

Optimal Recursive Estimation Techniques for Dynamic Medical Image Reconstruction

by

Joe Qranfal

Eng., Université Joseph Fourier, Grenoble, France

M.Sc., Université Joseph Fourier, Grenoble, France

A THESIS SUBMITTED IN PARTIAL FULFILLMENT
OF THE REQUIREMENTS FOR THE DEGREE OF

DOCTOR OF PHILOSOPHY

in the Department

of

Mathematics

© Joe Qranfal 2009

SIMON FRASER UNIVERSITY

Spring 2009

All rights reserved. This work may not be
reproduced in whole or in part, by photocopy
or other means, without the permission of the author.

APPROVAL

Name: Joe Qranfal
Degree: Doctor of Philosophy
Title of Thesis: Optimal Recursive Estimation Techniques
for Dynamic Medical Image Reconstruction

Examining Committee: Dr. JF Williams
Chair

Dr. Manfred Trummer
Senior Supervisor

Dr. Anna Celler
Supervisory Committee

Dr. Steve Ruuth
Supervisory Committee

Dr. Torsten Möller
Internal Examiner

Dr. Heinz Bauschke
External Examiner

Date of Defense: December 1, 2008



SIMON FRASER UNIVERSITY
LIBRARY

Declaration of Partial Copyright Licence

The author, whose copyright is declared on the title page of this work, has granted to Simon Fraser University the right to lend this thesis, project or extended essay to users of the Simon Fraser University Library, and to make partial or single copies only for such users or in response to a request from the library of any other university, or other educational institution, on its own behalf or for one of its users.

The author has further granted permission to Simon Fraser University to keep or make a digital copy for use in its circulating collection (currently available to the public at the "Institutional Repository" link of the SFU Library website <www.lib.sfu.ca> at: <<http://ir.lib.sfu.ca/handle/1892/112>>) and, without changing the content, to translate the thesis/project or extended essays, if technically possible, to any medium or format for the purpose of preservation of the digital work.

The author has further agreed that permission for multiple copying of this work for scholarly purposes may be granted by either the author or the Dean of Graduate Studies.

It is understood that copying or publication of this work for financial gain shall not be allowed without the author's written permission.

Permission for public performance, or limited permission for private scholarly use, of any multimedia materials forming part of this work, may have been granted by the author. This information may be found on the separately catalogued multimedia material and in the signed Partial Copyright Licence.

While licensing SFU to permit the above uses, the author retains copyright in the thesis, project or extended essays, including the right to change the work for subsequent purposes, including editing and publishing the work in whole or in part, and licensing other parties, as the author may desire.

The original Partial Copyright Licence attesting to these terms, and signed by this author, may be found in the original bound copy of this work, retained in the Simon Fraser University Archive.

Simon Fraser University Library
Burnaby, BC, Canada

Abstract

The focus of this thesis is to mathematically model and solve the inverse problem of reconstructing a dynamic medical image. We use a stochastic approach based on a Markov process to model the problem. We introduce a novel proximal approach based on a Bregman projection, and we apply it during the Kalman filter algorithm to ensure positivity and spatial regularization. We do not postulate precise a-priori information about the underlying dynamics of the physical process. We establish theoretical properties of our solution, and we test our method for the case of image reconstruction in time-dependent single photon emission computed tomography (SPECT). Static SPECT reconstruction algorithms assume that the activity does not vary in time. In many situations, however, physicians are interested in the dynamics of the underlying physiological process. For example, rate of uptake or wash-out of the pharmaceutical tracer will provide functional diagnosis information. Thus arises the need to explore time-varying SPECT which, mathematically, is an ill-posed inverse problem.

In this thesis, we investigate a projected Kalman reconstruction approach to estimate the dynamic activity. We give a brief overview of imaging in general, and medical imaging in particular. We then describe some important aspects of SPECT imaging, one of the two main imaging modalities in nuclear medicine.

We formulate a linear state-space model of the problem, and we introduce the optimal recursive Kalman filter (KF) and smoother. However, the Kalman output

image is unidentifiable because of the presence of negative components in the activity. Setting negative values of the activity to zero or taking their absolute value does not lead to an acceptable solution. We thus incorporate a proximal method to induce a positive estimator, and then we establish a number of mathematical and statistical properties of our estimator.

While KF does a temporal smoothing, it does not include a spatial regularization. We present spatial regularization schemes, and we give a detailed description on how to implement them. We provide numerical results to corroborate the effectiveness of our reconstruction method and to confirm our theoretical results. Finally, we summarize our findings and state directions of present and future work.

Dedication

To the memory of both my parents

To my wife Salwa

and our children

Samy, Ramy,

Sabrina, and Farris

Acknowledgements

First and foremost, all praise are due to the Almighty God. I am amazed at God's optimized ways of sorting things out gracefully under all constraints and circumstances. My research project, and indeed, my life experiences are no exceptions. I express my gratitude towards Simon Fraser University and the Department of Mathematics for their financial support and for their welcoming atmosphere to work. My supervisors Dr. Manfred Trummer and Dr. Anna Celler have introduced me to the field of medical imaging and kept providing me with insightful comments; I am really grateful to both of them. I appreciate the time and energy of my committee members Dr. JF Williams, Dr. Steve Ruuth, Dr. Torsten Möller. The work of Dr. Heinz Bauschke, a Canada Research Chair (CRC) in convex analysis and optimization, and of Dr. Charles Byrne as well as discussion with them were of great inspiration, I thank both of them. I was lucky to meet Dr. Germain Tanoh as a colleague who then became a wonderful friend of mine; I thank him for all his support. I encountered some people during my thesis research who directly or indirectly helped me out; I give thanks to them all. I treasure all my wife's efforts while standing by me; for without her my Ph.D. journey would have been less enjoyable.

Contents

Approval	ii
Abstract	iii
Dedication	v
Acknowledgements	vi
Contents	vii
List of Tables	x
List of Figures	xi
1 Dynamic Medical Image Reconstruction	1
1.1 Nuclear Medicine Imaging	2
1.2 Image Reconstruction	6
1.3 Dynamic SPECT Image Reconstruction	10
1.4 Thesis Overview	12
2 Temporal Recursive Filtering	14
2.1 Kalman Filter in SPECT	16
2.2 Stochastic Modeling	17
2.2.1 Problem setting	17
2.3 Kalman Filtering	20

2.4	Drawbacks of the Kalman Filter	23
3	Generalized Proximal Method	25
3.1	Convex Optimization	26
3.2	Basic Notions	26
3.3	Bregman Distance	30
3.4	Bregman Projection	33
3.5	Euclidian Proximal Approach	35
3.6	Generalized Proximal Operators	37
3.7	Nonnegative Minimization Method	37
3.8	Bregman-Legendre Functions	40
3.9	Boltzmann-Shannon Entropy	42
3.10	Nonnegativity Minimization Algorithm	43
4	Parameter Estimation	45
4.1	Overview	46
4.2	Parameter Estimation Properties	50
4.3	Positive Image	52
4.4	Properties of the Projected KF Estimator	53
4.4.1	Maximum Likelihood	54
4.4.2	Consistency	56
4.4.3	Unbiasedness	57
4.4.4	Optimality	60
4.5	Summary	63
5	Spatial Smoothness	64
5.1	Nonnegativity Constraint	64
5.2	Iterative Solver	65
5.3	Tikhonov Regularization	66

5.4	Energy Function and Approximation	71
5.5	Hölder Filter	77
5.6	Segmentation Regularization	78
6	Numerical Experiments	80
6.1	Procedure	80
6.2	Simulation	81
6.3	First Tests	83
6.3.1	Results	88
6.4	Positive Kalman	93
6.5	Tikhonov Regularization	99
6.5.1	Augmented	101
6.6	Median Regularization	103
6.7	Hölder Filter Regularization	106
6.8	Comparing with Improved dEM Algorithm	109
6.9	Various Sized Phantoms	113
6.10	Spatial Regularization via Segmentation	119
7	Conclusions	121
	Bibliography	123

List of Tables

6.1	CPU times and percentages.	99
6.2	Deviation error τ_{avg} of several reconstructions.	109
6.3	CPU time of several reconstructions	109
6.4	Tuning parameters.	113
6.5	Ratio of data to unknowns.	115
6.6	CPU time of various sizes.	116
6.7	Average deviation τ_{avg} of different sizes.	116

List of Figures

1.1	Photons radiating from the region of interest in SPECT.	3
3.1	Non convex function.	27
3.2	Non smooth convex function.	27
3.3	Non strictly convex function.	27
3.4	Strictly convex smooth function.	27
3.5	Convex set.	29
3.6	Non convex set.	29
3.7	Conjugate function.	30
3.8	Bregman distance.	32
3.9	Orthogonal projection.	34
3.10	Bregman projection.	34
3.11	Kullback-Leibler function.	41
4.1	Oblique projection	55
4.2	Projection inequality	62
5.1	First order neighborhood configuration.	72

5.2	log cosh function	73
5.3	$\tanh(\eta v)/(2v)$ function	74
6.1	Simulated annulus with its different ROI and their TACs	82
6.2	Sinogram	84
6.3	Reconstruction with different heads	86
6.4	Without and with positivity reconstructed images and TACs	88
6.5	All reconstructed TACs without positivity	89
6.6	All reconstructed TACs with positivity	90
6.7	Noisy sinogram	91
6.8	Reconstructed TACs and images with and without noise in data.	92
6.9	τ function with noisy data.	92
6.10	τ function for data without noise.	93
6.11	τ_{avg} as a function of $\log(\sigma_Q^2)$	94
6.12	τ_{avg} vs number of iterations	95
6.13	Image and TACs at different number of iterations	95
6.14	Images at different number of iterations	96
6.15	Nonnegative images at various times	96
6.16	Averaged TACs for each region with positivity	98
6.17	Tikhonov regularized images at various times	100
6.18	Augmented vs Bregman Tikhonov regularization	103
6.19	Median regularized images	104
6.20	Three reconstructed images	105

6.21	Mean-before and mean-after regularization	106
6.22	Median-before and median-after regularization	107
6.23	All reconstructed images	108
6.24	Kidney phantom	110
6.25	Improved dEM and Projected Kalman kidney images	111
6.26	Improved dEM and projected Kalman kidney TACs	112
6.27	Full annulus phantom	113
6.28	Size 31×31 digital phantom	114
6.29	Size 45×43 digital phantom	115
6.30	Size 64×64 digital phantom	116
6.31	Size 31×31 with Tikhonov regularization	117
6.32	Size 45×43 with Tikhonov regularization	117
6.33	Size 31×31 with Median regularization	118
6.34	Size 45×43 with Median regularization	118
6.35	Regularization via segmentation images	119
6.36	Regularization via segmentation TACs	120

Chapter 1

Dynamic Medical Image Reconstruction

Images and visualization have become increasingly important in many areas of science and technology. Advances in hardware and software have allowed computerized image processing to become a standard tool in many scientific applications. Applications include, although are not restricted to, remote sensing when imaging the earth or a planet, electrical resistivity imaging as a geophysical method to image the underground, SONAR as a sound navigation ranging imaging, Radar imaging, and medical imaging [32].

An image lifetime goes through three stages, acquisition, processing, and interpretation. In this thesis we focus on image reconstruction, which belongs to the second stage, processing. Processing involves contrast enhancing, denoising, deblurring, inpainting, coregistration, segmentation, or reconstruction of an image. To reconstruct an image, we first obtain or record data via some form of sensing. We need then to link the data to the object we aim to image via physical models, usually in a simplified form. We are often faced with additional challenges such as data scarcity and even the data available to us may be distorted or tarnished by noise. Mathematical tools

such as Fourier transforms, matrix theory, optimization techniques, probability and statistics are essential in imaging. Image analysis, for instance, uses mathematical tools such as geometry of curves and surfaces and bounded variations (BV) functions. It utilizes also elements from Bayesian statistical inference, wavelets, and iterative optimization techniques. Image modeling employs tools such as distributions, L^p and Sobolev $H^n(\Omega)$ spaces, Markov and Gibbs random fields and processes, level sets, PDEs, and Mumford-Shah free boundary [32, 38, 92].

Medicine is an area where image science has supplied many essential tools for diagnosis, treatment, and intervention to improve health care. Since the discovery and application of X-Rays by Wilhelm Röntgen (1895), medical imaging has grown and much improved [37]. This includes multi-dimensional modalities such as X-ray computed tomography, ultrasonic imaging (1942), and magnetic resonance imaging or MRI (1973). Different disciplines including physics, engineering computer science, mathematics, and medicine have contributed to the evolution of medical imaging. Since the introduction of the Gamma camera by Hal Anger (1957), nuclear medicine now provides two imaging modalities, positron emission tomography (PET) and single photon emission computed tomography (SPECT). The word “tomography” is derived from the Greek $\tau\omicron\mu\omicron\sigma$ (tomos), to cut or slice, and $\gamma\rho\alpha\phi\omicron\varsigma$ (graphy), to write. Some image modalities focus on revealing structures (X-ray tomography), others on revealing function (functional MRI or fMRI, PET, SPECT).

1.1 Nuclear Medicine Imaging

Emission tomography and transmission tomography are the two main families of medical imaging. Nuclear medicine belongs to the first family, where the radiation source is inside the patient, while X-ray belongs to the second family, where the radiation source is outside the patient. In case of SPECT and PET, a judiciously designed

chemical tagged with a radioactive tracer is administered to the patient, usually by intravenous injection. A radioactive tracer, or just radiotracer, is a special molecule carrying an unstable isotope named radionuclide. It is chosen to amass in a targeted organ or region of the body, the heart or the brain for instance. The unstable isotope emits γ rays (photons) and an external device, the gamma camera, detects the radioactivity originating from the patient, see Figure 1.1. As a consequence, nuclear medicine measures the function or metabolism of a targeted organ or region of the body [26]. MRI or ultrasound imaging assess functionality of an organ as well, however, nuclear medicine has the advantage that it provides much higher SNR (signal to noise ratio) than any other modality [92]. We get a considerable amount of data in both transmission and emission medical tomography although the dosage to the patient is limited in both cases, thus decreasing the SNR.

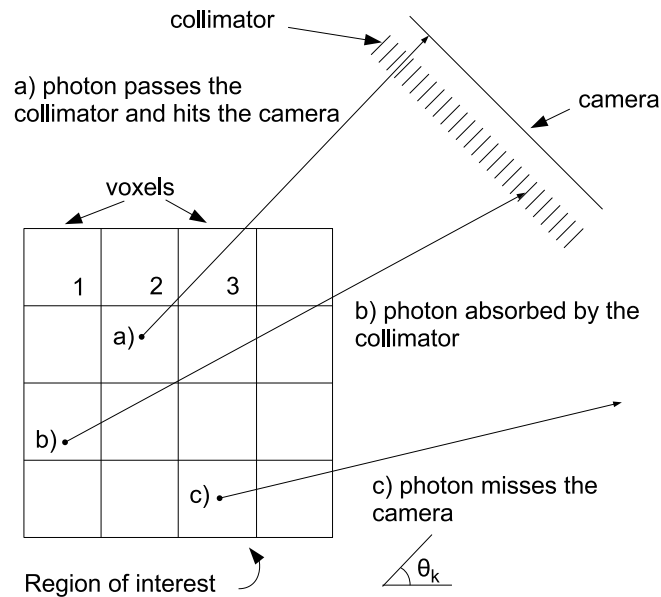


Figure 1.1: Photons radiating from the region of interest in SPECT.

We might have a camera with one or more heads. The camera could be stationary

or rotating so that we obtain one or several angles of view, respectively. Data from a rotating camera form what we call a sinogram which is a set of binned data representing one slice through all the projections, see Figure 6.2 for an example. At any angle of view, the data is the projection of a 3D activity distribution onto the 2D camera detectors. In transmission tomography, X-ray tomography for instance, the source's position is known so that every collected photon yields exact information about the projection line, that is the line joining the detection incidence and the source. This is not the case in emission tomography where the activity distribution, the emitting source of photons, is the unknown that we aim to solve for. To extract information about the spatial distribution of the activity, a collimator is used [38]. In SPECT, the collimator is a thick plate with hexagonal holes. This permits the detection of only these photons that are almost parallel to the axis of the holes, see Figure 1.1. Since most photons are absorbed by the collimator, the price is that this operation hurts the sensitivity. Even with knowledge of the direction of the emitted photons, information about the depth at which the disintegrations take place is lacking. It is as if we are asked to determine uniquely two or more numbers given only their sum, thus the need of a large number of angles of view around the patient. Indeed, the overlap that we get in the data informs us about the position of the object under scrutiny.

Another difference between X-ray imaging and nuclear medicine imaging is that we have a smaller amount of detected photons in the latter [92]. Thus noise is a big player in this process and must be taken care of; hence stochastic modeling becomes very useful. We cannot predict the exact moment at which the atom will disintegrate, however we know the decay probability as

$$\frac{dN(t)}{dt} = -\beta N(t)$$

where β is an isotope dependent decay constant and $N(t)$ is the activity at time t . This differential equation gives away the expected value as

$$N(t) = N(t_0)e^{-\beta(t-t_0)}$$

Since the process is statistical, we usually measure only an approximation to the true $N(t)$. The Poisson distribution describes the probability of rare events occurring in a fixed period of time but having very many opportunities to happen. It is a limiting case of the binomial distribution. It can be shown that it models well the probability of detection of photons in SPECT [45, 67, 86] and improves the reconstruction compared to the assumption of normality of the measurement noise. Let k be a nonnegative integer and suppose λ is the expected number of occurrences during a given interval of time, then the probability that there are exactly k detected photons is equal to

$$p_\lambda(k) = \frac{e^{-\lambda}\lambda^k}{k!} \quad (1.1.1)$$

The Gaussian distribution, with the same mean λ and standard deviation $\sqrt{\lambda}$, approximates very well this Poissonian probability when λ becomes large ($\lambda \geq 8$ in practice). However, with small values of λ , the Poissonian distribution is not even approximately symmetric.

Line integrals through the body could model the data in SPECT. Photons that travel through the body could be absorbed or could be scattered, thus not reaching the detector or deviating their trajectory from a straight path, respectively. Line integrals can not model these phenomena and more sophisticated physical models are called upon. A stochastic approach is one of them. Instead of calculating line integrals, we estimate rather the probability that an emitted photon from a certain location reaches a certain position in the camera taking into consideration factors such as attenuation, scattering, and blurring [86]. On the basis of the emitted photons that are registered by the gamma camera, the distribution of the radioactivity within the body is estimated; this mathematical operation is called *reconstruction*. A medical condition could then be indicated by an unusual lack of the radionuclide or more than usual presence of it in the targeted region.

Interested readers may wish to consult [26, 32, 37, 38, 92] for more about imaging in general and medical imaging in particular including nuclear medicine, especially

the SPECT modality, and the mathematics involved in all of these.

1.2 Image Reconstruction

The standard process is to discretize the problem by dividing the region(s) under inspection into small parts called pixels in 2D or voxels in 3D. If we denote by x the vector of these pixel/voxel values and by y the projections, then we are basically solving a large system of linear equations, $Cx = y$. Here, C is the system matrix, the mathematical model of how x and y are connected. This system does not need to be square. We could affiliate to it two square systems, $C^T Cx = C^T y$, referred to as the normal equation, and $CC^T z = y$. Direct methods exist to solve this system with a square matrix, however, they cannot be applied in our setting [26]. Not only are they computationally intensive, but a solution might not exist or may not be unique. In addition, the unknown x has many components, in the thousands and even in the millions, since it is the vectorization of a discrete approximation of a 2 or 3 dimensional continuous function. Thus direct methods are not widely applied and, for example, the use of Gauss elimination is precluded. The matrix C for our application of interest is often if not always a rectangular matrix, with fewer rows than columns. Thus we might have infinitely many solutions or no solutions at all. We have only noisy measured data, thereby we do not want to find an exact solution even if such a solution exists. An exact solution is not useful because it is the consequence of over-fitting the result to noisy data. We prefer a regularized approximate solution.

The mathematics behind image reconstruction has seen early developments. Radon published his famous paper in 1917 [80] linking any function f to its sinogram $p(s, \theta)$. It is called the Radon transform operator that associates with each line its line integral. It saw its first application in computed tomography (CT) in 1970 and later on in nuclear medicine imaging. The matrix C , in this case, is modeled as a transform

operator. The projected data can be represented as follows,

$$g(s, \theta) = \mathcal{R}f(r)$$

where \mathcal{R} is the Radon transform operator and f is the activity function, usually used in a discrete form as a vector x ; θ is the position angle of the camera head and s is the distance between the camera and the object. We need to find $f(r)$ given the sinogram $g(s, \theta)$; that is we must invert the Radon transform as follow

$$f(r) = \mathcal{R}^{-1}\{g(s, \theta)\}$$

A less used way to solve for f is by direct inversion of the Radon operator using the Hilbert transform operator [51]. A very important operator in signal processing is the Fourier transform. There exists a useful relationship between the Fourier transform and the Radon transform known as central slice theorem or the projection theorem. In the 2-dimensional case for instance, the projection-slice theorem states that the Fourier transform of the projection of a 2-dimensional function $f(r)$ onto a line is equal to a slice through the origin of the 2-dimensional Fourier transform of that function which is parallel to the projection line. As a consequence f can be reconstructed by performing first a 1D Fourier transform of g at different angles, followed by a 2D inverse Fourier transform. The Filtered back-projection algorithm (FBP) is the earliest reconstruction method tailored to medical imaging. It is the most widely used analytic reconstruction method in tomography and is a numerical implementation of the inversion formula of the Radon transform. Roughly speaking, the backprojection step produces a blurred version of f and the filtering step aims to reduce this blur; some useful references are in [71, 76]. If we had the line integrals for every line, then we could use that data to determine the activity, with some precision. In practice though, we have available only finitely many noisy line integral values, so using the central slice theorem gives us only approximate solutions. Nonetheless, the main advantage of this approach is its time efficiency.

In SPECT, the ideal setting is that photons do not interact with matter, however they do in practice. Take for instance the isotope ^{99m}Tc whose energy is 140 keV. Every 5 cm of tissue absorbs about 50% of the photons [92]. This information is required for the matrix C , thus we must correct for this effect in order to apply the central slice theorem; this is not an obvious task in SPECT. In addition, data suffers significantly from Poisson noise. Recently, an inversion formula for the attenuated radon transform was offered by Novikov [77]. Based on this inversion formula, Kuyansky [66] looked into a SPECT reconstruction.

An alternative to analytical approaches is to use iterative approaches. The maximum likelihood expectation maximization (MLEM) or just EM algorithm, introduced in 1982 [86], is the most popular algorithm; it is based on a Bayesian model. It yields automatically nonnegative solutions, a desirable feature in medical imaging which is lacking in FBP solutions. This statistical view sees the activity x_j as the expected number of emitted photons at the j^{th} location during the scanning time. Thereby the location values can be seen as parameters to be estimated. Thus the expected number of detected photons at the i^{th} detector is

$$\mathbb{E}(y_i) = \sum_j C_{ij} x_j$$

However, the actual count y_i replaces the expected count $\mathbb{E}(y_i)$. Hence we do not seek an exact solution, rather, an approximate one. We have only noisy data and as the number of iterations increases, we obtain projections closer and closer to this noisy data. This phenomenon has been observed in the EM algorithm. Thus EM is only semi-convergent; that is noise is amplified at high iteration numbers [26]. A remedy to this fallout is to stop the EM algorithm early on in the iteration. Another remedy is to ameliorate our comparison criterion between the measured data and the projections of the actual approximate. An improved criterion could be to have the projections of the actual approximate as close as possible to the measured data and the reconstructed image not being too noisy. Thereby we might introduce a prior

knowledge as a constraint into our optimization problem. This operation is called *regularization*. The prior, based on our assumption of what the true image should be like, is usually chosen to penalize noisy images. Maximizing both criteria has been done for EM and a one-step-late (OSL) algorithm was introduced by Green in 1990 [47].

We can say that there exist mainly two classes of reconstruction methods in tomography. The first class corresponds to noniterative methods that includes the analytical deterministic approaches like convolution techniques; FBP is one of these [6, 61, 76]. The second class corresponds to the iterative approaches that includes the stochastic methods based on bayesian analysis. Statistical criteria that have been utilized in devising these methods include the minimum mean squares error (MMSE), weighted least squares (WLS), maximum entropy (ME), maximum likelihood (ML), and maximum a posteriori (MAP). The algebraic reconstruction technique (ART) and multiplicative ART (MART) were first introduced by Gordon et al [46] (1970); although it was noticed later on that ART is but a particular case of Kaczmarz's algorithm [59] introduced earlier (1937). ART and MART are two examples of this second class of iterative methods. Other iterative approaches have been used such as Gauss-Seidel, conjugate gradient (CG), EM and OSEM (ordered subsets EM), a faster variant of EM [54]. Even though CG can be quicker than EM, it is still slow for the large problems that we face in image reconstruction [99]. It is also harder to find pre-conditioners for this ill-conditioned image reconstruction problem except when we deal with extremely structured matrices, which is not usually the case. Time-consuming convergence of EM has restricted its use clinically although it produces acceptable reconstructions early in the iterative process. The authors of [65] proposed a Newton and Conjugate Gradient based algorithms for both PET and SPECT using Bayesian estimators, while Jonhson et al [58] experimented with a nonlinear optimization method.

Standard SPECT imaging, where FBP and EM for instance are applied, assumes the distribution of the radioactive tracer is stationary or with little temporal change. The typical acquisition time in SPECT is 20 minutes thereby we get inconsistent data if the image does change with time.

1.3 Dynamic SPECT Image Reconstruction

Nuclear medicine is interested in the dynamics of the human body's physiological processes and biochemical function [92]. It uses a tracer that is showing time-dependent physiological effects because this time-dependency is of medical interest. Assume a triple-head camera is used while taking three angular measurements at each time frame, one measurement per head, and let the region of interest and each camera be discretized into 64^3 and 64^2 locations respectively. This will give us only 3×64^2 data for 64^3 unknowns. Hence we have a very under-determined problem to solve. The goal is to reconstruct a sequence of images (movie) from very few data. This is an *inverse* problem.

An inverse problem is the task where the values of some model parameters must be obtained from the observed data. In our setting, an inverse problem consists of estimating or reconstructing a target function from a limited number of measurements that are connected to the target function. We do not observe or measure the target function, a dynamic activity in our case, directly whose values would have made up a complete set of data. Hence, we have indirect or incomplete data of the function. In SPECT imaging, the target function is the radioactivity distribution. Standard SPECT imaging assumes the distribution of the radioactive tracer is stationary. Analytical reconstruction techniques such as FBP with an eventual 3 DRP (3D reprojection method) [39] and iterative ones such as OSEM can be used in the static case of tomography. However, these approaches break down when trying to solve a dynamic

SPECT problem since they are only suited for a static activity. Subsequently, we need different reconstruction approaches for non-static images.

Time-varying or dynamic SPECT reconstruction can be mathematically modeled as an optimization problem with a huge number of variables. A problem is called *well-posed* in the sense of Hadamard [15] if it obeys three conditions, the solution of the problem is unique, exists for any data, and depends continuously on the data. A problem is *ill-posed* if it fails to satisfy any of these conditions. Given the recorded data, the solution of the reconstruction problem may not be unique or may not even exist in the case of data inconsistency. The problem is then *ill-posed*. Inverse problems are also *ill-conditioned* in practice. As a way to diminish sensitivity to noise and other modeling errors, we call on regularization. Thus even if the large system $Cx = y$ had an exact solution, which is unlikely if the data contain noise, this solution is not sought after since it is a solution that is overfitted to the noisy sinogram y . Regularization assists in curing an ill-posed problem. This ill-posedness of the reconstruction problem is further amplified by physical degradation like camera blurring, photon scattering, or attenuation. The reconstructed image has to be a tradeoff between accuracy and damping of the noise within it. Thus arises the need for fast and robust algorithms and regularization can assist to make the solution less sensitive to noise and modeling errors.

Some authors have used nonlinear least squares method to fit sum of exponentials in one or two compartments modeling in the context of dynamic SPECT [70, 72, 73]. Instability could happen that results in an inadequate reconstruction; this was mentioned by Blondel et al [17]. Meanwhile, estimating the activity straight off from the observations while mixing spline and least squares was introduced by Reutter et al [81]. Bauschke et al [13] have offered what they call a dynamic EM approach by using the activity dynamics as linear constraints. Later on, Celler et al [29] investigated reconstructing dynamic images assuming monotone activity only. They

modeled the reconstruction problem as a constrained least squares problem. Activities are not always monotone in their behavior. Take for example the liver where activity is increasing, arrives at a peak, and then drops-off.

Research on reconstruction algorithms of time-varying SPECT imaging has been relatively scarce. Farncombe [40] implemented a time-varying adaptation of the EM algorithm in the same spirit of the work developed earlier in [13]. To cope with different behaviors of dixels/doxels (dynamic pixels/voxels) within different regions, they used a “mask”. The major inconvenience is that obtaining this mask can be a challenging problem in itself since it might require prior knowledge, or introduces additional variables. The same could be said about the work of Tanoh [93] who uses linear constraints for the primal-dual algorithm he proposed to solve the inverse problem of dynamic SPECT. Some of the existing dynamic SPECT reconstruction techniques use an optimization black box solver. Limber et al [70] employ the L-BFGS-B package to solve the least squares reconstruction problem; while Blinder et al [16] use the KNITRO package [55]. The drawback of these plans of attack is the lack of flexibility and the difficulty in finding good tuning parameters. Furthermore, methods as in [13, 16, 17, 29, 40, 72, 93], need an additional pre-processing step to build their mask. The pre-processing step may fail to estimate accurately the activity uptake time, and worse, the shape of the TAC, which could be misleading. The reconstructed TACs might also not be as smooth as they should be, refer for instance to [16]. As a consequence, these methods are costly in time and prone to introduce some bias.

1.4 Thesis Overview

The focus of this thesis is to mathematically model and solve the inverse problem of reconstructing a dynamic medical image. We introduce a Kalman-based algorithm

which is a *stochastic-recursive-iterative* hybrid approach while relying on a linear state-space model of the problem. Based on a stochastic model, we employ a time-recursive scheme (the Kalman filter) to get a solution. This solution might be unidentifiable, because it has negative components that appear to have a deleterious effect on other components. Setting negative values of the activity to zero or taking their absolute value does not result in an acceptable image. We then use a novel iterative proximal approach based on a Bregman projection, and we apply it during the Kalman filter algorithm to enforce nonnegativity and add spatial regularization. We do not postulate precise a-priori information about the underlying dynamics of the physical process. We establish theoretical properties of our solution, and we test our method for the case of image reconstruction in SPECT. While KF does a temporal smoothing, it does not include a spatial regularization. We present spatial regularization schemes, and we give a detailed description on how to implement them. We apply this approach numerically to dynamic SPECT reconstructing a digital phantom. Optimization is at the core of our technique. The Kalman algorithm has been applied to dynamic SPECT only seldom and not usually for image reconstruction purposes [19, 63].

Chapter 2

Temporal Recursive Filtering

As far as we can reckon from human history, we humans have been filtering things for almost all the time. Take the basic example of water filtering. A simple way to filter out undesired contents is by using our hands to take leaves off the top of the water. Another instance is when we instinctively filter noise from our milieu. It is fortunate this way, otherwise, our state of being could have been affected drastically had we cared about the small noises around us. We automatically disregard redundant sounds, such as the one of traffic and noisy neighbors, and we concentrate on significant sounds, such as the voice of an inspiring speaker for instance. Engineers have turned to filtering in their pursuit of bettering our life. They filter out noise from electromagnetic signals to improve radio communications [87]. Thus we receive clear useful information from corrupted signals.

A recursive filter is a kind of filter which reuses one or more of its outputs as a feedback input. The Kalman filter is one known example [62]. It has been introduced almost fifty years ago and has its roots in the least squares approach, also known as regression analysis, that goes as far back as Gauss. He is usually credited with developing the fundamentals of this analysis in 1795 when he was only 18 years old, eventhough he did not publish his work until 1809. The idea of least squares

analysis was also independently formulated by the Frenchman Adrien-Marie Legendre in 1805; the term “least squares method” is a direct translation from the French “*méthode des moindres carrés*”. Kalman published his famous paper [62] in 1960, describing a recursive solution to the discrete-data linear filtering problem; however, Thiele and Swerling had developed a similar algorithm earlier [68]. Since then, there are literally thousands of articles and dozens of books dealing with this technique; see [7, 27, 74, 87, 105] and references therein. Kalman filtering (KF) has been applied in diverse areas such as aerospace, marine navigation, nuclear power plant instrumentation, demographic modeling, electrical impedance, manufacturing, and many others. Recently, KF was applied to time-varying SPECT [63].

KF is an efficient recursive filter which estimates the state of a dynamic system from a series of incomplete and noisy measurements. The Kalman filter takes advantage of the dynamics of the activity, which regulate its time evolution, to take away the effects of noise and errors in order to get a good estimate of the activity at the present time (*filtering*), at a future time (*prediction*), or at a time in the past (*interpolation or smoothing*) [7]. KF is out of the ordinary since it is a purely time domain filter, which is more suitable to our time-varying SPECT case. Most filters, for example a low-pass filter or Wiener filter, are formulated in the frequency domain when we need to suppose that we have a stationary or periodic activity distribution; which is obviously not the case for dynamic SPECT.

KF is a temporal regularizer in its essence; however it does not take care of spatial regularization, which is a desirable qualitative feature in medical imaging. In addition, KF may produce negative activity; that is an activity with some or all of its components taking negative values. This is meaningless in medical imaging. These are two drawbacks of the KF algorithm. Our aim is to reconstruct an activity using the KF algorithm while remedying its two shortcomings.

2.1 Kalman Filter in SPECT

In applications, the most significant class of filters are the shift invariant filters. If the input is shifted in time and such a filter is implemented then the result is just shifted in time. Put differently, the response of the filter to an input does not depend on the time that the input arrives. Boulfefel et al [19] used a shift-variant KF for post-reconstruction restoration of SPECT slices. They showed that their approach performed better than shift-invariant KF, nonetheless they did not use KF to reconstruct the activity. Artemiev et al [9] established a general framework to apply KF to solve for any recursive tomographic image. To speed up the process and overcome the storage restriction, they introduced a pseudo KF. They were not concerned, for instance, with the reconstructed activity having a physical meaning, that is being nonnegative. The authors in [10, 60, 102] use an augmented system in implementing Tikhonov spatial regularization into KF. The method is an extension of Tikhonov in which the original observation model is replaced by an augmented one. This approach is expensive memory and time wise.

