



HAL
open science

Optimisation algorithms in non-standard Banach spaces for inverse problems in imaging

Marta Lazzaretti

► **To cite this version:**

Marta Lazzaretti. Optimisation algorithms in non-standard Banach spaces for inverse problems in imaging. Image Processing [eess.IV]. Université Côte d'Azur; Università degli studi (Gênes, Italie), 2024. English. NNT: 2024COAZ4009 . tel-04558943

HAL Id: tel-04558943

<https://theses.hal.science/tel-04558943>

Submitted on 25 Apr 2024

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

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

PHD THESIS

Optimisation algorithms in non-standard Banach spaces for inverse problems in imaging

Marta LAZZARETTI

Dipartimento di Matematica, Università di Genova

&

Laboratoire I3S, CNRS, Centre Inria d'Université Côte d'Azur

**Thesis presented to obtain the degree
of *docteur* in** Automatique, Traitement du Signal et des Images **of** Université Côte d'Azur
and of *dottore di ricerca* in Matematica e Applicazioni **of** Università di Genova

Supervisors : Luca CALATRONI, CR CNRS, Université Côte d'Azur, France
Claudio ESTATICO, Associate Professor, Università di Genova, Italy

Defense date : 05/04/24

In front of the jury :

Laure BLANC-FÉRAUD, DR CNRS, Université Côte d'Azur, France
Serena MORIGI, Full Professor, Università di Bologna, Italy
Ronny RAMLAU, Full Professor, JKU Linz, Austria
Nelly PUSTELNIK, DR CNRS, ENS Lyon, France
Silvia VILLA, Associate Professor, Università di Genova, Italy
Luca CALATRONI, CR CNRS, Université Côte d'Azur, France
Claudio ESTATICO, Associate Professor, Università di Genova, Italy

**OPTIMISATION ALGORITHMS IN NON-STANDARD BANACH SPACES
FOR INVERSE PROBLEMS IN IMAGING**

*Algorithmes d'optimisation dans des espaces de Banach non standard
pour problèmes inverses en imagerie*

Marta LAZZARETTI



Jury :

Referees

Serena MORIGI, Full Professor, Università di Bologna, Italy
Ronny RAMLAU, Full Professor, JKU Linz, Austria

Examinators

Laure BLANC-FÉRAUD, DR CNRS, Université Côte d'Azur, France
Nelly PUSTELNIK, DR CNRS, ENS Lyon, France
Silvia VILLA, Associate Professor, Università di Genova, Italy

Supervisors

Luca CALATRONI, CR CNRS, Université Côte d'Azur, France
Claudio ESTATICO, Associate Professor, Università di Genova, Italy

THÈSE DE DOCTORAT

Algorithmes d'optimisation dans des espaces de Banach non standard pour problèmes inverses en imagerie

Marta LAZZARETTI

Dipartimento di Matematica, Università di Genova

&

Laboratoire I3S, CNRS, Centre Inria d'Université Côte d'Azur (Équipe Morpheme)

Présentée en vue de l'obtention du grade de docteur en Automatique, Traitement du Signal et des Images **d'Université Côte d'Azur et de dottore di ricerca en** Matematica e Applicazioni **d'Università di Genova**

Dirigée par : Luca CALATRONI, CR CNRS, Université Côte d'Azur, France
Claudio ESTATICO, Associate Professor, Università di Genova, Italy

Soutenue le : 05/04/24

Devant le jury, composé de :

Serena MORIGI, Full Professor, Università di Bologna, Italy

Ronny RAMLAU, Full Professor, JKU Linz, Austria

Laure BLANC-FÉRAUD, DR CNRS, Université Côte d'Azur, France

Nelly PUSTELNIK, DR CNRS, ENS Lyon, France

Silvia VILLA, Associate Professor, Università di Genova, Italy

Luca CALATRONI, CR CNRS, Université Côte d'Azur, France

Claudio ESTATICO, Associate Professor, Università di Genova, Italy

**OPTIMISATION ALGORITHMS IN NON-STANDARD BANACH SPACES
FOR INVERSE PROBLEMS IN IMAGING**

*Algorithmes d'optimisation dans des espaces de Banach non standard
pour problèmes inverses en imagerie*

Marta LAZZARETTI



Jury :

Rapporteurs

Serena MORIGI, Full Professor, Università di Bologna, Italy
Ronny RAMLAU, Full Professor, JKU Linz, Austria

Examineurs

Laure BLANC-FÉRAUD, DR CNRS, Université Côte d'Azur, France
Nelly PUSTELNIK, DR CNRS, ENS Lyon, France
Silvia VILLA, Associate Professor, Università di Genova, Italy

Directeurs de thèse

Luca CALATRONI, CR CNRS, Université Côte d'Azur, France
Claudio ESTATICO, Associate Professor, Università di Genova, Italy

Acknowledgments

It seems unreal but my PhD journey has come to an end after 3 and a half incredible years. I shared this path and experience with many amazing people and I would like to express my gratitude to them for making it special and unique.

First of all, let me start with institutional, but no less genuine, acknowledgments. Thanks to the referees of my thesis Prof. Serena Morigi and Prof. Ronny Ramlau for their invaluable time, effort, and insightful comments in reviewing my thesis: it was an honor for me to receive your positive feedback and your suggestions have contributed to the improvement of my work. I would also like to thank the other members of my PhD defense jury, DR Laure Blanc-Féraud, DR Nelly Pustelnik and Ass. Prof. Silvia Villa for your availability: having such a high-quality committee is a privilege and a responsibility.

I sincerely feel no words would be enough to thank my PhD supervisors, Luca Calatroni and Claudio Estatico for their constant support, guidance and encouragements when I most needed it, for their brilliant intuitions, precious help and advice that guided me during my doctoral journey. I am very grateful to have shared the last for 4 years with you. Thank you, Luca, for encouraging me to take part in schools and conferences and research visits abroad: I had the possibility to meet excellent scientists and people from all over the world and I feel greatly enriched and lucky for this. Thank you, Claudio, for having your door always for me and for understanding how I feel with a glance and your kind words.

During my PhD studies, I had the privilege of spending almost one year at the I3S Lab in Sophia Antipolis: thanks to all the people who welcomed me nicely there and made me feel like at home. I would like to thank especially Laure Blanc-Féraud for the thoughtful discussions we had about off-the-grid methods and Bastien Laville for sharing his code: it has been immensely useful in my work.

I am also grateful to Prof. Carola-Bibiane Schönlieb, Dr. Yury Korolev, and the CIA group for hosting me during my research visit in Cambridge between May and June 2022. A special thank goes also to Dr. Leila Muresan (CAIC, Cambridge) and Dr. Jerome Boulanger (MRC-LMB, Cambridge), for welcoming me in their imaging laboratories, for showing me for the very first times microscopes and the acquisition

process I had worked on before only from an abstract point of view and for providing challenging real microscopy data. My PhD thesis would not be the same without your help.

I would like also to thank Dr. Zeljko Kereta (UCL, London): collaborating with you has been a pleasure and you taught me a lot.

Now I will thank all the amazing people I met during my PhD. I have to start from my friends at DIMA: Alessandro, my office mate (together with Antoine of course, if you know you know), Silvia S. for giving me the best nickname ever and Silvia B. for supporting me from even before my PhD. Special thanks also to Issa, for organising hikes (and never getting lost), Luca, Chiara, Danilo, Betta, Laura, and many more. All of you has made the last 3 years unforgettable. Thanks for all the fun moments, coffee breaks, dinners, hikes, trips and for all the amazing memories we now have together.

I also have great memories at the I3S lab, and for this I have to thank you, Vasilina. My time in France would not have been the same without you. Thanks for welcoming me in your friends group, for your constant support, for being the best company during so many conferences, winter schools and during our time in Cambridge. Thank you also to Gabriele, we met at I3S before starting our PhDs and we shared our journey from afar but your support and friendship have been precious to me. I am also grateful to Nadia for helping me navigating the complex french bureaucracy and for forcing me to speak in french with her. If I can speak a bit of french it's also your merit.

Last but not least, I want to thank my family, my parents, my sister Laura and my boyfriend Gabriele. You have supported me immensely in this journey, being there for me when I was abroad, helping me in all the ways you could.

Thank you very much. This works is also yours.

Abstract

Optimisation algorithms in non-standard Banach spaces for inverse problems in imaging

This thesis focuses on the modelling, the theoretical analysis and the numerical implementation of advanced optimisation algorithms for imaging inverse problems (e.g., image reconstruction in computed tomography, image deconvolution in microscopy imaging) in non-standard Banach spaces.

It is divided into two parts: in the former, the setting of Lebesgue spaces with a variable exponent map $L^{p(\cdot)}$ is considered to improve adaptivity of the solution with respect to standard Hilbert reconstructions; in the latter a modelling in the space of Radon measures is used to avoid the biases observed in sparse regularisation methods due to discretisation.

In more detail, the first part explores both smooth and non-smooth optimisation algorithms in reflexive $L^{p(\cdot)}$ spaces, which are Banach spaces endowed with the so-called Luxemburg norm. As a first result, we provide an expression of the duality maps in those spaces, which are an essential ingredient for the design of effective iterative algorithms. To overcome the non-separability of the underlying norm and the consequent heavy computation times, we then study the class of modular functionals which directly extend the (non-homogeneous) p -power of L^p -norms to the general $L^{p(\cdot)}$. In terms of the modular functions, we formulate handy analogues of duality maps, which are amenable for both smooth and non-smooth optimisation algorithms due to their separability. We thus study modular-based gradient descent (both in deterministic and in a stochastic setting) and modular-based proximal gradient algorithms in $L^{p(\cdot)}$, and prove their convergence in function values. The spatial flexibility of such spaces proves to be particularly advantageous in addressing sparsity, edge-preserving and heterogeneous signal/noise statistics, while remaining efficient and stable from an optimisation perspective. We numerically validate this extensively on 1D/2D exemplar inverse problems (deconvolution, mixed denoising, CT reconstruction).

The second part of the thesis focuses on off-the-grid Poisson inverse problems formulated within the space of Radon measures. Our contribution consists in the modelling of a variational model which couples a Kullback-Leibler data term with

the Total Variation regularisation of the desired measure (that is, a weighted sum of Diracs) together with a non-negativity constraint. A detailed study of the optimality conditions and of the corresponding dual problem is carried out and an improved version of the Sliding Frank-Wolfe algorithm is used for computing the numerical solution efficiently. To mitigate the dependence of the results on the choice of the regularisation parameter, an homotopy strategy is proposed for its automatic tuning, where, at each algorithmic iteration checks whether an informed stopping criterion defined in terms of the noise level is verified and updates the regularisation parameter accordingly. Several numerical experiments are reported on both simulated 2D and real 3D fluorescence microscopy data.

Keywords: non-smooth optimisation, imaging inverse problems, regularisation in Banach spaces, sparse regularisation, fluorescence microscopy.

Résumé

Algorithmes d'optimisation dans des espaces de Banach non standard pour problèmes inverses en imagerie

Cette thèse porte sur la modélisation, l'analyse théorique et l'implémentation numérique d'algorithmes d'optimisation pour la résolution de problèmes inverses d'imagerie (par exemple, la reconstruction d'images en tomographie et la déconvolution d'images en microscopie) dans des espaces de Banach non standard.

Elle est divisée en deux parties: dans la première, nous considérons le cadre des espaces de Lebesgue à exposant variable $L^{p(\cdot)}$ afin d'améliorer l'adaptabilité de la solution par rapport aux reconstructions obtenues dans le cas standard d'espaces d'Hilbert; dans la deuxième partie, nous considérons une modélisation dans l'espace des mesures de Radon pour éviter les biais dus à la discrétisation observés dans les méthodes de régularisation parcimonieuse.

Plus en détail, la première partie explore à la fois des algorithmes d'optimisation lisse et non lisse dans les espaces $L^{p(\cdot)}$ réflexifs, qui sont des espaces de Banach dotés de la norme dite de Luxemburg. Comme premier résultat, nous fournissons une expression des cartes de dualité dans ces espaces, qui sont un ingrédient essentiel pour la conception d'algorithmes itératifs efficaces. Pour surmonter la non-séparabilité de la norme sous-jacente et les temps de calcul conséquents, nous étudions ensuite la classe des fonctions modulaires qui étendent directement la puissance (non homogène) $p > 1$ des normes L^p au cadre $L^{p(\cdot)}$. En termes de fonctions modulaires, nous formulons des analogues des cartes duales qui sont plus adaptées aux algorithmes d'optimisation lisse et non lisse en raison de leur séparabilité. Nous étudions alors des algorithmes de descente de gradient (à la fois déterministes et stochastiques) basés sur les fonctions modulaires, ainsi que des algorithmes modulaires de gradient proximal dans $L^{p(\cdot)}$, dont nous prouvons la convergence en termes des valeurs de la fonctionnelle. La flexibilité de ces espaces s'avère particulièrement avantageuse pour la modélisation de la parcimonie et les statistiques hétérogènes du signal/bruit, tout en restant efficace et stable d'un point de vue de l'optimisation. Nous validons cela numériquement de manière approfondie sur des problèmes inverses exemplaires en une/deux dimension(s) (déconvolution, débruitage mixte, tomographie).

La deuxième partie de la thèse se concentre sur la formulation des problèmes inverses avec un bruit de Poisson formulés dans l'espace des mesures de Radon. Notre contribution consiste en la modélisation d'un modèle variationnel qui couple un terme de données de divergence de Kullback-Leibler avec la régularisation de la Variation Totale de la mesure souhaitée (une somme pondérée de Diracs) et une contrainte de non-négativité. Nous proposons une étude détaillée des conditions d'optimalité et du problème dual correspondant. Nous considérons une version améliorée de l'algorithme de Sliding Franke-Wolfe pour calculer la solution numérique du problème de manière efficace. Pour limiter la dépendance des résultats du choix du paramètre de régularisation, nous considérons une stratégie d'homotopie pour son ajustement automatique où à chaque itération algorithmique, on vérifie si un critère d'arrêt défini en termes du niveau de bruit est vérifié et on met à jour le paramètre de régularisation en conséquence. Plusieurs expériences numériques sont rapportées à la fois sur des données de microscopie de fluorescence simulées en 1D/2D et réelles en 3D.

Mots-clés: optimisation non lisse, problèmes inverses en imagerie, régularisation en espaces de Banach, parcimonie, microscopie à fluorescence.

Abstract

Algoritmi di ottimizzazione in spazi di Banach non standard per problemi inversi di ricostruzione di immagini

Questa tesi si concentra sulla modellizzazione, l'analisi teorica e l'implementazione numerica di algoritmi di ottimizzazione avanzati per problemi inversi di imaging (ad esempio, ricostruzione di immagini in tomografia computerizzata, deconvoluzione di immagini in microscopia) in spazi di Banach non standard.

È diviso in due parti: nella prima, si considera il contesto degli spazi di Lebesgue a esponente variabile $L^{p(\cdot)}$ per migliorare l'adattività della soluzione rispetto alle ricostruzioni standard di Hilbert; nella seconda, si utilizza una modellizzazione nello spazio delle misure di Radon per evitare le distorsioni osservate nei metodi di regolarizzazione sparsi a causa della discretizzazione.

Più in dettaglio, la prima parte esplora algoritmi di ottimizzazione lisci e non lisci in spazi $L^{p(\cdot)}$ riflessivi, che sono spazi di Banach dotati della cosiddetta norma di Luxemburg. Come primo risultato, forniamo un'espressione delle mappe di dualità in questi spazi, che sono un ingrediente essenziale per la definizione di algoritmi iterativi efficaci. Per superare la non separabilità della norma sottostante e i conseguenti pesanti tempi di calcolo, studiamo poi la classe di funzionali modulari che estendono direttamente la potenza p (non omogenea) delle norme L^p al caso generale $L^{p(\cdot)}$. In termini dei funzionali modulari, formuliamo degli analoghi delle mappe di dualità, che, grazie alla loro separabilità, sono utilizzabili per algoritmi di ottimizzazione sia lisci che non lisci. Studiamo quindi algoritmi di discesa del gradiente basati su funzioni modulari (sia in un contesto deterministico che stocastico) e algoritmi di gradiente prossimale basati su funzioni modulari in $L^{p(\cdot)}$, e dimostriamo la loro convergenza in termini dei valori delle funzioni. La flessibilità spaziale di questi spazi si rivela particolarmente vantaggiosa per gestire la sparsità, la conservazione dei bordi e le statistiche eterogenee di segnale/rumore, pur rimanendo efficiente e stabile dal punto di vista dell'ottimizzazione. Abbiamo validato numericamente questo metodo su problemi inversi esempi 1D/2D (deconvoluzione, denoising misto, ricostruzione TC).

La seconda parte della tesi si concentra su problemi inversi con rumore di Pois-

son formulati nello spazio delle misure di Radon. Il nostro contributo consiste nell'elaborazione di un modello variazionale che accoppia il termine di fedeltà della divergenza di Kullback-Leibler con il termine di regolarizzazione della variazione totale della misura desiderata (cioè una somma pesata di Diracs) insieme a un vincolo di non negatività. Viene effettuato uno studio dettagliato delle condizioni di ottimalità e del corrispondente problema duale e viene utilizzata una versione migliorata dell'algoritmo Sliding Franke-Wolfe per calcolare in modo efficiente la soluzione numerica. Per attenuare la dipendenza dei risultati dalla scelta del parametro di regolarizzazione, viene proposta una strategia di omotopia per la sua selezione automatica, in cui, a ogni iterazione dell'algoritmo, si verifica se un criterio di arresto stabilito in termini di livello di rumore è verificato e si aggiorna di conseguenza il parametro di regolarizzazione. Sono riportati diversi esperimenti numerici su dati di microscopia a fluorescenza 2D simulati e 3D reali.

Parole chiave: ottimizzazione non liscia, problemi inversi di ricostruzione di immagini, regolarizzazione in spazi di Banach, regolarizzazione sparsa, microscopia a fluorescenza.

Contents

List of Figures	xi
List of Tables	xv
List of Algorithms	xvi
1 Introduction	1
1.1 Inverse problems	3
1.1.1 Least-squares solutions and generalised inverse	5
1.2 Regularisation methods	6
1.2.1 Tikhonov regularisation	7
1.2.2 Iterative regularisation	8
1.3 Bayesian framework and variational approach	9
1.3.1 Choice of fidelity, penalty and solution space	11
1.4 Why Banach spaces?	13
1.4.1 Variable exponent Lebesgue spaces	17
1.4.2 Banach space of Radon measures	17
1.5 Optimisation methods	18
1.5.1 Proximal operator and resolvent	19
1.5.2 Smooth optimisation methods	20
1.5.3 Forward-backward splitting	20
1.5.4 Mirror descent	22
1.6 Outline and contribution	22
I Modular-based optimisation in variable exponents Lebesgue spaces	27
2 Variable Exponent Lebesgue Spaces	29
2.1 Modular and Luxemburg norm	30
2.1.1 Inequalities between norm and modular	32
2.1.2 Properties of $L^{p(\cdot)}(\Omega)$ and immersions	35
2.2 Dual space and duality mappings	36

2.2.1	Definition of dual and associate space	36
2.2.2	Duality mappings	37
2.2.2.1	Duality mappings in $L^{p(\cdot)}(\Omega)$	38
2.2.2.2	Inverse of duality mappings	42
2.3	Modular-based alternative to duality maps	42
2.3.1	Separability	44
2.4	Final discussion	46
3	Smooth optimisation in $L^{p(\cdot)}(\Omega)$	47
3.1	Proximal operators in Banach spaces	48
3.1.1	Definition of the p-norm proximal operator	49
3.1.2	Bregman-proximal operator	50
3.2	Landweber methods in Banach spaces	52
3.2.1	Hilbert spaces setting	52
3.2.2	Dual method in Banach spaces	54
3.2.3	Primal method in Banach spaces	56
3.2.4	Primal and dual method as proximal point algorithms	56
3.3	Modular-based dual method in $L^{p(\cdot)}(\Omega)$	57
3.3.1	Primal and dual methods in $L^{p(\cdot)}(\Omega)$: main issues	58
3.3.1.1	Approximation of the inverse of the duality map	58
3.3.1.2	Heavy computational times	59
3.3.2	Modular-based alternative to dual method in $L^{p(\cdot)}(\Omega)$	59
3.3.3	Comparison between Landweber and modular-based gradient descent	60
3.3.4	How to choose variable exponents	62
3.3.5	Numerical tests with modular-based gradient descent	64
3.4	A Modular-based Stochastic variant	66
3.4.1	Stochastic Gradient Descent in Banach spaces	67
3.4.2	Variable exponents modular-based SGD	69
3.4.3	Numerical results	70
3.4.3.1	Hyper-parameter selection	70
3.4.3.2	Simulated data	71
3.4.3.3	Real CT dataset	73
3.5	Final discussion	74
4	Proximal gradient algorithms in $L^{p(\cdot)}(\Omega)$	77
4.1	Norm-based proximal-gradient algorithms in Banach spaces	78
4.1.1	Proximal primal-gradient algorithm	79
4.1.2	Proximal dual-gradient algorithm	79
4.2	Modular-based proximal primal-gradient algorithm	80
4.2.1	Convergence analysis	83
4.3	Modular-based proximal dual-gradient algorithm	88
4.4	Sparse reconstruction and thresholding functions	91
4.5	Numerical tests	94

4.5.1	Spike reconstruction	94
4.5.2	Deconvolution of heterogeneous signals	96
4.5.3	1D and 2D mixed noise removal	97
4.5.4	A numerical study on convergence rates	100
4.6	Final discussion	104

II Sparse optimisation in the Banach space of Radon measures with Poisson noise 105

5 Sparse off-the-grid optimisation methods in imaging 107

5.1	Inverse problems in the space of measures	108
5.1.1	Going <i>off-the-grid</i> for sparse spikes deconvolution	108
5.1.2	The space of Radon measures	109
5.1.3	Inverse problems in $\mathcal{M}(\Omega)$	111
5.2	The BLASSO problem	113
5.2.1	Optimality conditions	115
5.2.2	Dual problem and extremality conditions	116
5.2.2.1	Convex conjugate of the fidelity	116
5.2.2.2	Convex conjugate of the TV norm	117
5.2.2.3	Dual problem and extremality conditions	117
5.3	Frank-Wolfe and Sliding Frank-Wolfe algorithms	118
5.3.1	Frank-Wolfe algorithm	118
5.3.2	Frank-Wolfe for the minimisation of BLASSO	120
5.3.2.1	Greedy approach	121
5.3.3	Sliding Frank-Wolfe algorithm	123
5.4	Final discussion	124

6 Off-the-grid regularisation for Poisson inverse problems 127

6.1	Off-the-grid Poisson inverse problems	128
6.2	Dual problem and optimality conditions	130
6.2.1	Convex conjugate of the Kullback-Leibler divergence	130
6.2.2	Convex conjugate of the sum of penalty and indicator function	131
6.2.3	Dual problem formulation and extremality conditions	133
6.2.3.1	Subdifferential of the indicator function of positive measures	135
6.2.3.2	Extremality conditions	136
6.2.4	Optimality conditions	137
6.3	Algorithmic choices	137
6.3.1	Boosted Sliding Frank Wolfe algorithm	138
6.4	Homotopy algorithm	139
6.4.1	Homotopy algorithm for off-the-grid methods	142
6.4.1.1	Starting value	143
6.4.1.2	Updating rule	144

6.4.1.3	Descent property of the homotopy algorithm	144
6.4.1.4	Regularisation path in the discrete setting	146
6.4.1.5	Homotopy and regularisation path of BLASSO	147
6.5	Numerical tests	150
6.5.1	Comparison between Gaussian and Poisson modelling	150
6.5.2	Homotopy algorithm: choice of σ_{target}	154
6.5.3	Numerical test with 3D real dataset	156
6.6	Final discussion	158
7	Conclusions	161
	Appendix	165
A	Computed Tomography: forward model and geometries	167
B	Fluorescence microscopy imaging	171
	Bibliography	177

List of Figures

1.1	Inverse problems consists in retrieving the underlying signal corresponding to a given acquired data, where the acquisition process can be mathematically modelled.	2
1.2	Small perturbations in the acquisitions leads to instability in the reconstructions if the inverse T^{-1} of the forward operator is considered.	4
1.3	Effect of the regularisation parameter λ	8
1.4	Tikhonov-type reconstructions for different penalties based on Hilbert and Banach spaces norm. Images from [204].	14
1.5	Top row: the exact solution (left) and the exact data y (in blue) and the noisy y^δ (red). Middle row: results obtained in $L^{1.5}(\Omega)$. Bottom row: result obtained in L^2 . Image from [37].	15
1.6	Exemplar deconvolution imaging problem on a test image: ground truth, acquired data and reconstructions obtained in L^2 and in L^p with $p = 1.3$ spaces.	16
1.7	Exemplar deconvolution imaging problem on a satellite image: ground truth, acquired data and reconstructions obtained in L^2 and in L^p with $p = 1.3$ spaces.	16
2.1	Test image: ground truth, acquired data, variable exponent.	34
2.2	Satellite image: ground truth, acquired data, variable exponent.	34
3.1	Deblurring of test image in $L^{p(\cdot)}(\Omega)$ spaces. The restored images, here shown, are attained in correspondence with the minimum of the reconstruction error.	61
3.2	Deblurring of satellite image in $L^{p(\cdot)}(\Omega)$ spaces. The shown restored images are attained in correspondence of the minimum of the reconstruction error.	62
3.3	Semi-convergence of Landweber method and modular-based GD in $L^{p(\cdot)}(\Omega)$ spaces.	63

3.4	Numerical test on simulated CT data with gradient descent algorithms. Comparison between Hilbert and non-Hilbertian Lebesgue spaces reconstructions with constant and variable exponents after 500 iterations.	64
3.5	Quality metrics of \mathbf{GD}_2 , \mathbf{GD}_p , $\mathbf{GD}_{p(\cdot),q(\cdot)}$ and $\mathbf{GD}_{p(\cdot),q(\cdot)}$ adaptive.	66
3.6	Partition of the observation $y \in \mathcal{Y}$ into $y_i \in \mathcal{Y}_i$ with $N_s = 30$	68
3.7	Simulated CT data: noisy acquisition, ground truth, exponent maps $p(\cdot)$ and $q(\cdot)$, reconstructions (after 40 epochs) with \mathbf{SGD}_2 , \mathbf{SGD}_p , $\mathbf{SGD}_{p(\cdot),q(\cdot)}$ and $\mathbf{SGD}_{p(\cdot),q(\cdot)}$ adaptive.	71
3.8	Quality metrics along the first 100 epochs of \mathbf{SGD}_2 ; $\mathbf{SGD}_{1.1}$; $\mathbf{SGD}_{p(\cdot),q(\cdot)}$ with and without adapting the exponent maps $p(\cdot)$, $q(\cdot)$	72
3.9	Real CT data: noisy sinogram and different reconstructions obtained with SGD strategies.	74
4.1	1D thresholding functions for proximal-gradient algorithms in Banach spaces. (a) $T^{\text{Alg.P}}(\cdot, s, t, p)$, $T^{\text{Alg.D}}(\cdot, s, t, p)$ and $T^{\text{ISTA}}(\cdot, s, t)$ with $p = 1.3$, $s = 0.3$, $t = 0.4$. (b) $T^{\text{Alg.P}}(\cdot, s, t, p)$ with $s = 0.3$, $t = 0.4$ and $p \in \{1.2, 1.4, 1.6, 1.8, 2\}$. (c) $T^{\text{Alg.D}}(\cdot, s, t, p)$ with $s = 0.3$, $t = 0.4$ and $p \in \{1.2, 1.4, 1.6, 1.8, 2\}$	93
4.2	Spike reconstruction in Lebesgue spaces: comparison between constant $L^{1.7}(\Omega)$ and variable $L^{p(\cdot)}(\Omega)$ exponent Lebesgue spaces. Parameters: $\tau_k \equiv 0.5$; $\lambda = 10^{-2}$. Stopping criterion based on the normalised relative change between x^k and x^{k+1} : $\ x^k - x^{k+1}\ _2 / \ x^k\ _2 < 10^{-4}$	95
4.3	Deconvolution of heterogeneous signals in Hilbert $L^2(\Omega)$, classical Banach $L^{1.5}(\Omega)$ and variable exponent Lebesgue space $L^{p(\cdot)}(\Omega)$ settings. Parameters: $\tau_k \equiv 0.5$; $\lambda = 5 * 10^{-3}$. Stopping criterion based on the relative distance between x^k and x^{k+1} : $\ x^k - x^{k+1}\ _2 / \ x^k\ _2 < 4 * 10^{-6}$	96
4.4	1D mixed noise removal: different choices for the solution and output spaces. Parameters: $\tau_k \equiv 0.1$, $\lambda = 2 * 10^{-2}$. Stopping criterion based on the relative distance between x^k and x^{k+1} : $\ x^k - x^{k+1}\ _2 / \ x^k\ _2 < 4 * 10^{-6}$	98
4.5	2D mixed denoising. Parameters: $\tau_k \equiv 0.1$, $\lambda = 0.1$. Stopping criterion based on the normalized relative change between x^k and x^{k+1} : $\ x^k - x^{k+1}\ _2 / \ x^k\ _2 < 10^{-4}$	100
4.6	Numerical study on convergence rates in function values $\phi(x^k) - \phi(\tilde{x})$ for proximal gradient algorithms in Banach spaces. (a) Relative rates along the first $4 * 10^4$ iterations. (b) Normalised relative rates along the first 60 seconds of CPU time.	102
4.7	Numerical study of the convergence rate in function values of Algorithm 7. Relative rates along the first 14000 iterations and different choices of variable exponents functions $p(\cdot)$ with $p_- = 1.7$	103

5.1	From on-the-grid to off-the-grid formulation. In red, the coarse grid of the acquisition $y \in \mathbb{R}^M$. In blue, the fine grid of the reconstruction $x \in \mathbb{R}^N$ with $N = L^2M$. In green, the point-sources to localise.	109
5.2	Comparison between conventional discrete (on-the-grid) and off-the-grid reconstructions. In black: the ground-truth spikes to retrieve. In Fig.5.2a: in blue, the acquired blurred and noisy signal. In Fig.5.2b: in red, discrete reconstruction with support constrained on a grid with M pixels. In Fig.5.2c: in red, discrete reconstruction with support constrained on a grid with $N > M$ pixels. In Fig.5.2d: in green, the off-the-grid reconstruction. The green spikes are the reconstruction without an a priori fixed grid, so they can move continuously on the line.	110
5.3	Discrete ground truth measure (black) and acquisition (blue) attained with the convolution with a Gaussian PSF (and Gaussian noise in Fig.5.3b).	112
5.4	Examples of biological fluorescent microscopy images. From left to right: molecules, cells, microtubules.	113
5.5	Conic particle gradient descent applied for 2D sparse spike deconvolution with Gaussian kernel. White dots are the source measure and red dots are the measures at iterations k with $k = 0$ in (a), $k = 150$ in (b) and $k = 1000$ in (c). The background image is the acquisition y . The black lines are the paths of the particles and constitute the gradient flow. Images from [141].	119
5.6	Reconstruction of 1D peaks from a blurred and noisy signal obtained using Frank-Wolfe algorithm for $\lambda = 0.5$ and $\lambda = 1000$, and with Sliding Frank-Wolfe with $\lambda = 0.5$	123
5.7	Frank-Wolfe and Sliding Frank-Wolfe algorithms for 2D sparse spike deconvolution with Gaussian kernel. White dots are the ground truth measure and red dots the reconstructed measure at iterate k . The background image is the acquisition y , obtained with a 2D Gaussian PSF with $\sigma = 0.1$ and Poisson noise. Results with $\lambda = 0.001$ on the domain $\Omega = [0, 1]^2$. The importance of the sliding step appears evident in this 2D example with 3 Diracs.	124
6.1	Reconstructions obtained using SFW in a 1D sparse deconvolution numerical example with Poisson noise. In black, the ground truth spikes and in green the reconstructed ones are shown. With λ close to 0, the number and intensities of spikes are overestimated, while with a much higher value they are underestimated.	139
6.2	Pareto frontier	140
6.3	Regularisation path for LASSO (discrete setting). Plot of $\lambda \mapsto x_\lambda(j)$ with different colours for each different $j \in J$	146

6.4	Regularisation path in the off-the-grid setting	148
6.5	1D comparison between Gaussian and Poisson noise modelling. In black: ground truth spikes. In green: reconstructed spikes. For both models, $\lambda = 8.82$	150
6.6	True Positives, False Positives, False Negatives spikes with respect to a tolerance radius $\delta > 0$	151
6.7	Mean values over 100 different randomly generated ground truths with 6 spikes and their corresponding reconstructions. Shaded area corresponds to standard deviation. Maximum number of iterations of SFW: $2N_{molecules}$. Tolerance radius $\delta = 0.05$	152
6.8	2D sparse ground truth image (white crosses) and its corresponding noisy blurred acquisition y^δ on $\Omega = [0, 1]^2$ (obtained with a 2D Gaussian PSF with $\sigma = 0.07$, constant background $b = 0.05$, Poisson noise). On the right, visualisation of $y^\delta _{\Omega_{bg}}$ corresponding to background noise, i.e. the external square-ring.	154
6.9	2D reconstruction with homotopy Alg.11 with parameters: max. number outer homotopy iterations 20, max. number inner SFW iterations 1, $c = 1$, $\gamma = 0.2$	155
6.10	3D Gaussian PSF	157
6.11	ERES 3D data. Values of λ_t , σ_{target} and cost functional along the homotopy iterations.	158
6.12	Sparse reconstruction of the 3D real volume acquisition. Visualisation from different angles. The colours corresponds to the depth along the z direction.	159
6.13	ERES 3D real data. From top to bottom: (a) acquired volume y^δ , (b) estimated Gaussian 3D PSF, (c) $y^\delta _{\Omega_b}$, (d) reconstructed volume with 130 iterations of the homotopy algorithm μ_{rec} , (e) blurred observation $\Phi\mu_{rec} + b$ corresponding to the reconstruction μ_{rec}	160
A.1	CT scanner. Image from [12]	167
A.2	2D Parallel Beam Geometry. Image from [12]	168
A.3	2D Fan Beam Geometry. Image from this link ¹	169
B.1	Resolution limit imposed by wave nature of light. Image from this link ²	172
B.2	Airy pattern	172
B.3	Overlapping of two point sources approaching near. The Rayleigh criterion: two points are considered as just resolved when the maximum of one diffraction pattern coincides with the first minimum of the other. Image from [209].	173
B.4	Fluorescent labelled cell. Each colour represents a different emission wavelenght of the emitted light. Image from this link ³	174

List of Tables

2.1	Comparison between the values of Luxembourg norm, modular and classical p -norms for x being the image in Figures 1.6a, 1.6b and Figures 1.7a, 1.7b.	33
3.1	Comparison between norm-based and modular-based gradient descent for images in Figure 3.1.	61
3.2	Comparison between norm-based and modular-based gradient descent for images in Figure 3.2.	62
3.3	CPU times after 3000 iterations. MAE, PSNR, and SSIM values after 500 iterations (before noise over-fitting).	66
3.4	Comparison of per iteration cost and total CPU times after 3000 iterations for deterministic algorithms and after 100 epochs for stochastic algorithms with $N_s = 30$. MAE, PSNR and SSIM values for stochastic algorithms are computed after 40 epochs (before noise over-fitting).	72
4.1	Deconvolution of heterogeneous signals in Hilbert $L^2(\Omega)$, classical Banach $L^{1.5}(\Omega)$ and variable exponent Lebesgue space $L^{p(\cdot)}(\Omega)$ settings. Quantitative results.	97
4.2	1D mixed noise removal in Hilbert $L^2(\Omega)$, classical Banach $L^{1.5}(\Omega)$ and variable exponent Lebesgue space $L^{p(\cdot)}(\Omega)$ settings. Quantitative results.	97
4.3	2D mixed denoising in Hilbert $L^2(\Omega)$, classical Banach $L^{1.5}(\Omega)$ and variable exponent Lebesgue space $L^{p(\cdot)}(\Omega)$ settings. Quantitative results.	99
4.4	Algorithmic comparison: iterations required and CPU time till convergence.	102
6.1	Homotopy algorithmic choices for BLASSO and for the Poisson off-the-grid models.	143
6.2	Parameters used with the homotopy algorithm (Alg.11) in the 1D simulated comparison tests	153

6.3	Homotopy algorithm: comparison between BLASSO and the Poisson off-the-grid modelling. Mean values over 100 different randomly generated ground truths with 6 spikes and their corresponding reconstructions.	154
6.4	Different estimates of σ_{target}	155
6.5	Parameters used for Algorithm 11 for the reconstruction of the 3D volume	157

List of Algorithms

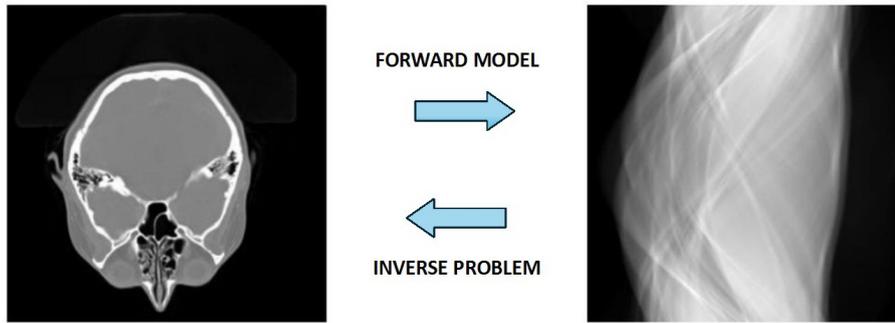
1	Landweber algorithm (dual method) in Banach spaces [207]	54
2	Primal method in Banach spaces	56
3	Modular-based Gradient Descent in $L^{p(\cdot)}(\Omega)$	59
4	Stochastic Modular-based Gradient Descent in $L^{p(\cdot)}(\Omega)$	70
5	Proximal primal-gradient algorithm in Banach spaces	78
6	Proximal dual-gradient algorithm in Banach spaces	79
7	Modular-based proximal primal gradient algorithm in $L^{p(\cdot)}(\Omega)$ spaces	81
8	Modular-based proximal dual gradient algorithm in $L^{p(\cdot)}(\Omega)$ spaces	89
9	Frank-Wolfe (FW) algorithm [97]	120
10	Sliding Frank-Wolfe (SFW) algorithm [75]	125
11	Homotopy algorithm for off-the-grid inverse problems	142

Introduction

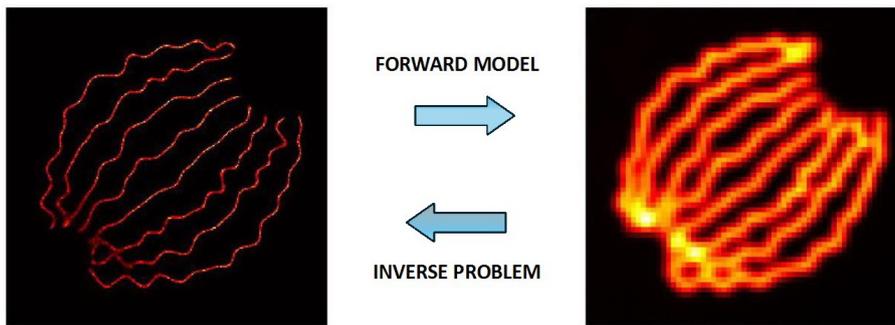
This chapter gives a brief introduction on the mathematical theory of inverse problems, focusing on regularisation methods and variational formulation, highlighting links and connections with optimisation techniques. In particular, the choice of solving inverse problems in Banach spaces will be motivated, providing some specific examples.

1.1	Inverse problems	3
1.1.1	Least-squares solutions and generalised inverse	5
1.2	Regularisation methods	6
1.2.1	Tikhonov regularisation	7
1.2.2	Iterative regularisation	8
1.3	Bayesian framework and variational approach	9
1.3.1	Choice of fidelity, penalty and solution space	11
1.4	Why Banach spaces?	13
1.4.1	Variable exponent Lebesgue spaces	17
1.4.2	Banach space of Radon measures	17
1.5	Optimisation methods	18
1.5.1	Proximal operator and resolvent	19
1.5.2	Smooth optimisation methods	20
1.5.3	Forward-backward splitting	20
1.5.4	Mirror descent	22
1.6	Outline and contribution	22

In many real-world situations, from medicine to geophysics, from neuroscience to biology, we encounter problems in which the quantities of interest are not directly observable but it is only possible to infer information on them from observations of



(a) On the left: section of a human head. On the right: the corresponding sinogram, i.e. the acquisition obtained with CT. Images from [18].



(b) On the right: acquisition of the ISBI SMLM 2013 dataset. On the left: super-resolved image from [144]. Dataset now available at [this link](#)¹.

Figure 1.1: Inverse problems consists in retrieving the underlying signal corresponding to a given acquired data, where the acquisition process can be mathematically modelled.

other measurable quantities, somehow related to the first ones. In applied mathematics these are known as inverse problems. The goal is to quantify unknown, but desired, parameters from relevant observed data, as sketched in Figure 1.1.

Inverse problems play a crucial role in modern medical diagnostics. Imaging techniques such as computed tomography (CT) scans, magnetic resonance imaging (MRI), and ultrasound are among the most known examples of inverse problems in medical imaging, see [168] for a complete analysis of these methods. Thanks to optimisation algorithms and computational modelling, detailed representations of internal anatomical structures are reconstructed from measurements of X-rays' attenuation, magnetic field generated by water dipoles in the human body, and scattered ultrasound waves, respectively.

Inverse problems have many other applications beyond the medical field. To name a few, in geophysics, a common inverse problem is subsurface imaging [132, 220], which involves inferring the composition and structure of the Earth's subsur-

¹https://github.com/KrakenLeaf/SPARCOM/blob/master/Example/EPFL/epfl_short_1.tif

face from electrical conductivity profiles reconstructed from electromagnetic induction measurements. In biology, fluorescence microscopy imaging techniques [76] are employed to reconstruct high-resolution images with fine scale details from noisy and blurred measurements. Images obtained through optical microscopes are inherently affected by blur caused by light diffraction. This phenomenon is unavoidable, no matter the quality of the lenses, and it causes distortions and a significant loss of resolution in the acquired data. Thus, acquired images need to be processed in order to enhance their resolution to enable scientists to observe biological structures with finer details, beyond the diffraction limit. Solving the inverse problem allows to image subcellular structures, such as proteins and cellular organelles, at a much finer scale than traditional light microscopy.

In neuroscience, electroencephalography (EEG) measures electrical activity on the scalp to determine the locations of the neural sources within the brain that generate these signals, i.e. EEG is a source localisation inverse problem [101]. EEG is used to study brain functionality, particularly in tasks like cognitive processing, emotion regulation, and identifying regions involved in various neurological conditions.

In some cases, the space of the solution (source space) and of the acquisition (measurement space) is the same, such as in deblurring of images (see for example Figure 1.1b), whilst in others solutions and acquisitions belong to different spaces. The latter is the classical setting of medical imaging problems, where the data is acquired outside of the human body to obtain visual information about inaccessible inside body parts, as in Figure 1.1a.

All the aforementioned scenarios have some common factors: availability of measured data, need to reconstruct not directly observable images or signals from the observations, and a known and mathematically modelled relation between event of interest and its indirect observations. Indeed, the term *inverse problem* makes sense only when there is an underlying *direct*, or *forward*, problem, describing the physics of the acquisition process, which can be modelled and parametrised into mathematical terms. A *forward operator* maps objects of interest into information collected about these objects, that is, the measurements or data. Solving an inverse problem means to obtain an estimate of the desired unknown parameters of interest, exploiting the measurements, and some knowledge of the forward model, when available.

1.1 Inverse problems

Let \mathcal{X} , \mathcal{Y} be two normed vector spaces and let $T : \mathcal{X} \rightarrow \mathcal{Y}$ a linear continuous operator between them, i.e. $T \in \mathcal{L}(\mathcal{X}, \mathcal{Y})$. The equation

$$Tx = y, \quad x \in \mathcal{X}, y \in \mathcal{Y} \tag{1.1}$$

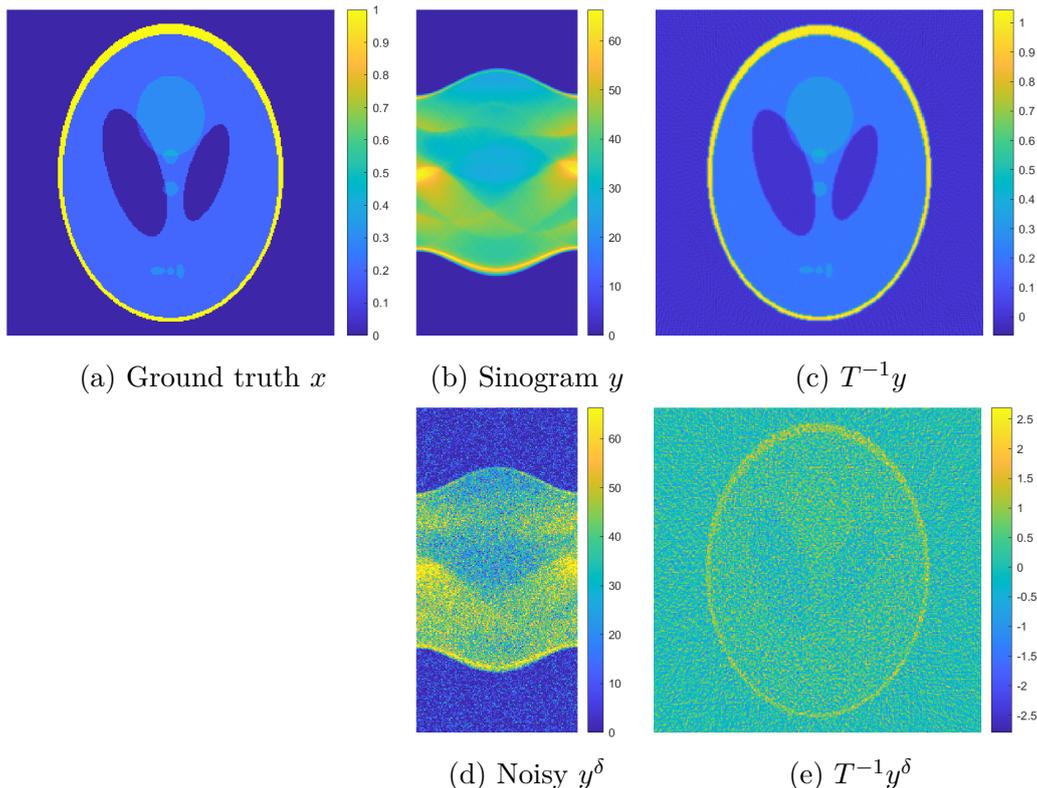


Figure 1.2: Small perturbations in the acquisitions leads to instability in the reconstructions if the inverse T^{-1} of the forward operator is considered.

implicitly defines the inverse problem [131] of finding x given y such that (1.1) is satisfied, for the operator T . An inverse problem is called *well-posed* in the sense of Hadamard [108] if the following conditions hold:

- the solution exists, i.e. for any $y \in \mathcal{Y}$ there exists $x \in \mathcal{X}$ such that $Tx = y$;
- the solution is unique, i.e. for any $x_1, x_2 \in \mathcal{X}$ such that $Tx_1 = Tx_2$ implies $x_1 = x_2$;
- the solution depends with continuity on the data: there exists a constant $C > 0$ such that $\|x_1 - x_2\|_{\mathcal{X}} \leq C\|y_1 - y_2\|_{\mathcal{Y}}$ where $Tx_1 = y_1$ and $Tx_2 = y_2$.

If at least one of the above conditions is not satisfied, the problem is said to be *ill-posed*.

The first condition means that the forward operator T has a full range $R(T) = \mathcal{Y}$ (T is surjective) and the second one is equivalent to require that the null space of T is trivial $N(T) = 0$ (T is injective). Thus, these two requirements guarantee the invertibility of the operator T . However, the existence of the inverse is not enough for the problem to be stable. The forward operator describes natural observable phenomena, such as the blur affecting image acquisition or the Radon transform for CT measurements (see Appendix A), and thus the measurement process inherently

produces small unavoidable errors. Indeed, it is nearly impossible in real applications to have a noise-free acquisition y . Instead of y , the ideal perfect data, only y^δ , a slightly perturbed version of y is often measured, is available, and (1.1) becomes:

$$Tx = y^\delta, \quad \|y^\delta - y\|_{\mathcal{Y}} \leq \delta \quad \text{for a small } \delta > 0. \quad (1.2)$$

Thus, it is reasonable to require that small perturbations of the data, quantified by δ , correspond to small perturbations of the solutions, that is the map $T^{-1} : \mathcal{Y} \rightarrow \mathcal{X}$ is continuous, so that a solution of (1.2) is not too far from a solution of (1.1). In the case of additive noise, i.e. $y^\delta = y + \delta$, by solving (1.2) just using the inverse T^{-1} , one gets

$$T^{-1}y^\delta = T^{-1}(y + \delta) = T^{-1}y + T^{-1}\delta = x + T^{-1}\delta.$$

Thus, if T^{-1} is unbounded the quantity $T^{-1}\delta$ might be very large compromising the quality of the reconstructed solution. This explains why the existence of the inverse T^{-1} does not suffice for the problem to be easily solvable. For a visual representation of this phenomenon see Figure 1.2.

1.1.1 Least-squares solutions and generalised inverse

The well-posedness of an inverse problem guarantees the existence (and uniqueness) of a solution but it does not necessarily entail that an explicit expression for T^{-1} is known. On the other hand, when the problem is ill-posed the definition of solution needs to be reformulated. For example, if the solution does not exist, i.e. given $y \in \mathcal{Y}$ for any $x \in \mathcal{X}$ equation (1.1) is not satisfied, one may still look for a *solution* that *almost* satisfies it. In the case of non-uniqueness of the solution, one may add some constraints to the solution set to ensure uniqueness, including *a priori* information of the expected solution (smoothness, sparsity with respect to some basis, etc.).

A possible way to relax the definition of solution when existence is not guaranteed, i.e. $y \notin R(T)$, is to look for

$$\bar{x} \in \mathcal{X} \text{ such that } \|T\bar{x} - y\|_{\mathcal{Y}} \leq \|Tu - y\|_{\mathcal{Y}} \quad \forall u \in \mathcal{X}. \quad (1.3)$$

Any such \bar{x} is called *least-squares solution* and it is not necessarily unique. The set of least-squares solutions can be found by computing

$$\operatorname{argmin}_{x \in \mathcal{X}} \frac{1}{2} \|Tx - y\|_{\mathcal{Y}}^2, \quad (1.4)$$

i.e. the set of minimisers of the smooth and convex functional $f : \mathcal{X} \rightarrow \mathbb{R}^+$ defined by $f(x) = \frac{1}{2} \|Tx - y\|_{\mathcal{Y}}^2$. To have uniqueness of the solution, one can choose in this set the *generalised solution* [91, 103] as

$$x^\dagger \in S \text{ s.t. } \|x^\dagger\|_{\mathcal{X}} \leq \|\bar{x}\|_{\mathcal{X}} \quad \forall \bar{x} \in S, \quad (1.5)$$

where S is the non-empty, closed and convex set of least-squares solutions, and thus the existence and uniqueness of x^\dagger is guaranteed.

Let now \mathcal{X} and \mathcal{Y} be Hilbert spaces and consider $T \in \mathcal{L}(\mathcal{X}, \mathcal{Y})$ with closed range $R(T)$. If $y \notin R(T)$, (1.1) does not admit any solution. Under these hypotheses, a definition of generalised inverse equivalent to (1.5) relies on relaxing the definition of solution by considering the orthogonal projection of the acquisition onto the range. In particular, by denoting with $P_{\bar{R}(T)} : \mathcal{Y} \rightarrow P_{\bar{R}(T)} \subseteq \mathcal{Y}$ the orthogonal projection onto the closure of $R(T)$, the following are equivalent [20, 21]:

- $T\bar{x} = P_{\bar{R}(T)}y$;
- $\|T\bar{x} - y\| \leq \|Tu - y\| \forall u \in \mathcal{X}$;
- $T^*T\bar{x} = T^*y$, being $T^* : \mathcal{Y} \rightarrow \mathcal{X}$ the adjoint operator.

In this setting, the least-squares solutions can be equivalently defined with one of the above. In particular, the generalised solution x^\dagger satisfies the equation

$$T^*Tx^\dagger = T^*y. \quad (1.6)$$

If $T : \mathcal{X} \rightarrow \mathcal{Y}$ is a linear continuous operator such that the dimension of the null space $N(T)$ coincides with the dimension of $R(T)^\perp$, the generalised solution can be equivalently defined as

$$x^\dagger = P_{N(T)^\perp} \tilde{T}^{-1}y$$

where

- $P_{N(T)^\perp} : \mathcal{X} \rightarrow N(T)^\perp \subseteq \mathcal{X}$ is the orthogonal projection onto $N(T)^\perp$, the orthogonal space of $N(T)$;
- $\tilde{T} : \mathcal{X} \rightarrow \mathcal{Y}$ is an invertible linear operator such that $T = P_{\bar{R}(T)}\tilde{T}$.

The operator $T^\dagger := P_{N(T)^\perp}\tilde{T}^{-1} : R(T) \oplus R(T)^\perp \subseteq \mathcal{Y} \rightarrow N(T)^\perp \subseteq \mathcal{X}$ is called *generalised inverse* [91]. If $R(T)$ is closed, then $R(T) \oplus R(T)^\perp = \mathcal{Y}$ and, thus, T^\dagger is well-defined everywhere on \mathcal{Y} and it is continuous. On the contrary, if the range of T is not closed, T^\dagger may not lead to a good solution. Indeed, being $R(T)$ non-closed, T^\dagger is not continuous and it is well-defined only on $R(T) \oplus R(T)^\perp \subsetneq \mathcal{Y}$. Having a non-continuous inversion process is undesirable, since it means that the solution produced does not depend with continuity on the data. This may lead to reconstructions far from the desired ones, when dealing with noisy data y^δ , as seen in Figure 1.2.

1.2 Regularisation methods

To tackle the aforementioned problem, a possible strategy is to consider regularisation methods. The idea behind regularisation is to stabilise the solution in

presence of noisy data by introducing a continuous inversion operator, so that the noise in the data is not reflected in the reconstruction.

Let $T \in \mathcal{L}(\mathcal{X}, \mathcal{Y})$, with \mathcal{X}, \mathcal{Y} Hilbert spaces, with non-closed range $R(T)$. Regularisation methods aim at approximating the discontinuous operator T^\dagger with a family of neighbouring continuous operators R_λ . Specifically, the family of operators $\{R_\lambda\}_{\lambda>0}$, with $R_\lambda : \mathcal{Y} \rightarrow \mathcal{X}$, is said a regularisation algorithm [102] if:

- R_λ is a linear and continuous (hence bounded) operator;
- $\forall y \in R(T) \oplus R(T)^\perp \lim_{\lambda \rightarrow 0} R_\lambda y = x^\dagger$ with $x^\dagger = T^\dagger y$.

The parameter $\lambda > 0$ is referred to as *regularisation parameter* and $x_\lambda = R_\lambda y$ is the *regularised solution*. Regularisation methods approximate T^\dagger with an operator R_λ which is continuous everywhere. In this way, a regularised solution of the noisy problem (1.2) is not too far from the generalised solution of the noiseless problem (1.1). This does not happen with the generalised inverse T^\dagger , which is not bounded. Regularisation methods play a crucial role in the presence of noisy data ensuring stability in their reconstruction. Denoting by x^\dagger a generalised solution of the ideal noise-free inverse problem (1.1), and by x_λ the regularised solution of the noisy problem (1.2), the reconstruction error is bounded by

$$\|x_\lambda - x^\dagger\| = \|R_\lambda y^\delta - x^\dagger\| \leq \|R_\lambda T x^\dagger - x^\dagger\| + \delta \|R_\lambda\|,$$

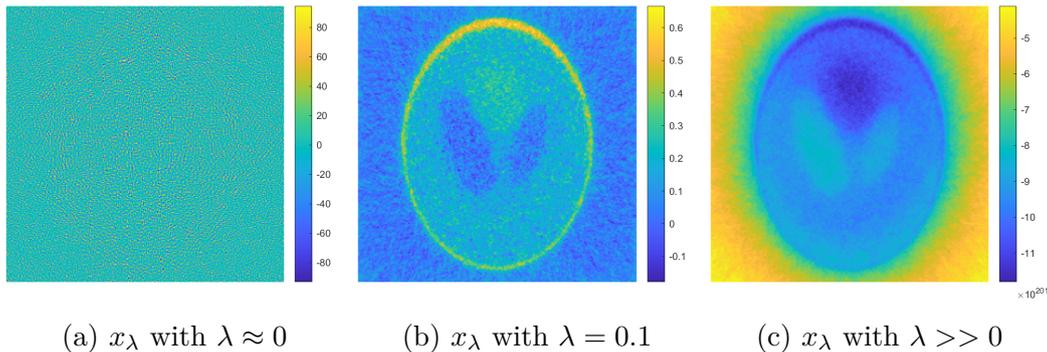
and it can be split into two parts. The first term is an approximation error due to the approximation of T^\dagger by R_λ and in general it is an increasing function of λ , whilst the second term measures the propagation error of the noise onto the solution and it is a decreasing function of λ . The choice of λ is therefore delicate: a good solution should compromise approximation and propagation errors at the same time.

1.2.1 Tikhonov regularisation

The most common regularisation method is the *Tikhonov regularisation* [212, 213] algorithm, which is defined by the following minimisation problem

$$x_\lambda := \operatorname{argmin}_{x \in \mathcal{X}} \|Tx - y\|_{\mathcal{Y}}^2 + \lambda \|x\|_{\mathcal{X}}^2, \quad \lambda > 0. \quad (1.7)$$

Using the notation of the previous section, the family of bounded linear operators $\{R_\lambda\}_{\lambda>0}$ defined as $R_\lambda = (T^*T + \lambda I)^{-1} T^*$ characterises Tikhonov regularisation. Indeed, the Tikhonov regularised solution satisfies the Euler-Tikhonov equation $(T^*T + \lambda I)x_\lambda = T^*y$ and it is an approximation of the generalised solution x^\dagger , which satisfies (1.6). The above functional is the sum of the residual $\|Tx - y\|_{\mathcal{Y}}^2$, that quantifies the mismatch between the observed data and the model's predictions, and the norm squared of the solution itself. By minimising the above functional, we look for an approximation of the solution with small energy, i.e. small norm, being this a characteristic associated with smooth, regular, noise free data. Thus,

Figure 1.3: Effect of the regularisation parameter λ

we seek for a solution with small norm that, at the same time, fits well enough the acquisition and is a good approximation of the generalised solution.

Tikhonov regularisation also refers to variational model as (1.7), where the penalty is of the form $\frac{1}{2}\|Dx\|^2$, with an operator $D : \mathcal{X} \rightarrow \tilde{\mathcal{X}}$. If D is the identity we retrieve exactly (1.7).

The regularisation parameter controls the trade off between having a small residual and having a solution with a small norm, that is between fitting the observations and having a regular and smooth solution. Choosing a small parameter might lead to overfitting of the reconstruction to the acquisition and thus to the noise, whilst a bigger value of λ might yield to a well-regularised and noise-free reconstruction that does not fit enough the data, as it is sketched in Figure 1.3.

1.2.2 Iterative regularisation

Another well-known approach to regularisation is the so-called *iterative regularisation*. One-step iterative algorithms, such as the Landweber method [110, 136] or the Conjugate Gradient method [117], represent the main class of regularisation schemes for functional equations, see e.g. [88, 174, 182], where the role of the regularisation parameter is played by the number of iterations of the given algorithm. It is now briefly presented the Landweber method in Hilbert spaces \mathcal{X} and \mathcal{Y} .

In order to approximate the generalised solution x^\dagger , consider (1.6). Given $\tau > 0$, it can be rewritten as a fixed-point iteration $x = G(x)$ with $G : \mathcal{X} \rightarrow \mathcal{X}$ such that $G(x) = x - \tau T^*(Tx - y)$. Given any starting point $x^0 \in \mathcal{X}$, its solution can be approximated by the sequence $(x^k)_{k \in \mathbb{N}}$ defined by $x^{k+1} = G(x^k)$, i.e.

$$x^{k+1} = x^k - \tau T^*(Tx^k - y). \quad (1.8)$$

The latter is referred to as *Landweber iteration* or *Landweber method* and it is a regularisation method [186]. It is important to observe that (1.8) can be also seen as an iteration of the gradient-descent algorithm [152] minimising the residual $\|Tx - y\|_{\mathcal{Y}}^2$, with $\tau \in \left(0, \frac{2}{\|T\|^2}\right)$ an appropriately chosen step-size.

To use the formalism about regularisation methods introduced before, we introduce the family of operators $\left\{R_{\frac{1}{k+1}}\right\}_{k \in \mathbb{N}}$, with $R_{\frac{1}{k+1}} : \mathcal{Y} \rightarrow \mathcal{X}$ defined by

$$R_{\frac{1}{k+1}}(\cdot) = (I - \tau T^* T)^{k+1} x^0 + \tau \sum_{i=0}^k (I - \tau T^* T)^i T^*(\cdot)$$

to rewrite (1.8) as $x^{k+1} = R_{\frac{1}{k+1}} y$. It can be proved that $\left\{R_{\frac{1}{k+1}}\right\}_{k \in \mathbb{N}}$ is a regularisation algorithm [102]. The more iterations are computed, the closer the regularised solution x^{k+1} will be to the generalised solution, at risk of over-fitting the noise. By stopping the regularisation algorithm early, this phenomenon is avoided and the obtained regularised reconstruction x^k is less influenced by the noise in the data. Indeed, an iterative method works as regulariser if an early-stopping strategy preventing over-fitting of the noise in the reconstructions is used according to the well-known semi-convergence property [167]. This is often being referred to as implicit regularisation [216], since no penalty term has to be introduced but regularisation is achieved by performing a relatively small number of iterations, and it is a very active field of research nowadays in disciplines such as machine and deep learning [7, 161, 164].

1.3 Bayesian framework and variational approach

In the mathematical theory of inverse problems, a classical approach to their resolution is the Bayesian one [123, 210]. In this framework, the unknown x and the data y (or y^δ if there is noise) are seen as realisations, respectively, of random variables \mathbf{X} and \mathbf{Y} distributed with probability density functions $\pi_{\mathbf{X}}(x)$ and $\pi_{\mathbf{Y}}(y)$.

Given an observation $y \in \mathbf{Y}$, let $\pi_{\mathbf{Y}|\mathbf{X}}(y|x)$ be the probability density of the fact that the observation y is produced by the underlying signal x . $\pi_{\mathbf{Y}|\mathbf{X}}(y|x)$ is called *likelihood*: it is a measure of how likely the realisation y of \mathbf{Y} comes from the realisation x of \mathbf{X} and it depends on the operator T . Within this setting, the problem (1.1) can be interpreted as the problem of maximising the likelihood function, that is

$$\hat{x} \in \operatorname{argmax}_{x \in \mathcal{X}} \pi_{\mathbf{Y}|\mathbf{X}}(y|x), \quad (1.9)$$

for a fixed realisation $y \in \mathbf{Y}$.

Instead of considering the maximum likelihood approach as above, another possible strategy is to maximise the so called *posterior probability* $\pi_{\mathbf{X}|\mathbf{Y}}(x|y)$, that is the probability of finding x given the observation y . Thanks to the Bayes formula,

$$\pi_{\mathbf{X}|\mathbf{Y}}(x|y) = \frac{\pi_{\mathbf{Y}|\mathbf{X}}(y|x)\pi_{\mathbf{X}}(x)}{\pi_{\mathbf{Y}}(y)}, \quad (1.10)$$

it is possible to express the *posterior probability* in terms of the *prior probability* $\pi_{\mathbf{X}}(x)$, the *conditional probability* or likelihood $\pi_{\mathbf{Y}|\mathbf{X}}(y|x)$ and the probability distribution $\pi_{\mathbf{Y}}(y)$ of the acquisitions, which serves as a normalisation constant. The prior

probability models *a priori* knowledge we have on the target solution x , such as its geometrical characteristics and/or sparsity features. It is independent of the data y . The likelihood depends on the forward model which maps the data $x \in \mathcal{X}$ with the observation $y \in \mathcal{Y}$ and on the probability distribution of the noise. According to the Maximum A Posteriori estimation, among all the possible realisation of the random variable \mathbf{X} , it is chosen as best approximation of the unknown x the one which maximises the posterior probability $\pi_{\mathbf{X}|\mathcal{Y}}(x|y)$ or, thanks to (1.10) (neglecting $\pi_{\mathcal{Y}}(y)$):

$$\operatorname{argmax}_{x \in \mathcal{X}} \pi_{\mathcal{Y}|\mathbf{X}}(y|x)\pi_{\mathbf{X}}(x),$$

which can be formulated as a minimisation problem by applying the logarithm function, as follows

$$\operatorname{argmin}_{x \in \mathcal{X}} \left[-\log \left(\pi_{\mathcal{Y}|\mathbf{X}}(y|x) \right) - \log \left(\pi_{\mathbf{X}}(x) \right) \right]. \quad (1.11)$$

Prior density functions are usually assumed to be of the form of Gibbs' potentials [98], that is:

$$\pi_{\mathbf{X}}(x) \sim e^{-\lambda R(x)}$$

where $R : \mathcal{X} \rightarrow \mathbb{R}_{\geq 0}$ is an energy functional, usually convex, and $\lambda > 0$ is a positive scalar parameter. A similar structure is assumed for the likelihood $\pi_{\mathcal{Y}|\mathbf{X}}(y|x)$ as well:

$$\pi_{\mathcal{Y}|\mathbf{X}}(y|x) \sim e^{-F(Tx,y)},$$

since most standard noise distributions have a negative exponential form. In this way, (1.11) turns out to be equivalent to the minimisation of the functional

$$\operatorname{argmin}_{x \in \mathcal{X}} J(x) := F(Tx, y) + \lambda R(x), \quad \lambda > 0. \quad (1.12)$$

- The first part of the functional is the *fidelity term* or *data-term*. It consists of a distance-like function $F : \mathcal{Y} \times \mathcal{Y} \rightarrow \mathbb{R}^+$, that measures the discrepancy between the model observations Tx corresponding to a possible solution x and the data y , and depends on the probability distribution of the noise.
- The second element is the *penalty* or *regularisation* term. It has to be chosen exploiting the *a priori* information about the desired solution. Its role is to stabilise the inversion process and to prevent overfitting or ill-posedness, so that to small perturbations of the data correspond small perturbations of the solution. This term has to promote solutions that are smooth or have certain properties for the sought solution.
- The *regularisation parameter* $\lambda > 0$ controls the trade-off between F and R , that is between fitting the observations and encouraging a *regular* solution, i.e. a solution with the properties enforced by the penalty term R .

This approach to the resolution of inverse problem is called *variational regularisation* [202]: it transforms the inverse problem into the minimisation problem of a structured functional $J : \mathcal{X} \rightarrow \mathbb{R}^+$. Under the assumption of Gaussian distribution for \mathbf{X} (with 0 mean and standard deviation $\sigma_{\mathbf{X}}$), i.e.

$$\pi_{\mathbf{X}}(x) = \frac{1}{\sqrt{2\pi\sigma_{\mathbf{X}}^2}} \exp\left(-\frac{1}{2} \frac{x^2}{\sigma_{\mathbf{X}}^2}\right), \quad (1.13)$$

and additive white Gaussian noise, that is

$$\pi_{\mathbf{Y}}(y|x, T) = \frac{1}{\sqrt{2\pi\sigma_{\mathbf{Y}}^2}} \exp\left(-\frac{1}{2} \frac{(Tx - y)^2}{\sigma_{\mathbf{Y}}^2}\right), \quad (1.14)$$

it can be proven that from (1.12) we retrieve exactly the Tikhonov regularisation functional (1.7).

Similarly, under the assumption of Gaussian distribution over \mathbf{X} (1.9) becomes (1.4), i.e. we retrieve the least-squares solutions.

1.3.1 Choice of fidelity, penalty and solution space

The choice of the different ingredients of the variational model (1.12) is crucial for obtaining a good reconstruction of the desired image or signal. The assumption of Gaussian distribution of the unknown signal (1.13) yielding to a penalty of the form $\|x\|_{\mathcal{X}}^2$ is convenient from a computational point of view, being it convex and smooth, but it is very often unrealistic. The *a priori* information about the desired solution is not well-described by such a penalty when, for example, the sought solution is sparse, i.e. equal to 0 in the vast majority of its domain $\Omega \subseteq \mathbb{R}^d$, $d \in \mathbb{N}$. The squared-norm might lead to an over-smoothed and not-sparse approximation of the solution. When the prior knowledge on the solution is sparsity-type information, different penalties have to be introduced. The most common one is the L^1 -norm [70, 82] of the solution

$$\|x\|_1 = \int_{\Omega} |x| dt$$

but many other choices are possible. Just to name a few, L^p penalties are proposed and studied in [70, 188], as clarified better in the next section, in [187] instead the use of a Besov norm as sparsity-promoting penalty in Tikhonov regularisation is proposed. In the discrete setting of compressed sensing [48, 81, 90], sparsity of the solution x with respect to some basis W is also sought, i.e. it is possible to write $x = Wz$ with z sparse, so that the penalty is $\|z\|_1$.

If the solution is known to have flat and constant regions with jumps and discontinuities, it means its gradient is sparse and a possible good choice for the penalty term is the Total-Variation norm, which for functions $x \in W^{1,1}(\Omega)$ reads [195]

$$\text{TV}(x) = \int_{\Omega} |\nabla x| dt.$$

Other studied penalty terms are the Non Local Total Variation [125], the Total Generalised Variation [38], the Hessian Schatten norm [150], the Elastic net regulariser [221] and many more.

A similar reasoning lies down the choice of the fidelity term. The assumption of Gaussian noise on the data (1.14) is often a good approximation of the real, and sometimes unknown, noise distribution by central limit theorem and leads to an L^2 -fidelity [53, 195]

$$\|Tx - y\|_2^2 = \int_{\Omega} (Tx - y)^2 dt.$$

However, in some circumstances this choice is too unrealistic and a more tailored choice for the fidelity leads to better quality in the reconstructions, for example, when dealing with Poisson noise. If the noisy data y is a realisation of a Poisson distributed random variable with mean Tx , i.e.

$$\pi_Y(y|x, T) = \frac{\exp(-Tx)}{y!} (Tx)^y,$$

where the operations are all intended point-wise, the Kullback-Leibler (KL) divergence [149, 201]

$$\mathcal{D}_{KL}(Tx, y) := \int_{\Omega} \left(Tx - y + y \log(y) - y \log(Tx) \right) dt \quad (1.15)$$

should be considered as fidelity. Alternatively, an approximation of the KL through Taylor expansions [199] can be considered. Computing the second order Taylor expansion of (1.15) with respect to the second variable centred in y yields a weighted- L^2 penalty

$$\|Tx - y\|_W^2 = \int_{\Omega} \frac{(Tx - y)^2}{y} dt,$$

where again the operations are all point-wise. We considered this fidelity in [143, 144, 147] coupled with a particular sparsity promoting ℓ^0 -type penalty, called WCEL0, that we proposed in [144]. This method has been used for the microscopy imaging reconstruction shown in Figure 1.1b. In presence of impulsive or salt-and-pepper noise, instead, the L^1 -norm of the residual [176, 177] is a good option

$$\|Tx - y\|_1 = \int_{\Omega} |Tx - y| dt,$$

since it well describes the sparse nature of the noise.

An ingredient that is sometimes overlooked and not properly taken into consideration is the space where solutions are sought for. Solving inverse problems in Hilbert spaces has many computational advantages but it can lead to over-smoothness of the solutions, bad reconstructions of edges and sparse patterns, such as small objects or impulse signals [37, 204]. For these reasons, it is often a good choice to consider inverse problems in Banach space setting, as it is sketched in the next section.

1.4 Why Banach spaces?

As well presented in [204], in a series of different applications models that use Hilbert spaces are not realistic or appropriate. To name a few, it is worth mentioning non-destructive testing, such as X-ray diffractometry [203, 205], phase retrieval [27, 28, 80, 119, 133, 134], parameter identification for special partial differential equations [62, 109], inverse problems in finance [113, 119, 120]. The nature of such applications requires Banach spaces modelling of the problem itself. In X-ray diffractometry, Banach spaces of continuous functions $\mathcal{C}([\alpha_1, \alpha_2])$ are needed and phase retrieval problems are defined in terms of a non-linear forward operator between Lebesgue space $L^p(\mathbb{R}^d)$ and its dual $L^{p'}(\mathbb{R}^d)$, with $p' = p/(p-1)$. For PDEs, again Lebesgue spaces are considered as solution and acquisition spaces, or more general non-Hilbertian Sobolev spaces. We refer to [204, Chapter 1] for an exhaustive description of all the above examples. Banach spaces not only can be used as solution and/or acquisition spaces \mathcal{X} and \mathcal{Y} in the definition of the problem (1.1), as for the above examples, but also in the formulation of variational models with specific fidelity and penalty terms defined in terms of Banach norms. Recently, the use of an L^p -norm to measure the fidelity term and of an L^q -norm to measure the regularisation term has received considerable attention [43, 122, 137].

Tikhonov-like regularisation with totally convex functions in reflexive Banach spaces has furthermore been recently studied, for instance, in [99] and in [64], in the framework of regularised learning schemes in feature Banach spaces.

Another well-known scenario where Banach spaces are known to perform better than Hilbert spaces is sparse inverse problems, that is when sparse solutions of ill-posed operator equations are to be determined. In the seminal work [70] by Defrise, Daubechies and De Mol, they propose the use of L^p -norms

$$\|x\|_p = \left(\int_{\Omega} |x|^p dt \right)^{1/p}, \quad 0 < p < 2,$$

as sparsity enforcing penalty, also studied in [130]. The effect of such a penalty is detailed with numerical examples in [204]. We report in Figure 1.4 a simulated numerical test from [204] on the problem of retrieving the velocity of a vehicle $x(t)$ given its position $y(t)$ for all $t \in [0, T]$:

$$y(t) = \int_0^t x(s) ds + y(0) + \delta = Tx + \delta,$$

where $\delta > 0$ is an additive noise. The following variational model is considered

$$\operatorname{argmin} \frac{1}{2} \|Tx - y\|_2^2 + \lambda \|x\|_p^p, \quad \lambda > 0,$$

for $p = 2$, $p = 1$ and $p = 1.1$. The ground truth velocity of Figure 1.4(a) presents three sharp peaks. However, due to the physical setting of the problem, the peaks have steep slopes but are not discontinuous.

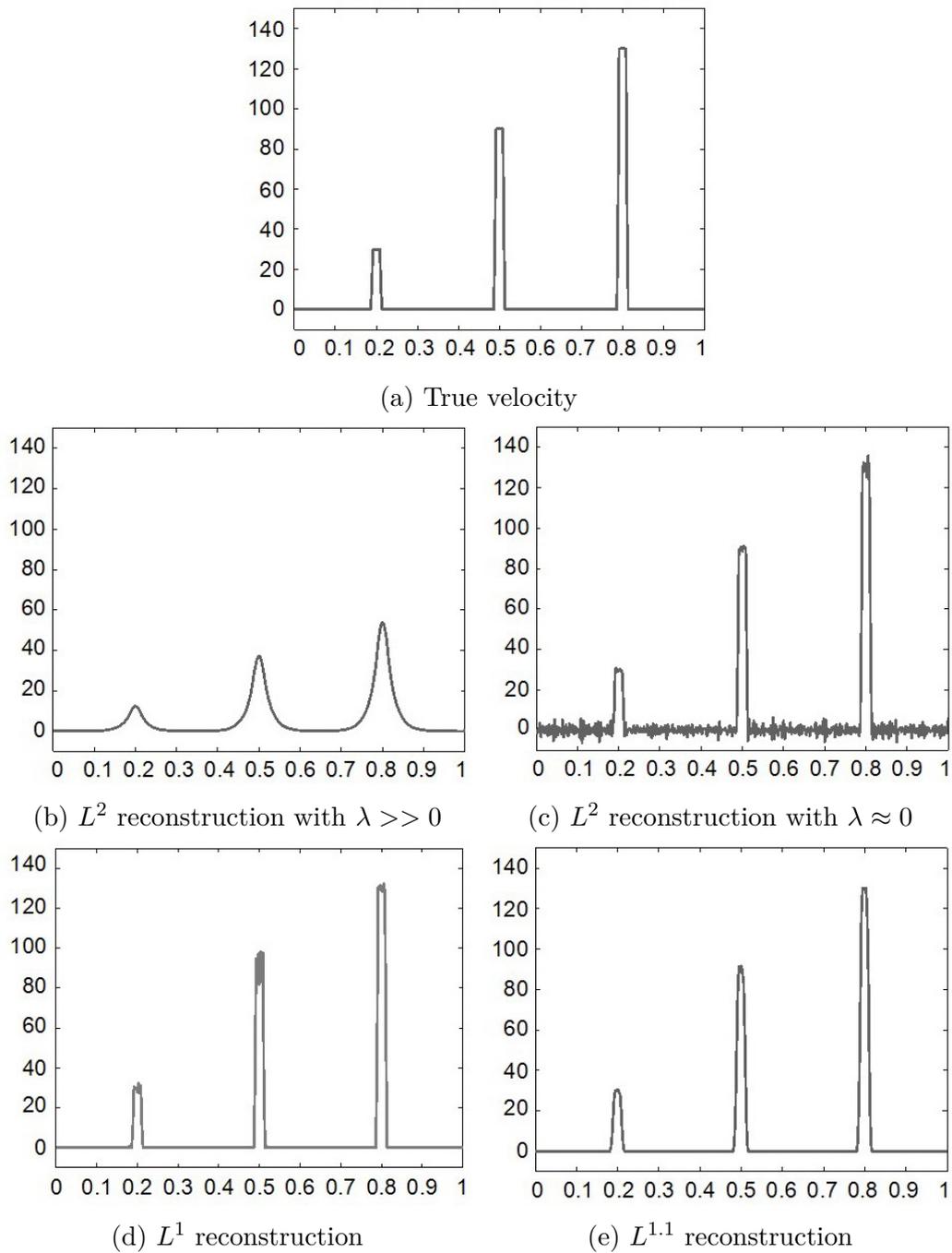


Figure 1.4: Tikhonov-type reconstructions for different penalties based on Hilbert and Banach spaces norm. Images from [204].

