from the prediction. With LASSO approach, this scheme could be efficient.

### 7.3.5 SCA for adaptive subtraction and the limit of dimensionality

We believe ICA and SCA are not contradictory, but *complementary*. Let us think one more time about the adaptive subtraction problem for multiple removal. For the sake of conciseness, we can think about ICA-based matching filter techniques as an englobing approach of all the methods using convex optimization (we could, indeed, always find a Bayesian framework associated).

In the perspective section 7.3, we propose to use SCA for checking the clustering of the data and the observations. If primaries and multiples are disjoint orthogonal (they cluster perfectly along two lines in the example of figure 7.3a), we have shown that ICA can be used for the separation (see section 3.2). There is no need for SCA in that case. However, if two clusters appear along with a diffuse correlation, then ICA

methods may fail. Actually, one could also keep only the data and observations points present in the clusters and perform ICA on this subsets only.

When ICA methods give satisfactory results, there is no need for other methods. Convex optimization is fast, and the methods can be used with more dimensions. In particular, using 3D data set (the data cube in standard 2D acquisitions) and 5D data set (the data 5-cube in modern 3D acquisitions) allow for more statistical diversity and more robust filter estimation. Using genetic algorithm for SCA with 5-cube may be limited by the computer memory, as the size of the population will be quite large to give accurate results. Also, curvelet transforms have been defined for 3D data cube, but not for more dimensions. In particular, the redundancy of the transform may explode and our limit of computer memory may be reached. A perfect curvelet transform may be difficult to use, but a filter-bank strategy, with a limited number of scales and directions may reduce the need of memory.

# Chapter 8

# Appendices

## 8.1 Some definitions related to norms

In this section, we would like to precise our use of the terms *norm*, *quasi-norm* and *pseudo-norm* in the manuscript. In particular, the term pseudo-norm is often used for referring to the $\ell_0$ pseudo-norm. However, strictly speaking, it is not a pseudo-norm. Several definitions are given in table 8.1. We consider a vector space $V$ with $\boldsymbol{u}, \boldsymbol{v} \in V$ and $\lambda, K \in \mathbb{R}$.

|  | Absolute homogeneity $p(\lambda\boldsymbol{u}) = |\lambda|p(\boldsymbol{u})$ | Sub homogeneity $p(\lambda\boldsymbol{u}) \leq |\lambda|p(\boldsymbol{u})$ | $p(\lambda\boldsymbol{u}) = p(\boldsymbol{u})$ | Triangle inequality $p(\boldsymbol{u}+\boldsymbol{v}) \leq p(\boldsymbol{u}) + p(\boldsymbol{v})$ | Mild triangle inequality $p(\boldsymbol{u}+\boldsymbol{v}) \leq K(p(\boldsymbol{u}) + p(\boldsymbol{v}))$ | Separation $p(\boldsymbol{u}) = 0 \Leftrightarrow \boldsymbol{u} = \boldsymbol{0}$ |
|---|---|---|---|---|---|---|
| norm | ✓ |  |  | ✓ |  | ✓ |
| semi-norm | ✓ |  |  | ✓ |  |  |
| quasi-norm | ✓ |  |  |  | ✓ | ✓ |
| pseudo-norm |  | ✓ |  | ✓ |  | ✓ |
| $\ell_0$ |  |  | ✓ | ✓ |  | ✓ |

**Table 8.1:** *Definitions related to norms. A mark ✓ indicates that the corresponding property is valid.*

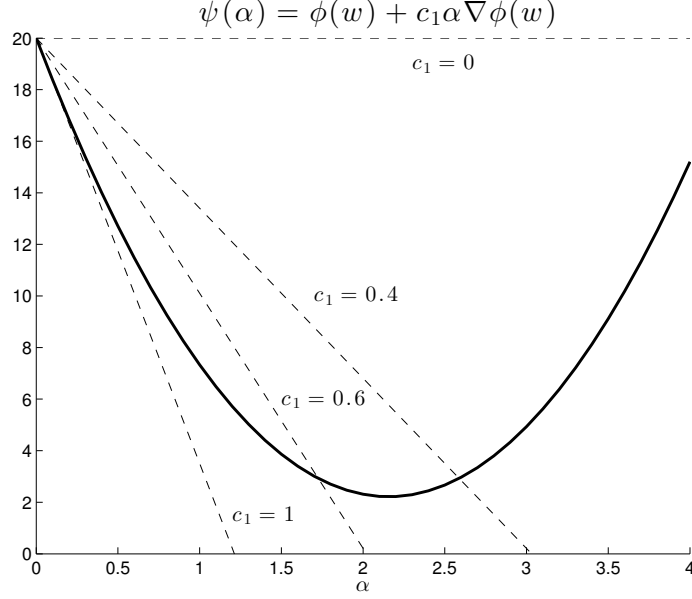$$\psi(\alpha) = \phi(w) + c_1\alpha\nabla\phi(w)$$

**Figure 8.1:** *Backtracking line search principle for obtaining the correct step $\alpha$ in gradient-based methods.*

## 8.2 Backtracking line search

We describe in this section a method for finding a correct step in gradient-based optimization methods. For non-quadratic functions, the quadratic approximation is not valid and an exact line search is impossible. Inexact line search such as backtracking has proven to be efficient in many cases [Boyd and Vandenberghe, 2004]. We consider a function $\phi(\boldsymbol{w})$ to be minimized with respect to the parameter $\boldsymbol{w}$. We consider that at the point $\boldsymbol{w}$ we have a direction of minimization $\boldsymbol{d}$ (for instance the opposite direction of the gradient). The line search problem is formulated as an optimization problem

$$\min_{\alpha} \quad \psi(\alpha) = \phi(\boldsymbol{w} + \alpha\boldsymbol{d}), \tag{8.1}$$

and backtracking search adaptively searches for a correct value $\alpha$. It starts with a fixed value $\alpha_0$ (generally small) that can be decreased if necessary by a factor $\rho \in [0, 1]$. The algorithm can be written as

$\alpha \leftarrow \alpha_0$
**while** $\phi(\boldsymbol{w} + \alpha\boldsymbol{d}) > \phi(\boldsymbol{w}) + c_1\alpha\boldsymbol{\nabla\phi^T d}$ **do**
 | $\quad \alpha \leftarrow \rho\alpha$
**end**

When the line search ends, the new current solution is obtained such that $\boldsymbol{w} \leftarrow \boldsymbol{w} + \alpha\boldsymbol{d}$. A new direction of descent is computed and the line search can be used again. This process is applied iteratively until convergence.