Lately, Kervinen et al. [63] used KF in dynamic SPECT, but they did not include spatial regularization. For the nonnegativity, they employed the Fast Non-Negativity-constrained Least Squares (FNNLS) developed earlier by Bro et al [22] in 1998. In the original article, Bro et al showed that their algorithm is 5 to 20 times faster than the NNLS (Non-Negativity-constrained Least Squares), first introduced by Lawson et al [69] in 1974. FNNLS uses an active-set-method optimization technique based on line search; which is known to be time consuming because of the Hessian and function evaluations. Van Benthem et al [100] improved FNNLS later on in 2004 via rearranging calculations, thus introducing the Fast Combinatorial NNLS (FC-NNLS). Kim et al [64] combine the strengths of gradient projection with a non-diagonal gradient scaling scheme to come up with a new algorithm in 2005, a Projected Quasi-Newton for the NNLS (PQN-NNLS). They showed that PQN-NNLS could outperform FNNLS

numerically hundreds times in certain cases. Unfortunately, their tests were mostly done on over-determined problems. To the best knowledge of the author, nothing is known regarding the performance of PQN-NNLS in the case of under-determined systems.

2.2 Stochastic Modeling

Notation. The following notation is used throughout the paper. We denote by \mathbb{R}^p and \mathbb{R}_+^p the p -dimensional Euclidean space and the nonnegative orthant, respectively. The set of all $n \times p$ matrices with real entries is denoted by $\mathbb{R}^{n \times p}$. I denotes the identity matrix; its size is always clear from the context. The operator $\text{Tr}(B)$ denotes the trace of the matrix B , which is the sum of its diagonal components. For a vector u , the Euclidean norm is denoted by $\|\cdot\|$ and u^\top denotes the transpose vector. The i^{th} component of a vector $u \in \mathbb{R}^p$ is denoted by u_i . Let x and y be two random vectors; $\mathbb{E}(x)$ and $\mathbb{E}(x|y)$ denote the expectation of x and the conditional expectation of x given y . The conditional expectation of x_k given y_1, \dots, y_s and its variance/covariance matrix $\mathbb{E}[(x_k - \hat{x}_{k|s})(x_k - \hat{x}_{k|s})^\top]$ are denoted $\hat{x}_{k|s} = \mathbb{E}(x_k|y_1, \dots, y_s)$ and $P_{k|s}$. We also refer to $\hat{x}_{k|k}$ as simply \hat{x}_k . We denote by $\text{int } C$ the interior of the set C .

2.2.1 Problem setting

We consider a physiological process where the distribution of the radioactive tracer in an organ or a specific region is time dependent. This region is divided into small parts called dynamic voxels in 3D or doxels and dynamic pixels in 2D or dixels; we also refer to them as locations. A SPECT camera, that could have one, two or three heads, is used to register the number of photons emitted by the patient. We assume three heads in our simulation, however, our method does not depend on the number of heads. The camera rotates around the patient through 180° with S stops during

a total acquisition time T , usually 20 minutes. The surface of each camera's head is composed of a set of detectors/bins. We assume that the activity is constant during the time interval of a single projection/stop, but is allowed to vary in time from one projection to the next one over the whole acquisition period. The distribution of the activity in the organ is denoted by a vector $x(t)$ whose dimension is equal to the total number of doxels/dixels. The data collected by the camera detectors at time t is denoted by a vector $y(t)$ whose dimension is equal to the total number of detectors. The vector $y(t)$ is also called measurement, projection, or observation vector. The activity in the organ is not directly observable and its dynamics is unknown. The goal is to reconstruct the emission object $x(t)$ from the measured data $y(t)$.

Let t_k , $k = 1, \dots, S$, be a sequence of acquisition times, N the total number of doxels/dixels and M the total number of bins. We denote by $x_k \in \mathbb{R}_+^N$ and $y_k \in \mathbb{R}_+^M$ the spatial distribution of the activity and the measured data at the k^{th} instant of time. The observations y_1, y_2, \dots, y_S are independent random vectors. Furthermore, each observation y_k depends on x_k only. In previous works on dynamic tomographic imaging [8, 63], a linear model is used to describe both the evolution of the activity and the link between the activity and data measurements. We describe next an optimal linear real-time reconstruction of dynamic images.

A stochastic process with the Markov property is called a Markov chain. Markov property means that the future depends only on the present since it is assumed that all past information is already fully captured in the present state. Hence future states are achieved from the present state via a probabilistic process. A random walk is a well known instance of a Markov process. It is related to a diffusion model which is the net motion of a substance from an area of high concentration to an area of low concentration [56]. The flow of a radioactivity could be seen as a diffusion model [81].

Although Bayes has proved a special case of Bayes' theorem more than 250 years

ago, the term Bayesian was coined only in the 1950s to refer to evidential interpretation of probability. This is the idea that probability should be viewed as a subjective degree of belief in a proposition in contrast to the frequentist view of the probability theory. To start off, we suppose that the image detection process is linear and includes additive noise. In addition, the radioactivity image is a random vector discrete in space and time. Therefore, the evolution model can be interpreted as a Markovian process. The Bayesian approach justifies the rationale for the Markovian stochastic estimation theory. This strategy results in Kalman filter based reconstruction in the time domain. We proceed as follows.

The activities sequence x_1, x_2, \dots, x_S satisfy the Markov property with unknown time varying transition/evolution matrix $A_k \in \mathbb{R}^{N \times N}$. That is

$$x_k = A_k x_{k-1} + \mu_k \quad (2.2.1)$$

where μ_k is the error random vector in modeling the transition from x_{k-1} to x_k with $\mathbb{E}(\mu_k)$ zero and covariance matrix Q_k . The argument we develop here is also valid for a homogeneous Markov sequence where the transition matrix is time invariant. Let $c_{\iota j}(k)$ be the probability that an emission from doxel/dixel ι during the acquisition time t_k will be detected in bin j . We call the projection or observation matrix the time varying matrix $C_k = [c_{\iota j}(k)]$. Finding a suitable system matrix C_k is an ongoing research topic [57] and some algorithms have been proposed [101]. We assume that this system matrix C_k is known. The observation and activity vectors are related by the following

$$y_k = C_k x_k + \nu_k \quad (2.2.2)$$

where ν_k is the noise vector in recording the data with $\mathbb{E}(\nu_k)$ zero and covariance matrix R_k . In practice, a high noise level makes the problem very challenging if no prior information is available. Each acquisition time constitutes a separate reconstruction problem. The sequential nature of the measurement is suitable for a recursive reconstruction method such as Kalman filtering. The sequence of activities x_k is a hidden

Markov chain (HMC) since it is observed through y_k using (2.2.2). We assume that we deal with white noise. Otherwise, our model should go first through the step of a pre-whitening process [49]. White noise, like white light, is a random signal with a flat power spectral density of all frequencies that constitute it.

2.3 Kalman Filtering

The Kalman filter is a Bayesian model that propagates the first and second moment of the conditional probability, namely the mean and the variance/covariance. We are solving for the vector activity x_k at every time k . We need to find a vector estimate \hat{x}_k which will be $\hat{x}_{k|k} = \mathbb{E}(x_k)$. The estimation error at the k^{th} sampling instant is computed via the covariance matrix $P_{k|k} = \mathbb{E}[(x_k - \hat{x}_{k|k})(x_k - \hat{x}_{k|k})^\top]$. The Kalman filter is an optimal estimator in the least square sense. We need then to find each $\hat{x}_{k|k}$ that minimizes $\text{Tr}(P_{k|k})$, that is the expected squared error, or $\mathbb{E}[(x_k - \hat{x}_{k|k})^\top (x_k - \hat{x}_{k|k})]$, using a subset of the projections y_1, \dots, y_k , where k indicates the index of the last available measurement.

We apply the *Kalman filter* algorithm to recursively find \hat{x}_k , which generates the conditional expectation $\hat{x}_{k|k}$. With the linear state-space equations (2.2.1) and (2.2.2), the Kalman filter propagates the activity estimate and its covariance. It proceeds in three steps, *predicting*, *correcting*, and *smoothing*; see [7, 62, 87, 105] and references therein. The first two steps, predicting and correcting, are usually referred to in the literature as the filtering portion of the algorithm.

Predicting Step Assume we have an initial estimate activity $\hat{x}_{0|0}$ and its covariance matrix $P_{0|0}$. For $k = 1, \dots, S$, compute the following steps that yield the predicted

variance $P_{k|k-1}$ and activity $\hat{x}_{k|k-1}$,

$$P_{k|k-1} = A_k P_{k-1|k-1} A_k^\top + Q_k \quad (2.3.1)$$

$$\hat{x}_{k|k-1} = A_k \hat{x}_{k-1|k-1} \quad (2.3.2)$$

where A_k is the transition/evolution matrix at time k , refer to equation 2.2.1. Recall that μ_k is the error random vector in modeling the transition from x_{k-1} to x_k with covariance matrix Q_k .

Correcting Step Then compute the following correcting steps that yield the filtered variance $P_{k|k}$ and activity $\hat{x}_{k|k}$,

$$K_k = P_{k|k-1} C_k^\top (C_k P_{k|k-1} C_k^\top + R_k)^{-1} \quad (2.3.3)$$

$$P_{k|k} = (I - K_k C_k) P_{k|k-1} (I - K_k C_k)^\top + K_k R_k K_k^\top \quad (2.3.4)$$

$$\hat{x}_{k|k} = \hat{x}_{k|k-1} + K_k (y_k - C_k \hat{x}_{k|k-1}) \quad (2.3.5)$$

where K_k is the Kalman gain. Recall that C_k is the projection or observation matrix and ν_k is the noise vector in recording the data with covariance matrix R_k , refer to equation 2.2.2.

Smoothing Step The recursive algorithm that calculates the estimate $\hat{x}_{k|S}$, where S denotes the total number of measurements, is called the *Kalman smoother*. We refer to $\hat{x}_{k|S}$ as \hat{x}_k too. To get smoothed values, run the following backward recursion for $k = S - 1, \dots, 1$:

$$J_k = P_{k|k} A_k^\top P_{k+1|k}^{-1} \quad (2.3.6)$$

$$P_{k|S} = P_{k|k} + J_k (P_{k+1|S} - P_{k+1|k}) J_k^\top \quad (2.3.7)$$

$$\hat{x}_{k|S} = \hat{x}_{k|k} + J_k (\hat{x}_{k+1|S} - \hat{x}_{k+1|k}) \quad (2.3.8)$$

where J_k is called the backward gain. This recursive procedure in KF goes in time from the evolution/transition matrix A_k , the variance matrix Q_k , and the calculated

matrix $P_{k-1|k-1}$ to predict the covariance $P_{k|k-1}$. The Kalman gain K_k is formed afterwards which helps to update the estimated activity $\hat{x}_{k|k}$, and the covariance matrix $P_{k|k}$. These in turn will contribute to perform the next step in the recursion. The difference $y_k - C_k \hat{x}_{k|k-1}$ in (2.3.5) is called the measurement *innovation* or the *residual* [105]; where $C_k \hat{x}_{k|k-1}$ is compared to the observation y_k . If this difference is zero the two will be in complete agreement; which means that our predicted estimated activity will be the corrected/filtered one as well.

On one hand, as R_k the observation covariance error draws near zero, think then of Kalman gain K_k as drawing near C_k^{-1} assuming that this inverse has some meaning. Hence the corrected estimated activity $\hat{x}_{k|k}$ will weigh the innovation, that is the observation, more heavily. On the other hand, as the predicted covariance matrix $P_{k|k-1}$ nears zero, the gain K_k nears zero and we will then weigh the innovation less heavily. Said otherwise, we remark that the weighting by the gain matrix K_k is that when we have high confidence in our measurements, that is the covariance R_k approaches zero, the present data y_k is “believed” more while the predicted estimated $\hat{x}_{k|k-1}$ is trusted less. By the same token, as we have high confidence in our evolution model, thus our covariance error $P_{k|k-1}$ approaches zero, the present data y_k is trusted less and emphasis is put more on the predicted estimated activity $\hat{x}_{k|k-1}$. The most impressive feature is that the KF technique is an on-line recursive algorithm in lieu of an off-line batch algorithm [105]. Hence, there is no need to store the past measurements in order to estimate the present activity. As a consequence, we deal with smaller size problems at each time recursion.

KF performs the conditional probability density propagation for our problem in which the system can be described through a *linear* model and in which system and measurement noises are *white* and *Gaussian*. The noise in dynamic SPECT is Poissonian (equation 1.1.1), however, the Poisson distribution is approximated by a Gaussian distribution when the number of detected photons goes to infinity, in practice

a number above 8, refer to section 1.1. Under these conditions, the mean, mode, median, and virtually any reasonable choice for an “optimal” estimate all coincide, including the maximum likelihood estimate; so there is in fact a unique “best” estimate of the value of *activity distribution* [7, 27, 74, 87, 105]. Under these three restrictions, the KF can be shown to be the best filter of any imaginable form [87, 105]. Some of the restrictions can be removed, granting a qualified optimal filter. For example, if the *Gaussian* assumption is removed, the KF can be shown to be the best (minimum error variance) filter out of the class of linear unbiased filters. That is, it is the BLUE (best linear unbiased estimator). These three assumptions are reasonable for our dynamic SPECT application, even though being Gaussian is not a must.

Derivation of the Kalman algorithm equations can be found in [7, 87] where optimization is yet another tool to obtain them. We wish to estimate the unknown x_k in both equations 2.2.1 and 2.2.2. First, the estimate \hat{x}_k must be a linear function of the data; that is it has to be of the form $\hat{x}_k = H_k^\top y_k$ where the matrix H_k is to be found. Second, the estimate \hat{x}_k is required to be unbiased or $\mathbb{E}(\hat{x}_k) = \mathbb{E}(x_k)$. Third, the matrix H_k has to be chosen as to minimize $\mathbb{E}(|\hat{x}_k - x_k|^2)$.

2.4 Drawbacks of the Kalman Filter

The Kalman algorithm has some drawbacks that we aim to correct for. First, the Kalman filter and smoother algorithm does not guarantee nonnegativity of the reconstructed activity. The update equations for \hat{x} in (2.3.5) at the correcting step and in (2.3.8) at the smoothing step can not guarantee the nonnegativity of \hat{x} . Both equations involve inversion of a matrix and subtraction of vectors, which may introduce negative elements. This is not feasible in nuclear medicine and also gives unidentifiable images. Setting negative values of the reconstructed activity to zero or taking the absolute value does not give an acceptable solution. Second, KF is based on temporal

assumptions on the object only. In particular, the smoothing of the object is done in the temporal domain. Thus, it is a very efficient temporal regularization technique for recovering an object. Since the problem we face is also spatially ill-posed, the spatial characteristics are poorly reconstructed by the Kalman filter. Another drawback of KF is the computational complexity. The cost of the algorithm grows like N^3 , refer to section 6.9.

Chapter 3

Generalized Proximal Method

We require a nonnegative approximate solution x^* while solving an *inverse problem* having the large real system $Cx = y$. We offer to use a projected temporal recursive filter. We saw in chapter 2 that the Kalman filter, a temporal recursive filter, produced the solution \hat{x} which might be nonpositive. Thus we want to find the weighted projection $x^* = \text{proj}_{\mathbb{R}_+^N}^W(\hat{x})$ with respect to a certain symmetric positive definite matrix W . Methods such as NNLS [69] and FNNLS [22] could be used to find the projection and we have covered their shortcomings in section 2.1. We propose a proximal method using a Bregman generalized distance based on a Bregman-Legendre function.

In 1967, Bregman [21] introduced what are now called *Bregman projections* with respect to generalized distances such as Kullback-Leibler distortion. We purport here a new alternative, to enforce the nonnegativity/positivity, using a *proximal* approach via an *entropy distortion*, which is an instance of the *generalized Bregman distance*. This approach not only achieves the desired nonnegativity and is very simple to implement, but it is also easily extendable to implement spatial regularization. Our method draws from *convex optimization*.

3.1 Convex Optimization

Optimization means seeking a minimum or a maximum of a real-valued function of one or several variables. When we put some restrictions on the acceptable solutions, such as being nonnegative, we have *constrained* optimization. Solving an optimization problem algebraically is usually very hard and we resort to iterative techniques. Convex optimization is a sub-class of optimization where interesting things take place. In general if a function has a minimum, this minimum is only a local minimum and not necessarily a global one. Yet for convex functions, a local minimum is also a global one. Looking for a global maximum of a concave downwards function f is equivalent to finding the global minimum of the concave upwards (convex) function $-f$.

Convex optimization plays a significant role in many applications, including ours. As in the case of general functions, we might have some constraints on the acceptable solutions and we usually use convex sets to describe the constraints. For more about convex optimization refer to Rockafellar's classic [82], where the subsequent fundamental definitions and theorems can be found.

3.2 Basic Notions

Definition 3.1.

(a) A function $f : \mathbb{R} \rightarrow \mathbb{R}$ is said to be convex if for each pair a and b with $a \neq b$, the chord with end points $(a, f(a))$ and $(b, f(b))$ is on or above the graph of $f(x)$.

That is $\forall a, b \in \mathbb{R}$ and $\lambda \in \mathbb{R}$ with $0 < \lambda < 1$

$$f(\lambda a + (1 - \lambda)b) \leq \lambda f(a) + (1 - \lambda)f(b)$$

(b) f is said to be strictly convex if it is convex and the strict inequality holds above whenever $a \neq b$.

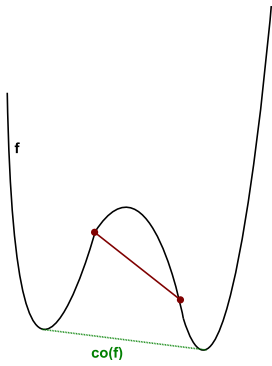


Figure 3.1: Non convex function.

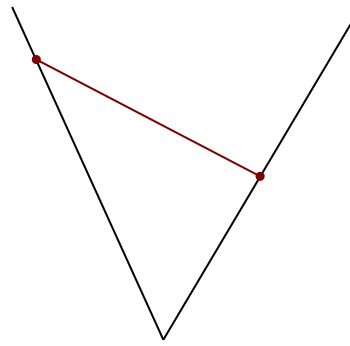


Figure 3.2: Non smooth convex function.

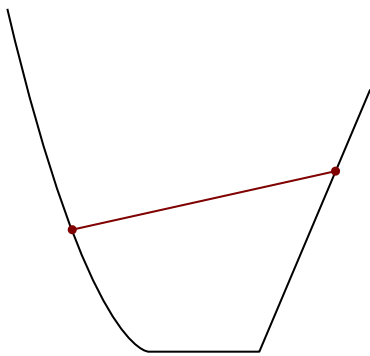


Figure 3.3: Non strictly convex function.

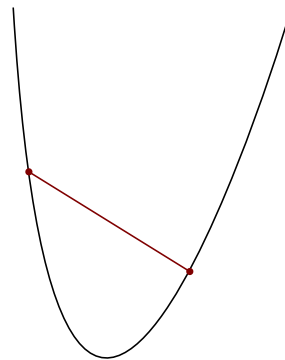


Figure 3.4: Strictly convex smooth function.

Figure 3.1 shows an example of a non convex function. It is evidenced by a chord, linking two points on the graph, that is below the graph of the function. The function $f(x) = x^4 - 2x^2 + 1$ is a simple example. Looking for a global minimum of a non convex function f is not an easy task. Thus we could work with a convex version of f called $co(f)$ or the convex hull of f . The advantage of working with $co(f)$ is that the global minimum of $co(f)$ is the same as the one of f , meanwhile $co(f)$ has the interesting property of being convex, refer to figure 3.1. Figure 3.2 exhibits a non differentiable convex function; $f(x) = |x|$ is an example. Figure 3.3 exhibits a non strictly convex function. Finally, figure 3.4 shows a strictly convex and smooth function; $f(x) = x^2$ is a very common example. The last class of functions proves to be useful for the generalized Bregman distances which is covered in the next section.

We have a similar definition of a convex function of several variables.

Definition 3.2.

- (a) *The function $f : \mathbb{R}^m \rightarrow \mathbb{R}$ is said to be convex if, for each pair of vectors a and b and for every λ such as $0 < \lambda < 1$, we have*

$$f(\lambda a + (1 - \lambda)b) \leq \lambda f(a) + (1 - \lambda)f(b)$$

- (b) *Let $C \subset \mathbb{R}^m$. C is said to be convex if $[a, b] \subset C$ whenever the vectors $a, b \in C$. That is $\forall \lambda \in [0, 1], \forall a, b \in C, \lambda a + (1 - \lambda)b \in C$*

Figure 3.5 shows an example of a convex set in \mathbb{R}^2 where it is evidenced by a chord, joining two points, as entirely belonging to the set. Figure 3.6 shows an example of a non-convex set where a chord is not entirely in the set. Constraints on x , such as when we aim to enforce nonnegativity or when we desire a regularized solution, often mean we impose on x to belong to certain convex sets. Measured data or incorporating some a priori knowledge about x usually contribute to defining these sets.

Another useful property in convex optimization is the notion of conjugate function [96].

Theorem 3.1. *A function $f : \mathbb{R} \rightarrow \mathbb{R}$ is convex if and only if there is a function $f^* : \mathbb{R} \rightarrow \mathbb{R} \cup \{+\infty\}$ such that $\forall x \in \mathbb{R}$*

$$f(x) = \sup_{x^* \in \mathbb{R}} [x^*x - f^*(x^*)]$$

The function f^* is called the *conjugate* of f . Hence f and f^* form a pair of functions such that for all $x, x^* \in \mathbb{R}$

$$f(x) + f^*(x^*) \geq xx^*$$

Figure 3.7 exhibits a geometrical interpretation of a conjugate function. For a convex function f , a line k with slope x^* and intercept $-s$ dwells nowhere above its graph is equivalent to saying that for any $z \in \mathbb{R}, x^*z - s \leq f(z)$ hence $s \geq x^*z - f(z)$. The lowest number s , which fulfills this inequality, is $\sup_{z \in \mathbb{R}} [x^*z - f(z)] = f^*(x^*)$. Moving k upwards as far as possible, we get a line $l(x^*)$ that crosses the graph of f and whose intercept is $-f^*(x^*)$. The graph of f is the envelope of the lines $l(x^*)$ ($x^* \in \mathbb{R}$) is equivalent to saying f is convex [96].

The conjugate function associated with a function of several variables f is the function

$$f^*(x^*) = \sup_z (\langle x^*, z \rangle - f(z))$$

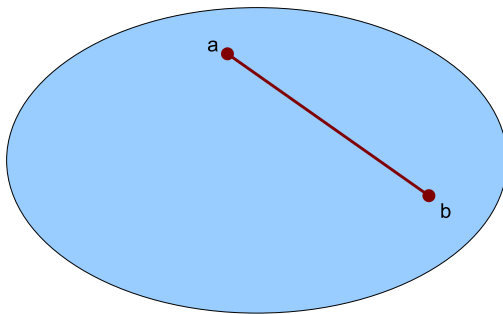


Figure 3.5: Convex set.

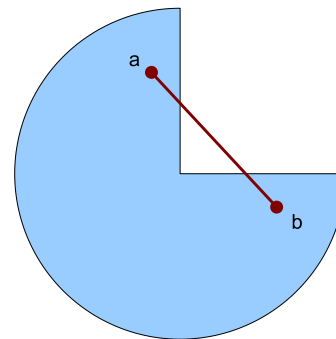


Figure 3.6: Non convex set.

where $\langle \cdot, \cdot \rangle$ is the usual scalar product between two vectors of the same dimension. We have a similar inequality as in the case of one dimension

$$f(x) + f^*(x^*) \geq \langle x, x^* \rangle$$

This is called *Fenchel's inequality*.

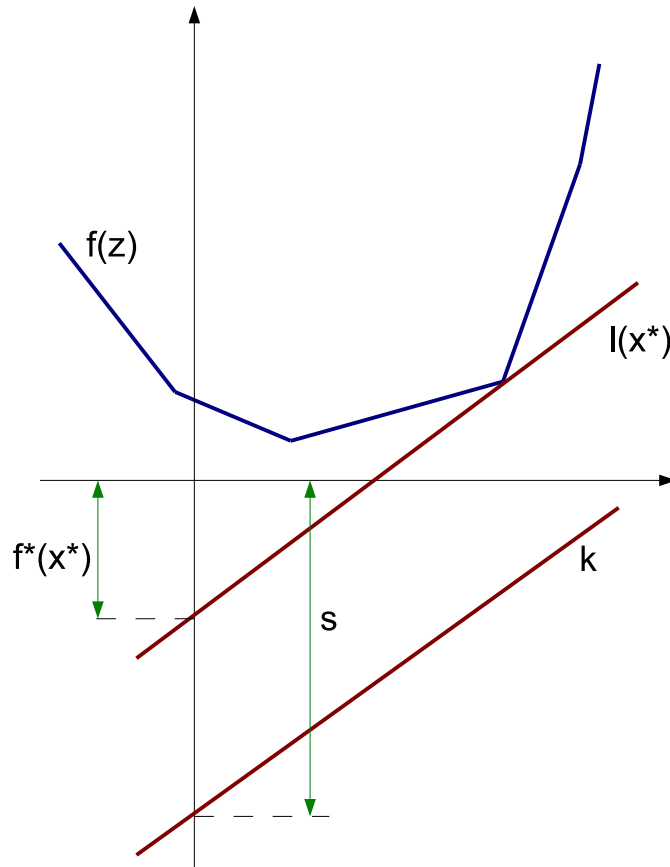


Figure 3.7: Conjugate function.

3.3 Bregman Distance

In order to define a Bregman distance, we need first additional definitions [82].

Definition 3.3.

(a) A function $f : S \subseteq \mathbb{R}^m \rightarrow [-\infty, +\infty]$ is convex if and only if its epigraph is convex. A convex function $f : S \subseteq \mathbb{R}^m \rightarrow [-\infty, +\infty]$ is proper if there is no $z \in S$ with $f(z) = -\infty$ and if there is some $z \in S$ with $f(z) < +\infty$.

(b) The essential domain of f is $\Delta_f = \{z \in S \mid f(z) < +\infty\}$.

(c) A proper convex function f is closed if it is lower semi-continuous, that is if $f(z) = \liminf f(y)$, as $y \rightarrow z$, $\forall z \in S$. Recall that

$$\liminf_{y \rightarrow z} f(y) = \inf_{\epsilon > 0} (\sup\{f(t) : t \in S \cap B(z; \epsilon) - \{z\}\})$$

where $B(z; \epsilon)$ denotes the metric ball of radius ϵ and centre z .

(d) A proper convex function f is said to be essentially smooth if

1. the interior of its domain $\text{int } \Delta_f \neq \emptyset$
2. f is differentiable on $\text{int } \Delta_f$, and
3. $\lim_{\ell \rightarrow +\infty} \|\nabla f(z^\ell)\| = +\infty$ whenever z^ℓ is a sequence in $\text{int } \Delta_f$ converging to a point on the boundary of $\text{int } \Delta_f$.

(e) The subdifferential of a function f at z is the set

$$\partial f(z) = \{z^* \mid \langle z^*, x - z \rangle \leq f(x) - f(z), \forall x\}$$

The domain of ∂f is the set $\text{dom } \partial f = \{z \mid \partial f(z) \neq \emptyset\}$

(f) A closed proper convex function f is essentially strictly convex if f is strictly convex on every convex subset of $\text{dom } \partial f$.

If f is differentiable at z , then the subdifferential is reduced to one element, namely the gradient

$$\partial f(z) = \{\nabla f(z)\} \tag{3.3.1}$$

We are ready to introduce the Bregman distance [31],

Definition 3.4. Let f be a closed proper convex function that is differentiable on $\text{int } \Delta_f \neq \emptyset$. For all $x \in \Delta_f$ and $y \in \text{int } \Delta_f$, the corresponding Bregman distance $D_f(x, y)$ is defined as follows

$$D_f(x, y) = f(x) - f(y) - \langle \nabla f(y), x - y \rangle \quad (3.3.2)$$

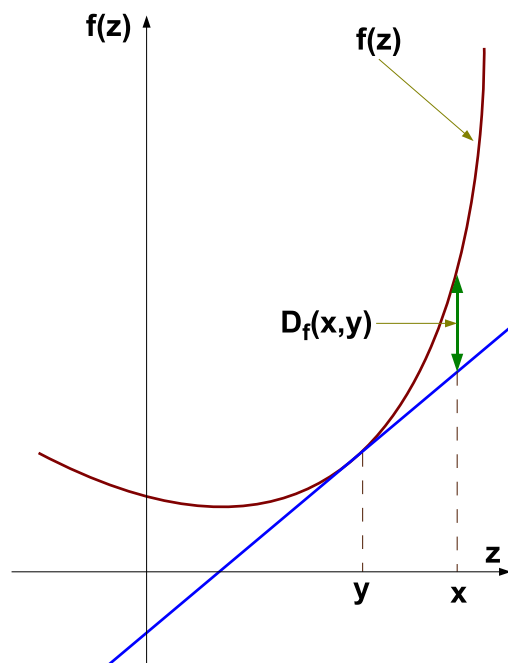


Figure 3.8: Bregman distance.

It is an easy exercise to check that $D_f(x, y) \geq 0$ is always true since a convex function is always on or above its tangents. Indeed, the minimum requirement for a distance is to be nonnegative. We could also have two vectors x and y with an infinite distance, that is $D_f(x, y) = +\infty$. In order to get the deduction $D_f(x, y) = 0$ implies $x = y$, it suffices to have f being essentially strictly convex.

Figure 3.8 shows the case of a real-valued function f of a single real variable. The Bregman distance associated with f is the vertical length between $f(x)$ and the tangent of f at y . The simplest and most widely utilized Bregman distance is

associated with the following function

$$f(x) = \langle x, x \rangle$$

so that (3.3.2) becomes

$$D_f(x, y) = \langle x, x \rangle - \langle y, y \rangle - \langle x - y, 2y \rangle = \|x - y\|^2$$

Hence we recover the usual Euclidian squared distance. We have all the ingredients to introduce the Bregman projection onto closed convex sets. For interested readers, a good reference is [31].

3.4 Bregman Projection

We associate to the Euclidian distance, a particular case of Bregman distance, the projection of a point or a vector onto a convex set referred to as the orthogonal projection. It is the shortest distance between the point and the convex set. The Bregman distance, which could be seen as the generalized distance, implies a generalized projection onto a convex set [25].

Definition 3.5. *Let $C \neq \emptyset$ be a closed convex set such that $C \cap \text{int } \Delta_f \neq \emptyset$. Choose $x \in \text{int } \Delta_f$. The Bregman projection of x onto C relatively to f is,*

$$\text{proj}_C^f(x) = \arg \min_{z \in C \cap \Delta_f} D_f(z, x)$$

It remains to prove that the above projection is well defined. This is where convex optimization plays an important role. It is useful for finding the projection into convex sets to guarantee feasibility where we commonly work with the argument of the function instead of the value of its minimum. We summarize here the discussion given in [11, 25]. If f is essentially smooth, then $\text{proj}_C^f(x)$ exists. If f is strictly convex on Δ_f , then $\text{proj}_C^f(x)$ is unique. Some restrictions on f help in uniquely finding this projection; more convex optimization notions are then in order [82].

Definition 3.6. A closed convex proper function f is Legendre, or a convex function of Legendre type, if

1. $\text{int } \Delta_f$ is nonempty
2. f is essentially smooth, and
3. f is essentially strictly convex.

It can be shown that [82],

Theorem 3.2. If f is a convex function of Legendre type then

$$\nabla f : \text{int } \Delta_f \rightarrow \text{int } \Delta_{f^*}$$

is a bijection, continuous in both directions, and $(\nabla f)^{-1} = \nabla f^*$, where the function f^* is the conjugate of f .

It entails that if $\text{int } \Delta_{f^*} = \mathbb{R}^m$ then the range of ∇f is \mathbb{R}^m and the equation $\nabla f(x) = z$ can be uniquely solved for every z in \mathbb{R}^m .

Definition 3.7. A function f is super-coercive or 1-coercive if

$$\lim_{\|x\| \rightarrow +\infty} \frac{f(x)}{\|x\|} = +\infty$$

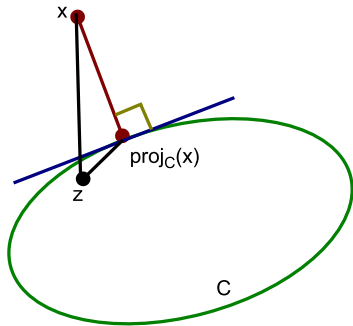


Figure 3.9: Orthogonal projection.

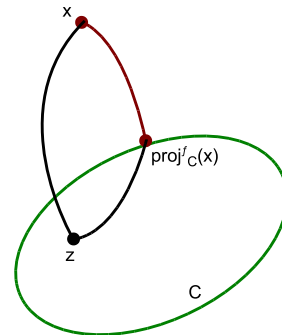


Figure 3.10: Bregman projection.

A Legendre function is super-coercive is equivalent to saying $\text{int } \Delta_{f^*} = \mathbb{R}^m$. The main consequence is that if f is Legendre, then $\text{proj}_C^f(x)$ is uniquely defined. Moreover, this projection is in $\text{int } \Delta_f$. This is cited in [31] as zone consistency.

Figure 3.9 shows the usual case, in \mathbb{R}^2 , of the orthogonal projection and figure 3.10 the one of the generalized projection. The orthogonal projection proj_C , that is using the euclidian distance, on a closed convex set C is a nonlinear operator unless C is a subspace. It has the interesting property under the form of an inequality

$$\|x - \text{proj}_C(x)\| \leq \|x - z\| \quad \forall z \in C$$

This inequality says that the shortest path from a point x to the set C is obtained by projecting x onto C and is a logical consequence of the generalized Pythagorean theorem

$$\|x - \text{proj}_C^f(x)\|^2 + \|\text{proj}_C^f(x) - z\|^2 \leq \|x - z\|^2$$

for all z in C . When f is Legendre we obtain an equivalent result when we think of a generalized distance rather as a squared than a squared root distance [25, 31]

$$D_f(z, \text{proj}_C^f(x)) + D_f(\text{proj}_C^f(x), x) \leq D_f(z, x)$$

for all z in C . This is called *Bregman's Inequality*. Thus we have

$$D_f(\text{proj}_C^f(x), x) \leq D_f(z, x) \quad \forall z \in C$$

It means that $\text{proj}_C^f(x)$ is the closest point, in the closed convex set C , to the point x when using the generalized distance D_f with respect to the Legendre function f . These generalized projections are the basis of proximal approaches in convex optimization.