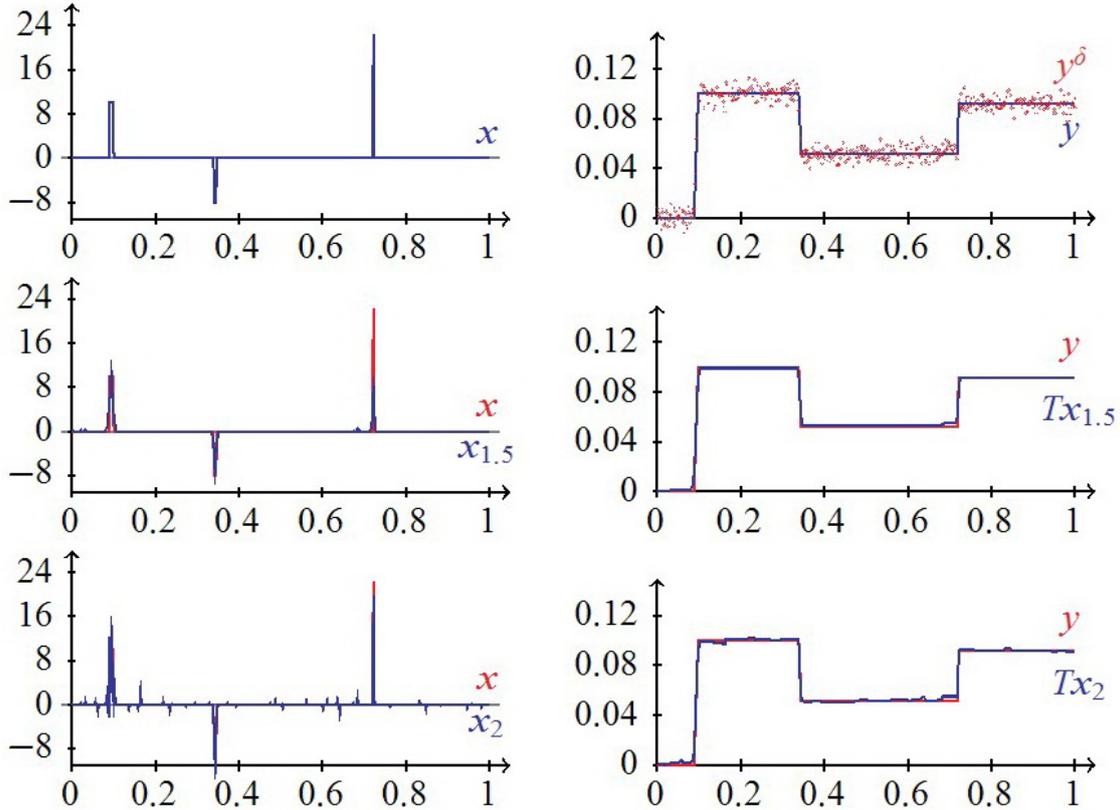


Figure 1.5: Top row: the exact solution (left) and the exact data y (in blue) and the noisy y^δ (red). Middle row: results obtained in $L^{1.5}(\Omega)$. Bottom row: result obtained in L^2 . Image from [37].

In the case of $p = 2$, reconstructions cannot recover well the signal: the reconstruction obtained with a big value of λ is over-smoothed (Figure 1.4(b)), whilst the one with $\lambda \approx 0$ is affected by the noise (Figure 1.4(c)). Since the true signal consists of peaks, it can be considered sparse, and thus by taking $p = 1$ we consider an L^1 sparsity promoting penalty. The corresponding reconstruction is shown in Figure 1.4(d), significantly improved with respect to the ones obtained with the L^2 Hilbertian penalty. A choice of $p = 1.1$, slightly larger than 1, gives even better results than $p = 1$, see Figure 1.4(e). This might be due to the fact that the $L^{1.1}$ norm is smooth, and thus it encourages not only continuous reconstructions but also smooth ones. We conclude from Figure 1.4 that the use of non-quadratic penalties and in particular penalties based on Banach space norms may improve the quality of the reconstructions [204].

In [37], Banach spaces are not used in the definition of a sparsity regulariser but instead as solution spaces and the numerical tests, reported in Figure 1.5 show the effectiveness of this choice. As a last example, we consider an image deblurring problem for images in Figure 1.6 and in Figure 1.7. We compare here reconstruction obtained by solving (1.2) in an Hilbert L^2 setting with Landweber method (1.8) and

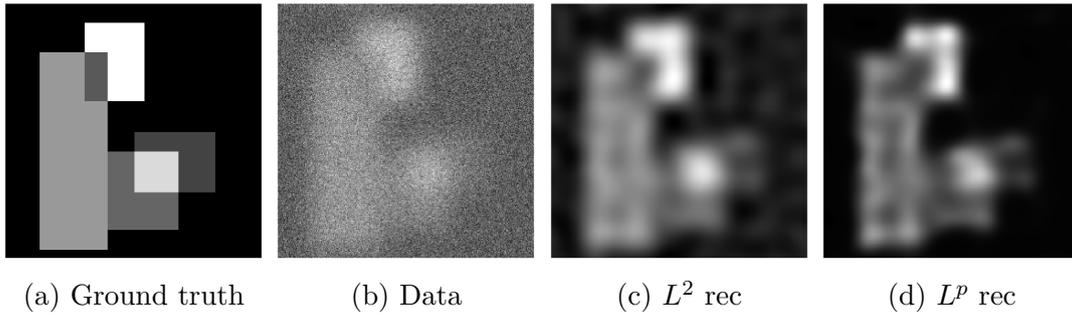


Figure 1.6: Exemplar deconvolution imaging problem on a test image: ground truth, acquired data and reconstructions obtained in L^2 and in L^p with $p = 1.3$ spaces.

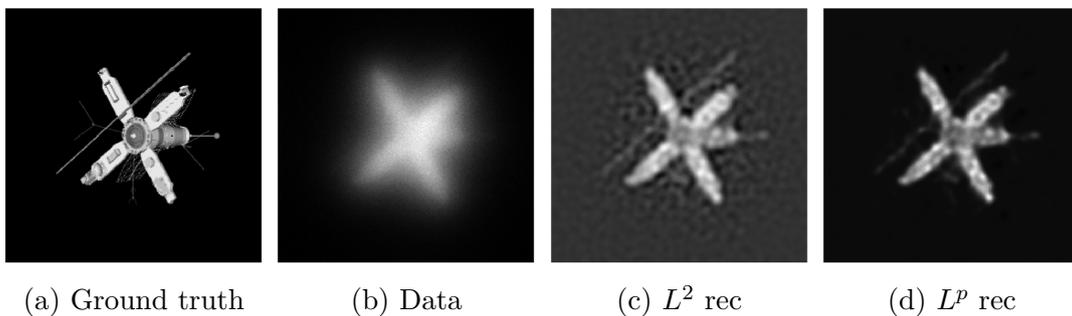


Figure 1.7: Exemplar deconvolution imaging problem on a satellite image: ground truth, acquired data and reconstructions obtained in L^2 and in L^p with $p = 1.3$ spaces.

in a Banach L^p setting ($p \in (1, 2)$) with the Banach Landweber method (that will be presented in Chapter 3). We can observe that in both images the Banach reconstructions present less artefacts in the background, which results cleaner. Moreover, especially in Figure 1.7, it is evident that considering a Banach space yields sharper solutions and better reconstructions of discontinuities.

Inverse problems in Banach spaces have been an extensive object of study. Several regularisation techniques and iterative methods, originally defined in a Hilbert setting, have been successfully extended and re-defined in Banach spaces, see for example [35, 41, 93, 112, 121, 151, 175, 192, 206, 207].

In this thesis, we often consider Banach spaces as solution space for inverse problems in imaging and study effective optimisation methods to minimise sparsity-promoting functionals arising from the variational formulation of the problem. In particular, this thesis focuses mainly on two particular Banach spaces, variable exponent Lebesgue spaces and the space of Radon measures, that are now briefly introduced.

1.4.1 Variable exponent Lebesgue spaces

Variable exponent Lebesgue spaces $L^{p(\cdot)}(\Omega)$, as name itself suggests, are Lebesgue spaces on $\Omega \subseteq \mathbb{R}^d$, $d \in \mathbb{N}$ defined in terms of a point-wise variable exponent $p(\cdot)$ instead of a constant one $p \in [1, +\infty]$. More precisely, they are defined in terms of a Lebesgue measurable function $p(\cdot) : \Omega \rightarrow [1, +\infty]$ that assigns coordinate-wise exponents to all points in the domain Ω . Under mild assumptions, they are Banach spaces intrinsically endowed with useful space-variant properties. A suitable variable exponent induces an adaptive regularisation in the reconstruction, having the possibility to enforce sparsity and preserve edges in certain parts of the domain, as when considering a constant exponent Lebesgue space $L^p(\Omega)$ with $1 < p < 2$ [204], and enforcing high L^2 regularity in smoother parts. Variable exponent Lebesgue spaces have proven to be useful in the design of adaptive regularisation with space-variant penalties terms to better deal with heterogeneous signals, i.e. signals with different properties in different parts of its domain. For example in [156] the so-called F-norm is defined in the discrete setting used as flexible penalty term that depends on a variable exponent $p(\cdot)$, point-wise defined on the domain. This penalty coincides to what we will refer to in the following as *modular* function. However, in [99, 156] a variable exponent is used in the definition of the penalty but $L^{p(\cdot)}(\Omega)$ spaces are not considered as solution spaces for the problem. This thesis proposes $L^{p(\cdot)}(\Omega)$ as solution spaces for inverse problems, showing their flexibility in the modelling of heterogeneous data and complex noise settings, to take advantage of the natural adaptivity of these spaces.

Previous works have considered $L^{p(\cdot)}(\Omega)$ spaces for applications. For instance, they have considered in the resolution of partial differential equations [36] and variational integrals with non-standard growth, see for example [2, 59, 66, 95, 103, 129]. In [153], the authors study a functional with variable exponent, which provides a model for image denoising, and its corresponding variable exponent heat equation. Variable exponent Lebesgue spaces have been used as solution spaces in electromagnetic non-linear inverse scattering, see [92], for noninvasive and nondestructive techniques to inspect materials, in [5] for microwave radiometer measurements, and in [26] for brain stroke microwave imaging.

1.4.2 Banach space of Radon measures

A research topic of particular interest in inverse problems is the recovery of sparse unknown signals, i.e. signals represented in terms of only a few non-zero components. This problem is commonly formulated in Hilbert and discrete settings by imposing an ℓ^0 -pseudonorm regularisation [47, 81, 83], defined as the number of non-zero components of x , $\|x\|_0 = \#\{x_i \text{ such that } x_i \neq 0\}$. However, this non-continuous and non-convex penalty makes variational problems NP-hard, and thus convex relaxations of ℓ^0 are considered in practice. The first one to mention is the ℓ^1 penalty, used in compressed sensing [90], which leads to the minimisation

of an $\ell^2 - \ell^1$ functional in the discrete setting [211]. To get rid of instabilities due to fine discretisations and to enhance the reconstruction precision without having to consider fine grids, the prior information on the sparsity of the signal can be characterised alternatively in the space of Radon measures $\mathcal{M}(\Omega)$, where a sparse signal $\mu \in \mathcal{M}(\Omega)$ can be modelled as a weighted sum of Dirac deltas, i.e.

$$\mu = \sum_{i=1}^N a_i \delta_{x_i}, \quad a_i \in \mathbb{R}, \quad x_i \in \Omega.$$

Sparse inverse problems in the space of Radon measures have been initially proposed in [39, 49, 72, 96] and studied further in many other works, see for example [34, 74, 85, 184]. In particular, in [39], the author establishes a framework for sparse inverse problems in general spaces of Radon measures and studies these spaces both as the solution space for linear inverse problems as well as the underlying space for numerical algorithms. The following variational model is therein proposed:

$$\operatorname{argmin}_{\mu \in \mathcal{M}(\Omega)} \frac{1}{2} \|T\mu - y\|^2 + \lambda |\mu|(\Omega), \quad (1.16)$$

where the penalty $|\mu|(\Omega)$ is the total-variation (TV) norm in the space of Radon measures $\mathcal{M}(\Omega)$, which is a Banach space endowed with such norm. The TV norm is a generalisation of the ℓ^1 one to the continuous setting, and, thus, (1.16) is considered an extension of the $\ell^2 - \ell^1$ model in the discrete setting. The framework of inverse problems in the space of Radon measures is often called *off-the-grid*, since the positions x_i of the Diracs can be anywhere in the domain Ω , which does not need to be discretised, as it happens when solving inverse problems in standard scenarios.

From an optimisation point of view, this setting is particularly challenging because $\mathcal{M}(\Omega)$ is a non-reflexive Banach space, and thus the scalar product and Riesz theorem's isomorphism are not available.

1.5 Optimisation methods

In the framework of variational models, the problem then becomes finding a minimiser of a structured functional as in (1.12), or of a smooth function as in (1.4).

In this section we refer to [54] and report some important definitions and optimisation algorithms used to tackle this minimisation problems in Hilbert spaces \mathcal{X} .

Let $f : \mathcal{X} \rightarrow \mathbb{R} \cup \{\pm\infty\}$ be an extended real valued function. It is said to be *convex* if and only if

$$\begin{aligned} x_1, x_2 \in \operatorname{dom}(f) &:= \{x \in \mathcal{X} \mid f(x) < +\infty\} \\ &\Downarrow \\ f(tx_1 + (1-t)x_2) &\leq tf(x_1) + (1-t)f(x_2) \quad \text{for all } t \in [0, 1]. \end{aligned} \quad (1.17)$$

If the above inequality is strict whenever $x_1 \neq x_2$ and $0 < t < 1$, f is *strictly convex*.

The function f is *proper* if it is not identically $+\infty$, i.e. there exists some $x \in \mathcal{X}$ such that $f(x) \neq +\infty$, and it is nowhere $-\infty$, i.e. for all $x \in \mathcal{X}$ we have $f(x) \neq -\infty$. In this case, f is convex if (1.17) holds for all $x_1, x_2 \in \mathcal{X}$.

The function f is *lower semi-continuous* (l.s.c.) if for all $x \in \mathcal{X}$

$$\text{if } x^k \rightarrow x \quad \Rightarrow \quad f(x) \leq \liminf_{k \rightarrow +\infty} f(x^k).$$

A well-known example of proper, convex and l.s.c. function is the indicator function of a convex set $C \subseteq \mathcal{X}$

$$\mathbb{1}_C(x) := \begin{cases} 0 & x \in C \\ +\infty & x \notin C \end{cases}.$$

It is really important in the framework of variational methods for inverse problems as it allows to easily model convex constraints, such as positivity constraints.

The set of $f : \mathcal{X} \rightarrow \mathbb{R} \cup \{+\infty\}$ proper, convex and l.s.c. extended real-valued functions is often denoted by $\Gamma_0(\mathcal{X})$.

For $f \in \Gamma_0(\mathcal{X})$, we recall the definition of *subdifferential* of f at $x \in \mathcal{X}$:

$$\partial f(x) := \{p \in \mathcal{X} \mid f(\tilde{x}) \geq f(x) + \langle p, \tilde{x} - x \rangle \text{ for all } \tilde{x} \in \mathcal{X}\}.$$

Any $p \in \partial f(x)$ is called *subgradient*. The above definition allows to generalise Fermat's stationary conditions for differentiable functions. Indeed, for non-smooth functions in $\Gamma_0(\mathcal{X})$ the following holds:

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

A function f is *strongly convex* with parameter μ if $x \mapsto f(x) - \frac{\mu}{2}\|x\|^2$ is convex. It is *coercive* if for all sequences $(x^k)_{k \in \mathbb{N}}$ such that $\|x^k\| \rightarrow +\infty$, we have $\lim_{k \rightarrow +\infty} f(x^k) = +\infty$.

1.5.1 Proximal operator and resolvent

A crucial role in optimisation is played by the so-called *proximal operator* [180] or *proximal map* of a convex function. Given $f \in \Gamma_0(\mathcal{X})$, it is defined as follows

$$\text{prox}_f(x) := \underset{u \in \mathcal{X}}{\text{argmin}} \frac{1}{2}\|x - u\|_{\mathcal{X}}^2 + f(u). \quad (1.18)$$

It is well-defined since the functional given by the sum of f and the squared norm term is strongly convex and proper, and thus has a unique minimiser. Subdifferential calculus computations [194] shows that

$$y = \text{prox}_f(x) \quad \Leftrightarrow \quad 0 \in \partial f(y) + (y - x) \quad \Leftrightarrow \quad y = (I + \partial f)^{-1}(x), \quad (1.19)$$

which shows that, from an operator perspective, $y = \text{prox}_f(x)$ is the *resolvent* of the maximal monotone operator ∂f at x [163].

1.5.2 Smooth optimisation methods

A first important class of optimisation methods is the family of smooth first-order methods, which are used to find a minimiser of a smooth convex function f , that is

$$\operatorname{argmin}_{x \in \mathcal{X}} f(x). \quad (1.20)$$

The most straightforward algorithm to solve (1.20) numerically is a gradient descent scheme

$$x^0 \in \mathcal{X}, \quad x^{k+1} = x^k - \tau \nabla f(x^k), \quad (1.21)$$

with a fixed step-size $\tau > 0$ such that $0 < \tau L < 2$ with L being the Lipschitz constant of $\nabla f(x)$, that is

$$\|\nabla f(x_1) - \nabla f(x_2)\|_{\mathcal{X}} \leq L \|x_1 - x_2\|_{\mathcal{X}} \quad \forall x_1, x_2 \in \mathcal{X}.$$

With $\tau < \frac{2}{L}$, the sequence $(f(x^k))_{k \in \mathbb{N}}$ is strictly decreasing, and, if in addition f is coercive it is possible to prove that $f(x^k)$ converges to a critical value of f . Moreover, in [173] the convergence rate in function values for the iterative scheme defined by (1.21) with $\tau < 2/L$ is proved and reads as

$$f(x^k) - f(x^*) \leq \frac{1}{2\tau k} \|x^* - x^0\|^2, \quad (1.22)$$

where x^* is any minimiser of f .

The rate of convergence (1.22) of order $\mathcal{O}(1/k)$ is sub-optimal and can be improved by considering Nesterov's acceleration [172]

$$\begin{aligned} t_{k+1} &= \frac{1 + \sqrt{1 + 4t_k^2}}{2} \\ y^k &= x^k + \frac{t_k - 1}{t_{k+1}}(x^k - x^{k-1}) \\ x^{k+1} &= y^k - \tau \nabla f(y^k). \end{aligned}$$

In this scenario, the rate of convergence with $\tau < 1/L$ reads

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

that is of order $\mathcal{O}(1/k^2)$. Such rate is interesting as it almost matches (up to constants) the worst case error bound of first order methods [172]. For more general results, see [107, 197].

1.5.3 Forward-backward splitting

Consider now minimisation problems in the more general form

$$\operatorname{argmin}_{x \in \mathcal{X}} f(x) + g(x), \quad (1.23)$$

where $g \in \Gamma_0(\mathcal{X})$ is possibly non-smooth and $f \in \Gamma_0(\mathcal{X})$ is a smooth function with L -Lipschitz gradient. To deal with such a minimisation problem, a possibility is to consider the forward-backward splitting algorithm, also called proximal gradient algorithm. They have been studied in a vast number of works, see for instance [13, 56, 63, 65, 155, 181]. The main idea of the forward-backward splitting scheme is to combine an explicit gradient descent step in the smooth part f , and an implicit step of descent on g . The iterations are defined by

$$x^0 \in \mathcal{X}, \quad x^{k+1} = \text{prox}_{\tau g} \left(x^k - \tau \nabla f(x^k) \right) \quad (1.24)$$

with $\tau < 2/L$. Using (1.19), iteration (1.24) can be equivalently defined as

$$x^{k+1} = \underset{x \in \mathcal{X}}{\text{argmin}} (I + \tau \partial g)^{-1} (I - \tau \nabla f)(x^k). \quad (1.25)$$

Forward-backward splitting has a rate of convergence of $\mathcal{O}(1/k)$. Similarly as in the case of gradient step, a rate of convergence of order $\mathcal{O}(1/k^2)$ has been proved in [13] with $\tau < 1/L$ for the modified scheme given by

$$\begin{aligned} t_{k+1} &= \frac{1 + \sqrt{1 + 4t_k^2}}{2} \\ y^k &= x^k + \frac{t_k - 1}{t_{k+1}} (x^k - x^{k-1}) \\ x^{k+1} &= \text{prox}_{\tau g} \left(y^k - \tau \nabla f(y^k) \right). \end{aligned}$$

This method is known with the name of FISTA, which stands for Fast Iterative Soft Thresholding Algorithm, being a fast version of ISTA, the forward-backward splitting (1.24) with $g(\cdot) = \|\cdot\|_1$, and the soft-thresholding being the proximal operator of $\|\cdot\|_1$.

Several other variants with/without accelerating strategies have been proposed. For instance in [173], the strong-convexity of the functions f and g is taken into account to improve convergence. In [217], strategies for inexact proximal computations are studied, since an explicit expression of the proximal operator is not always available. In [45], an adaptive backtracking strategy for the choice of the step-size of the algorithm studied in [55] is presented. In [191], the authors present a version of FISTA for strongly convex functions, with inexact proximal computations, the adaptive backtracking strategy proposed in [55], and a scaled metric is used in the definition of the proximal map. In [148], we applied this algorithm for some microscopy imaging problems. Proximal algorithms can be used in a myriad of situations. Extension of proximal algorithms to multicomponent signal and image recovery problems has been proposed in [42] and applied, for instance, to multispectral image denoising, and image decomposition into texture and geometry components. Proximal algorithms have also been defined to solve the minimisation problems of more than two terms that, i.e., arise when regularisation including several terms not necessarily acting in the same domain is considered to deal with complex noise models [185].

1.5.4 Mirror descent

A possible extension of proximal descent methods consist in replacing the distance induced by the L^2 norm in the definition of proximal operator (1.18) with other distances. There can be many reasons for this, such as the need to define a distance that blows up when approaching certain parts of the domain (and acts as a barrier or constraint), or the fact that the proximity operator of a function is not easily computable according to the standard definition (1.18) but it is simple in some non-quadratic metrics. Another possible reason, which we will investigate further in this thesis, is that the ambient space \mathcal{X} is neither Hilbertian nor Euclidean [171].

By using a different metric in the definition of (1.18), it is possible to characterise mirror descent, an algorithm suited for the minimisation problem (1.20). The distance function considered in mirror descent algorithms is the Bregman distance. Given a convex and smooth function $h : \mathcal{X} \rightarrow \mathbb{R}$, the Bregman distance $B_{\mathcal{X}}^h$ between $x, u \in \mathcal{X}$ is defined by

$$B_{\mathcal{X}}^h(u, x) = h(u) - \left(h(x) + \langle \nabla h(x), u - x \rangle \right), \quad \forall u \in \mathcal{X}.$$

The iterative scheme of mirror descent is defined in terms of the Bregman distance by

$$x^0 \in \mathcal{X}, \quad x^{k+1} \in \operatorname{argmin}_{x \in \mathcal{X}} B_{\mathcal{X}}^h(x, x^k) + \langle \tau \nabla f(x^k), x - x^k \rangle. \quad (1.26)$$

Iteration (1.26) can be equivalently formulated as

$$x^{k+1} \in \operatorname{argmin}_{x \in \mathcal{X}} B_{\mathcal{X}}^h(x, x^k) + f(x^k) + \langle \tau \nabla f(x^k), x - x^k \rangle \quad (1.27)$$

$$0 = \nabla h(x^{k+1}) - \nabla h(x^k) + \tau \nabla f(x^k)$$

$$\nabla h(x^{k+1}) = \nabla h(x^k) - \tau \nabla f(x^k). \quad (1.28)$$

From (1.27), we see that x^{k+1} corresponds to a minimiser of the linear approximation of f at x^k summed to the Bregman distance between the two points. Basically, mirror descent replaces the standard gradient descent (1.21) with (1.28), an unusual gradient step.

We refer the reader to the monograph [14] for many possible variants of mirror descent strategies, accompanied by their convergence rates.

1.6 Outline and contribution

In this thesis, we explore advanced smooth and non-smooth optimisation algorithms for imaging inverse problems in non-standard Banach spaces, which prove to be an effective and valid setting for the resolution of inverse problems. This thesis consists of two parts, corresponding to the two main different topics of research which has been investigated by the author during his Ph.D. studies, carried out as part of a joint PhD program between Università di Genova (Dipartimento di Matematica)

and Université Côte d’Azur (I3S Lab.-CNRS-INRIA (Morphéme team)) under the supervision of Luca Calatroni (UniCA) and Claudio Estatico (UniGE).

This thesis is based on the following publications:

- [147] [Stochastic Gradient Descent for Linear Inverse Problems in Variable Exponent Lebesgue Spaces](#). Lazzaretti M., Kereta Z., Estatico C., Calatroni L. In *Scale Space and Variational Methods in Computer Vision. SSVM 2023*. Lecture Notes in Computer Science, vol 14009. Springer, Cham. (2023)
- [31] [Dual descent regularization algorithms in variable exponent Lebesgue spaces for imaging](#). Bonino B., Estatico C., Lazzaretti M. Springer Numerical Algorithms Vol. 92. (2023)
- [145] [Modular-proximal gradient algorithms in variable exponent Lebesgue spaces](#). Lazzaretti M., Calatroni L., Estatico C. SIAM Journal on Scientific Computing Vol.44, Iss.6. (2022)
- [146] [Off-the-grid regularisation for Poisson inverse problems](#). Lazzaretti M., Estatico C., Melero A., Calatroni L. Submitted to Computational Optimization and Applications, Springer.

Other publications, not included in this thesis, are:

- [148] [A Scaled and Adaptive FISTA Algorithm for Signal-Dependent Sparse Image Super-Resolution Problems](#). Lazzaretti M., Rebegoldi S., Calatroni L., Estatico C. In *Scale Space and Variational Methods in Computer Vision - 8th International Conference, SSVM 2021*. Lecture Notes in Computer Science, vol 12679. Springer, Cham. (2021)
- [144] [Weighted-CEL0 sparse regularisation for molecule localisation in super-resolution microscopy with Poisson data](#). Lazzaretti M., Calatroni L., Estatico C. In *18th IEEE International Symposium on Biomedical Imaging, ISBI 2021*. (2021)
- [143] [A continuous, non-convex and sparse super-resolution approach for fluorescence microscopy data with Poisson noise](#). Lazzaretti M., Calatroni L., Estatico C. In *21st International Conference on Computational Science and its Applications, ICCSA 2021. IEEE CPS*. (2021)

The first part of this thesis focuses on optimisation strategies in general reflexive Banach spaces and proposes smooth and non-smooth algorithms in variable exponent Lebesgue spaces. The peculiarity of these spaces requires alternative algorithmic strategies, defined in terms of space specific functionals. The second part considers the non-reflexive Banach space of Radon measures as solution space for off-the-grid sparse imaging inverse problems, which are usually formulated with a Gaussian noise assumption. Here, a Poisson noise modelling is considered and a new variational formulation suited to this noise setting is studied both analytically and numerically.

Part I. Modular-based optimisation in variable exponents Lebesgue spaces.

This part is divided into three chapters. Chapter 2 gives a general introduction on variable exponent Lebesgue spaces $L^{p(\cdot)}(\Omega)$, reviewing the main definitions and properties and the role of the space-variant exponent $p(\cdot)$. A particular interest is given to the modular functions, since they are used to define the norm in $L^{p(\cdot)}(\Omega)$, and hence the space itself. The dual space of $L^{p(\cdot)}(\Omega)$ is analysed, underlying the differences with respect to the constant exponent case, and, in particular, an expression of the duality mappings in these spaces is explicitly computed. Our contribution is the proposal of a modular-based alternative to duality maps, motivated by the fact that duality maps in $L^{p(\cdot)}(\Omega)$ have some undesirable properties, i.e. they are not separable. A study on the separability of the modular concludes the chapter.

In Chapter 3, we present a review of gradient-based minimisation strategies in Banach spaces, namely primal and dual (Landweber) methods. We show how they can be equivalently formulated in terms of suitably defined proximal operators in Banach spaces, making a link between regularisation theory and convex optimisation. These algorithms, though, cannot be easily implemented in a variable exponent Lebesgue setting, due to the heavy-computations required by the non-separable norm and duality maps. Thus, we propose to replace the role played by duality maps by the modular-based alternative, analysed in Chapter 2, and we define a novel modular-based gradient descent algorithm in $L^{p(\cdot)}(\Omega)$. Numerical tests on simple image deblurring show the advantages of the proposed method with respect to the standard Landweber method for Banach spaces adapted to this scenario in terms of computational times. Another contribution in this chapter is the definition of a stochastic variant of the algorithm described above, which is again based on the modular, and it reduces significantly computational costs by Kaczmarz-type splitting of the problem. We validate the proposed methods with numerical experiments on CT reconstruction, both on simulated and real data.

Chapter 4 focuses on non-smooth optimisation in $L^{p(\cdot)}(\Omega)$. We first review forward-backward strategies in general Banach spaces, and then outline the reasons why they are hard to use in our variable exponent Lebesgue spaces modelling. We then propose two different modular-based proximal gradient algorithms and prove their convergence in function values, with rates. Exemplar 1D and 2D numerical tests prove that the spatial flexibility of $L^{p(\cdot)}(\Omega)$ spaces is particularly advantageous in addressing sparsity and heterogeneous signal/noise statistics, while remaining efficient and stable from an optimisation perspective.

Part II. Sparse optimisation in the Banach space of Radon measures with Poisson noise.

The second part of this thesis addresses optimisation in the space of Radon measures, a non-reflexive Banach space. This setting is particularly interesting for the so-called sparse *off-the-grid* methods in imaging. In Chapter 5, we review the main notions about the space of Radon measures $\mathcal{M}(\Omega)$. We recall the definition of Total-Variation norm of a measure and of its subdifferential, and present

the off-the-grid inverse problem formulation. In particular, we analyse the BLASSO variational problem, a generalisation of LASSO to this continuous setting, that is particularly suited to retrieve sparse signals, modelled in $\mathcal{M}(\Omega)$ as finite weighted sum of Diracs, under a Gaussian noise hypothesis. We conclude this chapter by reviewing some of the standard algorithms used for the resolution of the BLASSO problem, with a focus on Frank-Wolfe and Sliding Frank-Wolfe algorithms.

Our contribution for this second part of the thesis is outlined in Chapter 6. Since many off-the-grid methods are applied in microscopy imaging problems, assuming that the noise distribution is Gaussian is often unrealistic, due to the photon counting nature of the noise in this application. A more realistic assumption is Poisson noise, for which, following the Bayesian perspective, the Kullback-Leibler data term is the natural choice. We propose a novel variational model that couples the Kullback-Leibler divergence with the TV norm of measures and a non-negativity constraint, and provide a detailed theoretical analysis of its optimality conditions obtained by studying the corresponding dual problem. For its numerical resolution, we propose to consider the Sliding Frank-Wolfe algorithm, as for BLASSO, which is, however, quite sensitive to the choice of the regularisation parameter λ . To mitigate this undesirable effect, we propose an homotopy strategy for its automatic tuning. We validate the proposed Poisson off-the-grid model and compare it with BLASSO with several simulated numerical experiments, showing the effectiveness of the homotopy strategy considered. We conclude by showing the reconstruction of a real 3D fluorescence microscopy data obtained with the proposed Poisson off-the-grid model and the homotopy algorithm.

Part I

Modular-based optimisation in variable exponents Lebesgue spaces

Variable Exponent Lebesgue Spaces

In this chapter, the main definitions and properties about Lebesgue spaces with a variable exponent are introduced, with a particular focus on their norm and modular functions. In Section 2.2, the dual space of $L^{p(\cdot)}(\Omega)$ is discussed alongside the definition of duality mappings. In Section 2.3, we propose a modular-based alternative to duality mappings and study some of its properties, making a comparison with the standard definition.

2.1	Modular and Luxemburg norm	30
2.1.1	Inequalities between norm and modular	32
2.1.2	Properties of $L^{p(\cdot)}(\Omega)$ and immersions	35
2.2	Dual space and duality mappings	36
2.2.1	Definition of dual and associate space	36
2.2.2	Duality mappings	37
2.2.2.1	Duality mappings in $L^{p(\cdot)}(\Omega)$	38
2.2.2.2	Inverse of duality mappings	42
2.3	Modular-based alternative to duality maps	42
2.3.1	Separability	44
2.4	Final discussion	46

Let $\Omega \subseteq \mathbb{R}^d$, with $d \in \mathbb{N}$, $d \geq 1$, be a Lebesgue measurable subset with positive measure. In classical Lebesgue spaces $L^p(\Omega)$ with a constant exponent $p \in [1, +\infty)$, for any Lebesgue measurable function $x : \Omega \rightarrow \mathbb{R} \cup \{+\infty\}$, its norm is defined as

$$\|x\|_p = \left(\int_{\Omega} |x(t)|^p dt \right)^{1/p}$$

and, consequently, the Lebesgue space $L^p(\Omega)$ is naturally defined as the set of Lebesgue measurable functions with finite p -norm:

$$x \in L^p(\Omega) \iff \|x\|_p < +\infty.$$

In variable exponent Lebesgue spaces, as the name itself suggests, instead of a constant exponent $p \in [1, +\infty)$, a point-wise variable one is considered to measure the norm of elements and, hence, in the definition of the spaces. Any Lebesgue measurable function $p(\cdot) : \Omega \rightarrow [1, +\infty]$ can thus be taken as variable exponent. When considering a point-wise variable exponent $p(\cdot)$, the definition of the norm is not straightforward. Indeed, comparing to the above classical definition, it is only possible to compute the quantity

$$\int_{\Omega} |x(t)|^{p(t)} dt$$

but it is not clear which specific value of $p(\cdot)$ should be used to compute its radical. The computation of the radical, if possible, is crucial for the homogeneity property of the norm. In $L^{p(\cdot)}(\Omega)$ spaces, an alternative way to ensure homogeneity has to be considered. For this reason, the definition of norm in $L^{p(\cdot)}(\Omega)$ spaces has to be given in a different way; for that it is necessary to first introduce the so-called modular functions.

2.1 Modular and Luxemburg norm

Let the set of all possible exponents be

$$\mathcal{P}(\Omega) := \{p(\cdot) : \Omega \rightarrow [1, +\infty] \mid p(\cdot) \text{ is Lebesgue measurable}\}.$$

Given an exponent function $p(\cdot) \in \mathcal{P}(\Omega)$, the essential infimum and essential supremum of $p(\cdot)$ are denoted by

$$p_- := \operatorname{ess\,inf}_{u \in \Omega} p(u) \quad \text{and} \quad p_+ := \operatorname{ess\,sup}_{u \in \Omega} p(u).$$

Throughout this work, the following assumptions on the exponent are considered:

$$p_- > 1 \quad \text{and} \quad p_+ < +\infty. \tag{2.1}$$

As better specified in the next sections, under these hypotheses $L^{p(\cdot)}(\Omega)$ spaces satisfy important properties of Banach spaces. Moreover, if the set $\{t \in \Omega \mid p(t) = +\infty\}$ has positive measure, the definition of norm in $L^{p(\cdot)}(\Omega)$ spaces has a more complicated expression than the one provided in this chapter. Since in practical applications one always has $p_+ < +\infty$, we decided to limit the discussion on $L^{p(\cdot)}(\Omega)$ spaces to the hypothesis (2.1).

For the general case, see [68, 77] where a comprehensive dissertation about these spaces is carried out.

The characterisation of $L^{p(\cdot)}(\Omega)$ spaces is based on the key concept of modular function, whose definition is now reported. Let $\mathcal{F}(\Omega)$ be the set of all Lebesgue measurable functions $x : \Omega \rightarrow \mathbb{R} \cup \{+\infty\}$.

Definition 2.1.1. *Given an exponent $p(\cdot) \in \mathcal{P}(\Omega)$ with $p_+ < +\infty$, the functional $\rho_{p(\cdot)} : \mathcal{F}(\Omega) \rightarrow [0, +\infty]$ defined by*

$$\rho_{p(\cdot)}(x) = \int_{\Omega} |x(t)|^{p(t)} dt \quad (2.2)$$

is called modular associated to the exponent function $p(\cdot)$. An alternative definition of modular function consists in considering $\bar{\rho}_{p(\cdot)} : \mathcal{F}(\Omega) \rightarrow [0, +\infty]$ defined by

$$\bar{\rho}_{p(\cdot)}(x) = \int_{\Omega} \frac{1}{p(t)} |x(t)|^{p(t)} dt, \quad (2.3)$$

called normalised modular.

We will equivalently refer to (2.2) and (2.3) as modular functions, as needed.

Proposition 2.1.1. *[68, Proposition 2.7] Let $p(\cdot) \in \mathcal{P}(\Omega)$ and $x, y \in \mathcal{F}(\Omega)$. The following properties of the modular function $\rho_{p(\cdot)}(\cdot)$ holds true:*

- For all $x \in \mathcal{F}(\Omega)$, $\rho_{p(\cdot)}(x) \geq 0$ and $\rho_{p(\cdot)}(x) = \rho_{p(\cdot)}(|x|)$.
- $\rho_{p(\cdot)}(x) = 0$ if and only if $x(t) = 0$ for almost every $t \in \Omega$.
- If $\rho_{p(\cdot)}(x) < +\infty$, then $|x(t)| < +\infty$ for almost every $t \in \Omega$.
- $\rho_{p(\cdot)}$ is convex, that is

$$\rho_{p(\cdot)}(\alpha x + (1 - \alpha)y) \leq \alpha \rho_{p(\cdot)}(x) + (1 - \alpha) \rho_{p(\cdot)}(y), \quad 0 \leq \alpha \leq 1.$$

- $\rho_{p(\cdot)}$ is order preserving, i.e.

$$|x(t)| \geq |y(t)| \text{ a.e.} \Rightarrow \rho_{p(\cdot)}(x) \geq \rho_{p(\cdot)}(y).$$

As a consequence of the convexity property, it can be easily shown that given a scalar $\alpha > 0$, for any $x \in \mathcal{F}(\Omega)$:

- if $\alpha > 1$, then $\alpha \rho_{p(\cdot)}(x) \leq \rho_{p(\cdot)}(\alpha x)$;
- if $0 < \alpha < 1$, then $\rho_{p(\cdot)}(\alpha x) \leq \alpha \rho_{p(\cdot)}(x)$.

The above consideration shows that the modular does not satisfy the homogeneity property. Moreover, it is important to observe that the modular $\rho_{p(\cdot)}$ does not satisfy the triangle inequality either but instead it satisfies the following substitute. For any $x, y \in \mathcal{F}(\Omega)$:

$$\rho_{p(\cdot)}(x + y) \leq 2^{p_+ - 1} (\rho_{p(\cdot)}(x) + \rho_{p(\cdot)}(y)).$$

Notice that the modular $\rho_{p(\cdot)}(x)$ is the generalisation of the p -power of the norm $\|x\|_p^p = \int_{\Omega} |x(t)|^p dt$ in $L^p(\Omega)$ with constant exponent $p \in (1, +\infty)$. Similarly, the modular $\bar{\rho}_{p(\cdot)}(x)$ generalises the quantity $\frac{1}{p}\|x\|_p^p$. Modular functions are used to characterise the variable exponent space $L^{p(\cdot)}(\Omega)$ and to define its norm, in the general framework of the Luxemburg norms of Orlicz spaces [77].

Definition 2.1.2. *The space $L^{p(\cdot)}(\Omega)$ is the set of functions $x \in \mathcal{F}(\Omega)$ such that*

$$\rho_{p(\cdot)}\left(\frac{x}{\lambda}\right) \leq 1,$$

for some $\lambda > 0$. For any $x \in L^{p(\cdot)}(\Omega)$, we define $\|\cdot\|_{L^{p(\cdot)}} : L^{p(\cdot)}(\Omega) \longrightarrow \mathbb{R}$ as

$$\|x\|_{L^{p(\cdot)}} := \inf \left\{ \lambda > 0 : \rho_{p(\cdot)}\left(\frac{x}{\lambda}\right) \leq 1 \right\}. \quad (2.4)$$

Theorem 2.1.1. [68, Theorem 2.17] *The function $\|\cdot\|_{L^{p(\cdot)}}$ defined in (2.4) is a norm on the $L^{p(\cdot)}(\Omega)$. Moreover, the space $L^{p(\cdot)}(\Omega)$ endowed with such norm, that is the couple $(L^{p(\cdot)}(\Omega), \|\cdot\|_{L^{p(\cdot)}})$, is a Banach space.*

By extending the definition of the function $\|\cdot\|_{L^{p(\cdot)}}$ to any Lebesgue measurable function $x \in \mathcal{F}(\Omega)$ as follows

$$\|x\|_{L^{p(\cdot)}} = +\infty \quad \text{if} \quad \rho_{p(\cdot)}\left(\frac{x}{\lambda}\right) > 1 \text{ for any } \lambda > 0,$$

we retrieve the standard characterisation of Lebesgue spaces in terms of the norm, that is:

$$x \in L^{p(\cdot)}(\Omega) \iff \|x\|_{L^{p(\cdot)}} < +\infty.$$

The norm defined in (2.4) is often referred to as Luxemburg norm, after W. A. J. Luxemburg, who studied these concepts in his PhD thesis [157] in 1955. It can be considered as a more general definition of norm. Indeed, it is possible to retrieve the classical definition of norm $\|x\|_p$ in Lebesgue spaces $L^p(\Omega)$ with a constant exponent $p \in [1, +\infty)$. If $p(\cdot) \equiv p \in [1, +\infty)$, for any $\lambda > 0$

$$\rho_{p(\cdot)}\left(\frac{x}{\lambda}\right) = \rho_p\left(\frac{x}{\lambda}\right) = \frac{1}{\lambda^p} \rho_p(x) = \frac{1}{\lambda^p} \|x\|_p^p,$$

so that the infimum in (2.4) is equal to $\|x\|_p$.

2.1.1 Inequalities between norm and modular

It may seem that the Luxemburg norm and the modular are somewhat interchangeable but, if it is true to a certain extent that they have similar properties, it is also important to point out that they are truly different objects.

	$\rho_{p(\cdot)}(x)$	$\rho_{p(\cdot)}(x)^{1/p_-}$	$\rho_{p(\cdot)}(x)^{1/p_+}$	$\ x\ _{L^{p(\cdot)}}$	$\ x\ _{p_-}$	$\ x\ _{p_+}$	$\tilde{p}(x)$
Fig. 1.6a	1.3365	0.5635	0.0885	0.1275	0.6125	0.0864 · 10 ⁴	1.3238
Fig. 1.6b	1.2235	0.5200	0.0831	0.1325	0.59274	0.0785 · 10 ⁴	1.3092
Fig. 1.7a	0.0570	0.0739	0.1292	0.1152	0.4544	0.0844	1.3255
Fig. 1.7b	0.0635	0.0816	0.1395	0.1114	0.3962	0.0565	1.2561

Table 2.1: Comparison between the values of Luxemburg norm, modular and classical p -norms for x being the image in Figures 1.6a, 1.6b and Figures 1.7a, 1.7b.

A first very important difference is that the modular function does not satisfy homogeneity property, i.e. $\rho_{p(\cdot)}(\alpha x) \neq \alpha \rho_{p(\cdot)}(x)$ for any $\alpha \in \mathbb{R}$, that is satisfied by the norm $\|\alpha x\|_{L^{p(\cdot)}} = |\alpha| \|x\|_{L^{p(\cdot)}}$, as expected. In the classical case of constant Lebesgue spaces $L^p(\Omega)$, the computation of the p -radical of the modular $\rho_p(\cdot)$ ensures that the homogeneity property is satisfied by the p -norm. With a variable exponent, such computation is obviously not possible and in turn the one-dimensional (1D) minimisation problem (2.4) has to be solved. However, the Luxemburg norm is bounded by the p_- and p_+ radicals of the modular.

Lemma 2.1.1. [77, Lemma 3.2.5, Lemma 3.4.2] *Let $p(\cdot) \in \mathcal{P}(\Omega)$ with $p_+ < +\infty$.*

1. *If $\|x\|_{L^{p(\cdot)}} > 1$, then $\rho_{p(\cdot)}(x)^{1/p_+} \leq \|x\|_{L^{p(\cdot)}} \leq \rho_{p(\cdot)}(x)^{1/p_-}$.*
2. *If $0 < \|x\|_{L^{p(\cdot)}} \leq 1$, then $\rho_{p(\cdot)}(x)^{1/p_-} \leq \|x\|_{L^{p(\cdot)}} \leq \rho_{p(\cdot)}(x)^{1/p_+}$.*

At a certain extent, the norm can be viewed as $\tilde{p}(x)$ -radical (which depends on x) of the modular with $p_- \leq \tilde{p}(x) \leq p_+$, defined by

$$\tilde{p}(x) = \frac{\log(\rho_{p(\cdot)}(x))}{\log(\|x\|_{L^{p(\cdot)}})},$$

so that $\|x\|_{L^{p(\cdot)}} = (\rho_{p(\cdot)})^{1/\tilde{p}(x)}$.

As an example of Luxemburg norm computation, we report in Table 2.1 the values of the variable exponent norm for the ground truths and acquired images of Figures 2.1 and 2.2. In particular, the variable exponents used are shown in Figure 2.1c and Figure 2.2c respectively. The reason why they resemble the acquired solution will be explained in Chapter 3, Section 3.3.4 and it is not relevant for the comparison between norm and modular computation here discussed. In Table 2.1, it is also shown a comparison with the classical p -norms for $p = p_- = 1.1$ and $p = p_+ = 1.4$, and with the modular and its p_- and p_+ radicals. Moreover, the value $\tilde{p}(x)$ is reported for the tested images. It is important to consider the computational time needed for the computation of the Luxemburg norm. For the image in Figures 2.1a and 2.1b the CPU time amounts to 4.2s and 9.1s respectively, while for the

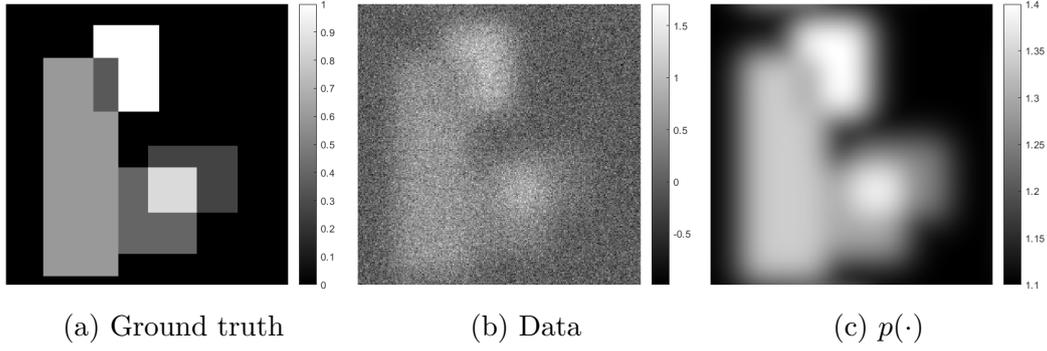


Figure 2.1: Test image: ground truth, acquired data, variable exponent.

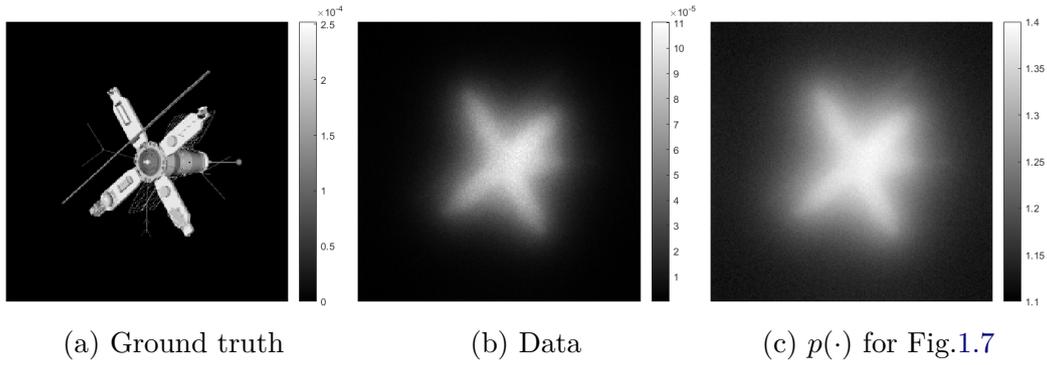


Figure 2.2: Satellite image: ground truth, acquired data, variable exponent.

satellite images in Figures 2.2a and 2.2b the CPU time is 2.2s and 10.4s respectively. On the other hand, the computation of the modular is very fast and amounts to almost 0s of CPU time.

Another property that is worth mentioning is the fact that the unit ball computed with respect to the $L^{p(\cdot)}(\Omega)$ norm is equivalent to the unit ball computed with the modular $\rho_{p(\cdot)}$.

Proposition 2.1.2. [77, Lemma 2.1.14] *Let $x \in L^{p(\cdot)}(\Omega)$ with $p(\cdot) \in \mathcal{P}(\Omega)$. Then:*

- $\|x\|_{L^{p(\cdot)}} < 1$ and $\rho_{p(\cdot)}(x) < 1$ are equivalent.
- $\|x\|_{L^{p(\cdot)}} = 1$ and $\rho_{p(\cdot)}(x) = 1$ are equivalent.

Moreover, as one can deduce from Lemma 2.1.1, the unit ball splits the space into two regions, one where the norm is smaller than the modular and the other where the vice-versa holds.

Proposition 2.1.3. [77, Lemma 3.2.4] *Let $x \in L^{p(\cdot)}(\Omega)$ with $p(\cdot) \in \mathcal{P}(\Omega)$. Then:*

1. If $\|x\|_{L^{p(\cdot)}} \leq 1$, then $\rho_{p(\cdot)}(x) \leq \|x\|_{L^{p(\cdot)}}$.
2. If $\|x\|_{L^{p(\cdot)}} > 1$, then $\rho_{p(\cdot)}(x) \geq \|x\|_{L^{p(\cdot)}}$.

2.1.2 Properties of $L^{p(\cdot)}(\Omega)$ and immersions

Some important properties about $L^{p(\cdot)}(\Omega)$ spaces are reported here. It is important to stress that the hypothesis made on the exponent function (2.1) are fundamental to prove the following properties for $L^{p(\cdot)}(\Omega)$.

Recall that for a Banach space \mathcal{X} , the notation $\langle x^*, x \rangle_{\mathcal{X}^* \times \mathcal{X}}$, or simply $\langle x^*, x \rangle$, is used to indicate the duality product, defined as

$$\langle x^*, x \rangle_{\mathcal{X}^* \times \mathcal{X}} = x^*(x), \quad \forall x^* \in \mathcal{X}^*, x \in \mathcal{X}.$$

It can be considered as a generalisation of the scalar product in Hilbert spaces. Moreover, a Banach space \mathcal{X} is:

- i) reflexive if $(\mathcal{X}^*)^* = \mathcal{X}$;
- ii) uniformly convex if for any $\varepsilon \in (0, 2]$, the inequalities $\|x\|_{\mathcal{X}} \leq 1, \|y\|_{\mathcal{X}} \leq 1$ and $\|x - y\|_{\mathcal{X}} \geq \varepsilon$ imply there exists a $\delta = \delta(\varepsilon) > 0$ such that $\|(x + y)/2\|_{\mathcal{X}} \leq 1 - \delta$;
- iii) strictly convex if for any $x, y \in \mathcal{X}$ such that $\|x\|_{\mathcal{X}} = \|y\|_{\mathcal{X}} = 1$ and $x \neq y$ there holds $\|(x + y)/2\|_{\mathcal{X}} < 1$;
- iv) smooth if, for every $x \neq 0$, there exists a unique $x^* \in \mathcal{X}^*$ such that $\|x^*\|_{\mathcal{X}^*} = 1$ and $\langle x^*, x \rangle = \|x\|_{\mathcal{X}}$.

Theorem 2.1.2. [77, Theorem 3.4.7] *Given $p(\cdot)$ such that $1 < p_- \leq p_+ < +\infty$, then $L^{p(\cdot)}(\Omega)$ is reflexive.*

Theorem 2.1.3. [77, Theorem 3.4.9] *Given $p(\cdot)$ such that $1 < p_- \leq p_+ < +\infty$, then $L^{p(\cdot)}(\Omega)$ is uniformly convex, and hence strictly convex.*

Theorem 2.1.4. [79, Lemma 1] *Given $p(\cdot)$ such that $1 < p_- \leq p_+ < +\infty$, then $L^{p(\cdot)}(\Omega)$ is smooth.*

As a last result of this section, we focus on immersions between any two Lebesgue spaces with variable exponents. To prove the following, one needs the boundedness of the domain Ω .

Proposition 2.1.4. [68, Corollary 2.48] *Given $p(\cdot) \in \mathcal{P}(\Omega)$ and $q(\cdot) \in \mathcal{P}(\Omega)$, with $p_+ < +\infty$ and $q_+ < +\infty$, if the domain Ω is bounded, then the following natural immersion holds*

$$L^{q(\cdot)}(\Omega) \hookrightarrow L^{p(\cdot)}(\Omega)$$

if and only if

$$p(t) \leq q(t) \quad \text{a.e. in } \Omega.$$

Moreover, it holds $\|x\|_{L^{p(\cdot)}} \leq (1 + |\Omega|) \|x\|_{L^{q(\cdot)}}$ being $|\Omega| < +\infty$ the finite Lebesgue measure of the bounded domain.

In particular, for any $p(\cdot) \in \mathcal{P}(\Omega)$ with $p_+ < +\infty$ the following inclusions follow from the previous proposition:

$$L^{p_+}(\Omega) \subseteq L^{p(\cdot)}(\Omega) \subseteq L^{p_-}(\Omega).$$

2.2 Dual space and duality mappings

In this section, the definition of dual spaces in the context of variable exponent Lebesgue spaces is analysed and a definition of duality mapping in $L^{p(\cdot)}(\Omega)$ spaces is provided, highlighting the differences between variable and constant exponent spaces. In particular, differing from the constant exponent case, an isometric isomorphism between $(L^{p(\cdot)}(\Omega))^*$ and $L^{p'(\cdot)}(\Omega)$ does not hold true in general.

For a comprehensive review of these arguments, see [77].

2.2.1 Definition of dual and associate space

Definition 2.2.1. Let $G : L^{p(\cdot)}(\Omega) \rightarrow \mathbb{R}$ be a linear functional. G is bounded if

$$\sup_{u \in L^{p(\cdot)}(\Omega), \|u\|_{L^{p(\cdot)}} \leq 1} |G(u)| < +\infty.$$

The dual space of $L^{p(\cdot)}(\Omega)$ can thus be defined as the set of linear and bounded functionals from $L^{p(\cdot)}(\Omega)$ to \mathbb{R} :

$$(L^{p(\cdot)}(\Omega))^* = \{G : L^{p(\cdot)}(\Omega) \rightarrow \mathbb{R} : G \text{ is linear and bounded}\},$$

which is a Banach space with the norm

$$\|G\|_{(L^{p(\cdot)}(\Omega))^*} := \sup_{u \in L^{p(\cdot)}(\Omega), \|u\|_{L^{p(\cdot)}} \leq 1} |G(u)|.$$

For $1 < p_- \leq p_+ < +\infty$, the Hölder conjugate of $p(\cdot)$ is a Lebesgue measurable function $p'(\cdot) \in \mathcal{P}(\Omega)$ such that

$$\frac{1}{p(t)} + \frac{1}{p'(t)} = 1 \text{ a.e. in } \Omega. \quad (2.5)$$

With the notation p' the conjugate of a constant exponent p will consistently be denoted. For a variable exponent $p(\cdot)$, the operations of taking the infimum and supremum do not commute with forming the conjugate exponent. In fact, the following relation holds:

$$(p'(\cdot))_- = (p_+)', \quad (p'(\cdot))_+ = (p_-)'. \quad (2.6)$$

To avoid ambiguous expressions, omitting one set of parentheses, the following notation will be adopted: $(p'(\cdot))_- = (p')_-$ and $(p'(\cdot))_+ = (p')_+$.

With a strong formal analogy to the constant exponent case, for any $z \in L^{p'(\cdot)}(\Omega)$, it can be shown, see [77], that there exists a unique $G_z \in (L^{p(\cdot)}(\Omega))^*$ such that

$$G_z(u) = \int_{\Omega} z(t)u(t)dt \quad \forall u \in L^{p(\cdot)}(\Omega). \quad (2.7)$$

Thus, the mapping $G_z \in (L^{p(\cdot)}(\Omega))^* \mapsto z \in L^{p'(\cdot)}(\Omega)$ is injective. In the following the duality pairing notation $G_z(x) = \langle z, x \rangle$ or $G_z(x) = \langle G_z, x \rangle$ will be thus adopted.

Definition 2.2.2. [77, Definition 2.7.1] *The associate space of $L^{p(\cdot)}(\Omega)$, denoted by $\mathcal{A}(L^{p(\cdot)}(\Omega))$, is the space of functions $z \in L^{p(\cdot)}(\Omega)$ such that*

$$\sup \left\{ \int_{\Omega} |z(t)||u(t)|dt : u \in L^{p(\cdot)}(\Omega), \|u\|_{L^{p(\cdot)}} \leq 1 \right\} < +\infty. \quad (2.8)$$

The function $\|\cdot\|'_{p(\cdot)} : \mathcal{A}(L^{p(\cdot)}(\Omega)) \rightarrow \mathbb{R}$ defined by

$$\|z\|'_{p(\cdot)} := \sup \left\{ \int_{\Omega} |z(t)||u(t)|dt : u \in L^{p(\cdot)}(\Omega), \|u\|_{L^{p(\cdot)}} \leq 1 \right\}$$

is a norm on $\mathcal{A}(L^{p(\cdot)}(\Omega))$.

First of all, observe that $\mathcal{A}(L^{p(\cdot)}(\Omega))$ might happen to be a proper subset of $L^{p(\cdot)}(\Omega)$. The supremum in (2.8) is finite if and only if the linear operator G_z is bounded according to Definition 2.2.1. Moreover, the operator norm of $G_z \in (L^{p(\cdot)}(\Omega))^*$ and the associate norm of $z \in \mathcal{A}(L^{p(\cdot)}(\Omega))$ coincides:

$$\begin{aligned} \|G_z\|_{(L^{p(\cdot)}(\Omega))^*} &= \sup \left\{ |G_z(u)| : u \in L^{p(\cdot)}(\Omega), \|u\|_{L^{p(\cdot)}} \leq 1 \right\} = \\ &= \sup \left\{ \int_{\Omega} |z(t)||u(t)|dt : u \in L^{p(\cdot)}(\Omega), \|u\|_{L^{p(\cdot)}} \leq 1 \right\} \\ &= \|z\|'_{p(\cdot)}. \end{aligned} \quad (2.9)$$

Thus, there exists an isometric embedding $\mathcal{A}(L^{p(\cdot)}(\Omega)) \hookrightarrow (L^{p(\cdot)}(\Omega))^*$ between the associate space and the dual space of $L^{p(\cdot)}(\Omega)$. However, the following proposition combined with (2.9) shows that $L^{p(\cdot)}(\Omega)$ and $(L^{p(\cdot)}(\Omega))^*$ are not isometrically isomorphic. Due to this, the expression of the inverse of duality mappings remains unknown, as it will be better explained in the following section.

Proposition 2.2.1. [77, Corollary 3.2.14] *For all $z \in \mathcal{A}(L^{p(\cdot)}(\Omega))$, there holds*

$$\frac{1}{2}\|z\|_{L^{p(\cdot)}} \leq \|z\|'_{p(\cdot)} \leq 2\|z\|_{L^{p(\cdot)}},$$

and the bounds are optimal.

2.2.2 Duality mappings

Before analysing the concept of duality mappings in $L^{p(\cdot)}(\Omega)$ spaces, the general definition of duality mapping for a Banach space \mathcal{X} is here reported.

Definition 2.2.3. [60] *Let \mathcal{X} be a Banach space and let a scalar $r > 1$. Then the duality map $\mathbf{J}_{\mathcal{X}}^r$ with gauge function $t \mapsto t^{r-1}$ is the operator $\mathbf{J}_{\mathcal{X}}^r : \mathcal{X} \rightarrow 2^{\mathcal{X}^*}$ such that*

$$\mathbf{J}_{\mathcal{X}}^r(x) = \left\{ x^* \in \mathcal{X}^* \mid x^*(x) = \langle x^*, x \rangle = \|x\|_{\mathcal{X}} \|x^*\|_{\mathcal{X}^*}, \|x^*\|_{\mathcal{X}^*} = \|x\|_{\mathcal{X}}^{r-1} \right\} \quad \forall x \in \mathcal{X}.$$

In general, duality maps are multi-valued operators. However, if (and only if) the Banach space \mathcal{X} is smooth the duality map $\mathbf{J}_{\mathcal{X}}^r$ is single valued, that is $\mathbf{J}_{\mathcal{X}}^r : \mathcal{X} \rightarrow \mathcal{X}^*$ a function between the space \mathcal{X} and its dual \mathcal{X}^* . In addition, it is important to point out that the only case when the duality map reduces to the identity operator is if \mathcal{X} is a Hilbert space \mathcal{H} and the gauge function has parameter $r = 2$. By virtue of the Riesz theorem, the duality map becomes $\mathbf{J}_{\mathcal{H}}^2(x) = x$, where the isometric isomorphism between \mathcal{H} and \mathcal{H}^* has been implicitly considered. On the other hand, when choosing $r \neq 2$, in the Hilbert setting, it means an unusual metric is considered.

In general, the following result gives a more intuitive characterisation of duality maps, providing a practical way to analytically compute them.

Theorem 2.2.1 ((Asplund) [60]). *Let \mathcal{X} be a Banach space, and let $r > 1$. The r -duality map $\mathbf{J}_{\mathcal{X}}^r$ is the subdifferential of the convex functional $h : \mathcal{X} \rightarrow \mathbb{R}$, $h(x) = \frac{1}{r} \|x\|_{\mathcal{X}}^r$:*

$$\mathbf{J}_{\mathcal{X}}^r = \partial h = \partial \left(\frac{1}{r} \|\cdot\|_{\mathcal{X}}^r \right).$$

Duality maps satisfy important properties, see [60] for details. A particular property is that being the subdifferential of a convex functional, $\mathbf{J}_{\mathcal{X}}^r$ is a monotone operator, that is

$$\langle \mathbf{J}_{\mathcal{X}}^r(x) - \mathbf{J}_{\mathcal{X}}^r(y), x - y \rangle \geq 0, \quad \forall x, y \in \mathcal{X}, \quad \forall r > 1$$

and it can be proven that $\mathbf{J}_{\mathcal{X}}^r(-x) = -\mathbf{J}_{\mathcal{X}}^r(x)$ and $\mathbf{J}_{\mathcal{X}}^r(\lambda x) = \lambda^{r-1} \mathbf{J}_{\mathcal{X}}^r(x)$ for any $x \in \mathcal{X}$ and for $\lambda \geq 0$.

2.2.2.1 Duality mappings in $L^{p(\cdot)}(\Omega)$

In [78, 79] the authors proved that the Luxemburg norm $\|\cdot\|_{L^{p(\cdot)}}$ is Gateaux-differentiable for any exponent $p(\cdot)$ such that $1 < p_- \leq p_+ < +\infty$, providing an analytical expression for its Gateaux derivative. This shows that hence the space $(L^{p(\cdot)}(\Omega), \|\cdot\|_{L^{p(\cdot)}})$ is smooth. In addition, in [159, 160], it is shown that the norm in $L^{p(\cdot)}(\Omega)$ is Fréchet differentiable too, for any $x \neq 0$. From these results, it follows that the functional $\frac{1}{r} \|\cdot\|_{L^{p(\cdot)}}^r$ for $r > 1$ is Fréchet differentiable for any $x \in \mathcal{X}$. Following arguments similar to those of [78, 79], we provided in [31] the analytical expression for the duality mapping $\mathbf{J}_{L^{p(\cdot)}}^r$, Gateaux derivative of $\frac{1}{r} \|\cdot\|_{L^{p(\cdot)}}^r$.

Theorem 2.2.2. *Let the exponent function $p(\cdot) \in \mathcal{P}(\Omega)$ be such that $1 < p_- \leq p_+ < +\infty$. Then, for each $x \in L^{p(\cdot)}(\Omega)$ and for any $r \in (1, +\infty)$, the duality mapping $\mathbf{J}_{L^{p(\cdot)}}^r : L^{p(\cdot)}(\Omega) \rightarrow (L^{p(\cdot)}(\Omega))^*$ is the linear operator with expression*

$$\langle \mathbf{J}_{L^{p(\cdot)}}^r(x), h \rangle = \frac{1}{\int_{\Omega} \frac{p(t)|x(t)|^{p(t)}}{\|x\|_{L^{p(\cdot)}}^{p(t)}} dt} \int_{\Omega} \frac{p(t) \operatorname{sign}(x(t)) |x(t)|^{p(t)-1}}{\|x\|_{L^{p(\cdot)}}^{p(t)-r}} h(t) dt \quad (2.10)$$

for any $h \in L^{p(\cdot)}(\Omega)$.

Proof. This proof is based on arguments similar to those of [78, 79], for the proof of Gateaux-differentiability of the Luxembourg norm in $L^{p(\cdot)}(\Omega)$ spaces. We presented this result in [31].

By Theorem 2.2.1, we know that $\mathbf{J}_{L^{p(\cdot)}}^r = \partial\left(\frac{1}{r}\|\cdot\|_{L^{p(\cdot)}}^r\right)$. Taking into account the smoothness of $(L^{p(\cdot)}(\Omega), \|\cdot\|_{L^{p(\cdot)}})$, in the following of the proof we will focus on the computation of the Gâteaux derivative the functional $x \in L^{p(\cdot)}(\Omega) \mapsto \|x\|_{L^{p(\cdot)}}^r$ (without the fixed scaling factor $\frac{1}{r}$ for simplicity), for any $x_0 \in L^{p(\cdot)}(\Omega)$.

We now first consider $x_0 \neq 0$. We have to prove that, for any possible direction $h \in L^{p(\cdot)}(\Omega)$, the real function $\sigma \mapsto \|x_0 + \sigma h\|_{L^{p(\cdot)}}^r$, with $\sigma \in \mathbb{R}$, is differentiable at $\sigma = 0$. We will use the implicit function theorem as follows.

Let $k > 1$ be a fixed real number, $D = (-1, 1) \times \left(\frac{1}{k}\|x_0\|_{L^{p(\cdot)}}^r, k\|x_0\|_{L^{p(\cdot)}}^r\right)$ and consider the function $\phi : D \rightarrow \mathbb{R}$ defined by means of the convex modular function $\rho_{p(\cdot)}$ which characterises the Luxembourg norm (2.4)

$$\phi(\sigma, \lambda) = \rho_{p(\cdot)}\left(\frac{x_0 + \sigma h}{\lambda^{1/r}}\right) - 1 = \int_{\Omega} \frac{|x_0(t) + \sigma h(t)|^{p(t)}}{\lambda^{p(t)/r}} dt - 1. \quad (2.11)$$

In the sequel, we will demonstrate the following statements, which are the hypothesis of the implicit function Theorem:

- i) $\phi \in C^1(D)$;
- ii) $\phi(0, \|x_0\|_{L^{p(\cdot)}}^r) = 0$;
- iii) $\frac{\partial \phi}{\partial \lambda}(0, \|x_0\|_{L^{p(\cdot)}}^r) < 0$.

Indeed, once proven i), ii) and iii), the implicit function Theorem guarantees that there exist neighbourhoods U of 0 and V of $\|x_0\|_{L^{p(\cdot)}}^r$ such that $U \times V \subset D$ and a unique C^1 -mapping $\lambda : U \rightarrow V$ which satisfies $\lambda(0) = \|x_0\|_{L^{p(\cdot)}}^r$, $\phi(\sigma, \lambda(\sigma)) = 0$ for any $\sigma \in U$, and

$$\lambda'(\sigma) = -\frac{\frac{\partial \phi}{\partial \sigma}(\sigma, \lambda(\sigma))}{\frac{\partial \phi}{\partial \lambda}(\sigma, \lambda(\sigma))}, \quad \forall \sigma \in U. \quad (2.12)$$

The equality $\phi(\sigma, \lambda(\sigma)) = 0 \forall \sigma \in U$, rewritten as

$$\rho_{p(\cdot)}\left(\frac{x_0 + \sigma h}{\lambda(\sigma)^{1/r}}\right) = 1, \quad \forall \sigma \in U,$$

together with Definition 2.1.2 of the norm in $L^{p(\cdot)}(\Omega)$, allows us to derive that

$$\lambda(\sigma) = \|x_0 + \sigma h\|_{L^{p(\cdot)}}^r, \quad \forall \sigma \in U. \quad (2.13)$$

Hence, from (2.12) and (2.13) we have that $\lambda'(0)$ exists and

$$\lambda'(0) = \lim_{\sigma \rightarrow 0} \frac{\|x_0 + \sigma h\|_{L^{p(\cdot)}}^r - \|x_0\|_{L^{p(\cdot)}}^r}{\sigma} = -\frac{\frac{\partial \phi}{\partial \sigma}(0, \|x_0\|_{L^{p(\cdot)}}^r)}{\frac{\partial \phi}{\partial \lambda}(0, \|x_0\|_{L^{p(\cdot)}}^r)}. \quad (2.14)$$

The functional $\|\cdot\|_{L^{p(\cdot)}}^r$ is thus Gâteaux differentiable at $x_0 \neq 0$, and the explicit computation of the ratio (2.14) will provide expressions (2.10) too.

We can now prove the statements i), ii) and iii).

i) To prove that $\phi \in C^1(D)$, let us consider the integrand $f : \Omega \times D \rightarrow \mathbb{R}$ of (2.11)

$$f(t; (\sigma, \lambda)) = \frac{|x_0(t) + \sigma h(t)|^{p(t)}}{\lambda^{p(t)/r}}, \quad t \in \Omega, \quad (\sigma, \lambda) \in D. \quad (2.15)$$

It is easy to show that, for any fixed $(\sigma, \lambda) \in D$, the map $t \mapsto f(t; (\sigma, \lambda))$ is integrable in Ω . Indeed, by definition of D , there hold $|\sigma| < 1$ and $\lambda \geq \frac{1}{k} \|x_0\|_{L^{p(\cdot)}}^r = \lambda_{\min} > 0$, which yields to

$$\frac{|x_0(t) + \sigma h(t)|^{p(t)}}{\lambda^{p(t)/r}} \leq \frac{k^{p(t)/r} (|x_0(t)| + |h(t)|)^{p(t)}}{\|x_0\|_{L^{p(\cdot)}}^{p(t)}} \leq \frac{k^{p_+/r}}{c} (|x_0(t)| + |h(t)|)^{p(t)}$$

with $c = \min \left(\|x_0\|_{L^{p(\cdot)}}^{p_-}, \|x_0\|_{L^{p(\cdot)}}^{p_+} \right)$ and $(|x_0(t)| + |h(t)|)^{p(t)}$ being integrable since $x_0, h \in L^{p(\cdot)}(\Omega)$ and $p_+ < +\infty$. Consequently, the function ϕ of (2.11) is well-defined.

We now show that for a.e. $t \in \Omega$, the map $(\sigma, \lambda) \mapsto f(t; (\sigma, \lambda))$, with $(\sigma, \lambda) \in D$, is a C^1 -mapping. By formal computation, the partial derivatives of (2.15) are

$$\frac{\partial f}{\partial \sigma}(t; (\sigma, \lambda)) = \frac{p(t) |x_0(t) + \sigma h(t)|^{p(t)-1} \text{sign}(x_0(t) + \sigma h(t)) h(t)}{\lambda^{p(t)/r}} \quad (2.16)$$

$$\frac{\partial f}{\partial \lambda}(t; (\sigma, \lambda)) = -\frac{p(t) |x_0(t) + \sigma h(t)|^{p(t)}}{s \lambda^{p(t)/r+1}}, \quad \forall (\sigma, \lambda) \in D. \quad (2.17)$$

Since $p_- > 1$ and $\lambda > 0$, it is evident from (2.16) and (2.17), that $(\sigma, \lambda) \mapsto \frac{\partial f}{\partial \sigma}(t; (\sigma, \lambda))$ and $(\sigma, \lambda) \mapsto \frac{\partial f}{\partial \lambda}(t; (\sigma, \lambda))$ are continuous mappings in D . Anyway, to explicitly compute both the numerator and the denominator of (2.12), that is, the partial derivatives of (2.11), we need to commute differentiation and integration operators. To this aim, we apply the Dominated Convergence Theorem, by searching for a function $g : \Omega \rightarrow \mathbb{R}$, integrable on Ω , such that

$$\left| \frac{\partial f}{\partial \sigma}(t; (\sigma, \lambda)) \right| \leq g(t), \quad \left| \frac{\partial f}{\partial \lambda}(t; (\sigma, \lambda)) \right| \leq g(t).$$

Similarly as before for the estimation of $|f(t; (\sigma, \lambda))|$, $(\sigma, \lambda) \in D$ implies that

$$\begin{aligned} \left| \frac{\partial f}{\partial \sigma}(t; (\sigma, \lambda)) \right| &\leq \frac{k^{p(t)/r} p(t) (|x_0(t)| + |h(t)|)^{p(t)}}{\|x_0\|_{L^{p(\cdot)}}^{p(t)}} \\ &\leq \frac{p_+ \cdot k^{p_+/r}}{c} (|x_0(t)| + |h(t)|)^{p(t)}, \end{aligned}$$

with $c = \min \left(\|x_0\|_{L^{p(\cdot)}}^{p_-}, \|x_0\|_{L^{p(\cdot)}}^{p_+} \right)$, and that

$$\left| \frac{\partial f}{\partial \lambda}(t; (\sigma, \lambda)) \right| \leq \frac{p_+ \cdot k^{p_+/r+1}}{c_1} (|x_0(t)| + |h(t)|)^{p(t)}$$

with $c_1 = \min \left(\|x_0\|_{L^{p(\cdot)}}^{p_-+r}, \|x_0\|_{L^{p(\cdot)}}^{p_++r} \right)$. Thus, we can now consider

$$g(t) = \max \left(\frac{p_+ \cdot k^{p_+/r}}{c}, \frac{p_+ \cdot k^{p_+/r+1}}{c_1} \right) (|x_0(t)| + |h(t)|)^{p(t)},$$

as dominating function, which is integrable on Ω , as already stated before. Hence differentiation and integration in ϕ commute, leading to

$$\begin{aligned} \frac{\partial \phi}{\partial \sigma}(\sigma, \lambda) &= \frac{\partial}{\partial \sigma} \left[\int_{\Omega} f(t; (\sigma, \lambda)) dt - 1 \right] = \int_{\Omega} \left[\frac{\partial}{\partial \sigma} f(t; (\sigma, \lambda)) \right] dt \\ &= \int_{\Omega} p(t) \frac{|x_0(t) + \sigma h(t)|^{p(t)-1} \text{sign}(x_0(t) + \sigma h(t))}{\lambda^{p(t)/r}} h(t) dt \end{aligned} \quad (2.18)$$

$$\begin{aligned} \frac{\partial \phi}{\partial \lambda}(\sigma, \lambda) &= \frac{\partial}{\partial \lambda} \left[\int_{\Omega} f(t; (\sigma, \lambda)) dt - 1 \right] = \int_{\Omega} \left[\frac{\partial}{\partial \lambda} f(t; (\sigma, \lambda)) \right] dt \\ &= - \int_{\Omega} \frac{p(t) \cdot |x_0(t) + \sigma h(t)|^{p(t)}}{r \lambda^{p(t)/r+1}} dt. \end{aligned} \quad (2.19)$$

From (2.18) and (2.19), the continuity of $\frac{\partial \phi}{\partial \sigma}$ and $\frac{\partial \phi}{\partial \lambda}$ is straightforward.

ii) By Definition 2.1.2 of Luxemburg norm,

$$\phi(0, \|x_0\|_{L^{p(\cdot)}}^r) = \int_{\Omega} \left| \frac{x_0(t)}{\|x_0\|_{L^{p(\cdot)}}} \right|^{p(t)} dt - 1 = 0.$$

iii) We have similarly

$$\begin{aligned} \frac{\partial \phi}{\partial \lambda}(0, \|x_0\|_{L^{p(\cdot)}}^r) &= - \int_{\Omega} p(t) \frac{|x_0(t)|^{p(t)}}{\|x_0\|_{L^{p(\cdot)}}^{p(t)+r}} dt \\ &\leq - \frac{p_-}{\|x_0\|_{L^{p(\cdot)}}^r} \int_{\Omega} \left| \frac{x_0(t)}{\|x_0\|_{L^{p(\cdot)}}} \right|^{p(t)} dt = - \frac{p_-}{\|x_0\|_{L^{p(\cdot)}}^r} < 0. \end{aligned}$$

We can now obtain formula (2.10) from (2.14), plugging $\sigma = 0$ into (2.18) and (2.19). Recall that we are computing the gradient of $\|\cdot\|_{L^{p(\cdot)}}^r$, thus a multiplication by $\frac{1}{r}$ is needed to obtain (2.10).

To conclude the proof, it remains to consider the case $x_0 = 0$. The real function $\sigma \mapsto \|x_0 + \sigma h\|_{L^{p(\cdot)}}^r$, with $\sigma \in \mathbb{R}$, becomes $\sigma \mapsto \|\sigma h\|_{L^{p(\cdot)}}^r$. We have

$$\lim_{\sigma \rightarrow 0} \frac{\|x_0 + \sigma h\|_{L^{p(\cdot)}}^r - \|x_0\|_{L^{p(\cdot)}}^r}{\sigma} = \lim_{\sigma \rightarrow 0} \frac{\|\sigma h\|_{L^{p(\cdot)}}^r}{\sigma} = \lim_{\sigma \rightarrow 0} \frac{|\sigma|^r \|h\|_{L^{p(\cdot)}}^r}{\sigma} = 0,$$

since $r > 1$, which proves the differentiability at the origin as well. \square

From (2.10), it is possible to retrieve the expression of the Gateaux derivative of the Luxemburg norm $\|\cdot\|_{L^{p(\cdot)}}$ of [78, 79], by plugging $r = 1$. In addition, it is easy to check that, if $p(\cdot) \equiv p$ is constant, with $1 < p < +\infty$, then $\mathbf{J}_{L^{p(\cdot)}}^r$ coincides with the duality map of constant exponent Lebesgue spaces $L^p(\Omega)$:

$$\langle \mathbf{J}_{L^p}^r(x), h \rangle = \|x\|_p^{r-p} \int_{\Omega} \text{sign}(x(t)) |x(t)|^{p-1} h(t) dt. \quad (2.20)$$

2.2.2.2 Inverse of duality mappings

In general, for any reflexive and strictly convex Banach space \mathcal{X} , $\mathbf{J}_{\mathcal{X}}^r$ is invertible and its inverse is given by

$$(\mathbf{J}_{\mathcal{X}}^r)^{-1} = \mathbf{J}_{\mathcal{X}^*}^{r'}, \quad (2.21)$$

where $\mathbf{J}_{\mathcal{X}^*}^{r'} : \mathcal{X}^* \rightarrow \mathcal{X}$ is the duality mapping of the dual space \mathcal{X}^* .