## 8.3 Sufficient condition for BSE

We want to retrieve the sufficient condition proposed by Duarte et al. [2011], starting from our theorem 9. We consider two sources $\boldsymbol{s_1}, \boldsymbol{s_2} \in \mathbb{R}^{N_x}$. We divide the support of indices such that $E = E^1 \cup E^2 \cup E^c$. Only $\boldsymbol{s_1}$ is active on $E^1$, only $\boldsymbol{s_2}$ is active
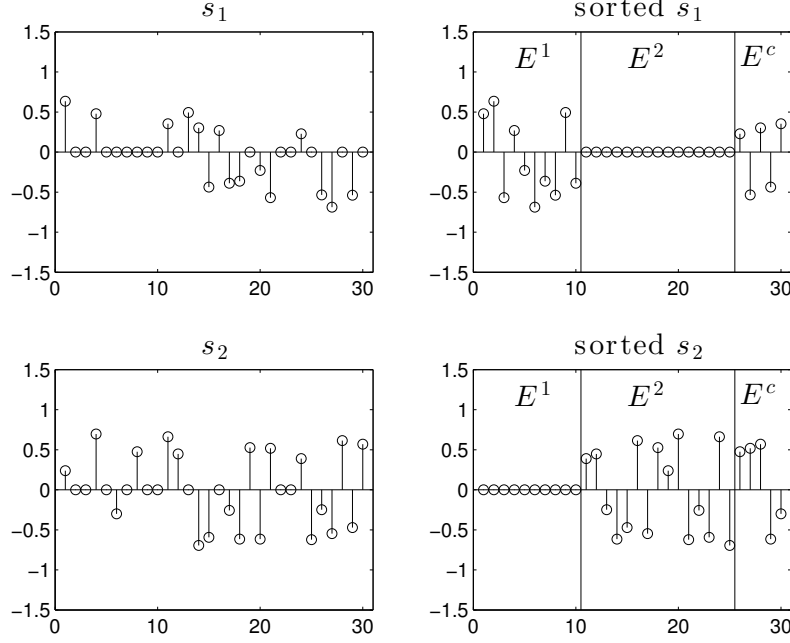
**Figure 8.2:** *Explanation of supports for a BSS problem with two sources. The coefficients of the sources are sorted without loss of generality.*

on $E^2$ and both sources are active on $E^c$ (see figure 8.2). The two sources $s_1$ and $s_2$ are respectively active on $N_1 + N_c$ and $N_2 + N_c$ indices. From the definition of inter-regressive processes, we have necessarily $E^* < N_c$.

**Theorem 14** (Sufficient condition)**.** *Let us consider two sources $s_1$ and $s_2$ for which $E^2$ and $E^c$ are of size $N_2$ and $N_c$. If*

$$N_2 > N_c, \tag{8.2}$$

*then $\ell_0$ pseudo-norm is a contrast for the extraction of $s_1$, no matter the size of $E^1$.*

*Proof.* From the condition we have $N_2 > N_c > E^*$. Also, if $N_c = 0$, then we have $E^* = 0$. Hence $N_c - N_2 \leq -1$ is a sufficient condition. It can be written $N_2 > N_c$.  ∎

If $\|s_1\|_0 < 0.5 \|s_2\|_0$, as proposed in Duarte et al. [2011], then

$$2(N_1 + N_c) \quad < \quad N_2 + N_c \tag{8.3}$$
$$2N_1 + N_c \quad < \quad N_2, \tag{8.4}$$

and so we have $N_2 > N_c$. The sufficiency is proved.

## 8.4   Separation $2 \times 2$

We want to proof our conjecture 1 for the case $N = M = 2$. We consider two sources $s_1, s_2 \in \mathbb{R}^{N_x}$. We divide the support of indices such that $E = E^1 \cup E^2 \cup E^c$. Only $s_1$ is active on $E^1$, only $s_2$ is active on $E^2$ and both sources are active on $E^c$ (see figure 8.2). We consider the following global mapping

$$\begin{bmatrix} \boldsymbol{u}_1 \\ \boldsymbol{u}_2 \end{bmatrix} = \begin{bmatrix} h_{11} & h_{12} \\ h_{21} & h_{22} \end{bmatrix} \begin{bmatrix} \boldsymbol{s}_1 \\ \boldsymbol{s}_2 \end{bmatrix}. \tag{8.5}$$

By definition, we denote the size of the support where the sources are active such that

$$\|\boldsymbol{s}_1\|_0 = N_1 \tag{8.6}$$

$$\|\boldsymbol{s}_2\|_0 = N_2 \tag{8.7}$$

$$\|\boldsymbol{s}_1 \odot \boldsymbol{s}_2\|_0 = N_{12}. \tag{8.8}$$

The value of the objective function is

$$\|\boldsymbol{u}_1 \odot \boldsymbol{u}_2\|_0 = \|\boldsymbol{h}_1 \boldsymbol{S} \odot \boldsymbol{h}_2 \boldsymbol{S}\|_0. \tag{8.9}$$

For most $\boldsymbol{H}$, the precedent equation is equal to $N$. We distinguish three cases leading to parasitic minima.

- An inter-regressive process of length $L$ on the support where both sources are active leads to a parasitic minimum. In the worse cases for which two inter-regressive processes exist (with arbitrarily $L_1 > L_2$), we have

$$\|\boldsymbol{u}_1 \odot \boldsymbol{u}_2\|_0 = N_x - L_1 - L_2, \tag{8.10}$$

and we end-up with the following condition

$$N_{12} < N_x - L_1 - L_2. \tag{8.11}$$

- One of the extracted vector can extract $\boldsymbol{s}_1$, while the other extracts a linear combination of the two sources. In the worse case with an existing inter-regressive process we have

$$\|\boldsymbol{u}_1 \odot \boldsymbol{u}_2\|_0 = N_x - (N_x - N_1) - L = N_1 - L, \tag{8.12}$$

leading to the following necessary and sufficient condition

$$N_{12} < N_1 - L \quad \Leftrightarrow \quad L < (N_1 - N_{12}). \tag{8.13}$$

The scalar $N_1 - N_{12}$ represents the size of the inactive support of $\boldsymbol{s}_2$.

- The symmetric case of the precedent one occurs when one of the extracted vector extracts $\boldsymbol{s}_2$, while the other extract a linear combination of the two sources. In the same way, we end up with

$$N_{12} < N_2 - L \quad \Leftrightarrow \quad L < (N_2 - N_{12}). \tag{8.14}$$

The scalar $N_2 - N_{12}$ represents the size of the inactive support of $\boldsymbol{s}_1$.