3.5 Euclidian Proximal Approach

Convex optimization is very useful for finding the projection onto a convex set to ensure feasibility. Recall that we are looking for a nonnegative solution to the Kalman

filter output \hat{x} . Thus $x^* = \text{proj}_{\mathbb{R}_+^N} \hat{x}$ is the best approximation in \mathbb{R}_+^N to \hat{x} , i.e. the point of \mathbb{R}_+^N nearest to \hat{x} . Imposing the nonnegativity on the solution implies we are seeking a certain regularized solution. Convex optimization is also useful for regularization of optimization problems, thereby we have a family of optimization approaches called *proximal* minimization algorithms. In the first place, these algorithms employed Euclidian distances then later on were generalized using Bregman distances.

Jean-Jacques Moreau (*1923) is one of the founding fathers of the convex analysis discipline together with Hermann Minkowski (1864-1909), Werner Fenchel (1905-1988), and Tyrrell Rockafellar (*1935). Moreau first introduced what has become known as Moreau's proximity operator when he published in 1963 a manuscript in French [75], "Inf-convolution des fonctions numériques sur un espace vectoriel". When we are interested in minimizing a proper closed convex function $\varphi : \mathbb{R}^m \rightarrow \mathbb{R} \cup \{+\infty\}$, we use techniques such as gradient descent when φ is differentiable. However, when φ is not, the problem is harder. Moreau's approach augments the function φ in order to make it differentiable, thus we have a Moreau's operator working as an envelope of the function φ .

$$\text{env}_\varphi(x) = \inf_z \left(\varphi(z) + \frac{1}{2} \|z - x\|^2 \right) \quad (3.5.1)$$

We have a proof in [82] that the infimum in (3.5.1) is uniquely achieved at a point referred to as $z = \text{prox}_\varphi(x)$. Note that the function env_φ is differentiable where $\nabla \text{env}_\varphi(x) = x - \text{prox}_\varphi(x)$ [26]. Hence, the function φ achieves its global minimum at $z = \text{prox}_\varphi(x)$. The differentiable convex envelope env_φ of φ can be seen as the "infimal convolution" of $\varphi(z)$ and $\frac{1}{2} \|z\|^2$ [52] and both minimizers of env_φ and φ are the same. This is due to Moreau's theorem [82] which generalizes the decomposition of vectors in \mathbb{R}^m relatively to a subspace.

3.6 Generalized Proximal Operators

Several authors, Moreau among them, took on this idea of proximity operators and developed it further. We could for instance take the fraction $\frac{1}{2}$ in 3.5.1 and replace it with $\frac{1}{\alpha}$, where we may vary the parameter α in the open real interval $]0, +\infty[$; this was done and is referred to as Moreau-Yosida regularization. Moreover, when we take α as a large number, we are then solving for a Tikhonov regularized problem where we are looking for a minimum norm solution. Other authors replaced the Euclidian distance with a generalized one [31, 95]. Hence they were interested in solving problems of the form

$$E_f(\varphi, x) = \arg \min_z (\varphi(z) + D_f(z, x))$$

The above problem is well defined and admits a unique solution when we impose some restrictions on φ and f . Bauschke et al [12] characterize this solution and associate to it a left proximal operator $\overleftarrow{\text{prox}}$. The operator $E_f(\varphi, x)$ has also properties similar to Moreau's proximity operator. More about proximal approaches can be found in [25] and references therein.

3.7 Nonnegative Minimization Method

We analyze a numerical algorithm to compute a nonnegative minimum of the convex function

$$\varphi(x) = \frac{1}{2} \|\hat{x} - x\|_W^2 + \alpha \psi(x) \tag{3.7.1}$$

In the case where $\psi = 0$, many authors [22, 63, 69, 100] have proposed nonnegative minimization techniques using active set method, Newton method or quasi-Newton method involving a line search strategy which is computationally expensive. Kervinen et al [63] use Fast Non-Negativity-Constrained Least Squares (FNNLS) developed earlier by Bro et al [22]. Recently, Kim et al [64] combined the strengths of gradient

projection with a non-diagonal gradient scaling scheme to come up with a new algorithm, a Projected Quasi-Newton for the NNLS (PQN-NNLS). We propose here an alternative based on a proximal approach.

Let us start with the following unconstrained optimization problem

$$\min_{x \in \mathbb{R}^N} \varphi(x) \tag{3.7.2}$$

The convex optimization problem (3.7.2) has a unique minimizer that we denote by $x^* = \text{prox}_\psi^\alpha(\hat{x})$; that is

$$\text{prox}_\psi^\alpha(\hat{x}) = \arg \min_{x \in \mathbb{R}^N} \left(\frac{1}{2} \|\hat{x} - x\|_W^2 + \alpha \psi(x) \right)$$

If $W = I$, prox_ψ^α is the Moreau's proximity operator [34] of index $\alpha \in]0, +\infty[$ of the function ψ . These operators generalize the projection onto convex sets. In the particular case when ψ is the indicator function of a convex set C , $\psi(v) = \chi_C(v)$ where it is zero if v is in the closed convex set C and $+\infty$ otherwise, prox_ψ^α is the weighted projection of \hat{x} onto the set C and the orthogonal one if $W = I$. When the set C is the nonnegative orthant \mathbb{R}_+^N , choosing $\hat{x}^+ = \max\{\hat{x}, 0\}$ as a projection of \hat{x} may give a good result and this is commonly used [15]. However, due to the weighted norm, such approach is not recommended.

Our approach could be referred to as a generalized proximal method. Instead of the original minimization problem, we are rather solving an augmented problem via an iterative algorithm. The name iterative comes from the Latin word *iterum*, which means “again”. In ancient times, Archimedes (c.250 BC) was using an iterative procedure when he was getting better approximations of the lower and upper bounds of the area of a circle by repeating the process of inscribing and circumscribing the circle with more regular polygons. The idea behind iterative algorithms is to find a solution by consecutive estimates. The limit of these successive approximates would converge, in the best case, to the solution we desire.

The generalized proximal method is an iterative approach for minimizing the convex and differentiable function φ over the closure of the essential domain $\overline{\Delta_f}$ of a second convex function $f : \mathbb{R}_+^N \rightarrow \mathbb{R}$ while we presume that such a minimizer x^* exists. We first augment the function $\varphi(x)$ with a generalized distance with respect to a function h , which has the same domain as f

$$\Phi(x, z) = \varphi(x) + \frac{1}{\gamma} D_h(x, z) \quad (3.7.3)$$

where $\gamma > 0$, and its significance would be seen shortly. When $\psi = 0$ in equation 3.7.1, we are just enforcing the nonnegativity, otherwise, we are asking for a spatial regularized nonnegative solution with α as a regularization parameter. From now on, we set $C = \mathbb{R}_+^N$, the nonnegative orthant. We are then more interested in the convex set $\text{int } C$ than in h itself. This leads us to limit our choices of h to functions that have $\text{int } C$ as their essential domain.

We start first with $x^0 \in \text{int } C$. Having found the iterate x^ℓ , we are then concerned with the following minimization problem

$$x^{\ell+1} = \arg \min_x \Phi(x, x^\ell) \quad (3.7.4)$$

We could characterize the next iterate $x^{\ell+1}$ by taking the sub-gradient of the right side of (3.7.3)

$$\partial h(x^\ell) - \partial h(x^{\ell+1}) \in \gamma \partial \varphi(x^{\ell+1})$$

We are interested more in differentiable functions φ in (3.7.1) and using the equality (3.3.1) we get

$$\nabla h(x^{\ell+1}) = \nabla h(x^\ell) - \gamma \nabla \varphi(x^{\ell+1}) \quad (3.7.5)$$

It is not apparent how we could isolate $x^{\ell+1}$ using (3.7.5). We are not interested in any specific function h , we are rather interested in its domain C . The idea is to let $h = f - \gamma \varphi$. The function f is differentiable and has the same domain C as h , but has the advantage in making calculating $x^{\ell+1}$ more obvious. With this choice of f ,

equation 3.7.5 reduces to the following algorithm

$$\nabla f(x^{\ell+1}) = \nabla f(x^\ell) - \gamma \nabla \varphi(x^\ell) \quad (3.7.6)$$

where the parameter $\gamma > 0$ is chosen so that the function $h(x) = f(x) - \gamma\varphi(x)$ is convex. C. Byrne [24] provides a similar algorithm when φ is a sum of convex functions with the particular case of $\gamma = 1$ when φ is a single convex function. Both functions f and φ , in the examples he took to illustrate his point, belong to the Kullback-Leibler divergence family; refer to section 3.9 about this family of functions. Our functions f and φ are, however, a mixture of a quadratic and a Kullback-Leibler, respectively. The function $f - \varphi$ will never be convex globally. Thus the introduction of the parameter γ to make $h = f - \gamma\varphi$ a convex function at least locally. C. Byrne [24] established a convergence theorem of his algorithm, which he refers to as IPA (*interior point algorithm*), by restricting the function f to the class of Bregman-Legendre functions.

3.8 Bregman-Legendre Functions

Censor and Lent [30] introduced Bregman functions in 1981 inspired by the key paper of Bregman [21] (1967). Neither a Bregman function alone nor a Legendre function alone possess the necessary requirements to establish convergence theorems while a function which is both Bregman and Legendre is too restrictive. Later on in 1995, Bauschke and Borwein [11] studied the Bregman projection method within the powerful framework of convex analysis. New perspectives open up as the rich class of Bregman-Legendre functions is brought in. Those functions are Legendre with some extra, but not to the extent of being Bregman as well. It is an in-between class of functions; it is included in the class of Legendre functions and it includes the class of Bregman and Legendre functions.

Definition 3.8. *A Legendre function f is a Bregman-Legendre function if:*

BL1: $\forall x \in \Delta_f$, $D_f(x, \cdot)$ is coercive

BL2: if $x \in \Delta_f$ and $x \notin \text{int } \Delta_f$, $\forall \ell$, $y^\ell \in \text{int } \Delta_f$ with $y^\ell \rightarrow y \in \text{bd}(\Delta_f)$ (boundary of Δ_f) and if $D_f(x, y^\ell)$ remains bounded, then $D_f(y, y^\ell) \rightarrow 0$, so that $y \in \Delta_f$

BL3: if $x^\ell, y^\ell \in \text{int } \Delta_f$, with $x^\ell \rightarrow x$, $y^\ell \rightarrow y$, $x, y \in \Delta_f$ but $\notin \text{int } \Delta_f$, and if $D_f(x^\ell, y^\ell) \rightarrow 0$ then $x = y$

Bregman-Legendre functions provide the proper context for the discussion of Bregman distances and Bregman projections onto convex sets, mainly [11]

Proposition 3.1. *Suppose f is a Bregman-Legendre function.*

If $D_f(x, y^\ell) \rightarrow 0 \forall y^\ell \in \text{int } \Delta_f$, then $y^\ell \rightarrow x$

This leads to the following convergence theorem [24] of the update formula (3.7.6). Recall that we choose $C = \overline{\Delta_f} = \mathbb{R}_+^N$.

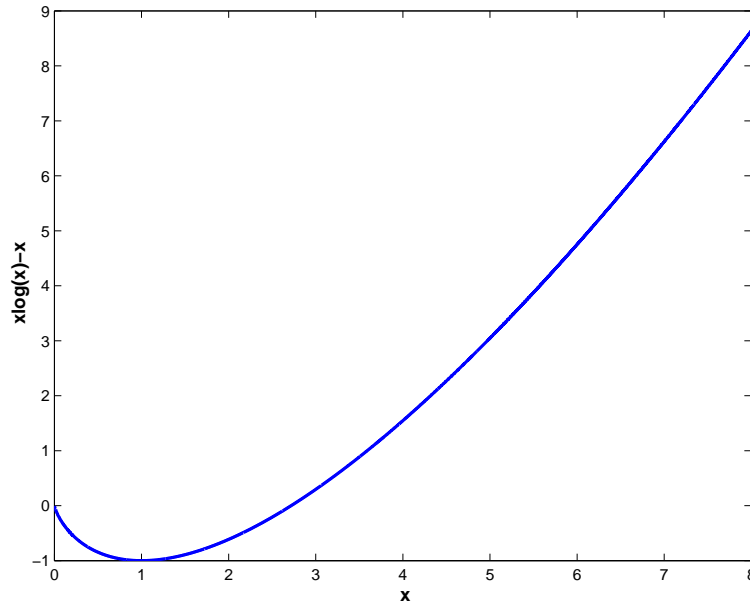


Figure 3.11: Kullback-Leibler function.

Theorem 3.3. *Let f be a Bregman-Legendre function. For any starting vector $x^0 \in \text{int } C$, the sequence x^ℓ in (3.7.6) converges to a minimizer of the function $\varphi + \chi_C$,*

where x is in the set C (the closure of the essential domain of f), assuming that such minimizers exist.

We know now that the IPA (3.7.6) converges to a unique solution, however, it is not evident how we can obtain explicitly $x^{\ell+1}$. The choice of the function of the generalized distance is very crucial indeed. We need a Bregman-Legendre function such that its domain is the nonnegative orthant. The next section introduces a typical example.

3.9 Boltzmann-Shannon Entropy

The similarity of the IPA and the projected gradient is clear, but the former avoids the computation of the projection required at each iterate. Hence IPA can be seen as an implicit projected gradient algorithm. It finds, as in the case of the EM algorithm, a nonnegative minimizer of φ . One of the powerful benefits of the IPA update formula is that all iterates remain in the essential domain of f if the initial iterate does. In our case, we choose f so that its essential domain is equal to the set of prior constraints. A well suited candidate is the Boltzmann-Shannon entropy function known also as the Kullback-Leibler (KL) distortion or distance associated with the convex differentiable function $f(x) = x \log(x) - x$, refer to figure 3.11. It is defined as follow,

Definition 3.9. *Let u and v be positive numbers*

$$KL(u, v) = u \log \frac{u}{v} + v - u$$

$$KL(u, 0) = +\infty \quad KL(0, v) = v \quad KL(0, 0) = 0$$

Extending to nonnegative vectors a and b

$$KL(a, b) = \sum_{j=1}^N KL(a_j, b_j) = \sum_{j=1}^N (a_j \log \frac{a_j}{b_j} + b_j - a_j)$$

A generalized distance D_f does not have to be symmetric unless f is quadratic [11]. The KL distance is one instance; note how in general

$$KL(a, b) \neq KL(b, a)$$

What appears to be a shortcoming gives actually rise to interesting features. Here is the tale of two algorithms. Although the EM algorithm was originally devised as a statistical parameter estimation and was not thought of as linked to any linear system, Byrne [23] showed convincingly that the EM method leads to a nonnegative minimizer of $KL(y, Cx)$. He showed also that the simultaneous MART (SMART) converges to a nonnegative minimizer of $KL(Cx, y)$. Like two faces of the same coin, EM and SMART are two algorithms deduced from the same distance.

3.10 Nonnegativity Minimization Algorithm

The KL function is a Bregman-Legendre function [11]. Its conjugate is the exponential function, $f^*(x^*) = \exp(x^*)$. Hence theorem 3.2 is very useful to invert the update formula (3.7.6). Recall theorem 3.2 that states for a Legendre function f

$$(\nabla f)^{-1} = \nabla f^*$$

We have then $(\nabla f)^{-1}(x) = \exp(x)$, which is a bijection. We also have $\nabla f(x) = \log(x)$. Recall equation (3.7.6)

$$\nabla f(x^{\ell+1}) = \nabla f(x^\ell) - \gamma \nabla \varphi(x^\ell)$$

which implies

$$\nabla^{-1} f(\nabla f(x^{\ell+1})) = \nabla^{-1} f(\nabla f(x^\ell) - \gamma \nabla \varphi(x^\ell))$$

Hence

$$\begin{aligned} x^{\ell+1} &= \exp(\nabla f(x^\ell) - \gamma \nabla \varphi(x^\ell)) \\ &= \exp(\log(x^\ell) - \gamma \nabla \varphi(x^\ell)) \\ &= x^\ell \exp(-\gamma \nabla \varphi(x^\ell)) \end{aligned}$$

knowing that

$$\varphi(x) = \frac{1}{2} \|\hat{x} - x\|_W^2 + \alpha \psi(x) \quad (3.10.1)$$

where ψ is a convex differentiable function. The minimization algorithm for this particular choice of f is the following

Algorithm 3.1. Choose $\gamma > 0$ and choose $\alpha \in]0, +\infty[$. Start with $x^0 \in \text{int } C$. For $\ell = 0, 1, \dots$ compute

$$x_i^{\ell+1} = x_i^\ell \exp(-\gamma (\nabla \varphi)_i(x^\ell)), \quad i = 1, \dots, N$$

until convergence.

In the limit, algorithm (3.1) finds an approximate solution $x^* = \text{prox}_\psi^\alpha(\hat{x})$. The convergence result for this algorithm is stated in theorem 3.3.

This is a generalized proximal approach that computes x^* , the weighted projection of \hat{x} onto \mathbb{R}_+^N . It is a proximal method, that generalizes the projection operator, and a distinctive iterative algorithm, that requires a relatively simple calculation executed repeatedly. The iterations give rise to a sequence of approximate answers that converges to the solution of the problem, x^* , regardless of the starting point $x^0 \in \text{int } C$. Resorting to a projector operator seems like an intuitive choice to get an “optimal” nonnegative solution. Can we justify this choice? This is the topic of the next chapter.

Chapter 4

Parameter Estimation

We are solving an inverse problem that involves the reconstruction of a dynamic nonnegative image x as a medical imaging application within the context of nuclear medicine. We showed how the formulation of this problem gave rise to the application of a temporal recursive filter, namely the Kalman filter. KF, having its origins in the well-known least squares (LS) techniques [7, 87], is more suitable for over-determined problems. Our dynamic image reconstruction we deal with is, however, an under-determined problem, and KF does not produce the desired image. As it was evidenced in section 2.4, we will probably get some negative components in the solution \hat{x} , which is not a feasible solution in nuclear medicine. We offered in section 3.10 a remedy to overcome this shortcoming. It is in essence a weighted projection of \hat{x} onto the nonnegative orthant \mathbb{R}_+^N while using a generalized proximal approach to achieve this goal. The nonnegative activity x gives rise to the solution \hat{x} as the BLUE (best linear unbiased estimator) when we do not enforce the positivity and by naturally projecting onto \mathbb{R}_+^N to have x^* , we intent to get back a “best” nonnegative estimate of the original activity x . How “good” is x^* ? Is it better than \hat{x} and in what sense? Is this projection the “optimal” one? Those are the questions we aim to answer in this chapter relying on parameter estimation theory. Parameter estimation is an active

field of research and many references can be consulted for more information; refer for instance to [2, 3, 18, 36, 90]. Let us first start with a summary.

4.1 Overview

An estimator seeks to approximate an unknown parameter of a physical model based on direct or indirect measurements. We are concerned with the accuracy and precision of the measurements and these are the precision and accuracy of the estimator as well [18]. An *estimator* is a vector function of the data made to evaluate the parameters. The *estimate* is a value taken by the estimator while the *estimand* is the quantity to be estimated; this will always be, for us, the true/simulated value of the parameter. Estimation theory, a branch of statistics and signal processing, concerns itself by studying desirable characteristics of an estimator and the consequences once we choose a certain one. When we are confronted with discrepancies in measuring the same quantities, we are then faced with the problem of getting the “best” estimation of these quantities while reducing the effect of the errors. Usually, we do not possess enough or exact knowledge about the errors or discrepancies, thus we treat them as random variables. Often, we have only limited or no knowledge about the nature of the randomness of the errors. Hence, the literature provides a number of so called “good” estimators, i.e., estimators with desirable properties based on certain assumptions on the error behaviour. So terms like “good” or “optimal” refer to specific underlying assumptions. Under certain conditions an estimator has all or almost all properties of a good estimator, and this is indeed the case for the Kalman filter. The LS, minimum mean square (MMS), maximum a posteriori (MAP), maximum likelihood (ML) are examples of criteria for choosing a “good” estimator [2, 3, 18, 36, 90]. Each criterion is, in its own right, concerned with achieving desirable properties. Hence, the need arises to codify and rationalize our choices and put them on firm ground.

Applications, particularly in astronomy and space, were the driving force of estimation. The Babylonians (300 BC) [90] were interested in knowing the positions of moving celestial bodies from various measurements, thus they were using the arithmetic mean as an estimator. Euler was interested (1750) in analyzing the irregularities in the motion of some planets taking advantage of the then recent development of the calculus of probabilities. Bernoulli (1777) was aware of the problem of outliers in the data when we use the arithmetic mean while the probabilistic Bayes' rule (1760) brought out the cornerstone of most methods in estimation theory. This leads the way for stochastic approaches to the estimation problem.

The LS method, a deterministic approach, is a very popular technique used to compute estimations of parameters to fit the measurements. This is the oldest technique and has its origins in the famous memoir, with application to astronomy, of the French mathematician Legendre, "Nouvelles Méthodes pour la Détermination des Orbites des Comètes" (1805) [90]. Then the German mathematician Gauss (1809) claimed that he had been using the LS technique as early as 1795. Historians later on could substantiate his claim and now LS is generally attributed to Gauss. This reminds us of the bitter anteriority dispute of the Leibniz-Newton controversy about the invention of Calculus. Several authors such as Cauchy contributed to the enrichment of the LS method.

Departure from the LS started early in the 20th century with Pearson, when he came up with the method of moments, and then with Fisher, when he introduced the ML approach in 1911. Contributions to parameter estimation theory was mainly done by statisticians until Wiener and Kolmogorov tried (1940) to bridge the gap between them and mathematicians, connecting the two views of the stochastic and the deterministic [90]. This led to the KF (1960).

The KF equations are not free from controversy either. The case with Legendre and Gauss contention of the antecedency of LS is repeated with Swerling who had

published a similar work (1958) as Kalman did, and claimed priority over Kalman but without success. The time was just ripe for the new method. KF [62], which could be seen as a sequential LS, saw its first application by NASA scientists to solve the problem of trajectory estimation for the Apollo program. This application to aerospace completed the full circle since the Babylonians' time 23 centuries ago.

LS has mainly two versions, the original and simple one is referred to as the ordinary LS (OLS), which could be seen as a determinist view of the estimation problem. LS has a more sophisticated version known as the weighted or generalized LS (WLS or GLS). The latter differs from the former in that it introduces weights which regulate the importance of each observation. Aitken [3], from New Zealand (1935), has demonstrated that with a proper choice of the weighting matrix equal to P^{-1} , where P is the covariance matrix, the LS approach performs as best as it can in the sense it is the BLUE. Hence, the GLS is also referred to as the Aitken estimator.

Suppose that we wish to have an estimate $\hat{\theta}_{LS}$ of an unknown θ from a data set y using LS. The idea is to minimize the incurred error $\tilde{\theta} = \theta - \hat{\theta}$, that is to say we aim to minimize $l_{LS}(\tilde{\theta}) = (\theta - \hat{\theta})^\top (\theta - \hat{\theta})$. The loss function l_{LS} of LS could be generalized by introducing a symmetric positive definite matrix W . We would then aim to minimize the cost function, $l_{GLS}(\tilde{\theta}) = (\theta - \hat{\theta})^\top W (\theta - \hat{\theta})$ instead, as we would recover the one of OLS by letting $W = I$.

In probability, the covariance matrix furnishes a measure of the behavior of the estimator. First, recall the well known property of the covariance matrix, see for instance [18].

Proposition 4.1. *The covariance matrix P of a stochastic real values vector is symmetric positive semi-definite. It is definite if and only if its stochastic components are linearly independent. In this case, its inverse P^{-1} is also symmetric positive definite.*

It comes at no surprise that we would shoot for an estimator that minimizes the

error variance. Hence the mean square error (MSE) was introduced and could be seen as the probabilistic equivalent to the deterministic LS cost function

$$\text{MSE}(\hat{\theta}) = L_{MS}(\tilde{\theta}) = \mathbb{E} \left((\theta - \hat{\theta})^\top (\theta - \hat{\theta}) \right) \quad (4.1.1)$$

The variance (Var) measures also the estimator's deviation from its expected value. It measures the variations of the estimates from try out to try out as a consequence of the variations in the observations and does so relative to the mean value of the estimator. This mean value needs not to be equal to the true value of the parameter. It can be proven that [90]

$$\text{MSE}(\hat{\theta}) = \mathbb{E}(\hat{\theta} - \mathbb{E}(\hat{\theta}))^2 + (\mathbb{E}(\hat{\theta}) - \theta)^2 = \text{Var}(\hat{\theta}) + \beta^2(\hat{\theta}) \quad (4.1.2)$$

where $\beta(\hat{x})$ is called the bias. The mean square error is then equal to the sum of the variance and the square of the bias. The bias of an estimator gives information about its systematic error while its covariance gives us its nonsystematic error. Unbiased estimator means that if we take a very large number of random observations, the average value of the parameter estimates will be theoretically exactly equal to the estimand. Minimum variance means that the estimator has the smallest variance, and thus the narrowest confidence interval, of all estimators of its type. The MSE evaluates the quality of an estimator in terms of its variation and unbiasedness. For an unbiased estimator, the MSE is the variance.

$$\text{MSE}(\hat{\theta}) = \text{Var}(\hat{\theta}) \quad (4.1.3)$$

The operators MSE and Var have the disadvantage of heavily weighting outliers. This is a consequence of the squaring of each term, which weighs large errors more heavily than small ones. Researchers have used alternatives to avoid this shortcoming. Hence loss functions based on the mean absolute error or the median have been introduced.

Like any other loss function L , notice how the loss function L_{MS} in (4.1.1), is chosen to satisfy three conditions [90].

1. $L(0) = 0$, to ensure that there is no loss with a zero error.
2. L is symmetric, to guarantee the independence of L of the sign of the error.
3. L is convex, to make sure that L is an increasing function of the error variable.

Hence the powerful tools of convex optimization, introduced in chapter 3, are very useful once more. For example, the KF estimator, a conditional expectation that has been shown acting as a projection operator, could be retrieved using convex optimization, refer for instance to [7, 87].

MSE is regarded as a way of measuring the performance of an estimator. It was proven [3] that the linear minimum MSE and GLS give the same estimator when the weighting matrix W is chosen to be the real symmetric positive definite P^{-1} , with P being the covariance matrix. Recall that, in section 2.3, an estimator which is unbiased, linear, and minimizes the MSE is called the BLUE. It is only natural that someone would use GLS with $W = P^{-1}$ even in the nonlinear case. This has been done in inverse problems, refer for example to [63, 102]. However, to the best knowledge of the author no property of optimality or other was formulated as yet; this is the theme of section 4.4.

4.2 Parameter Estimation Properties

The ML, introduced by Fisher (1911), is another estimator. When we have measurements y of the quantity z , which is only known to belong to a certain family of probability distributions $g(y|\theta)$, we want to estimate the unknown parameter θ that maximizes the likelihood of getting the data y . ML, LS, GLS estimators, and others have given rise to estimators of the form $\hat{\theta} = h(y)$ [2, 3, 18, 36, 90]. Based on the errors $\tilde{\theta} = \theta - \hat{\theta}$, they should be “good” estimators in verifying as many as possible

of the following requirements, although these are artificial and arbitrary; in the sense that other interested people may not agree with the rationale behind this choice.

1. **Consistency** The ideal case is to have $\tilde{\theta} = 0$, however, this is close to impossible because of the Cramér-Rao inequality that set a lower bound on $\text{Cov}(\tilde{\theta})$ [18]. Hence we must settle for less and require that at least this error converges, with probability 1, to 0 when we augment the size n of the data.

$$\lim_{n \rightarrow \infty} P(|\hat{\theta} - \theta| < \epsilon) = 1 \quad \forall \epsilon > 0 \quad (4.2.1)$$

2. **Sufficiency** A reasonable requirement is that the estimator should use all available data and extract all information carried in them.
3. **Acceptability** When we compute the estimator with a different set of data, it should still give an acceptable value of the parameter. The data vectors in the observed data set contain no missing elements.
4. **Unbiasedness** Although a biased estimator could be useful, an unbiased one, that is $\mathbb{E}(\hat{\theta}) = \mathbb{E}(\theta)$, is usually desired at least asymptotically.
5. **Efficiency** The error covariance, $\mathbb{E}\left((\theta - \hat{\theta})(\theta - \hat{\theta})^\top\right)$, should be minimal at least asymptotically.

The Kalman filter estimator was established [87] to satisfy the requirements above thus providing estimates that are optimal in the LS, ML and Bayesian sense for a Gauss-Markov model. Optimization techniques that are used to numerically compute estimators include, but are not limited to, Newton procedure, generalized Gauss-Newton, steepest descent, iteratively reweighted LS, conjugate gradient, and Polak-Polyak-Ribiere-Sorenson procedure [18].

We usually make the Gaussian assumption about the errors $\tilde{\theta} = \theta - \hat{\theta}$ based on the rationale of the Central Limit Theorem of probability theory. Given an unbiased

estimator $\hat{\theta}$ of θ and that the error belongs to the Gaussian family with a given covariance matrix P , the probability density function (pdf) is

$$g(\tilde{\theta}) = \frac{1}{(2\pi)^{\frac{N}{2}} \sqrt{\det(P)}} \exp\left(-\frac{1}{2}(\theta - \hat{\theta})^\top P^{-1}(\theta - \hat{\theta})\right) \quad (4.2.2)$$

Let

$$a^2 = (\theta - \hat{\theta})^\top P^{-1}(\theta - \hat{\theta}) \quad (4.2.3)$$

The quadratic form equation (4.2.3), which is also the useful part of the log-likelihood LL function, defines a surface and characterizes $g(\tilde{\theta})$ of (4.2.2) in the sense that for each a^2 , we have a constant pdf value. It can be shown that this surface is an ellipsoid in \mathbb{R}^n [90]. Its center is the estimate $\hat{\theta}$ and its semiaxes have magnitude $a^2 \lambda_i > 0$ ($i = 1, \dots, n$) and directions defined by ξ_i , λ_i and ξ_i being the eigenvalues and eigenvectors respectively of the symmetric positive definite covariance matrix P .

So far we have seen that two main methods of parameter estimation theory exist, the deterministic LS approach and its variants and the probabilistic Bayesian approaches such as ML. Different estimation criteria could give away different estimators stressing the point of arbitrariness of an optimality choice. Fortunately, there are times when both methods yield the same estimator as in the case of KF [7, 87], a linear, unbiased, GLS and MSE estimator, within the framework of Gaussian pdf and linear model. These last two assumptions prove useful to extract some properties about the estimator x^* .

4.3 Positive Image

The nonnegative activity x gives rise to the KF estimator \hat{x} . KF propagates \hat{x} and covariance matrix P at each step of the recurrences (2.3.3) and (2.3.6) with an error $\tilde{x} = \hat{x} - x$. Making the assumption that \tilde{x} follows a normal distribution of centre 0

and covariance matrix P , $\tilde{x} \sim \mathcal{N}(0, P)$, the pdf is

$$g(\tilde{x}) = \frac{1}{(2\pi)^{\frac{N}{2}} \sqrt{\det(P)}} \exp\left(-\frac{1}{2}(x - \hat{x})^\top P^{-1}(x - \hat{x})\right) \quad (4.3.1)$$

We have

$$\hat{x} = x + \tilde{x} \quad (4.3.2)$$

We reverse the process. Knowing or observing \hat{x} , we aim to select the nonnegative parameter value x^* which realizes the largest possible pdf $g(\tilde{x})$. In other words, we look for a constrained ML estimator of (4.3.2). We have seen that equation (4.2.3) characterizes completely the Gaussian pdf $g(\tilde{x})$. Hence we need to minimize

$$\frac{1}{2}(x - \hat{x})^\top P^{-1}(x - \hat{x}) \quad (4.3.3)$$

This is equivalent to minimizing $\|x - \hat{x}\|_{P^{-1}}^2$ in \mathbb{R}_+^N . The minimizer is an oblique/weighted projection onto the nonnegative orthant. We obtain the same quantity to minimize while using the likelihood terminology. The log-likelihood function, that is the log of the pdf $g(\tilde{x})$, to be maximized with respect to x is

$$LL(x) = -\frac{1}{2}(x - \hat{x})^\top P^{-1}(x - \hat{x}) + c$$

for some constant c . Because the logarithm is a continuous strictly increasing function over the range of the likelihood, the values which maximize the likelihood will also maximize its logarithm. Since maximizing the logarithm requires simpler algebra, it is the logarithm which is maximized.

4.4 Properties of the Projected KF Estimator

We deal with an operator that we could refer to as a projected KF and we would like to assess this estimator's goodness in terms of known properties of parameter estimation. Next, we establish that it is a ML estimator.

4.4.1 Maximum Likelihood

Maximum likelihood estimation is based on the assumptions that the distribution of the data is known and the expectation model is correct. ML methods have desirable mathematical and optimality properties. Recall the invariance property of ML estimators [90].

Lemma 4.1. *Suppose that $\hat{\theta}$ is the ML estimator of θ in \mathbb{R}^n . Consider the (not necessarily injective) vector mapping $\varrho : \mathbb{R}^n \rightarrow \mathbb{R}^m$. Then $\varrho(\hat{\theta})$ is the ML of $\varrho(\theta)$ in \mathbb{R}^m .*

We set $C = \mathbb{R}_+^N$, the nonnegative orthant where N is the size of the activity we are solving for. For the symmetric positive definite matrix P^{-1} , define

$$\begin{aligned} \text{proj}^{P^{-1}} : \mathbb{R}^N &\rightarrow C \\ \hat{z} &\mapsto z^* = \arg \min_{z \in C} \|z - \hat{z}\|_{P^{-1}}^2 \end{aligned} \quad (4.4.1)$$

where

$$\|z - \hat{z}\|_{P^{-1}}^2 = (z - \hat{z})^\top P^{-1} (z - \hat{z}) \quad (4.4.2)$$

Refer to figure 4.1 that illustrates this projection in a 2D setting. The estimator \hat{x} of the nonnegative activity x , being the KF estimator, is then a ML of x in \mathbb{R}^N within the framework of a Gaussian pdf and a linear model. Moreover, since C is a closed and convex set and the quadratic form $(z - \hat{z})^\top P^{-1} (z - \hat{z})$ is convex in the variable $z \in \mathbb{R}^N$, we conclude that $\text{proj}^{P^{-1}}(\hat{z})$ exists and is unique [1], so that the mapping 4.4.1 is well defined. Furthermore, $\text{proj}^{P^{-1}}(v) = v$ if and only if $v \in C$ because the matrix P^{-1} is positive definite. Recall that x is nonnegative and \hat{x} is its ML in \mathbb{R}^N . Apply lemma 4.1 with $\varrho = \text{proj}^{P^{-1}}$, we deduce that $x_{P^{-1}}^* = \text{proj}^{P^{-1}}(\hat{x})$ is the ML estimator of $x = \text{proj}^{P^{-1}}(x)$ with respect to the matrix P^{-1} . We have just proved the following theorem.

