

Distribution Agreement

In presenting this thesis or dissertation as a partial fulfillment of the requirements for an advanced degree from Emory University, I hereby grant to Emory University and its agents the non-exclusive license to archive, make accessible, and display my thesis or dissertation in whole or in part in all forms of media, now or hereafter known, including display on the world wide web. I understand that I may select some access restrictions as part of the online submission of this thesis or dissertation. I retain all ownership rights to the copyright of the thesis or dissertation. I also retain the right to use in future works (such as articles or books) all or part of this thesis or dissertation.

Signature:

Yunyi Hu

Date

Matrix Computations and Optimization for Spectral Computed Tomography

By

Yunyi Hu

Doctor of Philosophy

Mathematics

James G. Nagy

Advisor

Lars Ruthotto

Committee Member

Yuanzhe Xi

Committee Member

Accepted:

Lisa A. Tedesco, Ph.D.

Dean of the James T. Laney School of Graduate Studies

Date

Matrix Computations and Optimization for Spectral Computed Tomography

By

Yunyi Hu

B.S., Northeastern University (P. R. China), 2014

Advisor: James G. Nagy, Ph.D.

An abstract of

A dissertation submitted to the Faculty of the

James T. Laney School of Graduate Studies of Emory University

in partial fulfillment of the requirements for the degree of

Doctor of Philosophy

in Mathematics

2019

Abstract

Matrix Computations and Optimization for Spectral Computed Tomography

By Yunyi Hu

In the area of image science, the emergence of spectral computed tomography (CT) detectors highlights the concept of *quantitative imaging*, in which not only reconstructed images are offered, but also weights of different materials that compose the object are provided. In this thesis, we focus on optimization, preconditioning and model development of spectral CT. For simple energy discriminating detectors, a nonlinear optimization framework is built on a Poisson likelihood estimator and bound constraints. A nonlinear interior-point trust region method is implemented to compute the solution. For energy-windowed spectral CT, a nonlinear least squares approach is proposed to describe the problem and under bound constraints, a two-step method using the projected line search and the trust region approach, incorporated with an adaptive preconditioner, is used to solve the problem. In addition, a weighted least squares formulation is derived from the Gaussian noise assumption and another preconditioner that is based on rank-1 approximation is introduced to obtain robust reconstruction. The Fast Iterative Shrinkage-Thresholding Algorithm (FISTA), along with a projection step, is used to calculate the solution iteratively. Compared with a direct solver, a two-step model is developed using an auxiliary variable. With this two-step model, a row-wise computational method is proposed, which further reduces memory requirements and improves solution accuracy. Numerous numerical experiments are conducted to indicate the strength of methods and real-life examples are presented to show possible applications.

Matrix Computations and Optimization for Spectral Computed Tomography

By

Yunyi Hu

B.S., Northeastern University (P. R. China), 2014

Advisor: James G. Nagy, Ph.D.

A dissertation submitted to the Faculty of the
James T. Laney School of Graduate Studies of Emory University
in partial fulfillment of the requirements for the degree of
Doctor of Philosophy
in Mathematics

2019

Acknowledgments

I would like to express my deep gratitude to my advisor, Professor James G. Nagy. In academia, Jim helped me set up my current projects, taught me the necessary skills, introduced me to other great researchers, inspired me with new ideas and encouraged me when I met obstructions. Without him, I really doubt if I could move my research to the current point. He is not only a great coach in academia, but also a great educator in my life. He supported me when I was distracted by daily trivia and gave me priceless advices on my career development. All of these experiences build my maturity and profession, which are significant preparations for my future of life. I feel exceptionally lucky of being a student of him and words cannot fully reveal my gratefulness.

I also want to show my great thankfulness to Professor Martin S. Andersen who offered me the chance to visit Technical University of Denmark. I really enjoyed the life in Copenhagen where I immersed myself in work. Under his supervision, I gained a deeper insight into my research projects and we made great process on proposing new methods. This experience is an invaluable treasure in my life.

In addition, I would like to express appreciation to Professor Lars Ruthotto, Professor Michele Benzi and Professor Yuanzhe Xi, who guided me to explore the charm of mathematics and trained me with their great efforts. Moreover, I am grateful to other faculty and staff in the Department of Mathematics. They are always friendly and approachable when I need any helps.

Finally, I am indebted to my friends and family. I am fortunate to meet my friends who are from all around the world. I always remember the times when we shared our happiness and sadness with each other. I also want to express special thanks to my parents. Thank you for offering me the entire freedom and allowing me to choose whatever I like.

Contents

1	Introduction	1
1.1	Contributions of Work	2
1.2	Outline of Thesis	5
2	X-ray Computed Tomography	7
2.1	Physical Background	9
2.2	The Beer-Lambert's Law	10
2.3	Classical Reconstruction Methods	12
2.3.1	The Radon Transform and the Inverse Radon Transform . . .	12
2.3.2	Algebraic Reconstruction Technique	15
2.3.3	Statistical Reconstruction Methods	18
2.4	Ill-posed Inverse Problems	20
3	Nonlinear Optimization for the Energy Integrating Detector Model	22
3.1	The Energy Integrating Detector Model	24
3.2	Poisson Log-likelihood Function	28
3.3	Implementation of Nonlinear Interior Point Trust Region Method . .	33
3.4	Numerical Experiments	40
3.4.1	Full Angle Reconstruction	43
3.4.2	Limited Angle Reconstruction	46
3.5	Conclusions and Remarks	48

4	Nonlinear Optimization for Energy-windowed Spectral Computed Tomography	50
4.1	The Energy-windowed Spectral CT Model	52
4.2	Problem Set-up and Preconditioning	56
4.2.1	The Constrained Least Squares Problem	56
4.2.2	Preconditioning of the Hessian	58
4.3	Optimization and Regularization	63
4.3.1	Optimization with the Proposed Preconditioner	63
4.3.2	Regularization and Scaling	68
4.4	Numerical Experiments	70
4.5	Conclusions and Remarks	76
5	Preconditioning and Optimization for Energy-windowed Spectral Computed Tomography	78
5.1	The Weighted Least Squares Problem	81
5.2	Preconditioning and Regularization	87
5.2.1	Preconditioning	87
5.2.2	Regularization	91
5.3	FISTA and Projections	93
5.3.1	FISTA	93
5.3.2	Lipschitz Constant	94
5.3.3	Projections	95
5.4	Numerical Experiments	97
5.5	Conclusions and Remarks	103
6	A Two-Step Method for Energy-windowed Spectral Computed Tomography	105
6.1	The Two-step Method	107

6.1.1	The Framework of Two-step Model	107
6.1.2	A Solution to the Two-step Model	110
6.2	The Coupled Method	115
6.3	Numerical Experiments	119
6.4	Conclusions and Remarks	122
7	Conclusions and Future Works	124
	Bibliography	126

List of Figures

2.1	Illustration of a 2D CT imaging setup	10
2.2	Explanation of the Beer-Lambert's law	11
2.3	Explanation of the Radon transform	13
2.4	Explanation of the line integral over the j -th image basis function	16
3.1	Photon flux density versus photon energy	40
3.2	The linear attenuation curves for adipose, air and calcium	41
3.3	The true images for air, adipose and calcium	42
3.4	The reconstructed images for air, adipose and calcium for the full CT reconstruction	44
3.5	The plot of relative errors for air, adipose and calcium for the full CT simulation	45
3.6	The plot of decrease of the objective function value	45
3.7	The reconstructed images for air, adipose and calcium with 90 degrees projection	47
3.8	The plot of relative errors of air, adipose and calcium for the 90 degrees limited angle simulation	48
4.1	The comparison of eigenvalues before and after preconditioning.	62
4.2	The original images for plexiglass and PVC	71
4.3	The linear attenuation coefficients for plexiglass and PVC	71

4.4	Detector bins and photon flux density	72
4.5	The reconstructed images for plexiglass and PVC.	73
4.6	The relative errors for Plexiglass and PVC	74
4.7	The decay of norm of the gradient	75
5.1	The original material maps for plexiglass and PVC	97
5.2	The reconstructed images for plexiglass and PVC	99
5.3	The relative errors of each iteration for plexiglass and PVC	100
5.4	The MSE of each iteration for plexiglass and PVC	101
5.5	The PSNR of each iteration for plexiglass and PVC	101
5.6	The SSIM of each iteration for plexiglass and PVC	102
5.7	The decay of norm of the gradient for overall materials	102
5.8	The comparison of decay of related errors with different preconditioners	103
6.1	The soft shrinkage function with a bound constraint	118
6.2	The original material maps for plexiglass and PVC	119
6.3	The reconstructed material maps for plexiglass and PVC using the two-step method	120
6.4	The relative errors of two material maps solved by the two-step method	120
6.5	The reconstructed material maps for Plexiglass and PVC using the coupled method	121
6.6	The relative errors of two material maps solved by the coupled method	122

List of Tables

4.1	Geometry parameters of CT machine	72
4.2	The comparison of CG iterations	76
5.1	The comparison of condition numbers	91

Chapter 1

Introduction

An active area of interest in tomographic imaging is *quantitative imaging*, where in addition to producing an image, information about the material composition of the object is recovered. In order to obtain the information of material composition, it is necessary to better model of the image formation (i.e., forward) problem and/or to collect additional independent measurements. In x-ray computed tomography (CT), better modeling of the physics can be done by using the more accurate polyenergetic representation of source x-ray beams. In addition, recent advances in engineering have produced detectors that are made up of several energy windows and each energy window is assumed to detect a specific range of energy spectra. With this technique, a nonlinear matrix equation is formulated to represent the discretized process of attenuation of x-ray intensity

$$\mathbf{Y} = \exp(-\mathbf{A}\mathbf{W}\mathbf{C}^T) \mathbf{S} + \boldsymbol{\varepsilon}, \quad (1.1)$$

where \mathbf{Y} is a matrix that gathers the projected data of each energy window in the corresponding column and the exponential operator is applied element-wise (i.e., it is not a matrix function). \mathbf{A} is a matrix that is related to the quantitative information of ray trace and \mathbf{C} is a matrix that contains linear attenuation coefficients for par-

ticular (known) materials at specified energies. \mathbf{S} is the matrix that accumulates the spectrum energies for each energy window in the corresponding column. We assume that these data are known and the target is to solve the unknown weight matrix \mathbf{W} . \mathbf{W} real and is of size N_v by N_m , where N_v is the number of voxels (pixels if 2D) for each material map and N_m is the number of materials. The derivation of this equation is not straightforward and we will go into details later.

1.1 Contributions of Work

For different noise distributions, model (1.1) can be transformed into various formats and the methods used to compute solutions are diversified. If we assume that we only have one energy window, then the matrix equation (1.1) is reduced to a nonlinear system of the form

$$\mathbf{y} = \exp(-\mathbf{A}\mathbf{W}\mathbf{C}^T) \mathbf{s} + \boldsymbol{\eta}, \quad (1.2)$$

where \mathbf{y} is the vector of projected data, \mathbf{s} collects the information of energy spectrum and $\boldsymbol{\eta}$ represents the noise vector. For this problem, we transform it into a nonlinear optimization problem that involves both equality and inequality constraints. In order to keep the second order derivative positive semi-definite, we propose a modified Hessian. In addition to the modified Hessian, a problem-specific nonlinear interior-point trust region method is implemented to solve this problem. Moreover, total variation regularization is applied to stabilize the solutions. Both for full CT and limited angle cases, we can obtain images of high quality with this method.

If we take the assumption of multiple energy windows and each energy window can only detect a specific range of energy, then the basic model is the same as (1.1). Under the constraint of nonnegative weights, we can build a nonlinear optimization

problem as

$$\begin{aligned} \min_{\mathbf{w}} \quad & \frac{1}{2} \|\mathbf{y} - (\mathbf{S}^T \otimes \mathbf{I}) \exp\{-(\mathbf{C} \otimes \mathbf{A}) \mathbf{w}\}\|_2^2 \\ \text{subject to} \quad & \mathbf{w} \geq \mathbf{0}. \end{aligned} \tag{1.3}$$

In (1.3), \otimes represents the Kronecker product and the exponential operator is applied point-wise. $\mathbf{w} \geq \mathbf{0}$ indicates that each entry in \mathbf{w} should be bounded below by zero. Moreover, \mathbf{y} and \mathbf{w} are vectorization of \mathbf{Y} and \mathbf{W} , respectively. For problem (1.3), we formulate it as a nonlinear least squares problem, which is solved by a Gauss-Newton scheme. We show that if the object contains a mixture of materials with one known to be sparsely represented, then a combination of generalized Tikhonov and ℓ_1 regularization can be very effective in producing high quality quantitative reconstructions. Because the approximate Hessian system in the Gauss-Newton scheme is very ill-conditioned, we propose a preconditioner that effectively clusters eigenvalues and, therefore, accelerates convergence when the conjugate gradient method is used to solve the linear subsystems. To implement the preconditioner, a two-step method is used to guarantee faster convergence and to provide further stabilization. In particular, in the first step we compute the approximate Cauchy point, and then in the second step we set up and solve a quadratic programming problem. Numerical experiments illustrate the convergence, effectiveness, and significance of the proposed method.