Without lost of generality, we can assume that $\boldsymbol{s}_1$ is sparser than $\boldsymbol{s}_2$, i.e. $N_1 < N_2$ and so $(N_1 - N_{12}) < (N_2 - N_{12})$ . This means that the inactive support of $\boldsymbol{s}_1$ is larger than the one of $\boldsymbol{s}_2$ and equation 8.14 is verified if equation 8.13 is true. Also

$$L_1 + L_2 < L_1 + L_1 < (N_1 - N_{12}) + (N_2 - N_{12}) = N_x - N_{12}, \tag{8.15}$$

and equation 8.11 is verified if equation 8.13 is true. Our conjecture is proved for the case $N = M = 2$.

# Bibliography

Abbad, B., Ursin, B., and Porsani, M. J. (2011). A fast, modified parabolic Radon transform. *Geophysics*, 76(1):V11–V24.

Abma, R., Kabir, N., Matson, K., Michell, S., Shaw, S., and McLain, B. (2005). Comparisons of adaptive subtraction methods for multiple attenuation. *The Leading Edge*, 24(3):277–280.

Adali, T., Schreier, P. J., and Scharf, L. L. (2011). Complex-valued signal processing: the proper way to deal with impropriety. *IEEE Transactions on signal processing*, 59(11):5101–5125.

Ahmed, I. (2007). 2D wavelet transform domain adaptive subtraction for enhancing 3D SRME. In *87th SEG Annual Meeting*, pages 2490–2494.

Aki, K. and Richards, P. G. (2002). *Quantitative Seismology (second Edition)*. University Science Books.

Amari, S.-I. and Cichocki, A. (1998). Adaptive blind signal processing – neural network approaches. *Proceedings of the IEEE*, 86(10):2026–2048.

Anemüller, J. (2007). Maximization of component disjointness: a criterion for blind source separation. In *7th International Conference on Independent Component Analysis and Signal Separation (ICA)*, pages 325–332.

Arlot, S. (2010). A survey of cross-validation procedures for model selection. *Statistics Surveys*, 4:40–79.

Asano, F., Ikeda, S., Ogawa, M., Asoh, H., and Kitawaki, N. (2003). Combined approach of array processing and independent component analysis for blind separation of acoustic signals. *IEEE Transactions on Speech and Audio Processing*, 11(3):204–215.

Baraniuk, R. G. (2007). Compressive sensing [lecture notes]. *IEEE Signal Processing Magazine*, 24(4):118–121.

Batany, Y.-M., Donno, D., Tomazeli Duarte, L., Chauris, H., and Travassos Romano, J. M. (2016a). A necessary and sufficient condition for the blind extraction of the sparsest source in convolutive mixtures. In *European Signal Processing Conference (EUSIPCO)*.

Batany, Y.-M., Tomazeli Duarte, L., Donno, D., Travassos Romano, J. M., and Chauris, H. (2016b). Adaptive multiple subtraction: Unification and comparison of matching filters based on the $\ell_q$-norm and statistical independence. *Geophysics*, 81(1):V43–V54.

Batany, Y.-M., Tomazeli Duarte, L., Donno, D., Travassos Romano, J. M., and Chauris, H. (2016c). Comparison of matching filters for adaptive multiple subtraction: Lq-norm versus statistical independence. In *78th EAGE Conference and Exhibition*.

Bell, A. J. and Sejnowski, T. J. (1995). An information-maximization approach to blind separation and blind deconvolution. *Neural Computation*, 7(6):1129–1159.

Benichoux, A., Sudhakar, P., Bimbot, F., and Gribonval, R. (2012). Some uniqueness results in sparse convolutive source separation. In *10th International Conference on Latent Variable Analysis and Signal Separation (LVA/ICA)*, pages 196–203.

Berkhout, A. J. (2016). Utilization of multiple scattering: the next big step forward in seismic imaging. *Geophysical Prospecting*, pages 1–40.

Berkhout, A. J. and Verschuur, D. J. (1997). Estimation of multiple scattering by iterative inversion, part I: Theoretical considerations. *Geophysics*, 62(5):1586–1595.

Berkhout, A. J. and Verschuur, D. J. (2006). Imaging of multiple reflections. *Geophysics*, 71(4):SI209–SI220.

Bishop, T. N., Bube, K. P., Cutler, R. T., Langan, R. T., Love, P. L., Resnick, J. R., Shuey, R. T., Spindler, D. A., and Wyld, H. W. (1985). Tomographic determination of velocity and depth in laterally varying media. *Geophysics*, 50(6):903–923.

Bobin, J., Rapin, J., Larue, A., and Starck, J.-L. (2015). Sparsity and adaptivity for the blind separation of partially correlated sources. *IEEE Transactions on Signal Processing*, 63(5):1199–1213.

Bofill, P. (2003). Underdetermined blind separation of delayed sound sources in the frequency domain. *Neurocomputing*, 55(3–4):627–641.

Bofill, P. and Zibulevsky, M. (2001). Underdetermined blind source separation using sparse representations. *Signal Processing*, 81(11):2353–2362.

Boyd, S. and Vandenberghe, L. (2004). *Convex Optimization*. Cambridge University Press.

Brossier, R. (2011). Two-dimensional frequency-domain visco-elastic full waveform inversion: parallel algorithms, optimization and performance. *Computers & Geosciences*, 37(4):444–455.

Brossier, R., Operto, S., and Virieux, J. (2009). Seismic imaging of complex onshore structures by 2D elastic frequency-domain full-waveform inversion. *Geophysics*, 74(6):WCC105–WCC118.

Bunks, C., Saleck, F. M., Zaleski, S., and Chavent, G. (1995). Multiscale seismic waveform inversion. *Geophysics*, 60(5):1457–1473.

Caiafa, C. F. (2012). On the conditions for valid objective functions in blind separation of independent and dependent sources. *Journal on advances in signal processing*, 2012(1):1–17.

Campigotto, P., Passerini, A., and Battiti, R. (2014). Active learning of Pareto front. *IEEE transactions on neural networks and learning systems*, 25(3):506–519.

Candès, E. and Demanet, L. (2005). The curvelet representation of wave propagators is optimally sparse. *Communications on pure and applied mathematics*, LVIII:1473–1528.

Candès, E., Demanet, L., and Donoho, D. (2006a). Fast discrete curvelet transforms. *Multiscale Modeling and Simulation*, 5(3):861–899.