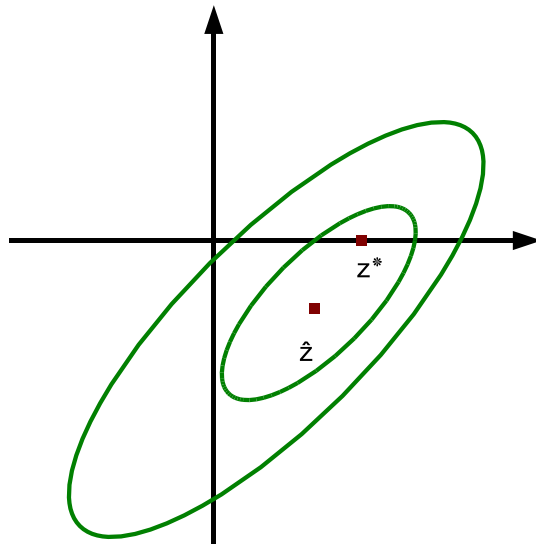


Figure 4.1: Illustrating the oblique projection for 2D. \hat{z} is the KF estimate and z^* is the nonnegative estimate as an oblique projection onto the first quadrant.

Theorem 4.1. *Let x be a nonnegative activity and \hat{x} be its KF estimator in \mathbb{R}^N . Let P^{-1} be the symmetric positive definite matrix inverse of the covariance matrix of the error $\tilde{x} = x - \hat{x}$. While \hat{x} is the ML estimator of x in \mathbb{R}^N , the projected KF estimator $x_{P^{-1}}^*$ is the constrained ML estimator of x in $C = \mathbb{R}_+^N$ within the framework of a linear model and Gaussian pdf given by equation 4.3.1.*

For ease of notation, we will drop from now on the subscript P^{-1} . Hence we would refer to $x_{P^{-1}}^*$ as x^* and it is assumed that the projection is done with respect to P^{-1} .

4.4.2 Consistency

Besides the invariance property, a ML estimator possesses a second property ([42] and references therein).

Theorem 4.2. *A ML estimator is consistent in the sense of equation 4.2.1.*

Thus we have another property of our estimator x^* giving the following result.

Corollary 4.1. *The estimator $x^* \in C$ is an asymptotic consistent ML estimator of the nonnegative activity x , that is to say*

$$\lim_{n \rightarrow \infty} P(|x^* - x| < \epsilon) = 1 \quad \forall \epsilon > 0 \quad (4.4.3)$$

where n is the size of the data.

The set $C = \mathbb{R}_+^N$ is not a subspace of \mathbb{R}^N and $\text{proj}^{P^{-1}}$ is a projection, the mapping 4.4.1 is a nonlinear operator that maps \hat{x} into x^* . Hence x^* is not a linear function of x since \hat{x} is already one. Usually, there will be no unbiased and optimal nonlinear estimator of x even in the event of normally distributed data. Nonetheless, ML estimators could exhibit asymptotic behavior, to the extent they could be unbiased and optimal for a fixed number of data. In addition, the convex set C , being a cone, has the salient property of “almost” linearity referred to as the nonnegative homogeneity [20], $\text{Proj}_C(\alpha z) = \alpha \text{Proj}_C(z) \quad \forall \alpha > 0$.

4.4.3 Unbiasedness

Theorem 4.3. *As its parent \hat{x} , x^* is an unbiased constrained estimator of the non-negative activity x , meaning*

$$\mathbb{E}(x^*) = \mathbb{E}(x) \quad (4.4.4)$$

Simon et. al. [88] proved the same property and the optimality one, that follows shortly, for Kalman filtering with state equality constraints $Dx_k = d_k$; that is when the state x_k is known to belong to a hyperplane. Simon et.al. [89] then proved both properties in the case of state variable inequality constraints $Dx_k \leq d_k$. They notice that almost all algorithms for solving such optimization problems belong to the active set methods. They base their argument on this fact assuming that the correct set of active constraints is known a-priori to them. We do not use an active set method to solve the constrained Kalman and we do not even know if there are any null components nor where they are located in the activity x_k . We are instead solving for lower bounds $x_k \geq 0$ constraints. Their arguments are therefore not useful to us.

Proof: We seek to find the oblique projection of \hat{x} on the positive orthant C ,

$$\min_z \frac{1}{2} \|z - \hat{x}\|_{P^{-1}} \quad \text{subject to } z \in C$$

The Lagrangian of the constrained problem is,

$$\mathcal{L}(z, \lambda) = \frac{1}{2} (z - \hat{x})^\top P^{-1} (z - \hat{x}) - \lambda^\top z \quad (4.4.5)$$

We formulate the first order Karush-Kuhn-Tucker conditions,

Stationarity:

$$\nabla \mathcal{L}(x^*, \lambda^*) = 0 \quad (4.4.6)$$

that is,

$$P^{-1}(x^* - \hat{x}) - \lambda^* = 0$$

or

$$\lambda^* = P^{-1}(x^* - \hat{x}) \quad (4.4.7)$$

Primal feasibility:

$$x^* \geq 0 \quad (4.4.8)$$

Dual feasibility:

$$\lambda^* \geq 0 \quad (4.4.9)$$

Complementary slackness:

$$(\lambda^*)^\top x^* = 0 \quad (4.4.10)$$

Even though we are interested in a general case of an oblique projection, the argument in the case of the orthogonal projection is interesting in itself. So let us consider the particular case when $P^{-1} = I$. Equation 4.4.7 combined with equation 4.4.9 gives,

$$\lambda^* = x^* - \hat{x} \geq 0$$

so that $x^* \geq \hat{x}$. But we have $x^* \geq 0$, equation 4.4.8, and $(x^* - \hat{x})^\top x^* = 0$, equation 4.4.10, which implies that $x^* = \max(\hat{x}, 0)$.

Recall Jensen's Inequality [20]. If x is a random variable such that $x \in \Delta_g$ with probability one, and g is convex, then we have

$$g(\mathbb{E}(x)) \leq \mathbb{E}(g(x)) \quad (4.4.11)$$

provided the expectations exist. We apply Jensen's Inequality 4.4.11 with the convex function $\max(y, 0)$ to get,

$$0 \leq \max(\mathbb{E}(\hat{x}), 0) \leq \mathbb{E}(\max(\hat{x}, 0))$$

$$0 \leq \mathbb{E}(\hat{x}) \leq \mathbb{E}(x^*) \quad (4.4.12)$$

Using equation 4.4.10 and the fact that the error and the state are uncorrelated [7], we have

$$\mathbb{E}(\lambda_i^* x_i^*) = \mathbb{E}(\lambda_i^*) \mathbb{E}(x_i^*) = 0 \quad \forall i \quad (4.4.13)$$

Recall that the Kalman output \hat{x} , as an estimate of the positive activity x , is unbiased, $\mathbb{E}(\hat{x}) = \mathbb{E}(x) \geq 0$. There are two cases to consider for equation 4.4.13. On one hand, if $\mathbb{E}(\lambda_i^*) = \mathbb{E}(x_i^* - \hat{x}_i) = 0$ for some i , then $\mathbb{E}(x_i^*) = \mathbb{E}(\hat{x}_i)$. On the other hand, if $\mathbb{E}(x_i^*) = 0$ for some i , then using inequality 4.4.12 we have $0 \leq \mathbb{E}(\hat{x}_i) \leq \mathbb{E}(x_i^*) = 0$; that is $\mathbb{E}(\hat{x}_i) = 0$. Both cases sum up to

$$\mathbb{E}(x^*) = \mathbb{E}(\hat{x})$$

We therefore conclude the proof of the unbiasedness in the case of orthogonal projection, $\mathbb{E}(x^*) = \mathbb{E}(x)$.

Let us now proceed to the more general case when the projection is oblique. The complementary slackness says that $(x^*)_i \lambda_i^* = 0 \quad \forall i$. Since the error and the state are uncorrelated [7], we have $\mathbb{E}(x_i^* \lambda_i^*) = \mathbb{E}(x_i^*) \mathbb{E}(\lambda_i^*) = 0$. There are two cases to consider. On one hand, if $\mathbb{E}(\lambda_i^*) = 0$ for some i , then $P^{-1} \mathbb{E}(x_i^* - \hat{x}_i) = 0$ so that $\mathbb{E}(x_i^*) = \mathbb{E}(\hat{x}_i)$. On the other hand, if $\mathbb{E}(x_i^*) = 0$ for some i , let \mathcal{I} be the set of these i .

Since $\forall i \in \mathcal{I}$, $\mathbb{E}(x_i^*) = 0$, then $x_i^* = 0$ since $x^* \geq 0$; which implies $\lambda_i^* = (P^{-1}(x^* - \hat{x}))_i > 0$. Since $\forall i \in \bar{\mathcal{I}}$, $x_i^* > 0$, then $(P^{-1}(x^* - \hat{x}))_i = 0$. The matrix P^{-1} is symmetric positive definite, therefore

$$\begin{aligned} \sum_{i \in \mathcal{I} \cup \bar{\mathcal{I}}} (x^* - \hat{x})_i (P^{-1}(x^* - \hat{x}))_i &= \sum_{i \in \mathcal{I}} (x^* - \hat{x})_i (P^{-1}(x^* - \hat{x}))_i + \sum_{i \in \bar{\mathcal{I}}} (x^* - \hat{x})_i (P^{-1}(x^* - \hat{x}))_i \\ &= \sum_{i \in \mathcal{I}} (x^* - \hat{x})_i (P^{-1}(x^* - \hat{x}))_i \\ &= (x^* - \hat{x})^\top P^{-1} (x^* - \hat{x}) \\ &\geq 0 \end{aligned}$$

Consequently,

$$\sum_{\iota \in \mathcal{I}} (x^* - \hat{x})_{\iota} (P^{-1}(x^* - \hat{x}))_{\iota} \geq 0$$

or

$$-\sum_{\iota \in \mathcal{I}} \hat{x}_{\iota} (P^{-1}(x^* - \hat{x}))_{\iota} \geq 0$$

Passing to the expectation we have,

$$-\sum_{\iota \in \mathcal{I}} \mathbb{E}(\hat{x}_{\iota}) \mathbb{E}((P^{-1}(x^* - \hat{x}))_{\iota}) \geq 0$$

We know that $\mathbb{E}(\hat{x}_{\iota}) \geq 0$ and $\mathbb{E}((P^{-1}(x^* - \hat{x}))_{\iota}) > 0$; this implies $\mathbb{E}(\hat{x}_{\iota}) = 0$ for all $\iota \in \mathcal{I}$. The two cases that we considered here allow us to conclude that $\mathbb{E}(x^*) = \mathbb{E}(\hat{x}) = \mathbb{E}(x)$. The projected KF x^* is an unbiased constrained estimator of the nonnegative activity x .

□

4.4.4 Optimality

Hitherto we have seen that $x^* \in C$ is an affine ML, consistent, and unbiased constrained estimator of the nonnegative activity x . At the end of subsection 4.4.1, we concluded that x^* , being ML, is an asymptotically optimal estimator with probability 1. However, the KF estimator \hat{x} is also an unbiased and optimal estimator of x while being linear. The estimator \hat{x} is not anymore optimal when we introduce the positivity constraint into the estimator. Does x^* do better than the KF solution \hat{x} since x^* is rather a constrained ML in C ? The answer is yes in the following sense.

Theorem 4.4. *The estimator $x^* \in C$ of x performs better than the estimator \hat{x} in the sense that the mean square error of x^* is smaller than the mean square error of \hat{x} ,*

$$\text{MSE}(x^*) \leq \text{MSE}(\hat{x}) \tag{4.4.14}$$

which is equivalent to say

$$\text{Tr}(\text{Cov}(x - x^*)) \leq \text{Tr}(\text{Cov}(x - \hat{x}))$$

Before giving a proof, we need a definition of nonexpansive mappings [15].

Definition 4.1. Let χ be a Hilbert space and $\|\cdot\|$ a vector norm. An operator T , not necessary linear, in χ is a nonexpansive mapping if $\forall z_1, z_2 \in \chi$, then

$$\|T(z_2) - T(z_1)\| \leq \|z_2 - z_1\| \quad (4.4.15)$$

Recollect a classical result in convex optimization about projection onto closed convex sets [91].

Proposition 4.2. Let $u, v \in \chi$, a Hilbert space with $\|\cdot\|$ as a vector norm, and let $\text{proj}_F(u), \text{proj}_F(v)$ be the corresponding projections onto any closed and convex set F , then

$$\|\text{proj}_F(u) - \text{proj}_F(v)\| \leq \|u - v\| \quad (4.4.16)$$

This property of the projection onto closed convex set states that the projection operator is nonexpansive (definition 4.1), see figure 4.2 for a geometrical intuition of this result.

Proof: Recall the definition of the operator MSE giving by equation (4.1.1).

$$\begin{aligned} \text{MSE}(x^*) &= \mathbb{E}((x - x^*)^\top (x - x^*)) \\ \text{MSE}(\hat{x}) &= \mathbb{E}((x - \hat{x})^\top (x - \hat{x})) \end{aligned}$$

As we have seen in the case of unbiased estimators, like both x^* and \hat{x} , equation (4.1.1) becomes equation (4.1.3) entailing

$$\begin{aligned} \text{MSE}(x^*) &= \text{Var}(x^*) \\ \text{MSE}(\hat{x}) &= \text{Var}(\hat{x}) \end{aligned}$$

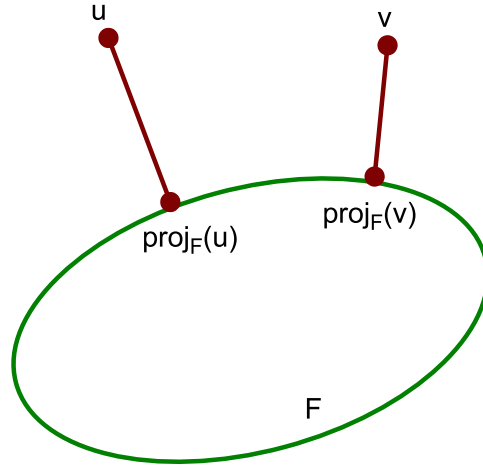


Figure 4.2: Illustrating the inequality (4.4.16) for projection onto closed convex set.

Apply proposition 4.2 to $\chi = \mathbb{R}^N, F = C = \mathbb{R}_+^N$ with proj_C the oblique projection onto C as defined in the operation (4.4.1), so that $x^* = \text{proj}_C(\hat{x})$. Since $x \in C$ means $x = \text{proj}_C(x)$, we have

$$\|x - x^*\|^2 \leq \|x - \hat{x}\|^2 \quad (4.4.17)$$

The expectation operator is positive, $Y \geq 0 \Rightarrow \mathbb{E}(Y) \geq 0$. Thus inequality (4.4.17) implies

$$\mathbb{E}((x - x^*)^\top (x - x^*)) \leq \mathbb{E}((x - \hat{x})^\top (x - \hat{x}))$$

that is

$$\text{MSE}(x^*) \leq \text{MSE}(\hat{x})$$

Note

$$\text{MSE}(Y) = \text{Tr}(\text{Cov}(Y))$$

Hence theorem 4.4 is also saying

$$\text{Tr}(\text{Cov}(x^* - x)) \leq \text{Tr}(\text{Cov}(\hat{x} - x)) \quad (4.4.18)$$

The estimator $x^* \in C$ performs better than \hat{x} in the minimum variance sense.

□

4.5 Summary

The KF estimator $\hat{x} \in \mathbb{R}^N$ is the optimal estimator of the activity x in the mean squares error sense. However, it is not necessarily nonnegative as it is the case for x . The matrix P is the covariance of the error $(x - \hat{x})$ which is updated at each step of the KF recurrence. We have seen that using the oblique projection x^* , with respect to the symmetric positive definite matrix P^{-1} , as an estimate in \mathbb{R}_+^N performs better than \hat{x} . Not only does the estimator x^* conserve the same properties as the KF \hat{x} of being unbiased, consistent, and ML, but in addition, the $\text{MSE}(x^*)$ is smaller than $\text{MSE}(\hat{x})$. Results in chapter 6 confirms these properties. The next chapter will tackle the lack of spatial regularization in KF.

Chapter 5

Spatial Smoothness

We have concluded chapter 2 by stating that while the KF estimator includes temporal smoothness, it lacks the spatial one due to two main reasons. First we are solving for a nonnegative solution, and second we have at hand a spatially ill-posed problem; KF does not handle well either one of these two challenges. Hence we need to impose a spatial regularization which is basically done by introducing a priori knowledge into the problem. Two main practical approaches are generally applied namely, including constraints into the problem and using iterative solvers.

5.1 Nonnegativity Constraint

We have already covered in chapter 3 how we impose nonnegativity while using an iterative algorithm to achieve that goal. We have also seen in chapter 4 how the nonnegative solution x^* performs better than the KF \hat{x} .

The solution x^* is a constrained ML in the nonnegative orthant since we showed in chapter 4 that we had to maximize the likelihood function over its domain of definition. We also established that this constrained ML was equivalent to a GLS with

an appropriate choice of the weighting matrix. Thus we injected a-priori information about the unknown x by restricting the likelihood function domain to be nonnegative. This is fundamentally a statistical cornerstone to some of the regularization approaches. Recall the cost function that we aimed to minimize

$$(x - \hat{x})^\top W (x - \hat{x}),$$

where W is a symmetric positive definite weighting matrix. This is equivalent to minimizing

$$\|x - \hat{x}\|_W^2.$$

Using the matrix $W = P^{-1}$, where P is the covariance matrix of the error $x - \hat{x}$, gives the optimal nonnegative estimator, refer to theorem 4.4. Letting

$$\begin{aligned} \varphi(x) &= \frac{1}{2} \|x - \hat{x}\|_{P^{-1}}^2 \\ &= \frac{1}{2} (x - \hat{x})^\top P^{-1} (x - \hat{x}) \end{aligned}$$

we have

$$\nabla \varphi(x) = P^{-1} (x - \hat{x})$$

and algorithm 3.1 reduces to

Algorithm 5.1. Choose $\gamma > 0$ and start with $x^0 \in \text{int } C$. For $\ell = 0, 1, \dots$ compute

$$x_i^{\ell+1} = x_i^\ell \exp(-\gamma (P^{-1}(x^\ell - \hat{x}))_i), \quad i = 1, \dots, N$$

until convergence.

The clustering point of algorithm 5.1 is the nonnegative solution we desire.

5.2 Iterative Solver

Iterative methods could serve as a regularization of ill-posed problems. The number of iterations plays the role of the regularization parameter [15] since semiconvergence

happens when we deal with noisy images as it was remarked for the EM algorithm; refer to section 1.2 for more details. A typical iterative algorithm involves a relatively simple calculation, performed repeatedly. An iterative method gives rise to a sequence of approximate answers that, in the best case, converges to the solution of the problem. We have at hand an inverse problem of image reconstruction from projections, which is ill-posed. Consequently, as the number of iterations increases the iterates get at first closer to the desired solution and then move away. An obvious remedy is to stop the iterations earlier. Our approach involves stopping iterative algorithms prior to convergence; we confirm this numerically in section 6.4. Next we see how we employ a very well known regularization technique.

5.3 Tikhonov Regularization

V. Ivanov and D. Phillips were the first to propose each a different cure for ill-posed problems (1962), even though one year later, it was A. Tikhonov who independently offered a general method called regularization [15] that unifies both approaches of Ivanov and Phillips. The central idea is to proceed by approximate solutions that depend on a so called regularization parameter α . In the absence of noise in the data and errors in modeling, the approximated regularized solution converges to the exact solution as α gets closer to 0. Otherwise, we get an optimal approximate solution associated with an optimal α . Furthermore, appropriate choices of α give back both methods of Ivanov and Phillips. Lagrange multipliers are used to incorporate constraints into a minimization problem in the optimization field. Tikhonov [98] built upon properties of these multipliers to develop his regularization algorithm. The original idea in Tikhonov regularization was based on the approximation of the operator that transforms the unknown x into the data y by a coercive operator, which had a bounded inverse. KF filters out the solution over time, and we utilize Tikhonov

regularization for a spatial filtering.

We mention in sections 1.2 and 1.3 that we deal with an inverse problem that is ill-posed. Recall that a problem is called well-posed in the sense of Hadamard [15] if it obeys three conditions: the solution of the problem is unique, exists for any data, and depends continuously on the data. A problem is ill-posed if it fails to satisfy any of these conditions. We are solving linear systems involving huge, in the thousands or even millions of entries, and not necessary sparse or of any known structure matrices. When a small change in the coefficients of a matrix results in a large change in the solution, we say that the matrix, say B , (or the problem) is ill-conditioned. This is captured in the so-called condition number; which is

$$\kappa(B) = \|B^{-1}\| \|B\|$$

A problem with a huge condition number is said to be ill-conditioned while a problem whose condition number is close to one is said to be well-conditioned. For instance, when this number is 10^6 , a relative error in the data of magnitude 10^{-6} could generate an error of magnitude 100% in the solution. The condition number is also a measure of the cooperativeness of the problem with digital computation. Thus a well-conditioned numerically problem is well-posed and hence the close connection between these two properties. While ill-posedness pertains to a continuous problem, ill-conditioning is related to a discrete problem. When we discretize our ill-posed problem, the condition number of the corresponding discrete problem can be very large.

We face the ill-posedness problem when we estimate the unknown activity x from the data pertaining to x only. This data is incomplete, that is the system does not convey complete information about the activity. Even if we had very accurate, that is, noise free data, we would have difficulties solving our inverse problem. We should not look for an exact solution and should focus instead only on an approximate one, since an exact solution would be matching the noisy measurements. This first basic approach might help to cure ill-posedness. This does, however, not address issues of

uniqueness and ill-conditioning.

Having the ability to use prior knowledge concerning x could stabilize the algorithm. Tikhonov regularization [97, 98], known as ridge regression in the statistical community, is our third remedy to help cure ill-posedness. It has been introduced in various settings. It is also known in the literature as Tikhonov-Phillips regularization due to the work of D. L. Phillips [78]. Recall that we are minimizing the function of 3.7.1

$$\varphi(x) = \frac{1}{2} \|\hat{x} - x\|_W^2 + \alpha\psi(x) \quad (5.3.1)$$

When we are interested in a nonnegative solution only, we set $\alpha = 0$. Since we aim for a regularized solution, we must impose $\alpha > 0$. To enforce a Tikhonov regularization type, we choose

$$\psi(x) = \frac{1}{2} \|L(x - \bar{x})\|^2 \quad (5.3.2)$$

where \bar{x} is some target value of x and L is some appropriately selected Tikhonov matrix. For instance, if we choose $\bar{x} = 0$ and $L = I$, we are then concerned with a minimum norm solution. If we take $\bar{x} = 0$ and L to be some differential operator, we are then interested in a spatially smooth outcome. Choosing α to be high implies we are relying more in our prior information while having it extremely small means that we are not really interested in a regularized answer. Hence there is a risk of ending up with a solution that is shaped more with our prior and there is also a risk of ending with an undesired solution in case we forsake our prior. Attaining an equilibrium between including prior information and working with the data only is our goal, while accomplishing it is not an easy endeavor.

Choosing the regularization parameter α is at the same time important and hard. Intensive research in both the mathematical and statistical communities has been going on and the perfect recipe has yet to be found. It is known that an optimal value of α exists in theory in the sense that the regularized solution is the closest to the activity we are solving for [84]. The parameter α is plotted against an energy function

and we obtain a L-curve plot. It is called so because, in most cases, its log-log plot has the shape of the letter L that has a global minimum $\alpha_{optimal}$. Nevertheless in practice, determining $\alpha_{optimal}$ involves knowledge of the unknown x . Methods have been developed to estimate $\alpha_{optimal}$ when we have no knowledge about the unknown such as prescribed energy, prescribed discrepancy, and Miller approaches [14, 83, 84, 103].

The behavior of the regularized solution as a function of the regularization parameter emphasizes again the *semiconvergence* property [15]. Thanks to this property, we know that an optimal value of the regularization parameter exists, even if its determination may be difficult. Optimal here means that among all regularized solutions, that corresponding to this value of the regularization parameter provides the best approximation of the unknown object. As we mentioned earlier, an interesting property of several iterative methods for the solution of a linear problem is that they can be viewed as regularization methods if the problem is ill-posed. In these cases the role of the regularization parameter is played by the number of iterations and the semiconvergence, which holds true again, implies the existence of an optimal value of α , in the sense specified above. An $\alpha_{optimal}$ corresponds to an optimal number of iterations.

Regularization and optimization are closely connected. We emphasized in chapter 3 the central role that convex optimization plays and how generalized distance, known also as Bregman distance (1967), is important. These generalized distances and projections associated with them give rise to useful iterative algorithms [31]. They are useful also for regularization of optimization problems, thus we have a class of optimization algorithms called proximal minimization algorithms where originally Euclidian distances were used. Consider the following optimization problem

$$\begin{aligned} \text{Minimize} \quad & \varphi(x) \\ \text{s. t.} \quad & x \in F \end{aligned} \tag{5.3.3}$$

where $\varphi : \mathbb{R}^N \rightarrow \mathbb{R}$ is a given convex function and $F \subseteq \mathbb{R}^N$ is a nonempty closed convex subset of \mathbb{R}^N .

The idea is based on converting (5.3.3) into a sequence of optimization problems by adding terms to $\varphi(x)$. These added terms measure the distance between the variable vector x and the current iterate x^ℓ , either in the Euclidian sense, or in a generalized sense, according to some D_f distance. The proximal minimization algorithm can be viewed within the general approach known as regularization. Suppose that the problem (5.3.3) is ill-posed in some sense. By this we mean simply that it lacks some desirable property, such as spatial smoothness, for example. In such a situation it is sometimes possible to modify the original objective function $\varphi(x)$ by adding a term $\alpha\psi(x)$ such that the perturbed function $\varphi(x) + \alpha\psi(x)$ has the desired property, which $\varphi(x)$ lacks. α is the regularization parameter. For the regularization method to work we usually look for conditions that will guarantee that as $\alpha \rightarrow 0^+$ the solution of the perturbed problem will converge to a solution of the original problem [31]. Regularization and optimization are indeed closely interconnected. Next we see our proximal approach at work to achieve regularization of the Tikhonov type.

Using equations (5.3.1)-(5.3.2), we have

$$\nabla\varphi(x) = P^{-1}(x - \hat{x}) + \alpha L(x - \bar{x})$$

and algorithm 3.1 becomes

Algorithm 5.2. Choose $\gamma > 0$ and choose $\alpha \in]0, +\infty[$. Start with $x^0 \in \text{int } C$. For $\ell = 0, 1, \dots$ compute

$$x_i^{\ell+1} = x_i^\ell \exp \left(-\gamma \left((P^{-1}(x^\ell - \hat{x}))_i + \alpha (L(x^\ell - \bar{x}))_i \right) \right), \quad i = 1, \dots, N$$

until convergence.

The authors in [10] analyze a Tikhonov based spatial regularization method. They claim that their work is the first to describe how to incorporate spatial regularization

in Kalman filtering. We have implemented their method and compared it to our approach; our performs better. We detail this issue more in subsection 6.5.1. The drawback of Tikhonov regularization, though, is that it has too strong a smoothing effect and does not preserve edges. As a consequence, the reconstruction produces blurred images. Alternative methods exist for regularization of imaging problems. In addition to Tikhonov regularization, we propose a regularization by a function with better edge preserving properties.

5.4 Energy Function and Approximation

A cost function that involves a 2-norm as a regularizer, à la Tikhonov-Philips, is usually unsatisfying because many images are not globally smooth. They have region boundaries across which the image values change sharply. The quadratic regularization causes the edges to become blurred. Little variations between neighboring locations are due to noise while large variations are due to the presence of edges. This premise is the basis of most edge preserving regularization schemes including applications to tomography [4, 5, 53]. We need a cost function that favors local smoothness with well defined boundaries. We propose to use a 1-norm instead of the 2-norm in the penalty cost function ψ . Both penalty cost functions are convex functions. However, the 1-norm based one has the advantage that it increases less rapidly than the quadratic function for sufficiently large arguments since it is a linear increase instead of a quadratic one. Thus large differences between neighboring locations are penalized less than with the quadratic penalty. This uses local information to detect if an edge is present or not.

Like variance, mean squared error has the disadvantage of heavily weighting outliers. An example of such a use is KF which is based on MSE. This is again a result of the squaring of each term, which effectively weights large errors more heavily than

small ones. This property, undesirable in many applications, has led researchers to use alternatives such as the mean absolute error, or those based on the median. In the language of probability theory, the value of c that minimizes $\mathbb{E}(|X - c|)$ is the median of the probability distribution of the random variable X .

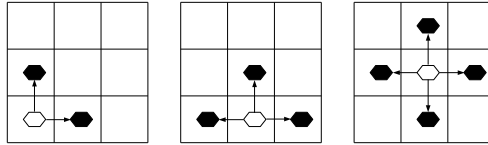


Figure 5.1: First order neighborhood configuration.

Let \mathcal{N}_i denote the set of indexes of voxels/pixels which are neighbors of pixel i . Define the energy function ψ by

$$\psi(x, m) = \sum_i \sum_{j \in \mathcal{N}_i} w_{ij} |x_i - m_j| \quad (5.4.1)$$

where $w_{ij} \geq 0$ are the neighborhood weights, m is a target image. From now on, we assume that $w_{ij} = 1 \forall j \in \mathcal{N}_i$ and is zero otherwise, in the sense that all neighboring locations have the same contribution. Therefore

$$\psi(x, m) = \sum_i \sum_{j \in \mathcal{N}_i} |x_i - m_j| \quad (5.4.2)$$

The function ψ is related to the Gibbs distribution in the Bayesian imaging context [43, 44]. This energy function does not penalize large differences between locations in the same neighborhood. We adopt a first order neighborhood, see Figure 5.1 for a 2D example. We refer the reader to [43, 50] for higher order neighborhoods. The absolute value function preserves the edges, e.g. abrupt changes in the image texture. This function penalizes deviations within uniform regions without necessarily penalizing the larger differences which occurs at the boundary between two different regions

of the image. This is an advantage over the Tikhonov based method. A connection exists between $|x|$ and the median as follow

$$\text{median}\{z_1, \dots, z_m\} = \arg \min_{s \in \mathbb{R}} \sum_i |s - z_i| \tag{5.4.3}$$

For instance, $\text{median}\{1, 1, 7\} = 1 = \arg \min_{s \in \mathbb{R}} (|s - 1| + |s - 1| + |s - 7|)$. The result (5.4.3) has been proven, that is, the median minimizes the sum of the absolute deviations [35, 48, 85]. Said otherwise, given a set of values z_1, z_2, \dots, z_m , the sum of absolute deviations is minimal when deviations are calculated from the median. For the continuous case, Cramér [35] considers the random variable ξ , the median μ , and an eventual in-between position θ . When $\theta > \mu$, he takes advantage of the relation

$$\mathbb{E}(|\xi - \theta|) = \mathbb{E}(|\xi - \mu|) + 2 \int_{\mu}^{\theta} (\theta - z) dF(z)$$

and that $\int_{\mu}^{\theta} (\theta - x) dF(z)$ is positive to prove that $\mathbb{E}(|\xi - \theta|)$ achieves its minimum value at $\theta = \mu$.

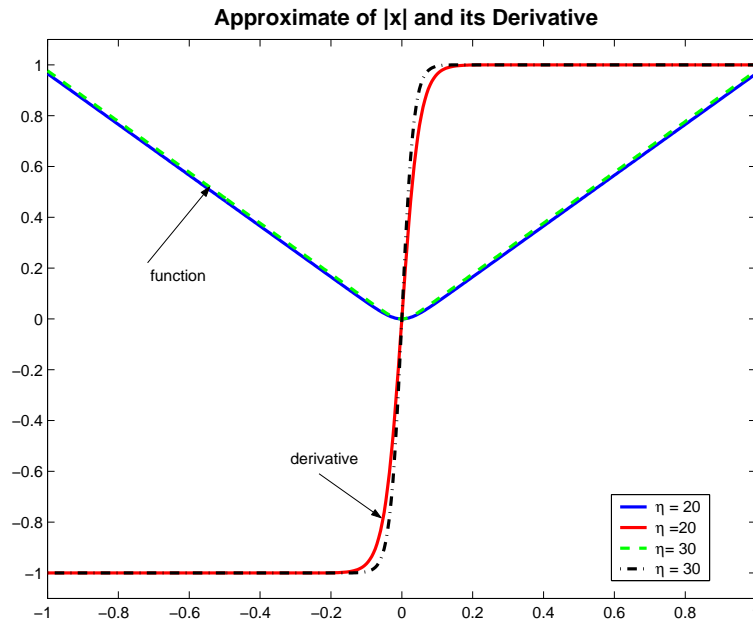


Figure 5.2: log cosh function and its derivative with two values of η .

We call “Median” the regularization involving the function ψ in (5.4.2). This regularization function is particularly suited to recover blocky images with sharp faces and

edges. Nonetheless, the absolute value function is convex but not differentiable where the voxel/pixel intensity is zero. Therefore, the optimization problem becomes non differentiable which is computationally impracticable. To circumvent this difficulty we approximate the absolute value with the function

$$\varphi_\eta(x) = \frac{1}{\eta} \log \cosh(\eta x) \quad (5.4.4)$$

which goes back to Green [47] as an extension of the Geman and McClure potential function [44]. There exists $\delta > 0$ such that when η is close to δ , then $\varphi_\eta(x) \rightarrow |x|$. Thus an appropriate choice of η may give a better approximation with numerical advantages in optimization. Note that $\varphi_\eta(x)$ is differentiable and its first derivative is given by $\varphi'_\eta(x) = \tanh(\eta x)$, see Figure 5.2. There exist other differentiable functions which approximate the absolute value quite well; take for instance $\varphi_\eta(x) = \sqrt{\eta^2 + x^2}$.