To avoid solving a nonlinear optimization problem, we might consider taking advantage of the Gaussian noise assumption and transform it into a weighted least squares problem. In this case, it is necessary to assume that \mathbf{S} is square and invertible. Moreover, for the noise term \mathcal{E} , we assume that $E_{il} \sim \mathcal{N}(0, y_{il})$ for each component E_{il} in \mathcal{E} and y_{il} in \mathbf{Y} . A linearization technique is used to transform the nonlinear equation (1.1) into an optimization problem that is based on a weighted

least squares term and a nonnegative bound constraint of the form

$$\begin{aligned} \min_{\tilde{\mathbf{w}}} \quad & \frac{1}{2} \|\mathcal{A}\tilde{\mathbf{w}} - \mathbf{b}\|_{\Sigma^{-1}}^2 \\ \text{subject to} \quad & (\mathbf{M} \otimes \mathbf{I}) \tilde{\mathbf{w}} \geq \mathbf{0}, \end{aligned} \tag{1.4}$$

where $\tilde{\mathbf{W}} = \mathbf{W}\mathbf{M}^{-T}$, $\tilde{\mathbf{w}} = \text{vec}(\tilde{\mathbf{W}})$, $\tilde{\mathbf{C}} = \mathbf{C}\mathbf{M}$, $\mathcal{A} = \tilde{\mathbf{C}} \otimes \mathbf{A}$, $\tilde{\mathbf{y}} = (\mathbf{S}^{-T} \otimes \mathbf{I}) \mathbf{y}$ and $\mathbf{b} = -\log(\tilde{\mathbf{y}})$. The matrix \mathbf{M} represents the corresponding preconditioner and $\|\cdot\|_{\Sigma^{-1}}^2$ is a weighted least squares term such that $\|\mathcal{A}\tilde{\mathbf{w}} - \mathbf{b}\|_{\Sigma^{-1}}^2 = (\mathcal{A}\tilde{\mathbf{w}} - \mathbf{b})^T \Sigma^{-1} (\mathcal{A}\tilde{\mathbf{w}} - \mathbf{b})$. To solve this optimization problem, we want to use a first order method to obtain fast and robust reconstruction. However, because of the ill-posedness, direct implementation of optimization methods does not offer us satisfactory results. Therefore, we propose a new Hessian preconditioner in order to significantly reduce the condition number, and with this preconditioner, we implement a highly efficient first order method, Fast Iterative Shrinkage-Thresholding Algorithm (FISTA), to achieve substantial improvements on convergence speed and image quality. We also use a combination of generalized Tikhonov regularization and ℓ_1 regularization to stabilize the solution. With the introduction of new preconditioner, a linear inequality constraint is also added. In each iteration, we decompose this constraint into small-sized problems that can be solved with fast optimization solvers. Furthermore, we conduct numerical experiments to indicate strengths of the proposed method and potential improvements.

Rather than solving \mathbf{W} directly in one step, we can use an auxiliary variable, $\mathbf{X} = \mathbf{A}\mathbf{W}$, and construct a two-step framework as

$$\begin{aligned} \mathbf{Y} &= \exp(-\mathbf{X}\mathbf{C}^T) \mathbf{S} + \mathbf{E}, \\ \mathbf{X} &= \mathbf{A}\mathbf{W}. \end{aligned} \tag{1.5}$$

For the first step, we can reuse the Gaussian assumption of noise and build a row-wise

model as

$$(\mathbf{x}_i)_{\text{ml}} = \underset{\mathbf{x}_i}{\operatorname{argmin}} \left\{ \left\| \mathbf{C}\mathbf{x}_i - \hat{\mathbf{b}}_i \right\|_{\Sigma_{\hat{\mathbf{b}}_i}^{-1}}^2 \right\}, \quad (1.6)$$

where $\Sigma_{\hat{\mathbf{b}}_i}$ is the noise covariance matrix related to $\hat{\mathbf{b}}_i$ and $(\mathbf{x}_i)_{\text{ml}}$ is the maximum likelihood estimator to the i -th row of \mathbf{X} . Even though this model is based on the maximum likelihood estimator, we can compute the solution corresponding to each row slice one at a time, and sum when they all are computed. Each small-sized problem does not depend on others, so these can also be solved in parallel. For the second step, we need to quantify the noise propagation and solve a least squares optimization problem under bound constraints. To solve the corresponding optimization problem, we use FISTA with projections onto the boundary to obtain faster convergence and higher accuracy. Instead of using the two-step method, we can build a coupled optimization problem based on (1.5) and solve \mathbf{W} directly. Numerical experiments show that the expense is lower and reconstructed images are of higher quality using the two-step framework compared with the coupled method.

1.2 Outline of Thesis

To reconstruct spectral CT images, we are required to solve a complicated nonlinear equation. Distinct assumptions might provide us with different models and various solutions. This thesis mainly focuses on how to build efficient models and find superior solutions.

The thesis is organized as follows. In Chapter 2, we first review the physical background and basics of computed tomography (CT) and present three classical methods to compute solutions. Since the traditional CT does not take energy spectra into consideration, we introduce the spectral CT model in Chapter 3. We also build a nonlinear optimization problem and implement a nonlinear interior-point trust region method to solve it. In Chapter 4, we switch to the energy-windowed spectral CT

model and build a nonlinear least squares problem based on it. An adaptive preconditioner based on the Gauss-Newton Hessian approximation is also included in Chapter 5. In Chapter 6, we use the assumption that noise is normally distributed to transform the energy-windowed spectral CT model into a weighted least squares problem. Another efficient preconditioner inspired by the interlacing of Kronecker products and diagonal matrices is also presented in Chapter 5. In Chapter 6, we move further to consider a two-step model as well as introduce a solution that is based on row slicing and weighted least squares optimization. As a comparison, the coupled method is also included in Chapter 6. Numerical experiments are conducted to show the strengths of proposed methods and conclusions and comparisons are drawn in the end of each chapter.

Chapter 2

X-ray Computed Tomography

In 1971 Godfrey Hounsfield [4] opened a new window for the world with his invention of the computed tomography (CT) machine. CT machines project x-ray beams from a known source through an object, which are then received by detectors. The energy of the source x-ray beams are known, and the detectors measure the energy after they pass through the object. The energy difference between the source and the detector depends on the *attenuation* properties of the object, and is modeled by the Beer's law [34]. If enough measurements are obtained, then by solving an inverse problem arising from the Beer's law, the inner structure of an object can be reconstructed. Because the Beer's law results in a very challenging nonlinear inverse problem, most image reconstruction algorithms are based on linear approximations. The linear approximations allow for very fast algorithms, and images they produce are often quite good. Thus CT immediately gained popularity for medical diagnostics, but the technique is also used widely in industry to, for example, inspect inanimate objects for defects, or in security to look for weapons and other dangerous materials.

Current CT machines mostly use single-energy tubes to conduct scanning. For single-energy CT, we assume that the x-ray tube only emits a uniform energy and this energy is used to estimate the attenuation properties of the object. Using the

Radon transform [14], a linear system is constructed to represent an approximation to this physical process. Under certain conditions, an image obtained from solving the corresponding linear system is clear enough to identify necessary information. However, the assumption of single uniform energy is only a simplification of the underlying physics. In actuality, the x-ray beams consist of a spectrum of energies, rather than a single energy. Because of this simplification, precise estimation of the attenuation properties of the object is very challenging, and often leads to appearance of so-called beam-hardening artifacts [45].

In 2006, the invention of dual-source CT machines refreshed this field [19]. Using two x-ray tubes with different voltages, two sets of projected data corresponding to two energy spectra are obtained to conduct image reconstruction. Basically, the attenuation properties of the object can be represented as a function of two variables, position and energy. With an extra energy, we can obtain more information about the object. On the one hand, the beam-hardening artifacts generated by single-energy CT can be dramatically reduced using dual-source CT, which offers images of higher quality. On the other hand, not only the inner structure but also the material separation of an object can be acquired, which provides the fundamental basics of quantitative spectral CT methods. Even if the idea of basis-material decomposition (BMD) was proposed by Alvarez and Macovski in 1976 [1], it is suddenly in the spotlight after the emergence of dual-source CT machines. If the contrast of different materials is small, it is hard to differentiate them by investigating the gray images generated by single-energy CT. For example, in breast imaging the object consists mainly of glandular and adipose tissues, which have similar densities and therefore are difficult to separate by traditional techniques. With spectral CT methods, it is likely that we can reconstruct the image of these two materials using different linear attenuation coefficients corresponding to distinct energies.

Nowadays, spectral CT still evolves both in theory and practice. If we assume

that the detector of CT machines has multiple energy windows and further assume each energy window can detect only a specific range of energy, then it is an energy-windowed spectral CT model. This model is more complicated, but on the other side, it can provide us with more flexibility on mathematical theories as well as more quantitative information of images. In practice, one step forward from dual-source CT machines, there are multi-energy CT machines that explore improvements over more than two sources. However, the real-life application of multi-energy CT is limited by the overlap of various energy spectra and this limitation might be overcome by photon counting detectors. On the one hand, we can obtain the material decomposition of an object and higher quality images with spectral CT. On the other hand, we want to lower the radiation dose to reduce risks of causing other health issues. To reach these goals, spectral CT is still an active research area and more promising results might show up in the near future.

2.1 Physical Background

In computed tomography, x-ray beams are emitted (usually in a cone) from a source at known energies and are directed to pass through an object under investigation, after which the remaining energy of the x-ray beams are measured at a detector. See Figure 2.1 for a 2D illustration. The amount of energy lost as the x-ray beams pass through the object is referred to as attenuation. The amount of attenuation depends on the energy of the x-ray beams, and on the material through which it penetrates; low dose energies are more easily attenuated, and denser materials have higher attenuation properties. The detector is typically partitioned into a grid of bins; in a full 3D model, these are often loosely referred to as detector pixels. If the 3D object is discretized into a grid of small volume elements (called voxels), then each voxel can be associated with a particular attenuation value, referred to as an

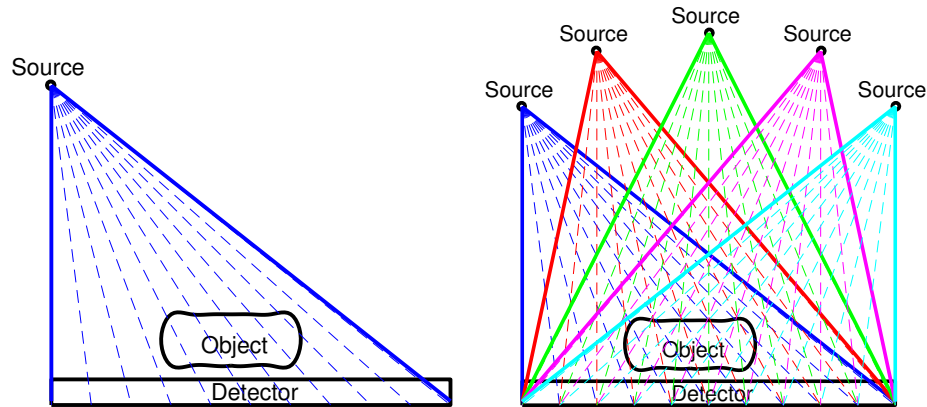


Figure 2.1: Illustration of a 2D CT imaging setup. The left illustrates how x-ray beams are emitted from the source in a cone, pass through an object of interest, and are then measured at the detector. The right illustrates how the source might be rotated around the object to collect additional data. In this illustration we assume the detector remains stationary, which is often the case in limited angle tomography applications such as tomosynthesis, but it should be obvious how the illustration would be modified if the detector rotates with the source.

attenuation coefficient. The problem of CT image reconstruction is to determine these attenuation coefficients from a sequence of measured projection data, which is obtained by rotating the source (at least partially) around the object; again, we refer to Figure 2.1 for a 2D illustration of the data collection process.