Candès, E., Romberg, J. K., and Tao, T. (2006b). Stable signal recovery from incomplete and inaccurate measurements. *Communications on pure and applied mathematics*, 59(8):1207–1223.

Candès, E. and Wakin, M. (2008). An introduction to compressive sampling: A sensing/sampling paradigm that goes against the common knowledge in data acquisition. *IEEE Signal Processing Magazine*, 25(2):21–30.

Cao, X.-R. and Liu, R.-W. (1996). General approach to blind source separation. *IEEE Transactions on Signal Processing*, 44(3):562–571.

Cardoso, J.-F. (1997). Infomax and maximum likelihood for blind source separation. *Signal Processing Letters, IEEE*, 4(4):112–114.

Cardoso, J.-F. (1999). High-order contrasts for independent component analysis. *Neural Computation*, 11(1):157–192.

Cardoso, J. F., Snoussi, H., and Delabrouille, J. (2002). Blind source separation of noisy Gaussian stationary sources. Application to cosmic microwave background imaging. In *11th European Signal Processing Conference (EUSIPCO)*, pages 1–4.

Castella, M. and Moreau, E. (2009). Generalized identifiability conditions for blind convolutive MIMO separation. *IEEE Transactions on Signal Processing*, 57(7):2846–2852.

Castella, M. and Pesquet, J.-C. (2004). An iterative blind source separation method for convolutive mixtures of images. In *5th International conference on independent component analysis and blind signal separation*, pages 922–929.

Chauris, H. and Nguyen, T. (2008). Seismic demigration/migration in the curvelet domain. *Geophysics*, 73(2):S35–S46.

Cocher, E., Chauris, H., and Lameloise, C. (2015). Imaging with surface-related multiples in the subsurface-offset domain. In *77th EAGE Conference and Exhibition*.

Comon, P. (1994). Independent component analysis, a new concept? *Signal Processing*, 36(3):287–314.

Comon, P. and Jutten, C. (2010). *Handbook of Blind Source Separation*. Academic Press, Oxford.

Costagliola, S., Mazzucchelli, P., and Bienati, N. (2011). Hybrid norm adaptive subtraction for multiple removal. *73rd EAGE Conference and Exhibition*, page P231.

Cover, T. M. and Thomas, J. A. (2006). *Elements of Information Theory (2nd Edition)*. Wiley.

Curry, W. and Shan, G. (2010). Interpolation of near offsets using multiples and prediction-error filters. *Geophysics*, 75(6):WB153–WB16.

Das, S. and Suganthan, P. N. (2001). Differential evolution: a survey of the state-of-the-art. *IEEE Transactions on Evolutionary Computation*, 15(1):4–31.

Deb, K. (2001). *Multi-objective Optimization using Evolutionary Algorithms*. Wiley.

Delfosse, N. and Loubaton, P. (1995). Adaptive blind separation of independent sources: A deflation approach. *Signal Processing*, 45(1):59–83.

Deville, Y. (2014). Sparse component analysis: A general framework for linear and nonlinear blind source separation and mixture identification. In *Blind Source Separation*, pages 151–196. Springer.

Diaz, E. and Sava, P. (2013). Data-domain and image-domain wavefield tomography. *The Leading Edge*, 32(9):1064–1072.

Donno, D. (2011). Improving multiple removal using least-squares dip filters and independent component analysis. *Geophysics*, 76(5):V91–V104.

Donno, D., Chauris, H., and Noble, M. (2010). Curvelet-based multiple prediction. *Geophysics*, 75(6):WB255–WB263.

Donoho, D. L. and Elad, M. (2003). Optimally sparse representation in general (nonorthogonal) dictionaries via $\ell_1$ minimization. *Proceedings of the National Academy of Sciences*, 100(5):2197–2202.

Dragoset, B., Verschuur, D. J., Moore, I., and Bisley, R. (2010). A perspective on 3D surface-related multiple elimination. *Geophysics*, 75(5):A245–A261.

Duarte, L. T., Batany, Y.-M., and Romano, J. M. T. (2015). Blind source separation: principles of independent and sparse component analysis. In *Signals and Images: Advances and results in speech, estimation, compression, recognition, filtering and processing*, chapter 1. CRC Press.

Duarte, L. T., Suyama, R., Attux, R., Romano, J. M. T., and Jutten, C. (2011). Blind extraction of sparse components based on $\ell_0$-norm minimization. In *IEEE Statistical Signal Processing (SSP) Workshop*, pages 617–620, Nice, France.

Duong, N. Q. K., Vincent, E., and Gribonval, R. (2009). Spatial covariance models for under-determined reverberant audio source separation. In *IEEE workshop on applications of signal processing to audio and acoustics*, pages 129–132.

Durrani, T. S. and Bisset, D. (1984). The Radon transform and its properties. *Geophysics*, 49(8):1180–1187.

Eriksson, J. and Koivunen, V. (2004). Identifiability, separability, and uniqueness of linear ICA models. *IEEE Signal Processing Letters*, 11(7):601–604.

Eriksson, J. and Koivunen, V. (2006). Complex random vectors and ICA models: Identifiability, uniqueness, and separability. *IEEE Transactions on Information Theory*, 52(3):1017–1029.

Feng, F., Wang, D.-L., Zhu, H., and Cheng, H. (2013). Estimating primaries by sparse inversion of the 3D curvelet transform and the l1-norm constraint. *Applied Geophysics*, 10(2):201–209.

Fiori, S. (2001). On blind separation of complex-valued sources by extended hebbian learning. *IEEE signal processing letters*, 8(8):217–220.

Fomel, S. (2007a). Local seismic attributes. *Geophysics*, 72(3):A29–A33.

Fomel, S. (2007b). Shaping regularization in geophysical-estimation problems. *Geophysics*, 72(2):R29–R36.

Fomel, S. (2009). Adaptive multiple subtraction using regularized nonstationary regression. *Geophysics*, 74(1):V25–V33.

Georgiev, P., Theis, F., and Cichocki, A. (2005). Sparse component analysis and blind source separation of underdetermined mixtures. *IEEE Transactions on Neural Networks*, 16(4):992–996.

Georgiev, P., Theis, F., and Ralescu, A. (2007). Identifiability conditions and subspace clustering in sparse BSS. In *Independent Component Analysis and Signal Separation*, volume 4666 of *Lecture Notes in Computer Science*, pages 357–364. Springer.

Gribonval, R. and Lesage, S. (2006). A survey of Sparse Component Analysis for blind source separation: principles, perspectives, and new challenges. In *14th European Symposium on Artificial Neural Networks (ESANN)*, pages 323–330.