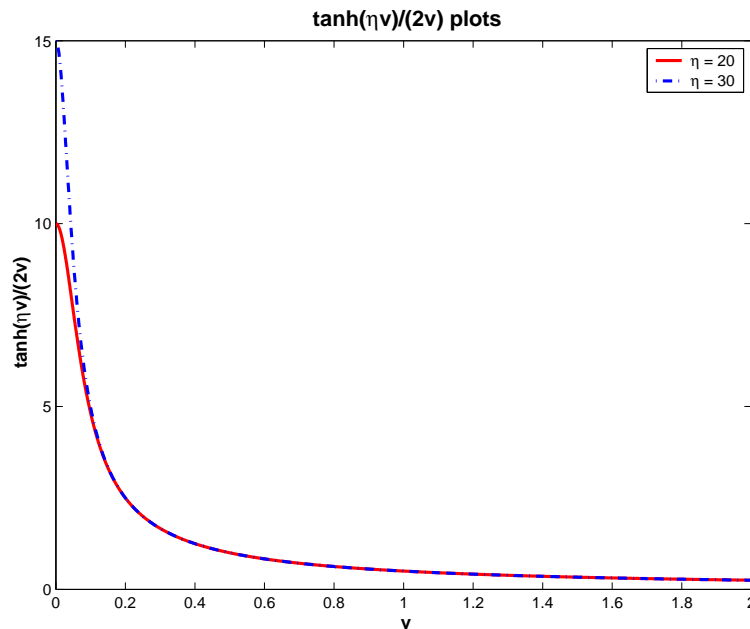


Figure 5.3: $\tanh(\eta v)/(2v)$ function with two values of η .

In order to encourage smoothing within a region and discourage smoothing across boundaries, Charbonnier et al [33] have suggested three conditions on the weighting function $\psi'(v)/(2v)$, namely

1. $0 < \lim_{v \rightarrow 0} \psi'(v)/(2v) = M$ to ensure isotropic smoothing in homogeneous areas.
2. $\lim_{v \rightarrow \infty} \psi'(v)/(2v) = 0$ to ensure preservation of edges.
3. $\psi'(v)/(2v)$ is strictly decreasing to avoid instabilities.

where M is finite. Tikhonov regularization is associated with $\psi(v) = v^2$ which is convex. However, $\psi'(v)/(2v) = 1$ does not satisfy the second and third conditions. Total variation is related to $\psi(v) = v$ which is convex, but $\psi'(v)/(2v) = 1/(2v)$ does not satisfy the first condition. Geman and McClure [44] use $\psi(v) = v^2/1 + v^2$ which verifies the three conditions but is not convex. To approximate the convex function we employ, $\psi(v) = |v|$, with the convex and differentiable function $\varphi_\eta(v) = \frac{1}{\eta} \log \cosh(\eta v)$. It has $\psi'(v)/(2v) = \tanh(\eta v)/(2v)$ which satisfies the three conditions with $M = \eta/2$, Figure 5.3 illustrates these facts.

Now we have assembled all the ingredients to study the numerical solution of problem (3.7.1). With ψ given in (5.4.2), we get

$$\varphi(x, m) = \frac{1}{2} \|\hat{x} - x\|_{P^{-1}}^2 + \alpha \sum_i \sum_{j \in \mathcal{N}_i} |x_i - m_j| \quad (5.4.5)$$

where \hat{x} is KF output activity and P is the covariance matrix of its error. We seek to minimize (5.4.5), hence

$$(x^*, m^*) = \arg \min_{x \geq 0, m} \varphi(x, m) \quad (5.4.6)$$

The function φ is continuous, nonnegative, convex, and coercive so that it has a global minimum (x^*, m^*) . This is a *joint estimation* of vectors x and m that we solve iteratively via an *alternating algorithm* as follow

$$x^{\ell+1} = \arg \min_{x \geq 0} \varphi(x, m^\ell) \quad (5.4.7)$$

$$m^{\ell+1} = \arg \min_m \sum_i \sum_{j \in \mathcal{N}_i} |x_i^{\ell+1} - m_j| \quad (5.4.8)$$

Applying the result (5.4.3) and rearranging the double sums in (5.4.8), we have

$$m_j^{\ell+1} = \text{median}\{x_i^{\ell+1}, i \in \mathcal{N}_j\} \quad j = 1, \dots, N$$

where \mathcal{N}_j is the set of indexes of locations which are neighbors of location m_j . Thus (5.4.7) becomes

$$x^{\ell+1} = \arg \min_{x \geq 0} \frac{1}{2} \|\hat{x} - x\|_{P^{-1}}^2 + \alpha \sum_i \sum_{j \in \mathcal{N}_i} |x_i - m_j^\ell| \quad (5.4.9)$$

In order to make the minimization problem differentiable, we employ the approximation (5.4.4), so that

$$x^{\ell+1} = \arg \min_{x \geq 0} \frac{1}{2} \|\hat{x} - x\|_{P^{-1}}^2 + \frac{\alpha}{\eta} \sum_i \sum_{j \in \mathcal{N}_i} \log \cosh(\eta(x_i - m_j^\ell)) \quad (5.4.10)$$

Let

$$\varphi(x) = \frac{1}{2} \|\hat{x} - x\|_{P^{-1}}^2 + \frac{\alpha}{\eta} \sum_i \sum_{j \in \mathcal{N}_i} \log \cosh(\eta(x_i - m_j^\ell)) \quad (5.4.11)$$

Fixing the index i and taking the partial derivative of φ w.r.t. x_i , we obtain

$$\frac{\partial \varphi}{\partial x_i}(x) = P^{-1}(x - \hat{x})_i + \alpha \sum_{j \in \mathcal{N}_i} \tanh(\eta(x_i - m_j^\ell)) \quad (5.4.12)$$

The general algorithm 3.1 yields the following alternating algorithm.

Algorithm 5.3. Choose $\gamma > 0$ and choose $\alpha \in]0, +\infty[$. Start with $x^0 \in \text{int } C$ and $m_j^0 = \text{median}\{x_i^0, i \in \mathcal{N}_j\}$ $j = 1, \dots, N$. For $\ell = 0, 1, \dots$ compute

$$\begin{aligned} x_i^{\ell+1} &= x_i^\ell \exp \left(-\gamma \left((P^{-1}(x^\ell - \hat{x}))_i - \alpha \sum_{j \in \mathcal{N}_i} \tanh(\eta(x_i^\ell - m_j^\ell)) \right) \right) \quad i = 1, \dots, N \\ m_j^{\ell+1} &= \text{median}\{x_i^{\ell+1}, i \in \mathcal{N}_j\} \quad j = 1, \dots, N \end{aligned}$$

until convergence.

Tikhonov regularization and median regularization differ only in the norm they use, the former uses the 2-norm and the latter uses the 1-norm. Next, we introduce a more generalized context of regularization.

5.5 Hölder Filter

Now, we describe some regularization schemes depending on a different choice of $\psi(v)$. When using Tikhonov regularization, we have $\psi(v) = \|Lv\|$, where L is a differential operator. The median regularization function is $\psi(v) = \sum_i \sum_{j \in \mathcal{N}_i} |v_i - \hat{x}_j|$; whose approximation is $\psi_\eta(v) = \sum_i \sum_{j \in \mathcal{N}_i} \varphi_\eta(v_i - \hat{x}_j)$. Another type of regularization function close to the Tikhonov family is the mean regularization, where $\psi(v) = \sum_i \sum_{j \in \mathcal{N}_i} |v_i - \hat{x}_j|^2$. The two former regularization operators belong to a more general family of filter operators that we call *Hölder filter*. Typically, the Hölder filter replaces each location by the Hölder mean of its neighborhood. That is,

$$\tilde{v}_i = \arg \min_{z \in \mathbb{R}} \sum_{j \in \mathcal{N}_i} |z - v_j|^p$$

where $1 \leq p < +\infty$. Choosing p such that $1 \leq p < +\infty$ makes the above functional convex and differentiable except the case $p = 1$ at the origin. This latter instance was dealt with in section 5.4. The Hölder mean filter transforms an image v to a new one given by $\tilde{v} = M^p(v)$ such that

$$M^p(v) = \arg \min_{u \in \mathbb{R}^N} \sum_i \sum_{j \in \mathcal{N}_i} |u_i - v_j|^p$$

It is straightforward to show that $\tilde{v}_i = M_i^p(v)$. For the sake of clarity, we write $\tilde{v}_i = M^p(v_j, j \in \mathcal{N}_i)$. We use this filter in two special cases, $p = 1$ where the Hölder mean coincides with the median, and $p = 2$ where we find the arithmetic mean. The median and mean regularization methods are based on the following alternating minimization formulation,

$$\min_{v, u \in \mathbb{R}_+^N} \left(\frac{1}{2} \|\hat{x} - v\|_{P^{-1}}^2 + \alpha \sum_i \sum_{j \in \mathcal{N}_i} |u_i - v_j|^p \right) \quad (5.5.1)$$

An optimal solution (v^*, u^*) is such that $v^* = \text{prox}_{\psi_p}^\alpha(u^*)$, and $u^* = M^p(v^*)$, which can be computed by using the following iterations

$$v^{m+1} = \text{prox}_{\psi_p}^\alpha(\hat{x}), \quad u^m = M^p(v^m), \quad m = 1, 2, \dots$$

where the starting guess is $v^0 > 0$, $\psi_p^m(v) = \sum_i \sum_{j \in \mathcal{N}_i} |u_i^m - v_j|^p$, and $p \in \{1, 2\}$. We could utilize our algorithm 3.1 to solve $v^{m+1} = \text{prox}_{\psi_p^m}^\alpha(\hat{x})$. Notice that the incorporated target value in the regularization function is not fixed; it is updated during the iteration. Foundations of this alternating minimization, including convergence, are provided in [12].

5.6 Segmentation Regularization

In medical imaging, we are sometimes not interested in individual intensities of each and every pixel/voxel but rather on some ROI (region of interest) intensities. We are then more concerned with a segmented reconstruction [81]. A CT scan for instance might give us an idea about the ROI. In case we have this prior knowledge about the selection of ROI before hand, we could include this constraint, reduce the size of our problem, and have by the same token a better spatial regularization. A commonly used approach is to proceed through a change of variable, see for instance [28]. Let ξ be the activity vector of the disjoint p ROI. Let E represent the $N \times p$ belonging matrix of each pixel to a unique ROI. It has therefore only one 1 in every column and row and the rest of the entries are zeros. A 1 in row i and column j implies the i^{th} pixel belongs to the j^{th} ROI. Consequently we have the following relation

$$x = E\xi \tag{5.6.1}$$

Instead of a transition and an evolution models for the activity x , we have rather similar ones for the activity ξ . Hence equations (2.2.1) and (2.2.2) write

$$\xi_k = \tilde{A}_k \xi_{k-1} + \tilde{\mu}_k \tag{5.6.2}$$

$$y_k = \tilde{C}_k \xi_k + \tilde{v}_k \tag{5.6.3}$$

where $\tilde{C}_k = C_k E$. Therefore, in lieu of solving for a bigger size x , we solve for a much smaller size ξ . Algorithm 5.1 becomes

Algorithm 5.4. Choose $\gamma > 0$ and start with $\xi^0 \in \text{int } \mathbb{R}_+^p$. For $\ell = 0, 1, \dots$ compute

$$\xi_i^{\ell+1} = \xi_i^\ell \exp\left(-\gamma(P^{-1}(\xi^\ell - \hat{\xi}))_i\right), \quad i = 1, \dots, p$$

until convergence to ξ^* . Then set

$$x^* = E\xi^*$$

Chapter 6

Numerical Experiments

Previously, we have stated the inverse problem of reconstructing a medical image in nuclear medicine. We have employed the KF to give us an initial estimate that we project onto the nonnegative orthant. We have proven some properties of our nonnegative estimator and introduced spatial regularization schemes. Next, we corroborate the effectiveness of the developed algorithms using a digital phantom. Some of the presented theoretical foundation and validating experiments are reported in [79].

6.1 Procedure

Assume we are provided with projection matrices C_1, \dots, C_S and projections y_1, \dots, y_S . We give a systematic method on how we can implement our approach.

Procedure 6.1.

step 1 Start with an initial guess vector $\hat{x}_{0|0}$ and an initial covariance matrix $P_{0|0}$.

For $k = 1, \dots, S$, execute step 2, 3, and 4

step 2 At the k^{th} recursion, choose the covariance matrices Q_k and R_k (section 2.3); examples of how are given in section 6.3

- step 3** Use the Kalman filter algorithm, Eqs. (2.3.1) to (2.3.5), to calculate $\hat{x}_{k|k-1}$ and $\hat{x}_{k|k}$
- step 4** Use algorithm 3.1, in one of its particular forms developed in chapter 5, to achieve nonnegativity and/or smoothness of $\hat{x}_{k|k}$ if necessary
- step 5** For $k = S - 1, \dots, 1$, use the Kalman smoothing algorithm, Eqs. (2.3.6) to (2.3.8), to calculate the estimate $\hat{x}_{k|S}$ and
- step 6** Use algorithm 3.1, in one of its particular forms developed in chapter 5, to achieve nonnegativity of $\hat{x}_{k|S}$ if necessary.

6.2 Simulation

Our phantom is composed of six regions of interest (ROI) or segments. Each ROI has a different time activity curve (TAC), see Figure 6.1. The example investigated in this work is based on the teboroxime dynamics in the body during first hour post injection. The choice of the time activity curves (TACs) is motivated by the behavior of liver, healthy myocardium, muscles, stenotic myocardium, and lungs, respectively. Only one slice is modeled; that is we simulate a 2D object. The star-like shape placed on the left ensures that the phantom is not entirely symmetrical. We simulate 120 projections over 360° , one projection for every 3° with attenuation and a 2D Gaussian detector response.

There are three camera heads consisting of 64 square detectors each measuring 0.625 cm in each side, see Figure 1.1. The distance from the annulus to the detector-rotation axis is 30 cm. We simulate 40 time instances for three heads; that is we have $3 \times 40 = 120$ projections for a camera rotating clock wise (CW) in a circular orbit. Head 1 starts at -60° , head 2 at 60° , and head 3 at 180° . A low energy high resolution (LEHR) collimator is used with a full width at half maximum (fwhm). We determine

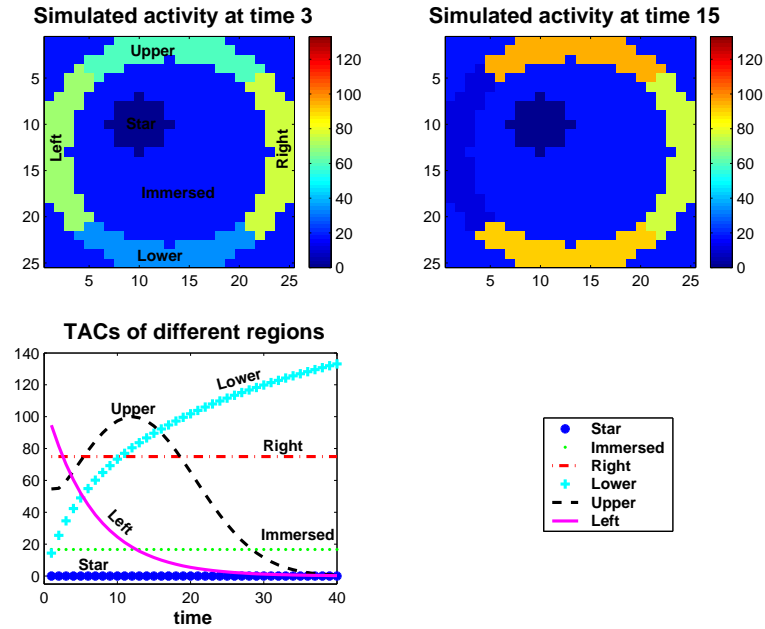


Figure 6.1: Simulated annulus with its different ROI and their TACs: (a) simulated activity at time 3, (b) at time 15, and (c) TACs of the 6 different ROI.

the blurred parallel strip/beam geometry system matrices for all projections with resolution recovery and attenuation correction [93, 94].

We have 64 projection values for each head, which amounts to a total of 192 observations at each time frame. The size of the image we aim to reconstruct is $625 = 25 \times 25$ dixels; this is an under-determined problem with a ratio of 1:3.25 of data to unknowns. It is an ill-posed problem. We have experimented with one-head as well as with two-head camera: This makes the problem even more ill-posed with a 1:9.76 and a 1:4.88 ratio respectively. We have six kinds of TACs that are very representative for clinical applications. The annulus has four arcs that we name “Left”, “Upper”, “Right”, and “Lower” according to their location. The activity is decreasing in the Left arc, increasing-decreasing in the Upper arc, constant in the Right arc, and increasing in the Lower arc; see Figure 6.1. The star-like shape has zero activity within it and is called the “Star” region; we refer to it as “Background” too. The annulus is immersed within a region that has the sky-blue color in our figure.

It is called “Immersed” and has a constant activity. We have six ROIs in total.

6.3 First Tests

We provide quantitative analysis of the reconstructed images in order to compare the simulated activity with the reconstructed one. We define the relative deviation error τ of the reconstructed activity v^* from the truth x , refer to (6.3.1) through (6.3.3). Hence we compare the simulated count $x_{i,k}$ with the corresponding reconstructed one $v_{i,k}^*$ at each time frame k for every location i . We sum over a ROI containing J dixels normalized by the total simulated/true counts in order to diminish the effect of statistical fluctuations. We have a $\tau_{ROI,k}$ for every sector. These indicators allow us to see how the method performs under different dynamic behaviors. We could compare, for instance, sectors with fast washout with those with slow one [17]. We calculate similar τ_k over the total number of doxels (dynamic voxels) N then we average them over the total number S of time acquisitions; so that we have τ_{avg} . This is an objective comparison of the quality of reconstruction for different sets of parameters such as iteration stopping criteria, noise levels, etc. The closer τ_{avg} is to zero, the better the reconstructed images should be.

$$\tau_{ROI,k}^2 = \frac{\sum_{i=1}^J (v_{i,k}^* - x_{i,k})^2}{\sum_{i=1}^J x_{i,k}^2} \quad (6.3.1)$$

$$\tau_k^2 = \frac{\sum_{i=1}^N (v_{i,k}^* - x_{i,k})^2}{\sum_{i=1}^N x_{i,k}^2} \quad (6.3.2)$$

$$\tau_{avg} = \frac{1}{S} \sum_{k=1}^S \tau_k \quad (6.3.3)$$

Preliminary tests are performed to choose initial covariance matrices, initial activity, number of iterations in the projection algorithm, as well as regularization tuning parameters. Our guide to distinguish between a variety of them is the τ value combined with visual inspection; this is a heuristic approach. The parameter that gives

us the least value of τ_{avg} in (6.3.3) and/or better smoothed images will be our choice.

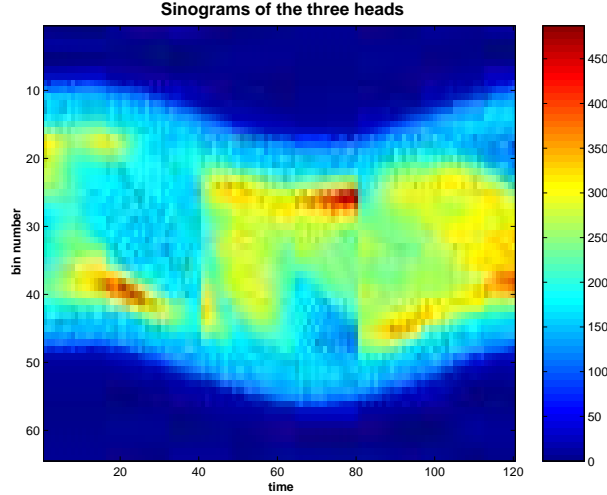


Figure 6.2: Sinogram or $2D$ projections: y -axis has the bin number and the x -axis has the 40 time instances of the 3 heads. Time instances from 1 to 40 are for head 1, 41 to 80 for head 2, and 81 to 120 for head 3. A color intensity of a pixel is the number of detected photons by a certain bin at a certain time.

To see the importance of the number of heads, we experimented with a camera that has one, two, and three heads to see if there is any improvement, see sinogram in Figure 6.2. We apply the procedure 6.1 with algorithm 5.1 in step 4 and 6 where we choose $\gamma = 1$, which is our choice from now on. The parameter γ is chosen to make the function $f - \gamma\varphi$ convex to ensure convergence of the algorithm 3.1. The choice $\gamma = 1$ is heuristic and the rationale behind it goes as follows.

The convex function f is the negative entropy and the convex function φ is quadratic. Consider the Hessian of $f - \gamma\varphi$ which is of the form

$$\text{diag}\left(\frac{1}{x}\right) - \gamma P^{-1} \quad (6.3.4)$$

where $\text{diag}(\frac{1}{x})$ is the diagonal matrix having the components $\frac{1}{x_i}$ in its main diagonal. Numerical experiments show that the matrix P^{-1} is almost diagonal with values of the order 10^{-3} or so. Recall that x is a nonnegative variable. When x_i is sufficiently large,

this Hessian is negative definite. That is any choice of γ will not render the function $f - \gamma\varphi$ convex for all values of x . However, we are interested only in convexity when the activity x takes some finite values. In our numerical experiments at hand, the maximal value of x_i is less than 150 or so. Thus the Hessian takes the approximate form

$$\text{diag}\left(\frac{1 - \gamma P_n^{-1} x_i}{x_i}\right) \quad (6.3.5)$$

The term $P_n^{-1} x_i$ is less than 1 and the choice of $\gamma = 1$ would make the Hessian positive definite.

With the choice of these parameters, the TACs look similar in shape to the true activity for the three reconstructed images; however this shape similarity phenomenon is more pronounced with three than with one head since we have less statistics, that is fewer photons detected by the camera when we have fewer heads; see Figure 6.3. The τ_{avg} value in the case of one head is 0.50. Images were also less noisy and TACs were closer to the true ones, for instance, with two than with one-head camera. TACs with two and three heads are almost the same; however, we have a slightly smaller τ_{avg} value of 0.37 with three compared to 0.40 with two. The results we present from now on are obtained with a triple-head camera.

We mentioned in section 2.2.1 that the flow of the radioactivity could be seen as a diffusion model which is related to a random walk model [56]. The latter model is very useful if we do not possess enough knowledge about the radioisotope substance flow. We assume that the system dynamics are unknown to us (2.2.1); therefore we use a random walk. In practical terms, we set $A_k = I$, for all $k = 1, \dots, S$. For the state transition linear model, we proceeded as follows. We are not interested in the background and we assume that we know the locations of these zero activities; this is a common practice [40, 63, 81]. We have run experiments without this assumption and results are very comparable to when we have run them with this assumption. One interesting way to deal with this assumption is as this. Set to zero the values of

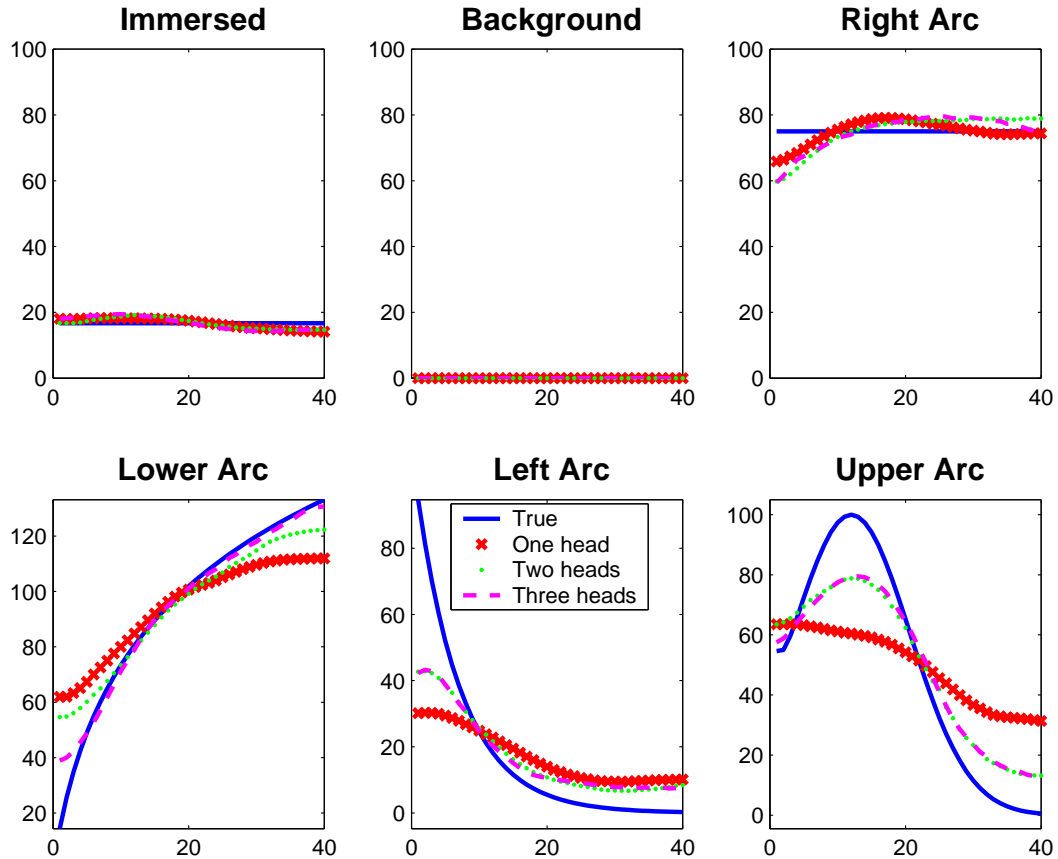


Figure 6.3: Reconstructed TACs with one, two, and three heads camera.

the corresponding positions of the matrix C_k and of the initial guess. The updating equations 2.3.1, 2.3.3, and 2.3.6 ensure that the updated activities will remain equal to zero; thus the KF reconstructs perfectly the star/background region(s).

We experiment with different initial guesses $\hat{x}_{0|0}$ such as $(10^{-6}, \dots, 10^{-6})^\top$ and $(1, \dots, 1)^\top$. We also start the algorithm with the static image given by OSEM where the background activity is set to zero; we call this initial guess *OSEM act*. We likewise run KF backward in time starting with an arbitrary starting image and use its final output at time 1 as another choice for $\hat{x}_{0|0}$. The average of the deviation error τ_{avg} combined with visual inspection show that there is a slight advantage in favor of the *OSEM act*, especially in rendering the edges of the background.

We do not have much confidence in our initial activity guess $\hat{x}_{0|0}$ so we choose the initial covariance matrix to be pretty large, $P_{0|0} = 10^5 I$. We have also tried the last covariance matrix P coming out from KF run backward in time; no improvement has been noticed though. In any case, KF has the interesting property that the effect of the initial values of $\hat{x}_{0|0}$ and $P_{0|0}$ diminishes over time, already after the first two or three steps of the recursion. The update formulas for the covariance matrix P ensure its symmetry in theory ((2.3.1), (2.3.3), and (2.3.6)) at each step. However, numerical calculations might introduce nonsymmetry into them; so we substitute $\frac{1}{2}(P + P^\top)$ for P at each step to ensure symmetry.

Recall that $\mathbb{E}(\mu_k) = 0$, and $\mathbb{E}(\mu_k \mu_k^\top) = Q_k$. We experimented with

$$Q_k = \sigma^2 \begin{pmatrix} 1 & \rho & \rho^2 & \rho^3 & \cdots \\ \rho & 1 & \rho & \rho^2 & \cdots \\ \rho^2 & \rho & 1 & \rho & \cdots \\ \rho^3 & \vdots & \vdots & \ddots & \rho \\ \vdots & \cdots & \rho^2 & \rho & 1 \end{pmatrix}$$

where ρ is some small positive value, for example, 10^{-1} . The higher the power associated with ρ the farther are pixels from each other. This choice made no difference compared to the simpler one by choosing $Q_k = Q = \sigma^2 I, \forall k \in \{1, \dots, S\}$.

Meanwhile $\mathbb{E}(\nu_k) = 0$ and $\mathbb{E}(\nu_k \nu_k^\top) = R_k$, so we take $R_k = \text{diag}(y_k)$, at each step k of the recursion, y_k being the projection data or number of detected photons at time k . Recall that, section 1.1, the data y_k is an instance of a Poisson random variable Y_k with mean $\mathbb{E}(Y_k) = y_k$ and standard deviation $\sigma(Y_k) = \sqrt{y_k}$.

We present here results obtained with $\hat{x}_{0|0} = \text{OSEM act}$ and $P_{0|0} = 10^5 I$ only. In contrast to other approaches, [81] for instance, we do not assume here the segments to be known exactly. We make use of these segments only to interpret the results. As a consequence, there are some differences in intensity between pixels within the same region. To assess the effectiveness of the method, we show the TACs averaged over

the pixels within the same ROI.

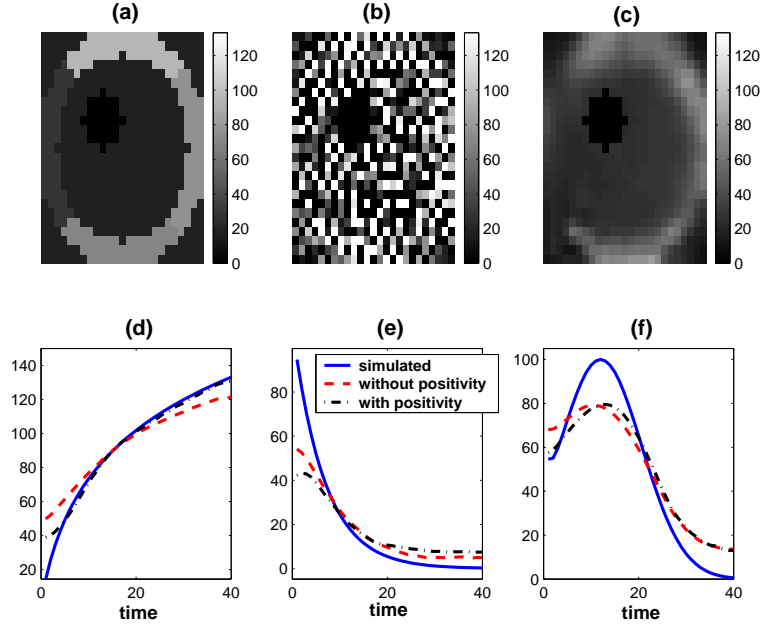


Figure 6.4: Without and with positivity reconstructed images and TACs: (a) simulated at time 9, (b) reconstructed without enforcing positivity, (c) reconstructed with enforcing positivity, (d), (e) and (f) TACs of Lower, Left, and Upper regions respectively.

6.3.1 Results

Right after each time step k of KF, we took the absolute value of the reconstructed activity, $\text{abs}(\hat{x}_k)$, and, in a second experiment, we set to zero its negative values; that is we took $\max(\hat{x}_k, 0)$ where \hat{x}_k is the Kalman output. We obtained a τ_{avg} equal to 2×10^{29} and 3×10^{17} respectively. We repeated the experiment in applying $\text{abs}(\hat{x})$ and $\max(\hat{x}, 0)$ only once at the end of the algorithm. We got high τ_{avg} as 2.12 and 1.80 respectively. More to the point, the images were unidentifiable in all these four cases. Recall that using $\max(\hat{x}_k, 0)$ means we apply the orthogonal projection. Consequently the fact that we did not get any meaningful image confirms theorem 4.1 that the projected KF estimator x_{p-1}^* is the constrained ML estimator of x in $C = \mathbb{R}_+^N$ within

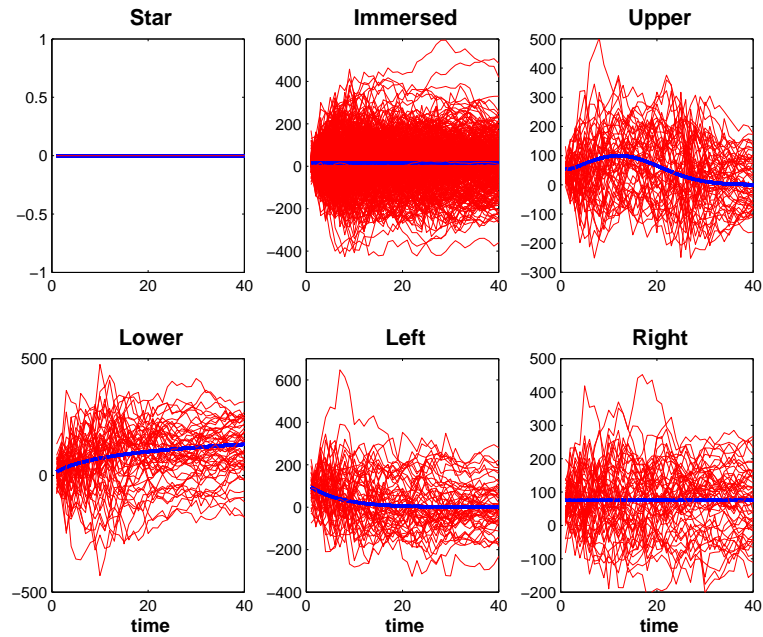


Figure 6.5: All reconstructed TACs without positivity. Blue — for true/simulated TAC, red — for reconstructed TACs. Note the scales.

the framework of a linear model and Gaussian pdf given by equation 4.3.1; refer also to section 3.7 for more details. We observe in Figure 6.4 that the algorithm reconstructs perfectly, as it should, the star region without or with enforcing any kind of positivity; refer to section 6.3. This is evidenced by the fact that the reconstructed TACs in the star region are exactly the same as the simulated one; refer for instance to Figure 6.5 and Figure 6.6. We ran KF without enforcing the positivity; that is without applying our proximal method. We then applied KF with positivity using algorithm 5.1 and compared both outputs. The averaged TACs in every region look similar without and with positivity which is in accordance with theorem 4.3 emphasizing the unbiasedness of x^* and \hat{x} , see Figure 6.4. This is explained by the fact that KF gives an optimal estimate on average at each time k for every region but not for every voxel. However, we got only meaningless images without positivity. This is more apparent when we plot all TACs over different regions for both without and with positivity, see Figure 6.5 and Figure 6.6. When you compare the figures, please note the difference in the scales.

Note for instance how some values are negative and go below and above the true ones by about 500 without the positivity; while values stay pretty close to the true ones and above zero with the positivity. This is reflected in the average over time τ_{avg} of 2.42 without positivity and 0.37 with positivity. Theoretical results about these observations have been established in section 4.4.4. It is clear that our approach of enforcing positivity in the output images, that come from the classical KF algorithm, is better than using the *abs* and *max* functions or than just doing nothing. The proximal approach to enforce nonnegativity is indeed an efficient tool to enforce some spatial regularization, refer to section 5.1 for more details including theorem 4.4 stating that the estimator x^* performs better in the MSE sense.