2.2 The Beer-Lambert's Law

As shown in Figure (2.1), the change of intensity before and after illumination can be described by the Beer's law or the Beer-Lambert's law [5, 38]:

$$\frac{dI}{dx} = -\mu(x)I, \quad (2.1)$$

where I represents the intensity and $\mu(x)$ is the linear attenuation coefficient at x . It shows that the change of intensity at a specific location equals to the product of the intensity at that location and the corresponding linear attenuation coefficient. The

physical meaning explained by Equation (2.1) is shown in Figure (2.2). In Figure (2.2), an object is located in the center and a x-ray beam with intensity I_0 penetrates it and the detector captures this x-ray beam with intensity I_1 . $I_1 < I_0$ and the attenuation characteristic of the object causes the reduction of intensity and this change is also captured by Equation (2.1). For non-homogeneous material $\mu(x)$, the

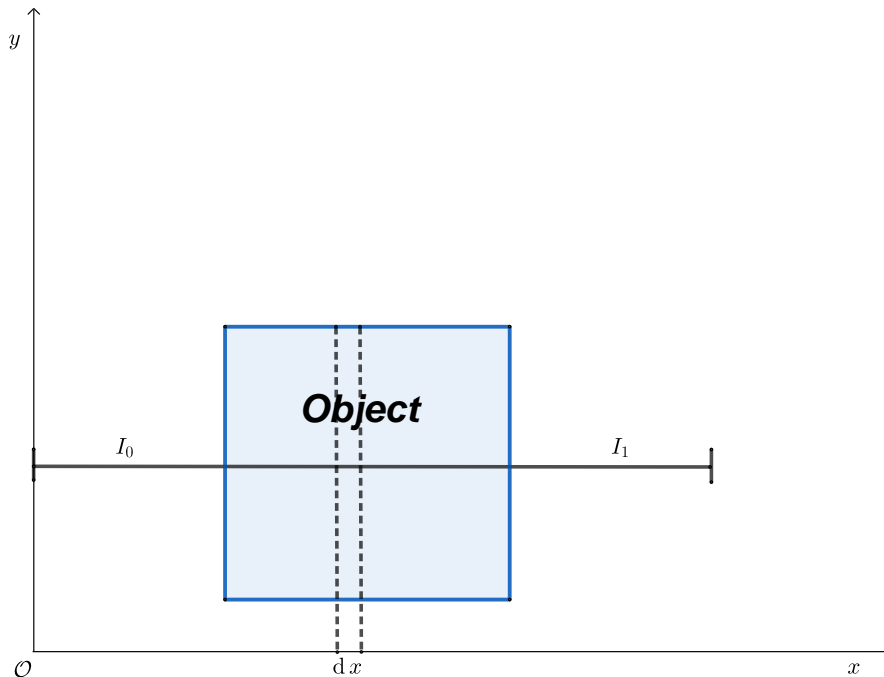


Figure 2.2: Explanation of the Beer-Lambert's law. The blue brick represents the object. There is a x-ray beam illuminating the object with the intensity I_0 and it is received by the detector with intensity I_1 .

ordinary differential equation (2.1) has an analytical solution

$$I = I_0 \exp \left\{ - \int_l \mu(x) dx \right\}. \quad (2.2)$$

As long as we know I_0 and I_1 , we can try to compute $\mu(x)$ using mathematical techniques. Basically, we have three mainstream methods to solve it: the analytical reconstruction using the Radon transform and the filtered back-projection, Algebraic Reconstruction Technique (ART) and statistical reconstruction methods.

Basically, Equation (2.1) has two underlying assumptions:

- x-ray beams pass through an object and they are not bent or diffracted by this object. They travel along a straight line.
- x-ray beams are monochromatic, which means they are of only one single energy.

The first assumption is realistic because x-rays are of high energy and thus of short wavelength. In most cases, the diffraction effect can be ignored and it is a reasonable approximation. However, the second assumption is not realistic, because in reality x-rays are made up of a spectrum of energies rather than a single energy. A more comprehensive equation might be written as

$$I = \int_E S(e) \exp \left\{ - \int_l \mu(x) dx \right\} de, \quad (2.3)$$

where $S(e)$ is the energy intensity at energy level e and E is the energy spectrum. Compared with Equation (2.2), Equation (2.3) has two integrals and the nonlinearity might cause more difficulties to compute the solution. In this chapter, we will focus on Equation (2.2) in most cases and discuss the classical methods to find solutions.

2.3 Classical Reconstruction Methods

2.3.1 The Radon Transform and the Inverse Radon Transform

The Radon transform was introduced by Johann Radon [50] in 1917. If we assume θ is expressed as $\theta = (\cos \theta, \sin \theta)^T$, then the Radon transform is an integral transform such that

$$\mathcal{R}\mu(s, \theta) = \int_{l_{s, \theta}} \mu(x) dx, \quad (2.4)$$

where $l_{s,\theta}$ is a signed line represented by the signed orthogonal distance s and the angle θ . $l_{s,\theta}$ can be represented as an equation

$$l_{s,\theta} = \{(x, y) | x \cos \theta + y \sin \theta = s\}. \quad (2.5)$$

The Radon transform is illustrated in Figure 2.3, where $l_{s,\theta}$ is drawn perpendicular

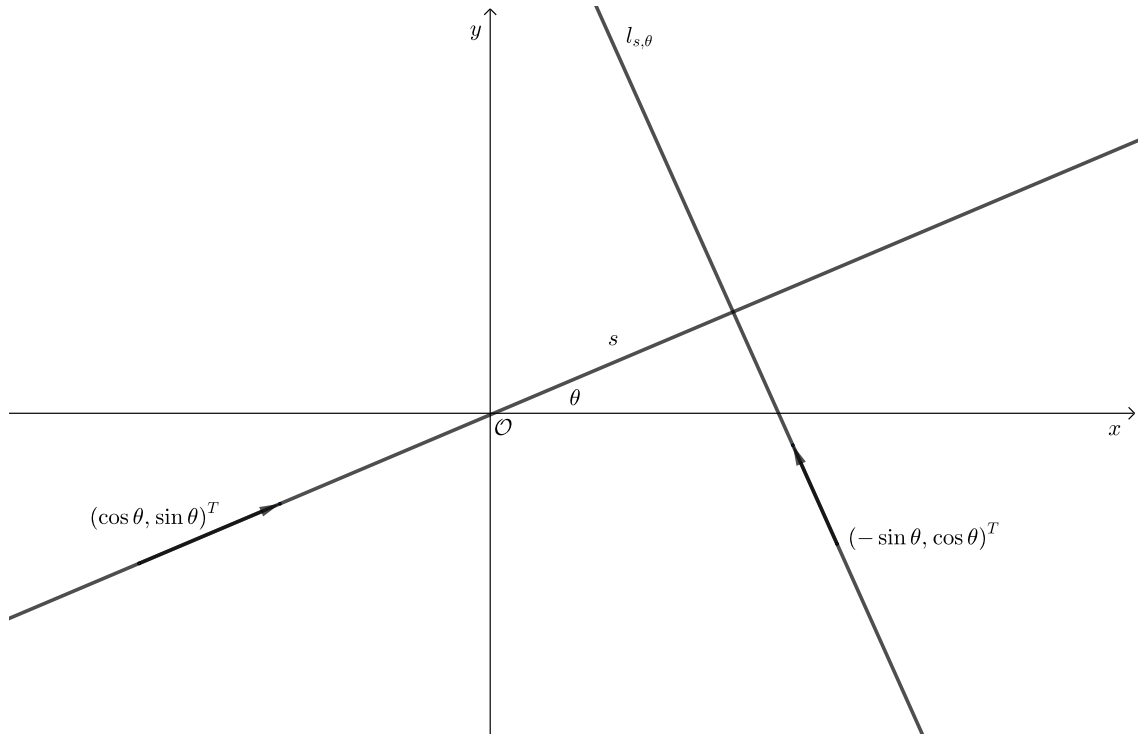


Figure 2.3: Explanation of the Radon transform. The line $l_{s,\theta}$ is perpendicular to the ray passing through the origin. The distance between these two lines is s .

to the original line.

Recall that from the Beer's law, the intensity I can be represented on behalf of s and θ , $I_{s,\theta}$. Then Equation (2.2) is equivalent to

$$I_{s,\theta} = I_0 \exp \left\{ - \int_{l_{s,\theta}} \mu(x) dx \right\}. \quad (2.6)$$

If we move all intensities in Equation (2.6) into one side and take a logarithm, then we can obtain the Radon transform. So the connection between the Beer-Lambert's

law and the Radon transform is

$$\mathcal{R}\mu(s, \theta) = \int_{l_{s,\theta}} \mu(x) dl = -\log\left(\frac{I_{s,\theta}}{I_0}\right). \quad (2.7)$$

Moreover, the Radon and Fourier transforms have a strong connection as

$$\widetilde{\mathcal{R}}\mu(r, \theta) = \int_{-\infty}^{\infty} \mathcal{R}\mu(s, \theta) e^{-irs} ds = \hat{\mu}(r\theta). \quad (2.8)$$

If we represent x in a polar coordinate form as $x = \rho n_\phi$, where $|\rho|$ is the distance between the current point and the origin, then the inverse Radon transform is expressed as

$$\mu(x) = \frac{1}{(2\pi)^2} \int_0^\pi \int_{-\infty}^{\infty} \widetilde{\mathcal{R}}\mu(r, \theta) |r| e^{ir\langle x, \theta \rangle} dr d\theta, \quad (2.9)$$

if $\mu(x)$ and $\widetilde{\mathcal{R}}\mu(r, \theta)$ are both absolutely integrable. The variable $|r|$ in Equation (2.9) can amplify the noise and the reconstructed images might not be satisfactory. This phenomenon also matches a slow decay of singular values of the Radon transform operator. Small singular values can exaggerate the influence of high-frequency noise so as to corrupt the reconstructed images. Compared with the inverse Radon transform, a more useful technique is called filtered back-projection (FBP) [12]. As the name indicates, it adds a filter to the inverse transform formula (2.9). In FBP, the radial integral is regarded as a filter used in the Radon transform

$$\mathcal{GR}(t, \theta) = \frac{1}{2\pi} \int_{-\infty}^{\infty} \widetilde{\mathcal{R}}\mu(r, \theta) |r| e^{irt} dr. \quad (2.10)$$

On the other hand, the angular integral is explained as a back-projection of the previous transform

$$\mu(x) = \frac{1}{2} \int_0^\pi \mathcal{GR}(\langle x, \theta \rangle, \theta) d\theta. \quad (2.11)$$

In contrast to the inverse Radon Transform, the filter of FBP can amplify high-

frequency components while suppressing low-frequency components so the process is more robust and reconstructed images are of higher quality and less noisy. Basically, it takes four steps to reconstruct images using FBP:

- For each angle θ , we collect measured data as (2.7).
- Compute the filter as a radial integral in (2.10).
- Calculate the back-projection as (2.11).
- Collect the reconstructions of all directions and synthesize the results.

2.3.2 Algebraic Reconstruction Technique

Rather than computing the linear attenuation coefficient μ analytically, we can also solve μ using iterative methods. If we use polar coordinates and discretize Equation (2.2) with respect to the j -th projection and the i -th detector pixel, we can obtain

$$I_i = I_0 \exp \left\{ - \int_{t \in l_{i,j}} \mu(\vec{r}(t)) dt \right\}. \quad (2.12)$$

To be convenient, we let $b_i = -\log \frac{I_i}{I_0}$ and can obtain that

$$b_i = \int_{t \in l_{i,j}} \mu(\vec{r}(t)) dt. \quad (2.13)$$

If we further decompose $\mu(\vec{r}(t))$ into an expansion

$$\mu(\vec{r}(t)) = \sum_j \mu_j \phi_j(\vec{r}(t)), \quad (2.14)$$

where $\phi_j(\vec{r}(t))$ is a basis function of the image representation. The line integral of the basis function, $a_{i,j}$, is the length of the x-ray beam through the j -th pixel of the target image, incident onto the i -th element of detector pixels. So $a_{i,j}$ can be

expressed as

$$a_{i,j} = \int_{t \in l_{i,j}} \phi_j(\vec{r}(t)) dt. \quad (2.15)$$

The geometric meaning of $a_{i,j}$ is shown in Figure 2.4. In Figure 2.4, a x-ray intersects

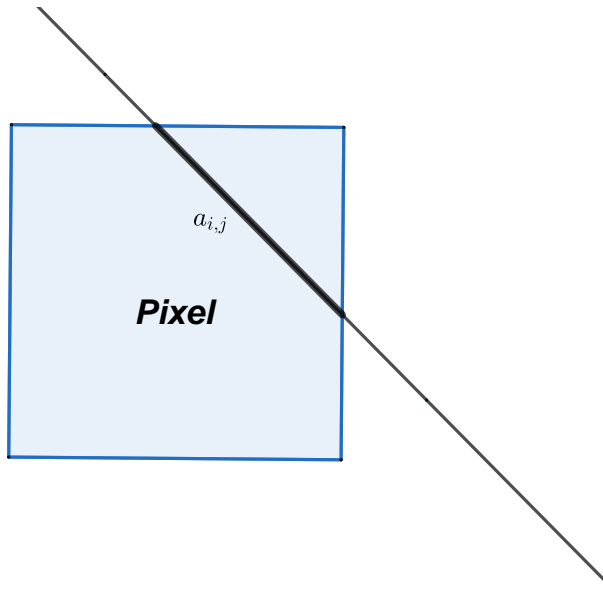


Figure 2.4: Explanation of the line integral over the j -th image basis function. The thick line inside the j -th pixel represents the distance $a_{i,j}$.

the j -th pixel of the target image and the line integral over the j -th basis function, $a_{i,j}$, corresponds to the distance of the x-ray inside the j -th pixel (the thick line).