Gribonval, R. and Nielsen, M. (2003). Ieee transactions on information theory. *Neurocomputing*, 49(12):3320–3325.

Guitton, A. (2003). Multiple attenuation with multidimensional prediction-error filters. In *73rd SEG Annual Meeting*, pages 1945–1948.

Guitton, A. (2005). Multiple attenuation in complex geology with a pattern-based approach. *Geophysics*, 70(4):V97–V107.

Guitton, A. (2006). A pattern-based approach for multiple removal applied to a 3D Gulf of Mexico data set. *Geophysical Prospecting*, 54(2):135–152.

Guitton, A., Brown, M., Rickett, J., and Clapp, R. (2001). Multiple attenuation using a $t - x$ pattern-based subtraction method. In *71th SEG Annual Meeting*, pages 1305–1308.

Guitton, A. and Cambois, G. (1999). Multiple elimination using a pattern-recognition technique. *The Leading Edge*, 18(1):92–98.

Guitton, A. and Verschuur, D. J. (2004). Adaptive subtraction of multiples using the l1-norm. *Geophysical Prospecting*, 52(1):27–38.

Guo, J. (2003). Adaptive multiple subtraction with a pattern-based technique. In *73rd SEG Annual Meeting*, pages 1953–1956.

Hansen, C. (2008). Inverse problem theory. Technical report, Technical University of Denmark.

Haykin, S. (2013). *Adaptive filter theory.* Prentice Hall, Inc., 5th edition.

Herrmann, F. and Moghaddam, P. (2004). Curvelet-based non-linear adaptive subtraction with sparseness constraints. In *74rd SEG Annual Meeting*, pages 1977–1980.

Herrmann, F. J., Böniger, U., and Verschuur, D. J. (2007). Non-linear primary-multiple separation with directional curvelet frames. *Geophysical Journal International*, 170(2):781–799.

Herrmann, F. J., Wang, D., and Verschuur, D. J. (2008). Adaptive curvelet-domain primary-multiple separation. *Geophysics*, 73(3):A17–A21.

Hung, B. and Wang, M. (2012). Internal demultiple methodology without identifying the multiple generators. In *82nd SEG Annual Meeting*, pages 1–5.

Hurley, N. and Rickard, S. (2009). Comparing measures of sparsity. *IEEE Transactions on Information Theory*, 55(10):4723–4741.

Hyvärinen, A. (1999). Fast and robust fixed-point algorithms for independent component analysis. *IEEE transaction on Neural Networks*, 10(3):626–634.

Hyvärinen, A., Karhunen, J., and Oja, E. (2001). *Independent Component Analysis.* John Wiley & Sons.

Hyvärinen, A. and Oja, E. (2000). Independent component analysis: algorithms and applications. *Neural Networks*, 13(4–5):411–430.

Ikelle, L. (2010). *Coding and decoding: seismic data.* Elsevier Science.

Jakubowicz, H. (1998). Wave equation prediction and removal of interbed multiples. In *68th SEG Annual Meeting*, pages 1527–1530.

Jourjine, A., Rickard, S., and Yilmaz, O. (2000). Blind separation of disjoint orthogonal signals: demixing $n$ sources from 2 mixtures. In *IEEE International Conference on Acoustics, Speech, and Signal Processing (ICASSP)*, volume 5, pages 2985–2988.

Jumah, B. and Herrmann, F. J. (2014). Dimensionality-reduced estimation of primaries by sparse inversion. *Geophysical Prospecting*, 62(5):972–993.

Kaplan, S. T. and Innanen, K. A. (2008). Adaptive separation of free-surface multiples through independent component analysis. *Geophysics*, 73(3):V29–V36.

Kearey, P., Brooks, M., and Hill, I. (2002). *An Introduction to Geophysical Exploration*. Wiley-Blackwell.

Kim, I. Y. and Weck, O. L. (2005). Adaptive weighted-sum method for bi-objective optimization: Pareto front generation. *Structural and multidisciplinary optimization*, 29(2):149–158.

Lambert, R. (1999). Difficulty measures and figures of merit for source separation. In *1st international symposium on independent component analysis and blind signal separation (ICA)*, pages 133–138.

Law, K., Fossum, R., and Do, M. (2009). Generic invertibility of multidimensional FIR filter banks and MIMO systems. *IEEE Transactions on Signal Processing*, 57(11):4282–4291.

Leon-Garcia, A. (2008). *Probability, Statistics, and Random Processes for Electrical Engineering (3rd Edition)*. Pearson.

Levin, S. A. (2008). Delft $\equiv$ inverse scattering surface-related multiple attenuation in three lines. In *78th SEG Annual Meeting*, pages 2512–2516.

Li, J. and Symes, W. W. (2007). Interval velocity estimation via NMO-based differential semblance. *Geophysics*, 72(6):U75–U88.

Li, Y. and Amari, S.-I. (2010). Two conditions for equivalence of 0-norm solution and 1-norm solution in sparse representation. *IEEE Transactions on Neural Networks*, 21(7):1189–1196.

Li, Y., Amari, S.-I., Cichocki, A., and Guan, C. (2006). Probability estimation for recoverability analysis of blind source separation based on sparse representation. *IEEE Transactions on Information Theory*, 52(7):3139–3152.

Li, Y., Cichocki, A., and Amari, S.-I. (2003a). Sparse component analysis for blind source separation with less sensors than sources. In *4th international symposium on independent component analysis and blind signal separation (ICA)*.

Li, Y., Cichocki, A., Amari, S.-I., Shishkin, S., Cao, J., and Gu, F. (2003b). Sparse representation and its applications in blind source separation. In *Seventeenth Annual Conference on Neural Information Processing Systems (NIPS)*.

Li, Z.-X., Li, Z.-C., and Lu, W.-K. (2016). Multichannel predictive deconvolution based on the fast iterative shrinkage-thresholding algorithm. *Geophysics*, 81(1):V17–V30.

Li, Z.-X. and Lu, W.-K. (2013). Adaptive multiple subtraction based on 3D blind separation of convolved mixtures. *Geophysics*, 78(6):V251–V266.

Li, Z.-X. and Lu, W.-K. (2014). Demultiple strategy combining Radon filtering and Radon domain adaptive multiple subtraction. *Journal of Applied Geophysics*, 103(4):1–11.

Lin, T. T. Y. and Herrmann, F. J. (2013). Robust estimation of primaries by sparse inversion via one-norm minimization. *Geophysics*, 78(3):R133–R150.