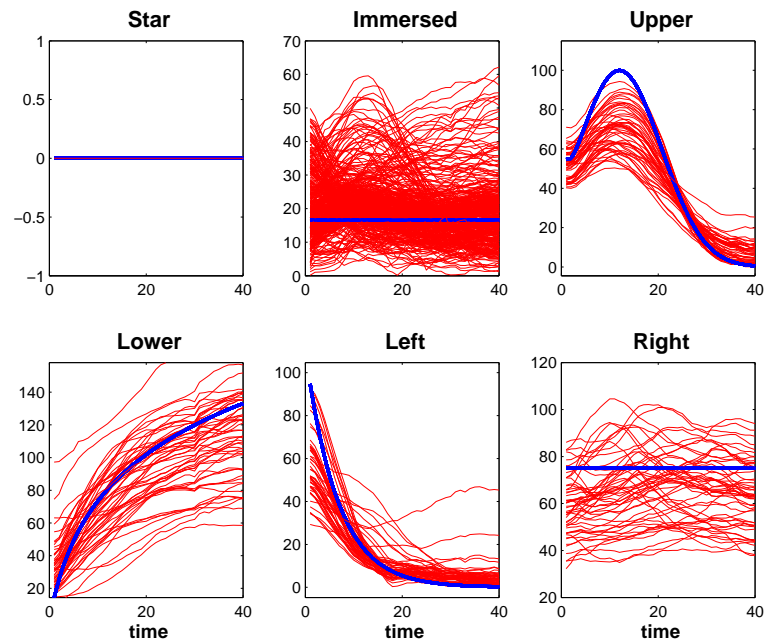


Figure 6.6: All reconstructed TACs with positivity. Blue — for true/simulated TAC, red — for reconstructed TACs. Note the scales.

We like to see how our algorithm behaves in the presence of good and bad signals, or in our setting, noiseless and noisy data, refer to Figure 6.7 and compare it to Figure 6.2. We run tests where we compare reconstructions with and without noise included into the data/observations. Instead of working with the observation y as is

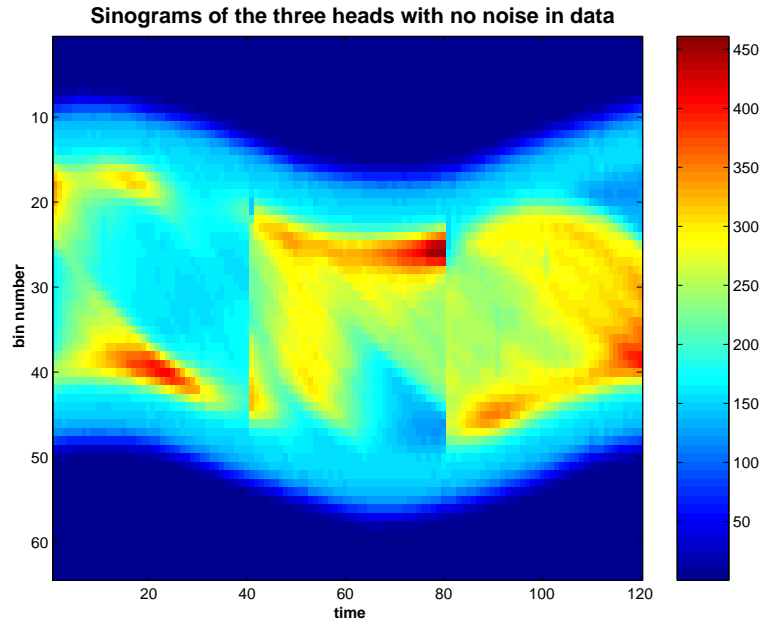


Figure 6.7: Noiseless sinogram or $2D$ projections without noise: y -axis has the bin number and the x -axis has the 40 time instances of the 3 heads. Time instances from 1 to 40 are for head 1, 41 to 80 for head 2, and 81 to 120 for head 3. A color intensity of a pixel is the number of detected photons by a certain bin at a certain time..

in the case of noiseless data, we took rather a Poisson random observation with mean y in the case of noisy data. We notice that there are very slight differences in the TACs and in the reconstructed images with and without noisy data, see Figure 6.8. We include, for reference, different plots for the τ function for both noisy and noiseless data, see Figure 6.9 and Figure 6.10. The τ_{avg} value at 0.36 without noise is slightly smaller than the one with noise, which is 0.37. This should come at no surprise since the less noisy the data the better the method should perform. However, the approach is not very sensitive to noise in the data; it filters out the noise from data very well. In general, the τ function is somehow “decreasing” over time (plot of “ τ in all” of Figure 6.9); which means that our algorithm improves in time as it should.

Theorem 4.1 posits that the estimator x^* is the ML of the nonnegative activity x in \mathbb{R}_+^N w.r.t. P^{-1} within the context of the Gaussian pdf, given by equation 4.3.1, using a linear model. To see the importance of the oblique projection through the

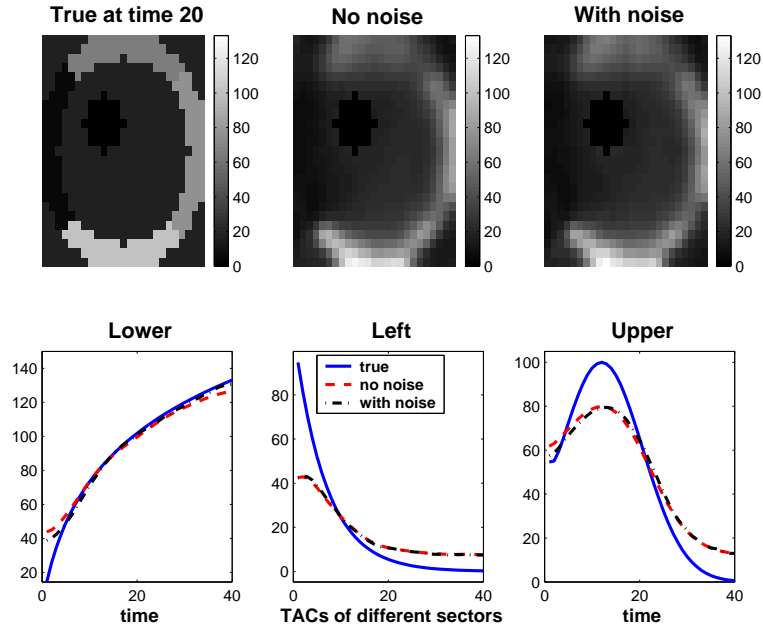


Figure 6.8: Reconstructed TACs and images with and without noise in data.

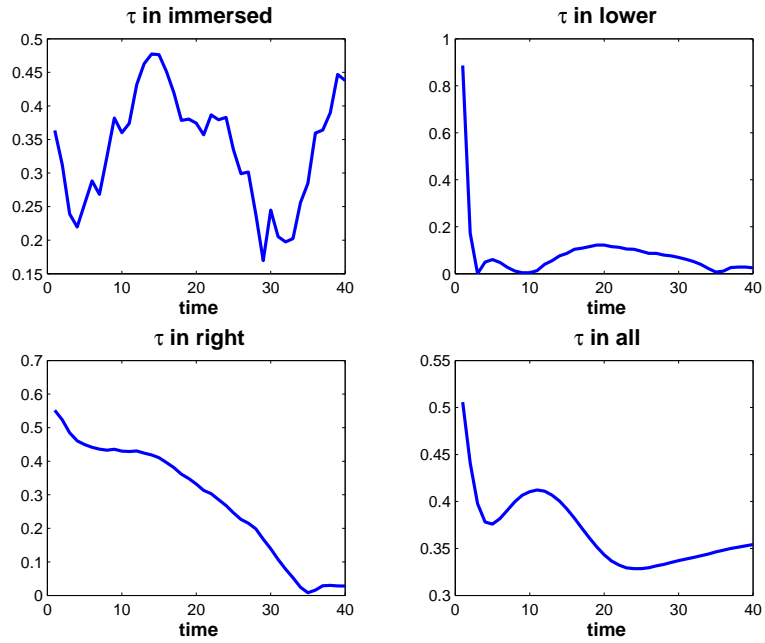
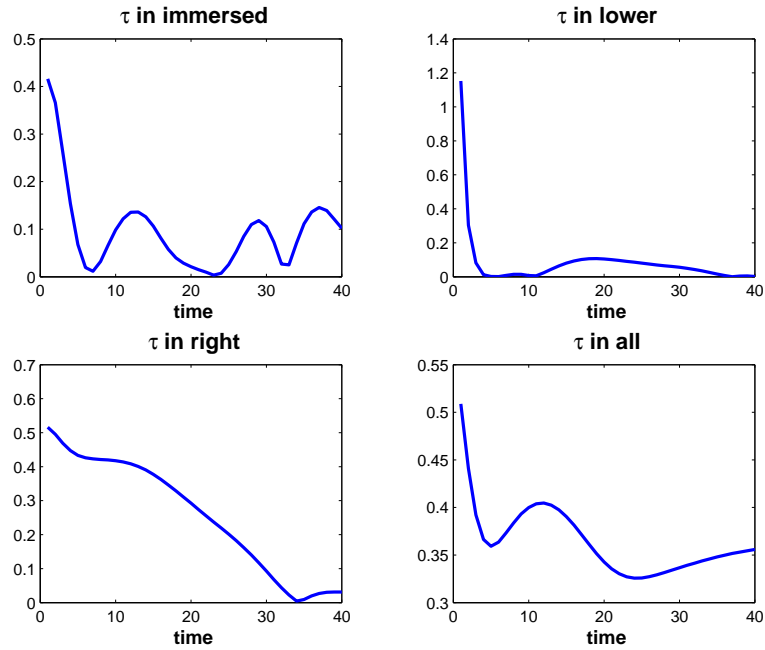


Figure 6.9: τ function with noisy data.

Figure 6.10: τ function for data without noise.

weighting matrix $W = P^{-1}$, we experimented with different matrices. We tried $W = I$, orthogonal projection, and $W = P$; neither choice leads to any meaningful reconstructed image. We also used $W = P^{-1} + \epsilon\Upsilon$ as a slightly perturbed P^{-1} . The matrix Υ was either a uniform random matrix within the interval $[0, 1]$, a Gaussian random matrix, or just the matrix having the number 1 in all its entries. We varied ϵ from 10^{-6} up to 10^{-3} . In all these cases the average deviation was greater than 0.37, the one with $W = P^{-1}$, starting from 0.37 up to 0.81. We did not get any meaningful image starting from $\epsilon = 10^{-3}$ and up. This suggests that the symmetric positive definite matrix P^{-1} is indeed the optimal weighting matrix.

6.4 Positive Kalman

Our first task was to determine the tuning parameters. Recall that in section 2.3 we choose to model the noise as white noise, that is $Q_k = \sigma^2 I$. The Kalman algorithm

gives two reconstructed images, one after the filtering step that we call “filtered”, another one after the smoothing step, that we call “smoothed”, refer to section 2.3. We are interested in the behavior of dynamic regions. Therefore, from now on we only show the TACs of the lower, left, and upper arcs. Using τ_{avg} as our primary criterion to discriminate for potential values and/or visual inspection while varying one parameter at the time thus keeping the rest constant, we found first that there is not much of a τ_{avg} gain when increasing or decreasing the value of σ^2 ; see Figure 6.11 for τ_{avg} as a *function* of σ^2 . We have a τ_{avg} around 0.36 for σ^2 between 4 and 4×10^8 ; thus the choice of a log scale for our x-axis. We have a minimum value of 0.37 when $Q \approx 40I$. The fact that our algorithm is not sensitive to the value of the error covariance matrix Q shows that it handles already “some” of the ill-posedness of our problem.

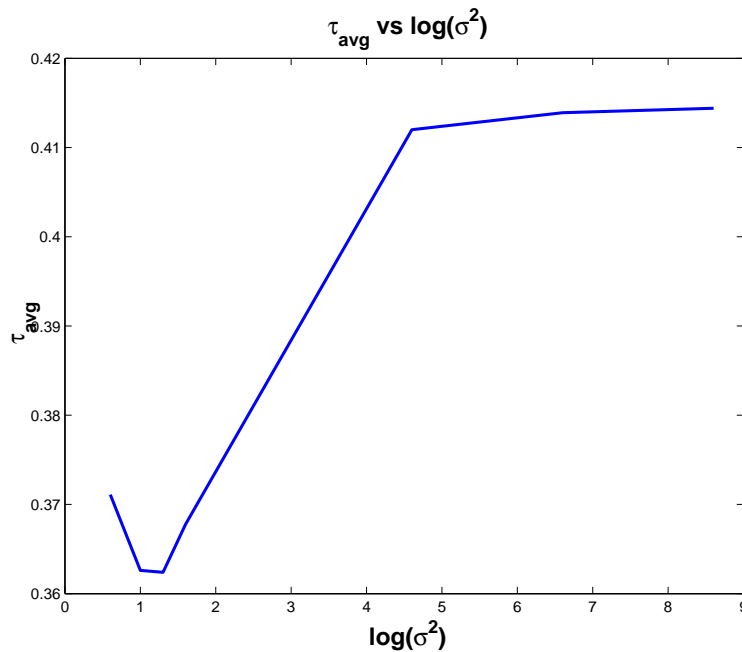


Figure 6.11: τ_{avg} as a function of $\log(\sigma_Q^2)$.

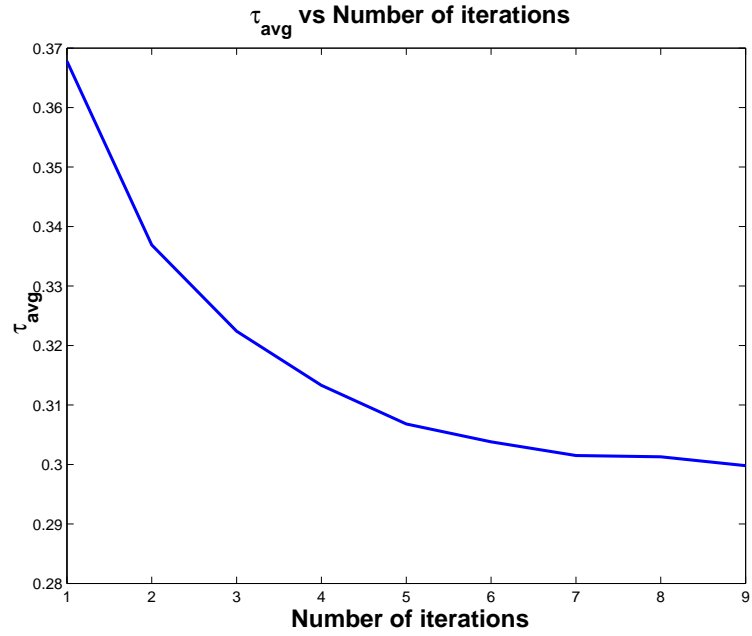


Figure 6.12: τ_{avg} as a function of the number of iterations.

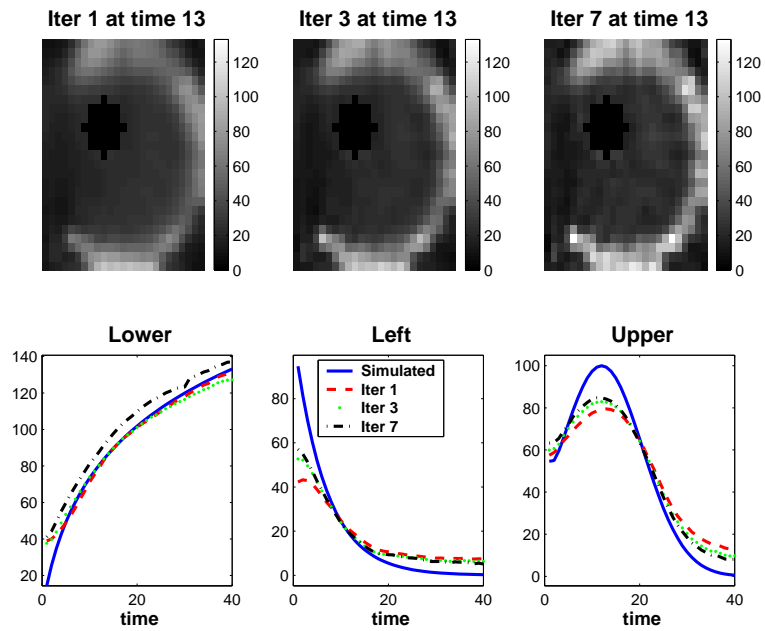


Figure 6.13: Image at time 13 and TACs at different number of iterations.

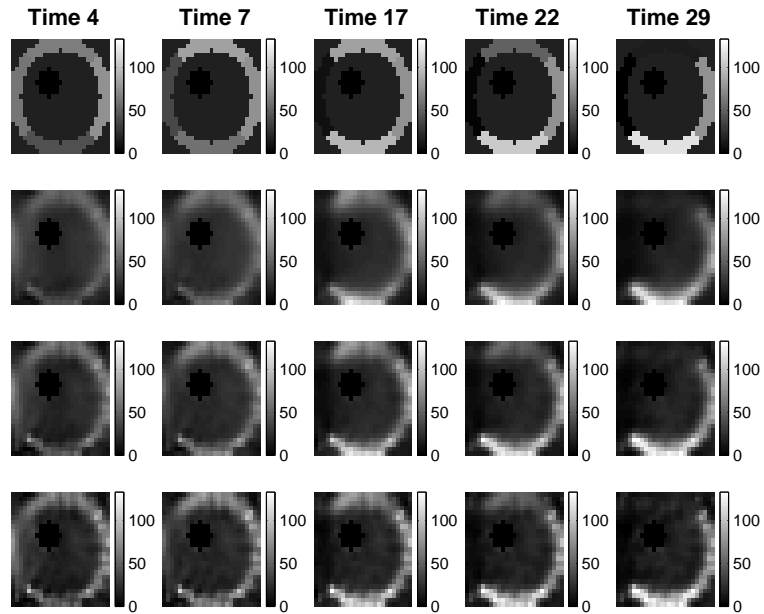


Figure 6.14: Images at different number of iterations: Truth in 1st row, images at 1st iteration in 2nd row, images at 3rd iteration in 3rd row, and images at 7th iteration in 4th row.

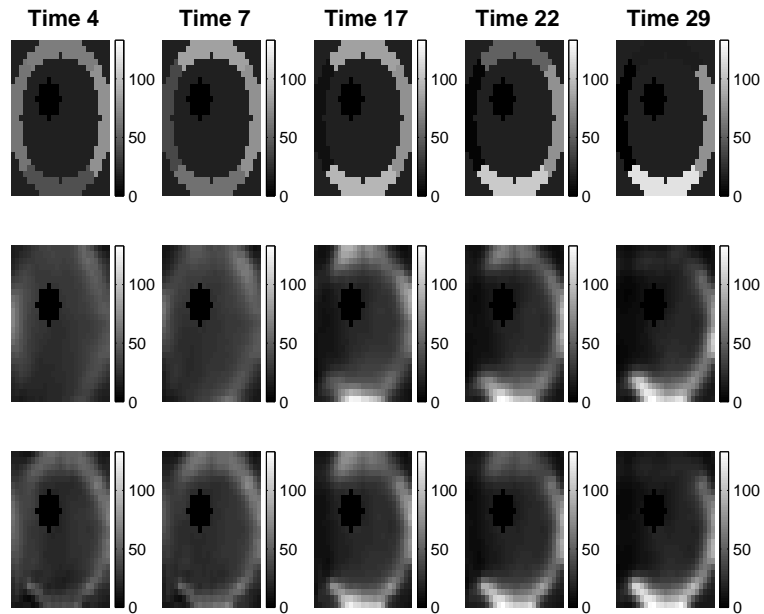


Figure 6.15: Nonnegative images at various times: Truth in 1st row, filtered image in 2nd row, smoothed image in 3rd row.

The projected Bregman algorithm to ensure nonnegativity is iterative in the time-domain. We use the same heuristic approach to find the number of iterations to be carried out. Iterative algorithms should give better results as we run them longer, that is, as we increase the number of iterations. In our case, we notice that τ_{avg} goes from 0.37 at the first iteration down to about 0.30 at the ninth iteration. The deviation τ_{avg} is indeed improving and Figure 6.12 illustrates this fact. However, images are getting noisier while the TACs are not improving much, Figure 6.13 testifies to that. Observe also how the reconstructed images in Figure 6.14 deteriorate as we increase the number of iterations. We are witnessing a pointwise but not a uniform convergence. A similar phenomenon is observed with the EM algorithm where we have semi-convergence instead of a normal “full” convergence. Indeed, it was noticed that we should stop EM earlier on in the iterations before the solution starts to fit to the noisy data. Section 5.2 details how an iterative scheme, as it is our case here, could be used as a tool to impose some spatial regularization. Subsequently, we only apply one iteration from now on as a trade off between reasonable τ_{avg} and smooth images. This is of course a subjective choice and one could run the algorithm for more than one iteration.

Figure 6.15 depicts images of the simulated/true annulus at various times together with the reconstructed filtered and smoothed ones when we enforce the positivity. Images look fine, however, the left arc looks noisy at the first few instances of time. This could be explained by the fact that this is where we have a rapidly decreasing activity and the algorithm takes longer before it catches on. Reconstructed smoothed images look smoother than the filtered ones as expected. We notice even an improvement intensity-wise since they are closer to the true ones; compare for instance the upper region at different times and how images look closer in color to the true ones than the filtered images. This is reflected in τ_{avg} where the smoothed reconstructed image has $\tau_{avg} = 0.37$ while the filtered one has $\tau_{avg} = 0.42$, see Table 6.2.

Figure 6.16 displays the averaged TACs over three different regions. We plot

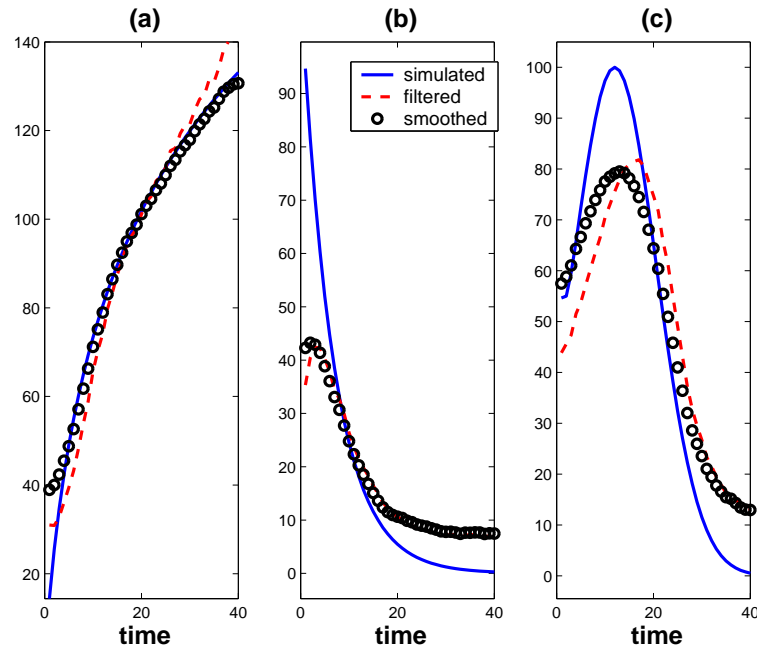


Figure 6.16: Averaged TACs for each region with positivity: Blue TACs for truth, red TACs for reconstructed filtered, and black TACs for smoothed. (a), (b), and (c) TACs for lower, left, and upper arcs respectively.

the true ones, shown in blue, the reconstructed filtered ones, shown in red, and the reconstructed smoothed ones, shown in black. Note that both reconstructed TACs look pretty close to the true ones shape-wise and in quantity/intensity/color. This is very interesting since we use only a basic approximation, namely first-order random walk, to describe the evolution model. Smoothed TACs look indeed “smoother” than the filtered ones as promised by the smoothing step in the Kalman algorithm. There are however some differences in intensity between pixels within the same region, see for instance the right and lower arcs. Smoothing is done over time with KF; that is we have a temporal regularization but not a spatial one. Tikhonov spatial regularization is the topic of the next section.

Table 6.1 lists CPU times taken to run our algorithm. We experimented on a P4 3.00 GHz desktop. All times are given in seconds (sec). It takes about 79.3 sec to run the whole algorithm; where 49.0 sec are spent in the filter step and the remaining

30.3 sec in the smoother step. We need to run the positivity stage during both steps. It takes about 8.4 sec in KF and 0.6 sec in the smoother; much less than KF or only 7.6%. This is explained by the fact that after the filtering step, images are already nonnegative and there is no need, most of the time, to enter the positivity stage during the smoothing step. The positivity takes about 9.1 sec in total, an average of 0.11 sec per recursion, less than 11.5% of the total running time. There are 79 recursions in total, 40 in the filtering step and 39 in the smoothing step.

Table 6.1: CPU times in seconds and percentages.

	Filter step	Smoother step	Both steps
total positivity time	8.43	0.64	9.07
positivity mean time per recursion	0.21	0.02	0.11
reconstruction total time	49.03	30.27	79.30
positivity time as percent of total	17.18%	2.11%	11.43%

6.5 Tikhonov Regularization

In chapter 5 we covered the spatial smoothness of the solution of our image reconstruction in nuclear medicine. Recall that two main practical approaches are generally applied: introducing constraints into the problem and using iterative solvers. In section 6.3.1 we have seen how our proximal method to enforce the nonnegativity constraint takes care of some of the spatial regularization. In section 6.4 we have also shown that the iterative scheme we employ is numerically efficient to help our solution to be more spatially regularized. In order to include more spatial constraints into the reconstruction, we experiment now with Tikhonov regularization known as ridge regression in the statistics community; refer to section 5.3 for more details about this method. The cost function that we are minimizing using a Bregman projection is

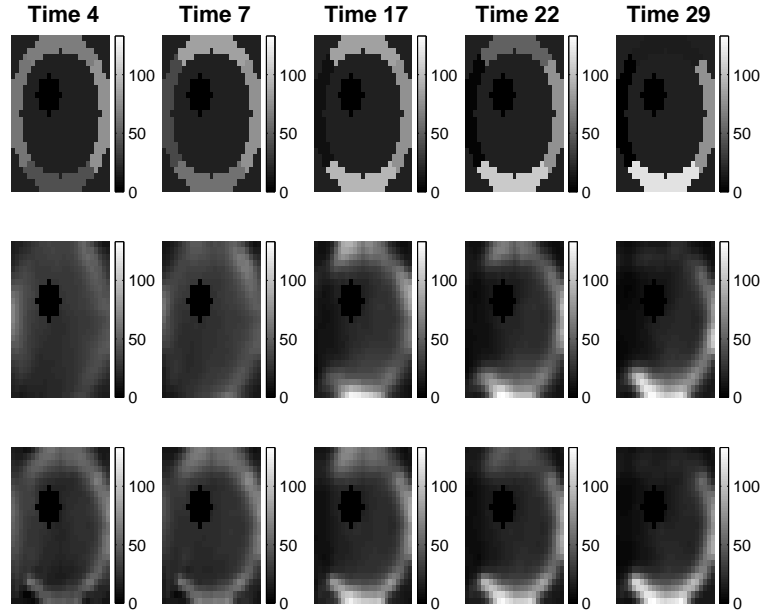


Figure 6.17: Tikhonov regularized images at various times: Truth in 1st row, filtered image in 2nd row, smoothed image in 3rd row.

$$f(v) = \frac{1}{2} \|v - \hat{x}\|_{P^{-1}}^2 + \frac{\alpha}{2} \|Lv\|^2 \quad (6.5.1)$$

where L is an appropriately chosen regularization operator, \hat{x} is the output activity of the Kalman algorithm. We follow the same systematic procedure 6.1 with the cost function from (6.5.1) and algorithm 5.2 in step 5. We tried $L = I$; that is we preferred a solution with smaller norm. We also chose L , see below for an example, to be the second order differential operator that we note as Diff2 where the neighboring system is shown in Figure 5.1. We did not notice any significant change or improvement from one setting to another. Images presented here are from the setting $L = \text{Diff2}$ and $Q = 40I$. We got, for instance, the following stencil for the Diff2 operator based on

the neighboring system of Figure 5.1,

$$\begin{pmatrix} -2 & 1 & 0 & 1 & 0 & 0 & 0 & 0 & 0 \\ 1 & -3 & 1 & 0 & 1 & 0 & 0 & 0 & 0 \\ 0 & 1 & -2 & 0 & 0 & 1 & 0 & 0 & 0 \\ 1 & 0 & 0 & -3 & 1 & 0 & 1 & 0 & 0 \\ 0 & 1 & 0 & 1 & -4 & 1 & 0 & 1 & 0 \\ 0 & 0 & 1 & 0 & 1 & -3 & 0 & 0 & 1 \\ 0 & 0 & 0 & 1 & 0 & 0 & -2 & 1 & 0 \\ 0 & 0 & 0 & 0 & 1 & 0 & 1 & -3 & 1 \\ 0 & 0 & 0 & 0 & 0 & 1 & 0 & 1 & -2 \end{pmatrix}$$

We observed in section 5.3 that an optimal value of α exists even though finding it is not an easy task and is an ongoing active research topic. Our take on this issue is heuristic. We thus experimented with various values of this parameter α . The smallest value of τ_{avg} happens when $\alpha_{optimal} \simeq 10^{-5}$; the smoothed reconstructed image has $\tau_{avg} = 0.36$ and of the filtered one has $\tau_{avg} = 0.42$. It takes about 79.7 sec to run the whole algorithm where nonnegativity and Tikhonov regularization takes about 11.5% of the total running time. It is as much time as when we impose the nonnegativity constraint only. We notice that there are some pixels' grouping if we compare the reconstructed image to the one done with enforcing the nonnegativity constraint only; compare for instance the lower and upper regions in Figure 6.15 and Figure 6.17. As it is known in the regularization literature, Tikhonov tends to over-smooth. We observe the same effect here where we see that images look blurred.

6.5.1 Augmented

We drew attention in section 5.3 to the fact that the authors in [10, 60, 102] use an augmented system in implementing Tikhonov regularization into KF. The method is an extension of Tikhonov in which the original observation model is replaced by

an augmented one; we therefore refer to this approach as “augmented”. The equation (2.2.2) is then replaced with

$$\tilde{y}_k = \tilde{C}_k x_k + \tilde{v}_k \quad (6.5.2)$$

where

$$\tilde{y}_k = \begin{pmatrix} y_k \\ \alpha L \bar{x}_k \end{pmatrix}, \tilde{C}_k = \begin{pmatrix} C_k \\ \alpha L \end{pmatrix}, \text{ and } \tilde{v}_k = \begin{pmatrix} \nu_k \\ \epsilon_k \end{pmatrix}$$

with L being the second-order difference or Diff2 matrix as we defined it before. The vector ϵ_k is a Gaussian zero-mean error of the fictitious noisy observations $\alpha L \bar{x}_k$ with covariance U_k and \bar{x}_k being some target value.

We implemented this method in order to compare it to ours, referred to as “Bregman” in Figure 6.18. We found the best tuning parameter α to be 10^{-2} via numerical trials with $\hat{x}_{0|0} = \text{OSEM act}$, $Q = 40I$, $P_{0|0} = 10^5 I$, and $U_k = I$. We used the same projected Bregman technique to ensure the positivity of the solution. So we keep the same conditions of the experiment except the way we implemented Tikhonov regularization. With $\bar{x} = 0$ we get $\tau_{avg} = 0.35$ and 0.36 with \bar{x} being the average of the true activity. Those values of deviation error are similar in range compared to Tikhonov regularization presented previously, which is 0.36 and TACs are almost the same, refer to Figure 6.18. It takes about 150.4 sec to find the solution using this method. It is about one and a half times as long as the computational time of the Tikhonov regularization using our implementation. Beside the usual Tikhonov over-smoothing effect, we do lose here, with this augmented system, some of the characteristics of the annulus; see for instance the region around the background in Figure 6.18. Furthermore, the augmented method is very demanding in terms of memory space because it is doubling the size of vectors and quadrupling one of the matrices in (6.5.2), especially the size of the matrix to invert in (2.3.3). Tikhonov regularization seems to behave better with our implementation.

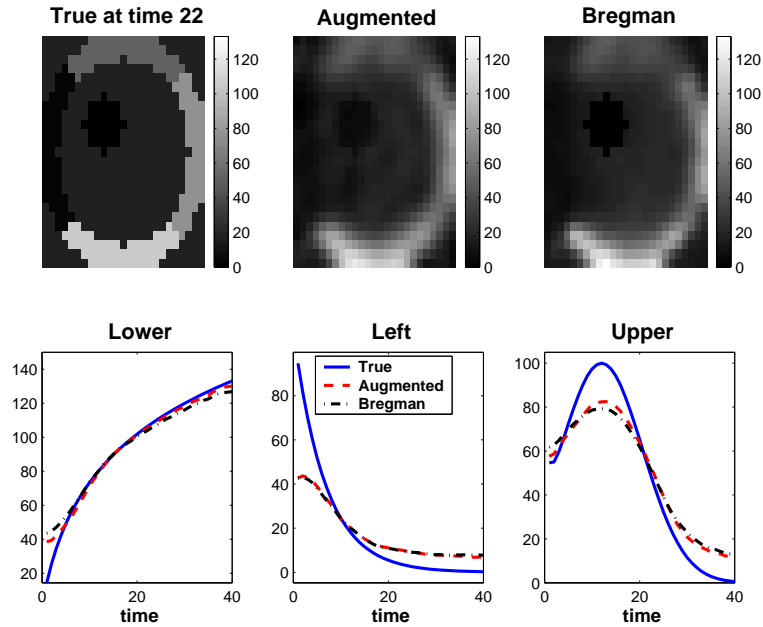


Figure 6.18: Augmented vs Bregman Tikhonov regularization at time 22: Blue for truth, red for augmented, and black for Bregman.