With the previous expansion, Equation (2.13) can be expressed as

$$b_i = \sum_j \mu_j \int_{t \in l_{i,j}} \phi_j(\vec{r}(t)) dt = \sum_j a_{i,j} \mu_j. \quad (2.16)$$

If we collect $a_{i,j}$, $b_{i,j}$ and μ_j with respect to their sub-indexes, we can obtain a linear system

$$\mathbf{b} = \mathbf{A}\mathbf{u}. \quad (2.17)$$

Therefore, given the vector \mathbf{b} and the matrix \mathbf{A} , our goal is to compute the unknown variable \mathbf{u} in a robust and efficient way. Since the matrix \mathbf{A} captures the geometry of CT machines, it is not square most of time. Moreover, it can be rank-deficient if

the data we obtained are based on limited angle projections. Rather than calculate it directly, we need to find alternative methods to handle these issues.

Even if the matrix \mathbf{A} might be rectangular, we can still build a least squares problem of the form

$$\min_{\mathbf{u}} \|\mathbf{A}\mathbf{u} - \mathbf{b}\|_2^2. \quad (2.18)$$

Many methods have been proposed to solve Equation (2.18) under different scenarios. For CT problems, one of the most famous methods is called Algebraic Reconstruction Technique (ART) [25], which is also a rediscovery of the Kaczmarz method [36]. We let \mathbf{a}_i represent the i -th row of \mathbf{A} , then the ART method is an iteration framework given by

$$\mathbf{u}_{k+1} = \mathbf{u}_k + \lambda_k \frac{b_i - \mathbf{a}_i^T \mathbf{u}_k}{\|\mathbf{a}_i\|^2} \mathbf{a}_i^T, \quad (2.19)$$

where λ_k is a step size parameter in the k -th iteration. Compared with other methods, ART does not require to save the full entries of the matrix \mathbf{A} and it only needs each row slice in the current iteration. Moreover, it has been observed for tomography problems that the initial convergence speed is often fast.

Mathematically speaking, it converges to the pseudoinverse solution, $\mathbf{A}^\dagger \mathbf{b}$. If \mathbf{A} has full column rank, then \mathbf{A}^\dagger can be expressed as

$$\mathbf{A}^\dagger = (\mathbf{A}^T \mathbf{A})^{-1} \mathbf{A}^T. \quad (2.20)$$

If \mathbf{A} is rank-deficient and has full column rank, Problem (2.18) has infinitely many solutions. Among these infinitely many solutions, $\mathbf{A}\mathbf{A}^\dagger \mathbf{b}$ is the orthogonal projection of \mathbf{b} onto the column space of \mathbf{A} . The orthogonal projection offers the “shortest” distance so the solution obtained by ART might be the “best” among all possible solutions. This is also a reason why ART is popular for limited angle CT reconstruction.

In addition to the previous iterative method, we can also build a regularized least

squares problem

$$\min_{\mathbf{u}} \|\mathbf{A}\mathbf{u} - \mathbf{b}\|_2^2 + \alpha R(\mathbf{u}), \quad (2.21)$$

where α is the regularization parameter and $R(\mathbf{u})$ is the selected regularization term. To smooth the edges, we can take ℓ_2 regularization, where $R(\mathbf{u}) = \|\mathbf{u}\|_2^2$. To strengthen the edges, we can take ℓ_1 regularization, where $R(\mathbf{u}) = \|\mathbf{u}\|_1$. In the case of ℓ_2 regularization we can build an augmented system,

$$\left\| \begin{bmatrix} \mathbf{A} \\ \sqrt{\alpha} \end{bmatrix} \mathbf{u} - \begin{bmatrix} \mathbf{b} \\ \mathbf{0} \end{bmatrix} \right\|_2^2, \quad (2.22)$$

which can then be solved by ART. The regularization parameter can be chosen using methods such as L-curve [31], the discrepancy principle [17] and generalized cross-validation [23].

2.3.3 Statistical Reconstruction Methods

We know that x-ray beams are made of x-ray photons, and rather than computing the deterministic solution directly from an equation, we can explore the randomness of x-ray photons. If we assume that the detector is a photon-counting detector, then we can use statistical tools to build the model.

In Equation (2.3), if we discretize the equation with respect to the j -th projection and the i -th detector pixel, then we can obtain

$$y_i = \int_E S(e) \exp \left\{ - \int_{t \in l_{i,j}} \mu(\vec{r}(t)) dt \right\} de, \quad (2.23)$$

where y_i represents the number of photons measured at the i -th pixel of detector. If we further discretize the equation over the energy spectrum, then it can be expressed

as

$$y_i = \sum_k s_k \exp \left\{ - \int_{t \in l_{i,j}} \mu(\vec{r}(t)) dt \right\}. \quad (2.24)$$

Repeating the same steps as (2.12), (2.13) and (2.14), we can obtain

$$y_i = \bar{s}_i \exp \{ - [\mathbf{A}\mathbf{u}]_i \}, \quad (2.25)$$

where \bar{s}_i is the average photon density of the j -th projection. If we both consider the errors introduced in the photon diffusion and counting, we can add a noise term

$$y_i = \bar{s}_i \exp \{ - [\mathbf{A}\mathbf{u}]_i \} + \eta_i. \quad (2.26)$$

One realistic assumption is that the projected data are Poisson distributed as

$$y_i \sim \text{Poisson}(\bar{s}_i \exp \{ - [\mathbf{A}\mathbf{u}]_i \}). \quad (2.27)$$

With this assumption, we can build a maximum likelihood function

$$f_{\mathbf{y}}(\mathbf{y}; \mathbf{u}) = \prod_i \frac{\bar{s}_i^{y_i} \exp \{ -y_i [\mathbf{A}\mathbf{u}]_i \} \exp (-\bar{s}_i \exp \{ - [\mathbf{A}\mathbf{u}]_i \})}{y_i!}. \quad (2.28)$$

If we ignore the constant terms, the corresponding likelihood function is

$$L(\mathbf{u}; \mathbf{y}) = \prod_i \exp \{ -y_i [\mathbf{A}\mathbf{u}]_i \} \exp (-\bar{s}_i \exp \{ - [\mathbf{A}\mathbf{u}]_i \}). \quad (2.29)$$

The log-likelihood function can be expressed as

$$l(\mathbf{u}; \mathbf{y}) = \sum_i (-y_i [\mathbf{A}\mathbf{u}]_i - \bar{s}_i \exp \{ - [\mathbf{A}\mathbf{u}]_i \}). \quad (2.30)$$

To maximize the log-likelihood function, it is equivalent to minimizing the negative

log-likelihood function. So the objective function is

$$l(\mathbf{u}; \mathbf{y}) = \sum_i (y_i [\mathbf{A}\mathbf{u}]_i + \bar{s}_i \exp\{-[\mathbf{A}\mathbf{u}]_i\}). \quad (2.31)$$

With this objective function, we can build an optimization problem as

$$\begin{aligned} \underset{\mathbf{u}}{\operatorname{argmin}} \quad & l(\mathbf{u}; \mathbf{y}) \\ \text{subject to} \quad & \mathbf{u} \geq 0. \end{aligned} \quad (2.32)$$

This optimization problem consists of a nonlinear objective function and bound constraints. Therefore, we cannot regard it as a linear system and compute the solution directly. Even if the nonlinearity provides more challenges for us, it also indicates numerous possibilities. To obtain the maximum likelihood estimators, nonlinear optimization skills are used to find the optimal solutions.

2.4 Ill-posed Inverse Problems

Using Hadamard's [28] definition, an inverse problem is well-posed as long as it satisfies three conditions:

- Existence: there should be at least one solution.
- Uniqueness: the solution should be unique.
- Stability: with initial conditions, the solution should depend continuously on the change.

Problems that fail to meet any of these three criteria are classified as ill-posed problems. For numerical analysis, it has corresponding explanations:

- The forward operator is compact, so it does not have a continuous inverse in the infinite-dimensional space, and the discrete problem inherits this property of

ill-posedness. Therefore, computing an accurate approximation is a challenging problem.

- The forward operator or linear system has small singular values. When we solve inverse problems, the inversion of small singular values will amplify the noise so the solution might be far away from the truth.

In our cases, the computed tomography problem is an ill-posed inverse problem since the forward operator has a decay of small singular values and the rank-deficient system has multiple solutions. To mitigate the effects, three classical methods use different techniques:

- The filtered back-projection uses a filter to emphasize the high-frequency components and reduce the influence of noise.
- Algebraic Reconstruction Technique takes orthogonal projections and regularizations to obtain a better solution. If the regularization operator $R(\mathbf{u})$ is omitted, the early termination of the iteration can be used to avoid noise amplification in the solution.
- Statistical reconstruction methods apply randomness of photons and compute the maximum likelihood estimator as evidence.

The solutions obtained using these classical methods are of higher quality. However, these results are based on the assumption of either a monochromatic x-ray source or an approximation of average of photon flux density. Furthermore, we cannot tell the composition of the target object we obtain as well as the percentages of different materials. In the following sections, we will focus on the spectral CT models, the material decomposition and optimization frameworks used to compute robust and efficient solutions.

Chapter 3

Nonlinear Optimization for the Energy Integrating Detector Model

In this chapter we focus on spectral computed tomography with a single energy discriminating detector. To reconstruct material maps of each composition, we need to find solutions to an inverse problem of the form

$$\mathbf{y} = \exp(-\mathbf{A}\mathbf{W}\mathbf{C}^T) \mathbf{s} + \boldsymbol{\eta}, \quad (3.1)$$

where \mathbf{y} is a vector of projection data, \mathbf{s} is a vector containing spectral information of the corresponding energies of the source x-ray, and $\boldsymbol{\eta}$ is the noise term. The exponentiation is done element-wise. \mathbf{A} is a matrix that is related to the quantitative information of ray trace, \mathbf{C} is a matrix of (known) material specific attenuation coefficients, and \mathbf{W} is a matrix of unknowns, whose columns correspond to the weights of each of known materials of the object being imaged.

To solve this problem, Elbakri and Fessler [15] and Chung et al. [11] suggest that we could use a 2-material model with a polyenergetic assumption on the source x-ray. Moreover, Mejia-Bustamante et al. [7] extended this idea, and provided a GPU implementation [6]. We remark that the problem can be reformulated in terms of mass

attenuation coefficients (see, e.g., [15]), which encodes density information into the mathematical model. This can be important in cases when materials have the same chemical composition but different densities. However, it does not fundamentally change our proposed optimization approach, so in the remainder of the chapter we use the model based on linear attenuation coefficients. To solve the nonlinear inverse problem in [6, 7], a maximum likelihood function was used to represent the evidence with respect to parameters \mathbf{W} . The gradient descent method is used to solve the resulting optimization problem, with an implicit enforcement of the constraint that in each voxel, the weights across all materials should sum to one.

In this chapter, we build an optimization problem with Poisson maximum likelihood function and use a problem-specific nonlinear interior point trust region method to solve it. Moreover, we provide a modified Hessian that is close to the true Hessian and is also positive semi-definite. The interior point method has proven to be efficient and stable for solving nonlinear optimization problems, and with the implementation of our modified Hessian, we can simultaneously obtain an accurate approximation of the true Hessian and avoid negative curvature. Numerical experiments have shown very promising results for this method.

In Section 3.1, we present the general polyenergetic model for spectral computed tomography and derive the specific mixed attenuation model from the general model. The discretization of this model is included in Section 3.1 as well. In Section 3.2, we revise this model to one that is more amenable to numerical implementation and formulate an optimization problem that is based on Poisson likelihood function. The standard form of this optimization problem and the method to solve it are discussed in Section 3.3. Moreover, numerical experiments for both full CT and limited angle reconstruction are presented in Section 3.4. In Section 3.5, we conclude with merits and limitations about the model and the optimization method.

3.1 The Energy Integrating Detector Model

The computed tomography (CT) process mainly consists of three parts: x-ray beams are emitted from a source with specific energies, an object is illuminated by x-ray beams and attenuated x-ray beams are received by a detector. In this process, the intensities of x-ray beams are reduced and using Beer's law [18], the energy integrating detector model can be written as

$$y_i = \int_E S(e) \exp \left(- \int_{t \in \ell} \mu(\vec{r}(t), e) dt \right) de + \eta_i, \quad i = 1, 2, \dots, N_d \times N_p, \quad (3.2)$$

where

- y_i is the x-ray intensity of the i -th pixel in the detector.
- E is the photon flux density. Figure 3.1 shows a curve of E versus photon energy with relative low potential (26 keV).
- N_d is the number of detector pixels. For a material map of size n by n , we assume $N_d = n$ and the number of projection rays for each angle is equivalent to $\lfloor \sqrt{2}N_d \rfloor$.
- N_p is the number of projections. For cone/fan beam CT, projections are distributed equally from 0 to 360 degrees.
- $S(e)$ represents the system spectral response, which is a product of x-ray energy with the number of incident photons at that energy.
- The outer integral is over all x-ray energies emitted from the source, and the inner integral is along lines that follow the x-ray beam paths through the object.
- $\mu(\vec{r}(t), e)$ denotes the attenuation coefficient, which is related to the position function $\vec{r}(t)$ and the energy level e .