Lindenbaum, O., Yeredor, A., Vitek, R., and Mishali, M. (2015). Blind separation of orthogonal mixtures of spatially-sparse sources with unknown sparsity levels and with temporal blocks. *Journal of Signal Processing Systems*, 79(2):167–178.

Linsker, R. (1989). An application of the principle of maximum information preservation to linear systems. In Touretzky, D. S., editor, *Advances in Neural Information Processing Systems 1*, pages 186–194. Morgan Kaufmann Publishers Inc.

Liu, J. and Lu, W. (2016). Adaptive multiple subtraction based on multiband pattern coding. *Geophysics*, 81(1):V69–V78.

Liu, K.-H. and Dragoset, W. H. (2013). Blind-source separation of seismic signals based on information maximization. *Geophysics*, 78(4):V119–V130.

Liu, K.-H. and Kostov, C. (2015). Multimodel adaptive subtraction with regularized parameter selection via generalized information criterion. *Geophysics*, 80(3):V33–V45.

Lu, W. (2006). Adaptive multiple subtraction using independent component analysis. *Geophysics*, 71(5):S179–S184.

Lu, W.-K. and Liu, L. (2009). Adaptive multiple subtraction based on constrained independent component analysis. *Geophysics*, 74(1):V1–V7.

Luo, Y., Kelamis, P. G., and Wang, Y. (2003). Simultaneous inversion of multiples and primaries: Inversion versus subtraction. *The Leading Edge*, 22(9):814–891.

Maeland, E. (2003). Disruption of seismic images by the parabolic Radon transform. *Geophysics*, 68(3):1060–1064.

Mallat, S. (1998). *A wavelet tour of signal processing – The sparse way.* Academic Press.

McGill, W. J. (1954). Multivariate information transmission. *Psychometrika*, 19(2):97–116.

Mei, T. and Mertins, A. (2008). Convolutive blind source separation based on disjointness maximization of subband signals. *IEEE Signal Precessing Letters*, 15:725–728.

Meles, G. A., L oer, K., Ravasi, M., Curtis, A., and da Costa Filho, C. A. (2015). Internal multiple prediction and removal using marchenko autofocusing and seismic interferometry. *Geophysics*, 80(1):A7–A11.

Métivier, L., Brossier, R., Mérigot, Oudet, E., and Virieux, J. (2016). Measuring the misfit between seismograms using an optimal transport distance: application to full waveform inversion. *Geophysical Journal International*, 205(1):345–377.

Mishali, M. and Eldar, Y. (2009). Sparse source separation from orthogonal mixtures. In *IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pages 3145–3148.

Mitra, S. K. and Kaiser, J. F. (1993). *Handbook for Digital Signal Processing.* J. Wiley & Sons.

Mohimani, G., Babaie-Zadeh, M., and Jutten, C. (2008). Complex-valued sparse representation based on smoothed $\ell_0$-norm. In *IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pages 3881–3884.

Mourad, N. and Reilly, J. P. (2010). Blind extraction of sparse sources. In *IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pages 2666–2669.

Mulder, W. A. and ten Kroode, A. P. E. (2002). Automatic velocity analysis by differential semblance optimization. *Geophysics*, 67(4):1184–1191.

Muller, M. E. (1959). A note on a method for generating points uniformly on *n*-dimensional spheres. *Comm. Assoc. Comput.*, 2:19–20.

Murillo-Fuentes, J. J. and Boloix-Tortosa, R. (2010). Strict separability and identifiability of a class of ICA models. *IEEE Signal Processing Letters*, 17(3):285–288.

Nadalin, E. Z. (2011). *Contribuições ao problema de separação cega de fontes, com ênfase no estudo de sinais esparsos*. PhD thesis, Faculdade de engenharia elétrica e de computação, Universidade de Campinas.

Naini, F., Mohimani, G., Babaie-Zadeh, M., and Jutten, C. (2007). Estimating the mixing matrix in sparse component analysis (SCA) based on multidimensional subspace clustering. In *IEEE International Conference on Telecommunications and Malaysia International Conference on Communications (ICT-MICC)*, pages 670–675.

Neelamani, R., Baumstein, A., and Ross, W. (2010). Adaptive subtraction using complex-valued curvelet transforms. *Geophysics*, 75(4):V51–V60.

Nguyen, T. T. and Chauris, H. (2010). Uniform discrete curvelet transform. *IEEE transactions on signal processing*, 58(7):3618–3634.

Nose-Filho, K. (2015). *Desconvolução e separação cega de sinais esparsos e aplicações em sísmica de reflexão*. PhD thesis, Faculdade de engenharia elétrica e de computação, Universidade de Campinas.

Ozerov, A. and Fevotte, C. (2010). Multichannel nonnegative matrix factorization in convolutive mixtures for audio source separation. *IEEE Transactions on Audio, Speech, and Language Processing*, 18(3):550–563.

Pang, T., Lu, W., and Ma, Y. (2009). Adaptive multiple subtraction using a constrained l1-norm method with lateral continuity. *Applied Geophysics*, 6(3):241–247.

Papadias, C. B. (2000). Globally convergent blind source separation based on a multiuser kurtosis maximization criterion. *IEEE Transactions on Signal Processing*, 48(12):3508–3519.

Peacock, K. L. and Treitel, S. (1969). Predictive deconvolution: theory and practice. *Geophysics*, 34(2):155–169.

Pedersen, M. S., Larsen, J., Kjems, U., and Parra, L. C. (2008). Convolutive blind source separation methods. In *Springer Handbook of Speech Processing*, pages 1065–1094. Springer.

Pham, M.-Q. (2015). *Seismic wave field restoration using spare representations and quantitative analysis*. Theses, Université Paris-Est.

Pham, M. Q., Chaux, C., Duval, L., and Pesquet, J.-C. (2015). Sparse adaptive template matching and filtering for 2D seismic images with dual-tree wavelets and proximal methods. In *International Conference on Image Processing (ICIP)*, pages 2339–2343, Québec City, Canada.

Pham, M. Q., Duval, L., Chaux, C., and Pesquet, J.-C. (2014). A primal-dual proximal algorithm for sparse template-based adaptive filtering: Application to seismic multiple removal. *IEEE Transactions on Signal Processing*, 62(16):4256–4269.

Pica, A., Poulain, G., David, B., Magesan, M., Baldock, S., Weisser, T., Hugonnet, P., and Herrmann, P. (2005). 3D surface-related multiple modeling, principles and results. In *75th SEG Annual Meeting*, pages 2080–2083.