6.6 Median Regularization

To avoid the over smoothing of Tikhonov type regularization, we introduced in section 5.4 an edge preserving regularization as an alternative. It is in essence based on the absolute value function $|x|$. The absolute value function is convex but not smooth; it is not differentiable at zero. We used the following approximation that is both convex and differentiable, $\varphi_\eta(x) = \frac{1}{\eta} \log \cosh(\eta x)$, see Figure 5.2. As mentioned before in section 5.4, a connection exists between the $|x|$ function and the median, thus the name of “Median” for this edge-preserving approach. Recall that the cost function is

$$f(v) = \frac{1}{2} \|\hat{x} - v\|_W^2 + \alpha \sum_i \sum_{j \in \mathcal{N}_i} |\hat{x}_i - v_j| \quad (6.6.1)$$

We apply the systematic procedure 6.1, with the cost function (6.6.1), using the algorithm 5.3 in step 5. Our criteria of selecting parameters are again the value of

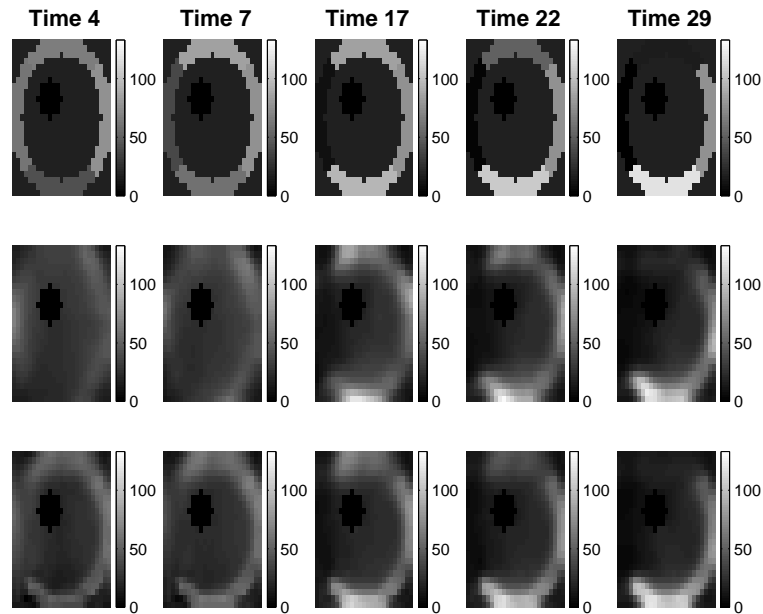


Figure 6.19: Median regularized images at different times: Truth in 1st row, filtered image in 2nd row, smoothed image in 3rd Row.

τ_{avg} and visual inspection. We notice that choosing any value of the parameter η greater or equal to 5 has not much bearing in improving the reconstructed images while we have an optimal value of the parameter α around 10^{-2} . We present results done with $\eta = 20$ and $\alpha = 10^{-2}$. Figure 6.19 exhibits the reconstructed filtered and smoothed images together with the simulated/true ones at different instances of time when we applied the median regularization. The smoothed reconstructed image has a τ_{avg} of 0.37 while the filtered one has 0.43. It takes about 83.8 sec to run the Median algorithm where nonnegativity and median regularization takes about 11.9% of the total running time, comparable to when we impose the nonnegativity only. Notice the blocky segments and edge-preserving at the borders of the regions. As in the previous approaches, reconstructed smoothed images are “smoother” and somehow “better” than the filtered ones. The median approach groups pixels together within a certain ROI. Furthermore, the too much smoothing effect of Tikhonov regularization, for instance in the middle of the images and at the borders of the arcs, has diminished

with the median regularization. Table 6.2 summarizes the deviation error τ of several reconstructions we presented up to now.

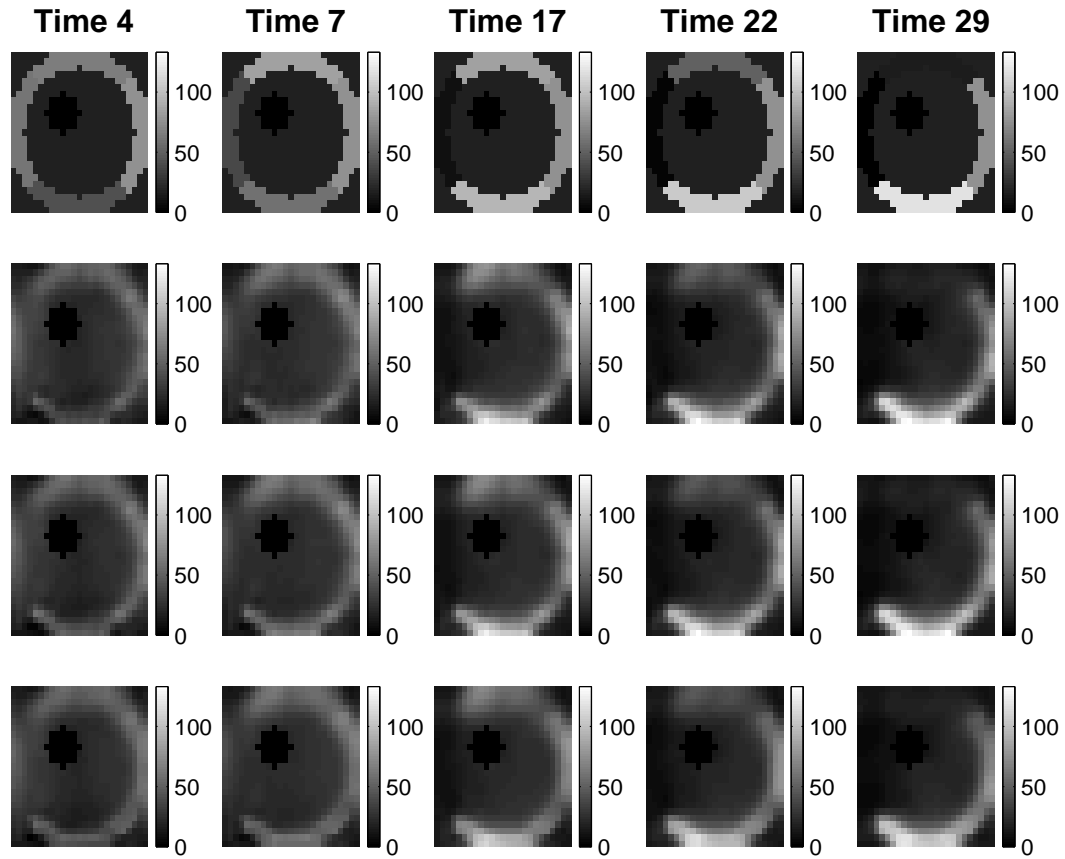


Figure 6.20: Three reconstructed images at different times: Truth in 1st row, just positivity in 2nd row, positivity and Tikhonov regularization in 3rd row, positivity and median regularization in 4th row.

So far we have reconstructed images using KF projected using Bregman distance to ensure nonnegativity. We have obtained quite good images and TACs. Temporal regularization is taken care of by KF itself. Enforcing the nonnegativity, stopping the iterative algorithm earlier on, implementing “Tikhonov” and “Median” all serve to impose spatial regularization and by the same token to minimize more the effect of the ill-posedness. See Figure 6.20 for the three reconstructed images put together one after the other. As we have already mentioned in this section, blocky segments

and edge-preserving at the borders of the regions are present when using the Median algorithm and the too much smoothing effect of the Tikhonov regularization is less present. We also mentioned in 6.5.1 that there are some pixels' grouping, when we use the Tikhonov algorithm, comparatively to with enforcing nonnegativity constraint only; compare for instance the lower and upper regions.

6.7 Hölder Filter Regularization

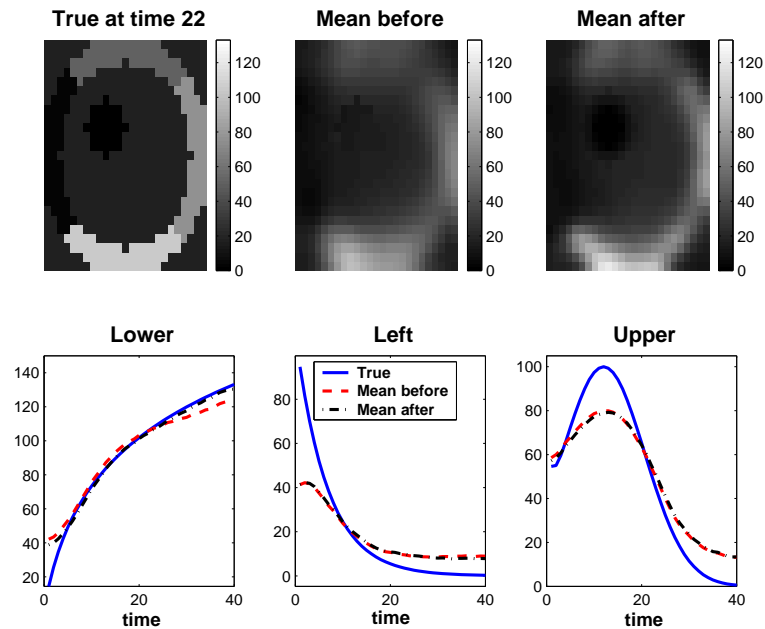


Figure 6.21: Mean-before and mean-after regularization at time 22 - TACs: Blue for truth, red for mean-before, and black for mean-after.

We implemented and tested more regularization approaches. Tikhonov and Median regularization use the 2-norm and 1-norm respectively to achieve their purpose of regularization. Both norms are associated essentially with the mean and median functions and we have put these norms in a more generalized Hölder norm form, see section 5.5. Therefore, we apply these functions to our \hat{x} , the output activity of KF.

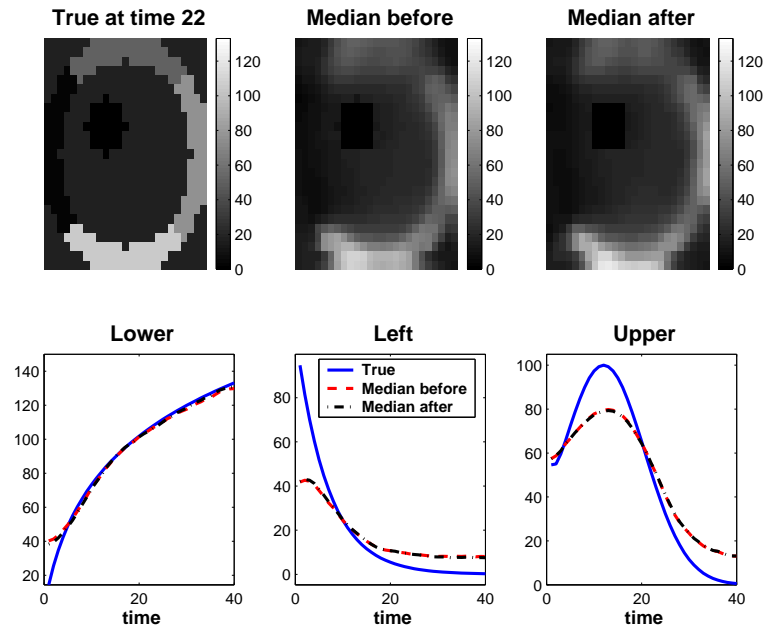


Figure 6.22: Median-before and median-after regularization at time 22 - TACs: Blue for truth, red for median-before, and black for median-after.

These two norms are rightfully referred to in the literature as filter window averaging and median filter. We referred to them before as Hölder filter of order 2 and 1. We used these two variants in two ways each. We applied one of these functions before entering the next step of the recurrence of the KF, meaning after applying our positivity algorithm. We also just applied one of them after applying the projected KF, that is at the end of the algorithm. We thus have another two variants for the mean and median functions, four in total. We name them “Mean-Before”, “Mean-After”, “Median-Before”, and “Median-After” respectively. Figure 6.21 presents the TACs and the reconstructed “Mean-Before” and “Mean-After” images at one instance of time while Figure 6.22 shows the ones of “Median-Before” and “Median-After”. These variants not only do a good job in edge-preserving, but also blur the images less. However, the four variants succeeded to an extent to group pixels within each one of the ROI. Notice, for instance, how the “star” region is even more blurred with this mean filter approach than when we use Tikhonov and how its shape looks like a

square with this median filter.

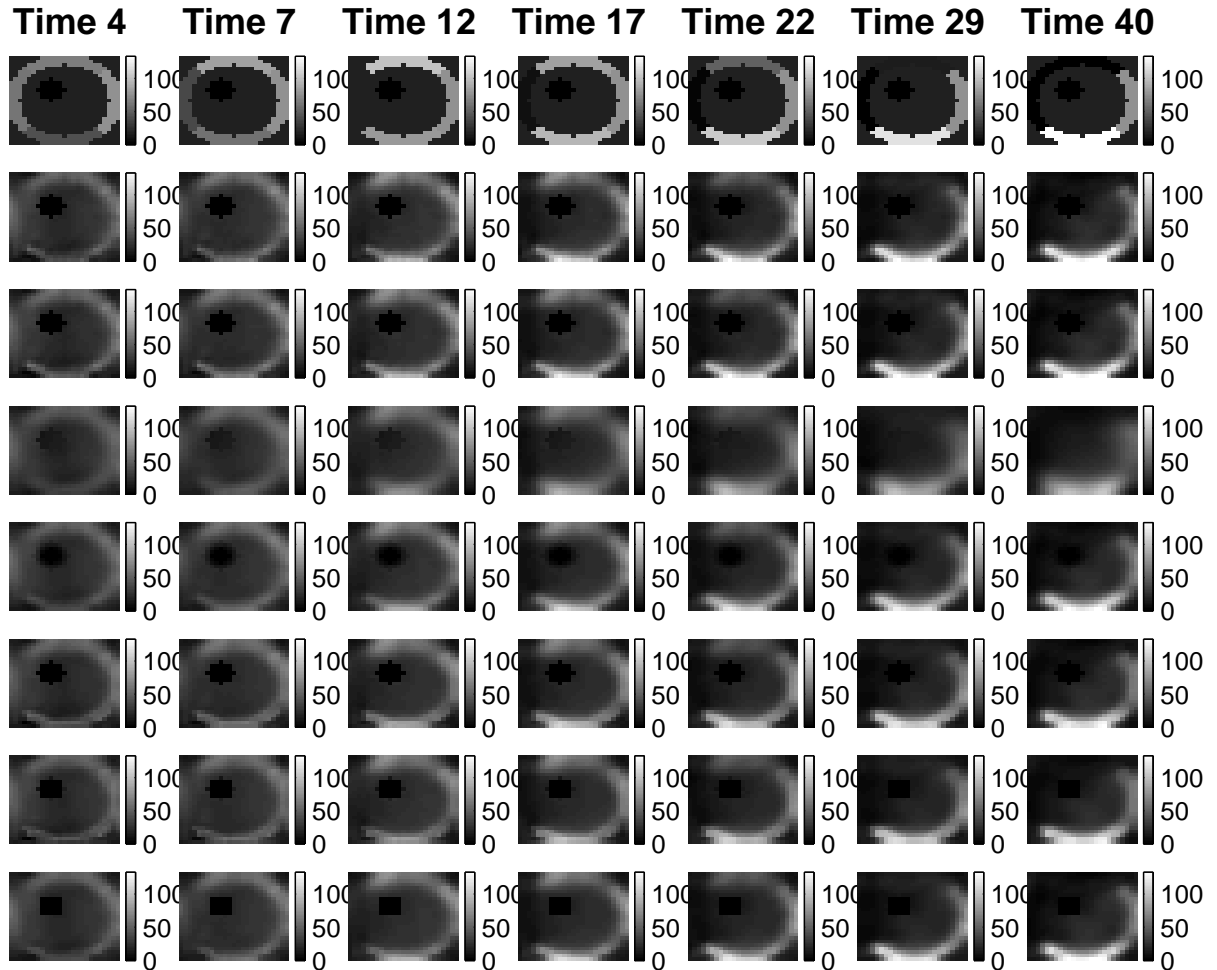


Figure 6.23: All reconstructed images: Truth in 1st row, just positivity in 2nd row, Tikhonov in 3rd row, mean-before in 4th row, mean-after in 5th row, median in 6th row, median-before in 7th row, median-after in 8th row.

For completeness sake we provide the seven sets of images at different instances of time that we reconstructed so far with and without regularization together with the simulated/true ones, see Figure 6.23. Notice that the three median-based approaches give better images. The “Median” ones look a little better even though it does not give the smallest τ value among the three median regularization-based ones. Table 6.2 summarizes the deviation error τ and Table 6.3 summarizes the total CPU times in

seconds of some reconstructions we presented so far.

Table 6.2: Deviation error τ_{avg} of several reconstructions.

	Filtering Step	Smoothing Step
No positivity	2.59	2.42
With positivity	0.42	0.37
Augmented	0.40	0.36
Tikhonov	0.42	0.36
Mean-before	0.49	0.46
Mean-after	0.43	0.39
Median-before	0.43	0.39
Median-after	0.43	0.38
Median	0.43	0.37

Table 6.3: CPU time in seconds of several reconstructions.

Positivity	Tikhonov	Median	Augmented
79.30	79.65	83.82	150.44

6.8 Comparing with Improved dEM Algorithm

We mentioned in section 1.3 that methods based in using a “mask” to incorporate inequality constraints on the variables could be costly in time and prone to introduce some bias. There is ongoing work to improve the dEM algorithm [40] mainly by updating the peaks of the TACs using techniques explored in [41]. We include here a comparison between this improved dEM and our approach that we refer to as projected Kalman.

The simulated 2D phantom is a combination of two slices of a 3D phantom, one containing the bladder and the other one containing the two kidneys. Five ROIs are simulated over 48 time frames knowingly, right kidney, left kidney, bladder, immersed,

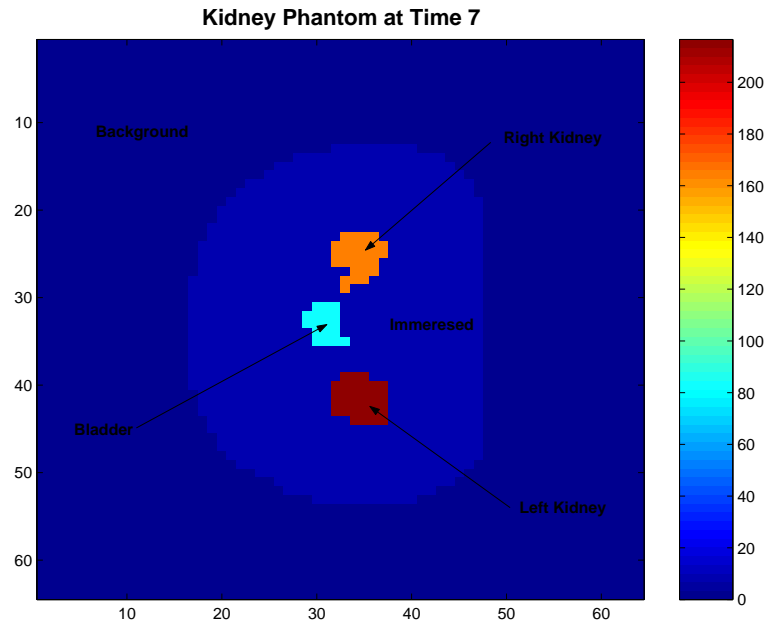


Figure 6.24: Simulated $64 \times 64 \times 48$ kidney phantom at time 7.

and background. The resulting 2D phantom, courtesy of Thomas Humphries, is a $64 \times 64 \times 48$ movie that we refer to as the kidney phantom; see Figure 6.24. The sinogram is done with two heads camera in L-mode, 64 detector bins in each head. One camera starts off behind the two kidneys and the other on the side closest to the right kidney. It then rotates clockwise behind the phantom over 360° . There is no Poisson noise added in the sinogram for dEM reconstruction while there is one included in the sinogram for our approach. The improved dEM reconstruction is better than just by using the classical dEM alone since the TACs are smoother than those, for instance, exhibited in [16]; see Figure 6.26. We observe that both reconstructions, improved dEM and projected Kalman, yield about the same quality level images and TACs; refer to Figure 6.25 and Figure 6.26.

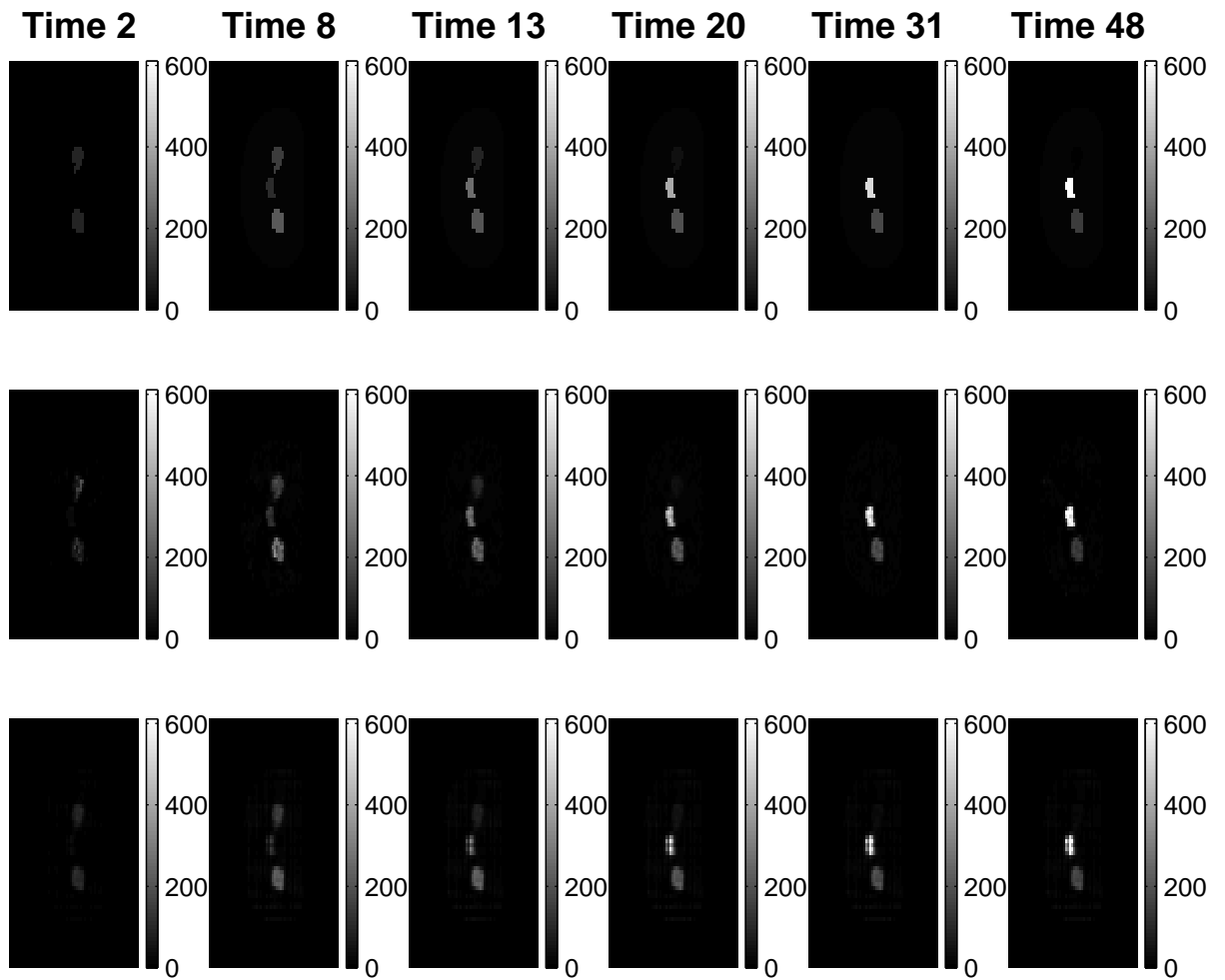


Figure 6.25: Improved dEM and Projected Kalman kidney images: Truth in 1st row, improved dEM in 2nd row, and projected Kalman in 3rd row.

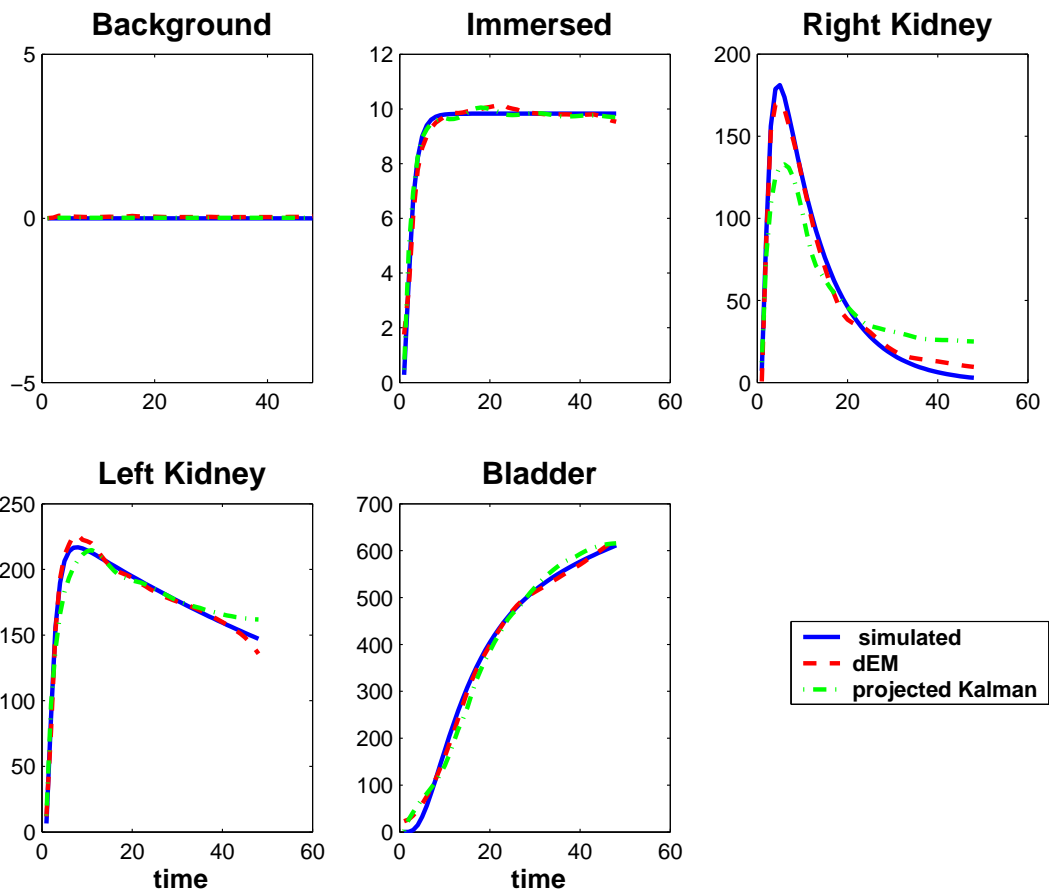


Figure 6.26: Reconstructed kidney TACs with improved dEM algorithm and projected Kalman approach.

6.9 Various Sized Phantoms

The 25×25 digital phantom, 625 pixels in total, we used for experiments is now set into a bigger 64×64 phantom, see Figure 6.27. We utilized the 25×25 one in order to determine heuristically the tuning parameters of the reconstruction. We are now ready to push for bigger phantoms with the same parameters that we summarize in table 6.4. The parameters might depend on the grid size and further experiments could help to clarify this issue. However, one heuristic way of getting in general the tuning parameters of any reconstruction is to start off with a coarser grid. Once we have them, we could proceed thereafter with finer grids.

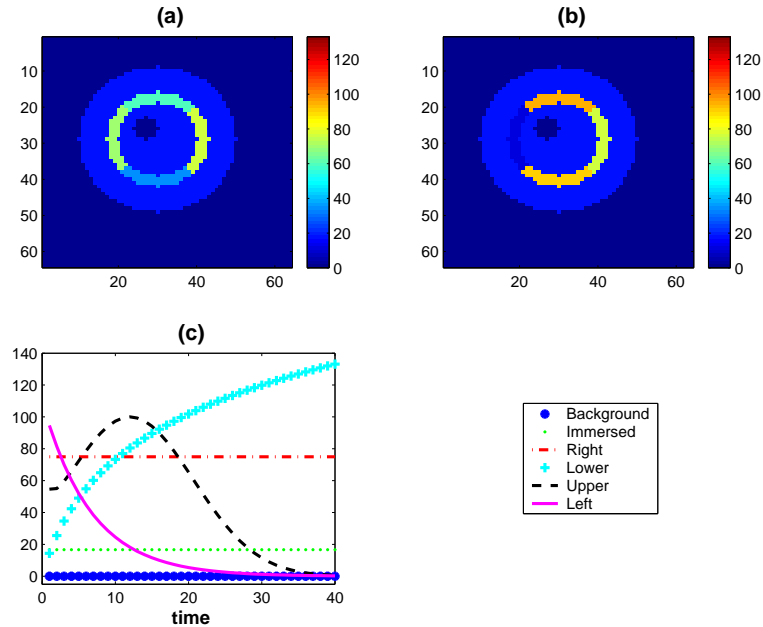


Figure 6.27: Simulated 64×64 annulus and TACs: (a) phantom at time 3, (b) at time 15, (c) TACs.

Table 6.4: Tuning parameters.

$\hat{x}_{0 0}$	$P_{0 0}$	Q_k	R_k	W	γ	# of iter	α_{Tikh}	α_{Median}
OSEM act	$10^5 I$	$40I$	$\text{diag}(y_k)$	P^{-1}	1	1	10^{-5}	10^{-2}

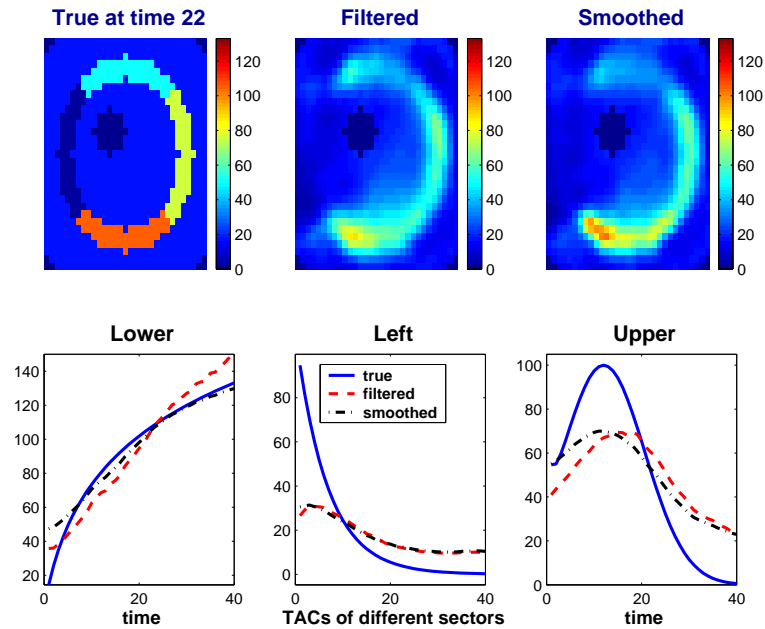


Figure 6.28: Size 31×31 digital phantom at time 22 - TACs: Blue for truth, red for filtered, and black for smoothed.

We conducted trials with sizes 31×31 , 961 pixels, 45×43 to include a zero-activity background, 1935 pixels, and 64×64 , 4096 pixels; Figures 6.28, 6.29, 6.30 give the reconstructed images and TACs respectively. The ratio of data to unknowns is summarized in table 6.5 and goes from 1:3.25 for the 25×25 size to 1:21.34 for the 64×64 size phantom.

We notice that the TACs of the three bigger sized phantoms look similar in shape to the true ones. The TAC of the “Upper” arc succeeds in catching the time of the peak as it was the case with 25×25 size, see Figure 6.16. As we increase the number of pixels from 625 to 4096, the intensity at the peak in the “Upper” arc, for instance, gets smaller; compare this intensity in Figure 6.16, 6.28, 6.29, and 6.30. The smoothed reconstructed images are again better than the filtered ones. For the static case, a reconstruction of a 3D image can take several hours using OSEM [104]. In our 2D dynamic SPECT reconstruction, an acceptable amount of CPU time would ideally be then few minutes. Table 6.6 sums up the CPU time taken to run each size of the

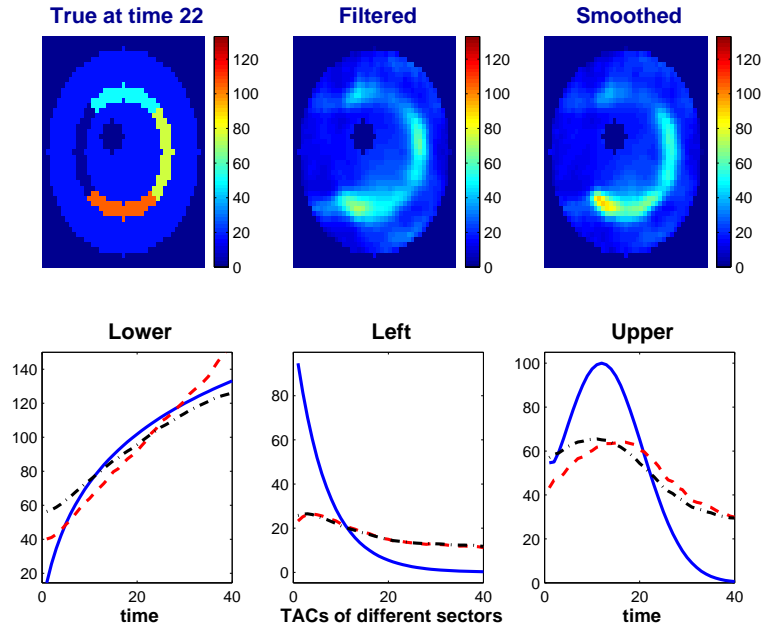


Figure 6.29: Size 45×43 digital phantom at time 22 - TACs: Blue for truth, red for filtered, and black for smoothed.

phantom. It goes from 80 seconds for 625 pixels to 4.5 hours for 4096 pixels. While the size 4096 is about 6.6 times the one of 625, the CPU time is almost 200 times as much. Therefore, we have a complexity of the order N^3 . This huge jump could be explained by the few matrix multiplications, of the same size as the one of the phantom, that are involved in (2.3.3) and (2.3.6) of the Kalman algorithm.

Table 6.5: Ratio of data to unknowns.