In $L^p(\Omega)$ spaces, the isometric isomorphism between $(L^p(\Omega))^*$ and $L^{p'}(\Omega)$ is fundamental to have an explicit expression for the inverse of the duality map $\mathbf{J}_{L^p}^r$. In fact, it leads to

$$(\mathbf{J}_{L^p}^r)^{-1} = \mathbf{J}_{(L^p)^*}^{r'} = \mathbf{J}_{L^{p'}}^{r'}, \quad (2.22)$$

so that we have an analytical expression of the inverse duality map.

In $L^{p(\cdot)}(\Omega)$ spaces, such isomorphism does not hold true, as shown in Proposition 2.2.1, so this fact cannot be used to obtain the inverse of $\mathbf{J}_{L^{p(\cdot)}}^r$, whose analytical expression remains unknown.

2.3 Modular-based alternative to duality maps

In this section, a modular-based alternative to duality maps is presented and some of its properties are shown and proved, see also [145]. The main idea is the fact that the computation of the Luxemburg norm is computationally expensive, since a one-dimensional minimisation problem has to be solved for each norm computation, and, moreover, the norm is not separable, in the meaning specified in the following. These undesirable characteristics are reflected also in the duality map $\mathbf{J}_{L^{p(\cdot)}}^r$ defined in Theorem 2.2.2, since its expression requires a norm computation. Inspired by Theorem 2.2.1 and recalling that the modular is a generalisation of the p -power of the p -norm in $L^p(\Omega)$, in [145] we propose to substitute duality maps with the subdifferential of the modular functions.

Proposition 2.3.1. *Let $p(\cdot) \in \mathcal{P}(\Omega)$ be such that $p_- > 1$ and $p_+ < +\infty$. For each $x \in L^{p(\cdot)}(\Omega)$, the modular functions $\rho_{p(\cdot)}(\cdot)$ and $\bar{\rho}_{p(\cdot)}(\cdot)$ are Gateaux differentiable. Their derivatives are respectively the linear operators $\mathbf{J}_{\rho_{p(\cdot)}} : L^{p(\cdot)}(\Omega) \rightarrow (L^{p(\cdot)}(\Omega))^*$ defined by*

$$\langle \mathbf{J}_{\rho_{p(\cdot)}}(x), h \rangle = \int_{\Omega} p(t) \operatorname{sign}(x(t)) |x(t)|^{p(t)-1} h(t) dt$$

and $\mathbf{J}_{\bar{\rho}_{p(\cdot)}} : L^{p(\cdot)}(\Omega) \rightarrow (L^{p(\cdot)}(\Omega))^*$ with expression

$$\langle \mathbf{J}_{\bar{\rho}_{p(\cdot)}}(x), h \rangle = \int_{\Omega} \operatorname{sign}(x(t)) |x(t)|^{p(t)-1} h(t) dt,$$

for any $h \in L^{p(\cdot)}(\Omega)$.

Proof. Let $x, h \in L^{p(\cdot)}(\Omega)$. The Gateaux derivative of $\rho_{p(\cdot)}$ at x along direction h is given by

$$\lim_{t \rightarrow 0} \frac{\rho_{p(\cdot)}(x + th) - \rho_{p(\cdot)}(x)}{t} = \left[\frac{\partial}{\partial t} \rho_{p(\cdot)}(x + th) \right]_{t=0}. \quad (2.23)$$

First, note that $\rho_{p(\cdot)}(x + th) = \int_{\Omega} |x(s) + th(s)|^{p(s)} ds < +\infty$. Indeed $x, h \in L^{p(\cdot)}(\Omega)$, then $x + th \in L^{p(\cdot)}(\Omega)$ and consequently $\|x + th\|_{L^{p(\cdot)}} < +\infty$ by definition of $L^{p(\cdot)}(\Omega)$ space. Since $p_+ < +\infty$ and by Lemma 2.1.1, $\|x + th\|_{L^{p(\cdot)}} < +\infty$ implies $\rho_{p(\cdot)}(x + th) < +\infty$. We can thus compute the partial derivative of $\rho_{p(\cdot)}(x + th)$ with respect to t by “differentiating under the integral sign.” To do so, we first need to verify the regularity of the integrand function. Let $f : \Omega \times (-1, 1) \rightarrow \mathbb{R}$ be defined by

$$f(s, t) := |x(s) + th(s)|^{p(s)}, \quad s \in \Omega, \quad t \in (-1, 1).$$

By direct computations, we obtain

$$\frac{\partial f}{\partial t}(s, t) = p(s)|x(s) + th(s)|^{p(s)-1} \operatorname{sign}(x(s) + th(s))h(s) \quad (2.24)$$

and, since $|t| < 1$,

$$\begin{aligned} \left| \frac{\partial f}{\partial t}(s, t) \right| &\leq p_+ |x(s) + th(s)|^{p(s)-1} |h(s)| \leq p_+ |x(s) + th(s)|^{p(s)-1} (|h(s)| + |x(s)|) \\ &\leq p_+ (|h(s)| + |x(s)|)^{p(s)} =: g(s), \end{aligned}$$

with $g(s)$ integrable.

Thanks to the dominated convergence theorem, we have

$$\frac{\partial}{\partial t} \rho_{p(\cdot)}(x + th) = \frac{\partial}{\partial t} \int_{\Omega} f(s, t) ds = \int_{\Omega} \frac{\partial}{\partial t} f(s, t) ds. \quad (2.25)$$

Thus, by combining (2.23), (2.24) and (2.25) we conclude the proof. \square

Remark 2.3.1. *We stress that although $\mathbf{J}_{\rho_{p(\cdot)}}$ and $\mathbf{J}_{\bar{\rho}_{p(\cdot)}}$ are not duality mappings, we nonetheless adopt a similar notation for consistency.*

It is interesting to observe the following property of the modular function and its gradient.

Lemma 2.3.1. *For any $x \in L^{p(\cdot)}(\Omega)$,*

$$\langle \mathbf{J}_{\bar{\rho}_{p(\cdot)}}(x), x \rangle = \rho_{p(\cdot)}(x). \quad (2.26)$$

Proof. By direct computation:

$$\langle \mathbf{J}_{\bar{\rho}_{p(\cdot)}}(x), x \rangle = \int_{\Omega} \operatorname{sign}(x(t)) |x(t)|^{p(t)-1} x(t) dt = \rho_{p(\cdot)}(x). \quad \square$$

Note that this is the analogue of a general property of duality mappings in Banach spaces. Indeed, by Definition 2.2.3 of duality mapping, if \mathcal{X} is a smooth Banach space, then for any $r > 1$ and $x \in \mathcal{X}$ there holds

$$\langle \mathbf{J}_{\mathcal{X}}^r(x), x \rangle = \|x\|_{\mathcal{X}} \|x^*\|_{\mathcal{X}^*} = \|x\|_{\mathcal{X}} \|x\|_{\mathcal{X}}^{r-1} = \|x\|_{\mathcal{X}}^r. \quad (2.27)$$

Differently from the duality maps $\mathbf{J}_{L^{p(\cdot)}}^r$, it is possible to compute explicitly the inverse of their modular-based alternative $\mathbf{J}_{\rho_{p(\cdot)}}$ and $\mathbf{J}_{\bar{\rho}_{p(\cdot)}}$, as we first showed in [147].

Proposition 2.3.2. *The functional $\mathbf{J}_{\bar{\rho}_{p(\cdot)}}$ is invertible. For all $G_z \in (L^{p(\cdot)}(\Omega))^*$ with $z \in \mathcal{A}(L^{p'(\cdot)}(\Omega))$, its inverse reads*

$$(\mathbf{J}_{\bar{\rho}_{p(\cdot)}})^{-1}(G_z) = |z|^{\frac{1}{p(\cdot)-1}} \text{sign}(z) \in L^{p(\cdot)}(\Omega). \quad (2.28)$$

Proof. By using the expression of $(\mathbf{J}_{\bar{\rho}_{p(\cdot)}})^{-1}$ (2.28), we show that both $(\mathbf{J}_{\bar{\rho}_{p(\cdot)}})^{-1}\mathbf{J}_{\bar{\rho}_{p(\cdot)}}$ and $\mathbf{J}_{\bar{\rho}_{p(\cdot)}}(\mathbf{J}_{\bar{\rho}_{p(\cdot)}})^{-1}$ are the identity operator.

First of all, observe that for any $x \in L^{p(\cdot)}(\Omega)$ the action of the linear operator $\mathbf{J}_{\bar{\rho}_{p(\cdot)}}$ can be expressed as

$$\mathbf{J}_{\bar{\rho}_{p(\cdot)}}(x) = G_z \in (L^{p(\cdot)}(\Omega))^*, \quad z = \text{sign}(x)|x|^{p(\cdot)-1} \in \mathcal{A}(L^{p'(\cdot)}(\Omega)),$$

with G_z defined as in (2.7). With a slight abuse of notation, we write here $\mathbf{J}_{\bar{\rho}_{p(\cdot)}}(x) = \text{sign}(x)|x|^{p(\cdot)-1}$. By straightforward computation, we have

$$\begin{aligned} \left(\mathbf{J}_{\bar{\rho}_{p(\cdot)}}\right)^{-1}\left(\mathbf{J}_{\bar{\rho}_{p(\cdot)}}(x)\right) &= \left|\mathbf{J}_{\bar{\rho}_{p(\cdot)}}(x)\right|^{\frac{1}{p(\cdot)-1}} \text{sign}\left(\mathbf{J}_{\bar{\rho}_{p(\cdot)}}(x)\right) \\ &= \left|\mathbf{J}_{\bar{\rho}_{p(\cdot)}}(x)\right|^{\frac{1}{p(\cdot)-1}-1} \mathbf{J}_{\bar{\rho}_{p(\cdot)}}(x) = \left|\mathbf{J}_{\bar{\rho}_{p(\cdot)}}(x)\right|^{\frac{2-p(\cdot)}{p(\cdot)-1}} \mathbf{J}_{\bar{\rho}_{p(\cdot)}}(x) \\ &= \left|x\right|^{p(\cdot)-1} \text{sign}(x) \left|x\right|^{\frac{2-p(\cdot)}{p(\cdot)-1}} \left|x\right|^{p(\cdot)-1} \text{sign}(x) = x. \end{aligned}$$

It can also be shown that $\mathbf{J}_{\bar{\rho}_{p(\cdot)}}\left(\left(\mathbf{J}_{\bar{\rho}_{p(\cdot)}}\right)^{-1}(G_z)\right) = G_z$:

$$\begin{aligned} \mathbf{J}_{\bar{\rho}_{p(\cdot)}}\left(\left(\mathbf{J}_{\bar{\rho}_{p(\cdot)}}\right)^{-1}(G_z)\right) &= \mathbf{J}_{\bar{\rho}_{p(\cdot)}}\left(\left|z\right|^{\frac{1}{p(\cdot)-1}} \text{sign}(z)\right) \\ &= \left\langle \left|\left|z\right|^{\frac{1}{p(\cdot)-1}} \text{sign}(z)\right|^{p(\cdot)-1} \text{sign}\left(\left|z\right|^{\frac{1}{p(\cdot)-1}} \text{sign}(z)\right), \cdot \right\rangle \\ &= \langle |z| \text{sign}(z), \cdot \rangle = \langle z, \cdot \rangle = G_z. \end{aligned}$$

□

Proposition 2.3.3. *The functional $\mathbf{J}_{\rho_{p(\cdot)}}$ is invertible. For all $G_z \in (L^{p(\cdot)}(\Omega))^*$ with $z \in \mathcal{A}(L^{p'(\cdot)}(\Omega))$, its inverse reads*

$$(\mathbf{J}_{\rho_{p(\cdot)}})^{-1}(G_z) = p(\cdot)^{\frac{1}{1-p(\cdot)}} |z|^{\frac{1}{p(\cdot)-1}} \text{sign}(z) \in L^{p(\cdot)}(\Omega).$$

2.3.1 Separability

As it will be more evident through the rest of this work and, in particular, better explained in Sections 3.4.1 and 4.4, it is handy having functionals and operators defined in Banach spaces that are *separable*, that is, their global computation can be decomposed into the sum of low-dimensional functionals. This fact usually allows component-wise computations, which are easier and more efficient. To be more precise, we consider the following definition of domain additive separability.

Definition 2.3.1. Let \mathcal{X} be a functional Banach space on Ω . An operator $S : \mathcal{X} \rightarrow \mathcal{X}^*$ or a functional $S : \mathcal{X} \rightarrow \mathbb{R}$ is domain additively separable, if, for any finite family of Lebesgue measurable subsets $(\Omega_i)_{i=1}^n$ of Ω such that $\dot{\Omega}_i \cap \dot{\Omega}_j = \emptyset$ for $i \neq j$, and $\Omega = \bigcup_{i=1}^n \Omega_i$, there holds $S(x) = \sum_{i=1}^n S(\chi_i x)$ for any $x \in \mathcal{X}$, where $\chi_i \in \mathcal{X}$ is the characteristic function of Ω_i , that is $\chi_i(t) = 1$ for $t \in \Omega_i$, and $\chi_i(t) = 0$ for $t \notin \Omega_i$.

In the following, domain additive separability will often be referred to simply as separability. It is quite evident that in Lebesgue spaces $L^p(\Omega)$ with a constant exponent, the norm functional $\|\cdot\|_p^p$, as well as the p -duality map $\mathbf{J}_p^p(\cdot)$, are domain additively separable, since, for any suitable family of subsets $(\Omega_i)_{i=1}^n$, there holds

$$\|x\|_p^p = \sum_{i=1}^n \|\chi_i x\|_p^p \quad \text{and} \quad \langle \mathbf{J}_p^p(x), u \rangle = \sum_{i=1}^n \langle \mathbf{J}_p^p(\chi_i x), u \rangle = \langle \sum_{i=1}^n \mathbf{J}_p^p(\chi_i x), u \rangle.$$

On the contrary, norms and duality maps in variable exponent spaces are not separable.

Lemma 2.3.2. The norm and the duality mapping in $L^{p(\cdot)}(\Omega)$ are not domain additively separable in the sense of Definition 2.3.1.

Proof. It is quite evident that the Luxemburg norm (2.4) requires the solution of a 1D minimization problem on the entire domain Ω , which, in general, cannot be divided into the solutions on single sets of the partition, that is, $\|x\|_{L^{p(\cdot)}} \neq \sum_{i=1}^n \|\chi_i x\|_{L^{p(\cdot)}}$. As the duality mapping is concerned, the two norms in the denominators of $\mathbf{J}_{L^{p(\cdot)}}$ (2.10) show that its computation cannot be decomposed into the computation of n integrals involving only the restriction of the function x onto single sets of the partition, or, in other words, $\mathbf{J}_{L^{p(\cdot)}(\Omega)}^r(x) \neq \sum_{i=1}^n \mathbf{J}_{L^{p(\cdot)}(\Omega)}^r(\chi_i x)$. \square

The modular functions introduced in Definition 2.2 as well as their gradients turn out instead to satisfy the separability property.

Lemma 2.3.3. The modular functions in $L^{p(\cdot)}(\Omega)$ and their gradients are domain additively separable, in the sense of Definition 2.3.1.

Proof. We consider the modular function $\bar{\rho}_{p(\cdot)}(x) = \int_{\Omega} \frac{1}{p(t)} |x(t)|^{p(t)} dt$ defined in (2.3) (for (2.2), the proof is similar). By direct computation, by the linearity property of the integral w.r.t. the integration domain, we have

$$\bar{\rho}_{p(\cdot)}(x) = \sum_{i=1}^n \int_{\Omega_i} \frac{1}{p(t)} |x(t)|^{p(t)} dt = \sum_{i=1}^n \int_{\Omega} \frac{1}{p(t)} |\chi_i(t)x(t)|^{p(t)} dt = \sum_{i=1}^n \bar{\rho}_{p(\cdot)}(\chi_i x).$$

Similarly, for $\mathbf{J}_{\bar{\rho}_{p(\cdot)}}$, we can write

$$\begin{aligned} \langle \mathbf{J}_{\bar{\rho}_{p(\cdot)}}(x), u \rangle &= \int_{\Omega} \text{sign}(x(t)) |x(t)|^{p(t)-1} u(t) dt = \sum_{i=1}^n \int_{\Omega_i} \text{sign}(x(t)) |x(t)|^{p(t)-1} u(t) dt \\ &= \sum_{i=1}^n \int_{\Omega} \text{sign}(x(t)) |\chi_i(t)x(t)|^{p(t)-1} u(t) dt = \sum_{i=1}^n \langle \mathbf{J}_{\bar{\rho}_{p(\cdot)}}(\chi_i x), u \rangle = \langle \sum_{i=1}^n \mathbf{J}_{\bar{\rho}_{p(\cdot)}}(\chi_i x), u \rangle \end{aligned}$$

which concludes the proof. \square

2.4 Final discussion

In this chapter, an overview on variable exponents Lebesgue spaces is given. $L^{p(\cdot)}(\Omega)$ spaces, as the name suggests, are defined in terms of a point-wise variable exponent function instead of a constant one, as in standard $L^p(\Omega)$ spaces. We mainly focused on the definition of the norm (2.4) in this context and of the so-called modular functions. The norm does not have a closed form expression and its computation requires the resolution of a 1-dimensional minimisation problem, and, thus, it is quite computationally expensive, as we have seen in the provided examples. Modular functions, defined in (2.2) and (2.3), on the other hand, maintain some of the properties of the norm but having an analytical closed form expression their computation is much faster than the norm's one.

Variable exponent Lebesgue spaces endowed with the Luxemburg norm are Banach spaces. In Section (2.2.2), we thus defined duality mappings, which allow to associate to any element of the primal space an element of the dual space (and vice-versa), firstly for a general Banach space \mathcal{X} , and then in the specific case of $L^{p(\cdot)}(\Omega)$. The duality map $\mathbf{J}_{L^{p(\cdot)}}^r$ is defined in (2.10): it requires norms computations and it is thus heavy to compute. Moreover, the lack of an isometric isomorphism between $L^{p(\cdot)}(\Omega)$ and $L^{p'(\cdot)}(\Omega)$ causes the impossibility to have an expression for the inverse of the duality map in $L^{p(\cdot)}(\Omega)$, using standard strategies. To the best of our knowledge, its expression remains unknown. Recalling that duality maps can be seen, thanks to the Asplund's Theorem 2.2.1, as the subdifferential of the norm (elevated to some constant r), we proposed a modular-based alternative to duality mappings in Section 2.3, defined as the derivatives of the modulars (instead of the norm). They do not require any norm computation and they are invertible with inverses given in Propositions 2.3.2 and 2.3.3.

In the last Section, we focused on the concept of separability (Definition 2.3.1) showing that the norm and duality maps are not additively separable in $L^{p(\cdot)}(\Omega)$ while modular functions and their derivatives are separable.

In the following, we will consider $L^{p(\cdot)}(\Omega)$ spaces as solution spaces for inverse problems and we will define optimisation strategies in this setting, making use of the modular functions and their derivatives instead of the norm and duality mappings, having the latter some undesirable properties to devise minimisation algorithms (non-separability, heavy-computational times, lack of an exact expression for the inverse of duality maps).

Smooth optimisation in $L^{p(\cdot)}(\Omega)$

In this chapter, we analyse gradient-based minimisation strategies, which depend on the definition of a gradient descent step, such as the Landweber algorithm, from a different point of view, making a link between regularisation theory and convex optimisation in Banach spaces. Specifically, we interpret Landweber algorithms in the context of proximal methods, defined in terms of an appropriate distance function. This novel reinterpretation allows a full understanding of the role of the geometrical properties of the Banach spaces \mathcal{X} and \mathcal{Y} . A modular-based version of gradient descent in variable exponent Lebesgue spaces is then presented, both in the deterministic and stochastic settings.

3.1	Proximal operators in Banach spaces	48
3.1.1	Definition of the p-norm proximal operator	49
3.1.2	Bregman-proximal operator	50
3.2	Landweber methods in Banach spaces	52
3.2.1	Hilbert spaces setting	52
3.2.2	Dual method in Banach spaces	54
3.2.3	Primal method in Banach spaces	56
3.2.4	Primal and dual method as proximal point algorithms	56
3.3	Modular-based dual method in $L^{p(\cdot)}(\Omega)$	57
3.3.1	Primal and dual methods in $L^{p(\cdot)}(\Omega)$: main issues	58
3.3.1.1	Approximation of the inverse of the duality map	58
3.3.1.2	Heavy computational times	59
3.3.2	Modular-based alternative to dual method in $L^{p(\cdot)}(\Omega)$	59
3.3.3	Comparison between Landweber and modular-based gradient descent	60
3.3.4	How to choose variable exponents	62
3.3.5	Numerical tests with modular-based gradient descent	64
3.4	A Modular-based Stochastic variant	66
3.4.1	Stochastic Gradient Descent in Banach spaces	67

3.4.2	Variable exponents modular-based SGD	69
3.4.3	Numerical results	70
3.4.3.1	Hyper-parameter selection	70
3.4.3.2	Simulated data	71
3.4.3.3	Real CT dataset	73
3.5	Final discussion	74

We start introducing the concept of proximal operators in Banach spaces, as they will be used to provide a new interpretation of gradient-based iterative strategies in this setting.

3.1 Proximal operators in Banach spaces

Let \mathcal{X} be a reflexive, smooth and strictly convex Banach space and let $\Gamma_0(\mathcal{X})$ be the set of proper, lower semi-continuous, convex functions $g : \mathcal{X} \rightarrow \mathbb{R} \cup \{+\infty\}$. In general, a proximal operator of a possibly non-smooth function $g \in \Gamma_0(\mathcal{X})$ depends on a chosen distance function d . Given $d : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$ such that $d(x, \cdot) \in \Gamma_0(\mathcal{X})$ for any $x \in \mathcal{X}$, the proximal operator of g with respect to d is defined as

$$\text{prox}_g^d(x) = \underset{u \in \mathcal{X}}{\text{argmin}} d(x, u) + g(u). \quad (3.1)$$

The functional $h(\cdot) = \inf_{u \in \mathcal{X}} d(\cdot, u) + g(u)$ can be seen as a regularised version of $g(\cdot)$, which share the same infimum as g , since:

$$\inf_{x \in \mathcal{X}} h(x) = \inf_{x \in \mathcal{X}} \inf_{u \in \mathcal{X}} \left(d(x, u) + g(u) \right) = \inf_{u \in \mathcal{X}} \inf_{x \in \mathcal{X}} \left(d(x, u) + g(u) \right) = \inf_{u \in \mathcal{X}} g(u).$$

This shows they play a huge role in devising minimisation algorithms, as better explained by the following results.

Proposition 3.1.1. *Let $g \in \Gamma_0(\mathcal{X})$. Then for all $x \in \mathcal{X}$*

$$g\left(\text{prox}_g^d(x)\right) \leq g(x).$$

Moreover, $g\left(\text{prox}_g^d(x)\right) = g(x) \iff x \in \text{prox}_g^d(x)$.

Proof. Let $x \in \mathcal{X}$ and $z \in \text{prox}_g^d(x)$. Define $h_x(u) := d(x, u) + g(u)$ so that $z \in \text{argmin}_{u \in \mathcal{X}} h_x(u)$ and, in particular, $h_x(z) = \inf_{u \in \mathcal{X}} h_x(u)$. Thus

$$g(z) \leq d(x, z) + g(z) = h_x(z) = \inf_{u \in \mathcal{X}} h_x(u) \leq h_x(x) = g(x).$$

Moreover, from the above inequality it is clear that

$$\begin{aligned} g(z) = g(x) &\iff h_x(z) = h_x(x) \iff \\ d(x, z) + g(z) = d(x, x) + g(x) &\iff d(x, z) = d(x, x) \iff x = z, \end{aligned}$$

that is $x \in \text{prox}_g^d(x)$. □

Lemma 3.1.1. *If \bar{x} minimises g , then $\bar{x} \in \text{prox}_g^d(\bar{x})$.*

Proof. Let $\bar{x} \in \text{argmin}_{x \in \mathcal{X}} g(x)$, then

$$g(\bar{x}) = d(\bar{x}, \bar{x}) + g(\bar{x}) \leq d(\bar{x}, u) + g(u) \text{ for all } u \in \mathcal{X}.$$

Hence $g(\bar{x}) \in \inf_{u \in \mathcal{X}} d(\bar{x}, u) + g(u)$, i.e. $\bar{x} \in \text{prox}_g^d(\bar{x})$. \square

As a last result of this section, the separability property of general proximal operators as defined in (3.1) is addressed, which depend on the separability property of the considered distance function d .

Lemma 3.1.2. *Let $\mathcal{X} = \bigcup_{i=1}^n \Omega_i$ such that $\overset{\circ}{\Omega}_i \cap \overset{\circ}{\Omega}_j = \emptyset$ for $i \neq j$. Suppose $g \in \Gamma_0(\mathcal{X})$ is additively separable in the sense specified by Definition 2.3.1. Suppose the distance function $d : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$ such that $d(x, \cdot) \in \Gamma_0(\mathcal{X})$ for any $x \in \mathcal{X}$ satisfies the property*

$$d(x, u) = \sum_{i=1}^n d(\chi_i x, \chi_i u),$$

where $\chi_i \in \mathcal{X}$ is the characteristic function of Ω_i . Then, the proximal operator defined by (3.1) is separable.

Proof. By hypothesis, g is separable thus $g(u) = \sum_{i=1}^n g(\chi_i u)$. Consider the definition of proximal operator (3.1) together with the separability of g and d :

$$\text{prox}_g^d(x) = \text{argmin}_{u \in \mathcal{X}} \sum_{i=1}^n d(\chi_i x, \chi_i u) + g(\chi_i u) = \sum_{i=1}^n \text{argmin}_{u \in \mathcal{X}} d(\chi_i x, \chi_i u) + g(\chi_i u) \quad (3.2)$$

We adopt now the following notation $u_i = \chi_i u \in \Omega_i$ and $g(\chi_i u) = g_i(u_i)$ with $g_i : \Omega_i \rightarrow \mathbb{R} \cup \{+\infty\}$. Similarly, for the distance function we write $d(\chi_i x, \chi_i u) = d_i(x_i, u_i)$ with $d_i : \Omega_i \times \Omega_i \rightarrow \mathbb{R}$, so that $d(x, u) = \sum_{i=1}^n d_i(x_i, u_i)$. Thus, (3.2) becomes

$$\text{prox}_g^d(x) = \sum_{i=1}^n \text{argmin}_{u \in \mathcal{X}} d_i(x_i, u_i) + g_i(u_i) = \sum_{i=1}^n \text{prox}_{g_i}^{d_i}(x_i),$$

which concludes the proof. \square

In the following, two instances of proximal operators in Banach spaces will be presented.

3.1.1 Definition of the p-norm proximal operator

Following [8], as a first instance of distance $d(\cdot, \cdot)$ we consider $d(x, u) = \frac{1}{p} \|u - x\|_{\mathcal{X}}^p$ for $p > 1$.

Definition 3.1.1 (p-norm proximal operator). *Given $g \in \Gamma_0(\mathcal{X})$ and $p > 1$, the p -proximal operator of g in the Banach space \mathcal{X} is*

$$\text{prox}_g^{1/p \|\cdot\|^p}(x) := \text{argmin}_{u \in \mathcal{X}} \frac{1}{p} \|u - x\|_{\mathcal{X}}^p + g(u). \quad (3.3)$$

A similar definition has been used in [6, 135] in metric spaces and a corresponding definition can be found in [11]. It is easy to show that the operator defined by (3.3) is single-valued, that is the minimiser of the functional appearing in (3.3) is unique.

Lemma 3.1.3. *For $p > 1$, the p -norm proximal operator (3.3) is well defined, i.e. the mapping $u \in \mathcal{X} \mapsto \frac{1}{p}\|u - x\|_{\mathcal{X}}^p + g(u)$ admits minimisers and the minimiser is unique for each $x \in \mathcal{X}$.*

Proof. Let $x \in \mathcal{X}$. Since g is convex, it is bounded from below by $g(u) \geq g(x) + \langle \partial g(x), u - x \rangle$. Thus,

$$\begin{aligned} \frac{1}{p}\|u - x\|_{\mathcal{X}}^p + g(u) &\geq \frac{1}{p}\|u - x\|_{\mathcal{X}}^p + (g(x) + \langle \partial g(x), u - x \rangle) \\ &= \|u\|_{\mathcal{X}} \left(\frac{1}{p} \frac{\|u - x\|_{\mathcal{X}}^p}{\|u\|_{\mathcal{X}}} + \frac{g(x) + \langle \partial g(x), u - x \rangle}{\|u\|_{\mathcal{X}}} \right) \xrightarrow{\|u\|_{\mathcal{X}} \rightarrow +\infty} +\infty, \end{aligned}$$

which means that the mapping is coercive. Moreover, it is strictly convex in u because $p > 1$ and thus it has a unique minimizer. \square

The following lemma and proposition further clarify the role played by the p -norm proximal operator in minimisation problems, and as well as in devising optimisation algorithms.

Lemma 3.1.4. *The following statements are equivalent:*

1. $\bar{x} = \text{prox}_g^{1/p\|\cdot\|^p}(x)$
2. $0 \in \mathbf{J}_{\mathcal{X}}^p(\bar{x} - x) + \partial g(\bar{x})$

Proof. Note that $\bar{x} = \text{prox}_g^{1/p\|\cdot\|^p}(x)$ if and only if it minimises $u \in \mathcal{X} \mapsto \frac{1}{p}\|u - x\|_{\mathcal{X}}^p + g(u)$ and thus $\bar{x} = \text{prox}_g^{1/p\|\cdot\|^p}(x)$ implies $0 \in \mathbf{J}_{\mathcal{X}}^p(\bar{x} - x) + \partial g(\bar{x})$, since for Theorem 2.2.1 it holds that $\partial\left(\frac{1}{p}\|\cdot\|_{\mathcal{X}}^p\right) = \mathbf{J}_{\mathcal{X}}^p(\cdot)$. \square

Lemma 3.1.5. *\bar{x} minimises g if and only if $\text{prox}_g^{1/p\|\cdot\|^p}(\bar{x}) = \bar{x}$.*

Proof. By Lemma 3.1.1, we know that \bar{x} minimises g implies $\text{prox}_g^{1/p\|\cdot\|^p}(\bar{x}) = \bar{x}$.

Consider now that $\text{prox}_g^{1/p\|\cdot\|^p}(\bar{x}) = \bar{x}$ holds. Thus \bar{x} minimises $u \in \mathcal{X} \mapsto \frac{1}{p}\|u - \bar{x}\|_{\mathcal{X}}^p + g(u)$ and therefore $0 \in [\mathbf{J}_{\mathcal{X}}^p(\cdot - \bar{x}) + \partial g(\cdot)](\bar{x}) = \partial g(\bar{x})$. Hence, \bar{x} minimises g . \square

3.1.2 Bregman-proximal operator

In this section, as distance function in the definition of proximal operator (3.1) we consider the Bregman distance associated to a convex functional. It is known that, in Banach spaces Bregman distances are more appropriate than norm distances, since they inherit the rich geometrical properties of the involved Banach space [44]. Indeed, Bregman distances are widely used to measure the distance between k -th

iteration and the generalised solution in many proofs of convergence and have been used a lot in the field of optimisation for defining proximal gradient algorithms under relaxed convergence assumptions [10, 29, 30, 166].

The Bregman distance is defined as the difference between the functional and its linear approximation as follows [40].

Definition 3.1.2. *Let $h : \mathcal{X} \rightarrow \mathbb{R}$ be a convex and continuously-differentiable functional on a Banach space \mathcal{X} . The Bregman distance $B_{\mathcal{X}}^h(\cdot, x) : \mathcal{X} \rightarrow [0, +\infty)$ of h at $x \in \mathcal{X}$ is defined as*

$$B_{\mathcal{X}}^h(u, x) = h(u) - \left(h(x) + \langle \nabla h(x), u - x \rangle \right), \quad \forall u \in \mathcal{X}. \quad (3.4)$$

For $h(\cdot) = \frac{1}{p} \|\cdot\|_{\mathcal{X}}^p$, with $p > 1$, the Bregman distance will be denoted simply as $B_{\mathcal{X}}^p$. Since in this case $\nabla h = \mathbf{J}_{\mathcal{X}}^p$ by Theorem 2.2.1, (3.4) becomes

$$B_{\mathcal{X}}^p(u, x) = \frac{1}{p} \|u\|_{\mathcal{X}}^p - \frac{1}{p} \|x\|_{\mathcal{X}}^p - \langle \mathbf{J}_{\mathcal{X}}^p(x), u - x \rangle,$$

or, equivalently,

$$B_{\mathcal{X}}^p(u, x) := \frac{1}{p} \|u\|_{\mathcal{X}}^p + \frac{1}{p'} \|x\|_{\mathcal{X}}^p - \langle \mathbf{J}_{\mathcal{X}}^p(x), u \rangle.$$

The latter equivalence can be shown by direct computations using (2.27), indeed:

$$-\frac{1}{p} \|x\|_{\mathcal{X}}^p - \langle \mathbf{J}_{\mathcal{X}}^p(x), u - x \rangle = -\frac{1}{p} \|x\|_{\mathcal{X}}^p + \|x\|_{\mathcal{X}}^p - \langle \mathbf{J}_{\mathcal{X}}^p(x), u \rangle = \frac{1}{p'} \|x\|_{\mathcal{X}}^p - \langle \mathbf{J}_{\mathcal{X}}^p(x), u \rangle.$$

In general, Bregman distances and norm-induced distances $\frac{1}{p} \|u - x\|_{\mathcal{X}}^p$ are different. Bregman distances do not satisfy either symmetry or the triangle inequality. Some basic results about Bregman distance in a smooth and uniformly convex Banach space \mathcal{X} can be found in [204, 207]. Note that, in a Hilbert space \mathcal{H} , $B_{\mathcal{H}}^2(u, x) = \frac{1}{2} \|u - x\|_{\mathcal{H}}^2$.

It is now possible to define the proximal operator in terms of the Bregman distance $B_{\mathcal{X}}^p$.

Definition 3.1.3 (Bregman-proximal operator). *Given $g \in \Gamma_0(\mathcal{X})$, the Bregman proximal operator of g in the Banach space \mathcal{X} is*

$$\text{prox}_g^{B_{\mathcal{X}}^p}(x) := \underset{u \in \mathcal{X}}{\text{argmin}} B_{\mathcal{X}}^p(u, x) + g(u). \quad (3.5)$$

A corresponding definition has been introduced in [3]. It is possible to show, similarly to Lemma 3.1.3, that the $\text{prox}_g^{B_{\mathcal{X}}^p}$ operator is single-valued and well-defined, being $\frac{1}{p} \|\cdot\|_{\mathcal{X}}^p$ strictly convex for $p > 1$. Moreover, one can easily show the following results.

Lemma 3.1.6. *The following statements are equivalent:*

1. $\bar{x} = \text{prox}_g^{B_{\mathcal{X}}^p}(x)$
2. $0 \in \mathbf{J}_{\mathcal{X}}^p(\bar{x}) - \mathbf{J}_{\mathcal{X}}^p(x) + \partial g(\bar{x})$
3. $\bar{x} = \left(\mathbf{J}_{\mathcal{X}}^p + \partial g \right)^{-1} \left(\mathbf{J}_{\mathcal{X}}^p(x) \right)$

Proof. Using the definition of $\text{prox}_g^{B_{\mathcal{X}}^p}$, from 1 of the statement we have:

$$\begin{aligned} \bar{x} &= \text{prox}_g^{B_{\mathcal{X}}^p}(x) = \underset{u \in \mathcal{X}}{\text{argmin}} B_{\mathcal{X}}^p(u, x) + g(u) \\ &= \underset{u \in \mathcal{X}}{\text{argmin}} \frac{1}{p} \|u\|_{\mathcal{X}}^p - \frac{1}{p} \|x\|_{\mathcal{X}}^p - \langle \mathbf{J}_{\mathcal{X}}^p(x), u - x \rangle + g(u) \\ 0 &\in \partial \left(\frac{1}{p} \|\cdot\|_{\mathcal{X}}^p - \frac{1}{p} \|x\|_{\mathcal{X}}^p - \langle \mathbf{J}_{\mathcal{X}}^p(x), \cdot - x \rangle + g(\cdot) \right) (\bar{x}) \\ 0 &\in \mathbf{J}_{\mathcal{X}}^p(\bar{x}) - \mathbf{J}_{\mathcal{X}}^p(x) + \partial g(\bar{x}), \end{aligned}$$

which is 2 of this lemma. Then, equivalently, we can write

$$\begin{aligned} \mathbf{J}_{\mathcal{X}}^p(x) &\in \mathbf{J}_{\mathcal{X}}^p(\bar{x}) + \partial g(\bar{x}) \\ \bar{x} &= \left(\mathbf{J}_{\mathcal{X}}^p + \partial g \right)^{-1} \left(\mathbf{J}_{\mathcal{X}}^p(x) \right) \end{aligned}$$

Observe that the operation $(\mathbf{J}_{\mathcal{X}}^p + \partial g)^{-1}$ is well defined, since $\mathbf{J}_{\mathcal{X}}^p + \partial g$ is maximal monotone operator being the subdifferential of a convex, lower semi-continuous function [194]. \square

Lemma 3.1.7. \bar{x} minimises g if and only if $\text{prox}_g^{B_{\mathcal{X}}^p}(\bar{x}) = \bar{x}$.

Proof. By Lemma 3.1.1, we know that \bar{x} minimises g implies $\text{prox}_g^{B_{\mathcal{X}}^p}(\bar{x}) = \bar{x}$.

Consider now that $\text{prox}_g^{B_{\mathcal{X}}^p}(\bar{x}) = \bar{x}$ holds. Thus \bar{x} minimises $u \in \mathcal{X} \mapsto \frac{1}{p} \|u\|_{\mathcal{X}}^p - \frac{1}{p} \|\bar{x}\|_{\mathcal{X}}^p - \langle \mathbf{J}_{\mathcal{X}}^p(\bar{x}), u - \bar{x} \rangle + g(u)$ and therefore $0 \in [\mathbf{J}_{\mathcal{X}}^p(\cdot) - \mathbf{J}_{\mathcal{X}}^p(\bar{x}) + \partial g(\cdot)](\bar{x}) = \partial g(\bar{x})$. Hence, \bar{x} minimises g . \square

3.2 Landweber methods in Banach spaces

In this section, we claim a connection between the popular Landweber iterative regularisation method and a suitable proximal operator, starting our analysis in Hilbert spaces. Then, we present the primal and dual methods in Banach spaces, both generalising gradient-descent strategies, with the dual method being considered the Landweber algorithm in Banach spaces.

3.2.1 Hilbert spaces setting

Recall that the Landweber method in Hilbert spaces is defined by the iteration

$$x^0 \in \mathcal{X}, \quad x^{k+1} = x^k - \tau T^*(Tx - y), \quad (3.6)$$

and is often used as a regularisation algorithm to solve (1.1), where $T : \mathcal{X} \rightarrow \mathcal{Y}$ is a linear and bounded operator between two Hilbert spaces \mathcal{X} and \mathcal{Y} [91, 110] and $T^* : \mathcal{Y}^* \rightarrow \mathcal{X}^*$ is its adjoint operator, defined by

$$\langle T^* y^*, x \rangle_{\mathcal{X}^* \times \mathcal{X}} = \langle y^*, Tx \rangle_{\mathcal{Y}^* \times \mathcal{Y}} \quad \forall y^* \in \mathcal{Y}^*, \forall x \in \mathcal{X}.$$

In an Hilbert spaces setting, the isometric isomorphism between \mathcal{X} and \mathcal{X}^* , and \mathcal{Y} and \mathcal{Y}^* , allows to consider the adjoint operator $T^* : \mathcal{Y}^* \rightarrow \mathcal{X}^*$ simply as $T^* : \mathcal{Y} \rightarrow \mathcal{X}$. As briefly mentioned in Section 1.2.2, the algorithm can be indeed interpreted as a gradient descent method (1.21) for the minimisation of the least square residual functional $f : \mathcal{X} \rightarrow \mathbb{R}$ defined as

$$f(x) = \frac{1}{2} \|Tx - y\|_{\mathcal{Y}}^2.$$

Indeed, since f is convex and differentiable, with $\nabla f(x^k) = T^*(Tx^k - y) \in \mathcal{X}$, iteration (3.6) is nothing but (1.21), that for simplicity we report here

$$x^{k+1} = x^k - \tau \nabla f(x^k). \quad (3.7)$$

Moreover, it is worth recalling that only in Hilbert spaces for all $x \in \mathcal{X}$ the element $\nabla f(x) \in \mathcal{X}^*$ is identified with a unique element in \mathcal{X} itself, up to the canonical isometric isomorphism, so that the design of gradient-type schemes is significantly simplified, being $\nabla f(x) \in \mathcal{X}$. Iteration (3.7) can be equivalently written as

$$\begin{aligned} x^{k+1} &= \operatorname{argmin}_{x \in \mathcal{X}} \frac{1}{2} \|x - (x^k - \tau \nabla f(x^k))\|_{\mathcal{X}}^2 \\ &= \operatorname{argmin}_{x \in \mathcal{X}} \left\{ \frac{1}{2} \|x - x^k\|_{\mathcal{X}}^2 + \tau \langle \nabla f(x^k), x \rangle \right\}, \end{aligned} \quad (3.8)$$

where the terms expanding on x^k and constant with respect to x have been neglected. It is quite evident that the latter minimisation problem is well defined, since the functional is strongly convex. Minimisation problem (3.8) can be thus recast in the framework of the theory of proximal operators as

$$x^{k+1} = \operatorname{prox}_{\tau \langle \nabla f(x^k), \cdot \rangle}(x^k),$$

where the prox operator in Hilbert spaces is defined by (1.18). This shows that iteration (3.6) corresponds to the computation of a point which decreases $\langle \nabla f(x^k), x \rangle$ and simultaneously is close (i.e., proximal) to the previous iteration. The step size τ can be here thought as the weight which balances between the two terms $\frac{1}{2} \|x - x^k\|_{\mathcal{X}}^2$ and $\langle \nabla f(x^k), x \rangle$.

The Landweber method in Hilbert spaces converges to the generalised inverse x^\dagger (1.5) of the problem (1.1) in the noisy free case.

Theorem 3.2.1. [182] *If $y \in R(T) \oplus R(T)^\perp$ and τ satisfies $0 < \tau < \frac{2}{\|T\|^2}$, the sequence $(x^k)_k$ given by (3.6) is strongly convergent for any initialisation $x^0 \in \mathcal{X}$:*

$$\lim_{k \rightarrow +\infty} x^k = x^\dagger + P_{N(T)} x^0,$$

Algorithm 1 Landweber algorithm (dual method) in Banach spaces [207]

Parameters: $p, r > 1$, $\{\tau_k\}_k$ specific real sequence with $\tau_k > 0$.

Initialisation: Start with $x_0 \in \mathcal{X}$.

repeat

$$x^{k+1} = \mathbf{J}_{\mathcal{X}^*}^{\mathbf{r}'} \left(\mathbf{J}_{\mathcal{X}}^{\mathbf{r}}(x^k) - \tau_k T^* \mathbf{J}_{\mathcal{Y}}^{\mathbf{p}}(Tx^k - y) \right) \quad (3.9)$$

until convergence

where $P_{N(T)}$ is the orthogonal projection onto $N(T)$. Moreover, with $x^0 = 0$ we have

$$\lim_{k \rightarrow +\infty} \|x^k - x^\dagger\|_{\mathcal{X}} = 0.$$

In the case of noisy data the method, in general, does not converge because y^δ does not satisfy the condition required by the theorem. In such a case, however, the method has a property of semi-convergence. In this case, it is interesting to study the reconstruction error

$$\|x^k - x\|_{\mathcal{X}}$$

where x is the unknown signal defined in equation (1.1). As k increases, the reconstruction error is decreasing first and increasing afterwards, meaning that the choice of the number of iteration is crucial in order to have a good reconstructed signal. A possible criterion for the stopping rule in case of noisy data is given by the discrepancy principle [91]

$$k(\delta, y^\delta, \tau) = \min \left\{ k \in \mathbb{N} \mid \|T(x^k)^\delta - y^\delta\|_{\mathcal{Y}} \leq \alpha\delta \right\} \quad \text{for a given } \alpha > 1,$$

where $\left((x^k)^\delta \right)_k$ is the sequence of Landweber iterates applied for the resolution of the noisy problem (1.2).

3.2.2 Dual method in Banach spaces

Consider now a different framework, where the forward operator $T : \mathcal{X} \rightarrow \mathcal{Y}$ acts between reflexive, strictly convex and smooth Banach spaces. In [207], a generalisation of the Landweber method to non-Hilbertian Banach spaces has been first proposed. Such a generalisation is not straightforward, because a non-Hilbertian Banach space is not always isometrically isomorphic to its dual, so that the iteration scheme (3.8) is no longer formally consistent, being $x^k \in \mathcal{X}$ summed to $\nabla f(x^k) \in \mathcal{X}^*$ with $\mathcal{X}^* \not\cong \mathcal{X}$ [41]. The key tool for the generalisation of the iteration scheme (3.8) to Banach spaces are the duality maps, which associate an element of a Banach space \mathcal{X} with an element of its dual \mathcal{X}^* .

On these grounds, for fixed parameters $p, r > 1$, the Landweber-type iteration scheme of the seminal paper [207] for the solution of (1.1) reads as in Algorithm 1, where $x_0 \in \mathcal{X}$ is the initial guess and $\tau_k > 0$ is a proper variable step-size. By analogy with Hilbertian Landweber iteration (3.6), we notice that the descent step

in (3.9) is performed in the dual space \mathcal{X}^* , since both $\mathbf{J}_{\mathcal{X}}^{\mathbf{r}}(x^k)$ and $T^*\mathbf{J}_{\mathcal{Y}}^{\mathbf{p}}(Tx^k - y)$ of (3.9) belong to \mathcal{X}^* . Hence, this algorithm will be referred to as *dual method* [204]. The result of the gradient-descent step computation in the dual space is then associated with its corresponding element of the primal space thanks to the inverse of $\mathbf{J}_{\mathcal{X}}^{\mathbf{r}}$, that is $(\mathbf{J}_{\mathcal{X}}^{\mathbf{r}})^{-1} = \mathbf{J}_{\mathcal{X}^*}^{\mathbf{r}'}$ by (2.21).

The dual method defined by (3.9) has been introduced by a pure formal approach. Anyway, its relation with gradient methods is evident, as sketched in [204]. Indeed, it is easy to show that it can be interpreted as suitable gradient-descent steps. By simple application of the chain rule for differentiation of the p -power residual

$$f(x) = \frac{1}{p} \|Tx - y\|_{\mathcal{Y}}^p, \quad (3.10)$$

we have

$$\begin{aligned} \nabla f(x) &= \left(\left(\nabla \left(\frac{1}{p} \|\cdot\|_{\mathcal{Y}}^p \right) \Big|_{Tx-y} \right)^* \nabla(Tx - y) \right)^* \\ &= \left(\left(\mathbf{J}_{\mathcal{Y}}^{\mathbf{p}}(Tx - y) \right)^* T \right)^* = T^* \mathbf{J}_{\mathcal{Y}}^{\mathbf{p}}(Tx - y), \end{aligned} \quad (3.11)$$

where $\left(\left(\mathbf{J}_{\mathcal{Y}}^{\mathbf{p}} \right)^* \right)^* = \mathbf{J}_{\mathcal{Y}}^{\mathbf{p}}$ because \mathcal{Y} is reflexive. This shows that iterative step (3.9) can be written as

$$\mathbf{J}_{\mathcal{X}}^{\mathbf{r}}(x^{k+1}) = \mathbf{J}_{\mathcal{X}}^{\mathbf{r}}(x^k) - \tau_k \nabla f(x^k),$$

which is a gradient descent step computed in the dual space \mathcal{X}^* , in analogy with the iterative step (3.7) of Hilbert setting. The duality mappings are crucial in the definition of the gradient-descent step in Banach spaces. The map $\mathbf{J}_{\mathcal{Y}}^{\mathbf{p}}$ of \mathcal{Y} appears only in the definition of the gradient of the residual function (3.10) and the parameter \mathbf{p} of the duality map is linked to the way the discrepancy between data y and model observations Tx is measured. On the other hand, the duality map $\mathbf{J}_{\mathcal{X}}^{\mathbf{r}}$ of \mathcal{X} and its inverse $\mathbf{J}_{\mathcal{X}^*}^{\mathbf{r}'}$ make the computation of the gradient descent step possible, moving the iterates x^k to the dual space \mathcal{X}^* where $\nabla f(x^k)$ lives and then the new iterate back to the primal space \mathcal{X} . The parameter \mathbf{r} refers only to the space \mathcal{X} and to how the distance between iterates is measured. If $\mathcal{X} = \mathcal{Y}$ it might be a good option to consider $\mathbf{p} = \mathbf{r}$. Moreover, in Lebesgue spaces with a constant exponent it is reasonable to take for \mathbf{p} and \mathbf{r} the same value of the constant exponent.

The dual method is also known as the generalisation of the Landweber regularisation method in Banach spaces. In this regard, we report here the result on the convergence of the Landweber algorithm in Banach spaces, see [204], for the resolution of the noise-free problem (1.1).

Theorem 3.2.2. *Algorithm 1 either stops after a finite number of iterations with x^\dagger or the sequence of the iterates $(x^k)_k$ converges strongly to x^\dagger .*

In the noisy case of (1.2), the following termination rule has to be applied

$$k(\delta, y^\delta, D) = \min \left\{ k \in \mathbb{N} \mid \|Ax^{k\delta} - y^\delta\|_{\mathcal{Y}} < \frac{\delta}{D} \right\} \quad \text{for some } D \in (0, 1). \quad (3.12)$$

This ensures that $(x^{k+1})^\delta$ is a better approximation to the generalised inverse x^\dagger than $(x^k)^\delta$ as long as $k < k(\delta, y^\delta, D)$. Algorithm 1 together with the discrepancy principle (3.12) as stopping rule is a regularisation method for problem (1.2). See [204] for more details.

3.2.3 Primal method in Banach spaces

Another possible way to consider gradient-based approaches in Banach spaces are the so-called *primal methods*, where the gradient step is directly computed in the primal space \mathcal{X} . In [204], the primal method is presented and analysed as a minimisation algorithm to minimise the residual functional (3.10). The iteration scheme is reported in Algorithm 2. Iteration (3.13) can be expressed as

Algorithm 2 Primal method in Banach spaces

Parameters: $p, r > 1$, $\{\tau_k\}_k$ specific real sequence with $\tau_k > 0$.

Initialisation: Start with $x_0 \in \mathcal{X}$.

repeat

$$x^{k+1} = x^k - \tau_k \mathbf{J}_{\mathcal{X}^*}^{r'} \left(T^* \mathbf{J}_Y^p (T x^k - y) \right) \quad (3.13)$$

until convergence

$$x^{k+1} = x^k - \tau_k \mathbf{J}_{\mathcal{X}^*}^{r'} \left(\nabla f(x^k) \right) \quad (3.14)$$

where it becomes evident that the primal method is a gradient-descent based method for f defined in (3.10), with gradient descent step computed in the primal space \mathcal{X} as for (3.7). As for the dual method, the gradient-descent step is possible thanks to the duality mappings. In particular, in (3.13) the map $\mathbf{J}_{\mathcal{X}^*}^{r'}$ allows to compute the descent step in the primal space \mathcal{X} , moving the gradient of the residual, that is an element of \mathcal{X}^* , to \mathcal{X} .

3.2.4 Primal and dual method as proximal point algorithms

The role of the duality maps $\mathbf{J}_{\mathcal{X}}^r$ and $\mathbf{J}_{\mathcal{X}^*}^{r'}$ in the definition of Algorithms 1 and 2 is related to the particular geometry induced by the r -norm of the Banach space \mathcal{X} and its dual \mathcal{X}^* with the dual r' -norm. To better clarify the role of the norms, we present a simple explanation of the dual method (3.9) and of the primal method (3.13), in terms of suitable proximal operators introduced in Sections 3.1.1 and 3.1.2.

We start considering the functional (3.10) and its subdifferential (3.11) in smooth, reflexive and strictly convex Banach spaces, inspired by (3.8) we write

$$x^{k+1} = \operatorname{argmin}_{x \in \mathcal{X}} \left\{ \frac{1}{r} \|x - x^k\|_{\mathcal{X}}^r + \tau_k \langle \nabla f(x^k), x \rangle \right\}, \quad (3.15)$$

which by differentiation leads to the optimality condition

$$\mathbf{J}_{\mathcal{X}}^r(x^{k+1} - x^k) + \tau_k \nabla f(x^k) = 0,$$

The latter can be solved explicitly, since $(\mathbf{J}_{\mathcal{X}}^{\mathbf{r}})^{-1} = \mathbf{J}_{\mathcal{X}^*}^{\mathbf{r}'}$ by (2.21), leading to

$$x^{k+1} = x^k - \tilde{\tau}_k \mathbf{J}_{\mathcal{X}^*}^{\mathbf{r}'}(\nabla f(x^k)),$$

where $\tilde{\tau}_k = \tau_k^{r^*-1}$, that is exactly the primal method iteration (3.14). Thus, the minimisation problem (3.15) allows to interpret the primal method as a proximal point algorithm for the functional (3.10) with the proximal operator defined by (3.3), that is

$$x^{k+1} = \text{prox}_{\tau_k \langle \nabla f(x^k), \cdot \rangle}^{1/r \|\cdot\|^r}(x^k).$$

Analogously, by using the Bregman distance $B_{\mathcal{X}}^{\mathbf{r}}$ instead of the norm distance $\frac{1}{r} \|\cdot\|_{\mathcal{X}}^r$ in (3.15), we can write

$$x^{k+1} = \underset{x \in \mathcal{X}}{\text{argmin}} \left\{ B_{\mathcal{X}}^{\mathbf{r}}(x, x^k) + \tau_k \langle \nabla f(x^k), x \rangle \right\} \quad (3.16)$$

where the objective is strictly convex and differentiable, so that

$$\nabla \left(\frac{1}{r} \|\cdot\|_{\mathcal{X}}^r + \frac{1}{r^*} \|x^k\|_{\mathcal{X}}^{r^*} - \langle \mathbf{J}_{\mathcal{X}}^{\mathbf{r}}(x^k), \cdot \rangle + \tau_k \langle \nabla f(x^k), \cdot \rangle \right) (x^{k+1}) = 0.$$

The latter equality leads to the following iterative gradient-type iteration

$$\mathbf{J}_{\mathcal{X}}^{\mathbf{r}}(x^{k+1}) - \mathbf{J}_{\mathcal{X}}^{\mathbf{r}}(x^k) + \tau_k \nabla f(x^k) = 0,$$

which can be written as

$$x^{k+1} = \mathbf{J}_{\mathcal{X}^*}^{\mathbf{r}'} \left(\mathbf{J}_{\mathcal{X}}^{\mathbf{r}}(x^k) - \tau_k \nabla f(x^k) \right).$$

We notice that this is exactly the Landweber dual method in Banach spaces (3.9) for the minimisation of the p -power residual functional (3.10). This novel interpretation allows us to write the dual method in Banach spaces in terms of Bregman proximal operators. More precisely, the basic Landweber iteration (3.9) can thus be written as

$$x^{k+1} = \text{prox}_{\tau_k \langle \nabla f(x^k), \cdot \rangle}^{B_{\mathcal{X}}^{\mathbf{r}}}(x^k),$$

which provides a useful interpretation of the dual algorithm in the context of convex optimisation.

3.3 Modular-based dual method in $L^{p(\cdot)}(\Omega)$

We present in this section the problems that arise when trying to use primal and dual methods in $L^{p(\cdot)}(\Omega)$ spaces and we propose a modular-based alternative to the dual method in this scenario.

3.3.1 Primal and dual methods in $L^{p(\cdot)}(\Omega)$: main issues

The primal and dual methods can be effectively used in Lebesgue spaces with constant exponents. Indeed, when $\mathcal{X} = L^s(\Omega)$ with $s \in (1, +\infty)$ and $\mathcal{Y} = L^q(\Omega)$ with $q \in (1, +\infty)$, the dual method (3.9) becomes

$$x^{k+1} = \mathbf{J}_{L^{s'}}^{r'} \left(\mathbf{J}_{L^s}^r(x^k) - \tau_k T^* \mathbf{J}_{L^q}^p(Tx^k - y) \right),$$

with duality maps defined by (2.20). Similarly, the primal method (3.13) reads as

$$x^{k+1} = x^k - \tau_k \mathbf{J}_{L^{s'}}^{r'} \left(T^* \mathbf{J}_{L^q}^p(Tx^k - y) \right) .$$

In the above equations, the inverse of $\mathbf{J}_{L^s}^r$, which by (2.21) is equal to $\mathbf{J}_{(L^s)^*}^{r'}$, is given by $\mathbf{J}_{L^{s'}}^{r'}$ thanks to (2.22) and to the isometric isomorphism between $(L^s(\Omega))^*$ and $L^{s'}(\Omega)$. Thus, as already introduced in Section 2.2.2, such isometric isomorphism is crucial to obtain an exact and explicit expression for the inverse of the duality maps given by (2.22).

We now consider Lebesgue spaces with a variable exponent $L^{p(\cdot)}(\Omega)$.

3.3.1.1 Approximation of the inverse of the duality map

As stated in Proposition 2.2.1, there is not an isometric isomorphism between $(L^{p(\cdot)}(\Omega))^*$ and $L^{p'(\cdot)}(\Omega)$ and thus the inverse of $\mathbf{J}_{L^{p(\cdot)}}^r$ cannot be obtained as in (2.22).

As a consequence of this fact, the inverse of $\mathbf{J}_{L^{p(\cdot)}}^r$ does not directly relate to the point-wise conjugate exponents of $p(\cdot)$. The following approximation is nevertheless considered in [5, 26, 31, 92]

$$\left(\mathbf{J}_{L^{p(\cdot)}}^r \right)^{-1} = \mathbf{J}_{(L^{p(\cdot)})^*}^{r'} \approx \mathbf{J}_{L^{p'(\cdot)}}^{r'} \quad (3.17)$$

to use the Landweber method (Algorithm 1) in $L^{p(\cdot)}(\Omega)$ spaces. It gives an inexact but explicit formula of the duality maps for the dual space $(L^{p(\cdot)}(\Omega))^*$. Proposition 2.2.1 states that the norms of $L^{p(\cdot)}(\Omega)$ and $(L^{p(\cdot)}(\Omega))^*$ are not isometric, and it provides optimal and finite bounds among them. Unfortunately, since the duality maps are the Fréchet derivative of the r -power of the norm, these bounds do not give any quantitative information about the goodness of the approximations. Anyway, due to continuity arguments, we can say that the approximation should be good for small ranges $[p_-, p_+]$ of exponent values, since for $p_- = p_+$, which coincides with the constant exponent case, the equality holds. Both the duality maps $\mathbf{J}_{L^{p(\cdot)}}^r$ and $\mathbf{J}_{L^{p'(\cdot)}}^{r'}$ can be computed by (2.10), so that Algorithm 1 with the approximation given by (3.17) is completely implementable in closed form.

We refer the reader to [31], where we performed some numerical tests on simple imaging deblurring problems solved in variable exponents Lebesgue spaces with the dual Landweber method (Alg. 1) with the approximation given by (3.17).

Algorithm 3 Modular-based Gradient Descent in $L^{p(\cdot)}(\Omega)$

Parameters: $\{\tau_k\}_k$ s.t. $0 < \bar{\tau} \leq \tau_k \leq \frac{pc(1-\delta)}{K}$ with $0 < \delta < 1$, for all $k \geq 0$.

Initialisation: $x^0 \in L^{p(\cdot)}(\Omega)$.

repeat

$$x^{k+1} = (\mathbf{J}_{\bar{\rho}_{p(\cdot)}})^{-1} \left(\mathbf{J}_{\bar{\rho}_{p(\cdot)}}(x^k) - \tau_k \nabla f(x^k) \right) \quad (3.18)$$

until convergence

3.3.1.2 Heavy computational times

However, the computation of the duality map $\mathbf{J}_{L^{p(\cdot)}}^r$ defined by (2.10) requires the computation of the Luxemburg norm (2.4) which, as previously discussed in Section 2.3, is not-separable and hence quite heavy to compute, since it requires the resolution of a one-dimensional minimisation problem. This fact makes the definition of gradient-descent iterative schemes rather inefficient in terms of computational times, since both Algorithm 1 and Algorithm 2 involve duality maps, and hence long norm computations. As a consequence, the expression (2.10) is not suited to be used in a computational optimisation framework.

3.3.2 Modular-based alternative to dual method in $L^{p(\cdot)}(\Omega)$

The computation of the Luxemburg norm (2.4) is quite heavy, as shown also in Section 2.1.1. On the other hand, the modular function, as seen in Definition 2.1.1, has a closed form expression, similar to the norm in the conventional constant case of L^p spaces, and thus its computation is more efficient than the Luxemburg norm's one. In Section 2.3, we observed that the modular is separable and differentiable and its gradient is invertible with expression provided in Propositions 2.3.2 and 2.3.3. We thus follow [147] and define in Algorithm 3 a more efficient modular-based gradient descent iteration in the general setting of variable exponent Lebesgue spaces for the minimisation of a proper, convex and smooth function $f : L^{p(\cdot)}(\Omega) \rightarrow \mathbb{R} \cup \{+\infty\}$. The following set of assumptions needs to hold:

A1. $\nabla f : L^{p(\cdot)}(\Omega) \rightarrow (L^{p(\cdot)}(\Omega))^*$ is $(p-1)$ -Hölder-continuous with exponent $1 < p \leq 2$ and constant $K > 0$, i.e.:

$$\|\nabla f(u) - \nabla f(v)\|_{(L^{p(\cdot)})^*} \leq K \|u - v\|_{L^{p(\cdot)}}^{p-1} \quad \forall u, v \in L^{p(\cdot)}(\Omega).$$

A2. There exists $c > 0$ such that, for all $u, v \in L^{p(\cdot)}(\Omega)$,

$$\langle \mathbf{J}_{\bar{\rho}_{p(\cdot)}}(u) - \mathbf{J}_{\bar{\rho}_{p(\cdot)}}(v), u - v \rangle \geq c \max \left\{ \|u - v\|_{L^{p(\cdot)}}^p, \|\mathbf{J}_{\bar{\rho}_{p(\cdot)}}(u) - \mathbf{J}_{\bar{\rho}_{p(\cdot)}}(v)\|_{(L^{p(\cdot)})^*}^p \right\}.$$

The latter bound was previously used in [105, 145]. It is a compatibility condition between the ambient space $L^{p(\cdot)}(\Omega)$ and the Hölder smoothness properties of the residual function to minimise to achieve algorithmic convergence.

The minimisation of the specific function f is achieved solving at each iteration (3.18) the following minimisation problem

$$x^{k+1} = \operatorname{argmin}_{u \in L^{p(\cdot)}(\Omega)} \mathbf{J}_{\bar{\rho}_{p(\cdot)}}(u) - \langle \mathbf{J}_{\bar{\rho}_{p(\cdot)}}(x^k), u \rangle + \tau_k \langle \nabla f(x^k), u \rangle \quad (3.19)$$

which can be equivalently formulated in terms of the Bregman distance of the modular $\mathbf{J}_{\bar{\rho}_{p(\cdot)}}$, that we will denote with $B_{\bar{\rho}_{p(\cdot)}}$. Indeed, by summing constant terms to (3.19), we obtain

$$\begin{aligned} x^{k+1} &= \operatorname{argmin}_{u \in L^{p(\cdot)}(\Omega)} \mathbf{J}_{\bar{\rho}_{p(\cdot)}}(u) - \mathbf{J}_{\bar{\rho}_{p(\cdot)}}(x^k) - \langle \mathbf{J}_{\bar{\rho}_{p(\cdot)}}(x^k), u - x^k \rangle + \tau_k \langle \nabla f(x^k), u \rangle \\ &= \operatorname{argmin}_{u \in L^{p(\cdot)}(\Omega)} B_{\bar{\rho}_{p(\cdot)}}(u, x^k) + \tau_k \langle \nabla f(x^k), u \rangle = \operatorname{prox}_{\tau_k \langle \nabla f(x^k), \cdot \rangle}^{B_{\bar{\rho}_{p(\cdot)}}}(x^k), \end{aligned}$$

which gives an interpretation of Algorithm 3 in terms of Bregman proximal operators, as carried out for Algorithm 1.

An important remark is that whenever $\nabla f(x_k) = 0$ at some $k \geq 0$, a stationary point $x^{k+1} = (\mathbf{J}_{\bar{\rho}_{p(\cdot)}})^{-1}(\mathbf{J}_{\bar{\rho}_{p(\cdot)}}(x^k)) = x^k$ is found, as expected. This cannot be obtained using Algorithm 1 with approximation (3.17), since the duality map $\mathbf{J}_{L^{p(\cdot)}}^r$ is not coupled with its exact inverse but only with an approximation.

Note that when not only the solution space is a variable exponent Lebesgue spaces $\mathcal{X} = L^{p(\cdot)}(\Omega)$, but also the measurement space is a variable exponent Lebesgue space $\mathcal{Y} = L^{q(\cdot)}(\Omega)$, a more natural and consistent choice for the objective function for the resolution of (1.1) is the modular of the discrepancy between the model observation and the data, i.e. $f(x) = \bar{\rho}_{q(\cdot)}(Tx - y)$. In this way, instead of (3.11), the gradient of f becomes $\nabla f(x^k) = T^* \mathbf{J}_{\bar{\rho}_{q(\cdot)}}(Tx^k - y)$: the heavy computations of the norm $\|\cdot\|_{L^{q(\cdot)}}$ are not required, making the iteration scheme faster.

3.3.3 Comparison between Landweber and modular-based gradient descent

To further motivate the choice of the modular to devise optimisation algorithms, we present here a brief analysis to compare Algorithm 1 with the approximation given by (3.17), i.e. norm-based gradient descent in $L^{p(\cdot)}(\Omega)$ with the needed approximation for the inverse of duality maps, and modular-based gradient descent in $L^{p(\cdot)}(\Omega)$ described in Algorithm 3.

To this aim, we consider the blurred and noisy images shown in Figures 3.1c and 3.2c and tackle the deblurring problem by norm-based and modular-based gradient descent in $L^{p(\cdot)}(\Omega)$. The forward operator $T : L^{p(\cdot)}(\Omega) \rightarrow L^{p(\cdot)}(\Omega)$ is a blurring operator given by a Gaussian Point-Spread-Function (PSF). The PSF for the test image in Figure 3.1c is shown in Figure 3.1b, while the one for the satellite image (Figure 3.2c) is shown in Figure 3.2c. The variable exponents used for the deblurring are shown in Figures 3.1d and 3.2d, respectively. We will comment on and motivate their selection in the next section.

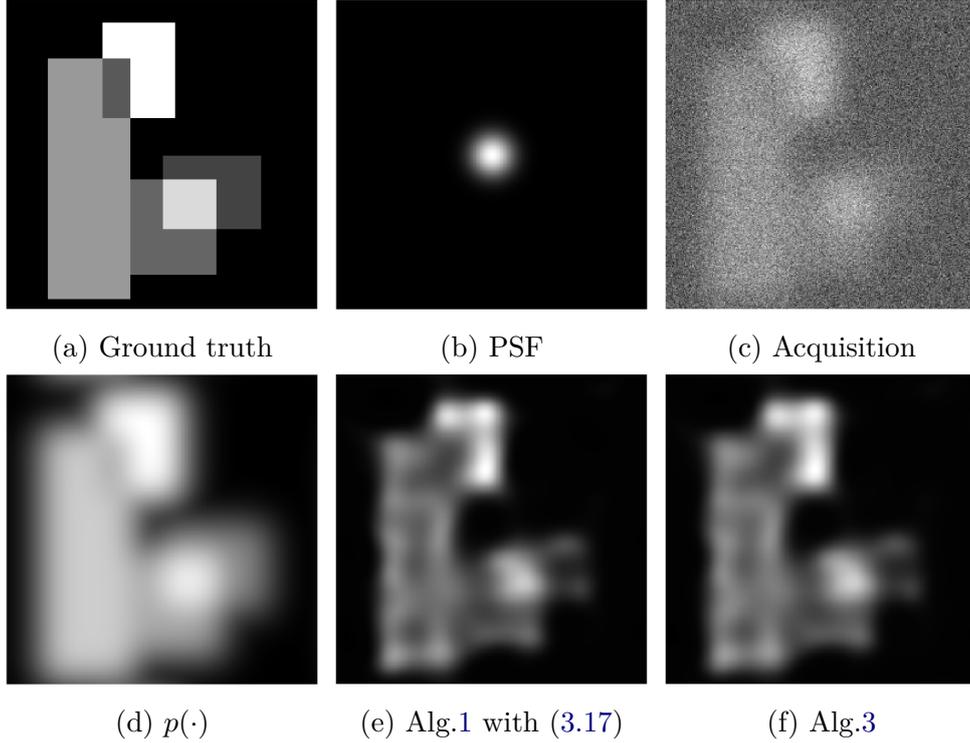


Figure 3.1: Deblurring of test image in $L^{p(\cdot)}(\Omega)$ spaces. The restored images, here shown, are attained in correspondence with the minimum of the reconstruction error.

	$\frac{\ x_{rec}-x_{gt}\ _{L^{p(\cdot)}}}{\ x_{gt}\ _{L^{p(\cdot)}}}$	$\frac{\rho_{p(\cdot)}(x_{rec}-x_{gt})}{\rho_{p(\cdot)}(x_{gt})}$	PSNR	SSIM	CPU time	#iter
Alg.1 approx	0.4836	0.4053	18.674	0.61975	9266	528
Alg.3	0.3602	0.2624	18.876	0.62435	48	618

Table 3.1: Comparison between norm-based and modular-based gradient descent for images in Figure 3.1.

We solve the minimisation of $\|Tx - y^\delta\|_{L^{p(\cdot)}}$ with Algorithm 1 with approximation (3.17) and the minimisation of $\bar{\rho}_{p(\cdot)}(Tx - y^\delta)$ with Algorithm 3. In Figure 3.1, the reconstruction error of both regularisation algorithms is shown and the semi-convergence behaviour is quite evident. We then consider to stop both algorithms when the minimum of the reconstruction error is attained and report it in Tables 3.1 and 3.2 the number of iterations needed to reach the minimum, alongside with the corresponding values of reconstruction error, residual, PSNR, SSIM, CPU time. The difference in terms of CPU time between the use of the norm and the modular is quite striking. It is interesting to point out that in terms of reconstruction quality, modular-based gradient descent yields better values of relative reconstruction error computed both in terms of the norm and of the modular, PSNR and SSIM. The approximation (3.17) introduced in Algorithm 1 to make it applicable in $L^{p(\cdot)}(\Omega)$ spaces introduces small errors slightly affecting the reconstructed images

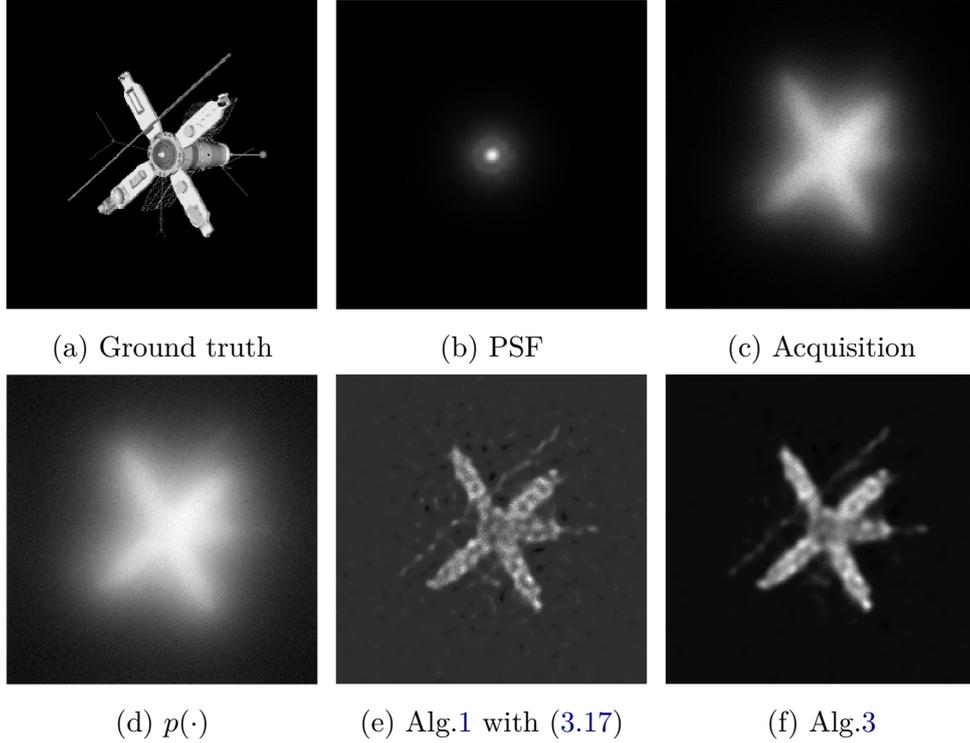


Figure 3.2: Deblurring of satellite image in $L^{p(\cdot)}(\Omega)$ spaces. The shown restored images are attained in correspondence of the minimum of the reconstruction error.

	$\frac{\ x_{rec}-x_{gt}\ _{L^{p(\cdot)}}}{\ x_{gt}\ _{L^{p(\cdot)}}}$	$\frac{\rho_{p(\cdot)}(x_{rec}-x_{gt})}{\rho_{p(\cdot)}(x_{gt})}$	SNR	PSNR	CPU time	#iter
Alg.1 approx	0.3133	0.1832	7.83	93.42	85345	1143
Alg.3	0.3168	0.1839	9.49	95.08	106.97	1606

Table 3.2: Comparison between norm-based and modular-based gradient descent for images in Figure 3.2.

and a slightly oscillating behaviour in the reconstruction error, as one can see from Figure 3.1.

We recall that in Chapter 1, in Figure 1.7 and in Figure 1.6 the reconstruction obtained in L^2 and in L^p of the test image (Fig.3.1a) and of the satellite image (Fig.3.2a) are shown.

3.3.4 How to choose variable exponents

In this section, some possible strategies to select variable exponents are detailed.

Consider the resolution of an inverse problems (1.1) in variable exponent Lebesgue spaces, with the forward operator $T : L^{p(\cdot)}(\Omega) \rightarrow L^{q(\cdot)}(\Omega)$ being a bounded linear operator between $\mathcal{X} = L^{p(\cdot)}(\Omega)$ and $\mathcal{Y} = L^{q(\cdot)}(\Omega)$ and with the acquisition $y \in \mathcal{Y} = L^{q(\cdot)}(\Omega)$. In general, there are two variable exponent maps to be chosen,

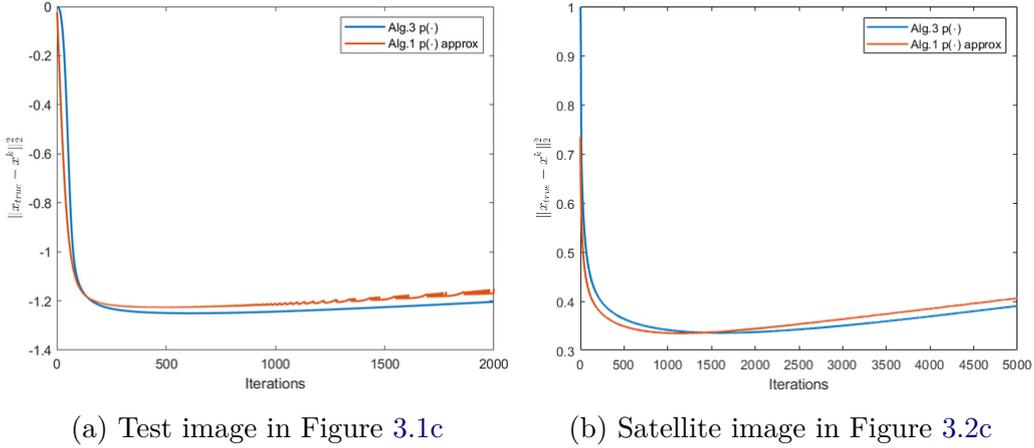


Figure 3.3: Semi-convergence of Landweber method and modular-based GD in $L^{p(\cdot)}(\Omega)$ spaces.

namely $p(\cdot)$ for the solution space $\mathcal{X} = L^{p(\cdot)}(\Omega)$ and $q(\cdot)$ for the acquisition space $\mathcal{Y} = L^{q(\cdot)}(\Omega)$.