Pope, K. and Bogner, R. (1996). Blind Signal Separation II. Linear, Convolutive Combinations. *Digital Signal Processing*, 6(1):17 – 28.

Porsani, M. J. and Ursin, B. (2007). Direct multichannel predictive deconvolution. *Geophysics*, 72(2):H11–H27.

Pratt, R., Song, Z., Williamson, P., and Warner, M. (1996). Two-dimensional velocity models from wide-angle seismic data by wavefield inversion. *Geophysical Journal International*, 124(2):323–340.

Prosser, R. T. (1966). A multidimensional sampling theorem. *Journal of mathematical analysis and applications*, 16:574–584.

Rajagopal, R. and Potter, L. C. (2003). Multivariate MIMO FIR inverses. *IEEE Transactions on image processing*, 12(4):458–465.

Rajendran, S. M. and Rajakumar, R. (2014). Curvelet transform for Boehmians. *Arab journal of mathematical sciences*, 20(2):264–279.

Reshef, M., Arad, S., and Landa, E. (2006). 3D prediction of surface-related and interbed multiples. *Geophysics*, 71(1):V1–V6.

Rickard, S. and Yilmaz, O. (2002). On the approximate W-disjoint orthogonality of speech. In *International Conference on Acoustics, Speech, and Signal Processing (ICASSP)*, volume 1, pages 529–532.

Rickett, J., Guitton, A., and Gratwick, D. (2001). Adaptive multiple subtraction with non-stationary helical shaping filters. *63rd EAGE Conference and Exhibition*, page P167.

Robein, E. (2010). *Seismic imaging: a review of the techniques, their principles, merits and limitations*. EAGE Publications BV.

Saab, R., Wang, D., Yilmaz, O., and Herrmann, F. J. (2007a). Curvelet-based primary-multiple separation from a bayesian perspective. In *87th SEG Annual Meeting*.

Saab, R., Yilmaz, O., McKeown, M., and Abugharbieh, R. (2007b). Underdetermined anechoic blind source separation via $\ell^q$-basis-pursuit with $q < 1$. *IEEE Transactions on Signal Processing*, 55(8):4004–4017.

Sacchi, M. D. and Ulrych, T. J. (1995). Improving resolution of radon operators using a model re-weighted least squares procedure. *Journal of Seismic Exploration*, 4:315–328.

Sava, P. and Biondi, B. (2004). Wave-equation migration velocity analysis. I. Theory. *Geophysical prospecting*, 52(6):593–606.

Savels, T., de Vos, K., and de Maag, J. W. (2011). Surface-multiple attenuation through sparse inversion: results for complex synthetics and real data. *First Break*, 29(1):55–64.

Sawada, H., Araki, S., and Makino, S. (2010). Underdetermined convolutive blind source separation via frequency bin-wise clustering and permutation alignment. *IEEE Transactions on Audio, Speech, and Language Processing*, 18(1):516–527.

Scott, R. (2007). The duet blind source separation algorithm. In *Blind Speech Separation*, pages 217–241. Springer.

Shannon, C. E. (1954). A mathematical theory of communication. *The Bell System Technical Journal*, 27(3,4):379–423,623–666.

Sheriff, R. E. and Geldart, L. P. (1995). *Exploration Seismology*. Cambridge University Press.

Snieder, R., Wapenaar, K., and Larner, K. (2006). Spurious multiples in seismic interferometry of primaries. *Geophysics*, 71(4):SI111–SI124.

Spitz, S. (1999). Pattern recognition, spatial predictability, and subtraction of multiple events. *The Leading Edge*, 18(1):55–58.

Spitz, S. (2000). Model-based subtraction of multiple events in the frequency-space domain. In *70th SEG Annual Meeting*, pages 1969–1972.

Storn, R. and Price, K. (1997). Differential evolution – a simple and efficient heuristic for global optimization over continuous spaces. *Journal of Global Optimization*, 11(4):341–359.

Sun, Y. and Xin, J. (2011). Underdetermined sparse blind source separation of nonnegative and partially overlapped data. *SIAM Journal on Scientific Computing*, 33(4):2063–2094.

T. Romano, J. M., Attux, R., C. Cavalcante, C., and Suyama, R. (2010). *Unsupervised Signal Processing: Channel Equalization and Source Separation*. J. Wiley & Sons.

Tarantola, A. (1987). *Inverse Problem Theory*. Elsevier Scientific Publishing Compnay.

Tu, N. and Herrmann, F. J. (2015). Fast imaging with surface-related multiples by sparse inversion. *Geophysical Journal International*, 201(1):304–317.

van der Neut, J., Wapenaar, K., Thorbecke, J., Slob, E., and Vasconcelos, I. (2015). An illustration of adaptive marchenko imaging. *The Leading Edge*, 34(7):818–822.

van Groenestijn, G. J. A. (2010). Estimation of primaries and multiples by sparse inversion. Master's thesis, Technische Universiteit Delft.

van Groenestijn, G. J. A. and Verschuur, D. J. (2009). Estimating primaries by sparse inversion and application to near-offset data reconstruction. *Geophysics*, 74(3):A23–A28.

van Groenestijn, G. J. A. and Verschuur, D. J. (2010). Estimation of primaries by sparse inversion from passive seismic data. *Geophysics*, 75(4):SA61–SA69.

Velis, D. R. (2003). Estimating the distribution of primary reflection coefficients. *Geophysics*, 68(4):1417–1422.

Ventosa, S., Le Roy, S., Huard, I., Pica, A., Rabeson, H., Ricarte, P., and Duval, L. (2012). Adaptive multiple subtraction with wavelet-based complex unary Wiener filters. *Geophysics*, 77(6):V183–V192.

Verschuur, D. J. (2013a). Estimation of primaries by sparse inversion including the ghost. In *75th EAGE Conference and Exhibition*.

Verschuur, D. J. (2013b). *Seismic Multiple Removal Techniques: Past, present and future*. EAGE Publications BV, revised edition.

Verschuur, D. J. and Berkhout, A. J. (1997). Estimation of multiple scattering by iterative inversion, part II: Practical aspects and examples. *Geophysics*, 62(5):1596–1611.

Verschuur, D. J., Berkhout, A. J., and Wapenaar, C. P. A. (1992). Adaptive surface-related multiple elimination. *Geophysics*, 57(9):1166–1177.

Vigario, R. and Oja, E. (2008). BSS and ICA in neuroinformatics: From current practices to open challenges. *IEEE Reviews in Biomedical Engineerng*, 1:50–61.