Size	25×25	31×31	45×43	64×64
Ratio	1 : 3.25	1 : 5.00	1 : 10.00	1 : 21.34

As far as the values of the average deviation τ_{avg} are concerned, we are still in the vicinity of 0.4 as it was the case with 25×25 size. Table 6.7 includes those values. We present reconstructions done with Tikhonov and Median spatial regularization for size 31×31 and 45×43 , see Figure 6.31, 6.32, 6.33, and 6.34. We notice again grouping

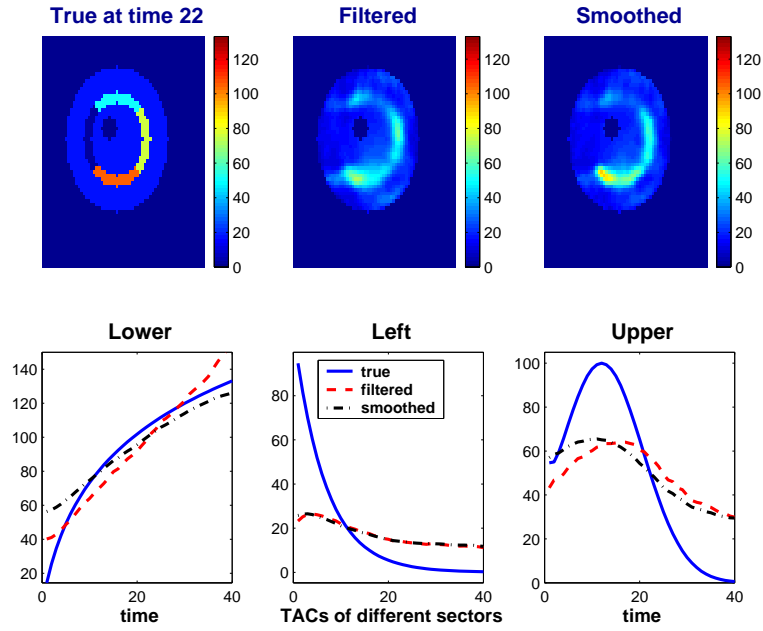


Figure 6.30: Size 64×64 digital phantom at time 22 - TACs: Blue for truth, red for filtered, and black for smoothed.

Table 6.6: CPU time of various sizes.

Size	25×25	31×31	45×43	64×64
Time	80 seconds	4.5 minutes	34 minutes	4.5 hours

of pixels due to Tikhonov, compare for instance Figure 6.29 with Figure 6.32. While Tikhonov regularization over smooths especially at the boundaries of the regions, the Median regularization preserves the edges, compare Figure 6.32 and Figure 6.34. TACs are once more similar in shape to the simulated ones in all these cases of regularization.

Table 6.7: Average deviation τ_{avg} of different sizes.

Size	25×25	31×31	45×43	64×64
τ_{avg}	0.37	0.43	0.46	0.46

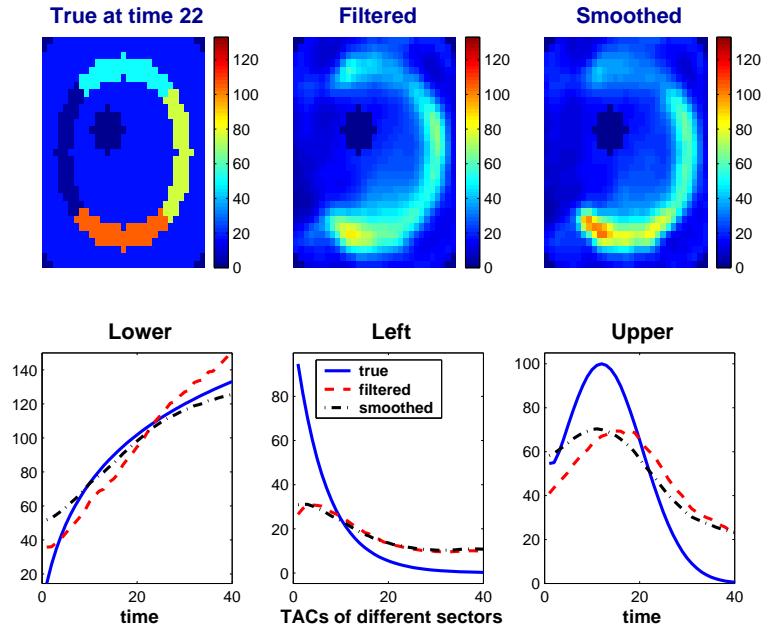


Figure 6.31: Size 31×31 digital phantom with Tikhonov regularization at time 22 - TACs: Blue for truth, red for filtered, and black for smoothed.

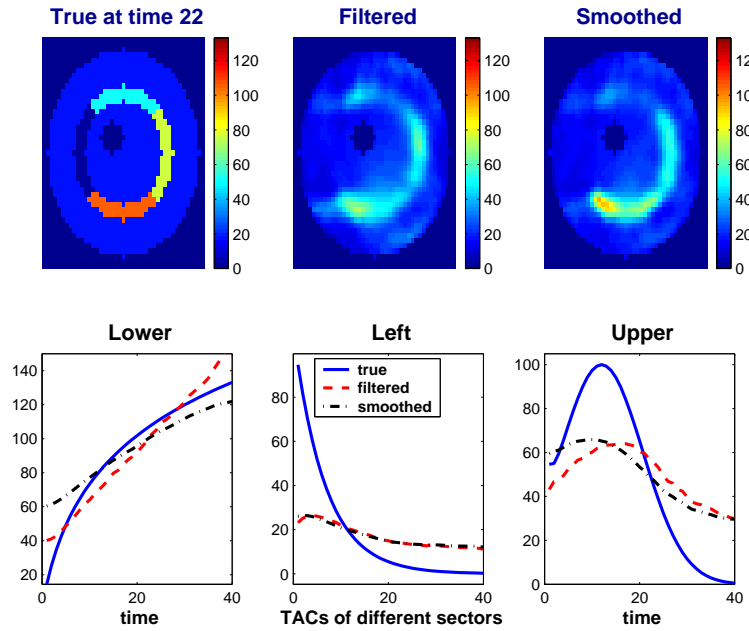


Figure 6.32: Size 45×43 digital phantom with Tikhonov regularization at time 22 - TACs: Blue for truth, red for filtered, and black for smoothed.

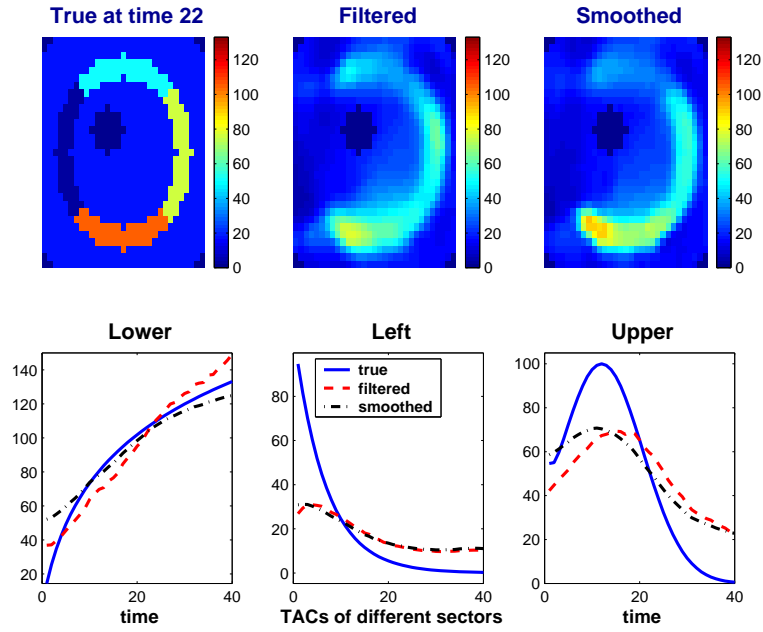


Figure 6.33: Size 31×31 digital phantom with Median regularization at time 22 - TACs: Blue for truth, red for filtered, and black for smoothed.

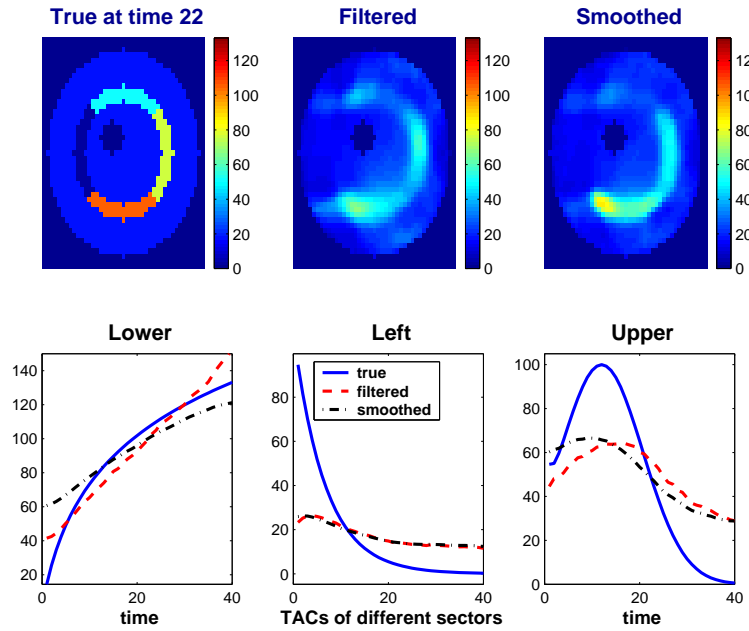


Figure 6.34: Size 45×43 digital phantom with Median regularization at time 22 - TACs: Blue for truth, red for filtered, and black for smoothed.

6.10 Spatial Regularization via Segmentation

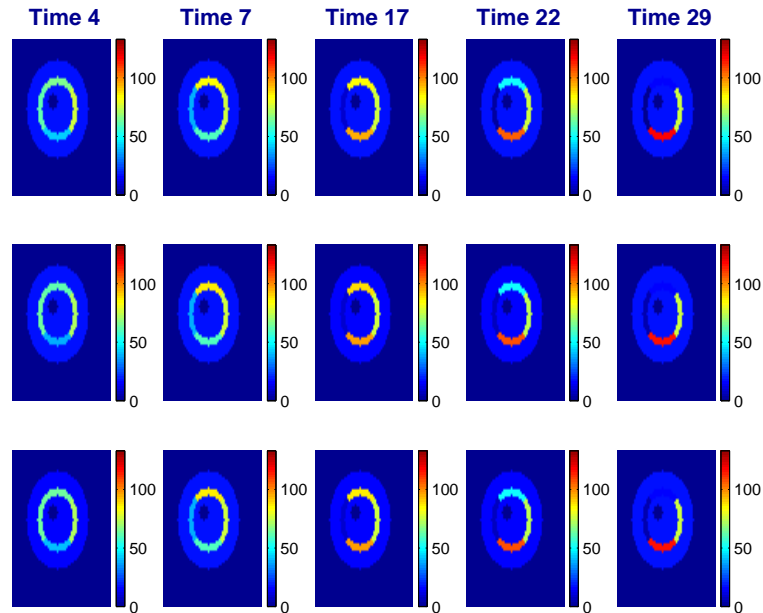


Figure 6.35: Regularization via segmentation images at different times: Truth in 1st row, filtered image in 2nd row, smoothed image in 3rd Row.

Finally we include spatial regularization via segmentation. When we know how to partition the digital phantom into nonintersecting regions, we can include this additional constraint into the problem and reduce its size at the same time. We offer here an experiment done when we have perfect information about the six ROIs. We utilize the procedure 6.1 where the variable is now ξ instead of x as we mentioned in section 6.10. We apply procedure 6.1 where we employ algorithm 5.4 in step 4 and 6 while the matrix $C_k E$ substitutes the matrix C_k in steps 4 and 6. The computation takes 1.72 seconds, 2% of the one of 25×25 phantom, while the average deviation is 0.06, 6.5 times better than the unsegmented. Figure 6.35 shows the reconstructed filtered and smoothed images together with the simulated one while Figure 6.36 exhibits the reconstructed TACs being very close to the true time activity curve. We get a much better reconstruction since we have more data. This sustains the asymptotic consistency property, as the observations number increases, of the

estimator x^* given in corollary 4.1. Recall the consistency property (4.2.1),

$$\lim_{n \rightarrow \infty} P(|x^* - x| < \epsilon) = 1 \quad \forall \epsilon > 0$$

where n is the size of the data. Furthermore, the estimator seems to confirm that it is sufficient, acceptable, unbiased, and efficient thus verifying the criteria stipulated in section 4.2,

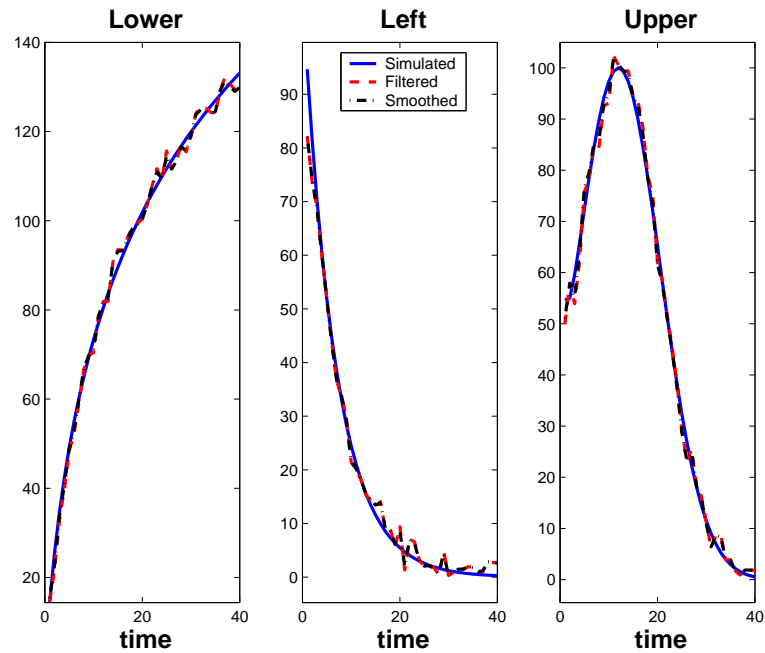


Figure 6.36: TACs of three regions: Blue TAC for truth, red TAC for reconstructed filtered, and black TAC for smoothed.

Chapter 7

Conclusions

In this thesis we have presented reconstruction schemes for a nuclear medicine problem. Our approach does not require a precise a-priori information about the underlying dynamic physical process. It consists of finding the best linear unbiased estimator solution of the dynamic SPECT reconstruction problem and then applies a generalized proximal approach. The initial solution was found using the classical Kalman algorithm. However, this solution is meaningless because it fails to be nonnegative. Setting negative values of the activity to zero or taking their absolute value do not work to get an acceptable solution. We remedied this Kalman algorithm's shortcoming by projecting this solution onto the nonnegative orthant via a Bregman approach. We showed that the projected estimator not only conserves the same properties as its parent Kalman filter estimator, but also performs better. Analysis of images and time activity curves showed a net improvement. The Kalman filter by nature takes care of temporal regularization. However, we do not have spatial smoothness, even among pixels within the same region. We propose then a few spatial regularization approaches. Numerical results confirm the effectiveness of our methods, especially with the "Median" approach that preserves the edges, while keeping the temporal smoothness feature. We also compared our Tikhonov regularization implementation with

the one implemented in earlier work and found that ours outperformed it in terms of running time and memory space needs, while producing images of comparable quality.

We used a heuristic method to determine the regularization parameters. Ongoing work consists of improving regularization by choosing a different operator L , neighboring system, and by automating the choice of the parameters. One way we are thinking of doing this is by using the L-curve [14, 83]. Cross validation [103] and discrepancy principle [84] are yet other approaches to experiment with. We are investigating other regularization avenues such as total variation.

First-order hidden Markov chain models very well a large portion of problems including those of tracer kinetics. We only experimented here with a first-order random walk; thus we chose the identity as the evolution matrix. We are looking into using a higher order random walk. Interesting aspects of the proposed method remain to be investigated, for example addressing the issue of computational complexity for large scale systems.

In this thesis we offered a mathematical model and numerical approach to solve the inverse problem of dynamic medical image reconstruction. Inverse problems are typically ill-posed and ill-conditioned. Our method performs well in handling both challenges. We believe that our approach could be utilized when the activity is static (time independent) and in other fields as well where nonnegativity or spatial regularization are required.

Bibliography

- [1] <http://cialab.ee.washington.edu/REPRINTS/1997-AlternatingProjections.pdf>,
(Last visited: October 1, 2008).
- [2] H. Abdi, *Least squares*, <http://www.utd.edu/~herve/Abdi-LeastSquares-pretty.pdf>,
2003 (Last visited: October 1, 2008).
- [3] A. C. Aitken, *On least squares and linear combinations of observations*, Proceedings of the Royal Society of Edinburgh, 55, 1935, pp. 42–48.
- [4] S. Alenius and U. Ruotsalainen, *Bayesian image reconstruction for emission tomography based median root prior*, E. J. Nucl. Med. **24** (1997), 258–265.
- [5] S. Alenius, U. Ruotsalainen, and J. Astola, *Using local median as the location of the prior distribution in iterative emission tomography image reconstruction*, IEEE Trans. Nucl. Sci. **45** (1998), 3097–3107.
- [6] B. Amini, M. Bjrkklund, R. Dror, and Nygren A., *Tomographic reconstruction of SPECT data*, <http://www.owl.net.rice.edu/~elec539/Projects97/cult/report.html>, 1997
(Last visited: October 1, 2008).
- [7] B. D. O. Anderson and J. B. Moore, *Optimal Filtering*, Printice-Hall, Englewood, Ciffs, NJ, 1979.

- [8] V. M. Artemiev, A. O. Naumov, and G. R. Tillack, *Dynamic image reconstruction: general estimation principles for dynamic tomography*, 15th WCNDT, Roma, 2000.
- [9] V. M. Artemiev, A. O. Naumov, and G. R. Tillack, *Recursive tomographic image reconstruction using Kalman filter approach in the time domain*, Jour. of Phys. D: Appl. Phys. **34** (2001), 2073–2083.
- [10] D. Baroudi, J. Kaipio, and E. Somersalo, *Dynamical electric wire tomography: a time series approach*, Inverse Problems **14** (1998), 799–813.
- [11] H. H. Bauschke and J. M. Borwein, *Legendre functions and the method of random Bregman projections*, J. Convex Anal. **4** (1997), no. 1, 27–67.
- [12] H. H. Bauschke, P. L. Combettes, and D. Noll, *Joint minimization with alternating Bregman proximity operators*, Pacific Journal of Optimization **2** (2006), 401–424.
- [13] H. H. Bauschke, D. Noll, A. Celler, and J. M. Borwein, *An EM-algorithm for dynamic SPECT tomography*, IEEE Trans. Med. Imag. **18** (1999), 252–261.
- [14] F. Benyah and L. S. Jennings, *The L-curve in regularisation of optimal control computation*, J. Austral. Math. **40** (1998), E138–E172.
- [15] M. Bertero and P. Boccacci, *Introduction to Inverse Problems in Imaging*, IOP Publishing, Bristol, 1998.
- [16] S. Blinder, D. Noll, X. Hatchondo, J.P. Celler, A. Esquerré, and P. Payoux, P. Gantet, *Fully 4D dynamic image reconstruction by nonlinear constrained programming*, Nuclear Science Symposium Conference Record, 5, 2003, pp. 3161–3165.

- [17] C. Blondel, D. Noll, J. Maeght, and T. Celler, A. Farncombe, *Comparison of different figure of merit functions for dynamic single photon computed tomography (dSPECT)*, IEEE Nuclear Science Symposium Conference Record, 2000, pp. 15115–15155.
- [18] A. V. Bos, *Parameter Estimation for Scientists and Engineers*, Cambridge University Press, Hoboken, NJ, 2007.
- [19] D. Boulfefel, L. J. Hahn, R. Kloiber, and G. R. Kuduvalli, *Two-dimensional restoration of single photon emission computed tomography images using the Kalman filter*, IEEE Trans. on Med. Imag. **13** (1994), 102–109.
- [20] S. Boyd and L. Vandenberghe, *Convex Optimization*, Cambridge, New York, NY, 2004.
- [21] L. M. Bregman, *The relaxation method of finding the common point of convex sets and its application to the solution of problems in convex programming*, USSR Comput. Math. Phys. **7/3** (1967), 200–217.
- [22] R. Bro and S. Jong, *A fast non-negativity-constrained least squares algorithm*, Jour. of Chemo. **11** (1997), 393–401.
- [23] C. L. Byrne, *Iterative image reconstruction algorithms based on cross-entropy minimization*, IEEE Trans. on Image Processing **2** (1993), 96–103.
- [24] ———, *Block-iterative interior point optimization methods for image reconstruction from limited data*, Inverse Problems **16** (2000), no. 5, 1405–1419.
- [25] ———, *Proximity function minimization using multiple Bregman projections, with applications to split feasibility and Kullback-Leibler distance minimization*, Annals of Operations Research **105** (2001), 77–98.

- [26] ———, *Signal Processing, A Mathematical Approach*, A K Peters, Wellesley, MA, 2005.
- [27] F. Campillo, *Filtrage particulaire et modèles de markov cachés*, <ftp://ftp.irisa.fr/local/sigma2/campillo/cours/2006-master2-toulon.pdf>, 2006 (Lat visited: October 1, 2008).
- [28] R. Carson, *A maximum likelihood method for region-of-interest evaluation in emission tomography*, *J. Comput. Assit. Tomogr.* **10** (1986), 654–663.
- [29] A. Celler, T. Farncombe, C. Bever, D. Noll, J. Maeght, and D. Harrop, R. Lyster, *Performance of the dynamic single photon emission computed tomography (dSPECT) method for decreasing or increasing activity changes*, *Phys. Med. Biol.* **45** (2000), 3525–3543.
- [30] Y. Censor and A. Lent, *An iterative row-action method for interval convex programming*, *J. Optim. Theory Appl.* **34** (1981), no. 3, 321–353.
- [31] Y. Censor and S. A. Zenios, *Parallel Optimization: Theory, Algorithms, and Applications*, Oxford University Press, New York, 1997.
- [32] T. F. Chan and J. Shen, *Image Processing and Analysis: Variational, PDE, Wavelet, and Stochastic Methods*, SIAM, Philadelphia, PA, 2005.
- [33] P. Charbonnier, L. Blanc-Féraud, G. Aubert, and M. Barlaud, *Deterministic edge-preserving regularization in computed imaging*, *IEEE Trans. on Image Processing* **5** (1997), 298–311.
- [34] P. L. Combettes and V. R. Wajs, *Signal recovery by proximal forward-backward splitting*, *Multiscale Model. Simul.* **4** (2005), no. 4, 1168–1200.
- [35] H. Cramér, *Mathematical Methods of Statistics*, University Press, Princeton, N.J., 1946.

- [36] A. C. Davidson, *Statistical Models*, Wiley and Sons Inc., Cambridge University Press, 2003.
- [37] A. P. Dhawan, *Medical Image Analysis*, John Wiley and Sons, Hoboken, NJ, 2003.
- [38] C. L. Epstein, *Introduction to the Mathematics of Medical Imaging*, second ed., SIAM, Philadelphia, PA, 2008.
- [39] K. Erlandsson, A.J. Reader, M.A. Flower, and R.J. Ott, *A new 3D backprojection and filtering method for PET using all detected events*, IEEE Trans. Nucl. Sci. **45** (1998).
- [40] T. Farncombe, *Functional Dynamic SPECT Imaging Using a Single Slow Camera Rotation*, Ph.D. thesis, University of British Columbia, 2000.
- [41] T. Farncombe, A. Celler, C. Bever, D. Noll, J. Maeght, and R. Harrop, *The incorporation of organ uptake into dynamic SPECT (dSPECT) image reconstruction*, Nuclear Science, IEEE Transactions on **48** (2001), 3–9.
- [42] C. H. Franklin, *Properties of ML estimators*, <http://www.polisci.wisc.edu/~franklin/Content/MLE/Lecs/MLELec04p4up.pdf>, 2005.
- [43] S. Geman and D. Geman, *Stochastic relaxation, Gibbs distributions and the Bayesian restoration of images*, IEEE Transactions on Pattern Analysis and Machine Intelligence **6** (1984), 721–741.
- [44] S. Geman and D. McClure, *Bayesian image analysis: An application to single photon emission tomography*, Statist. Comput. Sect., Amer. Statist. Assoc., Washington, DC, 1985, pp. 12–18.

- [45] ———, *Statistical methods for tomographic image reconstruction*, Bulletin of the International Statistical Institute **4** (1987), 5–21.
- [46] R. Gordon, R. Bender, and G. T. Herman, *Algebraic reconstruction technique (ART) for three-dimensional electron microscopy and X-ray photography*, J. Theoret. Biol. **29** (1970), 471–481.
- [47] P. G. Green, *Bayesian reconstruction from emission tomography data using a modified EM algorithm*, IEEE Trans. Med. Imaging **9** (1990), no. 1, 84–93.
- [48] J. A. Hanley, L. Joseph, R. W. Platt, M. K. Chung, and P. Béglise, *Visualizing the median as the minimum-deviation location*, Amer. Statist. **55** (2001), no. 2, 150–152.
- [49] S. R. Hare, *Prewhitening: A cure for the common correlation*, <http://www.iphc.washington.edu/Staff/hare/html/presentations/pices6/pices6.html>, 2006 (Last visited: October 1, 2008).
- [50] T. Hebert and R. Leahy, *A generalized EM algorithm for 3-D Bayesian reconstruction from Poisson data using Gibbs priors*, IEEE Transactions on Medical Imaging **8** (1989), no. 2.
- [51] D. Hilbert, *Grundzüge einer Allgemeinen Theorie der Linearen Integralgleichungen*, Chelsea Publishing Company, New York, N.Y., 1953.
- [52] J. B. Hiriart-Urruty and C. Lemaréchal, *Fundamentals of convex analysis*, Springer-Verlag, Berlin, 2001.
- [53] I. T. Hsiao, A. Rangarajan, and G. Gindi, *A new convex edge-preserving median prior with applications to tomography*, IEEE Trans. Med. Imaging **22** (2003), no. 5, 580–585.

- [54] A. Hudson and R. S. Larkin, *Accelerated image reconstruction using ordered subsets of projection data*, IEEE Transactions on Medical Imaging **13** (1994), 601–609.
- [55] Ziena Optimization Inc., <http://www.ziena.com/index.htm>, 2007 (Last visited: October 1, 2008).
- [56] S. Jin, Z Zador, and A.S. Verkman, *Random-walk model of diffusion in three dimensions in brain extracellular space: comparison with microfiberoptic photobleaching measurements*, IEEE Trans. Med. Imaging **95** (2008), no. 4, 1785–1794.
- [57] Q. Jinyi and R. H. Huesman, *Effect of errors in the system matrix on MAP image reconstruction*, Phys Med Biol. **50** (2005), no. 14, 3297–3312.
- [58] C. A. Johnson, J. Seidel, and A. Sofer, *Interior point methodology for 3-D PET reconstruction*, IEEE Trans. Med. Imag. **19** (2000), 271–285.
- [59] G. T. Kaczmarz, *Angenäherte Auflösung von Systemen Linearer Gleichungen*, Bulletin de l'Academie des Sciences et Lettres **A35** (1937), 355–357.
- [60] J. Kaipio and E. Somersalo, *Nonstationary inverse problems and state estimation*, J. Inv. Ill-Posed Problems **7** (1999), 273–282.
- [61] A. Kak and M. Slaney, *Principles of computerized tomographic imaging*, <http://www.slaney.org/pct/pct-toc.html>, 1999 (Last visited: October 1, 2008).
- [62] R. E. Kalman, *A new approach to linear filtering and prediction problems*, Trans. of the ASME-Jour. of Basic Eng. **82** (1960).
- [63] M. Kervinen, M. Vauhkonen, and P. A. Karjalainen, *Time-varying reconstruction in single photon emission computed tomography*, Int. J. of imaging syst. tech **14** (2004), 186–197.

- [64] D. Kim, S. Sra, and I. S. Dhillon, *A new projected quasi-Newton approach for the nonnegative least squares problem*, Tech. Report TR-06-54, Department of Computer Science, University of Texas, 2006.
- [65] H. Kudo and S. Sawada, *Newton-SOR method for fast statistical tomographic image reconstruction*, Systems and Computers in Japan **34** (2003), 1–11.
- [66] L. A. Kunyansky, *A new SPECT reconstruction algorithm based on the Novikov explicit inversion formula*, Inverse Problems **17** (2001), 293–306.
- [67] K. Lange and R. Carson, *EM reconstruction for emission and transmission tomography*, J. Comput. Assit. Tomogr. **8** (1984), 306–316.
- [68] S. L. Lauritzen, *Thiele: Pioneer in Statistics*, Oxford University Press, New York, NY, 2002.
- [69] C. L. Lawson and R. J. Hanson, *Solving Least Squares Problems*, SIAM, Philadelphia, PA, 1995.
- [70] M.A. Limber, A. Celler, J. Barney, M.N. Limber, and J.M. Borwein, *Direct reconstruction of functional parameters for dynamic SPECT*, IEEE Trans. Nuc. Sci. **42** (1995), 1249–1256.
- [71] A. K. Louis, *Medical imaging: state of the art and future development*, Inverse Problems **8** (1992), 709–738.
- [72] J. Maeght, D. Noll, and S. Boyd, *Dynamic emission tomography regularization and inversion*, Bulletin of the Canadian Math. Society, **27**, 2000, pp. 211–234.
- [73] J. Maltz, *Parsimonious basis selection in exponential spectral analysis*, Phys. Med. Biol. **47** (2002), 2341–2365.
- [74] P. S. Maybeck, *Stochastic Models, Estimation, and Control. Vol. I*, Academic Press, New York, NY, 1979.

- [75] J. J. Moreau, *Inf-convolution des fonctions numériques sur un espace vectoriel*, C. R. Acad. Sci. Paris **256** (1963), 5047–5049.
- [76] F. Natterer and F. Wübbeling, *Mathematical Methods in Image Reconstruction*, SIAM, Philadelphia, PA, 2001.
- [77] R. G. Novikov, *An inversion formula for the attenuated X-ray transformation*, Ark. Mat. **40** (2002), no. 1, 145–167.
- [78] D. L. Phillips, *A technique for the numerical solution of certain integral equations of the first kind*, J. Assoc. Comput. Mach. **9** (1962), 84–97.
- [79] J. Qranfal and G. Tanoh, *Regularized Kalman filtering for dynamic SPECT*, J. Phys.: Conf. Ser., 012042, vol. 124, 2008.
- [80] J. Radon, *Über die Bestimmung von Funktionen durch ihre Integralwärte längs Gewisser Männigfaltigkeiten*, Ber. Verh. Sächs. Akad. Wiss. (Leipzig), Math. Phys. Klass **69** (1917), 262–277.
- [81] B. W. Reutter, G. T. Gullberg, and Huesman R. H., *Direct least squares estimation of spatiotemporal distribution from dynamic SPECT projections using spatial segmentation and temporal B-splines*, IEEE trans. med. imaging **19** (2000), 434–450.
- [82] R. T. Rockafellar, *Convex Analysis*, Princeton University Press, Princeton, NJ, 1997, Reprint of the 1970 original.
- [83] G. Rodriguez and D. Theis, *An algorithm for estimating the optimal regularization parameter by the L-curve*, Rendiconti di Matematica **25** (2005), 69–84.
- [84] O. Scherzer, *The use of Morozov’s discrepancy principle for Tikhonov regularization for solving nonlinear ill-posed problems*, Rendiconti di Matematica **51** (1993), 45–60.

- [85] N. C. Schwertman, A. J. Gilks, and J. Cameron, *A simple noncalculus proof that the median minimizes the sum of the absolute deviations*, Amer. Statist. **44** (1990), 38–39.
- [86] L. A. Shepp and Y. Vardi, *Maximum likelihood reconstruction for emission tomography*, IEEE Trans. Med. Imag. **1** (1982), 113–122.
- [87] D. Simon, *Optimal State Estimation: Kalman, H Infinity, and Nonlinear Approaches*, Wiley-Interscience, 2006.
- [88] D. Simon and T. Chia, *Kalman filtering with state equality constraints*, IEEE Trans. Aero. and Electr. Syst. **39** (2002), 128–136.
- [89] D. Simon and D. L. Simon, *Kalman filtering with inequality constraints for turbofan engine health estimation*, IEE Proceedings, Vol. 153, no. 3, 9, 2006, pp. 371–378.
- [90] H. W. Sorenson, *Parameter Estimation, Principles and Problems*, Marcel Dekker Inc., New York, NY, 1980.
- [91] H. Stark, *Image Recovery: Theory and Applications*, Academic, New York, NY, 1987.
- [92] P. Suetens, *Fundamentals of Medical Imaging*, Cambridge Uni.Press, New York, NY, 2002.
- [93] G. Tanoh, *Algorithmes du point intérieur pour l’optimisation en tomographie dynamique et en mécanique du contact*, Ph.D. thesis, Université Paul Sabatier, 2004.
- [94] ———, Private communication, 2007.
- [95] M. Teboulle, *Entropic proximal mappings with applications to nonlinear programming*, Mathematics of Operations Research **17** (1992), no. 3, 670–690.

- [96] J. V. Tiel, *Convex Analysis: An Introductory Text*, John Wiley and Sons, Chichester, NY, 1984.
- [97] A. N. Tikhonov, *On the stability of inverse problems*, Dokl. Akad. Nauk SSSR **39** (1943), no. 5, 195–198.
- [98] A. N. Tikhonov and V.Y. Arsenin, *Solutions of Ill-Posed Problems*, Winston, New York, 1977.
- [99] B. Tsui, X. Zhao, E. Frey, and G. Gullberg, *Comparison between ML-EM and WLS-CG algorithms for SPECT image reconstruction*, IEEE Trans. Nucl. Sci. **38** (1991), 1766–1772.
- [100] M. H. Van Benthem and M. R. Keenan, *Fast algorithm for the solution of large-scale non-negativity-constrained least squares problems*, J.O.C., J. Chemometrics **18** (2004), 441–450.
- [101] E. Vandervoort, A. Celler, G. Wells, S. Blinder, K. Dixon, and Y. Pang, *Implementation of an analytically based scatter correction in SPECT reconstructions*, Nuclear Science Symposium Conference Record, 2003 IEEE, vol. 4, 2003, pp. 2647–2651.
- [102] A. Voutilainen, *Statistical Inversion Methods for the Reconstruction of Aerosol Size Distribution*, Ph.D. thesis, University of Kuopio, 2001.
- [103] G. Wahba, *Spline Models for Observational Data*, Society for Industrial and Applied Mathematics (SIAM), Philadelphia, PA, 1990.
- [104] J. Wei-Min, *Intra-and inter-iteration 3D OSEM PET image reconstruction*, Biomedical Engineering: Applications, Basis and Communications (BME) **19** (2007), no. 4, 239–249.

- [105] G. Welch and G. Bishop, *An introduction to the Kalman filter*, http://www.cs.unc.edu/~welch/media/pdf/kalman_intro.pdf/, 2006 (Last visited: October 1, 2008).