- η_i represents unknown errors in the measurements, which can include x-ray scatter and electronic noise.

The traditional methods to solve this inverse problem are mostly based on filtered back-projection (FBP) [49]. If we assume the source to be monoenergetic and the source energy is s_e , then we can build a linear inverse problem by dividing s_e on both sides and by applying the natural logarithm function to the data and to the model. Other approaches can be used to solve this linear inverse problem, such as incorporating different regularization schemes. These usually involve applying appropriate iterative optimization methods, such as the conjugate gradient method [45]. However, the images obtained from traditional methods on the simplified linear model might lead to significant beam hardening artifacts when the object is made up of several very distinct materials, such as bone and soft tissue. In addition, the linear models cannot be used to recognize the the actual types of materials from the results, nor can they separate different materials when they are mixed [33]. Moreover, the traditional methods are unstable when it comes to the limited angle cases. For these reasons, we consider the full nonlinear polyenergetic model.

In Model (3.2), the unknown linear attenuation coefficient $\mu(\vec{r}(t), e)$ is dependent on the position function $r(t)$ and energy levels e . If the object is assumed to be composed of several different materials, then a material expansion is introduced to further decompose the function $\mu(\vec{r}(t), e)$ [33]:

$$\mu(\vec{r}(t), e) = \sum_{m=1}^{N_m} u_{m,e} w_m(\vec{r}), \quad (3.3)$$

where

- N_m is the number of materials that form the object.
- $u_{m,e}$ is the linear attenuation coefficient for the m -th material at the energy level e .

- $w_m(\vec{r})$ is the unknown weight of the m -th material at the position \vec{r} .

With this decomposition, the unknown variable has been shifted from $\mu(\vec{r}(t), e)$ to weight fraction $w_m(\vec{r})$. If we further assume that $w_m(\vec{r})$ can be represented as a sum of product of the weight $w_{j,m}$ and the basis function $\phi_j(\vec{r})$, then another expansion can be expressed as

$$w_m(\vec{r}) = \sum_{j=1}^{N_v} w_{j,m} \phi_j(\vec{r}), \quad (3.4)$$

where

- N_v is the number of voxels (pixels if 2D) of images that compose the object.
- $w_{j,m}$ is the weight fraction of the m -th material in the j -th voxel (pixels if 2D).
- $\phi_j(\vec{r})$ is the basis function of image representation. The line integral of the basis function, $a_{i,j}$, is the length of the x-ray beam through the j -th voxel (pixel if 2D), incident onto the i -th element of the product of the detector pixels N_d and the number of projections N_p :

$$a_{i,j} = \int_{t \in l} \phi_j(\vec{r}(t)) dt. \quad (3.5)$$

By assumption, the attenuation coefficient for the m -th material and the energy level e , $\mu_{m,e}$, is already known and the only unknown variable is the weight $w_{j,m}$. Using this expression, we have transformed the goal of solving $\mu_{j,e}$ to the target of solving for $w_{j,m}$, which is dependent on the number of voxels and the number of materials. Usually, the number of materials is 2 or 3 and it is significantly fewer than the number of energies. To simplify the problem, we also assume that the sum of weights inside each voxel is equivalent to 1. That is to say,

$$\sum_{m=1}^{N_m} w_{j,m} = 1. \quad (3.6)$$

In this chapter, we limit the discussion to 2 to 3 materials, but note that for energy discriminating detectors, separating more materials is feasible. With (3.3), (3.4) and (3.5), the line integral in Model (3.2) is expressed by

$$\int_{t \in l} \mu(\vec{r}(t), e) dt = \sum_{m=1}^{N_m} \sum_{j=1}^{N_v} u_{m,e} w_{j,m} \int_{t \in l} \phi_j(\vec{r}(t)) dt = \sum_{j=1}^{N_v} \sum_{m=1}^{N_m} a_{i,j} w_{j,m} u_{m,e}. \quad (3.7)$$

If we also discretize the integral over energy E and ignore quadrature errors, then the discrete model of Equation (3.2) can be written as:

$$y_i = \sum_{e=1}^{N_e} s_e \exp\left(-\sum_{j=1}^{N_v} \sum_{m=1}^{N_m} a_{i,j} w_{j,m} u_{m,e}\right) + \eta_i, \quad (3.8)$$

where N_e is the number of discrete energies. If we collect $a_{i,j}$, $w_{i,j}$ and $u_{m,e}$ in matrix form and concatenate y_i , s_e , η_i into vectors, then the corresponding equation can be represented as

$$\mathbf{y} = \exp(-\mathbf{A}\mathbf{W}\mathbf{C}^T) \mathbf{s} + \boldsymbol{\eta}, \quad (3.9)$$

where

- \mathbf{y} is a vector of the size $N_d \cdot N_p$ that gathers x-ray photons.
- \mathbf{A} is a matrix of the size $(N_d \cdot N_p) \times N_v$ that collects the fan-beam geometry and each element corresponds to $a_{i,j}$.
- \mathbf{C} is a matrix of the size $N_e \times N_m$ that accumulates linear attenuation coefficients and each entry corresponds to $u_{e,m}$, the linear attenuation coefficient of the e -th energy and the m -th material.
- \mathbf{s} is a vector of the size N_e and collects the spectrum energy of a specific range.
- $\boldsymbol{\eta}$ is the noise vector that is of the size $N_d \cdot N_p$.

In Equation (3.9), the exponential function is point-wise exponential rather than

matrix function. In addition to Equation (3.9), we also require that weight fractions should be nonnegative and this can be illustrated by the constraint $\mathbf{W} \geq \mathbf{0}$. We have so far obtained the standard form of the polyenergetic multi-material model. Based on Equation (3.9), the goal is to solve for the unknown weight matrix \mathbf{W} such that $\mathbf{W} \in \mathbf{P}$, where $\mathbf{P} = \{\mathbf{W} \mid \mathbf{W}\mathbf{1}_{N_m} = \mathbf{1}_{N_v}, \mathbf{0} \leq \mathbf{W} \leq \mathbf{1}\}$. $\mathbf{1}_{N_m}$ and $\mathbf{1}_{N_v}$ are vectors of ones of the lengths N_m and N_v , respectively.

3.2 Poisson Log-likelihood Function

With different energy levels, the forward model is nonlinear and it is not possible to transform it into an equivalent linear model. Mejia-Bustamante et al. [7] use the assumption that each entry in \mathbf{y} follows a Poisson distribution,

$$\mathbf{y} \sim \text{Poisson}(\exp(-\mathbf{A}\mathbf{W}\mathbf{C}^T)\mathbf{s}). \quad (3.10)$$

With this assumption, one can formulate a maximum likelihood estimator (MLE) based on Poisson distribution. Before we show the expressions of objective function, gradient and Hessian, we first note that if the weights (unknowns) are stored as a matrix, then differentiation results in tensors, requiring tedious bookkeeping in the computations. In this work we instead rewrite the function to put the unknowns in vector form, and differentiate the objective function with respect to this vector. Notice that

$$\text{vec}(\mathbf{A}\mathbf{W}\mathbf{C}^T) = (\mathbf{C} \otimes \mathbf{A}) \text{vec}(\mathbf{W}), \quad (3.11)$$

where $\text{vec}(\cdot)$ reshapes a given matrix into a vector by stacking the columns on top of each other. Therefore, we can rewrite Equation (3.9) as

$$\mathbf{y} = (\mathbf{s}^T \otimes \mathbf{I}) \exp[-(\mathbf{C} \otimes \mathbf{A}) \text{vec}(\mathbf{W})] + \boldsymbol{\eta}. \quad (3.12)$$

If we let

$$\mathbf{w} = \text{vec}(\mathbf{W}) \quad \text{and} \quad K(\mathbf{w}) = \exp\{- (\mathbf{C} \otimes \mathbf{A}) \mathbf{w}\}, \quad (3.13)$$

then Equation (3.12) is equivalent to

$$\mathbf{y} = (\mathbf{s}^T \otimes \mathbf{I})K(\mathbf{w}) + \boldsymbol{\eta}. \quad (3.14)$$

In this problem, we assume that each element of measured data, \mathbf{y}_i , follows a Poisson distribution with mean $[(\mathbf{s}^T \otimes \mathbf{I})K(\mathbf{w})]_i$. That is to say,

$$y_i \sim \text{Poisson}([(\mathbf{s}^T \otimes \mathbf{I})K(\mathbf{w})]_i). \quad (3.15)$$

Based on this assumption, the corresponding probability density function can be expressed as

$$f_{\mathbf{y}}(\mathbf{y}; \mathbf{w}) = \prod_{i=1}^{N_d \times N_p} \frac{[(\mathbf{s}^T \otimes \mathbf{I})K(\mathbf{w})]_i^{y_i} \exp\{-[(\mathbf{s}^T \otimes \mathbf{I})K(\mathbf{w})]_i\}}{y_i!}. \quad (3.16)$$

If we ignore the constant terms, the corresponding likelihood function is

$$L(\mathbf{w}; \mathbf{y}) = \prod_{i=1}^{N_d \times N_p} [(\mathbf{s}^T \otimes \mathbf{I})K(\mathbf{w})]_i^{y_i} \exp\{-[(\mathbf{s}^T \otimes \mathbf{I})K(\mathbf{w})]_i\}. \quad (3.17)$$

Taking the logarithm, the log-likelihood function is expressed by

$$l(\mathbf{w}; \mathbf{y}) = \sum_{i=1}^{N_d \times N_p} \{y_i \log[(\mathbf{s}^T \otimes \mathbf{I})K(\mathbf{w})]_i - [(\mathbf{s}^T \otimes \mathbf{I})K(\mathbf{w})]_i\}. \quad (3.18)$$

To maximize the log-likelihood function, it is equivalent to minimizing the negative log-likelihood function. So the objective function is

$$\begin{aligned} l(\mathbf{w}; \mathbf{y}) &= \sum_{i=1}^{N_d \times N_p} \{[(\mathbf{s}^T \otimes \mathbf{I}) K(\mathbf{w})]_i - y_i \log[(\mathbf{s}^T \otimes \mathbf{I}) K(\mathbf{w})]_i\}. \\ &= \mathbf{1}_{N_p \times N_\theta}^T (\mathbf{s}^T \otimes \mathbf{I}) K(\mathbf{w}) - \mathbf{y}^T \log [(\mathbf{s}^T \otimes \mathbf{I}) K(\mathbf{w})], \end{aligned} \quad (3.19)$$

where $\mathbf{1}_{N_d \times N_p}$ is a vector of all ones of the length $N_d \times N_p$. The gradient to the objective function can be expressed as

$$\nabla l(\mathbf{w}) = -(\mathbf{C}^T \otimes \mathbf{A}^T) \text{diag}\{K(\mathbf{w})\} (\mathbf{s} \otimes \mathbf{I}) \{\mathbf{1}_{N_p \times N_\theta} - \mathbf{y} \odot [(\mathbf{s}^T \otimes \mathbf{I}) K(\mathbf{w})]\}. \quad (3.20)$$

Differentiating the gradient with respect to \mathbf{w} , the Hessian can be represented as a summation of two parts:

$$\nabla^2 l(\mathbf{w}) = H_1(\mathbf{w}) + H_2(\mathbf{w}), \quad (3.21)$$

where

$$\begin{aligned} H_1(\mathbf{w}) &= (\mathbf{C}^T \otimes \mathbf{A}^T) \text{diag}\{K(\mathbf{w})\} \text{diag}\{(\mathbf{s} \otimes \mathbf{I}) (\mathbf{1}_{N_p \times N_\theta} - \mathbf{y} \odot [(\mathbf{s}^T \otimes \mathbf{I}) K(\mathbf{w})])\} (\mathbf{C} \otimes \mathbf{A}) \\ H_2(\mathbf{w}) &= (\mathbf{C}^T \otimes \mathbf{A}^T) \text{diag}\{K(\mathbf{w})\} (\mathbf{s} \otimes \mathbf{I}) \text{diag}\{\mathbf{y} \odot [(\mathbf{s}^T \otimes \mathbf{I}) K(\mathbf{w})].^2\} \\ &\quad (\mathbf{s}^T \otimes \mathbf{I}) \text{diag}\{K(\mathbf{w})\} (\mathbf{C} \otimes \mathbf{A}). \end{aligned} \quad (3.22)$$

From the previous expressions, we can see that $H_2(\mathbf{w})$ is the Gauss-Newton approximation to the true Hessian and it is always positive semi-definite. Even if $H_1(\mathbf{w})$ is indefinite, we can transform it into a positive semi-definite matrix by setting a

threshold. Specifically, let $T(\mathbf{w})$ be defined as

$$T(\mathbf{w}) = \max [\mathbf{0}, (\mathbf{s} \otimes \mathbf{I}) (\mathbf{1}_{N_p \times N_\theta} - \mathbf{y} \oslash [(\mathbf{s}^T \otimes \mathbf{I})K(\mathbf{w})])] . \quad (3.23)$$

Then the new $\hat{H}_1(\mathbf{w})$ can be represented as

$$\hat{H}_1(\mathbf{w}) = (\mathbf{C}^T \otimes \mathbf{A}^T) \text{diag} \{K(\mathbf{w})\} \text{diag} \{T(\mathbf{w})\} (\mathbf{C} \otimes \mathbf{A}) . \quad (3.24)$$

Moreover, we define the modified Hessian as

$$\hat{H}(\mathbf{w}) = \hat{H}_1(\mathbf{w}) + H_2(\mathbf{w}) . \quad (3.25)$$

With the modified Hessian, we can include most information about the true Hessian as well as keep it positive semi-definite. Furthermore, both the bound constraint and the constraint (3.6) should be included. The bound constraint is equivalent to $\mathbf{0} \leq \mathbf{w} \leq \mathbf{1}$ and we can rewrite the constraint (3.6) in a matrix-vector form as

$$\mathbf{A}_{eq} \mathbf{w} = \mathbf{1}_{N_v}, \quad (3.26)$$

where \mathbf{A}_{eq} is a matrix of the form

$$\mathbf{A}_{eq} = \mathbf{1}_{N_m}^T \otimes \mathbf{I}_{N_v}. \quad (3.27)$$

$\mathbf{1}_{N_m}$ is a vector of ones of the length N_m . \mathbf{I}_{N_v} is an identity matrix of the size $N_v \times N_v$. In the following sections, we use \mathbf{I} to represent the identity matrix if the size of this matrix is clear to identify.