Vincent, E., Araki, S., Theis, F., Nolte, G., Bofill, P., Sawada, H., Ozerov, A., Gowreesunker, V., Lutter, D., and Duong, N. Q. K. (2012). The signal separation evaluation campaign (2007 – 2010): Achievements and remaining challanges. *Signal Processing*, 92(8):1928–1936.

Vincent, E., Gribonval, R., and Févotte, C. (2006). Performance measurement in blind audio source separation. *IEEE Transactions on audio, speech, and language processing*, 14(4):1462 – 1469.

Virieux, J. and Operto, S. (2009). An overview of full-waveform inversion in exploration geophysics. *Geophysics*, 74(6):WCC127–WCC152.

Wang, D., Saab, R., Yilmaz, O., and Herrmann, F. J. (2007). Recent results in curvelet-based primary-multiple separation: application to real data. In *87th SEG Annual Meeting*, pages 2500–2504.

Wang, G., Ma, R., Meng, Q., and Liu, W. (2015). Maximum non-gaussianity estimation revisit: Uniqueness analysis from the perspective of constrained cost function optimization. In *11th International Conference on Natural Computation (ICNC)*, pages 94–101.

Wang, Y. (2003). Multiple subtraction using an expanded multichannel matching filter. *Geophysics*, 68(1):346–354.

Wapenaar, K., Thorbecke, J., van der Neut, J., Broggini, F., Slob, E., and Snieder, R. (2014). Marchenko imaging. *Geophysics*, 79(3):WA39–WA57.

Watanabe, S. (1960). Information theoretical analysis of multivariate correlation. *IBM Journal of Research and Development*, 4:66–82.

Weglein, A., Gasparotto, F., Carvalho, P., and Stolt, R. (1997). An inverse-scattering series method for attenuating multiples in seismic reflection data. *Geophysics*, 62(6):1975–1989.

Weglein, A. B. (2015). Primaries – the only events that can be migrated and for which migration has meaning. *The Leading Edge*, 34(7):808–813.

Wiggins, J. W. (1988). Attenuation of complex water-bottom multiples by wave equation-based prediction and subtraction. *Geophysics*, 53(12):1527–1539.

Wiggins, J. W. (1999). Multiple attenuation by explicit wave extrapolation to an interpreted horizon. *The Leading Edge*, 18(1):46–54.

Wong, M. (2012). Introduction to multiple attenuation methods. University Lecture.

Wong, M., Biondi, B. L., and Ronen, S. (2015). Imaging with primaries and free-surface multiples by joint least-squares reverse time migration. *Geophysics*, 80(6):S223–S235.

Wu, X. and Hung, B. (2015). High-fidelity adaptive curvelet domain primary-multiple separation. *First Break*, 33(1):53–59.

Xie, X.-B., Jin, S., and Wu, R.-S. (2006). Wave-equation-based seismic illumination analysis. *Geophysics*, 71(5):S169–S177.

Xu, L. (2005). *One-Bit-Matching ICA Theorem, Convex-Concave Programming, and Combinatorial Optimization*, pages 5–20. Springer.

Yang, L., Zhang, H., and Tong, X. (2014). Blind signal separation of complex-valued sources based on gaussian mixture model for time-varying environment. In *6th International Conference on Wireless Communications and Signal Processing (WCSP)*, pages 94–101.

Yang, M., Li, C., Cai, Z., and Guan, J. (2015). Differential evolution with auto-enhanced population diversity. *IEEE Transactions on Cypernetics*, 45(2):302–315.

Yilmaz, O. (2001). *Seismic Data Analysis: Processing, Inversion, and Interpretation of Seismic data*. Society of Exploration Geophysicists.

Zhou, A., Qu, B.-Y., Li, H., Zhao, S.-Z., Suganthan, P. N., and Zhang, Q. (2011). Multiobjective evolutionary algorithms: a survey of the state of the art. *Swarm and Evolutionary Computation*, 1(1):32–49.

## Résumé

La séparation de signaux corrélés à partir de leurs combinaisons linéaires est une tâche difficile et possède plusieurs applications en traitement du signal. Nous étudions deux problèmes, à savoir la séparation aveugle de sources parcimonieuses et le filtrage adaptatif des réflexions multiples en acquisition sismique. Un intérêt particulier est porté sur les mélanges convolutifs : pour ces deux problèmes, des filtres à réponses impulsionnelles finies peuvent être estimés afin de récupérer les signaux désirés.

Pour les modèles de mélange instantanés et convolutifs, nous donnons les conditions nécessaires et suffisantes pour l'extraction et la séparation exactes de sources parcimonieuses en utilisant la pseudo-norme $\ell_0$ comme une fonction de contraste. Des équivalences entre l'analyse en composantes parcimonieuses et l'analyse en composantes disjointes sont examinées.

Pour la soustraction adaptative des réflexions sismiques, nous discutons les limites des méthodes basées sur l'analyse en composantes indépendantes et nous soulignons l'équivalence avec les méthodes basées sur les normes $\ell_p$. Nous examinons de quelle manière les paramètres de régularisation peuvent être plus décisifs pour l'estimation des primaires. Enfin, nous proposons une amélioration de la robustesse de la soustraction adaptative en estimant les filtres adaptatifs directement dans le domaine des curvelets. Les coûts en calcul et en mémoire peuvent être atténués par l'utilisation de la transformée en curvelet discrète et uniforme.

## Abstract

The recovery of correlated signals from their linear combinations is a challenging task and has many applications in signal processing. We focus on two problems that are the blind separation of sparse sources and the adaptive subtraction of multiple events in seismic processing. A special focus is put on convolutive mixtures: for both problems, finite impulse response filters can indeed be estimated for the recovery of the desired signals.

For instantaneous and convolutive mixing models, we address the necessary and sufficient conditions for the exact extraction and separation of sparse sources by using the $\ell_0$ pseudo-norm as a contrast function. Equivalences between sparse component analysis and disjoint component analysis are investigated.

For adaptive multiple subtraction, we discuss the limits of methods based on independent component analysis and we highlight equivalence with $\ell_p$-norm-based methods. We investigate how other regularization parameters may have more influence on the estimation of the desired primaries. Finally, we propose to improve the robustness of adaptive subtraction by estimating the extracting convolutive filters directly in the curvelet domain. Computation and memory costs are limited by using the uniform discrete curvelet transform.

## Mots Clés

traitement sismique • réflexions multiples • filtrage adaptatif • séparation aveugle de sources • analyse en composantes parcimonieuses • transformée en curvelet

## Keywords

seismic processing • seismic multiples • adaptive filtering • blind source separation • sparse component analysis • curvelet transform