It is difficult (and somehow undesirable) to have a unified configuration as the selection of both $p(\cdot)$ and $q(\cdot)$ is strictly problem-related. Parameters $q(\cdot)$ are related to the regularity of the measured data as well as the different noise distributions considered on the acquisitions. For instance, with impulsive noise, values of q_- and q_+ closer to 1 are preferred while for Gaussian noise values closer to 2 are more effective. Solution space parameters p_- and p_+ relate to the regularity of the solution to retrieve. As a consequence, their choice is intrinsically harder. A possible strategy for informed pixel-wise variable exponents consists in basing them on observed data for $q(\cdot)$ and an approximation of the reconstruction for $p(\cdot)$, as we did in [5, 31, 145]. In particular, given $z \in L^{p(\cdot)}(\Omega)$ an approximation of the desired solution, the exponent map $p(\cdot)$ is chosen in the following as

$$p(t) = p_- + (p_+ - p_-) \Upsilon \left(\frac{|z(t)|}{\max_{t \in \Omega} |z(t)|} \right), \quad t \in \Omega, \quad (3.20)$$

where $\Upsilon : \mathbb{R} \rightarrow \mathbb{R}$ is an interpolation function. Similarly, for the variable exponent $q(\cdot)$ of the measurement space, given $w \in L^{q(\cdot)}(\Omega)$ the exponent map $q(\cdot)$ is chosen as

$$q(t) = q_- + (q_+ - q_-) \Upsilon \left(\frac{|w(t)|}{\max_{t \in \Omega} |w(t)|} \right), \quad t \in \Omega, \quad (3.21)$$

where w can be the acquisition y directly, or Tz , where z is the approximation of the solution used in (3.20), or $T(p(\cdot))$ with $p(\cdot)$ given by (3.20). We refer the reader to [31], where a comparison between different choices for p_- , p_+ , q_- , q_+ and different interpolation strategies Υ is carried out for image deblurring with gradient descent (3.9) in $L^{p(\cdot)}(\Omega)$. It is shown that the interpolation strategy used does not affect the reconstruction quality significantly, but the latter is more sensitive to the choice of a good interval $[p_-, p_+]$. However, the tuning of p_- , p_+ , q_- , q_+ remains heuristic

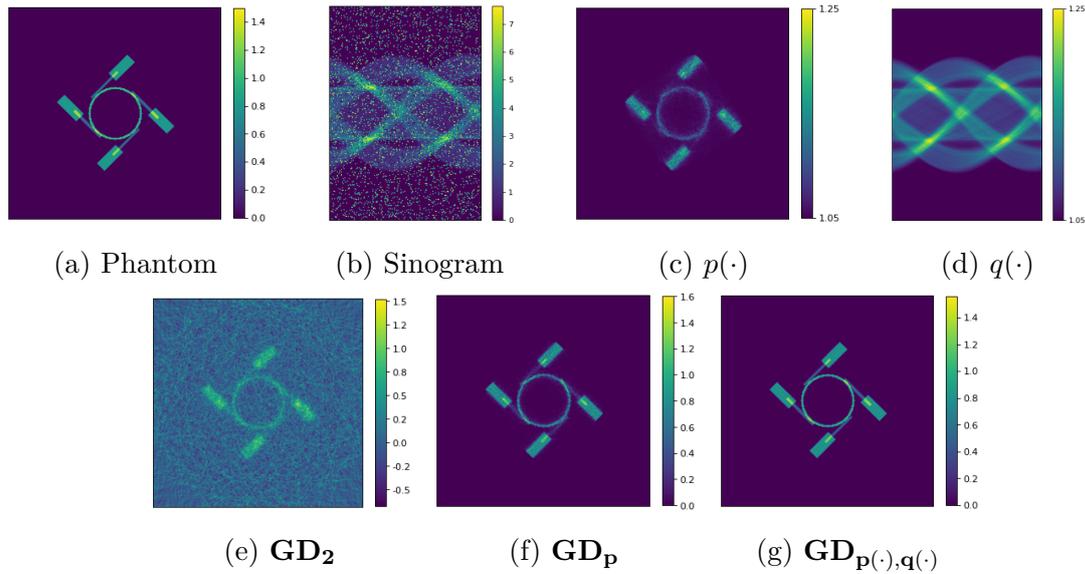


Figure 3.4: Numerical test on simulated CT data with gradient descent algorithms. Comparison between Hilbert and non-Hilbertian Lebesgue spaces reconstructions with constant and variable exponents after 500 iterations.

and specific to the particular problem and image considered. This will be object of further study.

3.3.5 Numerical tests with modular-based gradient descent

In Section 3.3.3, a numerical comparison between norm-based and modular-based gradient descent is presented and shows that the latter performs way better in terms of computational times.

Here, we present experimental results of the proposed modular-based strategy (Alg.3) on a Computed Tomography (CT) imaging problem in a simulated setting using a simulated dataset provided by the `python` CIL library [128]. The synthetic ground truth phantom is shown Figure 3.4a. We consider (1.2) where the forward operator $T \in \mathcal{L}(\mathcal{X}, \mathcal{Y})$ is given by the discrete Radon transform. Due to the nature of the CT problem, acquisitions and unknown images have two distinct spatial domains. For example, formulating the problem in a Hilbert setting we should consider $\mathcal{X} = L^2(\Omega_{\mathcal{X}})$ and $\mathcal{Y} = L^2(\Omega_{\mathcal{Y}})$. However, in the following for simplicity we denote with Ω both the spatial domain both over \mathcal{X} and \mathcal{Y} , with a slight abuse of notation. The acquisition is obtained using a 2D parallel beam geometry, with 180 projection angles on a 1 angle separation, 256 detector elements, and pixel size of 0.1. Further details on the mathematical modelling and on the 2D parallel beam geometry can be found in Appendix A. After applying the forward operator, a high level (15%) of salt-and-pepper noise is applied to the sinogram. The noisy sinogram is shown in Figure 3.4b. The goal of this section is to quantitatively compare the performance of Algorithm 3 with the corresponding Hilbert (3.7) and Banach space

versions (3.9):

- **GD₂**: $\mathcal{X} = \mathcal{Y} = L^2(\Omega)$, $f(x) = \frac{1}{2}\|Tx - y\|_2^2$ by GD (3.7);
- **GD_p**: $\mathcal{X} = \mathcal{Y} = L^p(\Omega)$, $f(x) = \frac{1}{p}\|Tx - y\|_p^p$ by Banach SGD (3.9) with $r = p$;
- **GD_{p(·),q(·)}**: $\mathcal{X} = L^{p(\cdot)}(\Omega)$, $\mathcal{Y} = L^{q(\cdot)}(\Omega)$ for appropriately chosen exponent maps, $f(x) = \bar{\rho}_{q(\cdot)}(Tx - y)$ with modular-based GD Algorithm 3.

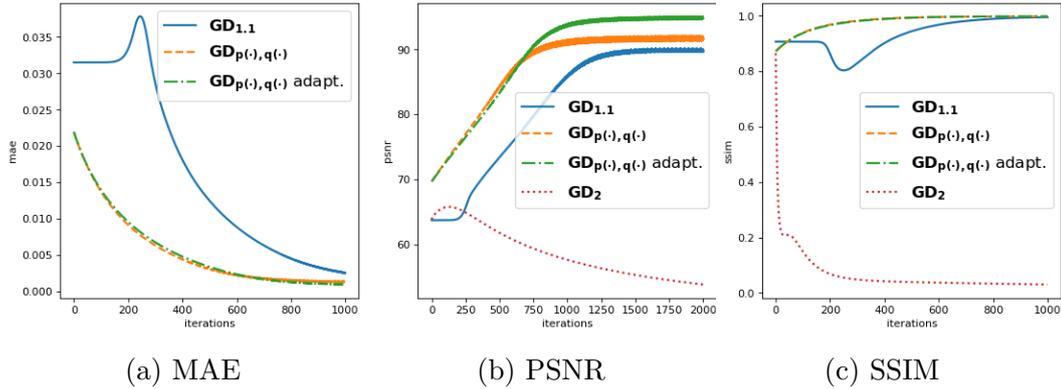
Recall that **GD₂** and **GD_p** are constant exponents methods while **GD_{p(·),q(·)}** depends on variable exponents.

The choice of the exponents is now discussed. In the constant case, we set $p = 1.1$ because of the salt-and-pepper noise present in the sinogram. The most suited fidelity for this noise setting is an L^1 -fidelity, which however is non-smooth and thus not-suited for our modelling. Thus, the choice of a constant value p close to 1 allows to approximate such non-smooth fidelity with a smooth one. In Lebesgue spaces with variable exponents, exponents are maps sensitive to local assumptions on both the solution and the measured data. The considered selection strategy is given by (3.20) for the exponent map $p(\cdot)$ of the solution space $\mathcal{X} = L^{p(\cdot)}(\Omega)$ and (3.21) for the exponent map $q(\cdot)$ of the acquisition space $\mathcal{Y} = L^{q(\cdot)}(\Omega)$. To this end, we first compute an approximate reconstruction \tilde{x} by running **GD_p** in $L^{1.1}(\Omega)$ for 70 iterations with a constant step-size regime. The map $p(\cdot)$ is then computed via a linear interpolation of \tilde{x} between $p_- = 1.05$ and $p_+ = 1.25$ and it is shown in Figure 3.4c. The map $q(\cdot)$ is chosen as the linear interpolation between $q_- = 1.05$ and $q_+ = 1.25$ of $T(p(\cdot))$ and it is reported in Figure 3.4d. The bounds p_-, p_+ and q_-, q_+ are chosen by prior assumptions on y (sparse phantom) and on the noise (impulsive).

The Hilbert reconstruction obtained with **GD₂** after 500 iterations is shown in Figure 3.4e and, as expected, it is very poor. In Figures 3.4f and 3.4f, reconstruction obtained after 500 iterations of **GD_p** and **GD_{p(·),q(·)}** respectively are reported. Considering a Banach space setting drastically improves the reconstruction quality and using a variable exponent yields particularly sharp solutions. This behaviour is particularly evident analysing quality metrics (MAE, PSNR and SSIM) computed after 500 iterations in Table 3.3.

The variable exponent setting also allows an heuristic, but effective, modification in the utilised solution space exponents. Assuming that the iterates are improving the quality of the reconstruction, an adaptive strategy would be to correspondingly update the solution space exponents $p(\cdot)$, based on the current iteration. To this end, we update $p(\cdot)$ based on the current solution estimate once every β_{updates} epochs to adapt the exponents along the iterations. We will refer to this strategy as **GD_{p(·),q(·)}** adaptive.

In Figure 3.5, we study the semi-convergence property and we report the mean absolute error (MAE), peak signal to noise ratio (PSNR) and structural similarity index (SSIM) of the iterates x^k w.r.t. the known ground-truth phantom along the

Figure 3.5: Quality metrics of \mathbf{GD}_2 , \mathbf{GD}_p , $\mathbf{GD}_{p(\cdot),q(\cdot)}$ and $\mathbf{GD}_{p(\cdot),q(\cdot)}$ adaptive.

Algorithm	It.	Tot.	MAE	PSNR	SSIM
\mathbf{GD}_2	0.44s	1324s	1.619e-1	61.81	0.0387
$\mathbf{GD}_{1.1}$	0.43s	1297s	1.247e-2	73.43	0.9452
$\mathbf{GD}_{p(\cdot),q(\cdot)}$	0.44s	1317s	3.101e-3	84.23	0.9947
$\mathbf{GD}_{p(\cdot),q(\cdot)}$ adapt.	0.47s	1403s	3.386e-3	83.33	0.9939

Table 3.3: CPU times after 3000 iterations. MAE, PSNR, and SSIM values after 500 iterations (before noise over-fitting).

first 1000 iterations. \mathbf{GD}_2 is omitted from MAE to improve the readability of the plots, due to its poor performance. We underline that Banach methods makes the algorithms more stable and less sensitive to the stopping rule considered, since the quality of the reconstructions obtained with \mathbf{GD}_p and $\mathbf{GD}_{p(\cdot),q(\cdot)}$ remains high increasing the number of iterations. This property to investigate further from a theoretical point of view is very useful in practice, since a definition of a univocal stopping rule is hard to attain.

3.4 A Modular-based Stochastic variant

The key challenge for the viability of many deterministic iterative methods for real-world image reconstruction problems is their scalability to data-size. The use of deterministic iterative algorithms, such as Algorithms 1 and 3, may be prohibitively expensive in large-size applications. For example, the highest per-iteration cost in emission tomography lies in the application of the entire forward operator at each iteration, whereas each image domain datum in computed tomography often requires several gigabytes of storage space. The same could thus be a bottleneck in the application of Algorithm 3. The stochastic gradient descent (SGD) paradigm addresses this issue [193].

Thus, following the strategy performed by the seminal work of Robbins and

Monro [193] we adapt a stochastic gradient descent (SGD) strategy to the non-standard setting of variable exponent Lebesgue space, in order to reduce the per-iteration complexity costs. In [193], to minimise a smooth function $f : \mathcal{H} \rightarrow \mathbb{R}$ in \mathcal{H} Hilbert space, the following SGD scheme is proposed:

$$x^{k+1} = x^k - \tau_k h(x^k, \xi_k), \quad (3.22)$$

where $h(\cdot, \xi)$ is an unbiased estimator of the gradient of f , i.e.

$$\mathbb{E}_\xi [h(x, \xi)] = \nabla f(x),$$

and the sequence of positive step-sizes $(\tau_k)_k$ satisfies the following conditions to ensure convergence:

$$\sum_{k \in \mathbb{N}} \tau_k^2 < +\infty, \quad \sum_{k \in \mathbb{N}} \tau_k = +\infty. \quad (3.23)$$

An unbiased estimator for the gradient of the residual function is usually obtained by defining a suitable decomposition of the original problem into (finitely many) sub-problems, and by implementing an iterative scheme where only a batch of data, typically one, is used to compute the current update. This greatly reduces the computational complexity per-iteration, and enjoys excellent scalability with respect to data size. Note that the use of SGD schemes is popular within the mathematical imaging community [126, 127] due to its applicability in large-scale applications such as medical imaging [116, 170, 214]. However, its extension to variable exponent Lebesgue setting is not trivial.

3.4.1 Stochastic Gradient Descent in Banach spaces

Let the forward operator of the model $T : \mathcal{X} \rightarrow \mathcal{Y}$ be a linear and bounded operator between the reflexive, strictly convex and smooth Banach spaces \mathcal{X} and \mathcal{Y} . We partition the forward operator $T : \mathcal{X} \rightarrow \mathcal{Y}$ into a finite number of block-type operators T_1, \dots, T_{N_s} with $T_i : \mathcal{X} \rightarrow \mathcal{Y}_i$, where $N_s \in \mathbb{N}$ is the number of subsets of data, and the family of Lebesgue measurable subsets $(\mathcal{Y}_i)_{i=1}^{N_s}$ of \mathcal{Y} is such that $\mathcal{Y}_i \cap \mathcal{Y}_j = \emptyset$ for $i \neq j$, and $\mathcal{Y} = \bigcup_{i=1}^{N_s} \mathcal{Y}_i$ (as in the hypothesis of Definition 2.3.1). This results in a partition of the forward model, i.e. the same partition is applied to the observations

$$T_i x = y_i, \quad y_i = \chi_i(y)$$

where χ_i is the characteristic function of \mathcal{Y}_i . Classical examples of this methodology include Kaczmarz methods [116, 170, 189], as sketched in Figure 3.6 in the case of Computed Tomography (CT). Though the vast majority of existing stochastic methods operate in Hilbert, and in particular Euclidean spaces, there has recently been a renewed interest in stochastic methods in Banach spaces and, in particular, $L^p(\Omega)$ spaces, see [126]. The SGD version of the iteration (3.9) in Banach spaces takes naturally the form

$$x^{k+1} = \mathbf{J}_{\mathcal{X}^*}^r \left(\mathbf{J}_{\mathcal{X}}^r(x^k) - \tau_k T_{i_k}^* \mathbf{J}_{\mathcal{Y}_{i_k}}^p (T_{i_k} x^k - y_{i_k}) \right) \quad (3.24)$$

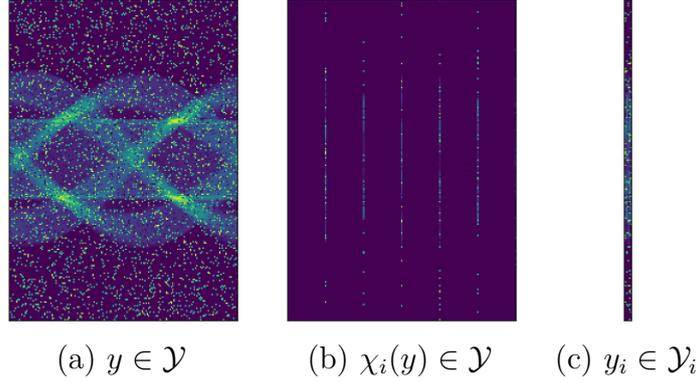


Figure 3.6: Partition of the observation $y \in \mathcal{Y}$ into $y_i \in \mathcal{Y}_i$ with $N_s = 30$

where the indices $i_k \in \{1, \dots, N_s\}$ are sampled uniformly at random.

Note that (3.24) is the standard form of SGD to minimise separable norm-based objectives as in (3.10). In particular, the partition of the forward operator results in the following splitting of the residual

$$f_i(x) = \|T_i x - y_i\|_{\mathcal{Y}_i}^p \Rightarrow f(x) = \frac{1}{N_s} \sum_{i=1}^{N_s} f_i(x) = \frac{1}{N_s} \|Tx - y\|_{\mathcal{Y}}^p. \quad (3.25)$$

The last equality is true only when the norm in \mathcal{Y} is separable. Moreover, by Theorem 2.2.1, decomposition (3.25) ensures that $\nabla f_i(x) = T_i^* \mathbf{J}_{\mathcal{Y}_i}^p (T_i x^k - y_i)$ is an unbiased estimator for the gradient of f and it shows that each step of (3.24) can thus be computed by simply taking a sub-differential of a single sum-function f_i :

$$x^{k+1} = \mathbf{J}_{\mathcal{X}^*}^{\mathbf{r}'} \left(\mathbf{J}_{\mathcal{X}}^{\mathbf{r}}(x^k) - \tau_k \nabla f_{i_k}(x^k) \right). \quad (3.26)$$

From (3.26), it becomes evident that this is a generalisation of the standard form of SGD (3.22). Indeed, if \mathcal{X} is an Hilbert space, taking $r = 2$, (3.26) reduces to:

$$x^{k+1} = x^k - \tau_k \nabla f_{i_k}(x^k).$$

Sampling reduces the per-iteration computational cost in \mathcal{Y} by a factor of N_s .

In [126] convergence of the iterates to a least-squares solution is shown, under conditions on the step-sizes similar to (3.23) in the standard SGD in Hilbert spaces. Before reporting the convergence result, we recall here the definition of p -convexity for Banach spaces.

Definition 3.4.1. *Let \mathcal{X} be a Banach space and $p > 1$. The space \mathcal{X} is p -convex if there exists a constant $c_p > 0$ such that*

$$\frac{1}{p} \|x - y\|_{\mathcal{X}}^p \geq \frac{1}{p} \|x\|_{\mathcal{X}}^p - \langle z, y \rangle + \frac{c_p}{p} \|y\|_{\mathcal{X}}^p$$

for all $x, y \in \mathcal{X}$ and all $z \in \mathbf{J}_{\mathcal{X}}^p(x)$.

Theorem 3.4.1. *Let \mathcal{X} be a p -convex and smooth Banach space and \mathcal{Y} a Banach space with separable norm. Let S be the set of least-squares solutions (1.3) of (1.1).*

If $(\mu_k)_{k \in \mathbb{N}}$ satisfy the following conditions

$$\sum_{k=1}^{\infty} \mu_k^{p^*} < \infty, \quad \sum_{k=1}^{\infty} \mu_k = \infty,$$

then the sequence $(x^k)_{k \in \mathbb{N}}$ given by (3.24) converges almost surely to a least-squares solution (1.3) of (1.1):

$$\mathbb{P} \left(\liminf_{k \rightarrow \infty} \inf_{\bar{x} \in S} \|x^k - \bar{x}\|_{\mathcal{X}} = 0 \right) = 1.$$

For noisy measurements and in standard $L^p(\Omega)$ spaces, the regularising property of SGD has been established in [126] by defining suitable stopping criteria. However, early stopping strategies are hard to use in practice and providing methods that are less sensitive to data over-fitting is crucial for their practical use.

3.4.2 Variable exponents modular-based SGD

We consider now the setting where both $\mathcal{X} = L^{p(\cdot)}(\Omega)$ and $\mathcal{Y} = L^{q(\cdot)}(\Omega)$ are variable exponents Lebesgue spaces. To define a suitable SGD in this scenario, we take as objective function the modular of the residual

$$f(x) = \frac{1}{N_s} \bar{\rho}_{q(\cdot)}(Tx - y).$$

The choice of a norm-based objective function, in this setting, is not possible because the Luxemburg norm $\|\cdot\|_{L^{p(\cdot)}}$ is not separable, as shown by Lemma 2.3.2, and, thus, splitting (3.25) cannot be achieved. Partitioning the space \mathcal{Y} into subsets \mathcal{Y}_i as described above, the same splitting of the exponent $q(\cdot)$ is implicitly considered as well, namely $q_i(\cdot) = \chi_i(q(\cdot))$. The modular-based objective is splitted into $N_s \geq 1$ sub-objectives

$$f_i(x) = \bar{\rho}_{q_i(\cdot)}(T_i x - y_i),$$

so that $\nabla f_i(x) = T_i^* \mathbf{J}_{\bar{\rho}_{q_i(\cdot)}}(T_i x - y_i)$ is an unbiased estimator of $\nabla f(x)$.

Then, at iteration k and a randomly sampled index $1 \leq i_k \leq N_s$, the corresponding stochastic iterates are given by (3.28). Moreover, similarly as before, one can show that (3.28) can be equivalently defined as

$$\begin{aligned} x^{k+1} &= \operatorname{argmin}_{u \in L^{p(\cdot)}(\Omega)} \bar{\rho}_{p(\cdot)}(u) - \langle \mathbf{J}_{\bar{\rho}_{p(\cdot)}}(x^k), u \rangle + \tau_k \langle \nabla f_{i_k}(x^k), u \rangle \\ &= \operatorname{argmin}_{u \in L^{p(\cdot)}(\Omega)} B_{\bar{\rho}_{p(\cdot)}}(u, x^k) + \tau_k \langle \nabla f_{i_k}(x^k), u \rangle = \operatorname{prox}_{\tau_k \langle \nabla f_{i_k}(x^k), \cdot \rangle}^{B_{\bar{\rho}_{p(\cdot)}}}(x^k). \end{aligned}$$

The pseudocode of the resulting stochastic modular-based gradient descent in $L^{p(\cdot)}(\Omega)$ is reported in Algorithm 4. We expect that through minimal modifications an analogous result to Theorem 3.4.1 can be proved in this setting too, as well the regularisation properties of Algorithm 4, as in [126] for Banach SGD defined by (3.24).

Algorithm 4 Stochastic Modular-based Gradient Descent in $L^{p(\cdot)}(\Omega)$

Parameters: τ_0 s.t. $0 < \bar{\tau} \leq \tau_0 \leq \frac{pc(1-\delta)}{K}$, $0 < \delta < 1$, $N_s \geq 1$, $\gamma > 0$, $\eta > 0$.**Initialisation:** $x^0 \in L^{p(\cdot)}(\Omega)$.**repeat** Select uniformly at random $i_k \in \{1, \dots, N_s\}$.

Set the step-size

$$\tau_k = \frac{\tau_0}{1 + \eta(k/N_s)^\gamma} \quad (3.27)$$

Compute

$$x^{k+1} = (\mathbf{J}_{\bar{p}_{p(\cdot)}})^{-1} \left(\mathbf{J}_{\bar{p}_{p(\cdot)}}(x^k) - \tau_k \nabla f_{i_k}(x^k) \right) \quad (3.28)$$

until convergence

A detailed convergence proof, however, is left for future research. The main issue in carrying out this analysis in $L^{p(\cdot)}(\Omega)$ lies in the fact that these spaces are not p -convex.

3.4.3 Numerical results

We now present experimental results of the proposed Algorithm 4 on two exemplar problems in CT (see Appendix A).

Similarly as in Section 3.3.5, we compare the performance of Algorithm 4 with the corresponding Hilbert and Banach space versions (3.24):

- **SGD₂**: $\mathcal{X} = \mathcal{Y} = L^2(\Omega)$, $f(x) = \frac{1}{2} \|Tx - y\|_2^2$ by SGD;
- **SGD_p**: $\mathcal{X} = \mathcal{Y} = L^p(\Omega)$, $f(x) = \frac{1}{p} \|Tx - y\|_p^p$ by Banach SGD (3.24) with $r = p$;
- **SGD_{p(\cdot), q(\cdot)}**: $\mathcal{X} = L^{p(\cdot)}(\Omega)$, $\mathcal{Y} = L^{q(\cdot)}(\Omega)$ for appropriately chosen exponent maps, $f(x) = \bar{\rho}_{q(\cdot)}(Tx - y)$ with modular-based SGD Algorithm 4.

The first set of experiments uses the same simulated setting of Section 3.3.5, whilst in the second set of experiments we consider the dataset of real-world CT scans of a walnut taken from <https://doi.org/10.5281/zenodo.4279549> [162], with a fan beam geometry. For this data, we utilise the insights from the first numerical tests on simulated data and apply Algorithm 4 in a setting with different noise modalities across the sinogram space. The experiments were conducted in python, using the open source package [128] for the tomographic back-end.

3.4.3.1 Hyper-parameter selection

In the following experiments, we employ a decaying step-size regime such that it satisfies (3.27), which has been shown to be a sufficient condition for the convergence

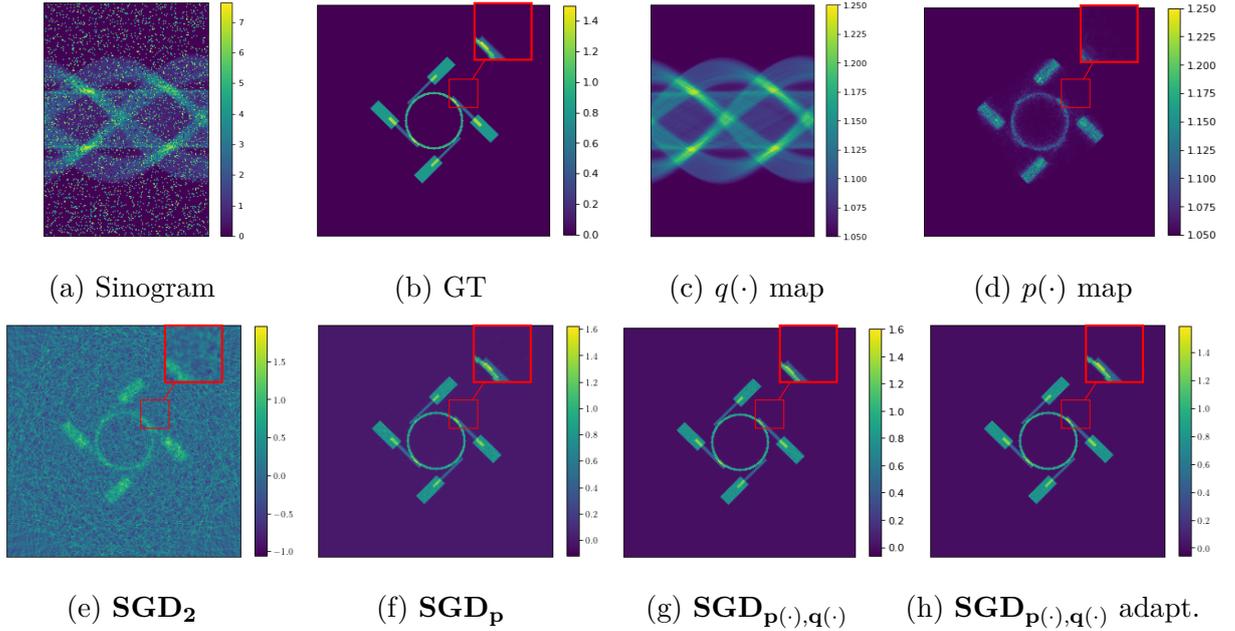


Figure 3.7: Simulated CT data: noisy acquisition, ground truth, exponent maps $p(\cdot)$ and $q(\cdot)$, reconstructions (after 40 epochs) with \mathbf{SGD}_2 , \mathbf{SGD}_p , $\mathbf{SGD}_{p(\cdot),q(\cdot)}$ and $\mathbf{SGD}_{p(\cdot),q(\cdot)}$ adaptive.

of SGD in Banach spaces [126]. A need for a decaying step-size regime is common for stochastic gradient descent to mitigate the effects of inter-iterate variance. Specifically, we use (3.27), where $\tau_0 > 0$ is the initial step-size, and $\gamma > 0$ and $c > 0$ control the decay speed. For the Hilbert space setting, \mathbf{SGD}_2 , initial step-size τ_0 is given by the Lipschitz constant of the gradient of the objective function, namely $\tau_0 = 0.95 / \max_i \|T_i\|^2$. For \mathbf{SGD}_p and $\mathbf{SGD}_{p(\cdot),q(\cdot)}$ the estimation of the respective Hölder continuity constant is more delicate and τ_0 has to be tuned to guarantee convergence. However, its tuning is rather easy and the employ of a decaying strategy makes the choice of τ_0 less critical.

3.4.3.2 Simulated data

We consider again the phantom of Figure 3.4a and the sinogram of Figure 3.4b, obtained using a 2D parallel beam geometry, with 180 projection angles on a 1 angle separation, 256 detector elements, and pixel size of 0.1, and with high level (15%) of salt-and-pepper noise.

To compute subset data T_i and y_i , the forward operator and the sinogram are pre-binned according to equally spaced views (w.r.t. the number of subsets) of the scanner geometry. Subsequent subset data are offset from one another by the subset index i , as sketched in Figure 3.6. We consider $N_s = 30$ batches.

As mentioned above, the goal is to compare the Hilbert, Banach and $L^{p(\cdot)}(\Omega)$ settings to see the role played by the choice of the solution space in the quality

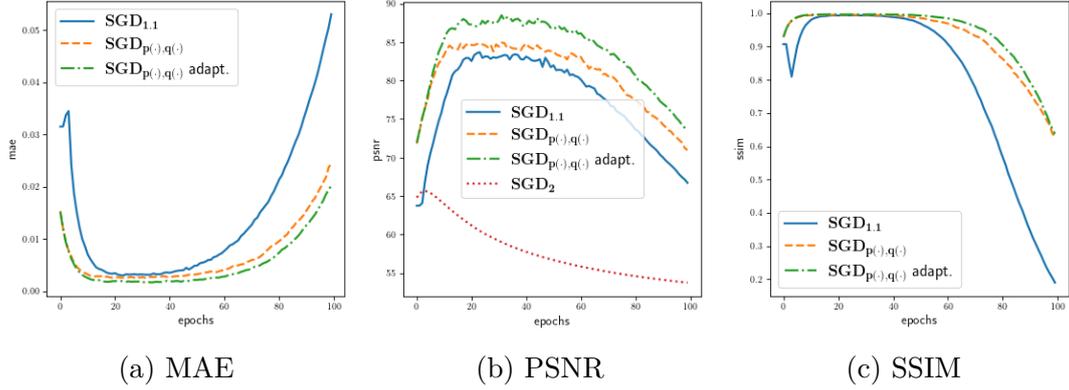


Figure 3.8: Quality metrics along the first 100 epochs of \mathbf{SGD}_2 ; $\mathbf{SGD}_{1.1}$; $\mathbf{SGD}_{p(\cdot),q(\cdot)}$ with and without adapting the exponent maps $p(\cdot)$, $q(\cdot)$.

Algorithm	Deterministic		Stochastic ($\cdot = \mathbf{S}$)					
	It.	Tot.	It.	Epoch	Tot.	MAE	PSNR	SSIM
$\cdot \mathbf{GD}_2$	0.44s	1324s	0.02s	0.74s	74.4 s	2.582e-1	57.89	0.0304
$\cdot \mathbf{GD}_{1.1}$	0.43s	1297s	0.03s	0.81s	81.3s	3.671e-3	82.64	0.9897
$\cdot \mathbf{GD}_{p(\cdot),q(\cdot)}$	0.44s	1317s	0.03s	0.91s	91.2s	2.887e-3	84.05	0.9927
$\cdot \mathbf{GD}_{p(\cdot),q(\cdot) \text{ adapt.}}$	0.47s	1403s	0.03s	0.96s	96.5s	1.777e-3	88.10	0.9965
Compute $p(\cdot)$, $q(\cdot)$	0.45s	16s	0.03s	0.8s	4.0s	-	-	-

Table 3.4: Comparison of per iteration cost and total CPU times after 3000 iterations for deterministic algorithms and after 100 epochs for stochastic algorithms with $N_s = 30$. MAE, PSNR and SSIM values for stochastic algorithms are computed after 40 epochs (before noise over-fitting).

of the reconstructions. The algorithmic choices are \mathbf{SGD}_2 , \mathbf{SGD}_p for $p = 1.1$ that will be denoted directly as $\mathbf{SGD}_{1.1}$ and $\mathbf{SGD}_{p(\cdot),q(\cdot)}$ for appropriately chosen variable exponents $p(\cdot)$ and $q(\cdot)$. The following choices for the step-sizes parameters τ_0 and γ are made depending on the specific algorithm. For \mathbf{SGD}_2 , τ_0 is set as $0.95/\max_i \|T_i\|^2$ and $\gamma = 0.51$. For \mathbf{SGD}_p and $\mathbf{SGD}_{p(\cdot),q(\cdot)}$, we use $\tau_0 = 0.015$ with $\gamma = (p - 1)/p + 0.01$ and $\gamma = (p_- - 1)/p_- + 0.01$ respectively.

The choice of the exponents has already been discussed in Section 3.3.5 for this same dataset. We just recall that as far as the variable exponents are concerned, we consider again the selection strategy given by (3.20) and (3.21) and we first compute an approximate reconstruction \tilde{x} by running \mathbf{SGD}_p in $L^{1.1}(\Omega)$ for 5 epochs with a constant step-size regime. The map $p(\cdot)$ is then computed via a linear interpolation of \tilde{x} between $p_- = 1.05$ and $p_+ = 1.25$ and it is shown in Figure 3.7(d). The map $q(\cdot)$ is chosen as the linear interpolation between $q_- = 1.05$ and $q_+ = 1.25$ of $T(p(\cdot))$. Recall that the parameters p_-, p_+ and q_-, q_+ are selected taking into account prior

assumptions on y (sparse phantom) and on the noise (impulsive).

In Figure 3.7 all the reconstructions of the noisy sinogram by the considered strategies are shown. In particular, as before, the \mathbf{SGD}_2 classical Hilbertian solution presents many artefacts and has very poor quality. All the other reconstructions obtained with Banach spaces methods appear to be very close to the ground truth phantom. In particular, $\mathbf{SGD}_{p(\cdot),q(\cdot)}$ (with and without adaptive strategy) yields solutions with a more truthful range than \mathbf{SGD}_p with a constant exponent $p = 1.1$. It is hard to notice differences visually but it is more interesting to have a look at the MAE, PSNR and SSIM of the iterates x^k w.r.t. the known ground-truth phantom along the first 100 epochs, Figure 3.8. Since PSNR favours smoothness, it is thus beneficial for \mathbf{SGD}_2 , whereas MAE promotes sparsity hence is beneficial for both \mathbf{SGD}_p and $\mathbf{SGD}_{p(\cdot),q(\cdot)}$. \mathbf{SGD}_2 is omitted from MAE and SSIM in Figure 3.8 to improve the readability of the plots, due to its poor performance. Figure 3.8 shows that Banach space algorithms provide better performance than \mathbf{SGD}_2 with respect to all three quality metrics. Note that all the results show the well-known semi-convergence behaviour with respect to the metrics considered. A theoretical result on the selection of an early stopping strategy is needed for stopping the iterations appropriately. We observe that not only the use of variable exponents improves all quality metrics, but also makes the algorithm more stable: the quality of the reconstructed solutions is significantly less sensitive to the number of epochs, making possible early stopping strategies more robust.

In Table 3.4, the CPU times for deterministic (\mathbf{GD}_2 , \mathbf{GD}_p and $\mathbf{GD}_{p(\cdot),q(\cdot)}$) approaches and stochastic ones (\mathbf{SGD}_2 , \mathbf{SGD}_p and $\mathbf{SGD}_{p(\cdot),q(\cdot)}$) are compared. It is worth observing that, as expected, the stochastic paradigm significantly reduces computational costs and CPU times. Moreover, the variable exponents setting has a similar CPU time than the others scenarios, thanks to the choice of a modular-based algorithm.

3.4.3.3 Real CT dataset

As real test, we now consider a fan beam CT dataset of a walnut taken from <https://doi.org/10.5281/zenodo.4279549> [162], from which we take a 2D fan beam sinograms from the centre plane of the cone. The fan beam data uses 0.5 angle separation over the range $[0, 360]$. The used sinogram is obtained by pre-binning the raw data by a factor of 8, resulting in 280 effective detector pixels. See Appendix A for a brief description of CT with 2D fan geometry.

We consider a more difficult noise setting that requires exponential maps which vary in the acquisition domain. Here, we assume that noise has a different effect on the background (zero entries) and the foreground (non-zero entries) of the clean sinogram. Namely, we apply 10% salt and pepper noise to the background, and speckle noise with mean 0 and variance 0.01 to the foreground, see Figure 3.9(a) for the resulting noisy sinogram. Notably, since this noise model has a non-uniform effect across the measurement data, it requires Banach space methods favouring the

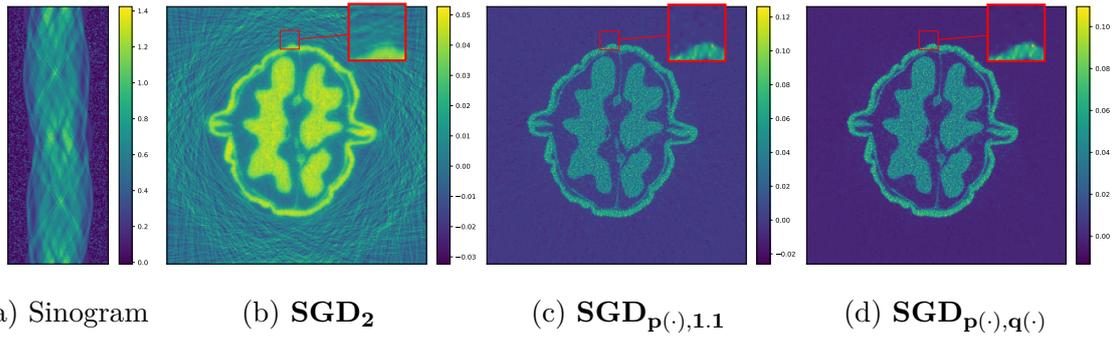


Figure 3.9: Real CT data: noisy sinogram and different reconstructions obtained with SGD strategies.

adjustment of the Lebesgue exponents are expected to perform better than those making use of a constant value.

Taking as a baseline the result obtained by \mathbf{SGD}_2 , shown in Figure 3.9(b), we compare here the effect of allowing variable exponents in the solution space only with the effect of allowing both maps $p(\cdot)$ and $q(\cdot)$ to be chosen. By choosing $p(\cdot)$ based on the initial image and interpolating it between $p_- = 1.2$ and $p_+ = 1.3$ we then compare $\mathbf{SGD}_{p(\cdot),1.1}$ (i.e., fixed exponent $q = 1.1$ in the measurement space) with $\mathbf{SGD}_{p(\cdot),q(\cdot)}$ where $p(\cdot)$ is as before while $q(\cdot)$ is chosen from the sinogram by interpolating between $q_- = 1.1$ and $q_+ = 1.9$. The reconstruction obtained with $\mathbf{SGD}_{p(\cdot),1.1}$ is reported in Figure 3.9(c), while in Figure 3.9(d) we see the one obtained with $\mathbf{SGD}_{p(\cdot),q(\cdot)}$. The results show that a flexible framework where both maps $p(\cdot)$ and $q(\cdot)$ adapt to local contents are more suited for dealing with this challenging scenario. Given the lack of ground-truth data, only a visual assessment is done here.

As step-size we used (3.27), with $N_s = 10$ subsets, and suitable τ_0 and γ . For \mathbf{SGD}_2 , $\tau_0 = 0.95 / \max_i \|T_i\|^2$, $\gamma = 0.51$. For $\mathbf{SGD}_{p(\cdot),q(\cdot)}$ we $\tau_0 = 0.001$, $\gamma = 0.58$.

3.5 Final discussion

In this chapter, we analysed regularisation methods in Banach spaces highlighting their links to optimisation theory and proximal operators. Indeed, we showed that the Landweber method in Banach spaces can be equivalently defined in terms of Bregman proximal operators as a proximal point algorithm for the residual function (3.10). The primal method is instead a proximal point algorithm for the smooth function (3.10) with respect to the p -norm proximal operator. We discuss the complications arising when trying to use such algorithms in $L^{p(\cdot)}(\Omega)$ spaces and propose an alternative method whose definition is based on the modular. The proposed algorithm consists of a modular-based gradient descent in $L^{p(\cdot)}(\Omega)$. We compared its performance with the approximation of the Landweber algorithm in $L^{p(\cdot)}(\Omega)$ given by (3.17), which validates the choice of the modular against the norm in the defin-

ition of optimisation strategies in this setting. In the last section, we presented a stochastic variant of the modular-based gradient descent algorithms and showed its performance in the reconstruction of both simulated and real CT data.

Proximal gradient algorithms in

$$L^{p(\cdot)}(\Omega)$$

In order to minimise functionals that result from the sum of a smooth part and a possibly non-smooth one, an effective strategy is the use of proximal gradient descent algorithms, also called forward-backward algorithms. After a brief review of these methods in Banach spaces, we present two possible way to define algorithms of this form in the setting of variable exponent Lebesgue spaces, making use of the modular instead of the norm in the definition of both the forward, i.e. gradient, and the backward, i.e. proximal, steps. We then focus on sparse reconstruction models and present some 1D and 2D numerical tests in mixed-noise scenarios or heterogeneous signals, in order to highlight the flexibility of the variable exponent spaces. To conclude, we present a numerical study on convergence rates.

4.1	Norm-based proximal-gradient algorithms in Banach spaces	78
4.1.1	Proximal primal-gradient algorithm	79
4.1.2	Proximal dual-gradient algorithm	79
4.2	Modular-based proximal primal-gradient algorithm	80
4.2.1	Convergence analysis	83
4.3	Modular-based proximal dual-gradient algorithm	88
4.4	Sparse reconstruction and thresholding functions	91
4.5	Numerical tests	94
4.5.1	Spike reconstruction	94
4.5.2	Deconvolution of heterogeneous signals	96
4.5.3	1D and 2D mixed noise removal	97
4.5.4	A numerical study on convergence rates	100
4.6	Final discussion	104

Algorithm 5 Proximal primal-gradient algorithm in Banach spaces

Parameters: $\{\tau_k\}_k$ s.t.

$$0 < \bar{\tau} \leq \tau_k \leq \frac{\mathbf{p}(1 - \delta)}{K} \quad (4.1)$$

with $\bar{\tau} > 0$, $0 < \delta < 1$, $\mathbf{p} \in (1, 2]$, $K > 0$.

Initialisation: $x^0 \in \mathcal{X}$.

repeat

$$x^{k+1} \in \operatorname{argmin}_{x \in \mathcal{X}} \frac{1}{\mathbf{p}} \|x - x^k\|_{\mathcal{X}}^{\mathbf{p}} + \tau_k \langle \nabla f(x^k), x \rangle + \tau_k g(x) \quad (4.2)$$

until convergence

In this chapter, we analyse and propose a study on proximal-gradient algorithms to solve the composite optimisation problem

$$\operatorname{argmin}_{x \in L^{p(\cdot)}(\Omega)} \phi(x) := f(x) + g(x) \quad (\text{P})$$

where $f : L^{p(\cdot)}(\Omega) \rightarrow \mathbb{R} \cup \{+\infty\}$ is a proper, convex, and Gateaux differentiable function while $g : L^{p(\cdot)}(\Omega) \rightarrow \mathbb{R} \cup \{+\infty\}$ is a lower semi-continuous, proper, convex, and possibly non-smooth one.

4.1 Norm-based proximal-gradient algorithms in Banach spaces

The generalisation of forward-backward strategies, initially formulated in Hilbert spaces, to the minimisation of a composite problem as in (1.23) and (P) but over a Banach space \mathcal{X} is not straightforward. As in the generalisation of Landweber method seen in Chapter 3 Section 3.2, there is the need to introduce duality maps, which link primal and dual spaces, to overcome the lack of isometric isomorphism between \mathcal{X} and \mathcal{X}^* . Duality maps allow to perform the forward gradient step either in the dual or in the primal space. The intrinsic non-linearity of their duality maps, however, introduces new challenges in the definition of a backward step in terms of the proximal operator of g . In [37, 105, 106], forward-backward algorithms have been proposed to solve composite minimisation problems in a general reflexive, strictly convex and smooth Banach space \mathcal{X} , by means of suitably defined notions of duality mappings and proximal operators.

We review in this section the main tools used in these works, which will be useful to extend these approaches to $L^{p(\cdot)}(\Omega)$ [145] in the next sections of this chapter, in which we highlight our contribution .

Algorithm 6 Proximal dual-gradient algorithm in Banach spaces

Parameters: $\{\tau_k\}_k$ s.t.

$$0 < \bar{\tau} \leq \tau_k \leq \frac{pc(1-\delta)}{K} \quad \text{with } 0 < \delta < 1. \quad (4.3)$$

 with $\bar{\tau} > 0$, $0 < \delta < 1$, $p > 1$, $K > 0$, $c > 0$.

Initialisation: $x^0 \in \mathcal{X}$.

repeat

$$x^{k+1} = \operatorname{argmin}_{x \in \mathcal{X}} \frac{1}{p'} \|x^k\|_{\mathcal{X}}^p + \frac{1}{p} \|x\|_{\mathcal{X}}^p - \langle \mathbf{J}_{\mathcal{X}}^p(x^k), x \rangle + \tau_k \langle \nabla f(x^k), x \rangle + \tau_k g(x) \quad (4.4)$$

until convergence

4.1.1 Proximal primal-gradient algorithm

In [37] Bredies introduced the iterative procedure defined in Algorithm 5 for smooth functions $f : \mathcal{X} \rightarrow \mathbb{R} \cup \{+\infty\}$ with $(p-1)$ -Hölder continuous gradient ∇f on bounded sets with $1 < p \leq 2$ with constant $K > 0$, i.e.

$$\|\nabla f(u) - \nabla f(v)\|_{\mathcal{X}^*} \leq K \|u - v\|_{\mathcal{X}}^{p-1} \quad \forall u, v \in \mathcal{X},$$

and lower semi-continuous, proper, convex functions $g : \mathcal{X} \rightarrow \mathbb{R} \cup \{+\infty\}$, with appropriately chosen step-sizes τ_k as in (4.1). Notice that, for $g \equiv 0$, the updating rule just (4.2) becomes (3.15) with $r = p$, retrieving hence the primal method defined in Algorithm 2, suited for smooth optimisation procedures. Algorithm 5 can thus be seen as a generalisation of the primal method. Such interpretation is not obvious from (4.2), where forward and backward steps are defined as one single minimisation problem, so that they cannot be distinguished. Moreover, note that iteration (4.2) can be interpreted in terms of the p -norm proximal operator of the functional $x \in \mathcal{X} \mapsto \tau_k \langle \nabla f(x^k), \cdot \rangle + \tau_k g(\cdot)$:

$$x^{k+1} \in \operatorname{prox}_{\tau_k \langle \nabla f(x^k), \cdot \rangle + \tau_k g}^{1/p \|\cdot\|_{\mathcal{X}}^p}(x^k),$$

according to the definition of proximal operator given by (3.3).

In the following section, this interpretation allows us to enhance the iteration by introducing a different proximal operator, namely the Bregman-proximal operator.

4.1.2 Proximal dual-gradient algorithm

In [105], the authors consider a different forward-backward splitting algorithm in Banach spaces, where the proximal step is defined in terms of a Bregman distance, as reported in Algorithm 6 for suitably chosen step-sizes (4.3). In this second case, differing from Algorithm 5, the algorithm requires the computation of the forward step in the dual space. This can be seen in (4.4) by putting together the third

and the fourth addenda, yielding to $-\langle \mathbf{J}_{\mathcal{X}}^{\mathbf{p}}(x^k) - \tau_k \nabla f(x^k), x \rangle$. Moreover, similarly as before, for $g \equiv 0$ (4.4) becomes exactly (3.16) with $r = \mathbf{p}$, which means that Algorithm 6 is a generalisation of the dual method defined in Algorithm 1. Another interesting fact to notice is that iteration (4.4) can be equivalently defined in terms of the Bregman proximal operator of the functional $x \in \mathcal{X} \mapsto \tau_k \langle \nabla f(x^k), \cdot \rangle + \tau_k g(\cdot)$ as follows

$$x^{k+1} \in \text{prox}_{\tau_k \langle \nabla f(x^k), \cdot \rangle + \tau_k g}^{B_{\mathcal{X}}^{\mathbf{p}}}(x^k),$$

according to the definition of Bregman-proximal operator given by (3.5). In Algorithm 5, it is not possible to split the gradient step from the proximal one as they are defined together in (4.2). Instead, in Algorithm 6, the use of a Bregman distance in the definition of the algorithm (indeed, its proximal step) allows to do so. Observe indeed that x^{k+1} defined by (4.4) satisfies

$$\begin{aligned} 0 &\in \partial \left(\frac{1}{\mathbf{p}'} \|x^k\|_{\mathcal{X}}^{\mathbf{p}} + \frac{1}{\mathbf{p}} \|\cdot\|_{\mathcal{X}}^{\mathbf{p}} - \langle \mathbf{J}_{\mathcal{X}}^{\mathbf{p}}(x^k), \cdot \rangle + \tau_k \langle \nabla f(x^k), \cdot \rangle + \tau_k g(\cdot) \right) (x^{k+1}) \\ 0 &\in \mathbf{J}_{\mathcal{X}}^{\mathbf{p}}(x^{k+1}) - \mathbf{J}_{\mathcal{X}}^{\mathbf{p}}(x^k) + \tau_k \nabla f(x^k) + \tau_k \partial g(x^{k+1}) \quad (\text{Theorem 2.2.1}) \\ \mathbf{J}_{\mathcal{X}}^{\mathbf{p}}(x^k) - \tau_k \nabla f(x^k) &\in \mathbf{J}_{\mathcal{X}}^{\mathbf{p}}(x^{k+1}) + \tau_k \partial g(x^{k+1}) \\ x^{k+1} &= \left(\mathbf{J}_{\mathcal{X}}^{\mathbf{p}} + \tau_k \partial g \right)^{-1} \left(\mathbf{J}_{\mathcal{X}}^{\mathbf{p}}(x^k) - \tau_k \nabla f(x^k) \right), \end{aligned} \quad (4.5)$$

thus generalising the standard forward-backward structure of (1.25).

4.2 Modular-based proximal primal-gradient algorithm

Algorithms 5 and 6 have been proposed for reflexive, strictly convex and smooth Banach spaces, hence they can a priori be applied to solve (P) in the variable exponent Lebesgue spaces $L^{p(\cdot)}(\Omega)$. However, as previously discussed in Chapter 2 at the beginning of Sections 2.3 and 2.3.1 and in Chapter 3 in Section 3.3.1, the definitions of $\|\cdot\|_{L^{p(\cdot)}}$ and of duality maps in $L^{p(\cdot)}(\Omega)$ make their use impracticable in real applications. To overcome this issue, we propose to replace the role of the norm, naturally appearing in the definition of both the gradient step and the proximal one, by the normalised modular function (2.3). Inspired by Algorithm 5, proposed in [37], and Algorithm 6, studied in [105], we devise two instances of modular-based proximal gradient algorithms in $L^{p(\cdot)}(\Omega)$ spaces, which we have proposed in [145].

First, we propose an iterative procedure to solve the minimisation problem (P), in which the gradient step is implicitly performed in the primal space.

We set $\bar{\phi} := \inf_{x \in L^{p(\cdot)}(\Omega)} \phi(x)$, and define

$$\text{Sol(P)} := \{x \in L^{p(\cdot)}(\Omega) : \phi(x) = \bar{\phi}\} \neq \emptyset. \quad (4.6)$$

We consider the following set of assumptions:

A1. The exponent function $p(\cdot)$ is such that $1 < p_- \leq p_+ \leq 2$.

Algorithm 7 Modular-based proximal primal gradient algorithm in $L^{p(\cdot)}(\Omega)$ spaces

Parameters: $\rho \in (0, 1)$, $\{\tau_k\}_k$ s.t.

$$0 < \bar{\tau} \leq \tau_k \leq \frac{p(1-\delta)}{K} \quad \text{with } 0 < \delta < 1. \quad (4.8)$$

Initialisation: Start with $x^0 \in L^{p(\cdot)}(\Omega)$.

repeat

repeat

1. Set $\tau_k = \rho^i \tau_k$.
2. Compute the next iterate as:

$$x^{k+1} = \underset{x \in L^{p(\cdot)}(\Omega)}{\operatorname{argmin}} \bar{\rho}_{p(\cdot)}(x - x^k) + \tau_k \langle \nabla f(x^k), x \rangle + \tau_k g(x). \quad (4.9)$$

3. $i = i + 1$.

until $\rho_{p(\cdot)}(x^k - x^{k+1}) < 1$

until convergence

A2. $\nabla f : L^{p(\cdot)}(\Omega) \longrightarrow (L^{p(\cdot)}(\Omega))^*$ is $(p-1)$ Hölder-continuous with exponent $p_+ \leq p \leq 2$ and constant $K > 0$, i.e.:

$$\|\nabla f(u) - \nabla f(v)\|_{(L^{p(\cdot)}(\Omega))^*} \leq K \|u - v\|_{L^{p(\cdot)}(\Omega)}^{p-1} \quad \forall u, v \in L^{p(\cdot)}(\Omega). \quad (4.7)$$

The first instance of modular-proximal gradient algorithm is reported in Algorithm 7. The inner loop over i is needed to select at each k -th iteration of the outer loop a sufficiently small step-size τ_k such that $\rho_{p(\cdot)}(x^k - x^{k+1}) < 1$, which is required in the following convergence analysis, as we will see in Proposition 4.2.3 and Lemma 4.2.2. It should be thought of as a backtracking-like procedure affecting just the first algorithmic iterations where the quantity $\rho_{p(\cdot)}(x^k - x^{k+1})$ is likely to be large.

Observe that the iteration step defined by (4.2) can be interpreted as the computation of a proximal operator defined with respect to the distance induced by the modular $\bar{\rho}_{p(\cdot)}$ of the functional $x \in L^{p(\cdot)}(\Omega) \mapsto \tau_k \langle \nabla f(x^k), \cdot \rangle + \tau_k g(\cdot)$:

$$x^{k+1} = \operatorname{prox}_{\tau_k \langle \nabla f(x^k), \cdot \rangle + \tau_k g}^{\bar{\rho}_{p(\cdot)}}(x^k).$$

The analysis of Algorithm 7 is related to the study of fixed points of iterations of (4.10). We have the following result.

Proposition 4.2.1. *The solutions of (P) coincide with the fixed points of the iteration defined by Algorithm 7.*

Proof. Suppose that, for some $k \geq 0$, x^k is a solution of (P), i.e. $x^k \in \text{Sol(P)}$ defined in (4.6). Then, by optimality condition, $-\nabla f(x^k) \in \partial g(x^k)$. Note that, in general, u solves (4.9) if and only if the following inclusion holds

$$-\tau_k \nabla f(x^k) \in \mathbf{J}_{\bar{\rho}_{p(\cdot)}}(u - x^k) + \tau_k \partial g(u),$$

where $\mathbf{J}_{\bar{\rho}_{p(\cdot)}}$ is the derivative of $\bar{\rho}_{p(\cdot)}$, according to Proposition 2.3.1. Clearly, $u = x^k$ is thus a solution of (4.9), hence $x^{k+1} = x^k$, so x^k is a fixed point of the iteration process.

Conversely, suppose now that x^k is a fixed point, that is $x^k = x^{k+1} \in L^{p(\cdot)}(\Omega)$, then

$$-\tau_k \nabla f(x^k) \in \mathbf{J}_{\bar{\rho}_{p(\cdot)}}(x^{k+1} - x^k) + \tau_k \partial g(x^{k+1}) = \tau_k \partial g(x^{k+1}) = \tau_k \partial g(x^k),$$

meaning that x^k is optimal, that is, $x^k \in \text{Sol(P)}$. \square

We start our analysis discussing the well-definition of the step (4.9), e.g. the existence and uniqueness of the minimizer of the functional defined by (4.9).

Proposition 4.2.2. *For each $x \in L^{p(\cdot)}(\Omega)$, $v^* \in (L^{p(\cdot)}(\Omega))^*$ and $\tau > 0$, the problem*

$$\operatorname{argmin}_{u \in L^{p(\cdot)}(\Omega)} \bar{\rho}_{p(\cdot)}(u - x) + \tau \langle v^*, u \rangle + \tau g(u) \quad (4.10)$$

has a unique solution.

Proof. Let $\tau > 0$. Note that when $\|u - x\|_{L^{p(\cdot)}} > 1$, by Lemma 2.1.1 we can write:

$$\bar{\rho}_{p(\cdot)}(u - x) \geq \frac{1}{p_+} \rho_{p(\cdot)}(u - x) \geq \frac{1}{p_+} \|u - x\|_{L^{p(\cdot)}}^{p_-}.$$

Let now $\bar{x} \in \text{Sol(P)}$. The optimality condition reads as $0 \in \nabla f(\bar{x}) + \partial g(\bar{x})$ or, equivalently, $\bar{\omega} := -\nabla f(\bar{x}) \in \partial g(\bar{x})$. By definition of subdifferential, there holds $g(u) \geq g(\bar{x}) + \langle \bar{\omega}, u - \bar{x} \rangle = g(\bar{x}) + \langle \bar{\omega}, u \rangle - \langle \bar{\omega}, \bar{x} \rangle$ for all $u \in L^{p(\cdot)}(\Omega)$. By combining such inequality with the Cauchy-Schwarz inequality, we get

$$\begin{aligned} & \bar{\rho}_{p(\cdot)}(u - x) + \tau \langle v^*, u \rangle + \tau g(u) \\ & \geq \frac{1}{p_+} \|u - x\|_{L^{p(\cdot)}}^{p_-} + \tau \langle v^* + \bar{\omega}, u \rangle + \tau g(\bar{x}) - \tau \langle \bar{\omega}, \bar{x} \rangle \\ & \geq \|u\|_{L^{p(\cdot)}} \left[\frac{1}{p_+} \frac{\|u - x\|_{L^{p(\cdot)}}^{p_-}}{\|u\|_{L^{p(\cdot)}}} + \tau \frac{\langle v^* + \bar{\omega}, u \rangle}{\|u\|_{L^{p(\cdot)}}} + \frac{\tau g(\bar{x}) - \tau \langle \bar{\omega}, \bar{x} \rangle}{\|u\|_{L^{p(\cdot)}}} \right] \\ & \geq \|u\|_{L^{p(\cdot)}} \left[\frac{1}{p_+} \frac{\|u - x\|_{L^{p(\cdot)}}^{p_-}}{\|u\|_{L^{p(\cdot)}}} - \tau \|v^* + \bar{\omega}\|_{(L^{p(\cdot)})^*} + \tau \frac{g(\bar{x}) - \langle \bar{\omega}, \bar{x} \rangle}{\|u\|_{L^{p(\cdot)}}} \right] \geq L \|u\|_{L^{p(\cdot)}} \end{aligned}$$

for some $L > 0$ and all $u \in L^{p(\cdot)}(\Omega)$ such that $\|u\|_{L^{p(\cdot)}}$ is large enough. Note, in particular, that $\frac{1}{p_+} \frac{\|u - x\|_{L^{p(\cdot)}}^{p_-}}{\|u\|_{L^{p(\cdot)}}} \rightarrow +\infty$ as $\|u\|_{L^{p(\cdot)}} \rightarrow +\infty$ since $p_- > 1$ and x is fixed. Hence, the functional in (4.10) is coercive. Moreover, it is convex, proper and lower semi-continuous in $L^{p(\cdot)}(\Omega)$, which, by Theorem 2.1.2, is reflexive, and thus at least one solution exists. Moreover, since $1 < p_- \leq p_+ \leq 2$, the functional is strictly convex, and hence the solution is unique. \square

4.2.1 Convergence analysis

We provide here a detailed convergence analysis of Algorithm 7 and provide an insight on its convergence speed in function values. Our analysis is inspired by the one conducted in [37], although it relies on different properties related to the modular function $\bar{\rho}_{p(\cdot)}$ rather than to $\|\cdot\|_{L^{p(\cdot)}}$.

Lemma 4.2.1. *For each $u \in L^{p(\cdot)}(\Omega)$, the following inequality holds true:*

$$\left\langle \frac{\mathbf{J}_{\bar{\rho}_{p(\cdot)}}(x^k - x^{k+1})}{\tau_k}, u - x^{k+1} \right\rangle \leq g(u) - g(x^{k+1}) + \langle \nabla f(x^k), u - x^{k+1} \rangle. \quad (4.11)$$

Moreover, denoting by $D(x^k, x^{k+1})$ the quantity

$$D(x^k, x^{k+1}) := g(x^k) - g(x^{k+1}) + \langle \nabla f(x^k), x^k - x^{k+1} \rangle, \quad (4.12)$$

we have that

$$\rho_{p(\cdot)}(x^k - x^{k+1}) \leq \tau_k D(x^k, x^{k+1}). \quad (4.13)$$

Proof. Note that, by optimality condition, x^{k+1} solves (4.9) if and only if

$$\begin{aligned} 0 \in \mathbf{J}_{\bar{\rho}_{p(\cdot)}}(x^{k+1} - x^k) + \tau_k \nabla f(x^k) + \tau_k \partial g(x^{k+1}) &\iff \\ \frac{\mathbf{J}_{\bar{\rho}_{p(\cdot)}}(x^k - x^{k+1})}{\tau_k} - \nabla f(x^k) \in \partial g(x^{k+1}). \end{aligned}$$

By definition of subdifferential, we thus have that, for all $u \in L^{p(\cdot)}(\Omega)$,

$$\left\langle \frac{\mathbf{J}_{\bar{\rho}_{p(\cdot)}}(x^k - x^{k+1})}{\tau_k} - \nabla f(x^k), u - x^{k+1} \right\rangle \leq g(u) - g(x^{k+1}),$$

which, by rearranging, coincides with (4.11). Choosing now $u = x^k$ above and recalling (4.12), we get $\left\langle \frac{\mathbf{J}_{\bar{\rho}_{p(\cdot)}}(x^k - x^{k+1})}{\tau_k}, x^k - x^{k+1} \right\rangle \leq D(x^k, x^{k+1})$. Applying now (2.26) entails

$$\left\langle \frac{\mathbf{J}_{\bar{\rho}_{p(\cdot)}}(x^k - x^{k+1})}{\tau_k}, x^k - x^{k+1} \right\rangle = \frac{\rho_{p(\cdot)}(x^k - x^{k+1})}{\tau_k},$$

by which (4.13) follows directly. \square

The following proposition shows that the iteration scheme (4.9) leads to a decrease of the functional ϕ of our minimization problem (P). This will be crucial for the following convergence analysis.

Proposition 4.2.3. *For every $k \geq 0$, if $\rho_{p(\cdot)}(x^k - x^{k+1}) < 1$, then the iteration defined by Algorithm 7 satisfies*

$$\phi(x^{k+1}) \leq \phi(x^k) - \left(1 - \frac{K\tau_k}{p}\right) D(x^k, x^{k+1}). \quad (4.14)$$

Proof. By (4.12), we have

$$\phi(x^k) - \phi(x^{k+1}) = D(x^k, x^{k+1}) + f(x^k) - f(x^{k+1}) - \langle \nabla f(x^k), x^k - x^{k+1} \rangle. \quad (4.15)$$

Considering the last three terms, we can write

$$\begin{aligned} & f(x^k) - f(x^{k+1}) - \langle \nabla f(x^k), x^k - x^{k+1} \rangle \\ &= \int_0^1 \langle \nabla f(x^k + t(x^{k+1} - x^k)) - \nabla f(x^k), x^k - x^{k+1} \rangle dt. \end{aligned} \quad (4.16)$$

By applying now the $(p-1)$ -Hölder continuity of ∇f (4.7), we can provide an estimate of the absolute value of the right-hand side of (4.16), since

$$\begin{aligned} & \left| \int_0^1 \langle \nabla f(x^k + t(x^{k+1} - x^k)) - \nabla f(x^k), x^k - x^{k+1} \rangle dt \right| \\ & \leq \int_0^1 \|\nabla f(x^k + t(x^{k+1} - x^k)) - \nabla f(x^k)\|_{(L^{p(\cdot)})^*} \|x^k - x^{k+1}\|_{L^{p(\cdot)}} dt \\ & \leq \int_0^1 K \|x^k - x^{k+1}\|_{L^{p(\cdot)}}^p t^{p-1} dt \leq \frac{K}{p} \|x^k - x^{k+1}\|_{L^{p(\cdot)}}^p \end{aligned}$$

Since we have $\rho_{p(\cdot)}(x^k - x^{k+1}) < 1$ by assumption, by Proposition 2.1.2 there holds $\|x^k - x^{k+1}\|_{L^{p(\cdot)}} < 1$ and $\|x^k - x^{k+1}\|_{L^{p(\cdot)}} < \rho_{p(\cdot)}(x^k - x^{k+1})^{1/p+} \leq \rho_{p(\cdot)}(x^k - x^{k+1})^{1/p}$ by Lemma 2.1.1. Hence, by (4.13) we obtain

$$\begin{aligned} & \left| \int_0^1 \langle \nabla f(x^k + t(x^{k+1} - x^k)) - \nabla f(x^k), x^k - x^{k+1} \rangle dt \right| \\ & \leq \frac{K}{p} \|x^k - x^{k+1}\|_{L^{p(\cdot)}}^p \leq \frac{K}{p} \rho_{p(\cdot)}(x^k - x^{k+1}) \leq \frac{K\tau_k}{p} D(x^k, x^{k+1}), \end{aligned}$$

which concludes the proof by combining this with (4.15) and (4.16). \square

For each $k \geq 1$, let us now define for simplicity the k -th residual

$$r_k := \phi(x^k) - \bar{\phi}, \quad \bar{\phi} = \inf_{x \in L^{p(\cdot)}(\Omega)} \phi(x). \quad (4.17)$$

Note that $r_k \geq 0$ by definition. We can thus rewrite (4.14) as

$$r_k - r_{k+1} \geq \left(1 - \frac{K\tau_k}{p}\right) D(x^k, x^{k+1}). \quad (4.18)$$

Thanks to the bounds on the step-sizes τ_k , there holds $r_k - r_{k+1} \geq 0$, hence the descent of the functional $\phi = f + g$ is guaranteed.

The following lemma shows that, by assuming the boundedness of the sequence $(x^k)_k$, an estimate of the right-hand side of (4.18) can be found. This estimate depends on the conjugate exponent $p'(\cdot) \in \mathcal{P}(\Omega)$ defined by (2.5). Since $1 < p(\cdot) \leq 2$ and $\frac{1}{p(t)} + \frac{1}{p'(t)} = 1$ a.e., there holds $2 \leq p'(\cdot) < +\infty$.

Lemma 4.2.2. *If $(x^k)_k$ is bounded, then $r_k - r_{k+1} \geq c_0 r_k^{(p-)'}$ with $c_0 > 0$.*

Proof. Since $(x^k)_k$ is bounded, for some $\bar{x} \in \text{Sol}(\text{P})$ there exists $C_1 > 0$ such that, for all $k \geq 1$, $\|x^k - \bar{x}\|_{L^{p(\cdot)}} \leq C_1$. The convexity of f as well as (4.11) with $u = \bar{x}$ gives

$$\begin{aligned}
 r_k &= f(x^k) + g(x^k) - f(\bar{x}) - g(\bar{x}) \leq g(x^k) - g(\bar{x}) + \langle \nabla f(x^k), x^k - \bar{x} \rangle \\
 &= \langle \nabla f(x^k), x^k - x^{k+1} \rangle + g(x^k) - g(\bar{x}) + \langle \nabla f(x^k), x^{k+1} - \bar{x} \rangle \\
 &= D(x^k, x^{k+1}) + g(x^{k+1}) - g(\bar{x}) + \langle \nabla f(x^k), x^{k+1} - \bar{x} \rangle \\
 &\leq D(x^k, x^{k+1}) + \left\langle \frac{\mathbf{J}_{\bar{\rho}_{p(\cdot)}}(x^k - x^{k+1})}{\tau_k}, x^{k+1} - \bar{x} \right\rangle \\
 &\leq D(x^k, x^{k+1}) + \tau_k^{-1} \|\mathbf{J}_{\bar{\rho}_{p(\cdot)}}(x^k - x^{k+1})\|_{(L^{p(\cdot)})^*} \|x^{k+1} - \bar{x}\|_{L^{p(\cdot)}} \\
 &\leq D(x^k, x^{k+1}) + \tau_k^{-1} \|\mathbf{J}_{\bar{\rho}_{p(\cdot)}}(x^k - x^{k+1})\|_{(L^{p(\cdot)})^*} C_1. \tag{4.19}
 \end{aligned}$$

Recalling now that $\langle \mathbf{J}_{\bar{\rho}_{p(\cdot)}}(u), v \rangle = \int_{\Omega} \text{sign}(u(t)) |u(t)|^{p(t)-1} v(t) dt$ for any $u, v \in L^{p(\cdot)}(\Omega)$, by (2.9) and Proposition 2.2.1 we have

$$\|\mathbf{J}_{\bar{\rho}_{p(\cdot)}}(u)\|_{(L^{p(\cdot)})^*} = \|\text{sign}(u) |u|^{p(\cdot)-1}\|'_{p(\cdot)} \leq 2 \|\text{sign}(u) |u|^{p(\cdot)-1}\|_{L^{p(\cdot)}}.$$

Observe now that the following equality holds true:

$$\begin{aligned}
 \rho_{p(\cdot)}(\text{sign}(u) |u|^{p(\cdot)-1}) &= \int_{\Omega} \left(\text{sign}(u(t)) |u(t)|^{p(t)-1} \right)^{p'(t)} dt \\
 &= \int_{\Omega} \left(\text{sign}(u(t)) |u(t)| \right)^{p(t)} = \int_{\Omega} |u(t)|^{p(t)} dt = \rho_{p(\cdot)}(u).
 \end{aligned}$$

Hence it follows that $\rho_{p(\cdot)}(\text{sign}(x^k - x^{k+1}) |x^k - x^{k+1}|^{p(\cdot)-1}) < 1$, since by construction we have $\rho_{p(\cdot)}(x^k - x^{k+1}) < 1$. Together with Lemma 2.1.1, this leads to $\|\text{sign}(x^k - x^{k+1}) |x^k - x^{k+1}|^{p(\cdot)-1}\|_{L^{p(\cdot)}} < 1$. Furthermore, keeping in mind that $(p')_+ = (p_-)'$ as in (2.6), Lemma 2.1.1 entails

$$\|\text{sign}(u) |u|^{p(\cdot)-1}\|_{L^{p(\cdot)}} \leq \left(\rho_{p(\cdot)}(\text{sign}(u) |u|^{p(\cdot)-1}) \right)^{1/(p_-)'}$$

which can now be evaluated in $u = x^k - x^{k+1}$ and combined with the previous inequalities to get from (4.19)

$$\begin{aligned}
 r_k &\leq D(x^k, x^{k+1}) + \tau_k^{-1} \|\mathbf{J}_{\bar{\rho}_{p(\cdot)}}(x^k - x^{k+1})\|_{(L^{p(\cdot)})^*} C_1 \\
 &\leq D(x^k, x^{k+1}) + \tau_k^{-1} C_1 2 \|\text{sign}(x^k - x^{k+1}) |x^k - x^{k+1}|^{p(\cdot)-1}\|_{L^{p(\cdot)}} \\
 &\leq D(x^k, x^{k+1}) + \tau_k^{-1} C_1 2 \left(\rho_{p(\cdot)}(\text{sign}(x^k - x^{k+1}) |x^k - x^{k+1}|^{p(\cdot)-1}) \right)^{1/(p_-)' } \\
 &= D(x^k, x^{k+1}) + \tau_k^{-1} C_1 2 \left(\rho_{p(\cdot)}(x^k - x^{k+1}) \right)^{1/(p_-)' },
 \end{aligned}$$

so that, by (4.13), we have $r_k \leq D(x^k, x^{k+1}) + 2C_1 \tau_k^{-1} \left(\tau_k D(x^k, x^{k+1}) \right)^{1/(p_-)'}$. The step-size constraints (4.8) together with (4.14) entail $r_k - r_{k+1} \geq \delta D(x^k)$ and $\tau_k \geq \bar{\tau}$.

Plugging these quantities into the last inequality above, we obtain

$$r_k \leq \frac{r_k - r_{k+1}}{\delta} + 2C_1 \bar{r}^{-\frac{(p_-)'-1}{(p_-)'}} \left(\frac{r_k - r_{k+1}}{\delta} \right)^{1/(p_-)'}$$

Note that since r_k is a nonnegative decreasing sequence, then $r_k - r_{k+1}$ is bounded, so that we can write

$$\delta r_k \leq (r_k - r_{k+1})^{1/(p_-)' } \left(R^{-\frac{(p_-)'-1}{(p_-)'}} + 2C_1 \bar{r}^{-\frac{(p_-)'-1}{(p_-)'}} \delta^{\frac{(p_-)'-1}{(p_-)'}} \right),$$

for a sufficiently large $R > 0$, which finally gives

$$r_k - r_{k+1} \geq \frac{\delta^{(p_-)'}}{\left(R^{-\frac{(p_-)'-1}{(p_-)'}} + 2C_1 \bar{r}^{-\frac{(p_-)'-1}{(p_-)'}} \delta^{\frac{(p_-)'-1}{(p_-)'}} \right)^{(p_-)'}} r_k^{(p_-)'}$$

□

Thanks to the previous lemma, we obtain the following convergence result.

Proposition 4.2.4. *If $(x^k)_k$ is bounded, then the following convergence rate in function values can be found for the iterates of Algorithm 7*

$$r_k \leq \eta \frac{1}{k^{p_- - 1}}, \quad (4.20)$$

where $\eta = \eta(\bar{r}, \delta, p_-, x^0, K, \bar{\phi})$ and the residual r_k is defined by (4.17).

Proof. The proof is based on analogous arguments as in [37, Proposition 4].

Apply the mean value theorem to get the identity

$$\frac{1}{r_{k+1}^{(p_-)'-1}} - \frac{1}{r_k^{(p_-)'-1}} = \frac{r_k^{(p_-)'-1} - r_{k+1}^{(p_-)'-1}}{(r_{k+1} r_k)^{(p_-)'-1}} = \frac{((p_-)' - 1) \vartheta^{(p_-)'-2} (r_k - r_{k+1})}{(r_{k+1} r_k)^{(p_-)'-1}}$$

with $r_{k+1} \leq \vartheta \leq r_k$. Thus, $\vartheta^{(p_-)'-2} \geq r_{k+1}^{(p_-)'-1} r_k^{-1}$ and, by Lemma 4.2.2,

$$\frac{1}{r_{k+1}^{(p_-)'-1}} - \frac{1}{r_k^{(p_-)'-1}} \geq \frac{((p_-)' - 1) c_0 r_{k+1}^{(p_-)'-1} r_k^{(p_-)'-1}}{(r_{k+1} r_k)^{(p_-)'-1}} = ((p_-)' - 1) c_0.$$

Summing up the above telescopic sequence, then yields

$$\frac{1}{r_k^{(p_-)'-1}} - \frac{1}{r_0^{(p_-)'-1}} = \sum_{i=0}^{k-1} \frac{1}{r_{i+1}^{(p_-)'-1}} - \frac{1}{r_i^{(p_-)'-1}} \geq k ((p_-)' - 1) c_0$$

and consequently,

$$r_k^{(p_-)'-1} \leq \left(r_0^{1-(p_-)'} + c_0 ((p_-)' - 1) k \right)^{-1} \Rightarrow r_k \leq C n^{1-p_-}$$

since $1/(1 - (p_-)') = 1 - p_-$.

□

Such convergence rate is related to the smoothness of the considered space $L^{p(\cdot)}(\Omega)$ and, in particular, to the infimum exponent value p_- appearing in the analytical proof of Lemma 4.2.2, where it is used in a majorisation as worst case analysis. It can thus be read as the worst-case convergence speed. It is highly expected that such convergence result can be improved and that, practically, faster convergence could be achieved, as our numerical tests will effectively show.

We now provide a result relative to the convergence of the sequence $(x_k)_k$ itself.

Proposition 4.2.5. *If $(x^k)_k$ is bounded, then the sequence $(x^k)_k$ has at least one accumulation point. All accumulation points belong to $\text{Sol}(P)$. If $\text{Sol}(P) = \{\bar{x}\}$, then $(x^k)_k$ converges weakly to \bar{x} .*

Proof. The proof is based on similar arguments to [37, Proposition 5] and [105, Proposition 3.4 (iii)].

Since $r_k \leq k^{1-p_-}$, the sequence is a minimising sequence, thus, due to the weak lower semi-continuity of the functional ϕ , each weakly convergent sub-sequence is a minimiser. Moreover, it also follows that $(x^k)_k$ is a bounded sequence in the reflexive Banach space $L^{p(\cdot)}(\Omega)$, meaning that there is a weakly-convergent sub-sequence. The statement that $(x^k)_k$ converges weakly to \bar{x} in case of uniqueness follows by the usual sub-sequence argument. \square

We conclude this section recalling the definition of totally convex and r -convex functions. Under the further hypothesis that f or g in (P) are totally convex or r -convex, it is possible to show that the sequence of iterates defined by Algorithm 7 converges strongly to a solution of the minimisation problem (P) and, moreover, that (P) has a unique solution, as we prove in Proposition 4.2.6.

Definition 4.2.1. *Let $h : L^{p(\cdot)}(\Omega) \rightarrow \mathbb{R} \cup \{+\infty\}$ be proper, convex and lower semi-continuous. The functional h is said to be totally convex in $\hat{u} \in L^{p(\cdot)}(\Omega)$ if for all $\omega \in \partial h(\hat{u})$ and for each $(u^n)_n$ such that*

$$h(u^n) - h(\hat{u}) - \langle \omega, u^n - \hat{u} \rangle \rightarrow 0,$$

there holds

$$\|u^n - \hat{u}\|_{L^{p(\cdot)}} \rightarrow 0, \quad \text{for } n \rightarrow +\infty.$$

We say that h is totally convex if it is totally convex in \hat{u} for all $\hat{u} \in L^{p(\cdot)}(\Omega)$.

Similarly, h is convex of power-type r (or r -convex) in $\hat{u} \in L^{p(\cdot)}(\Omega)$, with $r \geq 2$, if for all $M > 0$ and $\omega \in \partial h(\hat{u})$ there exists $\beta > 0$ such that for all $\|u - \hat{u}\|_{L^{p(\cdot)}} \leq M$

$$h(u) - h(\hat{u}) - \langle \omega, u - \hat{u} \rangle \geq \beta \|u - \hat{u}\|_{L^{p(\cdot)}}^r.$$

We say that h is convex of power-type r if it is convex of power-type r in \hat{u} for all $\hat{u} \in L^{p(\cdot)}(\Omega)$.