With the objective function (3.19), we can construct an optimization problem by

combining the regularization term and the constraints:

$$\begin{aligned} \min_{\mathbf{0} \leq \mathbf{w} \leq \mathbf{1}} \quad & f(\mathbf{w}) + \alpha R(\mathbf{w}) \\ \text{subject to} \quad & \mathbf{A}_{eq} \mathbf{w} = \mathbf{1}_{N_v}. \end{aligned} \quad (3.28)$$

In Problem (3.28), $f(\mathbf{w}) = l(\mathbf{w}; \mathbf{y})$ and $R(\mathbf{w})$ represents the regularization term, which is used to penalize the variable \mathbf{w} . α is the corresponding regularization parameter. For this problem, we choose the total variation regularization to preserve edges. Using the previous notations, we can express the bound constraint and the equality constraint as $\mathbf{Q} = \{\mathbf{w} \mid \mathbf{A}_{eq} \mathbf{w} = \mathbf{1}_{N_v}, \mathbf{0} \leq \mathbf{w} \leq \mathbf{1}\}$. So the notation of Problem (3.28) can be simplified as

$$\min_{\mathbf{w} \in \mathbf{Q}} f(\mathbf{w}) + \alpha R(\mathbf{w}). \quad (3.29)$$

We let \mathbf{w}_k be the k -th column in the matrix \mathbf{W} and \mathbf{W}_k be the corresponding reshaped image of \mathbf{w}_k . The total variation regularization for the k -th material can be represented as [30]:

$$R(\mathbf{W}_k) = \sum_{i=1}^n \sum_{j=1}^n \left((\mathbf{W}_k \mathbf{D}^T)_{ij}^2 + (\mathbf{D} \mathbf{W}_k)_{ij}^2 \right)^{1/2}, \quad (3.30)$$

where \mathbf{D} is either a forward, backward or central first order finite difference matrix and n is the dimension of the corresponding material map. For all m materials, the regularization term can be expressed as

$$R(\mathbf{W}) = \sum_{k=1}^m R(\mathbf{W}_k). \quad (3.31)$$

Other forms of regularization, such as the discrete Laplacian, can also be used with this framework. For simplicity, we use zero boundary conditions and assume that the regularization parameters for different materials are equal to α . With the regulariza-

tion terms, we can reduce the influence of noise and stabilize the solution.

With the modified Hessian and regularizations, we want to use a second order method to solve the optimization problem (3.28). However, direct implementation of the Newton's method is not effective. To guarantee the feasibility of each step, we need to project the current step onto the boundaries. The projected solution might be far away from the desired solution and hard to improve later. Furthermore, it is difficult to maintain the equality constraint in each step without changing the original model.

3.3 Implementation of Nonlinear Interior Point Trust Region Method

Recall that the optimization problem can be expressed as

$$\begin{aligned} \min_{\mathbf{0} \leq \mathbf{w} \leq \mathbf{1}} \quad & f(\mathbf{w}) + \alpha R(\mathbf{w}) \\ \text{subject to} \quad & \mathbf{A}_{eq} \mathbf{w} = \mathbf{1}, \end{aligned} \tag{3.32}$$

where $f(\mathbf{w}) = l(\mathbf{w}; \mathbf{y})$ and $R(\mathbf{w})$ is the regularization term.

To solve this constrained optimization problem, we use a nonlinear interior point trust region method, which combines sequential quadratic programming (SQP), a trust region dogleg method, and a projected conjugate gradient algorithm [8, 48]. To apply this method, we firstly establish a barrier problem based on (3.32):

$$\begin{aligned} \min_{\mathbf{w}, \mathbf{z}} \quad & f(\mathbf{w}) + \alpha R(\mathbf{w}) - \beta \sum_{i=1}^{2N_v} \ln(z_i) \\ \text{subject to} \quad & \mathbf{A}_{eq} \mathbf{w} - \mathbf{1} = \mathbf{0}, \\ & \mathbf{A}_{ieq} \mathbf{w} + \mathbf{y}_{ieq} + \mathbf{z} = \mathbf{0}, \end{aligned} \tag{3.33}$$

where

$$\mathbf{A}_{ieq} = \begin{bmatrix} -\mathbf{I} \\ \mathbf{I} \end{bmatrix} \quad \text{and} \quad \mathbf{y}_{ieq} = \begin{bmatrix} -\mathbf{1} \\ \mathbf{0} \end{bmatrix}. \quad (3.34)$$

β is the barrier parameter that should decrease to 0 and \mathbf{z} is the vector of slack variables that ensure all entries remain positive. The *permuted* KKT condition corresponding to (3.33) can be written as

$$\begin{aligned} \nabla f(\mathbf{w}) + \alpha \nabla R(\mathbf{w}) + \mathbf{A}_{eq}^T \boldsymbol{\lambda}_{eq} + \mathbf{A}_{ieq}^T \boldsymbol{\lambda}_{ieq} &= \mathbf{0}, \\ \mathbf{Z} \boldsymbol{\lambda}_{ieq} - \beta \mathbf{1} &= \mathbf{0}, \\ \mathbf{A}_{eq} \mathbf{w} - \mathbf{1} &= \mathbf{0}, \\ \mathbf{A}_{ieq} \mathbf{w} + \mathbf{y}_{ieq} + \mathbf{z} &= \mathbf{0}, \end{aligned} \quad (3.35)$$

where $\boldsymbol{\lambda}_{eq}$ and $\boldsymbol{\lambda}_{ieq}$ are the Lagrange multipliers corresponding to the equality and the inequality constraints, respectively. Furthermore, we should keep $\boldsymbol{\lambda}_{ieq}$ nonnegative. \mathbf{Z} is a diagonal matrix and $\mathbf{Z} = \text{diag} \{z_1, z_2, \dots, z_{2N_e}\}$. Compared with the original KKT system, this permuted KKT system is preferred because the matrix \mathbf{Z} is bounded when the entries of \mathbf{z} approach 0. We can construct an error function based on this system:

$$\begin{aligned} E(\mathbf{w}, \mathbf{z}, \boldsymbol{\lambda}_{eq}, \boldsymbol{\lambda}_{ieq}; \beta) &= \max \{ \|f(\mathbf{w}) + \alpha \nabla R(\mathbf{w}) \\ &\quad + \mathbf{A}_{eq}^T \boldsymbol{\lambda}_{eq} + \mathbf{A}_{ieq}^T \boldsymbol{\lambda}_{ieq}\|_{\infty}, \\ &\quad \|\mathbf{Z} \boldsymbol{\lambda}_{ieq} - \beta \mathbf{1}\|_{\infty}, \|\mathbf{A}_{eq} \mathbf{w} - \mathbf{1}\|_{\infty}, \\ &\quad \|\mathbf{A}_{ieq} \mathbf{w} + \mathbf{y}_{ieq} + \mathbf{z}\|_{\infty} \}. \end{aligned} \quad (3.36)$$

If the error function is less than a tolerance, for example, 10^{-8} , then we assume that we have solved this system accurately. The Newton system corresponding to the

permuted KKT system can be written as

$$D(\mathbf{w})\mathbf{p} = -g(\mathbf{w}), \quad (3.37)$$

where

$$D(\mathbf{w}) = \begin{bmatrix} \hat{H}(\mathbf{w}) + \alpha \nabla^2 R(\mathbf{w}) & \mathbf{0} & \mathbf{A}_{eq}^T & \mathbf{A}_{ieq}^T \\ \mathbf{0} & \mathbf{\Lambda}_{ieq} & \mathbf{0} & \mathbf{Z} \\ \mathbf{A}_{eq} & \mathbf{0} & \mathbf{0} & \mathbf{0} \\ \mathbf{A}_{ieq} & \mathbf{I} & \mathbf{0} & \mathbf{0} \end{bmatrix},$$

$$g(\mathbf{w}) = \begin{bmatrix} \nabla f(\mathbf{w}) + \alpha \nabla R(\mathbf{w}) + \mathbf{A}_{eq}^T \boldsymbol{\lambda}_{eq} + \mathbf{A}_{ieq}^T \boldsymbol{\lambda}_{ieq} \\ \mathbf{Z} \boldsymbol{\lambda}_{ieq} - \beta \mathbf{1} \\ \mathbf{A}_{eq} \mathbf{w} - \mathbf{1} \\ \mathbf{A}_{ieq} \mathbf{w} + \mathbf{y}_{ieq} + \mathbf{z} \end{bmatrix}, \quad (3.38)$$

$$\mathbf{p} = \begin{bmatrix} \mathbf{p}_x \\ \mathbf{p}_z \\ \Delta \boldsymbol{\lambda}_{eq} \\ \Delta \boldsymbol{\lambda}_{ieq} \end{bmatrix}.$$

In the matrix $D(\mathbf{w})$, $\mathbf{\Lambda}_{ieq}$ is a diagonal matrix with $\boldsymbol{\lambda}_{ieq}$ in the diagonal: $\mathbf{\Lambda}_{ieq} = \text{diag}\{\boldsymbol{\lambda}_{ieq}\}$. However, both direct and iterative methods to solve this system are computationally expensive, especially for large-scale problems. We therefore follow the strategy of Byrd et al. [9] by transforming the original problem into a sequential quadratic programming problem and solve by separating it into two subproblems. The first subproblem is called the normal subproblem, which is solved by trust-region dogleg method, while the second subproblem, the tangential subproblem, can be solved by the projected conjugate gradient method. By applying the idea of sequential

quadratic programming, we can construct an optimization problem as

$$\begin{aligned}
\min_{\mathbf{p}_x, \mathbf{p}_z} \quad & \nabla [f(\mathbf{w}) + \alpha R(\mathbf{w})]^T \mathbf{p}_x + \frac{1}{2} \mathbf{p}_x^T \left[\hat{H}(\mathbf{w}) + \alpha \nabla^2 R(\mathbf{w}) \right] \mathbf{p}_x \\
& - \beta \mathbf{1}^T \mathbf{Z}^{-1} \mathbf{p}_z + \frac{1}{2} \mathbf{p}_z^T \Sigma \mathbf{p}_z \\
\text{subject to} \quad & \mathbf{A}_{eq} \mathbf{p}_x + \mathbf{A}_{eq} \mathbf{w} - \mathbf{1} = \mathbf{r}_{eq}, \\
& \mathbf{A}_{ieq} \mathbf{p}_x + \mathbf{p}_z + \mathbf{A}_{ieq} \mathbf{w} + \mathbf{y}_{ieq} + \mathbf{z} = \mathbf{r}_{ieq}, \\
& \left\| \begin{bmatrix} \mathbf{p}_x^T & \mathbf{p}_z^T \mathbf{Z}^{-1} \end{bmatrix} \right\|_2 \leq \Delta, \\
& \mathbf{p}_z \geq -\tau \mathbf{z},
\end{aligned} \tag{3.39}$$

where $\Sigma = \mathbf{Z}^{-1} \mathbf{\Lambda}_{ieq}$. For this primal-dual system, we regard $\Sigma = \mathbf{Z}^{-1} \mathbf{\Lambda}_{ieq}$ as an approximation to the second order derivative $\beta \mathbf{Z}^{-2}$. \mathbf{r}_{eq} and \mathbf{r}_{ieq} are auxiliary variables to the linearized constraints. Moreover, the trust region constraint is used to guarantee a sufficient reduction in each step and stop the loop after a fixed number of iterations. The scale term \mathbf{Z}^{-1} in the trust region constraint is used to keep slack variables away from zero without enough iterations. Meanwhile, the bound constraint is applied to guarantee the positivity of the slack variable \mathbf{z} . That is to say, if the step is accepted, then we should have

$$\mathbf{z} + \mathbf{p}_z \geq \mathbf{z} - \tau \mathbf{z} = (1 - \tau) \mathbf{z} > 0. \tag{3.40}$$