Proposition 4.2.6. *If f or g is totally convex or r -convex with $r \geq 2$, the solution of problem (P) is unique. In this case, denoting it by $\bar{x} \in L^{p(\cdot)}(\Omega)$, we further have that the sequence $(x^k)_k$ defined by (4.9) converges strongly to \bar{x} .*

Proof. First, note that r -convexity implies total convexity. Then, focus on the case that f is totally convex. Suppose now by contradiction that there exists $\tilde{x} \in \text{Sol}(\text{P})$ with $\tilde{x} \neq \bar{x}$. Thus $\|\tilde{x} - \bar{x}\|_{L^{p(\cdot)}} > 0$. By defining

$$R(x) := f(x) - f(\bar{x}) - \langle \nabla f(\bar{x}), x - \bar{x} \rangle, \quad (4.21)$$

we have that $R(x) \geq 0$ for all $x \in L^{p(\cdot)}(\Omega)$. Moreover, by the optimality of \bar{x} and the subgradient inequality, there holds $\phi(x) - \phi(\bar{x}) \geq R(x)$ for all $x \in L^{p(\cdot)}(\Omega)$. Choosing $x = \tilde{x}$, we thus get $0 = \phi(\tilde{x}) - \phi(\bar{x}) \geq R(\tilde{x}) \geq 0$, whence $R(\tilde{x}) = 0$. Taking now $u^n = \tilde{x}$ for all $n \geq 1$, we find a contradiction as the total convexity property is violated, whence we deduce $\tilde{x} = \bar{x}$, that is the solution of (P) is unique.

To complete the proof, let us consider (4.21) once again. By optimality of \bar{x} and thanks to the subgradient inequality, there holds $r_k \geq R(x^k)$, whence, by Proposition 4.2.4, we deduce that $R(x^k) \leq \eta \frac{1}{k^p - 1}$. Letting now $k \rightarrow +\infty$ we thus infer that $R(x^k) \rightarrow 0$ and, by the total convexity of f , that $\|x^k - \bar{x}\|_{L^{p(\cdot)}} \rightarrow 0$, which completes the proof.

In the case that g is totally convex, the proof is analogous by defining $R(x) := g(x) - g(\bar{x}) + \langle \nabla f(\bar{x}), x - \bar{x} \rangle$. \square

4.3 Modular-based proximal dual-gradient algorithm

In this section, we introduce a different modular-based proximal gradient algorithm solving (P) where the proximal step is defined in terms of a modular Bregman-like distance. Our study is here inspired by the analysis carried out in [105] where Algorithm 6 is studied for a general Banach space \mathcal{X} . Similarly as above, we start this section by stating the required assumptions:

A1. $\nabla f : L^{p(\cdot)}(\Omega) \rightarrow (L^{p(\cdot)}(\Omega))^*$ is $(p-1)$ Hölder-continuous with $1 < p \leq 2$ with constant K , as in (4.7).

A2. There exists $c > 0$ such that for all $u, v \in L^{p(\cdot)}(\Omega)$

$$\langle \mathbf{J}_{\bar{\rho}_{p(\cdot)}}(u) - \mathbf{J}_{\bar{\rho}_{p(\cdot)}}(v), u - v \rangle \geq c \max \left\{ \|u - v\|_{L^{p(\cdot)}}^p, \|\mathbf{J}_{\bar{\rho}_{p(\cdot)}}(u) - \mathbf{J}_{\bar{\rho}_{p(\cdot)}}(v)\|_{(L^{p(\cdot)})^*}^{p'} \right\}. \quad (4.22)$$

Condition (4.22) links the geometrical properties of the space $L^{p(\cdot)}(\Omega)$ with the Hölder smoothness properties of f . It has to be interpreted as a sufficient compatibility condition between the ambient space $L^{p(\cdot)}(\Omega)$ and the function f for achieving the desired convergence result. Notice that Condition (4.22) is similar to Xu-Roach inequalities for p -convex spaces (see [204]). The pseudocode of the proposed algorithm is reported in Algorithm 8 with iteration step defined by (4.23) or, equivalently, by

$$x^{k+1} = \operatorname{argmin}_{x \in L^{p(\cdot)}(\Omega)} B_{\bar{\rho}_{p(\cdot)}}(u, x^k) + \tau_k \langle \nabla f(x^k), x \rangle + \tau_k g(x),$$

since $B_{\bar{\rho}_{p(\cdot)}}(u, x^k) = \bar{\rho}_{p(\cdot)}(u) - \bar{\rho}_{p(\cdot)}(x^k) - \langle \mathbf{J}_{\bar{\rho}_{p(\cdot)}}(x^k), u - x^k \rangle$, by neglecting constant terms with respect to u . From the last equation, it result evident that (4.23)

Algorithm 8 Modular-based proximal dual gradient algorithm in $L^{p(\cdot)}(\Omega)$ spaces

Parameters: $\{\tau_k\}_k$ s.t.

$$0 < \bar{\tau} \leq \tau_k \leq \frac{pc(1-\delta)}{K} \quad \text{with } 0 < \delta < 1.$$

Initialisation: Start with $x^0 \in L^{p(\cdot)}(\Omega)$.

repeat

Compute the next iterate as:

$$x^{k+1} = \underset{u \in L^{p(\cdot)}(\Omega)}{\operatorname{argmin}} \bar{\rho}_{p(\cdot)}(u) - \langle \mathbf{J}_{\bar{\rho}_{p(\cdot)}}(x^k), u \rangle + \tau_k \langle \nabla f(x^k), u \rangle + \tau_k g(u) \quad (4.23)$$

until convergence

can be interpreted by means of the proximal operator defined in terms of the Bregman distance induced by the modular $\bar{\rho}_{p(\cdot)}$ of the functional $x \in L^{p(\cdot)}(\Omega) \mapsto \tau_k \langle \nabla f(x^k), \cdot \rangle + \tau_k g(\cdot)$:

$$x^{k+1} = \operatorname{PROX}_{\tau_k \langle \nabla f(x^k), \cdot \rangle + \tau_k g}^{B_{\bar{\rho}_{p(\cdot)}}}(x^k).$$

Moreover, the iteration (4.23) can be equivalently formulated as

$$\begin{aligned} 0 &\in \mathbf{J}_{\bar{\rho}_{p(\cdot)}}(x^{k+1}) - \mathbf{J}_{\bar{\rho}_{p(\cdot)}}(x^k) + \tau_k \nabla f(x^k) + \tau_k \partial g(x^{k+1}) \\ \mathbf{J}_{\bar{\rho}_{p(\cdot)}}(x^k) - \tau_k \nabla f(x^k) &\in \mathbf{J}_{\bar{\rho}_{p(\cdot)}}(x^{k+1}) + \tau_k \partial g(x^{k+1}) \\ x^{k+1} &= \left(\mathbf{J}_{\bar{\rho}_{p(\cdot)}} + \tau_k \partial g \right)^{-1} \left(\mathbf{J}_{\bar{\rho}_{p(\cdot)}}(x^k) - \tau_k \nabla f(x^k) \right). \end{aligned} \quad (4.24)$$

This result is similar to the one shown by Lemma 3.1.6 for the proximal operator defined in terms of the norm and of the one in (4.5) for the proximal dual-gradient algorithm. It allows to better identify and split the forward gradient step and the backward proximal one of the iterative strategy. To clarify further this fact, it is useful to introduce the following notion which shows analogies to the standard scheme of the Moreau envelope (see, e.g., [9, Chapter 12]).

Definition 4.3.1. Given $h : L^{p(\cdot)}(\Omega) \rightarrow \mathbb{R} \cup \{+\infty\}$ smooth, convex, proper, lower semi-continuous, we define the Moreau-like envelope $e_h : (L^{p(\cdot)}(\Omega))^* \rightarrow \mathbb{R}$ and the modular-proximal mapping $\pi_h : (L^{p(\cdot)}(\Omega))^* \rightarrow L^{p(\cdot)}(\Omega)$ as follows

$$\begin{aligned} e_h(x^*) &:= \inf_{u \in L^{p(\cdot)}(\Omega)} \Delta(x^*, u) + h(u), \quad x^* \in (L^{p(\cdot)}(\Omega))^* \\ \pi_h(x^*) &:= \underset{u \in L^{p(\cdot)}(\Omega)}{\operatorname{argmin}} \Delta(x^*, u) + h(u), \quad x^* \in (L^{p(\cdot)}(\Omega))^* \end{aligned} \quad (4.25)$$

where $\Delta(x^*, u) = \bar{\rho}_{p(\cdot)}(u) - \langle x^*, u \rangle$ denotes the Bregman-like distance associated to $\bar{\rho}_{p(\cdot)}$.

Note that the minimum in (4.25) is uniquely attained as $\bar{\rho}_{p(\cdot)}$ is a strictly convex function. Moreover, the unique point $\pi_h(x^*)$ satisfies

$$\begin{aligned} 0 &\in \mathbf{J}_{\bar{\rho}_{p(\cdot)}}(\pi_h(x^*)) - x^* + \partial h(\pi_h(x^*)) \iff \\ x^* &\in \mathbf{J}_{\bar{\rho}_{p(\cdot)}}(\pi_h(x^*)) + \partial h(\pi_h(x^*)) \iff \\ \pi_h(x^*) &= \left(\mathbf{J}_{\bar{\rho}_{p(\cdot)}} + \partial h\right)^{-1}(x^*). \end{aligned} \quad (4.26)$$

This allows to write (4.24), which coincides with (4.23), equivalently as

$$x^{k+1} = \pi_{\tau_k g} \left(\mathbf{J}_{\bar{\rho}_{p(\cdot)}}(x^k) - \tau_k \nabla f(x^k) \right). \quad (4.27)$$

Thus, this shows that the iteration of Algorithm 8 can be equivalently formulated as (4.23), (4.24) or (4.27). The next proposition shows that we can interpret the proposed iteration as a fixed-point iteration scheme.

Proposition 4.3.1. *For any $\gamma > 0$, there holds*

$$\bar{x} \in \text{Sol}(P) \iff \bar{x} = \pi_{\gamma g} \left(\mathbf{J}_{\bar{\rho}_{p(\cdot)}}(\bar{x}) - \gamma \nabla f(\bar{x}) \right).$$

Proof. Since \bar{x} solves (P), we have

$$0 \in \nabla f(\bar{x}) + \partial g(\bar{x}) \iff 0 \in \gamma \nabla f(\bar{x}) - \mathbf{J}_{\bar{\rho}_{p(\cdot)}}(\bar{x}) + \mathbf{J}_{\bar{\rho}_{p(\cdot)}}(\bar{x}) + \gamma \partial g(\bar{x}),$$

that is $\mathbf{J}_{\bar{\rho}_{p(\cdot)}}(\bar{x}) - \gamma \nabla f(\bar{x}) \in \mathbf{J}_{\bar{\rho}_{p(\cdot)}}(\bar{x}) + \gamma \partial g(\bar{x}) = \left(\mathbf{J}_{\bar{\rho}_{p(\cdot)}} + \gamma \partial g\right)(\bar{x})$. By (4.26), since $\pi_{\gamma g} = \left(\mathbf{J}_{\bar{\rho}_{p(\cdot)}} + \gamma \partial g\right)^{-1}$, we thus deduce as required that $\bar{x} = \pi_{\gamma g} \left(\mathbf{J}_{\bar{\rho}_{p(\cdot)}}(\bar{x}) - \gamma \nabla f(\bar{x}) \right)$. \square

We now report the analogue of Proposition 4.2.2, omitting the proof since the reasoning is straightforwardly similar. It shows that the minimisation problem (4.23) has a unique solution at each iteration, and thus each new iteration x^{k+1} is well defined.

Proposition 4.3.2. *The problem*

$$\operatorname{argmin}_{u \in L^{p(\cdot)}(\Omega)} \bar{\rho}_{p(\cdot)}(u) - \langle \mathbf{J}_{\bar{\rho}_{p(\cdot)}}(x), u \rangle + \tau \langle v^*, u \rangle + \tau g(u)$$

has a unique solution for each $x \in L^{p(\cdot)}(\Omega)$, $v^ \in (L^{p(\cdot)}(\Omega))^*$ and $\tau > 0$.*

Similarly as in Proposition 4.2.4, it is possible to prove the convergence in function values for the iterates of Algorithm 8 and achieve a convergence rate. Here, we omit the proof, as it follows verbatim the one in [105] which is itself inspired by [37].

Proposition 4.3.3. *If $(x^k)_k$ is bounded, then the following convergence rate in function values can be found for the iterates of Algorithm 8*

$$r_k \leq \eta \frac{1}{k^{p-1}}, \quad (4.28)$$

where $\eta = \eta(\bar{\tau}, \delta, \mathbf{p}, x^0, K, c, \bar{\phi})$.

It is interesting to compare the rate (4.28) with the analogous one in (4.20) obtained for Algorithm 7. The dependence on the Hölder exponent \mathbf{p} in (4.28) links the speed of convergence to the smoothness of the smooth function f rather than to the one of the underlying $L^{p(\cdot)}(\Omega)$ space, which appears in (4.20). For reasonably smooth problems, we thus expect Algorithm 8 to show better performance than Algorithm 7.

Remark 4.3.1. *Algorithm 8 can be equivalently formulated in terms of the modular function $\rho_{p(\cdot)}$ maintaining the same convergence rate and convergence analysis. However, $\bar{\rho}_{p(\cdot)}$ intuitively better generalises the norm $\frac{1}{p}\|\cdot\|_{\mathcal{X}}$ used in the definition of Guan and Song’s algorithm [105]. On the other hand, we underline that Algorithm 7 is defined in terms of $\bar{\rho}_{p(\cdot)}$, and it cannot be defined with $\rho_{p(\cdot)}$, since the convergence analysis cannot be carried out otherwise. In particular, in the proof of Lemma 4.2.1, (2.26) is necessary and it requires the use of $\bar{\rho}_{p(\cdot)}$.*

Remark 4.3.2. *It is important to notice that Proposition 4.3.3 allows to obtain a convergence rate in function values for Algorithm 3 as well. Indeed, simply by considering $g \equiv 0$ in Algorithm 8 we retrieve Algorithm 3, which, thus, enjoys the same convergence rate, in the framework of smooth optimisation algorithms.*

4.4 Sparse reconstruction and thresholding functions

In this section, we consider a popular sparse reconstruction model used in a variety of signal/image inverse problems and discuss the application of the algorithms presented in the previous sections of this chapter for the computation of its numerical solution. Given a Lebesgue measurable map $p(\cdot) : \Omega \rightarrow (1, 2]$, we consider a bounded linear operator $A : L^{p(\cdot)}(\Omega) \rightarrow L^{p_+}(\Omega)$ and an observation $y \in L^{p_+}(\Omega)$, $p_+ \leq 2$. For $\lambda > 0$, we consider the Tikhonov-like functional

$$\operatorname{argmin}_{x \in L^{p(\cdot)}(\Omega)} \frac{1}{p_+} \|Ax - y\|_{p_+}^{p_+} + \lambda \|x\|_1$$

in the Banach space $L^{p(\cdot)}(\Omega)$, where $f(x) = \frac{1}{p_+} \|Ax - y\|_{p_+}^{p_+}$ is proper, convex and smooth, while $g(x) = \lambda \|x\|_1$ is proper, l.s.c, convex and non-smooth.

The gradient of f can be computed as $\nabla f(x) = A^* \mathbf{J}_{p_+}^{p_+}(Ax - y) \in (L^{p(\cdot)}(\Omega))^*$. In agreement with the assumption on the Hölder continuity of the gradient of f made for Algorithm 7 and for Algorithm 8, we start with the study of this property of ∇f . We claim that ∇f is $(p_+ - 1)$ -Hölder continuous. To show this, we recall the following useful definitions and properties.

Definition 4.4.1. [204] *A Banach space \mathcal{X} is called smooth of power-type r (or r -smooth) with $r \in (1, 2]$ if there exists a constant $C > 0$ such that for all $u, v \in \mathcal{X}$*

it holds that

$$\frac{\|v\|_{\mathcal{X}}^r}{r} - \frac{\|u\|_{\mathcal{X}}^r}{r} - \langle \mathbf{J}_{\mathcal{X}}^r(u), v - u \rangle \leq C\|v - u\|_{\mathcal{X}}^r,$$

where $\mathbf{J}_{\mathcal{X}}^r(\cdot)$ denotes the duality mapping between \mathcal{X} and \mathcal{X}^* .

Proposition 4.4.1. [37] *If a Banach space \mathcal{X} is smooth of power-type r , then $\frac{\|\cdot\|_{\mathcal{X}}^r}{r}$ is continuously differentiable with derivative $\mathbf{J}_{\mathcal{X}}^r$, which is $(r-1)$ -Hölder continuous.*

Furthermore, for $s \geq r$, the functional $\frac{\|\cdot\|_{\mathcal{X}}^s}{s}$ is continuously differentiable. Its derivative is given by $\mathbf{J}_{\mathcal{X}}^s$, which is still $(r-1)$ -Hölder continuous on each bounded subset of \mathcal{X} .

Lemma 4.4.1. [204] *Lebesgue spaces $L^r(\Omega)$ are $\min\{2, r\}$ -smooth.*

In the following, we will denote $\mathbf{J}_{L^p}^r$ defined by (2.20) simply as \mathbf{J}_p^r . The space $L^{p_+}(\Omega)$ is by Lemma 4.4.1 p_+ -smooth and $\mathbf{J}_{p_+}^{p_+}$ is $(p_+ - 1)$ -Hölder continuous, i.e.

$$\exists K_1 > 0 \quad \text{s.t.} \quad \forall y_1, y_2 \in L^{p_+}(\Omega) \quad \|\mathbf{J}_{p_+}^{p_+}(y_1) - \mathbf{J}_{p_+}^{p_+}(y_2)\|_{(p_+)^*} \leq K_1\|y_1 - y_2\|_{p_+}^{p_+-1}.$$

Using this combined with the linearity of A and the sub-multiplicativity of the norm, we can thus write

$$\begin{aligned} \|\nabla f(u) - \nabla f(v)\|_{(L^{p(\cdot)})^*} &= \left\| A^* \left[\mathbf{J}_{p_+}^{p_+}(Au - y) - \mathbf{J}_{p_+}^{p_+}(Av - y) \right] \right\|_{(L^{p(\cdot)})^*} \\ &\leq \|A^*\|_{(L^{p(\cdot)})^*} \|\mathbf{J}_{p_+}^{p_+}(Au - y) - \mathbf{J}_{p_+}^{p_+}(Av - y)\|_{(L^{p(\cdot)})^*} \leq \|A^*\|_{(L^{p(\cdot)})^*} K_1 \|A(u - v)\|_{p_+}^{p_+-1} \\ &\leq K_1 \|A^*\|_{(L^{p(\cdot)})^*} \|A\|_{p_+}^{p_+-1} \|u - v\|_{L^{p(\cdot)}}^{p_+-1} \leq K \|u - v\|_{L^{p(\cdot)}}^{p_+-1} \quad \forall u, v \in L^{p(\cdot)}(\Omega), \end{aligned}$$

showing that ∇f is $(p_+ - 1)$ -Hölder continuous, as needed for Algorithms 7 and 8.

We can thus focus now on the computation of the solutions of (4.9) in the discrete setting, where the domain Ω is discretised into the disjoint sum of n nonempty measurable subsets, i.e. $\Omega = \bigcup_{i=1}^n \Omega_i$ and $\overset{\circ}{\Omega}_i \cap \overset{\circ}{\Omega}_j = \emptyset$ for $i \neq j$. By considering a single real value on each subset Ω_i , with a slight abuse of notation we simply denote by $\ell^{p(\cdot)}(\mathbb{R}^n)$, the n -th dimensional subspace of the sequence space $\ell^{p(\cdot)}$ generated by the first n elements e_1, e_2, \dots, e_n of the canonical basis. To allow effective numerical resolution, we heavily exploit the separability property of the operators involved in the sense of Definition 2.3.1. By setting $\sigma^k := A^* \mathbf{J}_{p_+}^{p_+}(Ax^k - y) \in \mathbb{R}^n$, the iteration (4.9) of Algorithm 7 reads as

$$\begin{aligned} x^{k+1} &= \underset{u \in \ell^{p(\cdot)}(\mathbb{R}^n)}{\operatorname{argmin}} \left\{ \bar{\rho}_{p(\cdot)}(u - x^k) + \tau_k \langle \sigma^k, u \rangle + \tau_k \lambda \|u\|_1 \right\} \\ &= \underset{u \in \ell^{p(\cdot)}(\mathbb{R}^n)}{\operatorname{argmin}} \sum_{i=1}^n \left\{ \frac{1}{p_i} |u_i - (x^k)_i|^{p_i} + \tau_k \sigma_i^k u_i + \tau_k \lambda |u_i| \right\}, \end{aligned} \quad (4.29)$$

where the separability property of each term leads to the sum with respect to $i = 1, \dots, n$. In other words, thanks to the additive separability property, at each k -th iteration, with $k \geq 1$, the n -dimensional minimisation problem in (4.29) corresponds

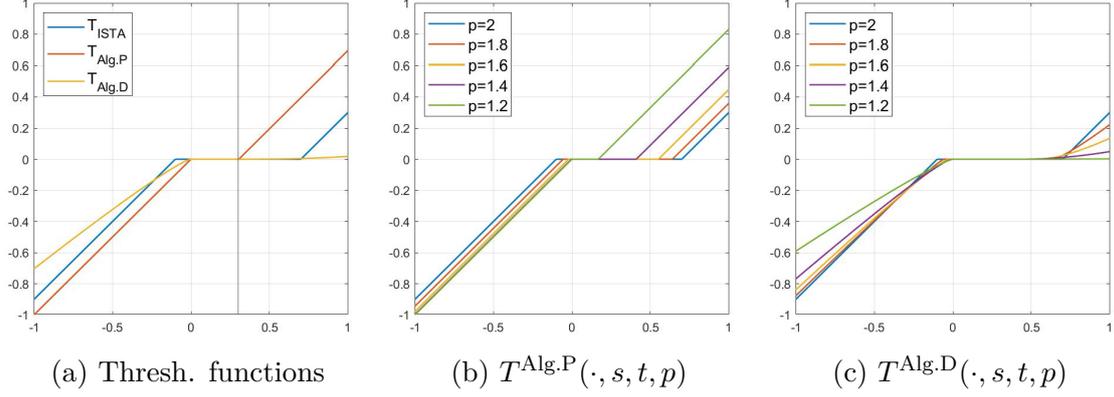


Figure 4.1: 1D thresholding functions for proximal-gradient algorithms in Banach spaces. (a) $T^{\text{Algo.P}}(\cdot, s, t, p)$, $T^{\text{Algo.D}}(\cdot, s, t, p)$ and $T^{\text{ISTA}}(\cdot, s, t)$ with $p = 1.3$, $s = 0.3$, $t = 0.4$. (b) $T^{\text{Algo.P}}(\cdot, s, t, p)$ with $s = 0.3$, $t = 0.4$ and $p \in \{1.2, 1.4, 1.6, 1.8, 2\}$. (c) $T^{\text{Algo.D}}(\cdot, s, t, p)$ with $s = 0.3$, $t = 0.4$ and $p \in \{1.2, 1.4, 1.6, 1.8, 2\}$.

to the sum of n 1D problems. Hence, each component can be treated independently, so one can consider the independent minimisation of the 1D functions:

$$\Psi_{x,s,t,p}^{\text{Algo.P}}(u) := \frac{1}{p}|u - x|^p + su + t|u|,$$

with $x = (x^k)_i$, $s = \tau_k \sigma_i^k$, $t = \tau_k \lambda$ and $p = p_i$, where the superscript Algo.P stands for primal algorithm. The minimisers of $\Psi_{x,s,t,p}^{\text{Algo.P}}$ can be computed by optimality as $(\partial \Psi_{x,s,t,p}^{\text{Algo.P}})^{-1}(0)$ and can be expressed in a compact form in terms of the thresholding function (see [37] for an analogous study in conventional $L^p(\Omega)$ spaces):

$$T^{\text{Algo.P}}(x, s, t, p) = \begin{cases} x - \text{sign}(s + t)|s + t|^{\frac{1}{p-1}} & \text{if } x > \text{sign}(s + t)|s + t|^{\frac{1}{p-1}} \\ x - \text{sign}(s - t)|s - t|^{\frac{1}{p-1}} & \text{if } x < \text{sign}(s - t)|s - t|^{\frac{1}{p-1}} \\ 0 & \text{otherwise.} \end{cases} \quad (4.30)$$

We can similarly focus on Algorithm 8, for which the Hölder continuity of the gradient of f (4.7) holds with $p = p_+$. The assumption defined by (4.22) is hard to verify in practice, although numerical tests show that it is not that challenging to find a step-size for which convergence is guaranteed. Proceeding similarly as above, we exploit again the separability of the modular appearing in the computation of the k -th iteration (4.23) of Algorithm 8, which leads to the computation of the minimisers of the following 1D function

$$\Psi_{x,s,t,p}^{\text{Algo.D}}(u) = \frac{1}{p}|u|^p - |x|^{p-1} \text{sign}(x)u + su + t|u|$$

where, as before, $x = (x^k)_i$, $s = \tau_k \sigma_i^k$, $t = \tau_k \lambda$, $p = p_i$ and the superscript Algo.D stands for dual algorithm. Similarly, such minimisers are now given by

$(\partial\Psi_{x,s,t,p}^{\text{Alg.D}})^{-1}(0)$ and correspond to the following thresholding function:

$$T^{\text{Alg.D}}(x, s, t, p) = \begin{cases} (|x|^{p-1} \text{sign}(x) - s - t)^{\frac{1}{p-1}} & \text{if } |x|^{p-1} \text{sign}(x) > s + t \\ -(s - t - |x|^{p-1} \text{sign}(x))^{\frac{1}{p-1}} & \text{if } |x|^{p-1} \text{sign}(x) < s - t \\ 0 & \text{otherwise.} \end{cases} \quad (4.31)$$

Figure 4.1a shows a comparison between the thresholding functions (4.30) and (4.31) with the classical soft-thresholding function

$$T^{\text{ISTA}}(x, s, t) = \begin{cases} x - s - t & \text{if } x > s + t \\ x - s + t & \text{if } x < s - t \\ 0 & \text{otherwise,} \end{cases} \quad (4.32)$$

which corresponds to the minimisation of the 1D function

$$\Psi_{x,s,t}^{\text{ISTA}}(u) = \frac{1}{2}(u - x + s)^2 + t|u|$$

appearing in the solution of the $\ell_2 - \ell_1$ LASSO optimisation problem in the Hilbert space ℓ_2 by means of the standard ISTA algorithm [70]. It is interesting to point out that differently from the ISTA thresholding function (4.32), both thresholding functions (4.30) and (4.31) are no longer symmetrical with respect to the vertical line $x = s$. Moreover, $T^{\text{Alg.P}}(\cdot, s, t, p)$ remains linear in x similarly to $T^{\text{ISTA}}(\cdot, s, t)$, as shown in Figure 4.1b for some exemplar values of p , while the thresholding function $T^{\text{Alg.D}}$ is, in general, nonlinear, as shown in Figure 4.1c. Note that for $p = 2$, both $T^{\text{Alg.P}}(\cdot, s, t, 2)$ and $T^{\text{Alg.D}}(\cdot, s, t, 2)$ coincide with the standard soft-thresholding $T^{\text{ISTA}}(\cdot, s, t)$ operator.

The thresholding functions $T^{\text{Alg.P}}$ and $T^{\text{Alg.D}}$ are respectively used in the implementation of Algorithms 5 and 6 for the resolution of (4.4) in constant exponent Lebesgue spaces $L^p(\Omega)$, since therein the p -power of the norm and the modular coincides.

4.5 Numerical tests

In this section, we provide some numerical tests showing how the proposed model adapts to deal with a variety of signal and image deconvolution and denoising problems. We include additional tests providing a numerical verification of the computational convergence properties of Algorithms 7 and 8.

4.5.1 Spike reconstruction

As a first example, we consider a 1D signal reconstruction problem where we seek for spikes defined on $\Omega = [0, 1]$ to be reconstructed from blurred measurements

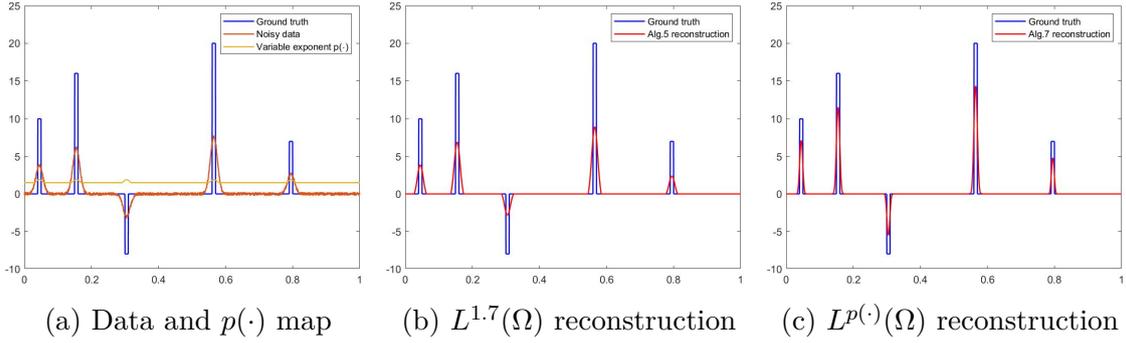


Figure 4.2: Spike reconstruction in Lebesgue spaces: comparison between constant $L^{1.7}(\Omega)$ and variable $L^{p(\cdot)}(\Omega)$ exponent Lebesgue spaces. Parameters: $\tau_k \equiv 0.5$; $\lambda = 10^{-2}$. Stopping criterion based on the normalised relative change between x^k and x^{k+1} : $\|x^k - x^{k+1}\|_2 / \|x^k\|_2 < 10^{-4}$.

corrupted with Gaussian noise, see Figure 4.2a. In order to favour sparse reconstructions and reduce possible over-smoothing, the formulation of a reconstruction model in a Banach space \mathcal{X} is considered, see, e.g., [37, 204]. Denoting by $A : \mathcal{X} \rightarrow L^2(\Omega)$ the blurring operator and by $y \in L^2(\Omega)$ the measured data, we thus aim to solve

$$\operatorname{argmin}_{x \in \mathcal{X}} \frac{1}{2} \|Ax - y\|_2^2 + \lambda \|x\|_1, \quad \lambda > 0 \quad (4.33)$$

where, in particular, we set $\mathcal{X} = L^{p(\cdot)}(\Omega)$ with $p_- = 1.6$ and $p_+ = 2$, as shown in orange in Figure 4.2a. The idea is to choose a higher value of $p(\cdot)$ where it is more likely to have a (spike) signal, while lower values are preferred elsewhere, so that, for these points, sparsity is enforced to a stronger extent. Note that the choice of the exponent map $p(\cdot)$ acts in fact as a prior model on the signal, together with the penalty term. To incorporate such prior knowledge, one can look directly at the shape of the data y , or, for instance, to the structure of a standard $\ell_2 - \ell_1$ reconstruction computed after a small number of iterations, in order to have a variable exponent $p(\cdot)$ consistent with an approximated (possibly over-smoothed) solution of the problem. The choice of $p(\cdot)$ has been already discussed in Section 3.3.5 of Chapter 3. Having chosen the exponent map $p(\cdot)$, we can thus solve (4.33) on $\mathcal{X} = L^{p(\cdot)}(\Omega)$ by means, e.g., of Algorithm 7.

In order to provide a comparison with existing models, we further consider problem (4.33) on $\mathcal{X} = L^p(\Omega)$ for $p = 1.7$ and solve it by means of Algorithm 5. We observe that using $L^{p(\cdot)}(\Omega)$ modelling improves the quality of the reconstruction with respect to a fixed $L^p(\Omega)$ modelling. In particular, we observe that the spikes look much more enhanced in the variable exponent reconstruction with respect to the constant one, which denoises well the acquisition but it does not manage to significantly deblur it.

In this test we have considered Algorithms 5 and 7, which both entail an implicit primal gradient-descent step. In the following, however, Algorithms 6 and 8 are considered, which are defined in terms of a Bregman proximal operator and

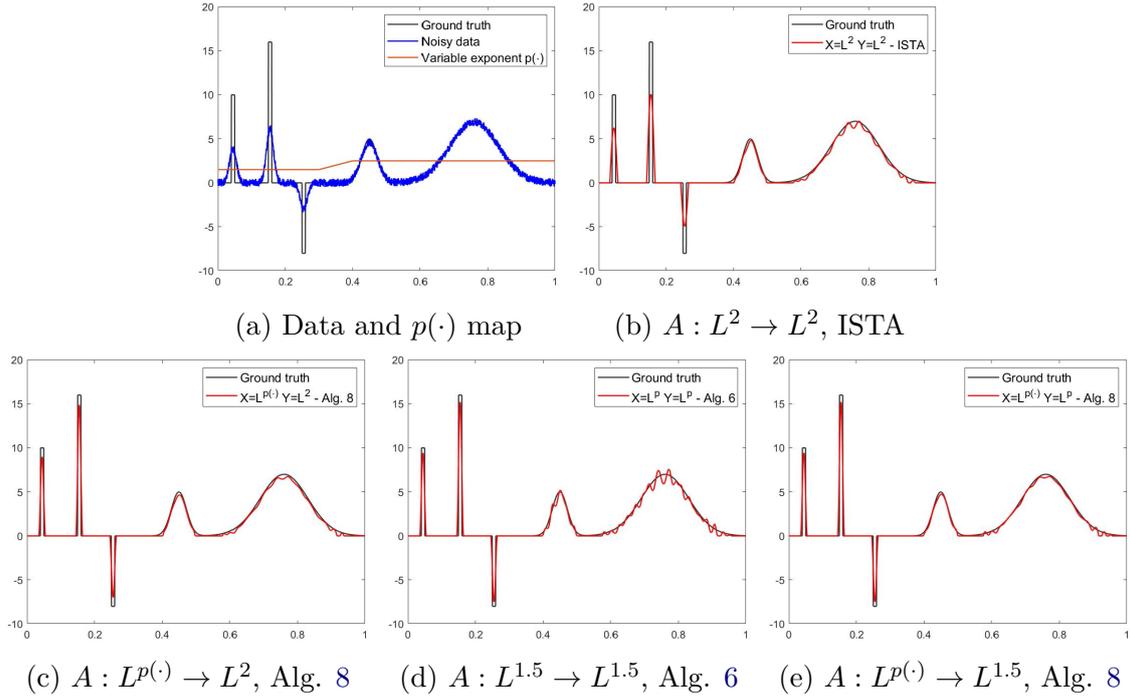


Figure 4.3: Deconvolution of heterogeneous signals in Hilbert $L^2(\Omega)$, classical Banach $L^{1.5}(\Omega)$ and variable exponent Lebesgue space $L^{p(\cdot)}(\Omega)$ settings. Parameters: $\tau_k \equiv 0.5$; $\lambda = 5 * 10^{-3}$. Stopping criterion based on the relative distance between x^k and x^{k+1} : $\|x^k - x^{k+1}\|_2 / \|x^k\|_2 < 4 * 10^{-6}$.

entail dual gradient-descent. We will see that the latter choices are faster in terms of convergence speed than the former primal algorithms and, thus, preferable in practice.

4.5.2 Deconvolution of heterogeneous signals

We now consider the signal deconvolution problem of a blurred and noisy data $y \in L^2(\Omega)$ corrupted by Gaussian noise of a 1D heterogeneous signal $x \in \mathcal{X}$ with $\Omega = [0, 1]$, that is sparse in some intervals and smooth in others, as shown in Figure 4.3a. As in the previous example, the blurring operator $A : \mathcal{X} \rightarrow L^2(\Omega)$ acts between \mathcal{X} , which will vary depending on the considered scenario, and the Hilbert space $L^2(\Omega)$. We consider the L^2 norm of the residual as data fitting term, so that model (4.33) is again used as reconstruction criterion. In Figures 4.3b and 4.3c we report the reconstructions obtained by solving (4.33) with solution spaces $\mathcal{X} = L^2(\Omega)$ using ISTA algorithm [70], and $\mathcal{X} = L^{p(\cdot)}(\Omega)$ using Algorithm 8, respectively, for the particular choice of exponent map $p(\cdot)$ having $p_- = 1.5$ and $p_+ = 2$ shown in orange in Figure 4.3a. We observe that the spikes in the left-hand side are better reconstructed when considering $L^{p(\cdot)}(\Omega)$ as solution space, see Figure 4.3c, due to its locally enhanced sparsifying property, while the Hilbert reconstruction of Figure 4.3b present loss of intensities in the spikes.

	$\frac{\ x_{rec}-x_{gt}\ _2}{\ x_{gt}\ _2}$	$\frac{\ Ax_{rec}-y^\delta\ _2}{\ y^\delta\ _2}$	SNR	SSIM
$A : L^2 \rightarrow L^2$, ISTA	0.3436	0.1169	9.28	0.64
$A : L^{p(\cdot)} \rightarrow L^2$, Alg. 8	0.2298	0.1152	12.77	0.70
$A : L^{1.5} \rightarrow L^{1.5}$, Alg. 6	0.2287	0.1023	12.81	0.54
$A : L^{p(\cdot)} \rightarrow L^{1.5}$, Alg. 8	0.2194	0.1031	13.18	0.70

Table 4.1: Deconvolution of heterogeneous signals in Hilbert $L^2(\Omega)$, classical Banach $L^{1.5}(\Omega)$ and variable exponent Lebesgue space $L^{p(\cdot)}(\Omega)$ settings. Quantitative results.

	$\frac{\ x_{rec}-x_{gt}\ _2}{\ x_{gt}\ _2}$	$\frac{\ Ax_{rec}-y^\delta\ _2}{\ y^\delta\ _2}$	SNR	SSIM
$A : L^2 \rightarrow L^2$, ISTA	0.5806	0.7989	4.72	$1.08 * 10^{-2}$
$A : L^2 \rightarrow L^{p(\cdot)}$, ISTA	0.4050	0.8226	7.85	$1.50 * 10^{-2}$
$A : L^{p(\cdot)} \rightarrow L^{p(\cdot)}$, Alg. 8	0.3813	0.8234	8.37	$1.67 * 10^{-2}$
$A : L^{1.4} \rightarrow L^{1.4}$, Alg. 6	0.7178	0.8294	2.88	$0.76 * 10^{-2}$
$A : L^{p(\cdot)} \rightarrow L^{1.4}$, Alg. 8	0.4089	0.8296	7.77	$1.58 * 10^{-2}$

Table 4.2: 1D mixed noise removal in Hilbert $L^2(\Omega)$, classical Banach $L^{1.5}(\Omega)$ and variable exponent Lebesgue space $L^{p(\cdot)}(\Omega)$ settings. Quantitative results.

Similarly as in [37], we could also consider an L^p , $1 < p < 2$, fidelity to better restore spikes. Consistently, we can consider $A : \mathcal{X} \rightarrow L^p(\Omega)$ with $p = 1.5$ and the following associated variational model

$$\operatorname{argmin}_{x \in \mathcal{X}} \frac{1}{p} \|Ax - y\|_p^p + \lambda \|x\|_1, \quad \lambda > 0.$$

In this case, Figures 4.3d and 4.3e show the reconstructions obtained with $\mathcal{X} = L^{1.5}(\Omega)$ using Algorithm 6 and $\mathcal{X} = L^{p(\cdot)}(\Omega)$ using Algorithm 8 respectively. We observe that smooth regions are better restored with reduced ringing effect artifacts in Figure 4.3e thanks to the flexible choice of the solution space. Similar observations can be made by looking at the quantitative results reported in Table 4.1, where for instance, we observe that the reconstruction of Figure 4.3e is the closest to the ground truth, among the computed reconstructions.

As a general comment, by these first numerical results we can say that working with a variable exponent allows us to deal with the different nature of the signal in a more flexible way.

4.5.3 1D and 2D mixed noise removal

We now focus on a mixed noise removal problem for blurred signals and images affected by Gaussian and impulsive (salt-and-pepper) noise, in different and disjoint parts of their spatial domain Ω , with $\Omega = [0, 1]$ and Ω being a compact of \mathbb{R}^2 , re-

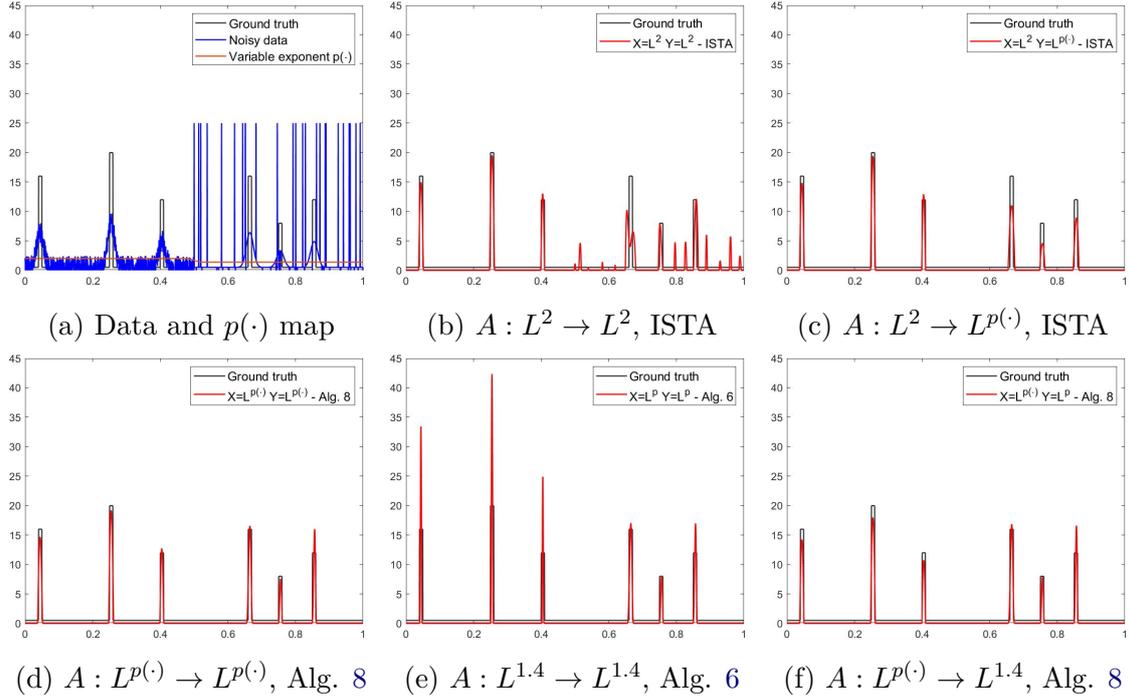


Figure 4.4: 1D mixed noise removal: different choices for the solution and output spaces. Parameters: $\tau_k \equiv 0.1$, $\lambda = 2 * 10^{-2}$. Stopping criterion based on the relative distance between x^k and x^{k+1} : $\|x^k - x^{k+1}\|_2 / \|x^k\|_2 < 4 * 10^{-6}$.

spectively. We exploit here the flexibility of Lebesgue spaces with variable exponent by effectively treating the different noise nature at the same time.

Let $\mathcal{X} = L^{p(\cdot)}(\Omega)$ and $\mathcal{Y} = L^{q(\cdot)}(\Omega)$ be two Lebesgue spaces with exponents $p(\cdot)$ and $q(\cdot)$, and let $A : L^{p(\cdot)}(\Omega) \rightarrow L^{q(\cdot)}(\Omega)$. To retrieve the sparse underlying signal, let us consider the following problem:

$$\operatorname{argmin}_{x \in L^{p(\cdot)}(\Omega)} \bar{\rho}_{q(\cdot)}(Ax - y) + \lambda \|x\|_1, \quad \lambda > 0. \quad (4.34)$$

For the 1D example, the measured blurred and noisy signal is shown in blue in Figure 4.4a: Gaussian noise is artificially added on the left part of the domain while impulsive noise is added on the right part. Choosing $p(\cdot) = q(\cdot) \equiv 2$ in (4.34) (and thus naturally setting $\mathcal{X} = \mathcal{Y} = L^2(\Omega)$) forces the fidelity term to reduce to $\|Ax - y\|_2^2$ which is well-known to be the most appropriate term to describe Gaussian noise degradation. Hence, in this case, an $L^2 - L^1$ variational model formulated in $L^2(\Omega)$ is obtained and the standard ISTA algorithm [70] is used for its numerical solution, shown in Figure 4.4b. We note that the reconstruction on the left part is good, while several artefacts can be observed on the right side, due to the poor adaptivity of the model to the different noise nature there. A data term more suited to describe the sparse nature of impulsive noise should be considered for the right part, such as, ideally, an L^1 fidelity; see [177]. In our modelling, however, exponents $p = 1$ cannot be chosen as they would correspond to a non-smooth data

	$\frac{\ x_{rec}-x_{gt}\ _2}{\ x_{gt}\ _2}$	$\frac{\ Ax_{rec}-y^\delta\ _2}{\ y^\delta\ _2}$	PSNR	SSIM	CPU time	# iterations
$L^2(\Omega)$	3.751	0.5115	21.82	0.78	249.84s	2292
$L^{1.4}(\Omega)$	1.143	0.5243	29.24	0.87	433.84s	2556
$L^{p(\cdot)}(\Omega)$	0.8554	0.6017	30.52	0.94	300.59s	2186

Table 4.3: 2D mixed denoising in Hilbert $L^2(\Omega)$, classical Banach $L^{1.5}(\Omega)$ and variable exponent Lebesgue space $L^{p(\cdot)}(\Omega)$ settings. Quantitative results.

term defined in a non-reflexive Banach space. However, exponents $p \approx 1$ can still be chosen. We thus consider as output space a Lebesgue space \mathcal{Y} with variable exponent map $q(\cdot) = p(\cdot)$ shown in orange in Figure 4.4a: it is equal to 2 in the part of the domain where the signal is corrupted by Gaussian noise, that is on the left side, and equal to 1.4 elsewhere. The corresponding space-variant modular-based fidelity is, therefore, differentiable and locally adapted to both natures of the noise thanks to the choice of $\mathcal{Y} = L^{p(\cdot)}$. In Figure 4.4c we show the ISTA reconstruction obtained by choosing the solution space as $\mathcal{X} = L^2(\Omega)$. In this case, although our choice of fidelity and measurement space \mathcal{Y} is better adapted to both noise distributions and favours the removal of impulsive noise from data, the spikes on the right-hand side suffer from some intensity loss. To improve upon this drawback, we propose a variable exponent in the definition of the solution space \mathcal{X} too, namely the same choice of variable exponent map $p(\cdot)$ of the output space is considered. In Figure 4.4d, we show the reconstruction obtained by solving (4.34) with $\mathcal{X} = \mathcal{Y} = L^{p(\cdot)}(\Omega)$ via Algorithm 8. The choice of an adaptive solution space improves the quality of the reconstruction and endows the model with enough flexibility to provide a good reconstruction of the signal over the whole domain. In this way, Gaussian noise and impulsive noise can be treated simultaneously.

A final comparison with $\mathcal{Y} = L^p(\Omega)$ with $p = 1.4$ again in the problem (4.34) is reported in Figures 4.4e and 4.4f. They show the reconstructions obtained by solving (4.34) with $\mathcal{Y} = L^{1.4}(\Omega)$ and solution spaces $\mathcal{X} = L^{p(\cdot)}(\Omega)$ and $\mathcal{X} = L^{1.4}(\Omega)$, respectively. Note that the reconstruction in Figure 4.4d is more accurate than the one in Figure 4.4f, in particular on the left-hand side since the spikes there are better reconstructed. In Table 4.2, we show a quantitative comparison of all the reconstructions obtained with the considered modellings. It is quite evident that, considering variable exponent Lebesgue spaces as solution spaces, the quality of the reconstructed signal improves significantly. In particular, considering the same Gaussian L^2 fidelity (that corresponds to $\mathcal{Y} = L^2(\Omega)$ with solution space $\mathcal{X} = L^{p(\cdot)}(\Omega)$ instead of $\mathcal{X} = L^2(\Omega)$) results in a significant improvement in terms of SNR and SSIM. The same consideration applies for the reconstructions obtained with the same acquisition space $\mathcal{Y} = L^{1.4}(\Omega)$ and different solution spaces $\mathcal{X} = L^{1.4}(\Omega)$ and $\mathcal{X} = L^{p(\cdot)}(\Omega)$. The reconstruction obtained using a variable exponent in both solution and acquisition space results the best among the ones we computed.

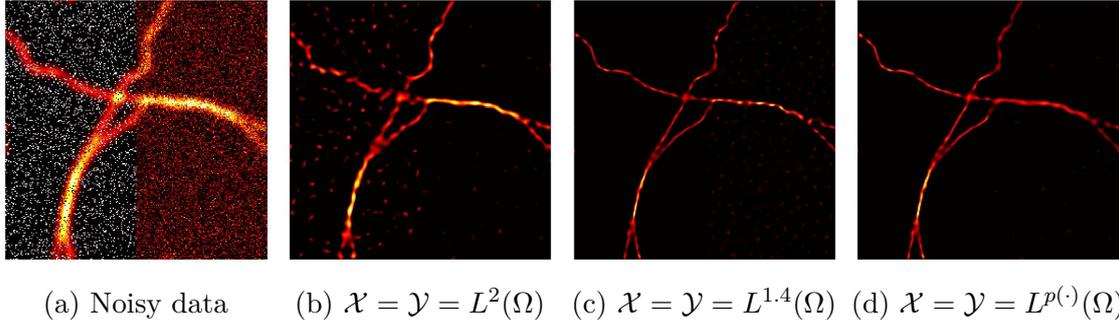


Figure 4.5: 2D mixed denoising. Parameters: $\tau_k \equiv 0.1$, $\lambda = 0.1$. Stopping criterion based on the normalized relative change between x^k and x^{k+1} : $\|x^k - x^{k+1}\|_2 / \|x^k\|_2 < 10^{-4}$.

As a similar test, we considered a mixed 2D denoising problem for the blurred and noisy image in Figure 4.5a, again with Gaussian noise on the left half and impulsive noise on the right half of the image domain. We solved again (4.34) for $p(\cdot) = q(\cdot) \equiv 2$ by the ISTA algorithm (Fig. 4.5b), for $p(\cdot) = q(\cdot) \equiv p = 1.4$ by Algorithm 6 (Fig. 4.5c), and for $p(\cdot) = q(\cdot) = P(\cdot)$, where $P(\cdot)$ is a variable exponent with $P_- = 1.4$ and $P_+ = 2$ such that it is equal to 2 on the Gaussian noisy half and on the impulsive noisy half it has value 1.4, by Algorithm 8 (Fig. 4.5d). Observations analogous to those discussed for the 1D case can still be made. Reconstruction artefacts are significantly reduced in the case of variable exponent modelling. From Table 4.3, we see that the values of the reconstruction error, PSNR and SSIM are better when considering a variable exponent setting. In particular, the reconstruction error drops significantly with $p = 1.4$ instead of $p = 2$ and again when taking a variable exponent $p(\cdot)$ instead of a constant one.

The flexibility of the model given by the choice of point-wise variable maps allows one to simultaneously reconstruct signals with different spatial properties on the whole domain. No domain decomposition techniques or domain splitting methods are required.

4.5.4 A numerical study on convergence rates

The numerical examples reported so far in this chapter show that the use of a variable exponent can help in improving reconstruction quality. It is thus natural to ask which algorithm – Algorithm 7 and Algorithm 8 – should be used in practice. As remarked already in [37], and as it can be observed from the convergence rate in function values (4.20), Algorithm 7 is expected to be very slow in practice, in particular, slower than a gradient-type algorithm whose well-known convergence speed is of the order $O(1/k)$.

In this section, we numerically compare the speed of convergence of different algorithms when used as solvers for the deblurring problem (4.33) in the different Hilbert and Banach scenarios discussed in Section 4.5.1. In particular, we compare

convergence speed of the ISTA algorithm used to solve (4.33) in $\mathcal{X} = L^2(\Omega)$, with Algorithms 5 and 6 for solving the same problem on $\mathcal{X} = L^p(\Omega)$ with $p = 1.7$ and with Algorithms 7 and 8 proposed in this work for $\mathcal{X} = L^{p(\cdot)}(\Omega)$. As exponent map $p(\cdot)$, we stick with the choice shown in orange Figure 4.2a.

Given $x^* \in \mathcal{X}$, solution of (4.33) in each different space \mathcal{X} considered, we recall the convergence rates in function values (4.20) and (4.28) for Algorithms 7 and 8, respectively. In principle, to compare the speed of convergence in a precise way, a pre-computation of (a suitable approximation of) x^* by means of benchmark algorithms should be done for all the different scenarios discussed above. However, since to our knowledge there is no existing algorithm for solving (4.33) in $\mathcal{X} = L^{p(\cdot)}(\Omega)$, instead of computing x^* we computed as a reference the value of \tilde{x} , solution of (4.33) with $\mathcal{X} = L^2(\Omega)$, by running ISTA for $2 * 10^4$ iterations. Note that by simple algebraic manipulations, from (4.20) we have for Algorithm 7

$$\phi(x^k) - \phi(\tilde{x}) = \phi(x^k) \pm \phi(x^*) - \phi(\tilde{x}) \leq \eta_1 \left(\frac{1}{k}\right)^{p-1} + c$$

and, similarly, from (4.28) for Algorithm 8

$$\phi(x^k) - \phi(\tilde{x}) \leq \eta_2 \left(\frac{1}{k}\right)^{p-1} + c$$

with $p = 2$ for the smooth part of (4.33), so that rates can still be compared up to an additive constant. We use \tilde{x} also for the computation of the convergence rates for Algorithms 5 and 6 having

$$\phi(x^k) - \phi(\tilde{x}) \leq \eta_{3,4} \left(\frac{1}{k}\right)^{p-1} + c',$$

again with $p = 2$ for (4.33). Finally, recall that for ISTA the convergence rate in function values is

$$\phi(x^k) - \phi(\tilde{x}) \leq \eta_5 \frac{1}{k},$$

where no additive constant to correct the rate is needed, being \tilde{x} computed with ISTA for $\mathcal{X} = L^2(\Omega)$. As stopping criterion, for all the tested algorithms we use the normalised relative rates with respect to \tilde{x} up to a tolerance parameter $\epsilon = 10^{-4}$:

$$|\phi(x^k) - \phi(\tilde{x})|/\phi(\tilde{x}) < \epsilon.$$

The results reported in Figure 4.6 and Table 4.4 show several interesting numerical convergence properties. First, we note that although Algorithm 5 (violet line) and Algorithm 6 (green line) are supposed to have the same convergence rate in theory for the specific problem at hand, they clearly have a very different behaviour. While the first needs more than $5 * 10^5$ iterations and more than 1000 seconds of CPU time to reach convergence, the second converges with a much faster speed. The same behaviour is observed also for the modular Algorithms 7 and 8 too, with the first (red line) being very slow and the second (yellow line) much faster. With

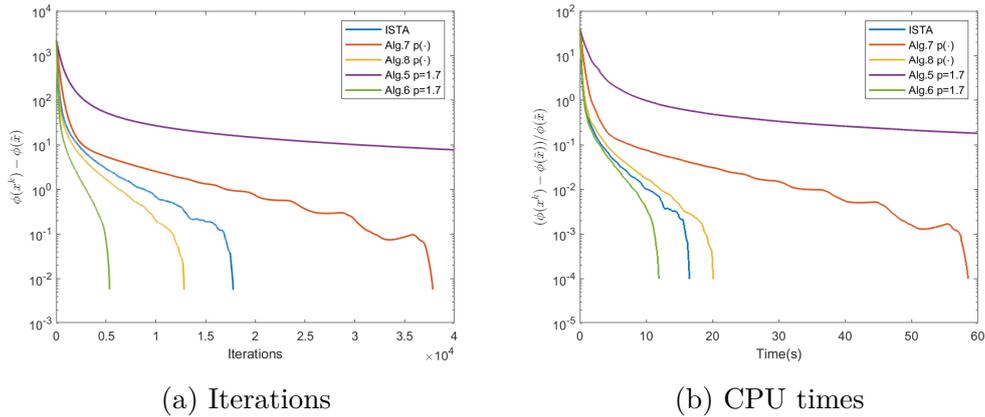


Figure 4.6: Numerical study on convergence rates in function values $\phi(x^k) - \phi(\tilde{x})$ for proximal gradient algorithms in Banach spaces. (a) Relative rates along the first $4 * 10^4$ iterations. (b) Normalised relative rates along the first 60 seconds of CPU time.

	ISTA	Alg. 5	Alg. 6	Alg. 7	Alg. 8
# iterations	17768	$5 * 10^5$	5352	37850	12849
CPU time (s)	16.5	1025.7	11.9	58.5	20.0

Table 4.4: Algorithmic comparison: iterations required and CPU time till convergence.

respect to standard ISTA (blue line), we observe that the modular-based (Alg. 8) (yellow line) and the norm-based proximal dual-gradient algorithms (Alg. 6) (green line) are faster in terms of number of iterations and comparable in terms of computational time, whilst the modular-based (Alg. 7) (red line) and norm-based (Alg. 5) (violet line) proximal primal-gradient algorithms are very slow. This fact suggests that the computation of the gradient step in the primal space is rather inefficient in terms of computational times, whilst performing it in the dual space speed the convergence of the algorithms. As a general comment, we can devise that this numerical study on convergence rates motivates the choice of Algorithm 8 instead of Algorithm 7 in Sections 4.5.2 and 4.5.3, as well as the choice of Algorithm 6 over Algorithm 5.

It is also important to address that the analytical convergence rate obtained for Algorithm 7 is a worst-case convergence speed, as already noticed after the proof of Proposition 4.2.4. Indeed, as expected, numerical tests show that the actual convergence speed does not depend just on p_- , the infimum value of $p(\cdot)$, but rather on the whole distribution of the exponent function. We studied the convergence behaviour of Algorithm 7 with respect to the choice of the exponent for problem (4.33) for the spike deconvolution problem presented in Section 4.5.1 and shown

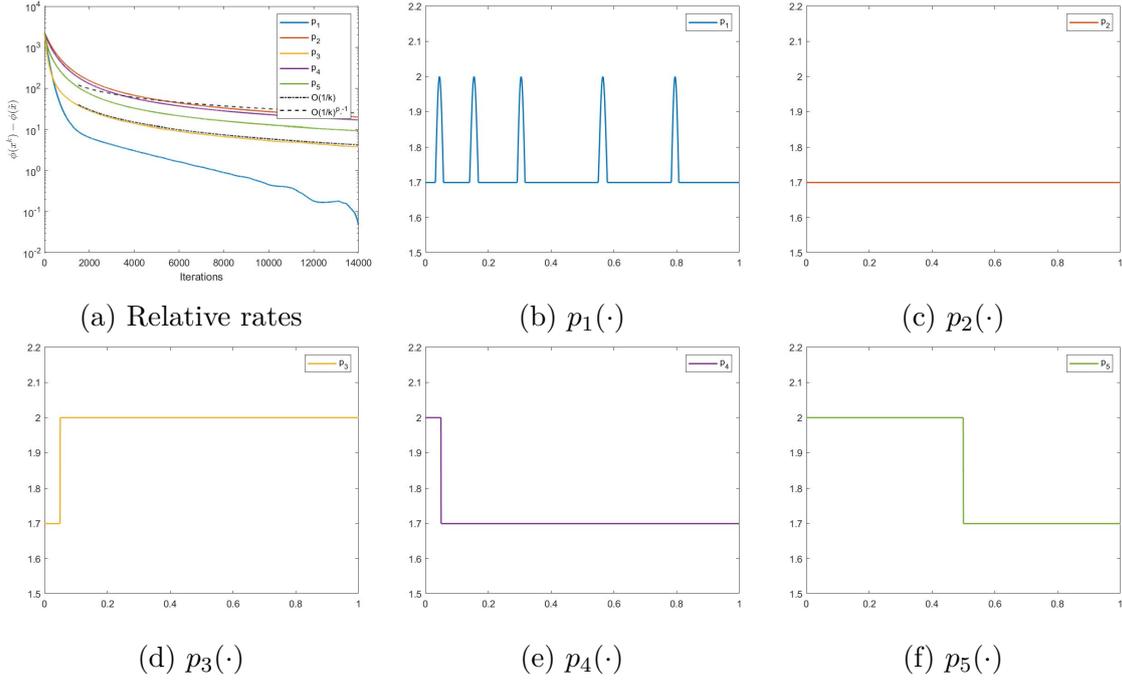


Figure 4.7: Numerical study of the convergence rate in function values of Algorithm 7. Relative rates along the first 14000 iterations and different choices of variable exponents functions $p(\cdot)$ with $p_- = 1.7$.

in Figure 4.2, with $\mathcal{X} = L^{p(\cdot)}(\Omega)$ and different choices of variable exponent $p(\cdot)$. We report in Figure 4.7 the convergence rate histories obtained with five different exponents, all having $p_- = 1.7$, together with ISTA convergence rate of $\mathcal{O}\left(\frac{1}{k}\right)$ and Algorithm 5 convergence rate of $\mathcal{O}\left(\frac{1}{k^{p_- - 1}}\right)$ with $p_- = 1.7$. We can observe that for the constant exponent $p_2(\cdot) \equiv 1.7$ empirically we obtain the estimated rate. A similar behaviour is observed for $p_4(\cdot)$, which is equal to $(p_4)_- = 1.7$ in the vast majority of the domain. However, different choices of the variable exponent leads to a better empirical convergence. For example with $p_3(\cdot)$ which is equal to 2 in almost all the interval $\Omega = [0, 1]$, we obtain an empirical rate comparable with the rate of ISTA, corresponding to the constant choice of $p = 2$. This shows that, even if $(p_3)_- = 1.7$, the convergence speed (which theoretically should be the same as the one obtained with $p_2(\cdot)$) is not affected by the small part of the domain where $p_3(\cdot) = 1.7$ and, instead, we retrieve ISTA convergence rate. For $p_5(\cdot)$, which takes the value of 1.7 in half of the domain, we observe a convergence speed which is in between the ISTA's and Algorithm 5's one. The choice of $p_1(\cdot)$ (which takes up the shape of the blurred data) significantly improves the convergence speed, even if it assumes the value p_- in the majority of the domain. This is a further motivation to such a choice for the variable exponent, that we detailed in Section 3.3.4.

These numerical validations suggest that the choice of the exponent map should take into account the whole structure, i.e. the regularity, of the data.

4.6 Final discussion

In this chapter, we presented two different ways to define forward-backward algorithms in general Banach spaces. A possible strategy is to follow [37] and define an algorithm that implicitly performs a gradient step in the primal space and a proximal step based on the p -norm proximal operator. Another definition is given in [105], where the gradient step is explicitly computed in the dual space and the proximal step is defined in terms of the Bregman proximal operator. Similarly as in the case of smooth optimisation, these algorithms are hard to use in practice in variable exponent Lebesgue spaces $L^{p(\cdot)}(\Omega)$, since they are all norm-based and the Luxemburg norm has many undesirable properties from a numerical point of view. Thus, similarly as for the definition of a gradient descent algorithms, we proposed to define forward-backward algorithms in $L^{p(\cdot)}(\Omega)$ by using the modular function and its derivative instead of the norm and the duality maps. We hence presented two instances of modular-based proximal gradient algorithms in $L^{p(\cdot)}(\Omega)$, adapting the conventional norm-based algorithms of general Banach spaces. We studied in both cases the convergence of the algorithms in function values and, under stronger hypothesis on the functional to minimise, to the minimiser of such functional.

Numerical tests to show the good performance of the proposed methods and to highlight the advantages that using a variable exponent versus a constant one yields have been computed and analysed. This is especially evident in mixed-noise scenarios and with heterogeneous signals.

To conclude, we showed with some simple 1-dimensional tests that the attained analytical convergence rates consist of worst-case estimates, as expected. Moreover, it resulted evident that Algorithm 8 is much faster than Algorithm 7, similarly to what happens to their constant exponents counterparts Algorithm 6 and Algorithm 5. It seems that computing the gradient descent step in the dual space improves the convergence speed and it is thus more desirable than doing it in the primal space, as well as for the simple Landweber primal and dual methods (Alg.2 and Alg.1) of [204]. Moreover, the choice of a good exponent appears to have an influence, not only in the reconstruction quality, but also in terms of convergence speed. This interesting behaviour will be object of further study.

Part II

Sparse optimisation in the
Banach space of Radon
measures with Poisson noise

Sparse off-the-grid optimisation methods in imaging

In this chapter, the mathematical theory of sparse off-the-grid optimisation is presented. Off-the-grid methods arise when formulating inverse problems in the continuous setting of the space of Radon measures $\mathcal{M}(\Omega)$, which is desirable in order to avoid discretisation biases and numerical instabilities of variational discretised approaches. In this context, the standard $L^2 - L^1$ LASSO regularisation takes the form of the so-called BLASSO problem, where an L^2 fidelity is coupled with a penalty consisting of the total-variation norm in $\mathcal{M}(\Omega)$. This penalty term promotes the reconstruction of discrete measures, i.e. finite linear combinations of weighted Diracs. We will review here optimality conditions for the BLASSO problem, obtained considering the corresponding dual problem, and the Sliding Frank Wolfe algorithm used for its resolution.

5.1	Inverse problems in the space of measures	108
5.1.1	Going <i>off-the-grid</i> for sparse spikes deconvolution	108
5.1.2	The space of Radon measures	109
5.1.3	Inverse problems in $\mathcal{M}(\Omega)$	111
5.2	The BLASSO problem	113
5.2.1	Optimality conditions	115
5.2.2	Dual problem and extremality conditions	116
5.2.2.1	Convex conjugate of the fidelity	116
5.2.2.2	Convex conjugate of the TV norm	117
5.2.2.3	Dual problem and extremality conditions	117
5.3	Frank-Wolfe and Sliding Frank-Wolfe algorithms	118
5.3.1	Frank-Wolfe algorithm	118
5.3.2	Frank-Wolfe for the minimisation of BLASSO	120
5.3.2.1	Greedy approach	121

5.3.3 Sliding Frank-Wolfe algorithm	123
5.4 Final discussion	124

In this second part of the thesis, we focus on imaging inverse problems defined in the space of Radon measures $\mathcal{M}(\Omega)$. This formulation is called *off-the-grid* or *gridless*, since images are not modelled anymore as vectors $x \in \mathbb{R}^N$ with $N \in \mathbb{N}$ being the number of pixels of the grid that discretises the domain Ω , but instead they are seen as measures in the continuous domain $\Omega \subseteq \mathbb{R}^d$, $d \geq 1$.

5.1 Inverse problems in the space of measures

Inverse problems in the space of measures and off-the-grid optimisation methods have been first proposed in [39, 72, 96] and since then they have been a topic of intense research for the mathematical community [34, 74, 75, 85, 184]. Off-the-grid methods are particularly useful for reconstructing fine-scale details from noisy acquisitions, with application to the localisation of spikes in a continuous domain Ω in astronomy or microscopy [75], to the parameter estimation for a super-positions of point sources in spectroscopy [84] and to density mixture estimation [52].

We start by briefly outlining discrete *on-the-grid* approaches and giving an intuitive description of their off-the-grid extensions. We sketch the motivation behind their first introduction taking the example of sparse spikes deconvolution.

5.1.1 Going *off-the-grid* for sparse spikes deconvolution

In off-the-grid approaches, the spatial domain is not discretised with a grid. The inverse problem is thus tackled directly in the continuous domain Ω and the desired solution is modelled as a Radon measure in $\mathcal{M}(\Omega)$ and not as a vectorised image $x \in \mathbb{R}^N$ with N pixels.

The application of interest in this thesis will be the sparse deconvolution problem of point-sources arising in fluorescence microscopy imaging, see Appendix B for more details on the topic. This problem has been widely studied in the conventional discrete setting, i.e. *on-the-grid*. The goal here is to estimate molecules' intensities and positions from blurred acquisitions of sparse samples of molecules, that is to reconstruct an image x from a blurred and noisy image y . To attain fine-scale details and a higher precision in the reconstructions two different grid are usually considered: the coarse grid of the acquisition $y \in \mathbb{R}^M$ and a finer grid for reconstruction $x \in \mathbb{R}^N$, as in Figure 5.1, with $N = L^2M$, where the parameter $L > 0$ controls how much finer the grid of the reconstruction is. Then, given $y \in \mathbb{R}^M$, the ill-posed inverse problem of reconstructing $x \in \mathbb{R}^N$ consists in solving

$$y = R_L H x + \omega,$$

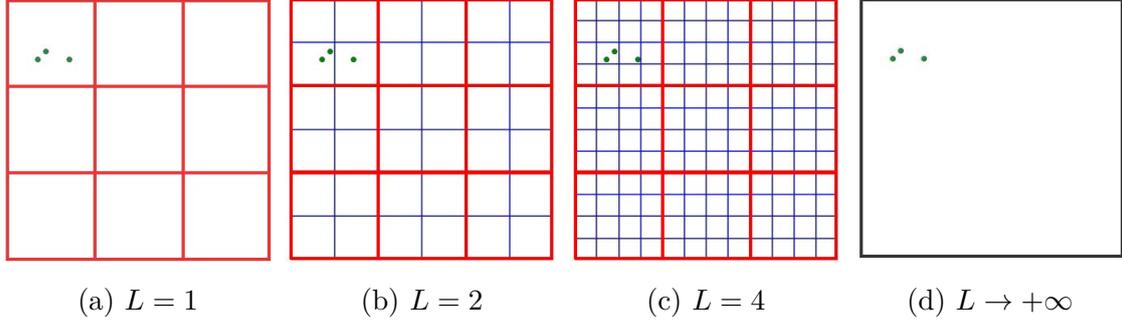


Figure 5.1: From on-the-grid to off-the-grid formulation. In red, the coarse grid of the acquisition $y \in \mathbb{R}^M$. In blue, the fine grid of the reconstruction $x \in \mathbb{R}^N$ with $N = L^2M$. In green, the point-sources to localise.

where $H \in \mathbb{R}^{N \times N}$ is the blurring operator corresponding to the microscope PSF, $R_L \in \mathbb{R}^{M \times N}$ is the under-sampling operator, mapping images in the finer grid to the coarse one, and $\omega \in \mathbb{R}^M$ is, in the basic scenario, an additive white Gaussian noise component, i.e. $\omega \in \mathcal{N}(0, \sigma^2 \text{Id})$. Nonetheless, other noise distributions are in general possible. Gaussian noise is often considered since it yields to simpler variational models than other hypotheses on the noise distribution. However, in the setting of fluorescence microscopy the Poisson distribution is more realistic, since in the measurement of light noise has a photon-counting nature [16]. The factor $L > 0$ controls the size of the finer grid: as the factor L increases, the grid is finer. In Figure 5.1 it is evident that with a coarse grid point-sources that are too close cannot be separated, thus a high value for L is desirable. In these cases, however, by increasing the size of the reconstruction grid the problem becomes more and more under-determined, causing instabilities in the reconstructions [87], and the dimensions of the inverse problem increase, making its resolution computationally expensive. One can thus think of *off-the-grid* approaches as the limit for L that goes to $+\infty$ of *on-the-grid* formulations [86, 141].

5.1.2 The space of Radon measures

We start now by introducing in a more formal way the space of Radon measures $\mathcal{M}(\Omega)$ and the off-the-grid inverse problem formulation. This section gathers some important definitions and properties of $\mathcal{M}(\Omega)$ from [75, 141, 183]. For more details on Radon measures see [61, 196].

Let $\Omega \subseteq \mathbb{R}^d$, with $d \in \mathbb{N}$, $d \geq 1$, be a compact subset of \mathbb{R}^d with non-empty interior. We denote by $\mathcal{C}_0(\Omega, \mathbb{R})$ the space of real continuous functions on Ω that vanish at infinity, namely all the continuous maps $\psi : \Omega \rightarrow \mathbb{R}$ such that:

$$\forall \varepsilon > 0, \exists K \subset \Omega \text{ compact}, \quad \sup_{x \in \Omega \setminus K} |\psi(x)| \leq \varepsilon.$$

It is now possible to give the following definition of the space of Radon measure through duality.

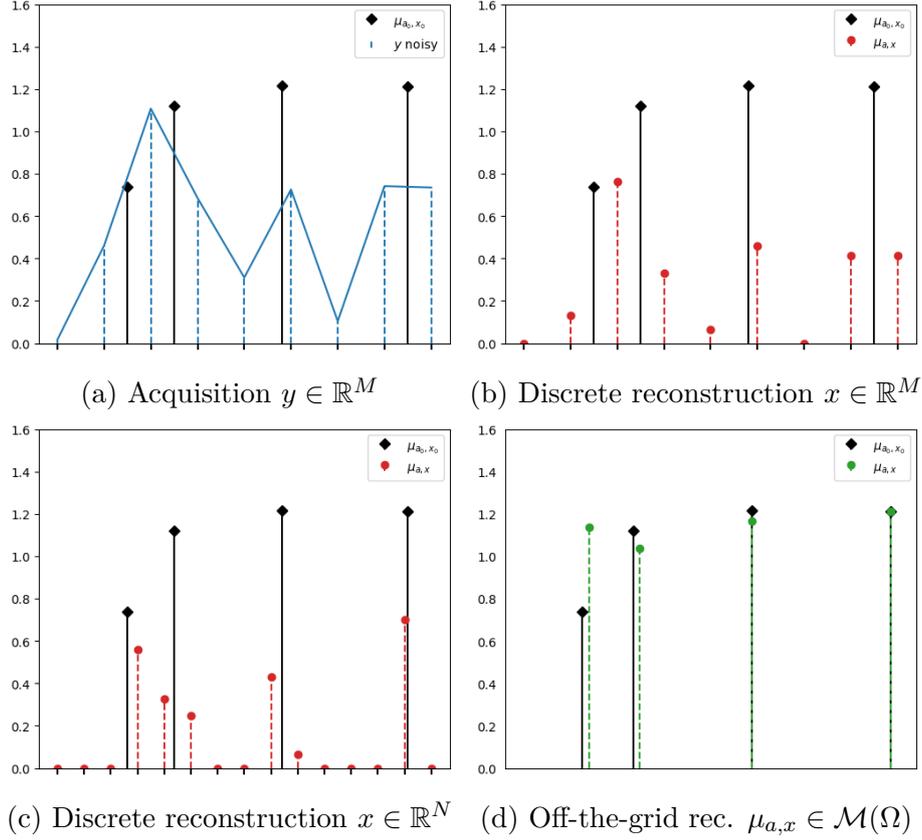


Figure 5.2: Comparison between conventional discrete (on-the-grid) and off-the-grid reconstructions. In black: the ground-truth spikes to retrieve. In Fig.5.2a: in blue, the acquired blurred and noisy signal. In Fig.5.2b: in red, discrete reconstruction with support constrained on a grid with M pixels. In Fig.5.2c: in red, discrete reconstruction with support constrained on a grid with $N > M$ pixels. In Fig.5.2d: in green, the off-the-grid reconstruction. The green spikes are the reconstruction without an a priori fixed grid, so they can move continuously on the line.

Definition 5.1.1. *The Banach space of real signed Radon measures on Ω $\mathcal{M}(\Omega)$ is the topological dual of $\mathcal{C}_0(\Omega, \mathbb{R})$ endowed with the supremum norm $\|\cdot\|_{\infty, \Omega}$, defined by $\|\psi\|_{\infty, \Omega} := \sup_{x \in \Omega} |\psi(x)|$.*

This interpretation allows to define a measure as a linear form on $\mathcal{C}_0(\Omega, \mathbb{R})$.

Definition 5.1.2. *A Radon measure $\mu \in \mathcal{M}(\Omega)$ is a continuous linear form evaluated on functions $\psi \in \mathcal{C}_0(\Omega, \mathbb{R})$, with duality pairing denoted by*

$$\langle \psi, \mu \rangle_{\mathcal{C}_0(\Omega, \mathbb{R}) \times \mathcal{M}(\Omega)} = \int_{\Omega} \psi d\mu. \quad (5.1)$$

A Radon measure $\mu \in \mathcal{M}(\Omega)$ is a positive measure if $\langle \psi, \mu \rangle_{\mathcal{C}_0(\Omega, \mathbb{R}) \times \mathcal{M}(\Omega)}$ is positive for any non-negative function ψ . This specifies the meaning of the term *signed* in the above definition: the quantity $\langle \psi, \mu \rangle_{\mathcal{C}_0(\Omega, \mathbb{R}) \times \mathcal{M}(\Omega)}$ can be also negative.

It is possible to define analogously the space of complex Radon measures $\mathcal{M}_{\mathbb{C}}(\Omega)$ as the dual of $\mathcal{C}_0(\Omega, \mathbb{C})$. On the other hand, the space of positive (non-negative) Radon measures $\mathcal{M}^+(\Omega)$ cannot be defined as the topological dual of $\mathcal{C}_0(\Omega, \mathbb{R}^+)$, since the latter is not a vector space.

Some classic examples of real Radon measures are:

- the Lebesgue measure of dimension $d \in \mathbb{N}$;
- the Dirac measure δ_z centred in $z \in \Omega$. For all $\psi \in \mathcal{C}_0(\Omega, \mathbb{R})$, one has $\langle \psi, \delta_z \rangle_{\mathcal{C}_0(\Omega, \mathbb{R}) \times \mathcal{M}(\Omega)} = \psi(z)$;
- discrete measures with $N \in \mathbb{N}, a \in \mathbb{R}^N, x \in \Omega^N$

$$\mu_{a,x} = \sum_{i=1}^N a_i \delta_{x_i}. \quad (5.2)$$

$\mathcal{M}(\Omega)$ is a non-reflexive Banach space endowed with its dual norm, called the total-variation (TV) norm, defined by:

$$|\mu|(\Omega) = \sup \left(\int_{\Omega} \psi d\mu \mid \psi \in \mathcal{C}_0(\Omega, \mathbb{R}), \|\psi\|_{\infty, \Omega} \leq 1 \right) \quad \forall \mu \in \mathcal{M}(\Omega).$$

The TV norm is non-differentiable but it is possible to consider its subdifferential [85]:

$$\partial|\mu|(\Omega) = \left\{ \psi \in \mathcal{C}_0(\Omega, \mathbb{R}) \mid \|\psi\|_{\infty, \Omega} \leq 1 \text{ and } \int_{\Omega} \psi d\mu = |\mu|(\Omega) \right\}. \quad (5.3)$$

In particular, for discrete measures (5.2), its TV norm coincides with the L^1 norm of the vector a :

$$|\mu_{a,x}|(\Omega) = \|a\|_1.$$

This explains why the TV norm is considered to be a generalisation of the L^1 norm. Moreover, in this special case the subdifferential has the following expression

$$\partial|\mu_{a,x}|(\Omega) = \{ \psi \in \mathcal{C}_0(\Omega, \mathbb{R}) \mid \|\psi\|_{\infty, \Omega} \leq 1, \forall i = 1, \dots, N \psi(x_i) = \text{sign}(a_i) \}. \quad (5.4)$$

5.1.3 Inverse problems in $\mathcal{M}(\Omega)$

We focus now on the formulation of inverse problems of the form (1.1) in $\mathcal{M}(\Omega)$. We will consider as acquisition space a Hilbert space \mathcal{H} .

Let $\mu \in \mathcal{M}(\Omega)$ be the continuous source measure, we call acquisition $y \in \mathcal{H}$ the result of the forward operator $\Phi : \mathcal{M}(\Omega) \rightarrow \mathcal{H}$ evaluated on μ , with a fixed measurement kernel $\varphi : \Omega \rightarrow \mathcal{H}$:

$$y := \Phi\mu = \int_{\Omega} \varphi(x) d\mu(x). \quad (5.5)$$

This last integral should not be confused with the concept of duality pairing (5.1), which is defined as the integral over Ω of a continuous real function with respect

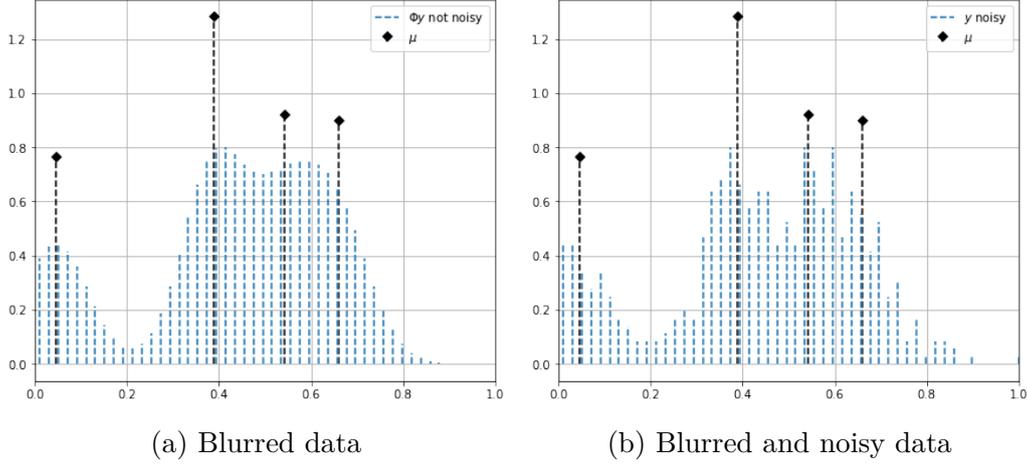


Figure 5.3: Discrete ground truth measure (black) and acquisition (blue) attained with the convolution with a Gaussian PSF (and Gaussian noise in Fig.5.3b).

to a measure $\mu \in \mathcal{M}(\Omega)$. In (5.5), the integral is then a Bochner integral [61] for vector-valued functions $\varphi(x) \in \mathcal{H}$, that is $\varphi(x)$ is not a real value but an element of \mathcal{H} , and it is well-defined if φ is continuous and bounded [58, 75].

The choice of the kernel φ and of the acquisition space \mathcal{H} depends on the physical process of acquisition considered. For the scope of our analysis, we consider here a convolution kernel, being the one of interest in fluorescence microscopy imaging [75]. In this setting, a natural choice for acquisition space is $\mathcal{H} = L^2(\Omega)$. The convolution kernel $\varphi : \Omega \rightarrow L^2(\Omega)$ with PSF $\tilde{\varphi} : \Omega \rightarrow L^2(\Omega)$ is defined as follows:

$$\varphi(x) := (s \mapsto \tilde{\varphi}(s - x)) \in L^2(\Omega). \quad (5.6)$$

Depending on the microscopy technique used one can have different PSFs. For instance, the Gaussian PSF, centred in $c \in \Omega$ with radius $\sigma > 0$, is defined by

$$s \mapsto \tilde{\varphi}(s - c) := 1/\sqrt[2]{2\pi\sigma^2} e^{-\|s-c\|_2^2/2\sigma^2}.$$

Other possible measurement kernels are

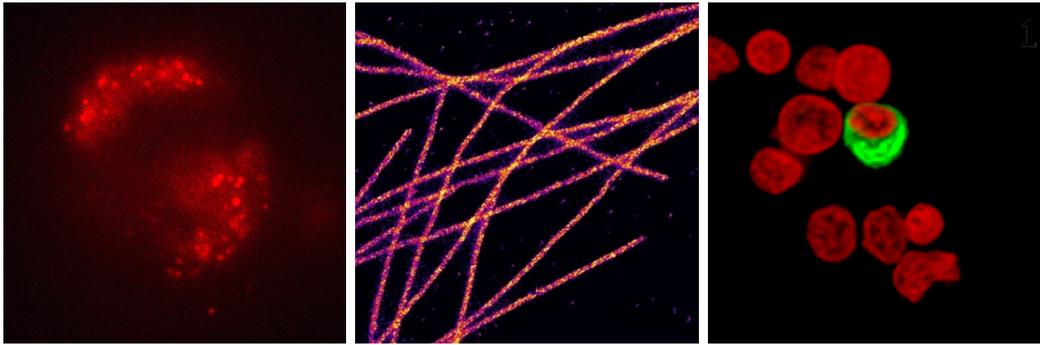
- the Fourier kernel, used i.e. in Nuclear Magnetic Resonance spectroscopy [84], where the kernel $\varphi : \Omega \rightarrow \mathcal{H}$ with acquisition space $\mathcal{H} = \mathbb{C}^{2f_c+1}$ and with cut-off frequency $f_c \in \mathbb{N}$ is

$$\varphi(x) = \left(e^{2i\pi kx} \right)_{|k| \leq f_c}. \quad (5.7)$$

- the Laplace kernel, used i.e. in MA-TIRF microscopy techniques [75], where the kernel is defined by

$$\varphi(x) = (s \mapsto \xi(x)e^{-sx}) \in \mathcal{H},$$

with respect to a non-negative weighting function $\xi \in \mathcal{C}(\Omega)$ specific to the physical acquisition process.



(a) Credit: Leila Muresan (b) Image from [114] (c) Image from [218]

Figure 5.4: Examples of biological fluorescent microscopy images. From left to right: molecules, cells, microtubules.

It is interesting to observe the action of the forward operator on finite linear combination of Dirac measures (5.2):

$$\Phi\mu_{a,x} = \int_{\Omega} \varphi(x) d\mu(x) = \sum_{i=1}^N a_i \varphi(x_i).$$

For simplicity, the following notation will be used $\Phi_x(a) = \sum_{i=1}^N a_i \varphi(x_i)$, instead of $\Phi\mu_{a,x}$.

In fluorescence microscopy imaging, the objects of interest are images of molecules, i.e. point-sources emitting fluorescent light. In this particular biological application, the unknown is well-described by discrete measures of the form (5.2), where $N \in \mathbb{N}$ is the number of molecules and any Dirac δ_{x_i} represents one molecule in the space Ω with position $x_i \in \Omega$ and amplitude $a_i \in \mathbb{R}$. Hence, a_i can only be positive (or null), since in a given position x_i either there is a source of light, and hence a_i is positive, or there is not, and hence $a_i = 0$. For this reason, the source measures $\mu_{a,x}$ are non-negative measures. Other possible objects of interest in microscopy images are microtubules, 1-dimensional curve structures, and cells (2-dimensional), that can be modelled as piece-wise constant functions, see Figure 5.4. In the following, we will focus only on sparse deblurring to recover point sources (0-dimensional), giving a brief description of the other two cases.

5.2 The BLASSO problem

The sparse spikes deconvolution problem, that we approach now, consists in recovering a (small) finite linear combination of Diracs $\mu_{a,x} = \sum_{i=1}^N a_i \delta_{x_i}$ from a blurred and noisy acquisition

$$y = \Phi\mu_{a,x} + \omega,$$

where $\Phi : \Omega \rightarrow L^2(\Omega)$ is a convolution operator, as defined in (5.6), and ω is an additive noise, typically white Gaussian noise.

A natural variational formulation of this problem is

$$\operatorname{argmin}_{\mu \in \mathcal{M}(\Omega)} \frac{1}{2} \|\Phi\mu - y\|^2 + \lambda |\mu|(\Omega), \quad \lambda > 0. \quad (L^2 - |\cdot|)$$

This convex functional is often called BLASSO, which stands for Beurling-LASSO after the work of the mathematician Beurling [19]. It is considered the generalisation of the LASSO $L^2 - L^1$ variational problem in the discrete setting, that is

$$\operatorname{argmin}_{x \in \mathbb{R}^N} \frac{1}{2} \|Tx - y\|^2 + \lambda \|x\|_1, \quad (5.8)$$

with $T = R_L H \in \mathbb{R}^{M \times N}$ and $y \in \mathbb{R}^M$, as in Section 5.1.1. It can be seen indeed as the functional limit of LASSO (5.8) on a finer and finer grid, as sketched in Figure 5.1. In Figure 5.2, we show a visual comparison between reconstructions obtained with on-the-grid approaches with LASSO and off-the-grid approaches with BLASSO.

In [39], the functional associated to $(L^2 - |\cdot|)$

$$T_\lambda(\mu) := \frac{1}{2} \|\Phi\mu - y\|^2 + \lambda |\mu|(\Omega) \quad (5.9)$$

is proved to be proper, convex and coercive. This fact also proves the existence of solutions of $(L^2 - |\cdot|)$. Under injectivity assumptions on the forward operator, uniqueness of the solution is also guaranteed [39].

Off-the-grid deblurring of curves. The reconstruction of images which are the superposition of a few 1D curves is an interesting task in fluorescence microscopy. Even if the images are sparse, it is not realistic to model them as sums of Diracs (5.2), since this yields to dotted reconstructions [142]. This task is investigated in the recent works [139, 140], which propose the penalty

$$R : \mathcal{M}(\Omega)^2 \rightarrow \mathbb{R} \cup \{+\infty\} \\ \mu \mapsto |\mu|(\Omega) + |\operatorname{div}(\mu)|(\Omega),$$

with $\operatorname{div}(\mu)$ being the divergence of μ (defined in the sense of distributions), and the associated CROC (Curves Represented On Charges) functional

$$\operatorname{argmin}_{\mu \in \mathcal{V}} \frac{1}{2} \|\Phi\mu - y\|^2 + \lambda R(\mu)$$

to implement the curve reconstruction problem, with $\mathcal{V} \subseteq \mathcal{M}(\Omega)^2$ space of charges, i.e. the set of $\mu \in \mathcal{M}(\Omega)^2$ with finite divergence $\operatorname{div}(\mu)$. In this way it is possible to obtain measures supported on a curve, i.e. μ_γ such that $\langle \psi, \mu_\gamma \rangle = \int_0^1 \psi(\gamma(t)) \gamma'(t) dt$ with $\psi \in \mathcal{C}_0(\Omega, \mathbb{R}^2)$ and $\gamma(\cdot)$ a regular curve onto $[0, 1]$.

Off-the-grid reconstructions of piece-wise constant functions. For completeness we briefly introduce here the imaging task in microscopy of precise localisation of cells, modelled as piece-wise constant functions, namely $u(\cdot) = a\mathbf{1}_E(\cdot)$ with $E \subset \mathbb{R}^2$, support regions of the cells. This problem has been studied in the off-the-grid setting in [50, 51, 71], where the authors propose to tackle the reconstruction of piece-wise constant images via total-variation regularisation with the functional

$$\operatorname{argmin}_{u \in L^2(\mathbb{R}^2)} \frac{1}{2} \|\Phi u - y\|^2 + \lambda |Du|(\mathbb{R}^2),$$

where $|Du|(\mathbb{R}^2)$ denotes the total variation of the gradient of u

$$|Du|(\mathbb{R}^2) = \sup \left\{ - \int_{\mathbb{R}^2} u \operatorname{div}(\psi) \mid \psi \in \mathcal{C}_c^\infty(\mathbb{R}^2, \mathbb{R}^2), \|\psi\|_\infty \leq 1 \right\}.$$

We stress that this problem is not defined in the space of Radon measures but instead in $L^2(\mathbb{R}^2)$, but it is interesting to mention since it is a continuous off-the-grid problem.

Returning to the analysis of the BLASSO ($L^2 - |\cdot|$) problem, in the following, we discuss its optimality conditions, which can be obtained by studying its corresponding dual problem. Then, we present some algorithmic strategies that can be used for its minimisation in the non-reflexive Banach space $\mathcal{M}(\Omega)$.

5.2.1 Optimality conditions

Since the functional T_λ is convex, μ minimises T_λ if and only if $0 \in \partial T_\lambda(\mu)$, i.e.

$$0 \in \Phi^*(\Phi\mu - y) + \lambda \partial|\mu|(\Omega) \iff \frac{1}{\lambda} \Phi^*(y - \Phi\mu) \in \partial|\mu|(\Omega),$$

which can be written as

$$\eta \in \partial|\mu|(\Omega), \tag{5.10}$$

by defining the so-called dual certificate η as

$$\eta := \frac{1}{\lambda} \Phi^*(y - \Phi\mu). \tag{5.11}$$

The dual certificate plays a crucial role in the characterisation of optimality conditions for ($L^2 - |\cdot|$) and in devising optimisation algorithms in this setting, as better specified in the following sections.

If μ is a finite linear combination of Dirac masses (5.2), the subdifferential of the TV norm, defined in (5.3), takes the form (5.4). Hence, (5.10) becomes simply as

$$\eta(x_i) = \operatorname{sign}(a_i) \quad \wedge \quad \|\eta\|_\infty \leq 1. \tag{5.12}$$

5.2.2 Dual problem and extremality conditions

The optimality conditions (5.12) can be also derived by studying the dual problem of $(L^2 - |\cdot|)$. This will give a further interpretation to the meaning of the dual certificate η . Before detailing the dual problem of $(L^2 - |\cdot|)$, we recall the basic definition of convex conjugate [194] and a standard result from [89].

Definition 5.2.1. *Given $f : \mathcal{X} \rightarrow \mathbb{R}$, with \mathcal{X} real Banach space, the convex conjugate of f is the function $f^* : \mathcal{X}^* \rightarrow \mathbb{R} \cup \{+\infty\}$ defined by*

$$f^*(x^*) := \sup_{x \in \mathcal{X}} \langle x^*, x \rangle_{\mathcal{X}^* \times \mathcal{X}} - f(x), \quad x^* \in \mathcal{X}^*. \quad (5.13)$$

Lemma 5.2.1. (Fenchel-Rockafellar duality.) *Let V, Y be two Banach spaces. For $\Lambda : V \rightarrow Y$ linear operator, $F : V \rightarrow \mathbb{R}$ and $G : Y \rightarrow \mathbb{R}$ convex functionals, we consider the following primal problem:*

$$\operatorname{argmin}_{u \in V} F(u) + G(\Lambda u). \quad (5.14)$$

The corresponding dual problem reads

$$\operatorname{argmax}_{p^* \in Y^*} -F^*(\Lambda^* p^*) - G^*(-p^*), \quad (5.15)$$

where $\Lambda^* : Y^* \rightarrow V^*$ is the adjoint operator of Λ and $F^* : V^* \rightarrow \mathbb{R} \cup \{+\infty\}$, $G^* : Y^* \rightarrow \mathbb{R} \cup \{+\infty\}$ are the convex conjugate of F and G .

Moreover, if $u \in V$ and $p^* \in Y^*$ are respectively solutions of the primal (5.14) and dual (5.15) problems, the following extremality conditions hold:

$$\begin{cases} \Lambda^* p^* \in \partial F(u) \\ -p^* \in \partial G(\Lambda u) \end{cases}. \quad (5.16)$$

In the formulation of the BLASSO $(L^2 - |\cdot|)$ problem, we have $V = \mathcal{M}(\Omega)$, $Y = L^2(\Omega)$, $F(\mu) = |\mu|(\Omega)$ and $G(p) = \frac{1}{2\lambda} \|y - p\|^2$. To write the dual problem, as defined in (5.15), it is necessary to determine the expressions for the convex conjugate of the TV norm and of G .

5.2.2.1 Convex conjugate of the fidelity

We start by detailing the computation of the convex conjugate for the Gaussian L^2 fidelity given by $G(p) = \frac{1}{2\lambda} \|y - p\|^2$. By Definition 5.2.1, we need to compute

$$G^*(p^*) = \sup_{p \in L^2(\Omega)} \langle p^*, p \rangle - \frac{1}{2\lambda} \|y - p\|^2.$$

Being the function $p \in L^2(\Omega) \mapsto \langle p^*, p \rangle - \frac{1}{2\lambda} \|y - p\|^2$ concave, its supremum is attained at p such that $0 \in \partial \left(\langle p^*, \cdot \rangle - \frac{1}{2\lambda} \|y - \cdot\|^2 \right) (p) = p^* - \frac{1}{\lambda}(y - p)$, i.e. for $p = \lambda p^* + y$. Thus, the convex conjugate has expression

$$G^*(p^*) = \langle p^*, y \rangle + \frac{\lambda}{2} \|p^*\|^2, \quad p^* \in L^2(\Omega). \quad (5.17)$$

5.2.2.2 Convex conjugate of the TV norm

To conclude and obtain the dual problem of $(L^2 - |\cdot|)$, we report now the computation of the convex conjugate of $F : V \rightarrow \mathbb{R}$, defined by $F(\mu) = |\mu|(\Omega)$ with $V = \mathcal{M}(\Omega)$. Thus, by Definition 5.1.1, the dual space is $V^* = \mathcal{C}_0(\Omega, \mathbb{R})$. In order to determine $F^* : V^* \rightarrow \mathbb{R} \cup \{+\infty\}$, let $\psi \in \mathcal{C}_0(\Omega, \mathbb{R})$:

$$\begin{aligned} F^*(\psi) &= \sup_{\mu \in \mathcal{M}(\Omega)} \langle \psi, \mu \rangle_{\mathcal{C}_0(\Omega, \mathbb{R}) \times \mathcal{M}(\Omega)} - |\mu|(\Omega) \\ &\geq \langle \psi, \mu \rangle_{\mathcal{C}_0(\Omega, \mathbb{R}) \times \mathcal{M}(\Omega)} - |\mu|(\Omega), \quad \forall \mu \in \mathcal{M}(\Omega). \end{aligned}$$

Let $x \in \Omega$, and $\mu = \alpha \delta_x$ with $\alpha > 0$. Then, we can write:

$$\sup_{\mu \in \mathcal{M}(\Omega)} \langle \psi, \mu \rangle_{\mathcal{C}_0(\Omega, \mathbb{R}) \times \mathcal{M}(\Omega)} - |\mu|(\Omega) \geq \alpha(\psi(x) - 1).$$

Taking the limit for $\alpha \rightarrow \infty$ of the latter inequality yields $F^*(\psi) \geq +\infty$ if $\psi(x) > 1$. A similar result for $\psi(x) < 1$ is obtained with the measure $\mu = -\alpha \delta_x$, with $\alpha > 0$. Thus, we have $F^*(\psi) = +\infty$ if $\|\psi\|_{\infty, \Omega} > 1$.

Assume now that $\|\psi\|_{\infty, \Omega} \leq 1$. First notice that

$$F^*(\psi) = \sup_{\mu \in \mathcal{M}(\Omega)} \langle \psi, \mu \rangle_{\mathcal{C}_0(\Omega, \mathbb{R}) \times \mathcal{M}(\Omega)} - |\mu|(\Omega) \geq 0$$

by considering $\mu = 0$. Moreover,

$$\begin{aligned} \langle \psi, \mu \rangle_{\mathcal{C}_0(\Omega, \mathbb{R}) \times \mathcal{M}(\Omega)} - |\mu|(\Omega) &\leq \|\psi\|_{\infty, \Omega} |\mu|(\Omega) - |\mu|(\Omega) \\ &\leq |\mu|(\Omega) (\|\psi\|_{\infty, \Omega} - 1) \leq 0 \end{aligned}$$

since $\|\psi\|_{\infty, \Omega} \leq 1$. By taking the sup on both sides of the last inequality, one finally gets $F^*(\psi) = 0$ if $\|\psi\|_{\infty, \Omega} \leq 1$. Altogether, this yields to

$$F^*(\psi) = \begin{cases} 0 & \|\psi\|_{\infty, \Omega} \leq 1 \\ +\infty & \|\psi\|_{\infty, \Omega} > 1 \end{cases}. \quad (5.18)$$

5.2.2.3 Dual problem and extremality conditions

We now have all the elements to write the dual problem formulation (5.15) for the BLASSO problem $(L^2 - |\cdot|)$, putting together the result in Lemma 5.2.1 with (5.17) and (5.18):

$$\operatorname{argmax}_{\|\Phi^* p^*\|_{\infty, \Omega} \leq 1} \langle y, p^* \rangle - \frac{\lambda}{2} \|p^*\|^2. \quad (5.19)$$

Moreover, from Lemma 5.2.1, the extremality conditions can be obtained. Given $\mu_\lambda \in \mathcal{M}(\Omega)$ solution of the primal problem $(L^2 - |\cdot|)$ with regularisation parameter $\lambda > 0$ and denoting by $p_\lambda^* \in L^2(\Omega)$ the solution of the dual problem (5.19), by (5.16) we get

$$\begin{cases} \Phi^* p_\lambda^* \in \partial |\mu_\lambda|(\Omega) \\ -p_\lambda^* = \frac{1}{\lambda} (\Phi \mu_\lambda - y) \end{cases}. \quad (5.20)$$

From the extremality conditions (5.20), we retrieve the optimality conditions (5.12), expressed in terms of the dual certificate (5.11). In addition, when the dual certificate satisfies (5.12) we can interpret it as $\eta = \Phi^* p_\lambda^*$ with p_λ^* being a solution of (5.19).

Optimality conditions characterise the solution(s) of the BLASSO problem ($L^2 - |\cdot|$) and they are crucial in devising algorithms for its minimisation and, in particular, in the definition of stopping rules, since a good stopping criterion stops at an iterate that satisfies the optimality conditions. In the next section, we present a possible algorithmic strategy to solve the BLASSO problem ($L^2 - |\cdot|$).

5.3 Frank-Wolfe and Sliding Frank-Wolfe algorithms

The BLASSO ($L^2 - |\cdot|$) is an optimisation problem over the space of Radon measures, an infinite dimensional and non-reflexive Banach space. Due to the non-reflexivity of $\mathcal{M}(\Omega)$, it is not clear how to define proximal strategies in this context. A preliminar result has recently been proposed in [215].

A possible solver for BLASSO has to take into account the infinite dimensional nature of $\mathcal{M}(\Omega)$. The most used algorithms in this setting are semi-definite programming approaches (for Fourier measurements (5.7)) [96], conditional gradient algorithms [75, 97] and particle gradient descent [57, 58], which is an optimal transport based algorithm. For a complete review of all these methods, we refer to [141].

A sketch of the Conic Particle Gradient Descent algorithm is shown in Figure 5.5. It proposes to solve BLASSO through optimal transport with gradient flows. Basically, starting from an initial over-parametrised measure, with particles that cover the domain Ω , it discretises the measure and, by performing a non-convex gradient descent on the positions and weights of the particles, approximates the gradient flow, which is proved to converge to a solution of BLASSO. Conic Particle Gradient Descent estimated the gradient flow through a gradient descent on both amplitudes and positions computed with respect to a specific cone metric.

In the next sections, we will present the conditional gradient method, initially proposed in 1956 in [97] by Frank and Wolfe, showing how it can be used as a solver for BLASSO ($L^2 - |\cdot|$). Then, we describe the Sliding Frank-Wolfe algorithm proposed in [75], for off-the-grid optimisation.

5.3.1 Frank-Wolfe algorithm

The Frank-Wolfe (FW) algorithm, or conditional gradient method, has been proposed in [97] to solve the following optimisation problem

$$\min_{m \in C} f(m) \tag{5.21}$$

where C is a weakly compact convex set of a Banach space, and f is a differentiable convex function. It relies on the iterative minimisation of a linearised version of

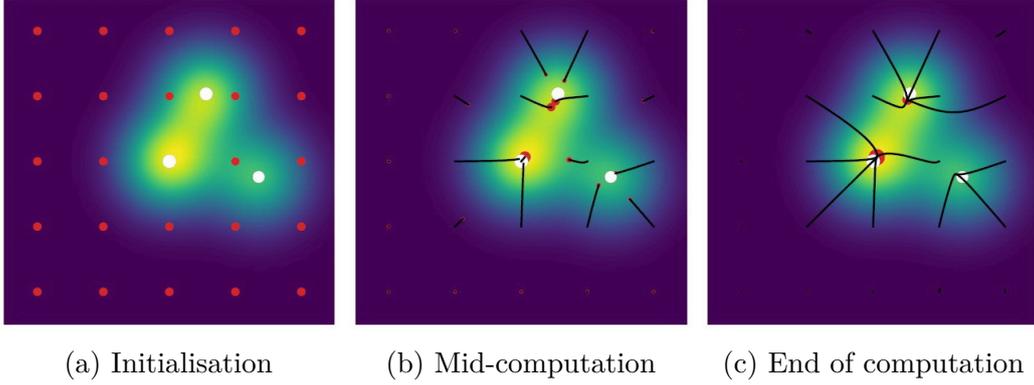


Figure 5.5: Conic particle gradient descent applied for 2D sparse spike deconvolution with Gaussian kernel. White dots are the source measure and red dots are the measures at iterations k with $k = 0$ in (a), $k = 150$ in (b) and $k = 1000$ in (c). The background image is the acquisition y . The black lines are the paths of the particles and constitute the gradient flow. Images from [141].

f without requiring any Hilbertian structure. This is a key advantage of FW with respect to most first order optimisation schemes, such as gradient descent or proximal splitting method, which on the contrary rely on an underlying Hilbertian structure. This fact makes the FW algorithm particularly interesting to work in $\mathcal{M}(\Omega)$. The pseudo-code of the FW algorithm is detailed in Algorithm 9.

Observe that the FW algorithm is naturally endowed with a stopping criterion for the iterates m^k which is equivalent to the standard optimality condition for constrained convex problems [73, 75], i.e.

$$\forall s \in C, \quad \mathrm{d}f(m^k)(s - m^k) \geq 0.$$

In this section, we denote with $\mathrm{d}f$ the directional derivative of f . Moreover, the stopping criterion of Algorithm 9 ensures that the method stops at a global minimiser. Indeed, since f is convex, for any $s, t \in C$ there holds

$$f(s) \geq f(t) + \mathrm{d}f(t)(s - t). \quad (5.22)$$

Then, since s^k is a minimiser of (5.25) the convexity property (5.22) implies that, for all $s \in C$,

$$f(s) \geq f(m^k) + \mathrm{d}f(m^k)(s - m^k) \geq f(m^k) + \mathrm{d}f(m^k)(s^k - m^k). \quad (5.23)$$

If the stopping criterion is satisfied by m^k , that is if

$$\mathrm{d}f(m^k)(s^k - m^k) = 0, \quad (5.24)$$

then (5.23) together with (5.24) ensures that m^k is a global minimiser of f , i.e. $f(s) \geq f(m^k)$ for all $s \in C$.

Another interesting property of this method is that it is possible to replace m^{k+1} in (5.27) by any element of $\tilde{m} \in C$ such that $f(\tilde{m}) \leq f(m^{k+1})$ without losing the convergence properties of this algorithm [33, 124].

Algorithm 9 Frank-Wolfe (FW) algorithm [97]

Initialisation: $m^0 \in C$.

repeat

Solve

$$s^k \in \operatorname{argmin}_{s \in C} f(m^k) + df(m^k)[s - m^k] \quad (5.25)$$

if $df(m^k)[s^k - m^k] = 0$
 m^k is a solution of (5.21) \Rightarrow **stop**
else

Step research:

$$\gamma^k = \frac{2}{k+2} \quad \vee \quad \gamma^k \in \operatorname{argmin}_{\gamma \in [0,1]} f(m^k + \gamma(s^k - m^k)) \quad (5.26)$$

Update:

$$m^{k+1} = m^k + \gamma^k(s^k - m^k) \quad (5.27)$$

until convergence

5.3.2 Frank-Wolfe for the minimisation of BLASSO

The FW algorithm cannot be applied straightforwardly to the BLASSO problem because $(L^2 - |\cdot|)$ is an optimisation problem over $\mathcal{M}(\Omega)$, which is not bounded, and the objective function (5.9) is not differentiable. Instead, following an idea of [111], in [75] the authors propose to consider an equivalent problem to the BLASSO by defining a differentiable epigraphical lift, which shares the same minimum as T_λ .

Lemma 5.3.1. *The BLASSO problem $(L^2 - |\cdot|)$ is equivalent to*

$$\operatorname{argmin}_{(t,\mu) \in C} \tilde{T}_\lambda(\mu, t) := \frac{1}{2} \|\Phi\mu - y\|^2 + \lambda t, \quad (5.28)$$

where

$$C := \left\{ (t, \mu) \in \mathbb{R}^+ \times \mathcal{M}(\Omega) \mid |\mu|(\Omega) \leq t \leq M \right\} \quad (5.29)$$

with $M := \frac{\|y\|^2}{2\lambda}$.

Proof. Let μ_λ be a minimiser of T_λ . Then,

$$|\mu_\lambda|(\Omega) \leq \frac{1}{\lambda} T_\lambda(\mu_\lambda) \leq \frac{1}{\lambda} T_\lambda(0) \leq \frac{\|y\|^2}{2\lambda}.$$

Therefore, it suffices to restrict the minimisation of BLASSO over the set of measures with $|\mu|(\Omega) \leq t \leq \frac{\|y\|^2}{2\lambda}$. \square

The equivalence stated in Lemma 5.3.1 is to be understood in the following sense: μ is a solution to $(L^2 - |\cdot|)$ if and only if (t, μ) is a solution to (5.28) for some $t \geq 0$. Moreover, if this is the case then $t = |\mu|(\Omega)$ and $\tilde{T}_\lambda(\mu, t) = T_\lambda(\mu)$. As a result, Lemma 5.3.1 allows to rewrite the minimisation of T_λ over $\mathcal{M}(\Omega)$ into the form (5.21), since the two problems are equivalent, C is bounded and \tilde{T}_λ is a differentiable functional on the Banach space $\mathbb{R} \times \mathcal{M}(\Omega)$, with differential

$$d\tilde{T}_\lambda(t, \mu) : (t', \mu') \longmapsto \int_{\Omega} \Phi^*(\Phi\mu - y) d\mu' + \lambda t'. \quad (5.30)$$

In [73], a linear rate of convergence in function values has been shown for any minimising sequence for the BLASSO.

Lemma 5.3.2. *Let $(t_k, \mu^k)_{k \in \mathbb{N}}$ be the sequence generated by Algorithm 9 applied to \tilde{T}_λ (5.28). Then, there exists $D > 0$ such that for any solution μ^* of BLASSO we have*

$$T_\lambda(\mu^k) - T_\lambda(\mu^*) \leq \frac{D}{k} \quad \forall k \in \mathbb{N}.$$

5.3.2.1 Greedy approach

Now we detail how the algorithmic steps of FW applied to (5.28) for the minimisation of BLASSO yields a greedy approach.

Consider the minimisation step (5.25) of Algorithm 9. First of all, note that (5.25), neglecting terms which are constant with respect to s , reads as

$$s^k \in \operatorname{argmin}_{s \in C} df(m^k)(s). \quad (5.31)$$

Moreover, since the function to minimise in (5.31) is a linear form and C is convex, it achieves a minimum at an extremal point of C .

Applying Algorithm 9 to (5.28), we have $f(s) = \tilde{T}_\lambda(s)$ with $s = (t, \mu)$. In our case, C is defined by (5.29) and has extremal points of the form $s = (M, \pm M\delta_x)$ with $x \in \Omega$. Thus, (5.31) (which is equivalent to the first step (5.25) of each iteration of Algorithm 9), by denoting $m^k = (t^k, \mu^k)$, becomes

$$\begin{aligned} \operatorname{argmin}_{(t, \mu) \in C} d\tilde{T}_\lambda(t^k, \mu^k)(t, \mu) &= \operatorname{argmin}_{x \in \Omega} \pm M(\Phi^*(\Phi\mu^k - y))(x) + \lambda M \\ &= \operatorname{argmin}_{x \in \Omega} \left(\pm(\Phi^*(\Phi\mu^k - y))(x) + 1 \right) \lambda M \\ &= \operatorname{argmin}_{x \in \Omega} \mp \eta^k(x) + 1 \quad \text{where } \eta^k = \Phi^*(y - \Phi\mu^k) \\ &= \operatorname{argmax}_{x \in \Omega} \left| \eta^k(x) \right|. \end{aligned}$$

This shows that

$$s^k = m^k + (M, \sigma M\delta_{x_*^k}) = (t^k, \mu^k) + (M, \sigma M\delta_{x_*^k}) \quad \text{with } \sigma = \operatorname{sign}(\eta^k(x_*^k)). \quad (5.32)$$

Therefore, at each iteration of the algorithm a new support point

$$x_*^k = \operatorname{argmax}_{x \in \Omega} |\eta^k(x)| \quad (5.33)$$

is introduced, resulting in a new spike at position x_*^k , with amplitude σM .

It is interesting to notice the similarity between η^k , introduced above, and the dual certificate (5.11). Recalling that $d\tilde{T}_\lambda$ is given by (5.30), the stopping criterion (5.24) here becomes

$$\begin{aligned} d\tilde{T}_\lambda(t^k, \mu^k)(s^k - m^k) &= 0 \quad \text{where by (5.32)} \quad s^k - m^k = (M, \sigma M \delta_{x_*^k}) \\ \sigma M \Phi^*(\Phi \mu^k - y)(x_*^k) + \lambda M &= 0 \\ -\sigma \eta^k(x_*^k) + 1 &= 0 \\ |\eta^k(x_*^k)| &= 1. \end{aligned}$$

By definition of x_*^k (5.33), it means that

$$|\eta^k(x)| \leq 1 \text{ for any } x \in \Omega,$$

and thus η^k is a dual certificate and satisfies optimality conditions (5.12). The algorithm iteratively constructs such a dual certificate.

Let analyse now the line-search step (5.26) of Algorithm 9. Without loss of generality, we can assume $\mu^k = \sum_{i=1}^k a_i^k \delta_{x_i^k}$ with $a^k = (a_1^k, \dots, a_k^k)$ and $t^k = \|a^k\|_1$, then

$$m^k + \gamma(s^k - m^k) = (t^k, \mu^k) + \gamma(M, \sigma M \delta_{x_*^k}) = (t^k + \gamma M, \mu_\gamma),$$

where $\mu_\gamma = \mu^k + \gamma \sigma M \delta_{x_*^k}$. Then, (5.26) reads as

$$\operatorname{argmin}_{\gamma \in [0,1]} \frac{1}{2} \|\Phi \mu_\gamma - y\|^2 + (1 - \gamma) \lambda \|a^k\|_1 + \gamma \lambda M. \quad (5.34)$$

Note that since this step can be replaced with any (t, μ) which improves the objective value [124], it seems sensible to simply perform a LASSO step in the form

$$a^{k+1} = \operatorname{argmin}_{a \in \mathbb{R}^{k+1}} \frac{1}{2} \|\Phi \mu_a - y\|^2 + \lambda \|a\|_1, \quad (5.35)$$

where

$$\mu_a = \sum_{i=1}^k a_i \delta_{x_i^k} + a_{k+1} \delta_{x_*^k} \quad (5.36)$$

This is a finite dimensional non-smooth convex optimisation problem and can be tackled using a variety of algorithms such as Forward Backward or FISTA [13].

Basically, by (5.36) at each step a new Dirac is added to the reconstruction. In this sense, FW can be described as a greedy approach that iteratively finds a better approximation to the desired solution adding just one spike per iteration.

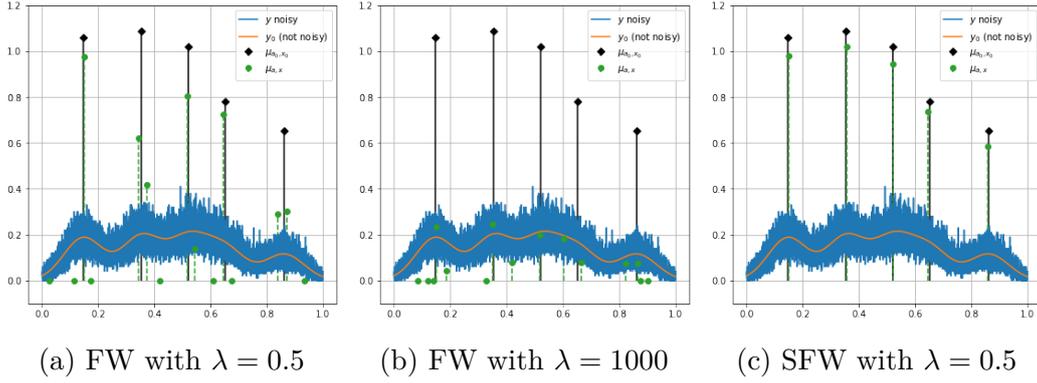


Figure 5.6: Reconstruction of 1D peaks from a blurred and noisy signal obtained using Frank-Wolfe algorithm for $\lambda = 0.5$ and $\lambda = 1000$, and with Sliding Frank-Wolfe with $\lambda = 0.5$.

5.3.3 Sliding Frank-Wolfe algorithm

Applying directly Algorithm 9 yields a sequence of measures $(\mu^k)_{k \in \mathbb{N}}$ which converges towards a solution in a greedy way. However, numerical tests show that the generated measures μ^k tend to be less sparse compared to the desired solution. Indeed, one can practically observe that each Dirac mass of the ground truth measure is approximated in μ^k by a cluster of Dirac masses with inexact positions, as shown in Figure 5.6 (and this undesired phenomenon does not disappear by changing λ). In [33, 39], the authors thus suggest to modify the Frank-Wolfe iterations for the resolution of the BLASSO and to let the Dirac positions slightly move.

To this aim, the fact that one can replace the update of step (5.27) in Frank-Wolfe algorithm (Alg. 9) by any value which improves the objective [124] is very useful. Before, this idea allowed to replace (5.34) with (5.35), for the estimation of the amplitudes. Similarly, one can further boost this step by optimising over the positions and the amplitudes simultaneously. This consideration allows to define the modified version of Frank-Wolfe in Algorithm 10 that has been proposed and studied in [75], namely the Sliding Frank-Wolfe (SFW) algorithm. The non-convex update step (5.38) of Algorithm 10 effectively decreases more the energy than the standard convex optimisation over the spikes amplitudes (5.37), as done in (5.35) with Frank-Wolfe, which in turn decreases the functional T_λ more than performing (5.34) using γ^k . Indeed,

$$T_\lambda(\mu^{k+1}) \leq T_\lambda(\mu^{k+1/2}) \leq T_\lambda((1 - \gamma^k)\mu^k + \gamma^k \sigma M \delta_{x_*^k}).$$

The non-convex step drastically improves the convergence property of the algorithm [32, 33]. Indeed, this tweak yields a theoretical convergence to a solution of BLASSO in a finite number of iterations [75]. Moreover, Figure 5.7 highlights the importance of the sliding step in the recovery of precise localised Diracs. Allowing the positions to slightly change at the end of each iteration corrects the position estimated in the insertion step with the argmax of the certificate, which does not necessarily

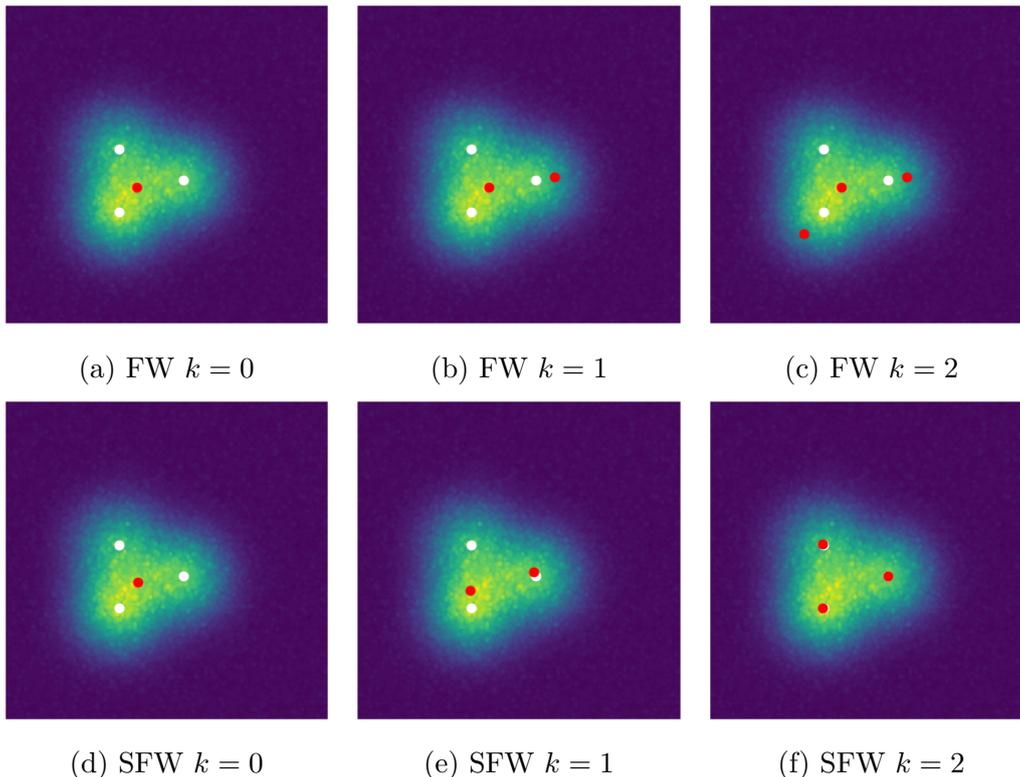


Figure 5.7: Frank-Wolfe and Sliding Frank-Wolfe algorithms for 2D sparse spike deconvolution with Gaussian kernel. White dots are the ground truth measure and red dots the reconstructed measure at iterate k . The background image is the acquisition y , obtained with a 2D Gaussian PSF with $\sigma = 0.1$ and Poisson noise. Results with $\lambda = 0.001$ on the domain $\Omega = [0, 1]^2$. The importance of the sliding step appears evident in this 2D example with 3 Diracs.

correspond to a spike position in the ground truth. Comparing the results given by FW and SFW in Figure 5.7, we observe that the localisation precision of SFW is way higher.

5.4 Final discussion

In this chapter, we presented an overview on the mathematical theory of sparse off-the-grid optimisation in imaging. The formulation of sparse inverse problems in an off-the-grid setting entails working in the space of Radon measures $\mathcal{M}(\Omega)$ and, in the case of a quadratic data term, yields to the BLASSO variational problem. BLASSO combines an L^2 fidelity with the TV norm of measures as penalty. Moreover, since the TV norm in $\mathcal{M}(\Omega)$ is a generalisation of the L^1 norm, BLASSO can be interpreted as a generalisation of the LASSO problem in \mathbb{R}^d , used in discrete settings, to the continuous scenarios of off-the-grid inverse problems. The TV norm enforces sparsity and, in particular, it favours the reconstruction of finite linear

Algorithm 10 Sliding Frank-Wolfe (SFW) algorithm [75]

Initialisation: $\mu^0 = 0$.

repeat

$\mu^k = \sum_{i=1}^{N^k} a_i^k \delta_{x_i^k}$, $a_i^k \in \mathbb{R}$, $x_i^k \in \Omega$, find $x_*^k \in \Omega$ s.t.:

$$x_*^k \in \operatorname{argmax}_{x \in \Omega} |\eta^k(x)| \quad \text{where} \quad \eta^k = \frac{1}{\lambda} \Phi^* (y - \Phi \mu^k)$$

if $|\eta^k(x_*^k)| \leq 1$

μ^k is a solution of $(L^2 - |\cdot|) \Rightarrow$ **stop**

else

– **insertion step:**

Add support for the new spike:

$$x^{k+1/2} = (x_1^k, \dots, x_{N^k}^k, x_*^k)$$

Estimation of amplitudes:

$$a^{k+1/2} \in \operatorname{argmin}_{a \in \mathbb{R}^{N^{k+1}}} \frac{1}{2} \|\Phi_{x^{k+1/2}} a - y\|_{\mathcal{H}}^2 + \lambda \|a\|_1 \quad (5.37)$$

Update:

$$\mu^{k+1/2} = \sum_{i=1}^{N^k} a_i^{k+1/2} \delta_{x_i^k} + a_{N^{k+1}}^{k+1/2} \delta_{x_*^k}$$

– **sliding step:**

Using a non-convex solver initialised with $(a^{k+1/2}, x^{k+1/2})$

$$(a^{k+1}, x^{k+1}) \in \operatorname{argmin}_{(a,x) \in \mathbb{R}^{N^{k+1}} \times \Omega^{N^{k+1}}} \frac{1}{2} \|\Phi_x a - y\|_{\mathcal{H}}^2 + \lambda \|a\|_1 \quad (5.38)$$

$$\mu^{k+1} = \sum_{i=1}^{N^{k+1}} a_i^{k+1} \delta_{x_i^{k+1}}$$

– **pruning:** eventually remove zero amplitudes Dirac masses from μ^{k+1}

until convergence

combinations of Diracs, which is particularly suited for applications. In the next chapter, we will consider fluorescence microscopy (see Appendix B) as a biological application for corresponding numerical tests.

From the optimisation point of view, BLASSO can be solved using conditional

gradient algorithms or particle gradient descent strategies. We discussed here only conditional gradient algorithms and, more precisely, the Frank-Wolfe and Sliding Frank-Wolfe algorithms.

Off-the-grid regularisation for Poisson inverse problems

Off-the-grid methods are usually formulated in the standard additive Gaussian noise setting, i.e. with an L^2 data fidelity term. However, in many scenarios other noise statistics better describe the underlying acquisition process. In this chapter, we aim to investigate the Poisson noise modelling for inverse problems in the space of Radon measures. We thus propose a variational model which couples a Kullback-Leibler data fitting term with the Total Variation of measures, as penalty term, together with a non-negativity constraint and study its optimality conditions by analysing the corresponding dual problem. The importance of the choice of a good regularisation parameter is addressed through an automatic selection strategy, based on the homotopy criterion. We conclude the chapter providing numerical experiments on both 1D/2D simulated and real 3D fluorescent microscopy data.

6.1	Off-the-grid Poisson inverse problems	128
6.2	Dual problem and optimality conditions	130
6.2.1	Convex conjugate of the Kullback-Leibler divergence	130
6.2.2	Convex conjugate of the sum of penalty and indicator function	131
6.2.3	Dual problem formulation and extremality conditions	133
6.2.3.1	Subdifferential of the indicator function of positive measures	135
6.2.3.2	Extremality conditions	136
6.2.4	Optimality conditions	137
6.3	Algorithmic choices	137
6.3.1	Boosted Sliding Frank Wolfe algorithm	138
6.4	Homotopy algorithm	139
6.4.1	Homotopy algorithm for off-the-grid methods	142

6.4.1.1	Starting value	143
6.4.1.2	Updating rule	144
6.4.1.3	Descent property of the homotopy algorithm	144
6.4.1.4	Regularisation path in the discrete setting	146
6.4.1.5	Homotopy and regularisation path of BLASSO	147
6.5	Numerical tests	150
6.5.1	Comparison between Gaussian and Poisson modelling	150
6.5.2	Homotopy algorithm: choice of σ_{target}	154
6.5.3	Numerical test with 3D real dataset	156
6.6	Final discussion	158

In this chapter, we study inverse problems in the space of Radon measures $\mathcal{M}(\Omega)$ under the hypothesis of signal dependent Poisson noise in the data. This choice is motivated by the particular biological application of interest, that is fluorescence microscopy. We refer to Appendix B for more details on this imaging problem. Due to the photon emission nature of the light, in microscopy imaging Poisson noise is better suited than the Gaussian one to describe the photon counts on acquired images [16, 144]. However, a Gaussian noise modelling is often preferred, as it is in general easier to work with, from both an analytical and a computational point of view. In the space of Radon measures $\mathcal{M}(\Omega)$, it leads to the well studied variational formulation of BLASSO, presented in the previous chapter. The contribution of this chapter is the definition and detailed analytical and numerical study of off-the-grid regularisation coupled with a Poisson data term.

6.1 Off-the-grid Poisson inverse problems

Similarly as in Section 5.2, we aim at solving the spike deconvolution problem under a Poisson noise hypothesis. To this aim, some adjustments to the inverse problem formulation are required and now discussed.

Let $\Phi : \mathcal{M}(\Omega) \rightarrow L^2(\Omega)$ the forward operator in (5.5). In the following, Φ is assumed to be a positive definite operator, i.e.

$$\mu \in \mathcal{M}^+(\Omega) \text{ positive measure} \Rightarrow \Phi\mu(x) \geq 0 \quad \forall x \in \Omega. \quad (6.1)$$

Observe that, with Φ defined as in (5.5), this is attained whenever the convolution kernel ψ is non-negative. We remark that in deblurring imaging problem the convolution kernel ψ is naturally non-negative, hence the forward operator Φ is always positive definite. We adopt the notation $L^2(\Omega)^+$ meaning

$$L^2(\Omega)^+ = \{f \in L^2(\Omega) \text{ such that } f(x) > 0 \text{ a .e. } x \in \Omega\}.$$

The Poisson sparse spike deconvolution problem in an off-the-grid setting consists in finding a positive discrete measure $\mu \in \mathcal{M}^+(\Omega)$ such that

$$y = \mathcal{P}(\Phi\mu + b),$$

where $b \in L^2(\Omega)^+$ is a strictly positive background term, and $y \in L^2(\Omega)$ is a realisation of a Poisson distributed random variable Y with mean $\Phi\mu + b > 0$.

We want to define a variational formulation well-suited to Poisson noise scenarios and thus the Kullback-Leibler divergence is now considered as fidelity term [4, 15, 17, 154, 200].

Definition 6.1.1. *The Kullback-Leibler divergence $\mathcal{D}_{KL} : L^2(\Omega)^+ \times L^2(\Omega)^+ \rightarrow \mathbb{R}$ is defined by*

$$\mathcal{D}_{KL}(s, t) := \int_{\Omega} s(x) - t(x) + t(x) \log(t(x)) - t(x) \log(s(x)) \, dx. \quad (6.2)$$

We would like to consider $\mathcal{D}_{KL}(\Phi\mu + b, y)$ as data fidelity term, but we observe that this is not well-defined since

- if μ is not a positive measure, then $\phi\mu + b$ might not be positive;
- although the assumption that Φ is a positive definite operator (6.1), together with the strict positivity of the background b , ensures that $\Phi\mu + b > 0$ almost everywhere when μ is positive, the noisy acquisition might still be null in a non negligible region of the domain Ω , that is $y \geq 0$. This is due to the fact that a Poisson random variable with mean α assumes the value of 0 with positive probability equal to $e^{-\alpha}$, and thus y might be equal to 0. For this reason, for (6.2) to be well defined it is required that $y > 0$ almost everywhere.

To solve the first issue, we introduce the function $\tilde{\mathcal{D}}_{KL} : L^2(\Omega) \times L^2(\Omega)^+ \rightarrow \mathbb{R} \cup \{+\infty\}$ which extends the definition of (6.2) as

$$\tilde{\mathcal{D}}_{KL}(s, t) = \begin{cases} \mathcal{D}_{KL}(s, t) & s \in L^2(\Omega)^+ \\ +\infty & s \notin L^2(\Omega)^+ \end{cases}. \quad (6.3)$$

Moreover, we just restrict our study to a positive acquisition y for being able to formulate a Poisson noise equivalent to BLASSO, namely we require

$$y \in L^2(\Omega)^+. \quad (6.4)$$

Under hypothesis (6.1) and (6.4) and using the Kullback-Leibler defined by (6.3) the quantity $\tilde{\mathcal{D}}_{KL}(\Phi\mu + b, y)$ is well-defined for all $\mu \in \mathcal{M}(\Omega)$.

Hence, we can now introduce the following variational problem

$$\operatorname{argmin}_{\mu \in \mathcal{M}(\Omega)} \tilde{\mathcal{D}}_{KL}(\Phi\mu + b, y) + \lambda |\mu|(\Omega) + \mathbf{1}_{\{\mathcal{M}^+(\Omega)\}}(\mu), \quad \lambda > 0, \quad (\tilde{\mathcal{D}}_{KL} - |\cdot|)$$

where the Poisson fidelity term, given by the Kullback-Leibler (6.3), is coupled with the TV norm, being a sparsity enforcing penalty, together with the indicator function of the positive measures $\mathcal{M}^+(\Omega)$, to enforce μ to be non-negative.

In the next section, we study the corresponding dual problem and we obtain optimality conditions for $(\tilde{\mathcal{D}}_{KL} - |\cdot|)$.

6.2 Dual problem and optimality conditions

Similarly as in Section 5.2.2, we analyse here the dual problem of $(\tilde{\mathcal{D}}_{KL} - |\cdot|)$ and we provide an analytical expression of the convex conjugate of the involved functions.

Following Lemma 5.2.1, for the definition of the dual problem we need to compute:

- $G^* : L^2(\Omega) \rightarrow \mathbb{R} \cup \{+\infty\}$, convex conjugate of $G : L^2(\Omega) \rightarrow \mathbb{R} \cup \{+\infty\}$ defined as $G(\cdot) := \frac{1}{\lambda} \tilde{\mathcal{D}}_{KL}(\cdot, y)$;
- $F^* : \mathcal{C}_0(\Omega, \mathbb{R}) \rightarrow \mathbb{R} \cup \{+\infty\}$, convex conjugate of the sum of the penalty and the indicator function of $\mathcal{M}^+(\Omega)$, i.e. of $F(\cdot) := |\cdot|(\Omega) + \mathbb{1}_{\{\mathcal{M}^+(\Omega)\}}(\cdot)$, with $F : \mathcal{M}(\Omega) \rightarrow \mathbb{R} \cup \{+\infty\}$.

6.2.1 Convex conjugate of the Kullback-Leibler divergence

To compute the convex conjugate of the Kullback-Leibler $G(\cdot) = \frac{1}{\lambda} \tilde{\mathcal{D}}_{KL}(\cdot, y)$, we start by considering the one-dimensional Kullback-Leibler function, defined by

$$g_t(s) = \frac{1}{\lambda} \left(s - t + t \log(t) - t \log(s) \right), \quad s, t > 0 \text{ and } \lambda > 0.$$

Applying the definition of convex conjugate (5.13) to g_t yields

$$\begin{aligned} g_t^*(s^*) &= \sup_{s>0} s s^* - g_t(s) = \sup_{s>0} s s^* - \frac{1}{\lambda} \left(s - t + t \log(t) - t \log(s) \right) = \\ &= \sup_{s>0} \underbrace{s \left(s^* - \frac{1}{\lambda} \right) + \frac{t}{\lambda} \log(s) + \frac{t}{\lambda} - \frac{t}{\lambda} \log(t)}_{h(s)}. \end{aligned}$$

We have two cases:

- If $s^* \geq \frac{1}{\lambda}$, then $\lim_{s \rightarrow +\infty} h(s) = +\infty$ implies $\sup_{s>0} h(s) = +\infty \Rightarrow g_t^*(s^*) = +\infty$.
- If $0 < s^* < \frac{1}{\lambda}$, then $\lim_{s \rightarrow \pm\infty} h(s) = -\infty$. Thus, being h a convex and differentiable function its supremum is attained at \hat{s} such that $h'(\hat{s}) = 0$, which can be computed

$$\begin{aligned} h'(\hat{s}) &= s^* - \frac{1}{\lambda} + \frac{t}{\lambda \hat{s}} = \frac{\lambda \hat{s} s^* - \hat{s} + t}{\lambda \hat{s}} = 0 \\ &\iff \lambda \hat{s} s^* - \hat{s} + t = 0 \iff \hat{s} = \frac{t}{1 - \lambda s^*}. \end{aligned}$$

Thus,

$$g_t^*(s^*) = h\left(\frac{t}{1 - \lambda s^*}\right) = -\frac{t}{\lambda} \log(1 - \lambda s^*).$$

Observe that $g_t^*(s^*)$ is well defined since $1 - \lambda s^* > 0 \iff s^* < \frac{1}{\lambda}$, which is exactly the case (ii) we are considering.

Hence, the convex conjugate g_t^* of g_t is defined by

$$g_t^*(s^*) = \begin{cases} +\infty & s^* \geq \frac{1}{\lambda} \\ -\frac{t}{\lambda} \log(1 - \lambda s^*) & s^* < \frac{1}{\lambda} \end{cases}. \quad (6.5)$$

Since $\tilde{\mathcal{D}}_{KL}(\cdot, t)$ is defined also for non-positive functions, its 1-dimensional counterpart is given by $\tilde{g}_t : \mathbb{R} \rightarrow \mathbb{R} \cup \{+\infty\}$ defined by

$$\tilde{g}_t(s) = \begin{cases} g_t(s) & s > 0 \\ +\infty & s \leq 0 \end{cases}.$$

Its convex conjugate coincides with (6.5),

$$\tilde{g}_t^*(s^*) = \sup_{s \in \mathbb{R}} s s^* - \tilde{g}_t(s) = \sup_{s > 0} s s^* - g_t(s) = \begin{cases} +\infty & s^* \geq \frac{1}{\lambda} \\ -\frac{t}{\lambda} \log(1 - \lambda s^*) & s^* < \frac{1}{\lambda} \end{cases}. \quad (6.6)$$

It is now possible to state the following lemma.

Lemma 6.2.1. *Consider the function $G : L^2(\Omega) \rightarrow \mathbb{R} \cup \{+\infty\}$ defined by*

$$G(\cdot) := \frac{1}{\lambda} \tilde{\mathcal{D}}_{KL}(\cdot, y),$$

where $\tilde{\mathcal{D}}_{KL}$ is given by (6.3).

The convex conjugate of G is given by $G^ : L^2(\Omega) \rightarrow \mathbb{R} \cup \{+\infty\}$ defined by*

$$\begin{aligned} G^*(p^*) &= \begin{cases} +\infty & p^* \geq \frac{1}{\lambda} \\ -\frac{y}{\lambda} \log(1 - \lambda p^*) & p^* < \frac{1}{\lambda} \end{cases} \\ &= -\frac{y}{\lambda} \log(1 - \lambda p^*) \mathbf{1}_{\{z < 1/\lambda\}}(p^*), \end{aligned} \quad (6.7)$$

where $-\frac{y}{\lambda} \log(1 - \lambda p^*) = \langle -\frac{y}{\lambda}, \log(1 - \lambda p^*) \rangle \in \mathbb{R}$.

Proof. The computation of (6.7) is straightforwardly analogous to the 1-dimensional case given by (6.6). \square

6.2.2 Convex conjugate of the sum of penalty and indicator function

We compute now the convex conjugate of the sum of the TV penalty and the indicator function of non-negative measures, that is the convex conjugate of

$$F : \mathcal{M}(\Omega) \rightarrow \mathbb{R} \cup \{+\infty\}, \quad F(\cdot) = |\cdot|(\Omega) + \mathbf{1}_{\{\mathcal{M}^+(\Omega)\}}(\cdot). \quad (6.8)$$

To simplify the notation, we denote in the following the TV norm and the indicator function of positive measures respectively as $A, B : \mathcal{M}(\Omega) \rightarrow \mathbb{R} \cup \{+\infty\}$

$$A(\cdot) = |\cdot|(\Omega), \quad B(\cdot) = \mathbf{1}_{\{\mathcal{M}^+(\Omega)\}}(\cdot).$$

It is possible to obtain the convex conjugate of the sum of functions as infimal convolution of the convex conjugate of each singular function. We report here a result from [94, 118].

Proposition 6.2.1. *Let $f_1, \dots, f_n : \mathcal{X} \rightarrow \mathbb{R}$ with \mathcal{X} real Banach space. Suppose there exists at least a point $x \in \mathcal{X}$ such that f_1, \dots, f_n are continuous in x . Then,*

$$(f_1 + \dots + f_n)^*(x^*) = \min_{x_1^* + \dots + x_n^* = x^*} f_1^*(x_1^*) + \dots + f_n^*(x_n^*).$$

Thus, following the proposition above we need to compute A^* , B^* first and then their infimal convolution:

$$F^*(\psi) = \min_{\psi_1 + \psi_2 = \psi} A^*(\psi_1) + B^*(\psi_2). \quad (6.9)$$

In Section 5.2.2, we reported the convex conjugate of $A(\cdot)$, the TV-norm, which is given by (5.18). Thus, we only need to compute the convex conjugate of $B(\cdot)$, the indicator function.

Lemma 6.2.2. *Let $B : \mathcal{M}(\Omega) \rightarrow \mathbb{R} \cup \{+\infty\}$, $B(\cdot) = \mathbb{1}_{\{\mathcal{M}^+(\Omega)\}}(\cdot)$, be the indicator function of $\mathcal{M}^+(\Omega)$. The convex conjugate of B is the function $B^* : \mathcal{C}_0(\Omega, \mathbb{R}) \rightarrow \mathbb{R} \cup \{+\infty\}$ given by*

$$B^*(\psi) = \begin{cases} 0 & \psi(x) \leq 0 \ \forall x \in \Omega \\ +\infty & \exists x \in \Omega \text{ s.t. } \psi(x) > 0 \end{cases}. \quad (6.10)$$

Proof. By Definition 5.2.1 of convex conjugate, for any $\psi \in \mathcal{C}_0(\Omega, \mathbb{R})$ we write

$$\begin{aligned} B^*(\psi) &= \sup_{\mu \in \mathcal{M}(\Omega)} \langle \psi, \mu \rangle_{\mathcal{C}_0(\Omega, \mathbb{R}) \times \mathcal{M}(\Omega)} - B(\mu) \\ &= \sup_{\mu \in \mathcal{M}^+(\Omega)} \langle \psi, \mu \rangle_{\mathcal{C}_0(\Omega, \mathbb{R}) \times \mathcal{M}(\Omega)} \\ &\geq \langle \psi, \mu \rangle_{\mathcal{C}_0(\Omega, \mathbb{R}) \times \mathcal{M}(\Omega)} \quad \forall \mu \in \mathcal{M}^+(\Omega). \end{aligned}$$

If there exists $\bar{x} \in \Omega$ such that $\psi(\bar{x}) > 0$, consider $\bar{\mu} = \alpha \delta_{\bar{x}}$ with $\alpha > 0$. Then,

$$B^*(\psi) \geq \langle \psi, \bar{\mu} \rangle = \alpha \psi(\bar{x}) \xrightarrow{\alpha \rightarrow +\infty} +\infty \Rightarrow B^*(\psi) = +\infty.$$

On the other hand, if $\psi(x) \leq 0 \ \forall x \in \Omega$, then $\langle \psi, \mu \rangle = \int_{\Omega} \psi d\mu \leq 0 \ \forall \mu \in \mathcal{M}^+(\Omega)$. Moreover, $\langle \psi, 0 \rangle = 0$. Thus, $B^*(\psi) = 0$ if $\psi(x) \leq 0 \ \forall x \in \Omega$. \square

We can now combine the two results.

Lemma 6.2.3. *Let $F : \mathcal{M}(\Omega) \rightarrow \mathbb{R} \cup \{+\infty\}$ be the function defined in (6.8). Its convex conjugate $F^* : \mathcal{C}_0(\Omega, \mathbb{R}) \rightarrow \mathbb{R} \cup \{+\infty\}$ for all $\psi \in \mathcal{C}_0(\Omega, \mathbb{R})$ is defined by*

$$F^*(\psi) = \begin{cases} 0 & \text{if } \psi(x) \leq 1 \ \forall x \in \Omega \\ +\infty & \text{otherwise} \end{cases}. \quad (6.11)$$

Proof. Recall that thanks to Proposition 6.2.1, the convex conjugate of F (6.8) is given by the infimal convolution (6.9) of A^* , defined in (5.18), with B^* , given by (6.10). Thus, we can write

$$\begin{aligned} F^*(\psi) &= \min_{\psi_1 + \psi_2 = \psi} A^*(\psi_1) + B^*(\psi_2) \\ &= \min_{\psi_1 + \psi_2 = \psi} \begin{cases} 0 & \|\psi_1\|_{\infty, \Omega} \leq 1 \\ +\infty & \|\psi_1\|_{\infty, \Omega} > 1 \end{cases} + \begin{cases} 0 & \psi_2(x) \leq 0 \ \forall x \in \Omega \\ +\infty & \exists x \in \Omega \text{ s.t. } \psi_2(x) > 0 \end{cases}. \end{aligned}$$

Notice that $F^*(\psi) \geq 0 \ \forall \psi \in \mathcal{C}_0(\Omega, \mathbb{R})$ and we have either $F^*(\psi) = 0$ or $F^*(\psi) = +\infty$. Thus, $F^*(\psi) = 0$ if and only if there exist $\psi_1, \psi_2 \in \mathcal{C}_0(\Omega, \mathbb{R})$ such that $\psi = \psi_1 + \psi_2$ with $A^*(\psi_1) = 0$ and $B^*(\psi_2) = 0$. On the contrary, $F^*(\psi) = +\infty$ if for all $\psi_1, \psi_2 \in \mathcal{C}_0(\Omega, \mathbb{R})$ such that $\psi = \psi_1 + \psi_2$ it holds $A^*(\psi_1) = +\infty$ or $B^*(\psi_2) = +\infty$.

- If $\|\psi\|_{\infty, \Omega} \leq 1$, consider $\psi_1 = \psi$, which implies $A^*(\psi_1) = 0$, and $\psi_2 = 0$, which implies $B^*(\psi_2) = 0$. Thus $F^*(\psi) = 0$.
- If $\|\psi\|_{\infty, \Omega} > 1$ and there exists \bar{x} such that $\psi(\bar{x}) > 1$, we have $F^*(\psi) = +\infty$. Indeed, assume by contradiction there exist $\psi_1, \psi_2 \in \mathcal{C}_0(\Omega, \mathbb{R})$ such that $\psi_1 + \psi_2 = \psi$ and $\|\psi_1\|_{\infty, \Omega} \leq 1$ and $\psi_2(x) \leq 0 \ \forall x \in \Omega$. We would have

$$1 < \psi(\bar{x}) = \underbrace{\psi_1(\bar{x})}_{\leq 1} + \underbrace{\psi_2(\bar{x})}_{\leq 0} \leq 1 \Rightarrow 1 < 1,$$

which is absurd. Then, $\forall \psi_1, \psi_2$ such that $\psi_1 + \psi_2 = \psi$ either $\|\psi_1\|_{\infty, \Omega} > 1$ or $\psi_2(x) > 0$ for some $x \in \Omega$. Hence, $F^*(\psi) = +\infty$.

- If $\|\psi\|_{\infty, \Omega} > 1$ and $\psi(x) \leq 1 \ \forall x \in \Omega$, consider $\psi_1 = \psi^+$ and $\psi_2 = \psi^-$, with

$$\psi^+(x) = \begin{cases} \psi(x) & \psi(x) \geq 0 \\ 0 & \psi(x) < 0 \end{cases} \quad \psi^-(x) = \begin{cases} 0 & \psi(x) > 0 \\ \psi(x) & \psi(x) \leq 0 \end{cases}.$$

It is evident that $\psi = \psi^+ + \psi^-$. Moreover, $\|\psi^+\|_{\infty, \Omega} \leq 1$, thus $A^*(\psi^+) = 0$, and $\forall x \in \Omega \ \psi^-(x) \leq 0$, thus $B^*(\psi^-) = 0$. We have proved that $F^*(\psi) = 0$.

□

6.2.3 Dual problem formulation and extremality conditions

The study of the dual problem with non-negativity constraints has been carried out in [169, 219] in the discrete setting of LASSO, with a Gaussian L^2 fidelity, and in [69] for the discrete counterpart of $(\tilde{\mathcal{D}}_{KL} - |\cdot|)$, with the Kullback-Leibler divergence as fidelity and the L^1 penalty. The analysis of $(\tilde{\mathcal{D}}_{KL} - |\cdot|)$, and of its optimality conditions and dual problem, that we present in this thesis, represents a novelty being in the continuous off-the-grid setting of $\mathcal{M}(\Omega)$ and having the Kullback-Leibler as Poisson noise fidelity term.

Following Lemma 5.2.1, we can obtain the dual problem corresponding to the primal problem $(\tilde{\mathcal{D}}_{KL} - |\cdot|)$ by plugging (6.7) and (6.11) into (5.15):

$$\begin{aligned}
 & \operatorname{argmax}_{p^* \in L^2(\Omega)} -F^*(\Phi^* p^*) - G^*(-p^*) \\
 &= \operatorname{argmax}_{p^* \in L^2(\Omega)} -F^*(\Phi^* p^*) + \begin{cases} -\infty & -p^* \geq \frac{1}{\lambda} \\ -\frac{b}{\lambda}(1 + \lambda p^*) + \frac{y}{\lambda} \log(1 + \lambda p^*) & -p^* < \frac{1}{\lambda} \end{cases} \\
 &= \operatorname{argmax}_{p^* \in L^2(\Omega)} -F^*(\Phi^* p^*) + \begin{cases} -\infty & p^* \leq -\frac{1}{\lambda} \\ -\frac{b}{\lambda}(1 + \lambda p^*) + \frac{y}{\lambda} \log(1 + \lambda p^*) & p^* > -\frac{1}{\lambda} \end{cases} \\
 &= \operatorname{argmax}_{p^* \in L^2(\Omega) \text{ s.t. } p^* > -\frac{1}{\lambda}} -F^*(\Phi^* p^*) - \frac{b}{\lambda}(1 + \lambda p^*) + \frac{y}{\lambda} \log(1 + \lambda p^*) \\
 &= \operatorname{argmax}_{p^* \in L^2(\Omega) \text{ s.t. } p^* > -\frac{1}{\lambda}} \begin{cases} 0 & \forall x \in \Omega \ \Phi^* p^*(x) \leq 1 \\ -\infty & \exists x \in \Omega \ \Phi^* p^*(x) > 1 \end{cases} - \frac{b}{\lambda}(1 + \lambda p^*) + \frac{y}{\lambda} \log(1 + \lambda p^*) \\
 &= \operatorname{argmax}_{p^* \in \mathcal{D}} -\frac{b}{\lambda}(1 + \lambda p^*) + \frac{y}{\lambda} \log(1 + \lambda p^*), \tag{6.12}
 \end{aligned}$$

where $\mathcal{D} = \{p^* \in L^2(\Omega) : p^* > -\frac{1}{\lambda} \text{ and } \forall x \in \Omega, \ \Phi^* p^*(x) \leq 1\}$.

Moreover, from Lemma 5.2.1, the extremality conditions can be obtained. Given $\mu_\lambda \in \mathcal{M}(\Omega)$ solution of the primal problem $(\tilde{\mathcal{D}}_{KL} - |\cdot|)$ with regularisation parameter $\lambda > 0$ and $p_\lambda^* \in L^2(\Omega)$ solution of the dual problem (6.12), the extremality conditions (5.16) in this case read as

$$\begin{cases} \Phi^* p_\lambda^* \in \partial F(\mu_\lambda) = \partial|\mu_\lambda|(\mathcal{X}) + \partial\mathbf{1}_{\mathcal{M}^+(\Omega)}(\mu_\lambda) \\ -p_\lambda^* \in \frac{1}{\lambda} \partial_1 \tilde{\mathcal{D}}_{KL}(\Phi \mu_\lambda + b, y) = \frac{1}{\lambda} \left(I - \frac{y}{\Phi \mu_\lambda + b} \right), \end{cases} \tag{6.13}$$

where $\partial_1 \tilde{\mathcal{D}}_{KL}(\Phi \mu_\lambda, y)$ denotes the subdifferential of $\tilde{\mathcal{D}}_{KL}(\cdot, \cdot)$ computed with respect to the first variable and evaluated in $(\Phi \mu_\lambda, y)$.

Remark 6.2.1. *If $\mu_\lambda \in \mathcal{M}(\Omega)$ is solution of the primal problem $(\tilde{\mathcal{D}}_{KL} - |\cdot|)$ and $p_\lambda^* \in L^2(\Omega)$ is solution of the dual (6.12), then from the extremality conditions (6.13) we have*

$$-p_\lambda^* = \frac{1}{\lambda} \left(I - \frac{y}{\Phi \mu_\lambda + b} \right) \Rightarrow p_\lambda^* = \frac{y - \Phi \mu_\lambda - b}{\lambda(\Phi \mu_\lambda + b)}.$$

It follows that $p_\lambda^ > -\frac{1}{\lambda} \iff y > 0$, which holds by hypothesis (6.4).*

In the extremality conditions (6.13), the meaning of $\partial|\mu_\lambda|(\Omega) + \partial\mathbf{1}_{\mathcal{M}^+(\Omega)}(\mu_\lambda)$ needs further explanation. First of all, observe that, in general, for two proper convex functions $f_1, f_2 : \mathcal{X} \rightarrow \mathbb{R} \cup \{+\infty\}$ it holds

$$\partial(f_1 + f_2) \subseteq \partial f_1 + \partial f_2,$$

and the equality holds if and only if $\operatorname{int}(\operatorname{dom}(f_1)) \cap \operatorname{int}(\operatorname{dom}(f_2)) \neq \emptyset$. In particular, we have $\partial F(\cdot) = \partial|\cdot|(\Omega) + \partial\mathbf{1}_{\mathcal{M}^+(\Omega)}(\cdot)$. Recall that the subdifferential of the TV

norm is given by (5.3), and it assumes the form (5.4) for finite linear combinations of Diracs. We need to study the subdifferential of the indicator function $\partial \mathbb{1}_{\mathcal{M}^+(\Omega)}(\cdot)$ of $\mathcal{M}^+(\Omega)$, and then analyse the sum of the two subdifferentials.

6.2.3.1 Subdifferential of the indicator function of positive measures

Recall that the indicator function of positive measures is defined by

$$B(\mu) = \mathbb{1}_{\{\mathcal{M}^+(\Omega)\}}(\mu) = \begin{cases} 0 & \mu \geq 0 \\ +\infty & \text{otherwise} \end{cases}.$$

Let $\mu \in \mathcal{M}^+(\Omega)$. By definition, the subdifferential of B at $\mu \in \mathcal{M}^+(\Omega)$ is

$$\begin{aligned} \partial B(\mu) &= \{\eta \in \mathcal{C}_0(\Omega, \mathbb{R}) \mid B(\tilde{\mu}) \geq B(\mu) + \langle \eta, \tilde{\mu} - \mu \rangle \forall \tilde{\mu} \in \mathcal{M}(\Omega)\} \\ &= \{\eta \in \mathcal{C}_0(\Omega, \mathbb{R}) \mid \langle \eta, \tilde{\mu} - \mu \rangle \leq 0 \forall \tilde{\mu} \in \mathcal{M}^+(\Omega)\}, \end{aligned} \quad (6.14)$$

that is the normal cone to $\mathcal{M}^+(\Omega)$ at μ .

Now, we consider (6.14) for some scenarios of interest.

- If $\mu = 0$, we have $\eta \in \partial B(0) \iff \langle \eta, \tilde{\mu} \rangle \leq 0 \forall \tilde{\mu} \in \mathcal{M}^+(\Omega)$. Since $\tilde{\mu}$ is a positive measure, we have that $\langle \eta, \tilde{\mu} \rangle = \int_{\Omega} \eta d\tilde{\mu} \leq 0$ for any η such that $\eta(x) \leq 0$ for all $x \in \Omega$. Thus,

$$\partial B(0) = \{\eta \in \mathcal{C}_0(\Omega, \mathbb{R}) \mid \forall x \in \Omega \quad \eta(x) \leq 0\}.$$

- If $\mu = \bar{a}\delta_{\bar{x}}$ with $\bar{a} > 0$, $\bar{x} \in \Omega$, we have that $\eta \in \partial B(\mu)$ if and only if $\forall \tilde{\mu} \in \mathcal{M}^+(\Omega)$ it holds

$$\langle \eta, \tilde{\mu} - \mu \rangle \leq 0 \iff \int_{\Omega} \eta d\tilde{\mu} = \langle \eta, \tilde{\mu} \rangle \leq \langle \eta, \mu \rangle = \int_{\Omega} \eta d\mu = \bar{a}\eta(\bar{x})$$

Consider $\tilde{\mu} = 2\mu$. Hence, $\int_{\Omega} \eta d\tilde{\mu} = \int_{\Omega} \eta d2\mu = 2\bar{a}\eta(\bar{x})$ and $2\bar{a}\eta(\bar{x}) \leq \bar{a}\eta(\bar{x})$ if and only if $\eta(\bar{x}) = 0$. Thus, $\eta \in \partial B(\mu)$ implies $\eta(\bar{x}) = 0$.

Let $\eta \in \partial B(\mu)$, which entails $\eta(\bar{x}) = 0$, and suppose $\exists \tilde{x} \in \Omega$ s.t. $\eta(\tilde{x}) > 0$. Then,

$$\forall \tilde{\mu} \in \mathcal{M}^+(\Omega) \quad \int_{\Omega} \eta d\tilde{\mu} \leq \int_{\Omega} \eta d\mu = \bar{a}\eta(\bar{x}) = 0.$$

Consider the above inequality with $\tilde{\mu} = \tilde{a}\delta_{\tilde{x}}$, $\tilde{a} > 0$:

$$0 \geq \int_{\Omega} \eta d\tilde{\mu} = \tilde{a}\eta(\tilde{x}) > 0,$$

which is absurd. Thus, $\eta(x) \leq 0 \forall x \in \Omega$ and

$$\partial B(\mu) = \{\eta \in \mathcal{C}_0(\Omega, \mathbb{R}) \mid \forall x \in \Omega \quad \eta(x) \leq 0 \text{ and } \eta(\bar{x}) = 0\}$$

- Similarly, if $\mu = \sum_{i=1}^N a_i \delta_{x_i}$ with $a_i > 0$ and $x_i \in \Omega$ for $i = 1, \dots, N$, we have:

$$\partial B(\mu) = \{\eta \in \mathcal{C}_0(\Omega, \mathbb{R}) \mid \forall x \in \Omega \quad \eta(x) \leq 0 \text{ and } \eta(x_i) = 0, \quad i = 1, \dots, N\}.$$

6.2.3.2 Extremality conditions

We now have an analytical expression for both $\partial|\cdot|(\Omega)$ and $\partial\mathbf{1}_{\mathcal{M}^+(\Omega)}(\cdot)$, at least for finite linear combinations of Diracs. Thus, we can now state the following proposition, which helps to understand the extremality conditions (6.13).