τ is a constant that is close to 1, for example, $\tau = 0.98$. If we let $\tilde{\mathbf{p}}_z = \mathbf{Z}^{-1} \mathbf{p}_z$, then

we can rewrite (3.39) as

$$\begin{aligned}
& \min_{\mathbf{p}_x, \tilde{\mathbf{p}}_z} \quad \nabla [f(\mathbf{w}) + \alpha R(\mathbf{w})]^T \mathbf{p}_x + \frac{1}{2} \mathbf{p}_x^T \left[\hat{H}(\mathbf{w}) + \alpha \nabla^2 R(\mathbf{w}) \right] \mathbf{p}_x \\
& \quad - \beta \mathbf{1}^T \tilde{\mathbf{p}}_z + \frac{1}{2} \tilde{\mathbf{p}}_z^T \mathbf{Z} \Sigma \mathbf{Z} \tilde{\mathbf{p}}_z \\
& \text{subject to} \quad \mathbf{A}_{eq} \mathbf{p}_x + \mathbf{A}_{eq} \mathbf{w} - \mathbf{1} = \mathbf{r}_{eq}, \\
& \quad \mathbf{A}_{ieq} \mathbf{p}_x + \mathbf{Z} \tilde{\mathbf{p}}_z + \mathbf{A}_{ieq} \mathbf{w} + \mathbf{y}_{ieq} + \mathbf{z} = \mathbf{r}_{ieq}, \\
& \quad \left\| \begin{bmatrix} \mathbf{p}_x^T & \tilde{\mathbf{p}}_z^T \end{bmatrix} \right\|_2 \leq \Delta, \\
& \quad \tilde{\mathbf{p}}_z \geq -\tau.
\end{aligned} \tag{3.41}$$

With the help of auxiliary variables, we can separate (3.41) into two subproblems, the normal subproblem and the tangential subproblem. The normal subproblem can be expressed as

$$\begin{aligned}
& \min_{\mathbf{v}_x, \mathbf{v}_z} \quad \left\| \mathbf{A}_{eq} \mathbf{v}_x + \mathbf{A}_{eq} \mathbf{w} - \mathbf{1} \right\|_2^2 + \left\| \mathbf{A}_{ieq} \mathbf{v}_x + \mathbf{Z} \mathbf{v}_z + \mathbf{A}_{ieq} \mathbf{w} + \mathbf{y}_{ieq} + \mathbf{z} \right\|_2^2 \\
& \text{subject to} \quad \left\| \begin{bmatrix} \mathbf{v}_x^T & \mathbf{v}_z^T \end{bmatrix} \right\|_2 \leq \zeta \Delta, \\
& \quad \mathbf{v}_z \geq -\frac{\tau}{2},
\end{aligned} \tag{3.42}$$

where ζ is a constant and $0 < \zeta < 1$, for example, $\zeta = 0.8$. Without the bound constraint, the normal subproblem is a standard form of trust region problem, which can be solved by the trust region dogleg method. So we solve it first by ignoring the bound constraint and test if the solution satisfies this constraint later. If not, we conduct backtracking to maintain feasibility [54]. After solving the normal subproblem approximately, we obtain the residuals as

$$\begin{aligned}
\mathbf{r}_{eq} &= \mathbf{A}_{eq} \mathbf{v}_x + \mathbf{A}_{eq} \mathbf{w} - \mathbf{1}, \\
\mathbf{r}_{ieq} &= \mathbf{A}_{ieq} \mathbf{v}_x + \mathbf{Z} \mathbf{v}_z + \mathbf{A}_{ieq} \mathbf{w} + \mathbf{y}_{ieq} + \mathbf{z}.
\end{aligned} \tag{3.43}$$

By substituting the residuals for the same terms in (3.41), the original optimization

problem can be represented as

$$\begin{aligned}
\min_{\mathbf{p}_x, \tilde{\mathbf{p}}_z} \quad & \nabla [f(\mathbf{w}) + \alpha R(\mathbf{w})]^T \mathbf{p}_x + \frac{1}{2} \mathbf{p}_x^T \left[\hat{H}(\mathbf{w}) + \alpha \nabla^2 R(\mathbf{w}) \right] \mathbf{p}_x \\
& - \beta \mathbf{1}^T \tilde{\mathbf{p}}_z + \frac{1}{2} \tilde{\mathbf{p}}_z^T \mathbf{Z} \Sigma \mathbf{Z} \tilde{\mathbf{p}}_z \\
\text{subject to} \quad & \mathbf{A}_{eq}(\mathbf{p}_x - \mathbf{v}_x) = 0, \\
& \mathbf{A}_{ieq}(\mathbf{p}_x - \mathbf{v}_x) + \mathbf{Z}(\tilde{\mathbf{p}}_z - \mathbf{v}_z) = 0, \\
& \left\| \begin{bmatrix} \mathbf{p}_x^T & \tilde{\mathbf{p}}_z^T \end{bmatrix} \right\|_2 \leq \Delta, \\
& \tilde{\mathbf{p}}_z \geq -\tau.
\end{aligned} \tag{3.44}$$

If we ignore the last two constraints, this optimization problem is a standard form of quadratic programming problem under linear equality constraints, which can be solved by the projected conjugate gradient method. In this case, we ignore the bound constraint at first and stop the iteration when the desired tolerance is attained or the current step crosses the trust region boundary. If the solution does not satisfy the bound constraint, we backtrack and choose the last feasible step as the solution.

After we have obtained \mathbf{p}_x and \mathbf{p}_z , we need to decide if we should accept them and update the current step as well as the size of the trust region. To realize this idea, we can construct a merit function based on the objective function and constraints from the original barrier problem to decide the actual reduction. For example, a merit function can be expressed as

$$\phi_\nu(\mathbf{w}, \mathbf{z}) = f(\mathbf{w}) + \alpha R(\mathbf{w}) - \beta \sum_{i=1}^{2N_v} \ln(z_i) + \nu \|\mathbf{A}_{eq} \mathbf{w} - \mathbf{1}\|_2 + \nu \|\mathbf{A}_{ieq} \mathbf{w} + \mathbf{y}_{ieq} + \mathbf{z}\|_2, \tag{3.45}$$

where $\nu > 0$ is a penalty parameter. The actual reduction can be represented as

$$\text{ared}(\mathbf{p}) = \phi_\nu(\mathbf{w}, \mathbf{z}) - \phi_\nu(\mathbf{w} + \mathbf{p}_x, \mathbf{z} + \mathbf{p}_z). \tag{3.46}$$

The predicted reduction can be constructed in many ways, most of which are based on the SQP problem and its constraints. For example, we can set up a function as

$$\begin{aligned}
q_\nu(\mathbf{p}) = & \nabla [f(\mathbf{w}) + \alpha R(\mathbf{w})]^T \mathbf{p}_x + \frac{1}{2} \mathbf{p}_x^T \left[\hat{H}(\mathbf{w}) + \alpha R(\mathbf{w}) \right] \mathbf{p}_x - \beta \mathbf{1}^T \mathbf{Z}^{-1} \mathbf{p}_z \\
& + \frac{1}{2} \mathbf{p}_z^T \Sigma \mathbf{p}_z + \left\| \begin{bmatrix} \mathbf{A}_{eq} \mathbf{p}_x + \mathbf{A}_{eq} \mathbf{w} - \mathbf{1} \\ \mathbf{A}_{ieq} \mathbf{p}_x + \mathbf{p}_z + \mathbf{A}_{ieq} \mathbf{w} + \mathbf{y}_{ieq} + \mathbf{z} \end{bmatrix} \right\|_2.
\end{aligned} \tag{3.47}$$

For $q_\nu(\mathbf{p})$, the predicted reduction is the difference between not taking any step and taking the obtained step \mathbf{p} , which can be indicated as

$$\text{pred}(\mathbf{p}) = q_\nu(\mathbf{0}) - q_\nu(\mathbf{p}), \tag{3.48}$$

where the variable \mathbf{p} is a concatenation of \mathbf{p}_x and \mathbf{p}_z . For a tiny constant $\eta = 10^{-8}$, if $\text{ared}(\mathbf{p}) \geq \eta \text{pred}(\mathbf{p})$, we accept \mathbf{p} and update the current step. We will also update the trust region with a standard criterion based on the ratio $\text{ared}(\mathbf{p}) / \text{pred}(\mathbf{p})$.

In conclusion, we implement this problem-specific nonlinear interior point trust region method to solve the corresponding optimization problem (3.33). Problem (3.33) is a nonlinear optimization problem under linear and bound constraints. The objective function contains a nonlinear log-likelihood term and a regularization term. For the log-likelihood term, we calculate the gradient and the modified Hessian as (3.20) and (3.25). For the regularization term, the total variation regularization is chosen to stabilize the solution. The modified Hessian is close to the true Hessian and it is positive semidefinite so solutions to the augmented Newton system are robust. Furthermore, the problem is prone to large-scale application since it is unnecessary to save the modified Hessian. We only need the Hessian-vector multiplication when we implement the Newton-CG method to solve the augmented Newton system. The cost of memory in each conjugate gradient iteration is close to an iteration of gradient descent.

3.4 Numerical Experiments

To test the method, we generate a 2D image of the size 128 by 128 and assume that the object is made up of three simulated materials that arise in polyenergetic image reconstruction – adipose, air and bones. For bones, we use the main component, calcium, to represent it. One application of polyenergetic image reconstruction is breast imaging, which requires low dose radiation for patients. To realize this application, we generate an energy spectrum with potential 26 keV with the help of function “spektrSpectrum” [53]. We also select a low radiation dose of $1e5$ total photons for the x-ray energy spectrum. The corresponding spectrum is shown in Figure 3.1. From Figure 3.1, we can find that the photon flux density is above zero when the energy is between 3 keV and 28 keV. Based on this observation, the discrete energies for the simulated source x-ray beam are chosen from 3 keV to 28 keV, with an interval of 1 keV.

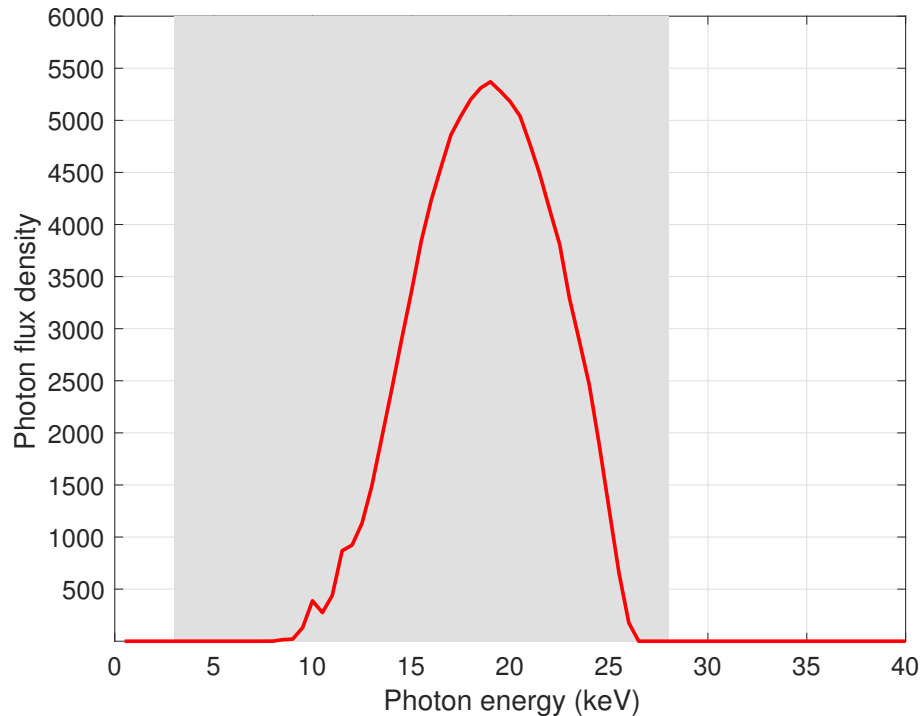


Figure 3.1: Photon flux density versus photon energy.