Proposition 6.2.2. *If $\mu = \sum_{i=1}^N a_i \delta_{x_i}$ with $a_i > 0$, $x_i \in \Omega$ we have*

$$\partial|\mu|(\Omega) + \partial\mathbf{1}_{\mathcal{M}^+(\Omega)}(\mu) = \{\eta \in \mathcal{C}_0(\Omega, \mathbb{R}) : \forall x \in \Omega \quad \eta(x) \leq 1 \text{ and } \eta(x_i) = 1, i = 1, \dots, N\}. \quad (6.15)$$

Proof. If $\mu = \sum_{i=1}^N a_i \delta_{x_i}$ with $a_i > 0$, $x_i \in \Omega$ we have

$$\partial|\mu|(\Omega) = \{\eta_1 \in \mathcal{C}_0(\Omega, \mathbb{R}) : \|\eta_1\|_\infty \leq 1 \text{ and } \eta_1(x_i) = 1, i = 1, \dots, N\}$$

$$\partial\mathbf{1}_{\mathcal{M}^+(\Omega)}(\mu) = \{\eta_2 \in \mathcal{C}_0(\Omega, \mathbb{R}) : \forall x \in \Omega \quad \eta_2(x) \leq 0 \text{ and } \eta_2(x_i) = 0, i = 1, \dots, N\}$$

Thus,

$$\partial|\mu|(\Omega) + \partial\mathbf{1}_{\mathcal{M}^+(\Omega)}(\mu) \subseteq \{\eta \in \mathcal{C}_0(\Omega, \mathbb{R}) : \forall x \in \Omega \quad \eta(x) \leq 1 \text{ and } \eta(x_i) = 1, i = 1, \dots, N\}.$$

We consider now the opposite inclusion. Let $\eta \in \mathcal{C}_0(\Omega, \mathbb{R})$ such that for all $x \in \Omega$ $\eta(x) \leq 1$ and $\eta(x_i) = 1$, $i = 1, \dots, N$. Our aim is to prove that there exist $\eta_1 \in \partial|\mu|(\Omega)$ and $\eta_2 \in \partial\mathbf{1}_{\mathcal{M}^+(\Omega)}(\mu)$ such that $\eta = \eta_1 + \eta_2$.

- If $\|\eta\|_\infty \leq 1$, one can take $\eta_1 = \eta$, $\eta_2 = 0$.
- If $\|\eta\|_\infty > 1$, consider

$$\eta_1(x) = \begin{cases} \eta(x) & -1 \leq \eta(x) \leq 1 \\ -1 & \eta(x) \leq -1 \end{cases} \quad \eta_2(x) = \begin{cases} 0 & -1 \leq \eta(x) \leq 1 \\ \eta(x) + 1 & \eta(x) \leq -1 \end{cases}.$$

In both cases, we have found $\eta_1 \in \partial|\mu|(\Omega)$ and $\eta_2 \in \partial\mathbf{1}_{\mathcal{M}^+(\Omega)}(\mu)$ such that $\eta = \eta_1 + \eta_2$. This concludes the proof. \square

Thanks to (6.15) of the latter proposition, it is now possible to better characterise the extremality conditions (6.13), under the assumption that μ_λ , solution of $(\tilde{\mathcal{D}}_{KL} - |\cdot|)$, is a discrete measure. Indeed, if $\mu_\lambda = \sum_{i=1}^{N_\lambda} (a_\lambda)_i \delta_{(x_\lambda)_i}$ we have that

$$\forall x \in \Omega \quad \Phi^* p_\lambda^*(x) \leq 1 \quad \text{and} \quad \Phi^* p_\lambda^*((x_\lambda)_i) = 1, i = 1, \dots, N_\lambda, \quad (6.16)$$

with p_λ^* solution of the dual problem (6.12).

6.2.4 Optimality conditions

It is possible to obtain a result similar to (6.16), analysing the optimality conditions of $(\tilde{\mathcal{D}}_{KL} - |\cdot|)$. In $(\tilde{\mathcal{D}}_{KL} - |\cdot|)$, the functional

$$T_\lambda^{KL}(\mu) := \tilde{\mathcal{D}}_{KL}(\Phi\mu + b, y) + \lambda|\mu|(\Omega) + \mathbf{1}_{\{\mathcal{M}^+(\Omega)\}}(\mu)$$

is convex. Hence μ minimises T_λ^{KL} if and only if $0 \in \partial T_\lambda^{KL}(\mu)$, i.e.

$$\begin{aligned} 0 &\in \Phi^* \partial_1 \tilde{\mathcal{D}}_{KL}(\Phi\mu, y) + \lambda \partial|\mu|(\Omega) + \partial \mathbf{1}_{\mathcal{M}^+(\Omega)}(\mu) \\ 0 &\in \Phi^* \left(I - \frac{y}{\Phi\mu + b} \right) + \lambda \partial|\mu|(\Omega) + \partial \mathbf{1}_{\mathcal{M}^+(\Omega)}(\mu) \\ \Phi^* \left(\frac{y}{\Phi\mu + b} - I \right) &\in \lambda \partial|\mu|(\Omega) + \partial \mathbf{1}_{\mathcal{M}^+(\Omega)}(\mu) \\ \frac{1}{\lambda} \Phi^* \left(\frac{y}{\Phi\mu + b} - I \right) &\in \partial|\mu|(\Omega) + \frac{1}{\lambda} \partial \mathbf{1}_{\mathcal{M}^+(\Omega)}(\mu) \end{aligned}$$

which can be written as

$$\eta \in \partial|\mu|(\Omega) + \frac{1}{\lambda} \partial \mathbf{1}_{\mathcal{M}^+(\Omega)}(\mu), \quad (6.17)$$

by defining the dual certificate

$$\eta := \frac{1}{\lambda} \Phi^* \left(\frac{y}{\Phi\mu + b} - I \right). \quad (6.18)$$

If μ is a finite linear combination of Dirac masses, Proposition 6.2.2 yields an expression for the sum of subdifferentials in (6.15). Hence, (6.17) becomes

$$\forall x \in \Omega \quad \eta(x) \leq 1 \quad \text{and} \quad \eta(x_i) = 1, i = 1, \dots, N. \quad (6.19)$$

Note the similarity between (6.19) and (6.16), and of the dual certificate (6.18) and the quantity $\Phi^* p_\lambda^*$ in the extremality conditions (6.13).

6.3 Algorithmic choices

For the minimisation of $(\tilde{\mathcal{D}}_{KL} - |\cdot|)$, we propose to consider the Frank-Wolfe algorithm (Algorithm 9), endowed with a similar epigraphic lift as for BLASSO (see Section 5.3.2).

Lemma 6.3.1. *The problem $(\tilde{\mathcal{D}}_{KL} - |\cdot|)$ is equivalent to*

$$\operatorname{argmin}_{(t, \mu) \in C} \tilde{T}_\lambda^{KL}(\mu, t) := \tilde{\mathcal{D}}_{KL}(\Phi\mu + b, y) + \lambda t + \mathbf{1}_{\{\mathcal{M}^+(\Omega)\}}(\mu), \quad (6.20)$$

where

$$C := \left\{ (t, \mu) \in \mathbb{R}^+ \times \mathcal{M}(\Omega); |\mu|(\Omega) \leq t \leq M \right\} \quad \text{with} \quad M := \frac{\tilde{\mathcal{D}}_{KL}(b, y)}{2\lambda}.$$

Some small changes with respect to the Gaussian case have to be considered. Indeed, step (5.31) of Algorithm 9 becomes

$$x_*^k = \operatorname{argmax}_{x \in \Omega} \eta^k(x), \quad \eta^k = \frac{1}{\lambda} \Phi^* \left(\frac{y}{\Phi \mu^k + b} - I \right). \quad (6.21)$$

This results in a change of the stopping criterion (5.24), that reads

$$\eta^k(x) \leq 1 \text{ for any } x \in \Omega,$$

in accordance with the optimality conditions (6.19).

Moreover, the line-search step (5.26) of Algorithm 9 applied to the minimisation of T_λ^{KL} (6.20) becomes

$$a^{k+1} = \operatorname{argmin}_{a \in \mathbb{R}^{k+1}} \tilde{\mathcal{D}}_{KL}(\Phi \mu_a + b, y) + \lambda \|a\|_1 + \mathbb{1}_{\geq 0}(a) \quad (6.22)$$

where $\mu_a = \sum_{i=1}^k a_i \delta_{x_i^k} + a_{k+1} \delta_{x_*^k}$, and x_*^k is the support, defined by (6.21), of the new spike to be added to $\mu^k = \sum_{i=1}^k a_i^k \delta_{x_i^k}$.

Similarly, the Sliding Frank-Wolfe algorithm can be adapted to the resolution of $(\tilde{\mathcal{D}}_{KL} - |\cdot|)$ with the same modifications. The sliding step (5.38) of Algorithm 10 now reads as

$$(a^{k+1}, x^{k+1}) \in \operatorname{arg min}_{(a,x) \in \mathbb{R}^{N_{k+1}} \times \Omega^{N_{k+1}}} \tilde{\mathcal{D}}_{KL}(\Phi_x a + b, y) + \lambda \|a\|_1 + \mathbb{1}_{\geq 0}(a). \quad (6.23)$$

6.3.1 Boosted Sliding Frank Wolfe algorithm

Sliding Frank Wolfe computational complexity increases significantly throughout the iterations as more spikes are added. This affects especially the non-convex sliding step on both amplitudes and positions, since it is a minimisation problem over $\mathbb{R}^{(d+1)k}$ at iteration k , with $d \in \mathbb{N}$ being the dimension of Ω . In [67], it is observed that most computational time is spent in the sliding step to finely optimise all the parameters, which are again updated and modified afterwards in the next iteration, for all iterations except the last one. Based on this, they propose an accelerated version of SFW, called Boosted SFW, which removes most of the costly non-convex minimisation steps. Boosted SFW core structure is the same as for SFW, but the following changes are made to accelerate the algorithm preserving its convergence properties.

1. The certificate η^k is computed to obtain the position x_*^k of a potential new spike.
2. If the stopping condition - $\|\eta^k\|_\infty \leq 1$ for BLASSO $(L^2 - |\cdot|)$, and $\|(\eta^k)_+\|_\infty \leq 1$ for the Poisson model $(\tilde{\mathcal{D}}_{KL} - |\cdot|)$ - of SFW is not met, then the insertion step - (5.37) for $(L^2 - |\cdot|)$, and (6.22) for $(\tilde{\mathcal{D}}_{KL} - |\cdot|)$ - is performed and the new amplitudes are estimated. Then, the algorithm proceeds again with (1) and a new iteration of BSFW commences. Note that, in this case, no sliding step is performed.

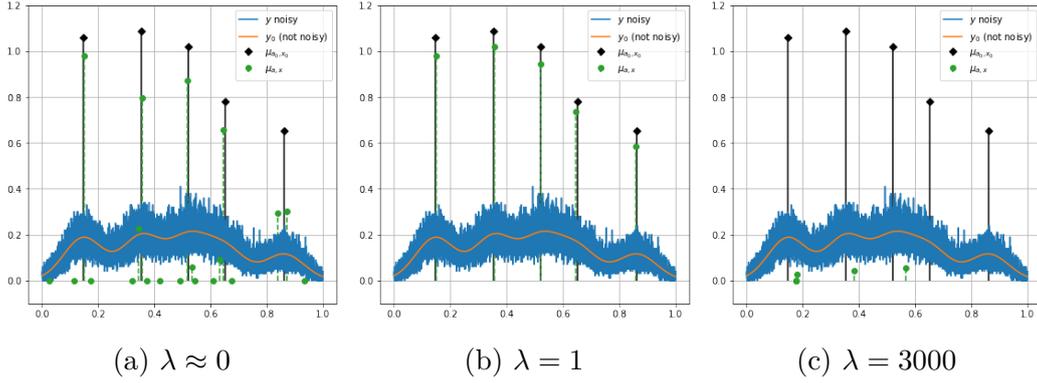


Figure 6.1: Reconstructions obtained using SFW in a 1D sparse deconvolution numerical example with Poisson noise. In black, the ground truth spikes and in green the reconstructed ones are shown. With λ close to 0, the number and intensities of spikes are overestimated, while with a much higher value they are underestimated.

3. Otherwise the sliding step - (5.38) for $(L^2 - |\cdot|)$, and (6.23) for $(\tilde{\mathcal{D}}_{KL} - |\cdot|)$ - finely re-estimates amplitudes and positions. Then, the new dual certificate η^{k+1} is computed to check if the stopping condition is still met, after the performed sliding step. Namely,

- If $\|\eta^{k+1}\|_\infty \leq 1$ for $(L^2 - |\cdot|)$ and $\|(\eta^{k+1})_+\|_\infty \leq 1$ for $(\tilde{\mathcal{D}}_{KL} - |\cdot|)$, the algorithm stops.
- Otherwise the algorithm continues again from (1).

Thus, BSFW performs the sliding step when needed and not systematically, thanks to the insight given by the dual certificate. In the worst case scenario, the sliding step is computed at each iteration of BSFW retrieving exactly the SFW algorithm. In the best case scenario, instead, the sliding step is performed only once at the very last iteration. In [67], the finite termination of BSFW is proved.

6.4 Homotopy algorithm

Note that in Frank-Wolfe and Sliding Frank-Wolfe algorithms, as well as in the just described Boosted SFW, the choice of the regularisation parameters is crucial. Indeed, it plays a fundamental role in the characterisation of the stopping criterion, defined as

$$\operatorname{argmax}_{x \in \Omega} |\eta^k(x)| \leq 1, \quad \eta^k = \frac{1}{\lambda} \Phi^*(y - \Phi \mu^k)$$

for BLASSO $(L^2 - |\cdot|)$, and as

$$\operatorname{argmax}_{x \in \Omega} \eta^k(x) \leq 1, \quad \eta^k = \frac{1}{\lambda} \Phi^*\left(\frac{y}{\Phi \mu^k + b} - I\right)$$

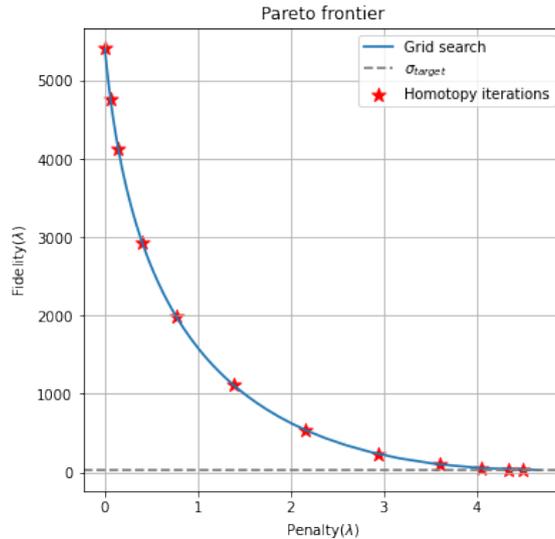


Figure 6.2: Pareto frontier

for the proposed Poisson model $(\tilde{\mathcal{D}}_{KL} - |\cdot|)$. The choice of high value for λ results in only few iterations of the algorithms, before the dual certificate η^k satisfies the above stopping rules. On the other hand, small values of λ yields to more algorithmic iterations before convergence. Moreover, the regularisation parameter has an impact on the estimation step of FW and SFW and on the sliding step of SFW, being associated with the L^1 penalty of the amplitudes vector. Altogether, with high values of λ the reconstruction usually presents fewer spikes with smaller intensities values, whilst small λ s provide a better data fit, with more spikes and higher intensities, see Figure 6.1.

In [67], the authors propose a method to search for the best regularisation parameter, for the resolution of BLASSO $(L^2 - |\cdot|)$ exploiting the idea of homotopy. Homotopy algorithms have been first proposed in [178, 179] to solve LASSO problem, the discrete counterpart of BLASSO, for some known (or estimated) noise value $\delta > 0$ for a sequence of decreasing regularisation parameters λ , appropriately chosen, until a stopping rule defined in terms of δ is met.

We provide more detail by introducing the concept of Pareto frontier, which helps to better visualise the meaning of the homotopy strategy.

Definition 6.4.1. *Let \mathcal{X} be a Banach space and let $f, g : \mathcal{X} \rightarrow \mathbb{R} \cup \{+\infty\}$ such that f and g are proper, lower semi-continuous and convex functions. Consider the minimisation problem*

$$x_\lambda \in \operatorname{argmin}_{x \in \mathcal{X}} f(x) + \lambda g(x), \quad \lambda > 0. \quad (6.24)$$

The set defined by

$$\left\{ (g(x_\lambda), f(x_\lambda)) \in \mathbb{R} \times \mathbb{R} \mid \lambda > 0 \right\}$$

is called *Pareto frontier* or *Pareto front*.

In the context of variational approaches for inverse problems (1.12), the Pareto Frontier is the plot of the values of the fidelity against the ones of the penalty for each solution (6.24) obtained by varying λ , as shown in blue in Figure 6.2. Basically, it is equivalent to the concept of L-curve, studied for instance in [46] to define a way to select the regularisation parameter, that is the plot of the values of the penalty against the ones of the fidelity, i.e. it is the L-curve can be obtained from the Pareto frontier with symmetry with respect to the bisector of the axes.

A possible way to select an optimal regularisation parameter λ is to perform a fine discretisation of the Pareto front and then select as best λ the one minimising the reconstruction error $\|x_\lambda - x^\dagger\|$, or an estimate of this quantity. This approach is time consuming since it requires minimising the variational problem for an elevated number of parameters λ . Moreover, this strategies fails whenever a good estimate of $\|x_\lambda - x^\dagger\|$ is hard to find.

Instead of exploring the Pareto frontier by grid search, the idea behind homotopy algorithms is to explore the Pareto frontier iteratively [67, 179], only for a few values of $\lambda > 0$. Starting from an initial and overestimated value $\lambda_1 > 0$ used at the first iteration, the solution x_{λ_1} of the corresponding variational problem is computed. Then, at each homotopy iteration, λ is updated if the solution does not fit well the data up to some tolerance $\sigma_{\text{target}}(\delta)$, which depends on the noise value δ . A new solution x_{λ_2} for the new value λ_2 of the regularisation parameter has thus to be computed in the next homotopy step. Homotopy algorithms thus explore the Pareto frontier for a small set of λ and select the biggest value for which the solution of the variational formulation, computed with respect to that value, meets a convergence criterion typically defined in terms of the noise value δ . This process is sketched in Figure 6.2, where one can see in red the discrete discontinuous jumps produced by the homotopy strategy we are going to describe, which stops when the fidelity term goes under the value of $\sigma_{\text{target}}(\delta)$, plotted in grey.

Each homotopy iteration $t \in \mathbb{N}$ performs the following steps:

- Compute $x_{\lambda_t}^t$, solution of (6.24) with $\lambda = \lambda_t$.
- Check if $x_{\lambda_t}^t$ is a good data fit, i.e. if

$$f(x_{\lambda_t}^t) \leq \sigma_{\text{target}}(\delta), \quad (6.25)$$

where the target value for the fidelity $\sigma_{\text{target}}(\delta)$ depends on the noise level $\delta > 0$.

- If not, decrease λ

$$\lambda_{t+1} = h(\lambda_t) < \lambda_t$$

with $h : \mathbb{R} \rightarrow \mathbb{R}$ pre-assigned decreasing function, until (6.25) is satisfied.

Observe that when computing $x_{\lambda_{t+1}}^{t+1}$ one already has computed $x_{\lambda_t}^t$, which can be used as warm start to reduce computations. Thus, for the design of a good homotopy

Algorithm 11 Homotopy algorithm for off-the-grid inverse problems

Input: $y \in L^2(\Omega)$, $b \in L^2(\Omega)$, $b \geq 0$, $\Phi \in \mathcal{L}(\mathcal{M}(\Omega), L^2(\Omega))$

Parameters: $\gamma \in (0, 1)$, $c > 0$, $\sigma_{\text{target}} > 0$

Output: estimation $\hat{\mu} \in \mathcal{M}(\Omega)$

Initialisation: $\hat{\mu}_0 \in \mathcal{M}(\Omega)$ and $\lambda_1 = \gamma \|\eta(1, \hat{\mu}_0)\|_\infty$

repeat

1. Compute $\hat{\mu}_t$ solution of $(\mathcal{P}(\lambda))$ with $\lambda = \lambda_t$ with warm start $\mu_t^{[0]} = \hat{\mu}_{t-1}$.
2. Compute σ_t from the residual:

$$\sigma_t = f_{y^\delta, b}(\Phi \hat{\mu}_t) \tag{6.26}$$

3. **if** $\sigma_t < \sigma_{\text{target}}$

$\hat{\mu}_t$ is a solution \Rightarrow **stop**

4. **else if** $\sigma_t \geq \sigma_{\text{target}}$

$$\text{Update } \lambda_{t+1} = \frac{\lambda_t \|\eta(\lambda_t, \hat{\mu}_t)\|_\infty}{c + 1} \tag{6.27}$$

until $\sigma_t < \sigma_{\text{target}}$

algorithm the choice of the starting value for λ , its updating rule for λ , how to take into account for past knowledge and the choice of a suitable noise representative value σ_{target} have to be discussed.

6.4.1 Homotopy algorithm for off-the-grid methods

We follow [67], where an homotopy version of Sliding Frank-Wolfe is proposed for the resolution of BLASSO ($L^2 - |\cdot|$) problem for automatic parameter selection, and extend its definition to, in principle, a general fidelity term in an off-the-grid setting.

Consider the noisy data $y^\delta \in L^2(\Omega)$ and the inverse problem $y^\delta = \Phi\mu + b$ with $b \in L^2(\Omega)$ such that $b \geq 0$. In this section, we allow b to be equal to 0 to describe scenarios where the background is not needed, or where it can be zero with positive measure. Anyway, we remark that for the formulation $(\tilde{\mathcal{D}}_{KL} - |\cdot|)$ a strictly positive background is necessary, as well as $y^\delta > 0$.

Given a convex and smooth fidelity $f_{y^\delta, b} : \mathcal{M}(\Omega) \rightarrow \mathbb{R}$ we define the following

	$(L^2 - \cdot)$	$(\tilde{\mathcal{D}}_{KL} - \cdot)$
λ_1	$\gamma \ \Phi^* y^\delta\ _\infty$	$\gamma \left\ \left(\Phi^* \left(\frac{y-b}{b} \right) \right)_+ \right\ _\infty$
σ_t	$\frac{1}{2} \ \Phi \hat{\mu}_t - y^\delta\ ^2$	$\tilde{\mathcal{D}}_{KL}(\Phi \hat{\mu}_t, y^\delta)$
$\sigma_{\text{target}}(\delta)$	$\frac{1}{2} \ y - y^\delta\ ^2 = \frac{\delta^2}{2}$	$\tilde{\mathcal{D}}_{KL}(y, y^\delta)$

Table 6.1: Homotopy algorithmic choices for BLASSO and for the Poisson off-the-grid models.

variational problem in the space of Radon measures:

$$\operatorname{argmin}_{\mu \in \mathcal{M}(\Omega)} f_{y^\delta, b}(\Phi \mu) + \lambda |\mu|(\Omega) + \alpha \mathbf{1}_{\mathcal{M}^+(\Omega)}(\mu), \quad \lambda > 0, \quad \alpha \in \{0, 1\}. \quad (\mathcal{P}(\lambda))$$

The parameter $\alpha \in \{0, 1\}$ allows to consider, at the same time, problems with and without positivity constraints. We propose to solve $(\mathcal{P}(\lambda))$ by Algorithm 11. The dual certificate of $(\mathcal{P}(\lambda))$ is a function $\eta(\lambda, \mu) \in L^2(\Omega)$ defined by

$$\eta(\lambda, \mu) = \frac{1}{\lambda} \tilde{\eta}(\mu) \quad \text{with} \quad \tilde{\eta}(\mu) := \begin{cases} -\Phi^* \frac{\partial}{\partial \mu} f_{y^\delta, b}(\Phi \mu) & \alpha = 0 \\ \left(-\Phi^* \frac{\partial}{\partial \mu} f_{y^\delta, b}(\Phi \mu) \right)_+ & \alpha = 1 \end{cases}. \quad (6.28)$$

Optimality conditions for $(\mathcal{P}(\lambda))$ read as

$$\|\eta(\lambda, \mu)\|_\infty \leq 1, \quad (6.29)$$

and, under the hypothesis that the solution is a finite linear combination of Diracs, the dual certificate of the solution satisfies $\eta(\lambda, \mu)(x_i) = 1$ where the points $x_i \in \Omega$ are the support of μ .

We will consider Algorithm 11 for the BLASSO problem $(L^2 - |\cdot|)$ and for the Poisson off-the-grid problem $(\tilde{\mathcal{D}}_{KL} - |\cdot|)$. BLASSO can be retrieved by means of $(\mathcal{P}(\lambda))$ by considering $b = 0$, $f_{y^\delta, b}(\Phi \mu) = \frac{1}{2} \|\Phi \mu - y\|^2$ and $\alpha = 0$. Instead, the Poisson off-the-grid problem $(\tilde{\mathcal{D}}_{KL} - |\cdot|)$ is obtained with $f_{y^\delta, b}(\Phi \mu) = \tilde{\mathcal{D}}_{KL}(\Phi \mu + b, y)$ and $\alpha = 1$. In Table 6.1, we outline the different algorithmic choices made applying the homotopy algorithm (Alg.11) to $(L^2 - |\cdot|)$ and $(\tilde{\mathcal{D}}_{KL} - |\cdot|)$.

Observe that $\eta(\lambda, \mu)$ is crucial in the definition of Algorithm 11 as it allows to define a good starting value for λ and the updating rule.

6.4.1.1 Starting value

We propose to initialise the homotopy strategy with the initial value

$$\lambda_1 := \gamma \|\eta(1, \hat{\mu}_0)\|_\infty = \gamma \|\tilde{\eta}(\hat{\mu}_0)\|_\infty, \quad \gamma \in (0, 1) \quad (6.30)$$

where $\hat{\mu}_0 \in \mathcal{M}(\Omega)$ is the initialisation of the solution and η is defined by (6.28). This choice is motivated by the optimality conditions (6.29). Indeed, if one takes at

the first iteration $\lambda_1 \geq \|\eta(1, \hat{\mu}_0)\|_\infty$ then $\hat{\mu}_0$ is an optimal solution for $(\mathcal{P}(\lambda))$ with $\lambda = \lambda_1$:

$$\|\eta(\lambda_1, \hat{\mu}_0)\|_\infty = \left\| \frac{1}{\lambda_1} \tilde{\eta}(\hat{\mu}_0) \right\|_\infty = \frac{1}{\lambda_1} \|\eta(1, \hat{\mu}_0)\|_\infty \leq 1 \iff \lambda_1 \geq \|\eta(1, \hat{\mu}_0)\|_\infty.$$

In this case the algorithm does not improve upon the initialisation $\hat{\mu}_0$ since it does not perform any iteration. On the contrary, choosing $\lambda_1 < \|\eta(1, \hat{\mu}_0)\|_\infty$ ensures that the initial measure $\hat{\mu}_0$ will be updated since it is not optimal for $(\mathcal{P}(\lambda))$ with the value $\lambda = \lambda_1$. Indeed, at the first iteration of the homotopy algorithm, the dual certificate computed with respect to the initialisation $\hat{\mu}_0$ and $\lambda_1 > 0$ given by (6.30) has supremum norm that satisfies

$$\|\eta(\lambda_1, \hat{\mu}_0)\|_\infty = \frac{1}{\lambda_1} \|\eta(1, \hat{\mu}_0)\|_\infty = \frac{1}{\gamma \|\eta(1, \hat{\mu}_0)\|_\infty} \|\eta(1, \hat{\mu}_0)\|_\infty = \frac{1}{\gamma} > 1,$$

since $\gamma \in (0, 1)$. Note that a similar starting point for λ is used in [158, 179] in the discretised case of LASSO and, here, generalised for a generic fidelity f and to the off-the-grid setting.

6.4.1.2 Updating rule

It is now worth discussing the updating rule (6.27). Indeed, (6.27) together with the choice of a strictly positive parameter $c > 0$ ensures that the measure $\hat{\mu}_t$, by which the algorithm initialises $(\mathcal{P}_{\lambda_{t+1}})$ as $\mu_{t+1}^{[0]} = \hat{\mu}_t$, is not already an optimal solution for the problem. Indeed, the dual certificate with respect to λ_{t+1} and the initialisation $\hat{\mu}_t$ reads

$$\eta(\lambda_{t+1}, \hat{\mu}_t) = \frac{\lambda_t}{\lambda_{t+1}} \eta(\lambda_t, \hat{\mu}_t) = \frac{1+c}{\|\eta(\lambda_t, \hat{\mu}_t)\|_\infty} \eta(\lambda_t, \hat{\mu}_t), \Rightarrow \|\eta(\lambda_{t+1}, \hat{\mu}_t)\|_\infty = 1+c > 1.$$

Thus, $\mu_{t+1}^{[0]} = \hat{\mu}_t$ does not satisfy the optimality condition (6.29) for $(\mathcal{P}(\lambda))$ with $\lambda = \lambda_{t+1}$ and hence the homotopy step $t+1$ computes a new $\hat{\mu}_{t+1}$. This is consistent with the choice of rejecting $\hat{\mu}_t$ at the previous step t .

6.4.1.3 Descent property of the homotopy algorithm

In this section, we discuss the descent properties of the proposed homotopy algorithm, showing that it produces a strictly decreasing sequence of residual errors $(\sigma_t)_t$. This properties gives insight on the good convergence of the algorithm.

Proposition 6.4.1. *The homotopy algorithm (Alg.11) for the resolution of $(\mathcal{P}(\lambda))$ produces a strictly decreasing sequence $(\sigma_t)_t$ provided that $f_{y^\delta, b}$ is twice differentiable with $\partial^2 f_{y^\delta, b} > 0$.*

Proof. Let us denote

$$\mu^*(\lambda) \in \underset{\mu}{\operatorname{argmin}} C(y, \lambda, \mu) = f_{y^\delta, b}(\Phi\mu) + \lambda|\mu|(\Omega),$$

where the argmin is computed over $\mathcal{M}(\Omega)$ if $\alpha = 0$ in $(\mathcal{P}(\lambda))$ and over $\mathcal{M}^+(\Omega)$ if $\alpha = 1$. Being $\mu^*(\lambda)$ the minimiser, then we have

$$\left. \frac{\partial C(y, \lambda, \mu)}{\partial \mu} \right|_{\mu=\mu^*} = \Phi^* \partial f_{y^\delta, b}(\Phi \mu^*) + \lambda = 0, \quad (6.31)$$

and thus deriving again with respect to λ yields

$$\begin{aligned} 0 &= \frac{\partial}{\partial \lambda} \left(\left. \frac{\partial C(y, \lambda, \mu)}{\partial \mu} \right|_{\mu=\mu^*} \right) = \frac{\partial}{\partial \lambda} \left(\Phi^* \partial f_{y^\delta, b}(\Phi \mu^*) + \lambda \right) \\ 0 &= \frac{\partial}{\partial \mu} \left(\Phi^* \partial f_{y^\delta, b}(\Phi \mu) \right) \Big|_{\mu=\mu^*} \cdot \frac{\partial \mu^*(\lambda)}{\partial \lambda} + 1 \\ 0 &= \Phi^* \Phi \cdot \partial^2 f_{y^\delta, b}(\Phi \mu^*) \cdot \frac{\partial \mu^*(\lambda)}{\partial \lambda} + 1. \end{aligned}$$

Since Φ is positive definite and by hypothesis on $f_{y^\delta, b}$, from the computations above it follows that $\frac{\partial \mu^*(\lambda)}{\partial \lambda} < 0$. Consider now the following:

$$\frac{\partial}{\partial \lambda} f_{y^\delta, b}(\Phi \mu^*(\lambda)) = \Phi^* \partial f_{y^\delta, b}(\Phi \mu^*) \cdot \frac{\partial \mu^*(\lambda)}{\partial \lambda}.$$

Returning to (6.31), we can see that $\frac{\partial}{\partial \mu} f_{y^\delta, b}(\Phi \mu^*) = -\lambda$. Moreover, since $\frac{\partial \mu^*(\lambda)}{\partial \lambda} < 0$, then $\frac{\partial}{\partial \lambda} f_{y^\delta, b}(\Phi \mu^*(\lambda)) > 0$. Recalling that for any t we have $\lambda_{t+1} < \lambda_t$ and denoting by $\hat{\mu}_t = \mu^*(\lambda_t)$, we can write

$$\sigma_t = f_{y^\delta, b}(\Phi \hat{\mu}_t) < \sigma_{t+1} = f_{y^\delta, b}(\Phi \hat{\mu}_{t+1}),$$

which concludes the proof. \square

Remark 6.4.1. For BLASSO ($L^2 - |\cdot|$), the fidelity $f_{y^\delta, b}(\Phi \mu) = \frac{1}{2} \|\Phi \mu - y\|^2$ satisfies the hypothesis required to prove Proposition 6.4.1. Indeed,

$$\partial^2 f_{y^\delta, b} = \partial^2 \left(\frac{1}{2} \|\cdot - y^\delta\|^2 \right) = 1.$$

For the Poisson off-the-grid problem ($\tilde{\mathcal{D}}_{KL} - |\cdot|$) the fidelity $f_{y^\delta, b}(\Phi \mu) = \tilde{\mathcal{D}}_{KL}(\Phi \mu + b, y)$ is twice differentiable with positive second derivative

$$\partial^2 f_{y^\delta, b} = \partial \tilde{\mathcal{D}}_{KL}(\cdot + b, y) = \frac{y^\delta}{\|\cdot + b\|^2}.$$

Remark 6.4.2. It is interesting to observe that in the proof of Proposition 6.4.1, the explicit expression of the updating rule (6.27) is not needed. It is only necessary to have a strictly decreasing sequence of $(\lambda_t)_t$ in order to have a strictly decreasing sequence of $(\sigma_t)_t$. However, with a different updating rule there is no guarantee that at iteration $t + 1$ of Algorithm 11 the initialisation measure given by warm start $\mu_{t+1}^{[0]} = \hat{\mu}_t$ might be already an optimal solution for the problem $\mathcal{P}(\lambda_{t+1})$, requiring to immediately update again λ . For this reason the updating rule 6.27 is preferable.

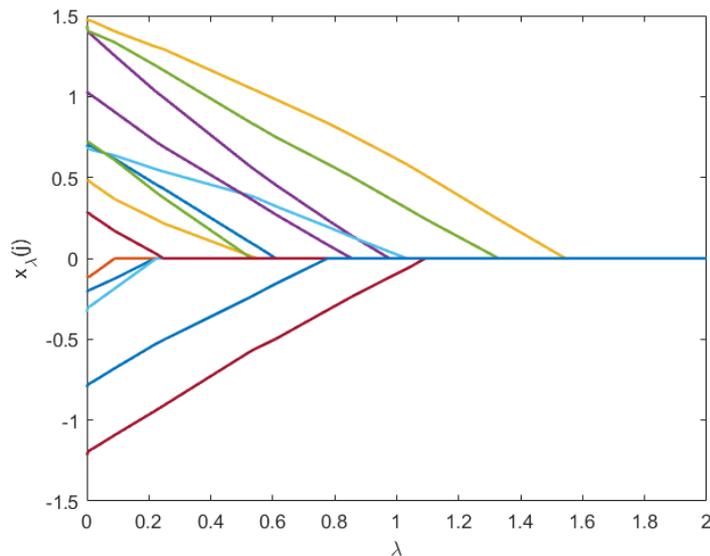


Figure 6.3: Regularisation path for LASSO (discrete setting). Plot of $\lambda \mapsto x_\lambda(j)$ with different colours for each different $j \in J$.

6.4.1.4 Regularisation path in the discrete setting

The homotopy algorithm has a link with the concept of regularisation path, defined in the discrete setting of Section 5.1.1 as follows.

Definition 6.4.2. *Let $x \in \mathbb{R}^N$ be sparse. Let x_λ be a reconstruction of x obtained with a regularisation algorithm with parameter $\lambda > 0$ given an acquisition $y \in \mathbb{R}^M$. Then, it is called component-wise regularisation path the set of functions*

$$\lambda \in \mathbb{R} \mapsto x_\lambda(j) \in \mathbb{R}, \quad j \in J,$$

where $J \subset \{1, \dots, N\}$ is the set of indices corresponding to the non-zero components of the desired solution x .

The regularisation path has been mainly studied in the discrete setting with an L^2 fidelity and L^1 penalty. In particular, in [158] authors considered the LASSO problem and proved that its regularisation path is piece-wise linear, as shown visually in Figure 6.3. The piece-wise linearity of the regularisation path has to be intended as follows. Let $x_{\lambda_1} \in \mathbb{R}^N$ and $x_{\lambda_2} \in \mathbb{R}^N$ solutions of LASSO with $\lambda_1 \neq \lambda_2$. By optimality, the (discrete) dual certificate $\eta(\lambda_1, x_{\lambda_1}) \in \mathbb{R}^N$ satisfies

$$\begin{aligned} |\eta(\lambda_1, x_{\lambda_1})(i)| &\leq 1 \quad \text{for all } i = 1, \dots, N \\ \eta(\lambda_1, x_{\lambda_1})(j) &= \text{sign}(x_{\lambda_1}(j)) \quad \text{for } j \text{ such that } x_{\lambda_1}(j) \neq 0 \end{aligned}$$

and the same applies to $\eta(\lambda_2, x_{\lambda_2}) \in \mathbb{R}^N$. If we assume $\eta(\lambda_1, x_{\lambda_1}) \in \mathbb{R}^N$ and $\eta(\lambda_2, x_{\lambda_2}) \in \mathbb{R}^N$ to be such that

- for all i , there holds

$$|\eta(\lambda_1, x_{\lambda_1})(i)| < 1 \iff |\eta(\lambda_2, x_{\lambda_2})(i)| < 1 \quad (6.32)$$

- for all $j \in \{1, \dots, N\}$ such that (6.32) is not met it holds

$$\eta(\lambda_1, x_{\lambda_1})(j) = \eta(\lambda_2, x_{\lambda_2})(j),$$

which implies $\text{sign}(x_{\lambda_1}(j)) = \text{sign}(x_{\lambda_2}(j))$,

then, for each $\theta \in [0, 1]$ the image $x^\theta = \theta x_{\lambda_1} + (1 - \theta)x_{\lambda_2}$ is a solution for $(\mathcal{P}(\lambda))$ with $\bar{\lambda} = \theta\lambda_1 + (1 - \theta)\lambda_2$, i.e.

$$x_{\bar{\lambda}} = x_{(\theta\lambda_1 + (1-\theta)\lambda_2)} = \theta x_{\lambda_1} + (1 - \theta)x_{\lambda_2}.$$

Moreover, its dual certificate has the same properties as: for all i it satisfies

$$|\eta(\bar{\lambda}, x_{\bar{\lambda}})(i)| < 1 \iff |\eta(\lambda_1, x_{\lambda_1})(i)| < 1 \iff |\eta(\lambda_2, x_{\lambda_2})(i)| < 1$$

and for all j such that $\eta(\bar{\lambda}, x_{\bar{\lambda}})(j) = 1$

$$\eta(\bar{\lambda}, x_{\bar{\lambda}})(j) = \eta(\lambda_1, x_{\lambda_1})(j) = \eta(\lambda_2, x_{\lambda_2})(j),$$

which entails $\text{sign}(x_{\bar{\lambda}}) = \text{sign}(x_{\lambda_1}(j)) = \text{sign}(x_{\lambda_2}(j))$.

Using an homotopy algorithm with an updating rule analogous to (6.27), in the discrete setting, every homotopy step thus produces a solution x_{λ_t} on a different linear segment of the regularisation path [67]. For this problem, homotopy techniques explore the Pareto frontier jumping between different linear segments of the regularisation path.

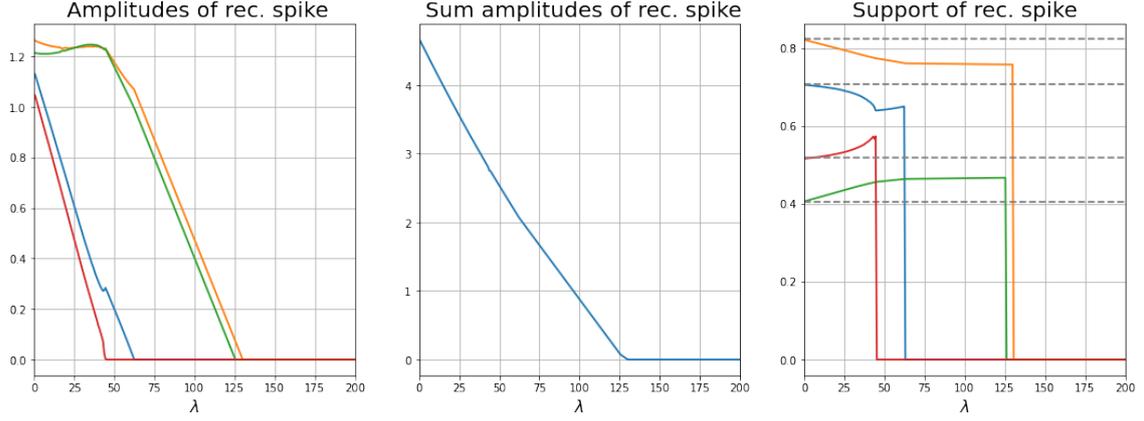
In the discrete setting, the study of the regularisation path has been investigated for more general variational problems in [25]. More specifically, when a smooth and convex fidelity term f is coupled with the L^1 norm, it is proven that the regularisation path of the general $(f - L^1)$ problem is piece-wise smooth. Its precise geometrical structure depends on the fidelity f .

6.4.1.5 Homotopy and regularisation path of BLASSO

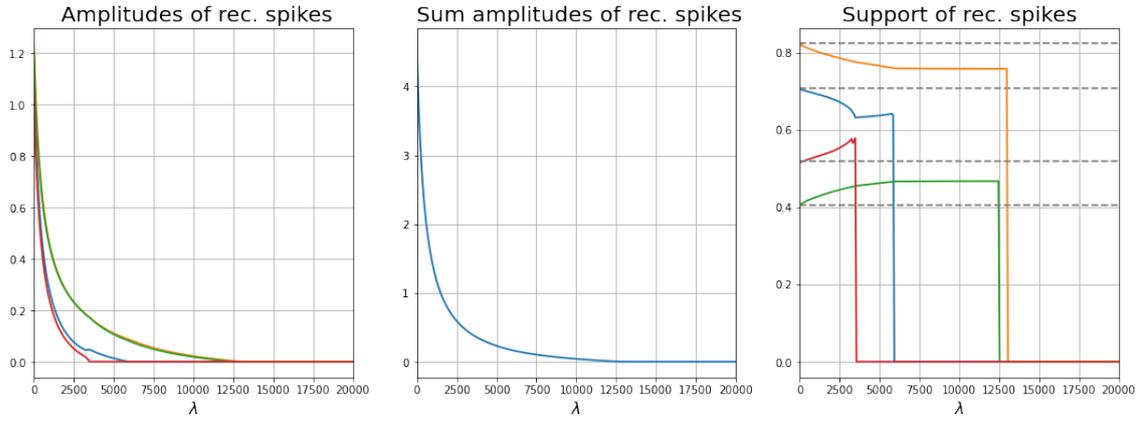
The study of regularisation paths in the continuous setting of off-the-grid problems is a bit more delicate. We start our analysis of the link between homotopy and regularisation paths with BLASSO $(L^2 - |\cdot|)$ problem.

Similarly as in the discrete setting, it is possible to show that the regularisation path of BLASSO is piece-wise linear. For $\mu_\lambda = \sum_{i=1}^N a_i^\lambda \delta_{x_i^\lambda}$ solution of BLASSO $(L^2 - |\cdot|)$, we denote with $J(\mu_\lambda)$ the support of the Diracs of μ_λ , that is $J(\mu_\lambda) = \{x_i^\lambda \mid i = 1, \dots, N\}$ and we recall that, by (5.12), its dual certificate $\eta(\lambda, \mu_\lambda) = \frac{1}{\lambda} \Phi^*(y - \Phi \mu_\lambda)$ satisfies $\eta(\lambda, \mu_\lambda)(x_j^\lambda) = \text{sign}(a_j^\lambda)$ for all $x_j^\lambda \in J(\mu_\lambda)$.

Let μ_{λ_1} and μ_{λ_2} be two solutions of BLASSO with $\lambda_1 \neq \lambda_2$ such that



(a) Regularisation path of BLASSO



(b) Regularisation path for the Poisson off-the-grid problem

Figure 6.4: Regularisation path in the off-the-grid setting

- $J(\mu_{\lambda_1}) = J(\mu_{\lambda_2})$;
- $\eta(\lambda_1, \mu_{\lambda_1})(x_j) = \eta(\lambda_2, \mu_{\lambda_2})(x_j)$, and thus $\text{sign}(a_j^{\lambda_1}) = \text{sign}(a_j^{\lambda_2})$, for all $x_j \in J(\mu_{\lambda_1}) = J(\mu_{\lambda_2})$.

Then, the measure defined by $\mu^\theta = \theta\mu_{\lambda_1} + (1 - \theta)\mu_{\lambda_2}$ is a solution of BLASSO with regularisation parameter $\bar{\lambda} = \theta\lambda_1 + (1 - \theta)\lambda_2$ and it moreover satisfies

- $J(\mu^\theta) = J(\mu_{\lambda_1}) = J(\mu_{\lambda_2})$;
- $\eta(\bar{\lambda}, \mu^\theta)(x_j) = \eta(\lambda_1, \mu_{\lambda_1})(x_j) = \eta(\lambda_2, \mu_{\lambda_2})(x_j)$ for all $x_j \in J(\mu^\theta)$, which implies that $\text{sign}(a_j^{\bar{\lambda}}) = \text{sign}(a_j^{\lambda_1}) = \text{sign}(a_j^{\lambda_2})$ for all $x_j \in J(\mu^\theta)$.

Indeed, since μ^θ is obtained as a convex combination of two measures that are weighted sum of the same Diracs, by construction $J(\mu^\theta)$ coincides with the support of μ_{λ_1} and μ_{λ_2} . Then, consider

$$\Phi^*(y - \mu^\theta) = \Phi^*(y - \Phi[\theta\mu_{\lambda_1} + (1 - \theta)\mu_{\lambda_2}]) = \theta\Phi^*(y - \mu_{\lambda_1}) + (1 - \theta)\Phi^*(y - \mu_{\lambda_2})$$

evaluated at $x_j \in J(\mu^\theta)$:

$$\begin{aligned}
\Phi^*(y - \mu^\theta)(x_j) &= \theta \Phi^*(y - \mu_{\lambda_1})(x_j) + (1 - \theta) \Phi^*(y - \mu_{\lambda_2})(x_j), \quad x_j \in J(\mu^\theta) \\
&= \theta \lambda_1 \operatorname{sign}(a_j^{\lambda_1}) + (1 - \theta) \lambda_2 \operatorname{sign}(a_j^{\lambda_2}) \\
&= (\theta \lambda_1 + (1 - \theta) \lambda_2) \operatorname{sign}(a_j^{\bar{\lambda}}) = \bar{\lambda} \operatorname{sign}(a_j^{\bar{\lambda}}) \\
\Rightarrow \eta(\bar{\lambda}, \mu^\theta)(x_j) &= \operatorname{sign}(a_j^{\bar{\lambda}}).
\end{aligned}$$

The latter shows that μ^θ is a solution of $(L^2 - |\cdot|)$ with $\bar{\lambda} = \theta \lambda_1 + (1 - \theta) \lambda_2$ and that it satisfies the properties listed above.

In Figure 6.4 (a), the regularisation path for a 1-dimensional spike deconvolution problem is presented in the case of BLASSO. Being in a continuous domain, the analytical illustrations have to be analysed carefully. Namely, changing λ does not only cause a change in the intensities and number of spikes in the reconstruction (or number of non-zero pixels) but also in their positions. On the left, the functions $\lambda \mapsto a_j(\lambda)$, with a_j being the amplitudes corresponding to the j -th spike of the reconstruction, are shown. Note that this is slightly different from the regularisation path that corresponds to the functions that map λ to the intensities of the spikes in the positions of the spikes in the ground truth. In the centre, we report the plot of the TV norm of the reconstructed spikes. On the right, we plot instead the functions $\lambda \mapsto x_j(\lambda)$, being x_j the position of the j -th spike of the reconstruction: as just described, the spikes are allowed to move in the continuous domain. The piece-wise linearity of the regularisation path is quite evident from the plot of $\lambda \mapsto a_j(\lambda)$ on the left.

In Figure 6.4 (b), we report the same plots in the case of a Poisson off-the-grid sparse deconvolution problem $(\tilde{\mathcal{D}}_{KL} - |\cdot|)$. Studying analytically the regularisation path in this scenario is challenging due to the presence of the Kullback-Leibler fidelity, which introduces non-linearity in the problem. Nonetheless, we can observe numerically that it is not piece-wise linear but it seems to have a piece-wise hyperbolic geometrical structure. This interesting claim has to be further investigated.

Recalling Algorithm 11, we thus observe that each homotopy iteration t produces a solution $\hat{\mu}_t$ belonging to a different linear segment of the piece-wise linear regularisation path. Thanks to the updating rule (6.27), the dual certificates $\eta(\lambda_t, \hat{\mu}_t)$ and $\eta(\lambda_{t+1}, \hat{\mu}_{t+1})$ cannot satisfy the conditions, defined at the beginning of this section, that characterise a linear segment of the path.

Without an analytical expression for the regularisation path in the Poisson $(\tilde{\mathcal{D}}_{KL} - |\cdot|)$ case, we cannot come to the same conclusion. However, the discussion about the updating rule (6.27) in Section 6.4.1.2 ensures that at each iteration we are improving the reconstruction and the numerical tests of the next section shows the good performance of Algorithm 11 in the resolution of $(\tilde{\mathcal{D}}_{KL} - |\cdot|)$. Moreover, Figure 6.2 refers to a 1D sparse deblurring problem with Poisson noise. We computed the Pareto frontier with SFW and the homotopy iterations with Algorithm 11 for the minimisation of $(\tilde{\mathcal{D}}_{KL} - |\cdot|)$. It is worth mentioning that the homotopy

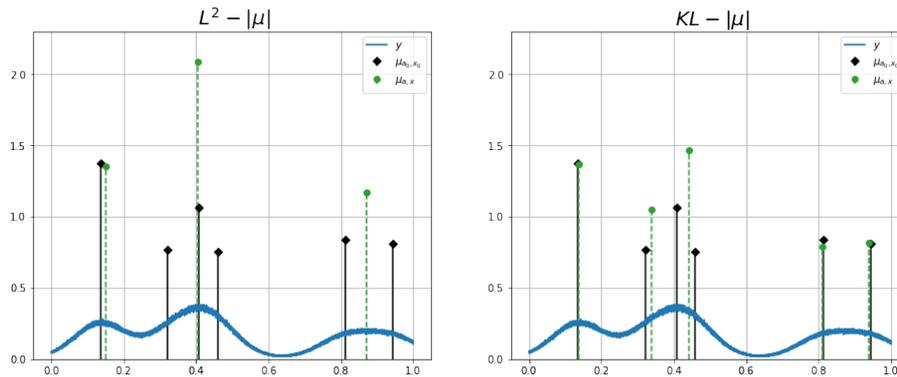


Figure 6.5: 1D comparison between Gaussian and Poisson noise modelling. In black: ground truth spikes. In green: reconstructed spikes. For both models, $\lambda = 8.82$.

algorithm in this case produces a sequence of iterations that moves along the Pareto frontier towards an optimal solution, as expected.

6.5 Numerical tests

In this section, we present numerical results of the proposed SFW algorithm plus homotopy step, on Poisson off-the-grid sparse deconvolution problems. We consider simulated 1-dimensional and 2-dimensional data and a real 3-dimensional dataset of fluorescent microscopy. In the following numerical experiments, the results computed by solving the $(\tilde{\mathcal{D}}_{KL} - |\cdot|)$ model were compared only with the ones obtained by solving the BLASSO model $(L^2 - |\cdot|)$ as we are not aware of any other off-the-grid model relevant for comparisons in the non-Gaussian noise setting.

As a first study, we consider a simulated dataset of 1-dimensional blurred acquisitions with Poisson noise and reconstruct the 1-dimensional sparse signals with the off-the-grid models $(\tilde{\mathcal{D}}_{KL} - |\cdot|)$ and BLASSO $(L^2 - |\cdot|)$, which is better suited for Gaussian noise, to compare the results. We consider in both cases the Sliding Frank-Wolfe algorithm (Alg. 10). Then, we discuss results obtained by the homotopy algorithm (Alg. 11) with SFW as an inner solver. In a 2-dimensional setting, we consider a blurred acquisition with Poisson noise of a sparse signal to study a practical way to choose a good value for σ_{target} and a strategy to estimate the background as spatially constant. We conclude our tests with the deconvolution of a 3D real image dataset where $(\tilde{\mathcal{D}}_{KL} - |\cdot|)$ is solved by means of the Boosted Sliding Frank Wolfe with the homotopy algorithm.

6.5.1 Comparison between Gaussian and Poisson modelling

The aim of the first set of experiments on simulated 1D blurred signals corrupted with Poisson noise is to validate our model $(\tilde{\mathcal{D}}_{KL} - |\cdot|)$ and to compare its performance with BLASSO $(L^2 - |\cdot|)$.

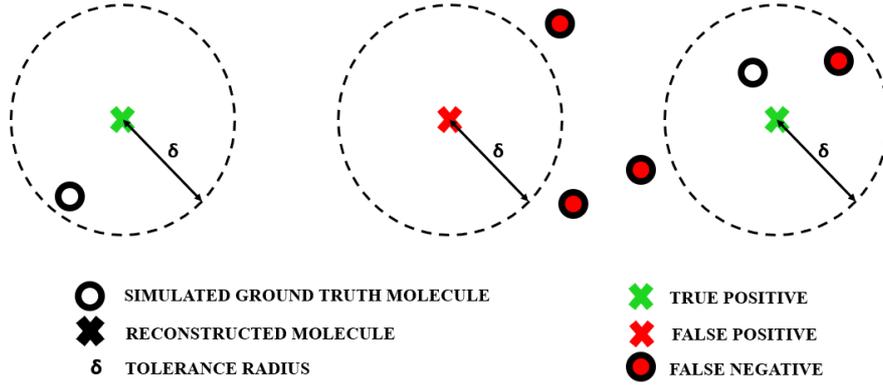


Figure 6.6: True Positives, False Positives, False Negatives spikes with respect to a tolerance radius $\delta > 0$.

We simulate 10 different ground truths with 6 randomly located spikes in the 1D domain $\Omega = [0, 1]$. The position of each spike of each simulated ground truth signal is sampled from a uniform distribution over Ω . The spikes' amplitudes are themselves sampled from a uniform distribution over $[1 - d, 1 + d]$ with $d = 0.4$. The corresponding acquisitions are blurred by a Gaussian 1D PSF with $\sigma = 0.07$, a spatially constant background $b = 0.01$ is considered, and then Poisson noise realisations are generated as acquired data. In Figure 6.5, one of the simulated ground truth signals μ_{gt} is shown (black) with the corresponding Poisson noisy and blurred data (blue).

We then reconstruct all 1D signals using both the Poisson noise model ($\tilde{\mathcal{D}}_{KL} - |\cdot|$) and the Gaussian noise model BLASSO ($L^2 - |\cdot|$) with $\lambda \in (0, 10]$. Figure 6.5 shows an example of the two reconstruction μ_{rec} (green) for $\lambda = 8.82$. It is evident that, in this particular case, the Poisson model separates better the spikes in the ground truth. In particular, with BLASSO only 3 out of 6 spikes are reconstructed whilst our proposed model ($\tilde{\mathcal{D}}_{KL} - |\cdot|$) manages to retrieve 5 spikes, with an accurate localisation.

To evaluate the goodness of the reconstructions, we consider the Jaccard index defined in terms of the number of True Positive, False Positive and False Negative spikes as follows

$$\text{Jac}_{\delta}(\mu_{\text{gt}}, \mu_{\text{rec}}) = \frac{\#\text{TP}}{\#\text{TP} + \#\text{FP} + \#\text{FN}} \in [0, 1]$$

with tolerance radius $\delta > 0$. In Figure 6.6, the meaning of True Positive, False Positive and False Negative is visually explained. We call TP the reconstructed spikes that are at a distance less than δ from a spike in the ground truth, while reconstructed spikes that are more than δ distant from each ground truth spikes are called FP. FN are spikes in the ground truth which have not been associated to any TP. Another good metric for the reconstruction quality is the RMSE of the

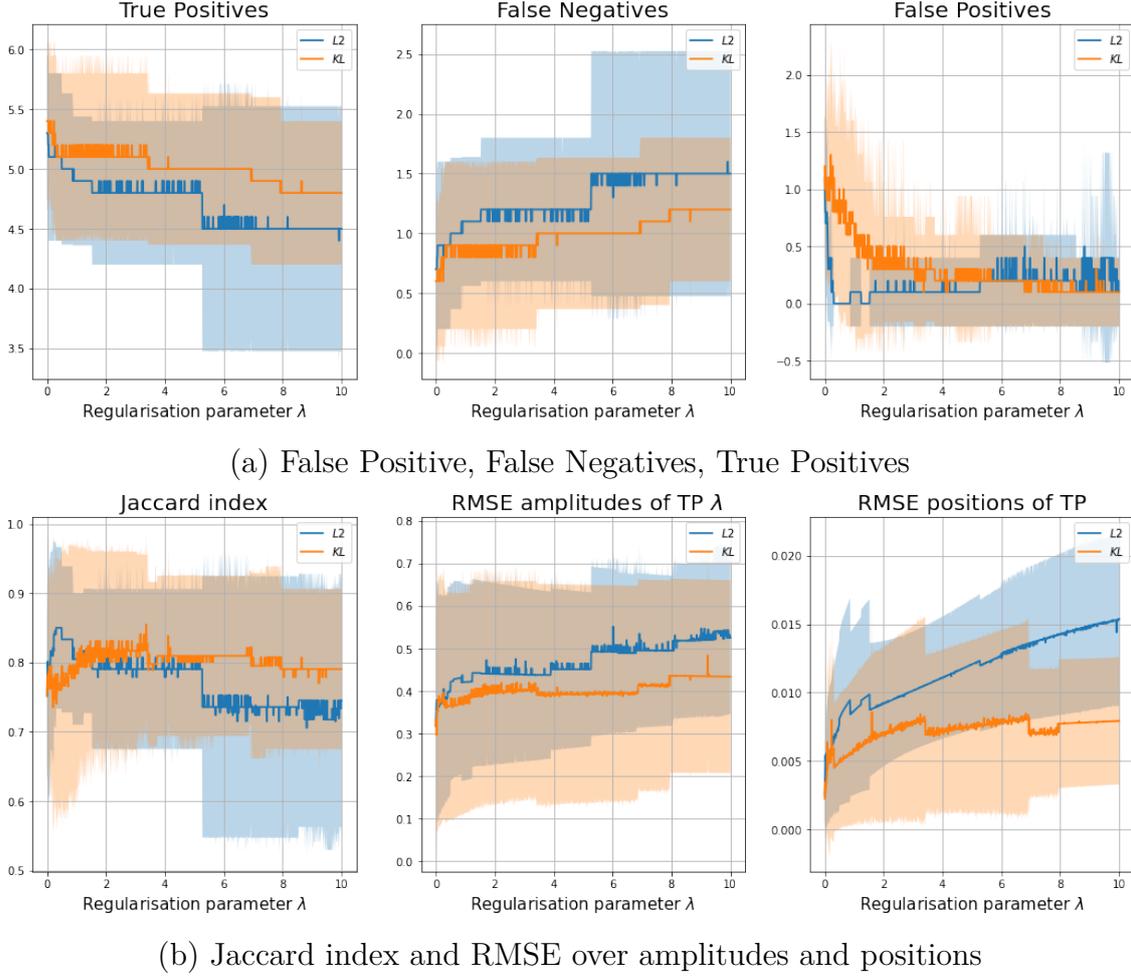


Figure 6.7: Mean values over 100 different randomly generated ground truths with 6 spikes and their corresponding reconstructions. Shaded area corresponds to standard deviation. Maximum number of iterations of SFW: $2N_{molecules}$. Tolerance radius $\delta = 0.05$.

amplitudes a and positions x of the TP spikes:

$$\text{RMSE}_x(\mu_{gt}, \mu_{rec}) = \sqrt{\frac{1}{\#\text{TP}} \sum_{i \in \text{TP}} ((x_{rec})_i - (x_{gt})_i)^2}$$

$$\text{RMSE}_a(\mu_{gt}, \mu_{rec}) = \sqrt{\frac{1}{\#\text{TP}} \sum_{i \in \text{TP}} ((a_{rec})_i - (a_{gt})_i)^2}.$$

In Figure 6.7(a), In Figure 6.7(a), we plot TP, FN and FP with respect to λ and in Figure 6.7(b) the Jaccard index (computed with $\delta = 0.05$) and the RMSE of amplitudes and positions. Note that the Poisson model ($\tilde{\mathcal{D}}_{KL} - |\cdot|$) has a better performance in terms of TP and FN spikes and in terms of Jaccard index and RMSE of amplitudes and positions. Only for the number of FP, BLASSO ($L^2 - |\cdot|$) results slightly better than ($\tilde{\mathcal{D}}_{KL} - |\cdot|$) for small values of λ . This is due to the fact that

Parameters	$L^2 - \mu $	$\tilde{\mathcal{D}}_{KL} - \mu $
Max. number of homotopy iterations	$2N_{\text{molecules}}$	$2N_{\text{molecules}}$
Max. number of inner SFW iterations	1	1
Homotopy parameter c	15	40
Homotopy parameter γ	0.9	0.9
Choice of σ_{target}	$1.5 \times \frac{1}{2} \ \Phi\mu_{\text{gt}} + b - y\ ^2$	$1.5 \times \mathcal{D}_{KL}(\Phi\mu_{\text{gt}} + b, y)$

Table 6.2: Parameters used with the homotopy algorithm (Alg.11) in the 1D simulated comparison tests

$(\tilde{\mathcal{D}}_{KL} - |\cdot|)$ usually requires more iterations of SFW before reaching convergence. It results in a better estimation of the number of molecules, with TP being closer to the actual number of spikes in the ground truth and it may cause an overestimation of the number of spikes with a consequently higher value of FP. However, this behaviour of $(\tilde{\mathcal{D}}_{KL} - |\cdot|)$ has overall a positive impact on the quality of the reconstructions, with all the other considered indices showing its better performance with respect to BLASSO $(L^2 - |\cdot|)$.

With the same dataset, we then compare the results obtained after running the homotopy algorithm (Alg. 10) for the automatic selection of the regularisation parameter λ for both methods $(L^2 - |\cdot|)$ and $(\tilde{\mathcal{D}}_{KL} - |\cdot|)$, with the parameters specified in Table 6.2. Algorithm 11 requires that at each homotopy iteration the problem with $\lambda = \lambda_t$ is solved to convergence with SFW. However, to avoid excessive computational time, we observed that one iteration of SFW has been always enough, since the estimated spikes at each homotopy iteration are then updated again in the next one. Then, for the same reason we set the maximum number of homotopy outer iterations to be equal to twice the number of peaks in the ground truth. As far as σ_{target} is concerned, we observe that, being in a simulated environment, it is possible to compute exactly the value of (6.26) as $f_{y^\delta, b}(\Phi\mu_{\text{gt}})$. However, since this is not possible in real situations, we will discuss a possible strategy to estimate σ_{target} in the next section. In Table 6.3, we report the values of TP, FN, FP, Jaccard index and RMSE of the reconstructed signals obtained using homotopy for $(L^2 - |\cdot|)$ and $(\tilde{\mathcal{D}}_{KL} - |\cdot|)$. The final estimated λ is also reported in Table 6.3, together with the number of performed homotopy iterations and of the value σ_{target} . First, we observe that using homotopy we retrieve values which are comparable with the best ones obtained using SFW with grid search. This shows the effectiveness of the homotopy strategy for the selection of a good λ . Similarly as above, $(\tilde{\mathcal{D}}_{KL} - |\cdot|)$ yields better results than $(L^2 - |\cdot|)$ in the presence of Poisson noise, with a reduction of the number of FN and an improvement of the accuracy in terms of Jaccard index. These first results confirm the advantage of $(\tilde{\mathcal{D}}_{KL} - |\cdot|)$ model for Poisson noisy acquisitions, which aimed our proposal.

	$L^2 - \mu $	$\tilde{\mathcal{D}}_{KL} - \mu $
Jaccard index	0.74	0.76
Number of TP	4.50	4.80
Number of FN	1.50	1.20
Number of FP	0.10	0.40
RMSE on amplitudes of TP	0.41	0.44
RMSE on positions of TP	0.014	0.015
Final estimated λ	6.09	40.21
Number of homotopy iterations	4.55	3.93
Value of σ_{target}	4.09	77.16

Table 6.3: Homotopy algorithm: comparison between BLASSO and the Poisson off-the-grid modelling. Mean values over 100 different randomly generated ground truths with 6 spikes and their corresponding reconstructions.

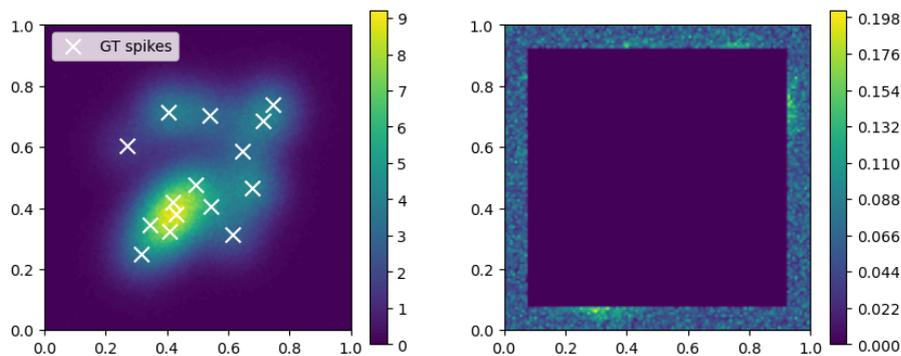
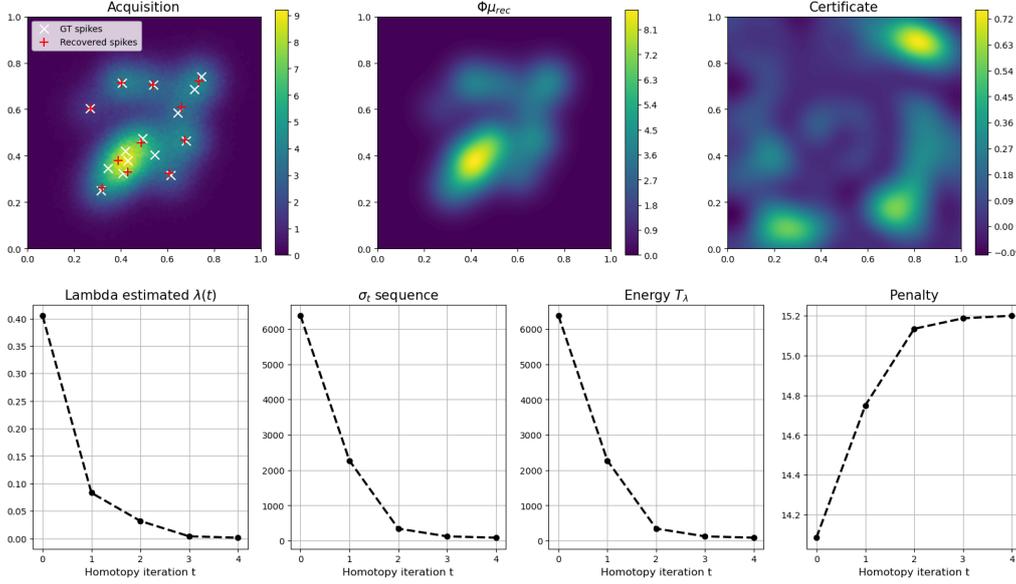


Figure 6.8: 2D sparse ground truth image (white crosses) and its corresponding noisy blurred acquisition y^δ on $\Omega = [0, 1]^2$ (obtained with a 2D Gaussian PSF with $\sigma = 0.07$, constant background $b = 0.05$, Poisson noise). On the right, visualisation of $y^\delta|_{\Omega_{\text{bg}}}$ corresponding to background noise, i.e. the external square-ring.

6.5.2 Homotopy algorithm: choice of σ_{target}

As mentioned above, in the case of real data, the choice of σ_{target} obviously must not require the knowledge of the ground truth image. In this section, we discuss a way to estimate a reasonable value σ_{target} that only relies on the acquisition and on the assumption that the signal is sparse. We make this discussion in a 2D fluorescence microscopy simulated setting, which is better for visualisation purposes. Under a suitable sparsity level of the ground truth image, it is safe to assume that its corresponding noisy and blurred acquisition y^δ presents regions containing only background noise, that we denote with $\Omega_{\text{bg}} \subset \Omega$. In Figure 6.8, we show on the left y^δ and on the right $y^\delta|_{\Omega_{\text{bg}}}$, i.e. the acquisition restricted to the area of background noise in the external square-ring. We propose to estimate the value of σ_{target} as

Theoretical value (based on the ground truth)	$\tilde{\mathcal{D}}_{KL}(\Phi\mu_{gt} + b, y^\delta)$	842.3
Poisson discrepancy principle (6.35) [17]	$\frac{ \Omega }{2}$	8192
Estimation based only on y^δ and Ω_{bg} (6.33)	$\mathcal{D}_{KL}(b, y_{bg}) \frac{ \Omega }{ \Omega_{bg} }$	846.4

Table 6.4: Different estimates of σ_{target} .Figure 6.9: 2D reconstruction with homotopy Alg.11 with parameters: max. number outer homotopy iterations 20, max. number inner SFW iterations 1, $c = 1$, $\gamma = 0.2$.

follows

$$\sigma_{\text{target}} = f_{y^\delta, b}(0) \Big|_{\Omega_{bg}} \frac{|\Omega|}{|\Omega_{bg}|}, \quad (6.33)$$

where the restriction of the fidelity term to $\mu = 0$ is due to the fact we assume the desired image μ to be null in Ω_{bg} , i.e. $\mu|_{\Omega_{bg}} = 0$. If we considered $\sigma_{\text{target}} = f_{y^\delta, b}(0) \Big|_{\Omega_{bg}}$, it would be equivalent to assuming the noise to be null in $\Omega \setminus \Omega_{bg}$, which obviously is not true. Then we adjust (6.33) in order to account for noise not only in Ω_{bg} but on the whole domain Ω . Thanks to Ω_{bg} , it is also possible to approximately estimate the background $b \in L^2(\Omega)$ when it is not known. Indeed, we propose a spatially constant estimate of b as mean integral value

$$b = \frac{1}{|\Omega_{bg}|} \int_{\Omega_{bg}} y^\delta(t) dt. \quad (6.34)$$

For the image shown in Figure 6.8, we compare in Table 6.4 different choices of σ_{target} for $(\tilde{\mathcal{D}}_{KL} - |\cdot|)$ with $f_{y^\delta, b}$ being the Kullback-Leibler fidelity term. In particular, we consider in the second row the estimate proposed in [17], where a discrepancy principle for Poisson data is studied under the following approximation

$$\tilde{\mathcal{D}}_{KL}(\Phi\mu_{gt} + b, y^\delta) \approx \frac{|\Omega|}{2}. \quad (6.35)$$

This value is obtained by computing the expected value for Kullback-Leibler fidelity and by approximating it with a first order Taylor expansion. As observed also in [17], the estimate (6.35) might not be optimal and, indeed, one should consider

$$\tilde{\mathcal{D}}_{KL}(\Phi\mu_{gt} + b, y^\delta) \approx \frac{1 - \epsilon}{2} |\Omega|,$$

where ϵ is a small positive or negative number. When Ω is big, this might lead to bad estimates even if ϵ is small. Indeed, Table 6.4 shows that in our modelling (6.35) does not give an accurate estimation of σ_{target} . On the contrary, the estimation given by (6.33), on the last row of Table 6.4, is very close to the real value (in the first row), computed exploiting the knowledge of the ground truth. In [22, 23, 24] other strategies have been considered and, in particular, in [23] a similar masking approach is studied.

By using the homotopy algorithm (Alg.11) with SFW for the reconstruction of the image in Figure 6.8 with background and σ_{target} estimated by (6.34) and (6.33) respectively, we obtain the results shown in Figure 6.9.

6.5.3 Numerical test with 3D real dataset

In this last section, we consider a 3D real blurred and noisy volume acquired by means of a widefield microscope. I have been given access to this real dataset by Jerome Boulanger (MRC Laboratory of Molecular Biology, Cambridge) during my visiting period at Cambridge Image Analysis group at DAMPT, University of Cambridge, between May and June 2022. The image was acquired in widefield microscopy with a 100x/1.49 objective, a Hamamatsu Flash 4 camera with 6.5um pixels, by Alejandro Melero Carrillo from Liz Miller's group at the MRC-LMB and was used in [100]. It is an acquisition of yeasts expressing fluorescent proteins (*SEC16-sfGFP* and a *SEC24-sfGFP*) localised at the Endoplasmic Reticulum exit sites (ERES). The acquired volume y^δ is shown in Figure 6.13 (first row) with maximum intensities projections over the xz , yz , yx planes. The 3D volume blurred and noisy acquisition has $190 \times 190 \times 17$ voxels with voxel size of 65nm in yz and 250nm in z . Signal dependency of the noise is observed.

To reconstruct a sparse volume from the acquisition y^δ , we use the homotopy algorithm (Alg. 11) with the Boosted SFW as an inner solver to minimise $(\tilde{\mathcal{D}}_{KL} - |\cdot|)$. We prefer the accelerated version of SFW, since in 3D the computational costs are significantly higher than in lower dimensions.

As forward operator, we consider a 3D convolution kernel (5.6) estimated as a 3D Gaussian PSF (as in Figure 6.10),

$$\varphi((x, y, z)) = \frac{1}{\sqrt{(2\pi)^3 \sigma_x \sigma_y \sigma_z}} \exp\left[-\frac{x^2}{2\sigma_x^2}\right] \exp\left[-\frac{y^2}{2\sigma_y^2}\right] \exp\left[-\frac{z^2}{2\sigma_z^2}\right].$$

The standard deviation σ_x, σ_y of the 3D PSF can be estimated from the Full Width

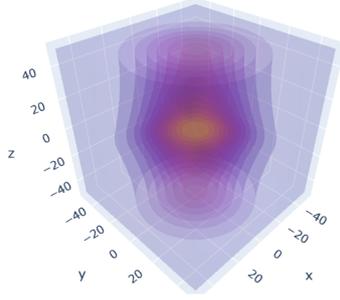


Figure 6.10: 3D Gaussian PSF

Parameters	$(\tilde{\mathcal{D}}_{KL} - \cdot)$
Max. number of homotopy iterations	10
Max. number of inner SFW iterations	50
Homotopy parameter c	0.5
Homotopy parameter γ	0.9
σ_{target} given by (6.33)	1102067.75
σ_{target} given by (6.35)	306850
Constant background estimate (6.34)	$b = 337.77$

Table 6.5: Parameters used for Algorithm 11 for the reconstruction of the 3D volume

Half Maximum, which is given by

$$\text{FWHM} = 0.61 \frac{\lambda_{\text{wavelength}}}{\text{NA}},$$

where $\lambda_{\text{wavelength}}$ is the emission wavelength of the fluorescent proteins and NA is the numerical aperture of the microscope (see Appendix B for more details). If the FWHM is known, then it is possible to retrieve information about the variance parameters of the PSF since $\text{FWHM} = 2.355 \cdot \sigma_x$ and $\text{FWHM} = 2.355 \cdot \sigma_y$. For the standard deviation in the z -axis, we assume $\sigma_z = 2 \cdot \sigma_x$. Since the value of the numerical aperture is known, $\text{NA} = 1.49$, and $\lambda_{\text{wavelength}} = 508\text{nm}$ for the green fluorescent proteins used, we obtain a PSF estimation with $\sigma_x = \sigma_y = 89\text{nm}$ and $\sigma_z = 178\text{nm}$ shown in the second row of Figure 6.13, which appears to be a good estimate for the underlying blur model.

To use the homotopy algorithm for the regularisation parameter automatic selection, we need to choose some important parameters that are reported in Table 6.5. In particular, for σ_{target} we observe that (6.35) and (6.33) give very different results and we consider the estimate given by (6.33), as in the previous section it proved to be more accurate. The background, whose restriction to Ω_b is shown in the third row (c) of Figure 6.13, is estimated by (6.34) with $b = 337.77$. In Figure 6.13 (forth row) and Figure 6.12, we show the reconstructed volume obtained with Algorithm 11. We computed this reconstruction up to 50 homotopy iterations fixing

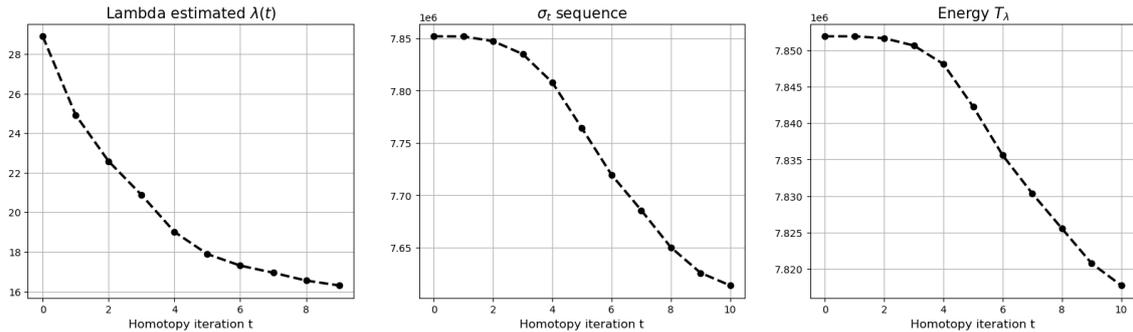


Figure 6.11: ERES 3D data. Values of λ_t , σ_{target} and cost functional along the homotopy iterations.

to 10 the maximum number of inner loops for the BSFW algorithm, using Google Colab CPUs for about 10 hours. We obtained a reconstructed volume μ_{rec} with 274 spikes. We assume this number to be underestimated, comparing the acquisition y^δ (top row) with the blurred $\Phi\mu_{\text{rec}} + b$ (bottom row) in Figure 6.13, and comparing the value of σ_{target} of Table 6.5 with the one attained along the iterations (see Figure 6.11). However, this first result is promising, since the use of Algorithm 11 yields very precise localisation of spikes automatically, with no need of estimating the regularisation parameter, and no information about the data.

6.6 Final discussion

In this chapter, we considered off-the-grid Poisson inverse problems in the space of Radon measures. Our contribution is the study of a variational model combining a Kullback-Leibler fidelity term, a TV norm of measures as penalty and non-negativity constraints. We presented a detailed study of the optimality conditions and of the corresponding dual problem. Then, we discussed how to implement the Sliding Frank-Wolfe algorithm to solve the Poisson off-the-grid model and the importance of the choice of a good regularisation parameter λ . For this last reason a tailored homotopy algorithm is presented. It provides a way to automatically select the parameter by decreasing it subsequently at each iteration if a certain informed condition is not met. We discussed the properties of the homotopy algorithm, which are linked to the concepts of Pareto frontier and regularisation path. Finally, we presented numerical experiments on simulated data to validate the proposed Poisson model regularisation, to compare its performance with BLASSO and to verify the proposed homotopy strategy yields good reconstructions. Concerning the estimation of important parameters of Algorithm 11, we proposed to estimate the background and the noise level on a masked region of the acquisition y^δ , where it is safe to assume no signal has been measured. To conclude, we reported a numerical experiment on a 3D real data of fluorescence microscopy, that we reconstructed using the homotopy algorithm with an accelerated version of SFW, the Boosted SFW, as inner solver.

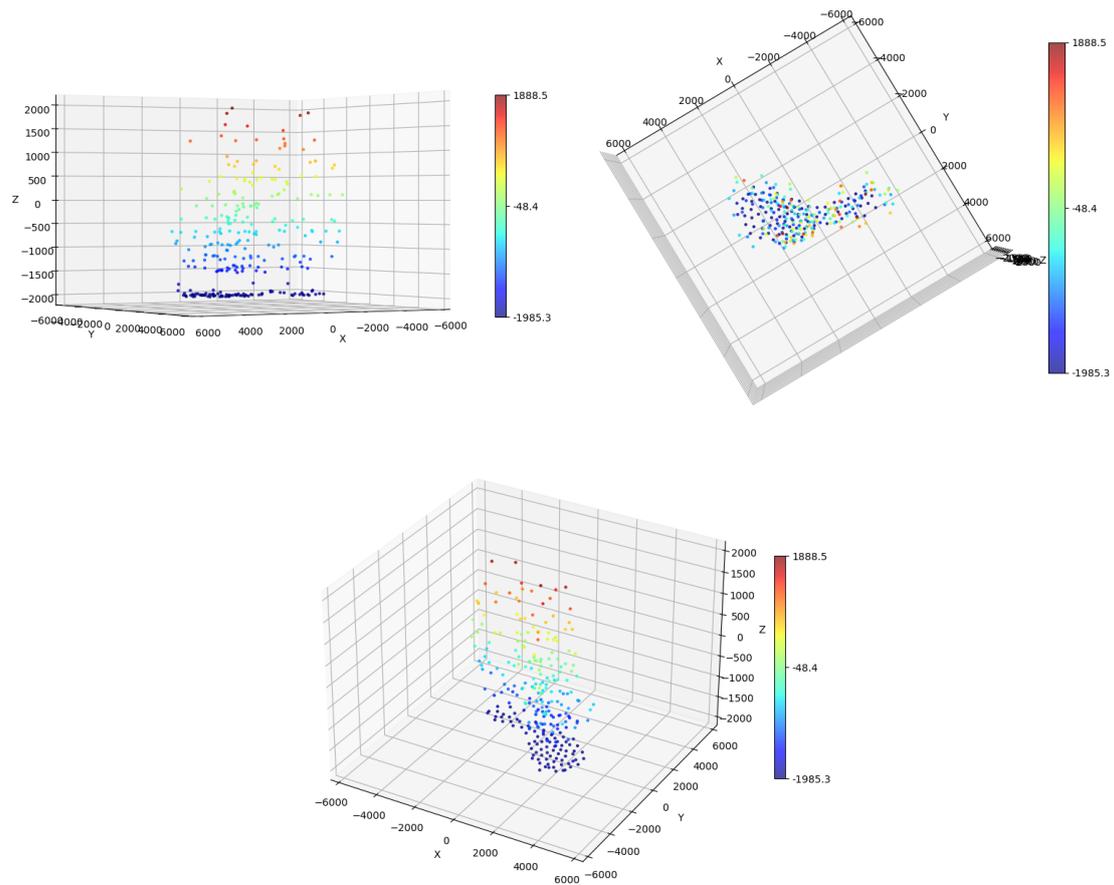


Figure 6.12: Sparse reconstruction of the 3D real volume acquisition. Visualisation from different angles. The colours corresponds to the depth along the z direction.

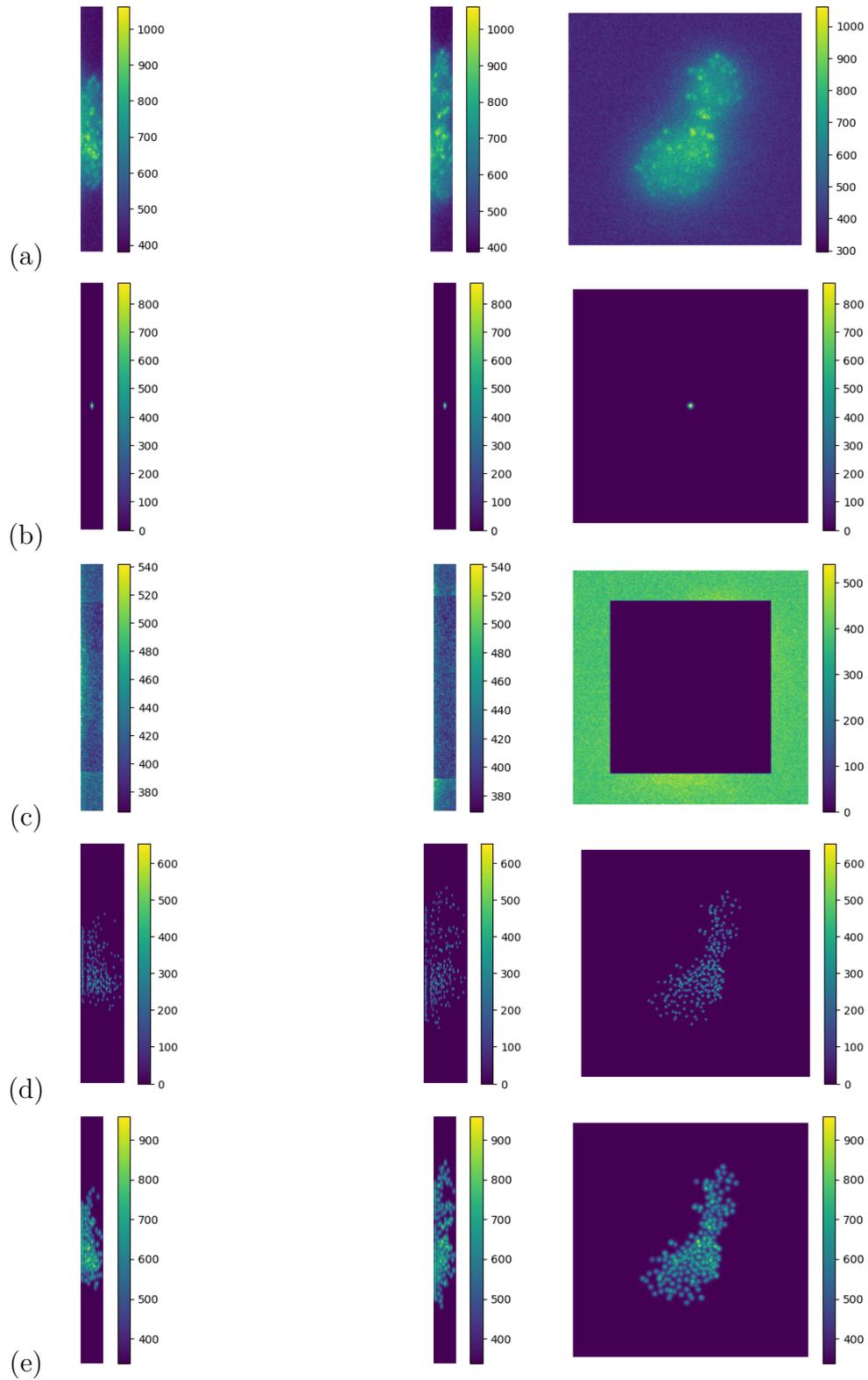


Figure 6.13: ERES 3D real data. From top to bottom: (a) acquired volume y^δ , (b) estimated Gaussian 3D PSF, (c) $y^\delta|_{\Omega_b}$, (d) reconstructed volume with 130 iterations of the homotopy algorithm μ_{rec} , (e) blurred observation $\Phi\mu_{\text{rec}} + b$ corresponding to the reconstruction μ_{rec} .

Conclusions

This thesis focused on the resolution of inverse problems in Banach spaces via variational approaches. Tailored optimisation techniques have been defined, since Banach spaces lack of the Riesz isomorphism and scalar product, which are often explicitly or implicitly used in the design of classical minimisation algorithms in Hilbert spaces. The first part of this thesis is concerned with Lebesgue spaces with variable exponents $L^{p(\cdot)}(\Omega)$ and on the use of gradient-based algorithms in this setting for smooth and non-smooth optimisation. In Chapter 2, we presented the main tools that we used for these methods in $L^{p(\cdot)}(\Omega)$, namely the modular functions and their derivatives, which we proposed to use instead of the Luxemburg norm and the duality maps in $L^{p(\cdot)}(\Omega)$, respectively. We showed both analytically and numerically that modulars allow much faster computations, having a closed form explicit expression and being separable. The same applies to their derivatives, that we proposed as modular-based alternatives to duality maps. In Chapter 3, in particular, we defined novel modular-based gradient-descent algorithms in $L^{p(\cdot)}(\Omega)$, where the descent step is computed in the dual space, and we further design its stochastic variant whose advantages are validated on CT data. Chapter 4 is, instead, focused on non-smooth optimisation proximal gradient methods in $L^{p(\cdot)}(\Omega)$. We proposed two instances of modular-based proximal gradient algorithms, where in the first gradient descent is performed in the primal space and in the second it is performed in the dual, as for Chapter 3. Extensive numerical validations confirmed that the choice of employing these algorithms tailored for the specific Banach spaces at hand is particularly suited to deal with mixed noise statistics, heterogeneous or sparse signals and to reconstruct discontinuities with more precision.

In the second part of this thesis, inverse problems in the Banach space of Radon measures $\mathcal{M}(\Omega)$ are considered. Chapter 5 presented a literature review on sparse off-the-grid problems, on the BLASSO variational problem proposed for Gaussian noise settings, and on standard optimisation algorithms to solve BLASSO in $\mathcal{M}(\Omega)$. Chapter 6 outlined our contribution on this topic. Motivated by fluorescence micro-

scopy applications, where the noise observed is rarely purely Gaussian, we considered a Poisson noise hypothesis on the data and proposed a variational approach where the total-variation norm is coupled with the Kullback-Leibler fidelity term and a positivity constraint. We analytically derived its optimality conditions and we studied its corresponding dual problem. We showed that Sliding Frank-Wolfe algorithm can be adapted to solve the minimisation of the new Poisson variational functional and we considered an homotopy strategy to automatically and iteratively select a good regularisation parameter for our problem. Numerical tests demonstrated that the use of the Poisson fidelity effectively improves the reconstruction quality and that the homotopy algorithm reduces the dependence of the reconstruction on the choice of regularisation parameter, providing a good estimate for the best λ , according to the noise intensities.

In both parts, some interesting questions remain open. In particular, concerning Part I:

- In Section 3.3.4, we detailed some possible strategies for an informed variable exponent $p(\cdot)$ selection for the resolution of inverse problems in $L^{p(\cdot)}(\Omega)$. The exponent map $p(\cdot)$ can be obtained interpolating an approximation of the desired solution between two fixed values p_- and p_+ . This method presents two main problems. It requires another reconstruction method to obtain an approximate solution to select $p(\cdot)$, and thus this strategy is not self-sufficient. Moreover, parameters p_- and p_+ have to be manually tuned and their choice is strictly problem dependent. Even though this selection method has been proved quite effective both in terms of reconstruction quality and in terms of convergence speed, we lack of theoretical guarantees. Since $p(\cdot)$ and the unknown solution x are related, a possibility could be to consider a minimisation problem on both exponent and solution, and to solve it with bi-level optimisation techniques.
- In Chapter 3, we proposed modular-based gradient descent and stochastic modular-based gradient descent algorithms in $L^{p(\cdot)}(\Omega)$. For the deterministic one, we observed in Chapter 4 that it is possible to retrieve a result of convergence in function values because it follows from Proposition 4.3.3, in the framework of smooth optimisation. However, for both deterministic and stochastic modular-based gradient descent it would be meaningful to also prove convergence to x^\dagger , in the noiseless case. We attempted to adapt standard descent lemmas to the variable exponent setting, but the lack of p -convexity of $L^{p(\cdot)}(\Omega)$ makes it almost impossible to obtain the same bounds and majorisations. If $L^{p(\cdot)}(\Omega)$ were p -convex, then it would be possible to conclude the convergence proofs. Thus, we suggest the following strategies: either to prove the p -convexity for some p , or to formulate some versions of the descent lemma in a more problem-specific way.

Concerning Part II:

-
- As far as the homotopy strategy presented in Chapter 6 is concerned, we recall that the proposed updating scheme for λ along the homotopy steps ensures that they happen on different linear segment of the solution path for the BLASSO problem. Differently from BLASSO, we observe that the solution path for $(\tilde{\mathcal{D}}_{KL} - |\cdot|)$ is not piece-wise linear but has a piece-wise structure too, then again the homotopy steps happen each on a different piece of the frontier. However, we have not been in able to prove analytically which is the exact shape of the regularisation path. Our claim is that in the Poisson fidelity scenario it has a piece-wise hyperbolic structure. By further analysing this aspect, one would have more insights on the nature of the problem and on alternative updating rule for λ , more adapted to the structure of the problem at hand.
 - Referring to the homotopy algorithm, in Chapter 6 we showed that the sequence of σ_t , which is the residual at each homotopy step, is strictly decreasing. We would like to show that the algorithm has a finite termination rule. A sufficient condition is $\sigma_t \rightarrow 0$. Although this is something that we observe in practice, from a theoretical standpoint it can happen only in absence of noise, hence leaving an open question.

Appendix

Computed Tomography: forward model and geometries

In Chapter 3, numerical tests on Computed Tomography imaging problems have been carried out to validate the proposed Algorithms 3 and 4, that are modular-based strategies for smooth optimisation in $L^{p(\cdot)}(\Omega)$. In particular, experiments in a 2D parallel beam geometry simulated setting and in a 2D fan beam geometry real case are shown. Here, we briefly detail the functioning and the forward model of CT, describing in particular the aforementioned geometries.

Computed Tomography (CT) is a medical imaging technique that uses X-rays to create detailed cross-sectional images of the body, scanned with X-ray beams. CT scanners consist of an X-ray source and a detector array mounted on opposite sides of a rotating structure. They make two measurements: the initial intensity of each X-ray beam at the radiation source and the final intensity of each beam, at the radiation detector, as sketched in Figure A.1. The changes in intensity for a single beam are dependent on the internal density of the body along the line the X-ray

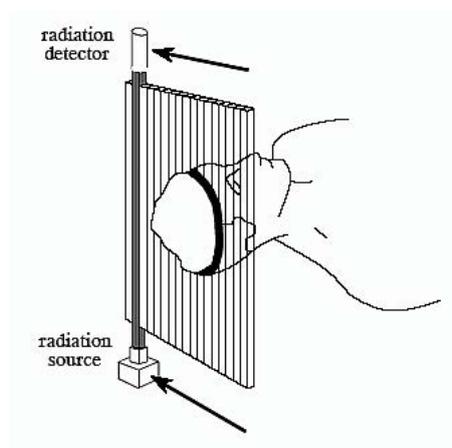


Figure A.1: CT scanner. Image from [12]

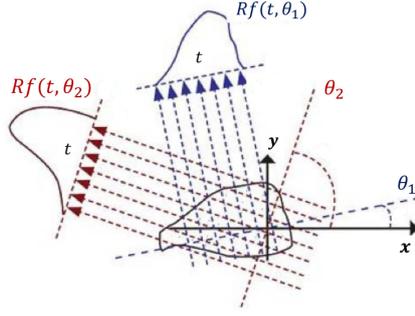


Figure A.2: 2D Parallel Beam Geometry. Image from [12]

passes through. The goal is to determine the attenuation coefficients of the X-rays caused by body absorption, which being related to the body's density, provide an image of the scanned region. The rotation of both the X-ray source and the detector array is a key aspect of this process, as it allows to measure attenuation coefficients from multiple angles around the body. For each considered angle, CT scanners collect multiple measurements, each corresponding to different beams.

A forward operator based on the Radon transform mathematically describes how X-rays interact with the body and how the acquired data are formed. Given a compact support $\Omega \subset \mathbb{R}^d$, consider an object of interest with density $f \in L^2(\Omega)$. The trajectory of the X-ray beam through the object usually depends on two parameters, the angle θ of rotation of the beam and the offset t , as shown in Figures A.2 and A.3. For each angle $\theta \in \tilde{\Omega} = [0, 2\pi]^{d-1}$ and for each offset $t \in I \subset \mathbb{R}$, the trajectory is a curve (generally a straight line) $\gamma_{t,\theta} : \mathbb{R} \rightarrow \Omega$, $s \mapsto \gamma_{t,\theta}(s) \in \Omega$. Then, the Radon transform of f is defined for all $t \in \mathbb{R}$ and $\theta \in \tilde{\Omega}$ as

$$\mathcal{R}f(t, \theta) = \int_{\mathbb{R}} f(\gamma_{t,\theta}(s)) ds. \quad (\text{A.1})$$

In CT imaging, two common geometries for data acquisition are the 2D parallel beam and the 2D fan beam geometries.

2D Parallel Beam Geometry. In this configuration, a collimated X-ray beam is projected through the patient, and detectors are arranged in a straight line opposite the X-ray source, see Figure A.2. The X-ray source and detectors remain parallel during the rotation. As the X-ray source and detectors rotate around the patient, measurements are taken at each angle $\theta \in \tilde{\Omega} \subset [0, 2\pi]$. The detectors, placed at positions $t \in I \subset \mathbb{R}$, record the attenuation of X-rays as they straightly pass through the body at various positions. As a result of this process, the acquired data is an element of $L^2(I \times \tilde{\Omega})$ obtained by (A.1) with

$$\gamma_{t,\theta}(s) = (t \cos(\theta) - s \sin(\theta), t \sin(\theta) + s \cos(\theta)), \quad s \in \mathbb{R}.$$

By considering the simulated setting of Section 3.4.3, we have $\tilde{\Omega} = \left\{ \frac{k\pi}{180}, k = 1, 2, \dots, 180 \right\}$ and $I = \{1, 2, \dots, 256\}$.

2D Fan Beam Geometry. In fan beam geometry, the X-ray source emits a diverging beam, and the detectors are arranged in an arc or fan shape opposite the source, as in Figure A.3. This geometry is closer to the natural divergence of X-rays as they pass through the body. Similar to parallel beam geometry, measurements are taken at various angles as the source and detectors rotate together. The detectors then record the attenuation of X-rays along fan-shaped straight paths through the body.

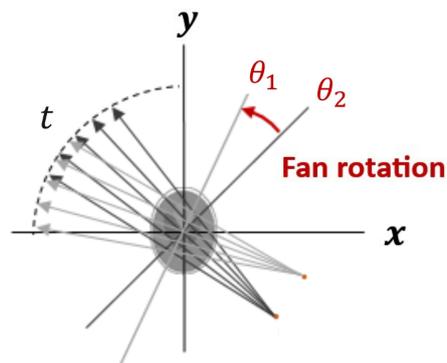


Figure A.3: 2D Fan Beam Geometry. Image from [this link](#)¹.

¹<https://it.mathworks.com/help/images/ref/fan2para.html>

Fluorescence microscopy imaging

In the second part of this thesis, we often refer to fluorescence microscopy imaging, which is our chosen application of interest for off-the-grid methods, namely for BLASSO ($L^2 - |\cdot|$) detailed in Chapter 5 and for $(\tilde{\mathcal{D}}_{KL} - |\cdot|)$ proposed in Chapter 6 for sparse Poisson deconvolution. We sketch here some general concepts about microscopy imaging. We focus on the description of fluorescence microscopy techniques providing some useful insight. Fluorescence microscopy is a popular imaging technique that allows the study of living cells and cellular organelles. However, the resolution of images obtained by fluorescence microscopes is physically limited due to the diffraction of visible light.

In biology, the observation of very small structures, such as cells or proteins, is crucial for the study of the behaviour of a bacteria or of a disease development. In this scenario, the use of conventional optical microscopes has played a key role in better understanding the phenomena involved at this scale. Over the years, the quality of the images produced by means of optical microscopes has improved dramatically thanks to technological advances and manufacturing breakthroughs. An optical microscope contains one or a series of lenses which create an enlarged image of a sample that is placed in the *focal plane* of the lens, that is the vertical plane in which the *focal point* lies, that is the point behind the lens at which the light from a far away object is brought to. However, despite the computer-aided optical design and automated methodology utilised to produce modern lens components, glass-based microscopes are still affected by an intrinsic limit in spatial resolution, due to the well-known physical barriers imposed by diffraction of visible light. In consequence, the highest achievable resolution that can be obtained in fluorescence microscopy is governed by some fundamental physical laws that cannot be overcome by physical means. These diffraction barriers restrict the ability of the optical instrument to distinguish between two objects which are too close to each other.

Diffraction limit. Passing through an optical system, the light wavefronts are distorted, which results in the perturbation of a point source into a diffraction figure.

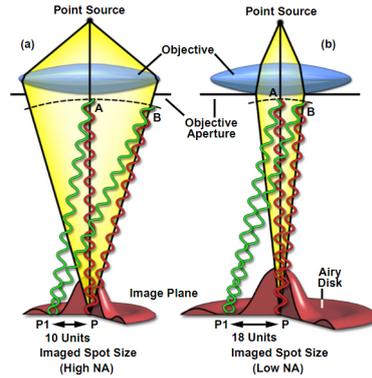


Figure B.1: Resolution limit imposed by wave nature of light. Image from [this link](#)¹.

The transmitted light wavefronts emanating from a point in the specimen plane of the microscope become diffracted at the edges of the objective aperture. Basically, the lens spreads the wavefronts to produce an image of the point source that is broadened into a diffraction pattern having a central disk of finite, but larger size than the original point. In Figure B.1, the two waves in red (green, respectively) give rise to constructive (destructive) interference, leading to the diffraction pattern. Therefore, due to diffraction, the image of a specimen never perfectly represents the real specimen, because there is a lower bound below which the microscope optical system cannot resolve structural details. In the case of a perfect optical instrument,

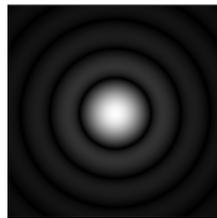


Figure B.2: Airy disk. The grayscale intensities have been adjusted to enhance the brightness of the outer rings of the diffraction pattern.

a single wavelength light from a point sources generates a diffraction pattern, shown in Figure B.2, called *Airy disk* (after Sir George B. Airy, a nineteenth century English astronomer), characteristic of this phenomenon, with radius

$$r = \frac{\lambda}{2} \text{NA}$$

where λ is the wavelength of the light and NA the characteristic numerical aperture of the microscope. The intensity profile of the Airy disk is called Point Spread Function (PSF). The ideal PSF is the diffraction light pattern emitted by a point source in the specimen and transmitted to the image plane through the objective of the microscope. It represents the impulse response of the microscope to a

¹<https://www.microscopyu.com/techniques/super-resolution/the-diffraction-barrier-in-optical-microscopy>

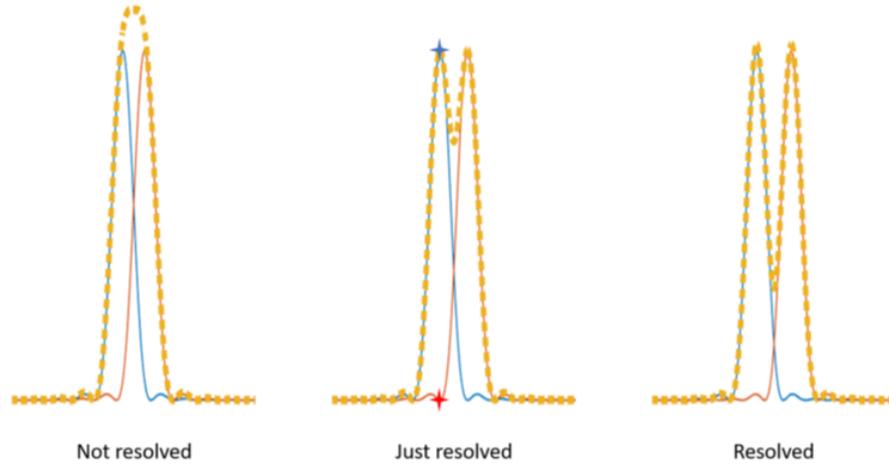


Figure B.3: Overlapping of two point sources approaching near. The Rayleigh criterion: two points are considered as just resolved when the maximum of one diffraction pattern coincides with the first minimum of the other. Image from [209].

single bright pixel (or a single fluorescent protein in fluorescence microscopy) and it is characteristic of any specific optical system. Due to the fast vanishing of the central disk of the Airy pattern, which represents the PSF, we usually estimate it by a Gaussian function. The PSF is the fundamental unit of an image in theoretical models of image formation, as it models the blur in acquisitions caused by the diffraction phenomenon. Rational alternations in objective lenses and/or changes to the numerical aperture design cannot overcome the physical laws governing optical microscopy. Thus, the highest achievable point-to-point resolution that can be obtained has severe limitations, often referred to as the *diffraction barrier*, which restricts the ability of optical instruments to distinguish between two objects. The physicist Ernst Abbe advanced the diffraction-limited resolution theory in 1873 [1], and Lord Rayleigh established a standard formula to characterise the spatial resolution of an optical device [[190]. According to Rayleigh, the resolution limit is equal to the minimum resolvable distance of two point sources, and two point sources are considered just resolved when the maximum of one diffraction pattern coincides with the first minimum of the other, as shown in Figure B.3. This results to a lateral resolution:

$$d = 0.61 \frac{\lambda}{\text{NA}},$$

where λ is the emission wavelength and NA is the numerical aperture of the objective. The numerical aperture is given by the formula $\text{NA} = n \sin \theta$ and depends on the refractive index n of the objective immersion medium and the half-angle θ of the cone of light collected by the lens [115]. The lateral resolution d is also identical to the Full-Width at Half-Maximum (FWHM) of the microscope's PSF. According to Abbe's and Rayleigh's theory, the images obtained by optical microscopes consist of many overlapping diffraction-limited spots with different intensities. The only way to

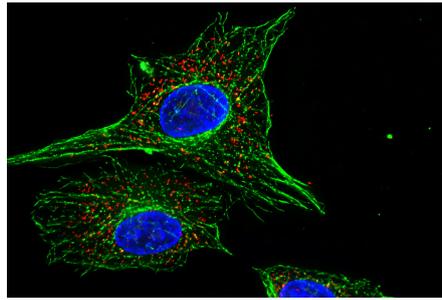


Figure B.4: Fluorescent labelled cell. Each colour represents a different emission wavelength of the emitted light. Image from [this link](#) ².

improve the lateral resolution is to minimise the size of the diffraction-limited spots either by increasing the NA of the objective lens or by decreasing the wavelength of the emitted light. However, even under ideal conditions, when imaging with visible light, the lateral resolution cannot drop under the level of 200 *nm*.

Fluorescence microscopy. A very powerful technique developed in the early 20th century is fluorescence microscopy, that allows the observation of cells, tissues and dynamics of structures with high precision in 3D and *in vivo*. This strategy takes advantage of the fluorescent nature of certain molecules which, when excited by a photon of a certain energy, return to a stable state by emitting a lower energy photon, and allows the acquisition of spatial and temporal information about objects that are either intrinsically fluorescent or coupled to extrinsic fluorescent molecules in samples [115, 198]. In practice, fluorescent molecules absorb some specific light wavelengths (or colours) and emit some others with longer wavelengths. It is thus possible to dissociate emission light from absorption light by filtering the wavelengths. The main function of fluorescence microscopes is, therefore, to deliver excitation energy to the fluorescent species in the sample and to separate the much weaker emitted fluorescence light from the brighter excitation light so that it reaches the detector and, finally, a high contrast image is generated. The light separation is usually achieved by optical filters.

Certain substances within cells are naturally fluorescent (fluorochromes), others, to benefit from this property, must be combined with a fluorescent protein. In this latter case, it is necessary to introduce fluorophores into the samples under observation, i.e. fluorescent chemical compounds that can re-emit light upon light excitation. They are physically attached to the structures of interest using certain labelling methods (such as antibody antigen pair, modification genomics, cell tracers...). The marked structures appear coloured by a different intensity, depending on the wavelength of the fluorescence emitted, as shown in Figure B.4. This is particularly useful to mark and locate specific structures of interest.

Over the past several decades, fluorescence microscopy has become an essential

²www.leica-microsystems.com/science-lab/the-fundamentals-and-history-of-fluorescence-and-quantum-dots/

tool for examining a wide variety of biological molecules, pathways and dynamics in living cells and tissues. Currently, modern and well-established techniques can resolve a variety of features in isolated cells and tissues, such as the nucleus, mitochondria, Golgi complex, cytoskeleton and endoplasmic reticulum, as in Figure 6.13 of Chapter 6. Various imaging modes in fluorescence are often used to dynamically track proteins and signal peptides, as well as for monitoring other interactions in living cells. However, the spatial resolution limitations due to light diffraction precludes the ability to resolve very small structures (such as synaptic vesicles, ribosomes...) and the investigation of biological processes close to the molecular scale, which lie beneath the limits of detection.

However, a sample generally contains many thousands of fluorescent molecules, which inevitably overlap and are hard or impossible to locate from the acquisition obtained, without any further post-processing. The localisation of the underlying molecules can be tackled by the design of a precise localisation method, capable to deal with medium-to-high density scenarios, and the challenge then is to find a good numerical method (in both time and memory complexity), capable of dealing with high-density scenarios improving the quality of the obtained images.

Different microscopy techniques can be used to enhance the visualisation and contrast of an image depending on the application. Each method has advantages and disadvantages, but all use the same fluorescence mechanism to observe a biological process. The most well-known fluorescence microscopy techniques are widefield microscopy [208], confocal microscopy [104], Total Internal Reflection Fluorescence (TIRF) microscopy [165], light sheet microscopy or Selective Plane Illumination Microscopy (SPIM) [138]. A brief overview only of widefield modality is now given, as we use it in the 3D real data experiments of Chapter 6.

Widefield microscopy is the most common fluorescence microscopy modality. It took its name because the whole (wide) field of view is illuminated. Observation and excitation of fluorescence are done from only one side of the sample and the entire sample is exposed at the same time to light, making it the simplest and fastest fluorescence modality. However, because fluorescence signals from all focal planes are detected, contrast is poor in thick samples, while thinner samples, such as adherent cells, are preferred.

Bibliography

- [1] E. Abbe. Beiträge zur theorie des mikroskops und der mikroskopischen wahrnehmung. *Arch. Mikr. Anat.*, 9(1):413–468, 1873. (Cited on page 173.)
- [2] E. Acerbi and G. Mingione. Regularity results for stationary electro-rheological fluids. *Arch. Ration. Mech. Anal.*, 164(3):213–259, 2002. (Cited on page 17.)
- [3] Y. I. Alber, R. S. Burachik, and A. N. Iusem. A proximal point method for nonsmooth convex optimization problems in Banach spaces. *Abstr. Appl. Anal.*, 2(1-2):97 – 120, 1997. (Cited on page 51.)
- [4] F. S. Alessandro Lanza, Serena Morigi and Y.-W. Wen. Image restoration with poisson–gaussian mixed noise. *Comput. Methods Biomech. Biomed. Eng. Imaging Vis.*, 2(1):12–24, 2014. (Cited on page 129.)
- [5] M. Alparone, F. Nunziata, C. Estatico, F. Lenti, and M. Migliaccio. An adaptive l^p -penalization method to enhance the spatial resolution of microwave radiometer measurements. *IEEE Trans. Geosci. Remote Sens.*, 57(9):6782–6791, 2019. (Cited on pages 17, 58, and 63.)
- [6] L. Ambrosio. Gradient flows in metric spaces and in the spaces of probability measures, and applications to fokker-planck equations with respect to log-concave measures. *Bollettino UMI*, 1(1):223–240, 2 2008. (Cited on page 50.)
- [7] V. Apidopoulos, T. A. Poggio, L. Rosasco, and S. Villa. Iterative regularization in classification via hinge loss diagonal descent. *ArXiv*, abs/2212.12675, 2022. (Cited on page 9.)
- [8] M. Bacák and U. Kohlenbach. On proximal mappings with young functions in uniformly convex banach spaces. *arXiv: Functional Analysis*, 2017. (Cited on page 49.)
- [9] H. Bauschke and P. Combettes. *Convex Analysis and Monotone Operator Theory in Hilbert Spaces*. CMS Books in Mathematics. Springer, 2017. (Cited on page 89.)

- [10] H. H. Bauschke, J. Bolte, and M. Teboulle. A descent lemma beyond lipschitz gradient continuity: First-order methods revisited and applications. *Math. Oper. Res.*, 42:330–348, 2017. (Cited on page 51.)
- [11] H. H. Bauschke and P. L. Combettes. Convex analysis and monotone operator theory in hilbert spaces. In *CMS Books in Mathematics*, 2011. (Cited on page 50.)
- [12] J. Beatty. *The Radon Transform and the Mathematics of Medical Imaging*. Honors thesis. Digital Commons @ Colby, 2012. (Cited on pages xiv, 167, and 168.)
- [13] A. Beck and M. Teboulle. A fast iterative shrinkage-thresholding algorithm for linear inverse problems. *SIAM J. Imaging Sci.*, 2(1):183–202, 2009. (Cited on pages 21 and 122.)
- [14] A. Ben-Tal, T. Margalit, and A. Nemirovski. The ordered subsets mirror descent optimization method with applications to tomography. *SIAM J. Optim.*, 12(1):79–108, 2001. (Cited on page 22.)
- [15] M. Bertero, P. Boccacci, G. Desidera, and G. Vicidomini. Image deblurring with poisson data: from cells to galaxies. *Inverse Probl.*, 25:123006, 2009. (Cited on page 129.)
- [16] M. Bertero, P. Boccacci, and V. Ruggiero. *Inverse Imaging with Poisson Data*. 2053-2563. IOP Publishing, UK, 2018. (Cited on pages 109 and 128.)
- [17] M. Bertero, P. Boccacci, G. Talenti, R. Zanella, and L. Zanni. A discrepancy principle for poisson data. *Inverse Probl.*, 26:105004, 2010. (Cited on pages 129, 155, and 156.)
- [18] M. Bertero and M. Piana. *Inverse problems in biomedical imaging: modeling and methods of solution*, pages 1–33. Springer Milan, Milano, 2006. (Cited on page 2.)
- [19] A. Beurling. Sur les intégrales de fourier absolument convergentes et leur application à une transformation fonctionnelle. In *Proc. Ninth Scand. Math. Congr.*, pages 345–366, Helsinki, Finland, 1938. (Cited on page 114.)
- [20] F. J. Beutler. The operator theory of the pseudo-inverse i. bounded operators. *J. Math. Anal. Appl.*, 10(3):451–470, 1965. (Cited on page 6.)
- [21] F. J. Beutler. The operator theory of the pseudo-inverse ii. unbounded operators with arbitrary range. *J. Math. Anal. Appl.*, 10(3):471–493, 1965. (Cited on page 6.)

-
- [22] F. Bevilacqua, A. Lanza, M. Pragliola, and F. Sgallari. Nearly exact discrepancy principle for low-count poisson image restoration. *J. Imaging*, 8(1):1, 2021. (Cited on page 156.)
- [23] F. Bevilacqua, A. Lanza, M. Pragliola, and F. Sgallari. Masked unbiased principles for parameter selection in variational image restoration under poisson noise. *Inverse Probl.*, 39(3):034002, jan 2023. (Cited on page 156.)
- [24] F. Bevilacqua, A. Lanza, M. Pragliola, and F. Sgallari. Whiteness-based parameter selection for poisson data in variational image processing. *Appl. Math. Model.*, 117:197–218, 2023. (Cited on page 156.)
- [25] K. Bieker, B. Gebken, and S. Peitz. On the treatment of optimization problems with l1 penalty terms via multiobjective continuation. *IEEE Trans. Pattern Anal. Mach. Intell.*, 44(11):7797–7808, Nov 2022. (Cited on page 147.)
- [26] I. Bisio, C. Estatico, A. Fedeli, F. Lavagetto, M. Pastorino, A. Randazzo, and A. Sciarrone. Variable-exponent lebesgue-space inversion for brain stroke microwave imaging. *IEEE Trans. Microw. Theory Tech.*, 68(5):1882–1895, 2020. (Cited on pages 17 and 58.)
- [27] B. Blaschke, H. W. Engl, W. Grever, and M. V. Klibanov. An application of tikhonov regularization to phase retrieval. *Nonlinear World*, 1996. (Cited on page 13.)
- [28] B. Blaschke-Kaltenbacher and H. W. Engl. Regularization methods for non-linear ill-posed problems with applications to phase reconstruction. In *Inverse Problems in Medical Imaging and Nondestructive Testing*, 1997. (Cited on page 13.)
- [29] J. Bolte, S. Sabach, and M. Teboulle. Proximal alternating linearized minimization for nonconvex and nonsmooth problems. *Math. Program.*, 146(1):459–494, 08 2014. (Cited on page 51.)
- [30] J. Bolte, S. Sabach, M. Teboulle, and Y. Vaisbourd. First order methods beyond convexity and lipschitz gradient continuity with applications to quadratic inverse problems. *SIAM J. Optim.*, 28(3):2131–2151, 2018. (Cited on page 51.)
- [31] B. Bonino, C. Estatico, and M. Lazzaretti. Dual descent regularization algorithms in variable exponent Lebesgue spaces for imaging. *Numer. Algorithms*, 92(6), 2023. (Cited on pages 23, 38, 39, 58, and 63.)
- [32] N. Boyd, G. Schiebinger, and B. Recht. The alternating descent conditional gradient method for sparse inverse problems. *Proc. IEEE CAMSAP 2015*, pages 57–60, 2015. (Cited on page 123.)

- [33] N. Boyd, G. Schiebinger, and B. Recht. The alternating descent conditional gradient method for sparse inverse problems. *SIAM J. Optim.*, 27(2):616–639, 2017. (Cited on pages 119 and 123.)
- [34] C. Boyer, Y. De Castro, and J. Salmon. Adapting to unknown noise level in sparse deconvolution. *Inf. Inference J. IMA*, 2017. (Cited on pages 18 and 108.)
- [35] R. I. Boş and T. Hein. Iterative regularization with a general penalty term—theory and application to l1 and tv regularization. *Inverse Probl.*, 28(10), 2012. (Cited on page 16.)
- [36] T. Brander and D. Winterrose. Variable exponent Calderón’s problem in one dimension. *Ann. Acad. Sci. Fenn. Math.*, 44(2):925–943, 2019. (Cited on page 17.)
- [37] K. Bredies. A forward–backward splitting algorithm for the minimization of non-smooth convex functionals in Banach space. *Inverse Probl.*, 25(1):015005, 2008. (Cited on pages xi, 12, 15, 78, 79, 80, 83, 86, 87, 90, 92, 93, 95, 97, 100, and 104.)
- [38] K. Bredies, K. Kunisch, and T. Pock. Total generalized variation. *SIAM J. Imaging Sci.*, 3(3):492–526, 2010. (Cited on page 12.)
- [39] K. Bredies and H. K. Pikkarainen. Inverse problems in spaces of measures. *ESAIM: COCV.*, 2013. (Cited on pages 18, 108, 114, and 123.)
- [40] L. Bregman. The relaxation method of finding the common point of convex sets and its application to the solution of problems in convex programming. *USSR Comput. Math. Math. Phys.*, 7(3):200–217, 1967. (Cited on page 51.)
- [41] P. Brianzi, F. Di Benedetto, and C. Estatico. Preconditioned iterative regularization in banach spaces. *Comput Optim Appl*, 54(2):263–282, 2013. (Cited on pages 16 and 54.)
- [42] L. M. Briceño-Arias, P. L. Combettes, J. C. Pesquet, and N. Pustelnik. Proximal algorithms for multicomponent image recovery problems. *J. Math. Imaging Vis.*, 41(1):3–22, 2011. (Cited on page 21.)
- [43] A. Buccini and L. Reichel. Generalized cross validation for ℓ_p - ℓ_q minimization. *Numer. Algorithms*, 88(4):1595–1616, December 1 2021. (Cited on page 13.)
- [44] D. Butnariu and E. Resmerita. Bregman distances, totally convex functions, and a method for solving operator equations in banach spaces. *Abstr. Appl. Anal.*, 2006. (Cited on page 50.)

-
- [45] L. Calatroni and A. Chambolle. Backtracking strategies for accelerated descent methods with smooth composite objectives. *SIAM J. Optim.*, 29(3):1772–1798, 2019. (Cited on page 21.)
- [46] D. Calvetti, S. Morigi, L. Reichel, and F. Sgallari. Tikhonov regularization and the l-curve for large discrete ill-posed problems. *J. Comput. Appl. Math.*, 123(1):423–446, 2000. Numerical Analysis 2000. Vol. III: Linear Algebra. (Cited on page 141.)
- [47] E. J. Candès, J. K. Romberg, and T. Tao. Stable signal recovery from incomplete and inaccurate measurements. *Comm. Pure Appl. Math.*, 59:1207–1223, 2006. (Cited on page 17.)
- [48] E. J. Candes and M. B. Wakin. An introduction to compressive sampling. *IEEE Signal Processing Mag.*, 25(2):21–30, 2008. (Cited on page 11.)
- [49] E. J. Candès and C. Fernandez-Granda. Towards a mathematical theory of super-resolution. *Commun. Pure Appl. Math.*, 2013. (Cited on page 18.)
- [50] Y. D. Castro, V. Duval, and R. Petit. Exact recovery of the support of piecewise constant images via total variation regularization, 2023. (Cited on page 115.)
- [51] Y. D. Castro, V. Duval, and R. Petit. Towards off-the-grid algorithms for total variation regularized inverse problems. *J. Math. Imaging Vis.*, 65(1):53–81, 2023. (Cited on page 115.)
- [52] Y. D. Castro, S. Gadat, C. Marteau, and C. Maugis-Rabusseau. SuperMix: Sparse regularization for mixtures. *Ann. Statist.*, 49(3):1779 – 1809, 2021. (Cited on page 108.)
- [53] A. Chambolle and P.-L. Lions. Image recovery via total variation minimization and related problems. *Numer. Math.*, 76:167–188, 1997. (Cited on page 12.)
- [54] A. Chambolle and T. Pock. An introduction to continuous optimization for imaging. *Acta Numer.*, 25:161–319, 2016. (Cited on page 18.)
- [55] A. Chambolle and T. Pock. An introduction to continuous optimization for imaging. *Acta Numer.*, 25:161–319, 05 2016. (Cited on page 21.)
- [56] G. H.-G. Chen and R. T. Rockafellar. Convergence rates in forward–backward splitting. *SIAM J. Optim.*, 7(2):421–444, 1997. (Cited on page 21.)
- [57] L. Chizat. Sparse optimization on measures with over-parameterized gradient descent. *Math. Program.*, 194(1–2):487–532, jul 2022. (Cited on page 118.)

- [58] L. Chizat and F. Bach. On the global convergence of gradient descent for over-parameterized models using optimal transport. In *Proc. 32nd Int. Conf. Neural Inf. Process. Syst., NIPS'18*, page 3040–3050, 2018. (Cited on pages 112 and 118.)
- [59] P. G. Ciarlet, G. Dinca, and P. Matei. Operator equations and duality mappings in sobolev spaces with variable exponents. *Chin. Ann. Math., B*, 34(5):639–666, 2013. (Cited on page 17.)
- [60] I. Cioranescu. Geometry of Banach spaces, duality mappings and nonlinear problems. *Springer*, 1990. (Cited on pages 37 and 38.)
- [61] D. L. Cohn. *Measure Theory*. Springer New York, New York, 2013. (Cited on pages 109 and 112.)
- [62] F. Colonius and K. Kunisch. Output least squares stability in elliptic systems. *Appl. Math. Opt.*, 19:33–63, 1989. (Cited on page 13.)
- [63] P. L. Combettes and J.-C. Pesquet. *Proximal Splitting Methods in Signal Processing*, pages 185–212. Springer New York, New York, NY, 2011. (Cited on page 21.)
- [64] P. L. Combettes, S. Salzo, and S. Villa. Regularized learning schemes in feature banach spaces. *Anal. Appl.*, 16(01):1–54, 2018. (Cited on page 13.)
- [65] P. L. Combettes and V. R. Wajs. Signal recovery by proximal forward-backward splitting. *Multiscale Model Simul.*, 4(4):1168–1200, 2005. (Cited on page 21.)
- [66] A. Coscia and G. Mingione. Hölder continuity of the gradient of $p(x)$ -harmonic mappings. *C. R. Acad. Sci. Paris Sér.*, 328(4):363–368, Feb. 1999. (Cited on page 17.)
- [67] J.-B. Courbot and B. Colicchio. A fast homotopy algorithm for gridless sparse recovery. *Inverse Probl.*, 37(2):025002, jan 2021. (Cited on pages 138, 139, 140, 141, 142, and 147.)
- [68] D. V. Cruz-Uribe and A. Fiorenza. Variable Lebesgue spaces. *Springer Birkhäuser Basel*, 2013. (Cited on pages 30, 31, 32, and 35.)
- [69] C. F. Dantas, E. Soubies, and C. Févotte. Safe screening for sparse regression with the Kullback-Leibler divergence. In *Proc. - ICASSP IEEE Int. Conf. Acoust.*, Toronto, Canada, June 2021. (Cited on page 133.)
- [70] I. Daubechies, M. Defrise, and C. D. Mol. An iterative thresholding algorithm for linear inverse problems with a sparsity constraint. *Commun. Pure Appl. Math.*, 57:1413–1457, 2003. (Cited on pages 11, 13, 94, 96, and 98.)

-
- [71] Y. De Castro, V. Duval, and R. Petit. Towards off-the-grid algorithms for total variation regularized inverse problems. In *Proc. SSVM 2021*, pages 553–564, Cham, 2021. Springer International Publishing. (Cited on page 115.)
- [72] Y. De Castro and F. Gamboa. Exact reconstruction using beurling minimal extrapolation. *J. Math. Anal. Appl.*, 2012. (Cited on pages 18 and 108.)
- [73] V. F. Demyanov, A. M. Rubinov, and G. M. Kranc. *Approximate methods in optimization problems*. American Elsevier Pub. Co., New York, [2. ed.]. edition, 1970. (Cited on pages 119 and 121.)
- [74] Q. Denoyelle, V. Duval, and G. Peyré. Support recovery for sparse super-resolution of positive measures. *J. Fourier Anal. Appl.*, 2017. (Cited on pages 18 and 108.)
- [75] Q. Denoyelle, V. Duval, G. Peyré, and E. Soubies. The sliding frank–wolfe algorithm and its application to super-resolution microscopy. *Inverse Probl.*, 36(1):014001, dec 2019. (Cited on pages xvii, 108, 109, 112, 118, 119, 120, 123, and 125.)
- [76] A. Diaspro, editor. *Confocal and Two-Photon Microscopy: Foundations, Applications and Advances*. Wiley, November 2001. (Cited on page 3.)
- [77] L. Diening, P. Harjulehto, P. Hästö, and M. Ruzicka. *Lebesgue and Sobolev Spaces with Variable Exponents*. Lecture Notes in Math. Springer-Verlag, Germany, 2011. (Cited on pages 30, 32, 33, 34, 35, 36, and 37.)
- [78] G. Dinca and P. Matei. Geometry of Sobolev spaces with variable exponent: smoothness and uniform convexity. *C. R. Math.*, 347(15):885–889, 2009. (Cited on pages 38, 39, and 41.)
- [79] G. Dinca and P. Matei. Geometry of Sobolve spaces with variable exponent and a generalization of the p -Laplacian. *Anal. Appl.*, 07(04):373–390, 2009. (Cited on pages 35, 38, 39, and 41.)
- [80] D. C. Dobson. Phase reconstruction via nonlinear least-squares. *Inverse Probl.*, 8:541 – 558, 1992. (Cited on page 13.)
- [81] D. L. Donoho. Compressed sensing. *IEEE Trans. Inf. Theory*, 52:1289–1306, 2006. (Cited on pages 11 and 17.)
- [82] D. L. Donoho. For most large underdetermined systems of linear equations the minimal ℓ_1 -norm solution is also the sparsest solution. *Commun. Pure Appl. Math.*, 59(6):797–829, 2006. (Cited on page 11.)
- [83] D. L. Donoho, M. Elad, and V. N. Temlyakov. Stable recovery of sparse over-complete representations in the presence of noise. *IEEE Trans. Inf. Theory*, 52:6–18, 2006. (Cited on page 17.)

- [84] C. Dossal, V. Duval, and C. Poon. Sampling the fourier transform along radial lines. *SIAM J. Numer. Anal.*, 55(6):2540–2564, 2017. (Cited on pages 108 and 112.)
- [85] V. Duval and G. Peyré. Exact support recovery for sparse spikes deconvolution. *Found. Comput. Math.*, 15:1315 – 1355, 2013. (Cited on pages 18, 108, and 111.)
- [86] V. Duval and G. Peyré. Sparse regularization on thin grids I: the Lasso. *Inverse Probl.*, 33(5):055008, May 2017. (Cited on page 109.)
- [87] V. Duval and G. Peyré. Sparse spikes super-resolution on thin grids II: the continuous basis pursuit. *Inverse Probl.*, 33(9):095008, Sept. 2017. (Cited on page 109.)
- [88] B. Eicke. Iteration methods for convexly constrained ill-posed problems in hilbert space. *Numer. Funct. Anal. Optim.*, 13(5-6):413–429, 1992. (Cited on page 8.)
- [89] I. Ekeland and R. Témam. *Convex Analysis and Variational Problems*. SIAM, Philadelphia, PA, USA, 1999. (Cited on page 116.)
- [90] Y. Eldar and G. Kutyniok. *Compressed Sensing: Theory and Applications*. Cambridge University Press, 11 2012. (Cited on pages 11 and 17.)
- [91] H. Engl, M. Hanke, and A. Neubauer. *Regularization of Inverse Problems*. Kluwer Academic Publishers Group, 1996. (Cited on pages 5, 6, 53, and 54.)
- [92] C. Estatico, A. Fedeli, M. Pastorino, and A. Randazzo. Quantitative Microwave Imaging Method in Lebesgue Spaces With Nonconstant Exponents. *IEEE Trans. Antennas Propag.*, 66(12):7282–7294, 2018. (Cited on pages 17 and 58.)
- [93] C. Estatico, S. Gratton, F. Lenti, and D. Titley-Peloquin. A conjugate gradient like method for p-norm minimization in functional spaces. *Numer. Math.*, 137:895–922, 2017. (Cited on page 16.)
- [94] M. D. Fajardo, J. Vicente-Pérez, and M. M. L. Rodríguez. Infimal convolution, c-subdifferentiability, and fenchel duality in evenly convex optimization. *Inverse Probl.*, 20(2):375–396, jul 2012. (Cited on page 131.)
- [95] X.-L. Fan and D. Zhao. The quasi-minimizer of integral functionals with $m(x)$ growth conditions. *Nonlinear Anal.*, 39:807–816, 2000. (Cited on page 17.)
- [96] C. Fernandez-Granda. Support detection in super-resolution. 2013. (Cited on pages 18, 108, and 118.)

-
- [97] M. Frank and P. Wolfe. An algorithm for quadratic programming. *Nav. Res. Logist. Q.*, 3(1-2):95–110, 1956. (Cited on pages [xvii](#), [118](#), and [120](#).)
- [98] S. Geman and D. Geman. Stochastic relaxation, gibbs distributions, and the bayesian restoration of images. *IEEE Trans. Pattern Anal. Mach. Intell.*, 6:721–741, 1984. (Cited on page [10](#).)
- [99] D. Ghilli, D. A. Lorenz, and E. Resmerita. Nonconvex flexible sparsity regularization: theory and monotone numerical schemes. *Optimization*, 0(0):1–33, 2021. (Cited on pages [13](#) and [17](#).)
- [100] N. Gomez-Navarro, A. Melero, X.-H. Li, J. Boulanger, W. Kukulski, and E. A. Miller. Cargo crowding contributes to sorting stringency in COPII vesicles. *J. Cell. Biol.*, 219(7), 05 2020. (Cited on page [156](#).)
- [101] R. Grech, T. Cassar, J. Muscat, K. P. Camilleri, S. G. Fabri, M. Zervakis, P. Xanthopoulos, V. Sakkalis, and B. Vanrumste. Review on solving the inverse problem in eeg source analysis. *J. Neuroeng. Rehabil.*, 5(1):25, November 7 2008. (Cited on page [3](#).)
- [102] C. Groetsch. *The theory of Tikhonov regularization for Fredholm equations of the first kind*. Longman Higher Education, 01 1984. (Cited on pages [7](#) and [9](#).)
- [103] C. W. Groetsch. *Generalized Inverse of Linear Operators*. Marcel Dekker, 1977. (Cited on pages [5](#) and [17](#).)
- [104] R. Gräf, J. Rietdorf, and T. Zimmermann. Live cell spinning disk microscopy. *Adv. Biochem. Eng. Biotechnol.*, 95:57–75, 2005. (Cited on page [175](#).)
- [105] W.-B. Guan and W. Song. The Generalized Forward-Backward Splitting Method for the Minimization of the Sum of Two Functions in Banach Spaces. *Numer. Funct. Anal. Optim.*, 36(7):867–886, 2015. (Cited on pages [59](#), [78](#), [79](#), [80](#), [87](#), [88](#), [90](#), [91](#), and [104](#).)
- [106] W.-B. Guan and W. Song. The forward–backward splitting method and its convergence rate for the minimization of the sum of two functions in Banach spaces. *Optim. Lett.*, 15:1735–1758, 2021. (Cited on page [78](#).)
- [107] O. Güler. New proximal point algorithms for convex minimization. *SIAM J. Optim.*, 2(4):649–664, 1992. (Cited on page [20](#).)
- [108] J. Hadamard. Sur les problèmes aux dérivés partielles et leur signification physique. *Princeton University Bulletin*, 13:49–52, 1902. (Cited on page [4](#).)
- [109] M. Hanke. A regularizing levenberg - marquardt scheme, with applications to inverse groundwater filtration problems. *Inverse Probl.*, 13(1):79, feb 1997. (Cited on page [13](#).)

- [110] M. Hanke, A. Neubauer, and O. Scherzer. A convergence analysis of the landweber iteration for nonlinear ill-posed problems. *Numer. Math.*, 72(1):21–37, 1995. (Cited on pages 8 and 53.)
- [111] Z. Harchaoui, A. Juditsky, and A. Nemirovski. Conditional gradient algorithms for norm-regularized smooth convex optimization. *Math. Program.*, 152(1-2):75–112, 2015. (Cited on page 120.)
- [112] T. Hein. Tikhonov regularization in banach spaces—improved convergence rates results. *Inverse Probl.*, 25(3):035002, jan 2009. (Cited on page 16.)
- [113] T. Hein and B. Hofmann. On the nature of ill-posedness of an inverse problem arising in option pricing. *Inverse Probl.*, 19(6):1319–1338, Dec. 2003. (Cited on page 13.)
- [114] D. A. Helmerich, G. Beliu, and M. Sauer. Multiple-labeled antibodies behave like single emitters in photoswitching buffer. *ACS Nano*, 14(10):12629–12641, 2020. (Cited on page 113.)
- [115] B. Herman. Fluorescence microscopy. *Curr. Protoc. Cell Biol.*, 00(1):4.2.1–4.2.10, 1998. (Cited on pages 173 and 174.)
- [116] G. Herman and L. Meyer. Algebraic reconstruction techniques can be made computationally efficient (positron emission tomography application). *IEEE Trans. Med. Imaging*, 12(3):600–609, 1993. (Cited on page 67.)
- [117] M. R. Hestenes and E. Stiefel. Methods of conjugate gradients for solving linear systems. *J. Res. Natl. Bur. Stand.*, 49(6):409–436, 1952. (Cited on page 8.)
- [118] J.-B. Hiriart-Urruty. A note on the legendre-fenchel transform of convex composite functions. In *Nonsmooth Mechanics and Analysis*, pages 35–46, Boston, MA, 2006. Springer US. (Cited on page 131.)
- [119] B. Hofmann, B. Kaltenbacher, C. Pöschl, and O. Scherzer. A convergence rates result for Tikhonov regularization in Banach spaces with non-smooth operators. *Inverse Probl.*, 23(3):987–1010, June 2007. (Cited on page 13.)
- [120] B. Hofmann and R. Krämer. On maximum entropy regularization for a specific inverse problem of option pricing. *J. Numer. Math.*, 13(1):41–63, 2005. (Cited on page 13.)
- [121] B. Hofmann and P. Mathé. Parameter choice in banach space regularization under variational inequalities. *Inverse Probl.*, 28(10):104006, 2012. (Cited on page 16.)

-
- [122] G. Huang, A. Lanza, S. Morigi, L. Reichel, and F. Sgallari. Majorization–minimization generalized krylov subspace methods for ℓ_p – ℓ_q optimization applied to image restoration. *BIT Numer. Math.*, 57(2):351–378, June 1 2017. (Cited on page 13.)
- [123] J. Idier. *Bayesian Approach to Inverse Problems*. John Wiley & Sons, 2013. (Cited on page 9.)
- [124] M. Jaggi. Revisiting Frank-Wolfe: Projection-free sparse convex optimization. volume 28 of *Proc. Mach. Learn. Res.*, pages 427–435, Atlanta, Georgia, USA, 17–19 Jun 2013. (Cited on pages 119, 122, and 123.)
- [125] P. Jidesh and S. H. K. Non-local total variation regularization models for image restoration. *Comput. Electr. Eng.*, 67:114–133, 2018. (Cited on page 12.)
- [126] B. Jin and Z. Kereta. On the convergence of stochastic gradient descent for linear inverse problems in banach spaces. *SIAM J. Imaging Sci.*, 16(2):671–705, 2023. (Cited on pages 67, 68, 69, and 71.)
- [127] Q. Jin, X. Lu, and L. Zhang. Stochastic mirror descent method for linear ill-posed problems in banach spaces. *Inverse Probl.*, 39(6):065010, may 2023. (Cited on page 67.)
- [128] J. S. Jørgensen and et al. Core Imaging Library - Part I: a versatile Python framework for tomographic imaging. *Phil. Trans. R. Soc. A*, 2021. (Cited on pages 64 and 70.)
- [129] A. Kaltenbach and M. Růžička. Variable exponent bochner–lebesgue spaces with symmetric gradient structure. *J. Math. Anal. Appl.*, 503(2):125355, 2021. (Cited on page 17.)
- [130] K. S. Kazimierski, P. Maass, and R. Strehlow. Norm sensitivity of sparsity regularization with respect to p. *Inverse Probl.*, 28(10):104009, 2012. (Cited on page 13.)
- [131] J. B. Keller. Inverse problems. *Am. Math. Mon.*, 83(2):107–118, 1976. (Cited on page 4.)
- [132] P. C. Kendall. Geo-electromagnetism. *Geophys. J. Int.*, 74:639–640, 1983. (Cited on page 2.)
- [133] M. V. Klibanov and P. E. Sacks. Phaseless inverse scattering and the phase problem in optics. *J. Math. Phys.*, 33(11):3813–3821, Nov. 1992. (Cited on page 13.)
- [134] M. V. Klibanov, P. E. Sacks, and A. V. Tikhonravov. The phase retrieval problem. *Inverse Probl.*, 11(1):1–28, Feb. 1995. (Cited on page 13.)

- [135] K. Kuwae. Resolvent flows for convex functionals and p-harmonic maps. *Analysis and Geometry in Metric Spaces*, 3(1):000010151520150004, 2015. (Cited on page 50.)
- [136] L. Landweber. An iteration formula for fredholm integral equations of the first kind. *Am. J. Math.*, 73(3):615–624, 1951. (Cited on page 8.)
- [137] A. Lanza, S. Morigi, L. Reichel, and F. Sgallari. A generalized krylov subspace method for ℓ_p - ℓ_q minimization. *SIAM J. Sci. Comput.*, 37(5):S30–S50, 2015. (Cited on page 13.)
- [138] Z. Lavagnino, F. C. Zanacchi, and A. Diaspro. *Selective Plane Illumination Microscopy (SPIM)*, pages 2307–2308. Springer Berlin Heidelberg, Berlin, Heidelberg, 2013. (Cited on page 175.)
- [139] B. Laville, L. Blanc-Féraud, and G. Aubert. Off-the-grid charge algorithm for curve reconstruction in inverse problems. In *Proc. SSVM 2023*, pages 393–405, Cham, 2023. Springer International Publishing. (Cited on page 114.)
- [140] B. Laville, L. Blanc-Féraud, and G. Aubert. Off-the-grid curve reconstruction through divergence regularization: An extreme point result. *SIAM J. Imaging Sci.*, 16(2):867–885, 2023. (Cited on page 114.)
- [141] B. Laville, L. Blanc-Féraud, and G. Aubert. Off-the-grid variational sparse spike recovery: Methods and algorithms. *J. Imaging.*, 7(12), 2021. (Cited on pages xiii, 109, 118, and 119.)
- [142] B. Laville, L. Blanc-Féraud, and G. Aubert. Off-the-grid covariance-based super-resolution fluctuation microscopy. In *Proc. ICASSP 2022*, pages 2315–2319, 2022. (Cited on page 114.)
- [143] M. Lazzaretti, L. Calatroni, and C. Estatico. A continuous, non-convex and sparse super-resolution approach for fluorescence microscopy data with poisson noise. In *Proc. ICCSA 2021*, pages 80–86, 2021. (Cited on pages 12 and 23.)
- [144] M. Lazzaretti, L. Calatroni, and C. Estatico. Weighted-celo sparse regularisation for molecule localisation in super-resolution microscopy with poisson data. In *IEEE ISBI 2021*, pages 1751–1754, 2021. (Cited on pages 2, 12, 23, and 128.)
- [145] M. Lazzaretti, L. Calatroni, and C. Estatico. Modular-proximal gradient algorithms in variable exponent Lebesgue spaces. *SIAM J. Sci. Comput.*, 44(6), 2022. (Cited on pages 23, 42, 59, 63, 78, and 80.)
- [146] M. Lazzaretti, C. Estatico, A. Melero, and L. Calatroni. Off-the-grid regularisation for poisson inverse problems, 2024. Submitted to *Comput. Optim. Appl.* (Cited on page 23.)

-
- [147] M. Lazzaretti, Z. Kereta, C. Estatico, and L. Calatroni. Stochastic gradient descent for linear inverse problems in variable exponent lebesgue spaces. In *Proc. SSVM 2023*, pages 457–470, Cham, 2023. Springer International Publishing. (Cited on pages [12](#), [23](#), [43](#), and [59](#).)
- [148] M. Lazzaretti, S. Rebegoldi, L. Calatroni, and C. Estatico. A scaled and adaptive fista algorithm for signal-dependent sparse image super-resolution problems. In *Proc. SSVM 2021*, pages 242–253, Cham, 2021. Springer International Publishing. (Cited on pages [21](#) and [23](#).)
- [149] T. T. Le, R. Chartrand, and T. J. Asaki. A variational approach to reconstructing images corrupted by poisson noise. *J. Math. Imaging Vis.*, 27(3):257–263, 2007. (Cited on page [12](#).)
- [150] S. Lefkimmiatis, J. P. Ward, and M. Unser. Hessian schatten-norm regularization for linear inverse problems. *IEEE Trans. Image Process.*, 22(5):1873–1888, 2013. (Cited on page [12](#).)
- [151] A. Leitão and M. Marques Alves. On landweber–kaczmarz methods for regularizing systems of ill-posed equations in banach spaces. *Inverse Probl.*, 28(10):104008, 2012. (Cited on page [16](#).)
- [152] C. Lemarechal. Cauchy and the gradient method. *Doc Math Extra*, 251–254:10, 2012. (Cited on page [8](#).)
- [153] F. Li, Z. Li, and L. Pi. Variable exponent functionals in image restoration. *Appl. Math. Comput.*, 216(3):870–882, 2010. (Cited on page [17](#).)
- [154] J. Li, Z. Shen, R. Yin, and X. Zhang. A reweighted l^2 method for image restoration with poisson and mixed poisson-gaussian noise. *Inverse Probl. Imaging.*, 9(3):875–894, 2015. (Cited on page [129](#).)
- [155] P. L. Lions and B. Mercier. Splitting algorithms for the sum of two nonlinear operators. *SIAM J. Numer. Anal.*, 16(6):964–979, 1979. (Cited on page [21](#).)
- [156] D. A. Lorenz and E. Resmerita. Flexible sparse regularization. *Inverse Probl.*, 33(1):014002, 2016. (Cited on page [17](#).)
- [157] W. Luxemburg. *Banach Function Spaces*. PhD thesis, T.U. Delft, 1955. (Cited on page [32](#).)
- [158] J. Mairal and B. Yu. Complexity analysis of the lasso regularization path. In *Proc. ICML’12*, page 1835–1842, Madison, WI, USA, 2012. Omnipress. (Cited on pages [144](#) and [146](#).)
- [159] P. Matei. A nonlinear eigenvalue problem for the generalized Laplacian on Sobolev spaces with variable exponent. *Rom. J. Math. Comput. Sci.*, 2(2):70–82, 2012. (Cited on page [38](#).)

- [160] P. Matei. On the Fréchet differentiability of Luxemburg norm in the sequence spaces $\ell^{p(n)}$ with variable exponents. *Rom. J. Math. Comput. Sci.*, 4(2):167–179, 2014. (Cited on page 38.)
- [161] S. Matet, L. Rosasco, S. Villa, and B. L. Vu. Don’t relax: early stopping for convex regularization. *ArXiv*, abs/1707.05422, 2017. (Cited on page 9.)
- [162] A. Meaney. X-ray dataset of walnut (2020-11-11). Zenodo, Nov. 2020. (Cited on pages 70 and 73.)
- [163] G. J. Minty. Monotone (nonlinear) operators in Hilbert space. *Duke Math. J.*, 29(3):341 – 346, 1962. (Cited on page 19.)
- [164] C. Molinari and M. Massias. Iterative regularization for low complexity regularizers, 1 2022. (Cited on page 9.)
- [165] P. P. Mondal and A. Diaspro. Simultaneous multilayer scanning and detection for multiphoton fluorescence microscopy. *Sci. Rep.*, 1:149, 2011. (Cited on page 175.)
- [166] M. C. Mukkamala, J. Fadili, and P. Ochs. Global convergence of model function based bregman proximal minimization algorithms. *J. Glob. Optim.*, 83(4):753–781, 08 2022. (Cited on page 51.)
- [167] F. Natterer. The mathematics of computerized tomography. *John Wiley*, 1986. (Cited on page 9.)
- [168] F. Natterer and F. Wübbeling. *Mathematical Methods in Image Reconstruction*. SIAM, 2001. (Cited on page 2.)
- [169] E. Ndiaye, O. Fercoq, A. Gramfort, and J. Salmon. Gap safe screening rules for sparsity enforcing penalties. *J. Mach. Learn. Res.*, 18(1):4671–4703, jan 2017. (Cited on page 133.)
- [170] D. Needell, R. Zhao, and A. Zouzias. Randomized block Kaczmarz method with projection for solving least squares. *Linear Algebra Appl.*, 484:322–343, 2015. (Cited on page 67.)
- [171] A. S. Nemirovsky, D. B. Yudin, and E. R. Dawson. *Problem Complexity and Method Efficiency in Optimization*. Wiley-Interscience Series in Discrete Mathematics. John Wiley Chichester, Chichester, 1983. (Cited on page 22.)
- [172] Y. Nesterov. A method for solving the convex programming problem with convergence rate $o(1/k^2)$. *Proc. USSR Acad. Sci.*, 269:543–547, 1983. (Cited on page 20.)
- [173] Y. Nesterov. Introductory lectures on convex optimization. *Appl. Optimizat.*, 2004. (Cited on pages 20 and 21.)

-
- [174] A. Neubauer. Tikhonov-regularization of ill-posed linear operator equations on closed convex sets. *J. Approx. Theory*, 53(3):304–320, 1988. (Cited on page 8.)
- [175] A. Neubauer. On enhanced convergence rates for tikhonov regularization of nonlinear ill-posed problems in banach spaces. *Inverse Probl.*, 25(6):065009, may 2009. (Cited on page 16.)
- [176] M. Nikolova. Minimizers of cost-functions involving nonsmooth data-fidelity terms. application to the processing of outliers. *SIAM J. Numer. Anal.*, 40(3):965–994, 2003. (Cited on page 12.)
- [177] M. Nikolova. A variational approach to remove outliers and impulse noise. *J. Math. Imaging Vision*, 20, 2004. (Cited on pages 12 and 98.)
- [178] M. Osborne, B. Presnell, and B. Turlach. A new approach to variable selection in least squares problems. *IMA J. Numer. Anal.*, 20(3):389–403, 07 2000. (Cited on page 140.)
- [179] M. R. Osborne, B. Presnell, and B. A. Turlach. On the lasso and its dual. *J. Comput. Graph. Stat.*, 9(2):319–337, 2000. (Cited on pages 140, 141, and 144.)
- [180] N. Parikh and S. Boyd. Proximal algorithms. *Found. Trends Optim.*, 1(3):127–239, 2013. (Cited on page 19.)
- [181] N. Parikh and S. Boyd. Proximal algorithms. *Found. Trends Optim.*, 1(3):127–239, jan 2014. (Cited on page 21.)
- [182] M. Piana and M. Bertero. Projected Landweber method and preconditioning. *Inverse Probl.*, 13(2):441–463, apr 1997. (Cited on pages 8 and 53.)
- [183] C. Poon. An introduction to sparse spikes recovery via the blasso. Lecture notes, 2019. (Cited on page 109.)
- [184] C. Poon and G. Peyré. Multi-dimensional sparse super-resolution. *SIAM J. Math. Anal.*, 2019. (Cited on pages 18 and 108.)
- [185] N. Pustelnik, C. Chaux, and J.-C. Pesquet. Parallel proximal algorithm for image restoration using hybrid regularization. *IEEE Trans. Image Process.*, 20(9):2450–2462, 2011. (Cited on page 21.)
- [186] R. Ramlau. A modified landweber method for inverse problems. *Numer. Funct. Anal. Optim.*, 20(1-2):79–98, 1999. (Cited on page 8.)
- [187] R. Ramlau. Regularization properties of tikhonov regularization with sparsity constraints. *ETNA*, 30:54–74, 2008. (Cited on page 11.)

- [188] R. Ramlau and E. Resmerita. Convergence rates for regularization with sparsity constraints. *ETNA*, 2010. (Cited on page 11.)
- [189] R. Ramlau and M. Rosensteiner. An efficient solution to the atmospheric turbulence tomography problem using kaczmarz iteration. *Inverse Probl.*, 28(9):095004, jul 2012. (Cited on page 67.)
- [190] L. Rayleigh. Investigations in optics, with special reference to the spectroscope. *London Edinburgh Philos. Mag. & J. Sci.*, 8(49):261–274, 1879. (Cited on page 173.)
- [191] S. Rebegoldi and L. Calatroni. Scaled, inexact, and adaptive generalized fista for strongly convex optimization. *SIAM J. Optim.*, 32(3):2428–2459, 2022. (Cited on page 21.)
- [192] S. Reich, T. Bonesky, K. S. Kazimierski, P. Maass, F. Schöpfer, and T. Schuster. Minimization of tikhonov functionals in banach spaces. *Abstr. Appl. Anal.*, 2008:192679, 2008. (Cited on page 16.)
- [193] H. Robbins and S. Monro. A Stochastic Approximation Method. *Ann. Math. Stat.*, 22(3):400 – 407, 1951. (Cited on pages 66 and 67.)
- [194] R. T. Rockafellar. *Convex Analysis*. Princeton University Press, Princeton, NJ, USA, 1970. (Cited on pages 19, 52, and 116.)
- [195] L. I. Rudin, S. Osher, and E. Fatemi. Nonlinear total variation based noise removal algorithms. *Phys. D: Nonlinear Phenom.*, 60(1):259–268, 1992. (Cited on pages 11 and 12.)
- [196] W. Rudin. *Real and complex analysis, 3rd ed.* McGraw-Hill, Inc., USA, 1987. (Cited on page 109.)
- [197] S. Salzo and S. Villa. Inexact and accelerated proximal point algorithms. *J. Convex Anal.*, 19:1167–1192, 01 2012. (Cited on page 20.)
- [198] M. J. Sanderson, I. Smith, I. Parker, and M. D. Bootman. Fluorescence microscopy. *Cold Spring Harb. Protoc.*, 2014(10), 2014. (Cited on page 174.)
- [199] A. Sawatzky. *(Nonlocal) Total Variation in Medical Imaging*. PhD thesis, University of Münster, 2011. (Cited on page 12.)
- [200] A. Sawatzky, C. Brune, T. Kösters, F. Wübbeling, and M. Burger. *EM-TV Methods for Inverse Problems with Poisson Noise*, pages 71–142. Springer International Publishing, Cham, 2013. (Cited on page 129.)
- [201] A. Sawatzky, C. Brune, J. Müller, and M. Burger. Total variation processing of images with poisson statistics. In *Comput. Anal. Images Patterns*. Springer, 2009. (Cited on page 12.)

-
- [202] O. Scherzer, M. Grasmair, H. Grossauer, M. Haltmeier, and F. Lenzen. *Variational Methods in Imaging*. Springer Publishing Company, Incorporated, 1 edition, 2008. (Cited on page 11.)
- [203] T. Schuster. A stable inversion scheme for the laplace transform using arbitrarily distributed data scanning points. *J. Numer. Math.*, 11(3):263–287, 2003. (Cited on page 13.)
- [204] T. Schuster, B. Kaltenbacher, B. Hofmann, and K. S. Kazimierski. Regularization methods in Banach spaces. *De Gruyter*, 2012. (Cited on pages xi, 12, 13, 14, 15, 17, 51, 55, 56, 88, 91, 92, 95, and 104.)
- [205] T. Schuster, J. Plöger, and A. K. Louis. Depth-resolved residual stress evaluation from x-ray diffraction measurement data using the approximate inverse method. *IJMR*, 94:934 – 937, 2003. (Cited on page 13.)
- [206] T. Schuster, A. Rieder, and F. Schöpfer. The approximate inverse in action: Iv. semi-discrete equations in a banach space setting. *Inverse Probl.*, 28(10):104001, 2012. (Cited on page 16.)
- [207] F. Schöpfer, A. K. Louis, and T. Schuster. Nonlinear iterative methods for linear ill-posed problems in Banach spaces. *Inverse Probl.*, 22(1):311–329, 2006. (Cited on pages xvii, 16, 51, and 54.)
- [208] A. Stemmer, M. Beck, and R. Fiolka. Widefield fluorescence microscopy with extended resolution. *Histochem. Cell Biol.*, 130(5):807–817, 2008. (Cited on page 175.)
- [209] V. Stergiopoulou. *Learning and optimization for 3D super-resolution in fluorescence microscopy*. Theses, Université Côte d’Azur, Jan. 2023. (Cited on pages xiv and 173.)
- [210] A. M. Stuart. Inverse problems: A bayesian perspective. *Acta Numer.*, 19:451–559, 2010. (Cited on page 9.)
- [211] R. Tibshirani. Regression shrinkage and selection via the lasso. *J. R. Stat. Soc. Series B Stat. Methodol.*, 58(1):267–288, 1996. (Cited on page 18.)
- [212] A. N. Tikhonov. On the stability of inverse problems. *Proc. USSR Acad. Sci.*, 39:195–198, 1943. (Cited on page 7.)
- [213] A. N. Tikhonov. Solution of incorrectly formulated problems and the regularization method. *Soviet Math. Dokl.*, 4:1035–1038, 1963. (Cited on page 7.)
- [214] R. Twyman, S. Arridge, and et al. An investigation of stochastic variance reduction algorithms for relative difference penalized 3D PET image reconstruction. *IEEE Trans. Med. Imaging*, 42(1):29–41, 2023. (Cited on page 67.)

- [215] T. Valkonen. Proximal methods for point source localisation. *J. Nonlinear Anal. Optim.*, Volume 4, Sept. 2023. (Cited on page 118.)
- [216] S. Villa, S. Matet, B. C. Vũ, and L. Rosasco. Implicit regularization with strongly convex bias: Stability and acceleration. *Anal. Appl.*, 21(01):165–191, 2023. (Cited on page 9.)
- [217] S. Villa, S. Salzo, L. Baldassarre, and A. Verri. Accelerated and inexact forward-backward algorithms. *SIAM J. Optim.*, 23(3):1607–1633, 2013. (Cited on page 21.)
- [218] V. Vshyukova, A. Meleshko, N. Mihal, and O. Aleinikova. Changing of ikzf1 genotype during philadelphia-negative precursor-b acute lymphoblastic leukemia progression: A short clinical report. *Leuk. Res. Rep.*, 6, 06 2016. (Cited on page 113.)
- [219] J. Wang and J. Ye. Two-layer feature reduction for sparse-group lasso via decomposition of convex sets. In *Proc. 27th Int. Conf. Neural Inf. Process. Syst. - Vol. 2, NIPS'14*, page 2132–2140, Cambridge, MA, USA, 2014. MIT Press. (Cited on page 133.)
- [220] S. H. Ward and G. W. Hohmann. Electromagnetic Theory for Geophysical Applications. In *Electromagnetic Methods in Applied Geophysics: Volume 1, Theory*. Society of Exploration Geophysicists, 01 1987. (Cited on page 2.)
- [221] H. Zou and T. Hastie. Regularization and variable selection via the elastic net. *J. R. Stat. Soc. Series B Stat. Methodol.*, 67(2):301–320, 2005. (Cited on page 12.)