The plots of linear attenuation coefficients to materials adipose, air and calcium are shown in Figure 3.2. In Figure 3.2, the red, blue and black curves represent

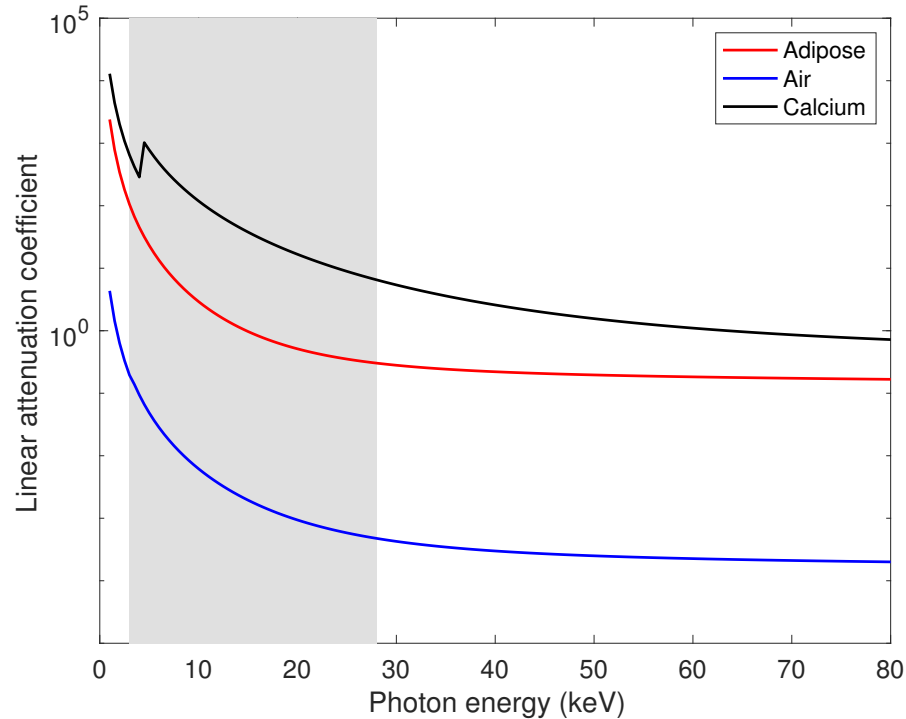


Figure 3.2: The linear attenuation curves for adipose (red), air (blue) and calcium (black).

adipose tissue, air and calcium, respectively, and the gray patch corresponds to the area of energy flux that is not equivalent to zero. From Figure 3.2, we can see that the curvatures of air and adipose are similar, while the curve of calcium has a K-edge [52]. The similarity of curvatures between adipose and air might cause the collinearity of linear attenuation coefficient matrix \mathbf{C} and so as the ill-conditioning of Hessian, while the K-edge might result in difficulty for reconstruction.

The simulations of the true object, shown in Figure 3.3, contain four distinct regions: 100% adipose, 0% air, 0% calcium; 0% adipose, 100% air, 0% calcium; 0% adipose, 0% air, 100% calcium; 50% adipose, 50% air, 0% calcium¹. In Figure 3.3, the

¹We actually tested many different combinations of mixed materials, for example, 20% adipose, 60% air and 20% calcium. The results are very similar to the one case considered in this experiment, thus to conserve space, we omit the results.

yellow color represents regions that contain 100% of the corresponding material, the turquoise color indicates regions that contain 50% of the adipose and air materials, the blue color indicates that the corresponding material does not exist in this area. Since we only have three materials, the weights corresponding to these three plots in Figure 3.3 should add to one. Moreover, since we use a 2-dimensional object for the

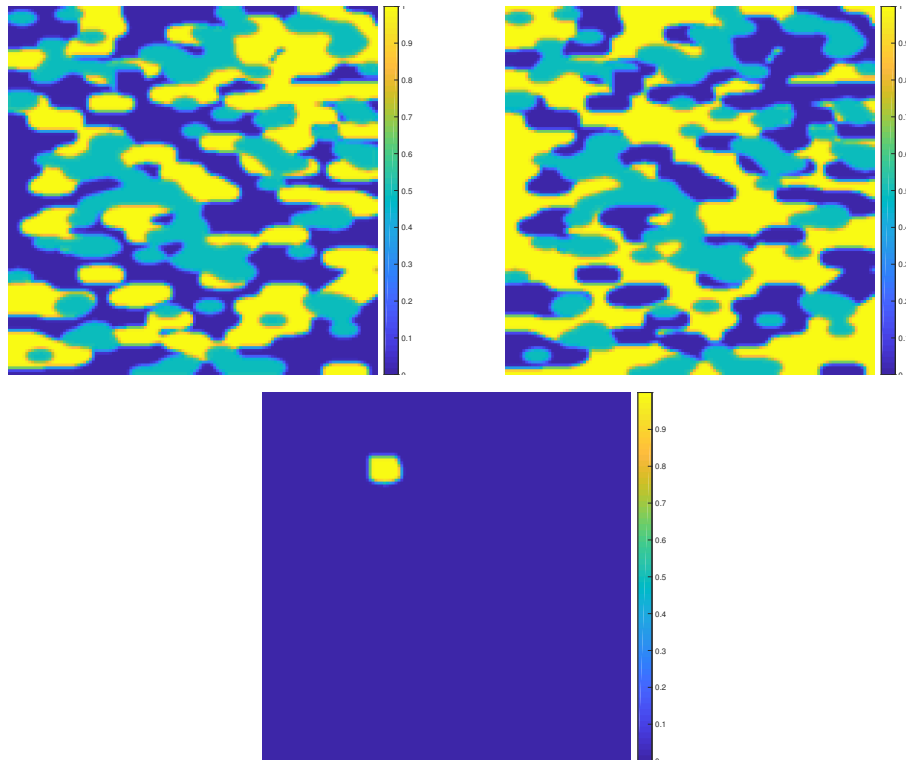


Figure 3.3: The true images for air (upper left), adipose (upper right) and calcium (middle). The turquoise colored regions are those areas in the object with a mixture of 50% glandular and 50% adipose tissue.

simulation, we use fan-beam (instead of cone beam) tomography model to generate a projection matrix \mathbf{A} using the AIR Tools software [32]. The distance between the source and the detector is 70 cm, with 2.5 cm air gap between the object and the detector. In order to keep the gauge of projection matrix the same under different size of images, we scale the projection matrix by the grid size and the dimension of images. For example, we choose the grid size as 2 cm and the dimension as 128 pixels so the scaling results in a pixel size of $2/128$ cm/pixel.

To avoid the inverse crime, we use spectral energies discretized on a finer grid and images with higher resolution to build the forward problem, but then use a coarser grid and lower resolution when solving the inverse problem. In particular, we collect the photon flux density corresponding to energies from 3 keV to 28 keV, with an interval of 0.5 keV for the forward problem, but then use an interval of 1 keV when we solve the inverse problem. Moreover, the resolution of object is initially 256×256 when we build the forward problem, and for the inverse problem, we solve (using a function included in the package IR Tools [22]) on a 128×128 grid. The number of x-rays used in building the ray trace matrix \mathbf{A} is scaled to match the projected data generated with higher resolution images. Both full CT and limited angle reconstructions are presented in the following sections.

3.4.1 Full Angle Reconstruction

First, we only consider the full CT case, where the range of projection angle is from 0 to 179 degrees in one degree increment. We use Poisson distribution to generate the measurements as (3.10). The initial guess is a random vector whose entries are between 0 and 1 and it is not required to satisfy the equality constraint. Moreover, a total variation regularization is introduced to preserve the edges; specifically, we use the forward difference operator and zero boundary conditions for this regularization term. The regularization parameter is chosen among the set $\{10^{-1}, 10^{-2}, 10^{-3}, 10^{-4}, 10^{-5}\}$ and the most effective parameter is used to compute the result. With our (un-optimized) MATLAB implementation on a laptop computer, we need around 20 minutes to finish 30 Newton iterations where the stationary point is achieved. For the full CT case, the reconstructed images are presented in Figure 3.4. From Figure 3.4, we can see that the reconstructed images are of high quality in general. It successfully separates the areas corresponding to adipose, air and calcium as well as mixture of adipose and air. Edges of the reconstructed images are clear, which might be con-

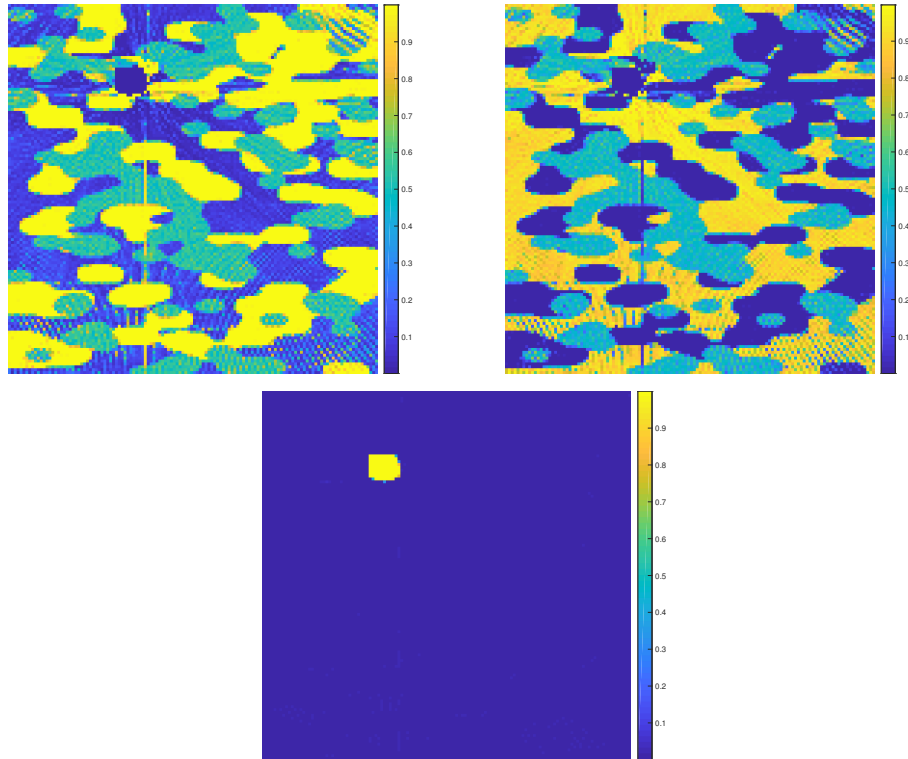


Figure 3.4: The reconstructed images for air (upper left), adipose (upper right) and calcium (middle) for the full CT reconstruction.

tributed from the total variation regularization. On the other hand, we can also find several artifacts that appear as blurred spots concentrating in the upper right corner, as well as other small artifacts scattered around the image. This results from measured data that are generated by Poisson distribution. With lower radiation dose, the relative noise level is higher compared with higher dose. Moreover, we can illustrate the convergence behavior by investigating the curves of relative errors. This plot is shown in Figure 3.5.

From Figure 3.5, we can find that the relative errors of materials air and adipose decrease in a similar way while the relative error of material calcium drops much faster. It is likely that calcium is only composed of a small part of the area and it can achieve faster convergence and higher accuracy. We also observe that the relative errors of three materials decrease to a particular level and then stagnate. The relative error between the last step and the true solution is about 19% for each material. Note

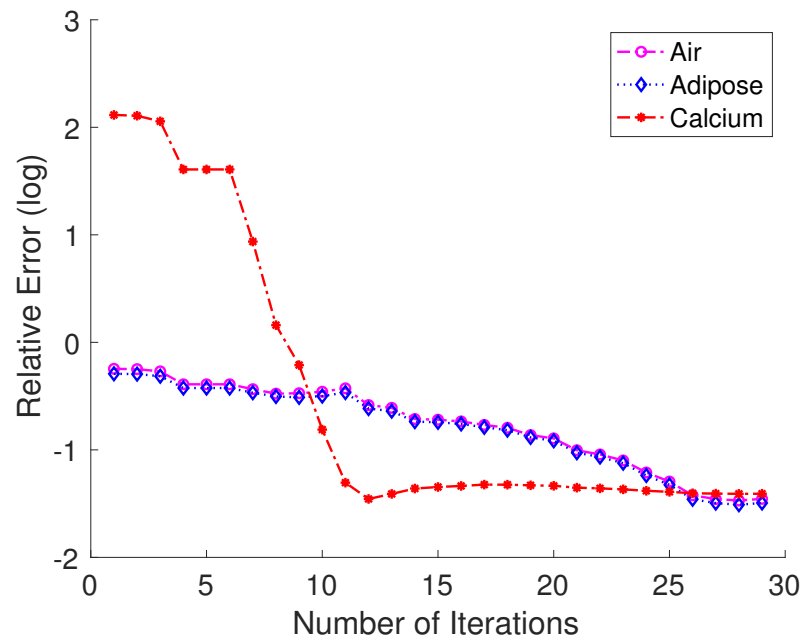


Figure 3.5: The plot of relative errors of air (magenta), adipose (blue) and calcium (red) for the full CT simulation.

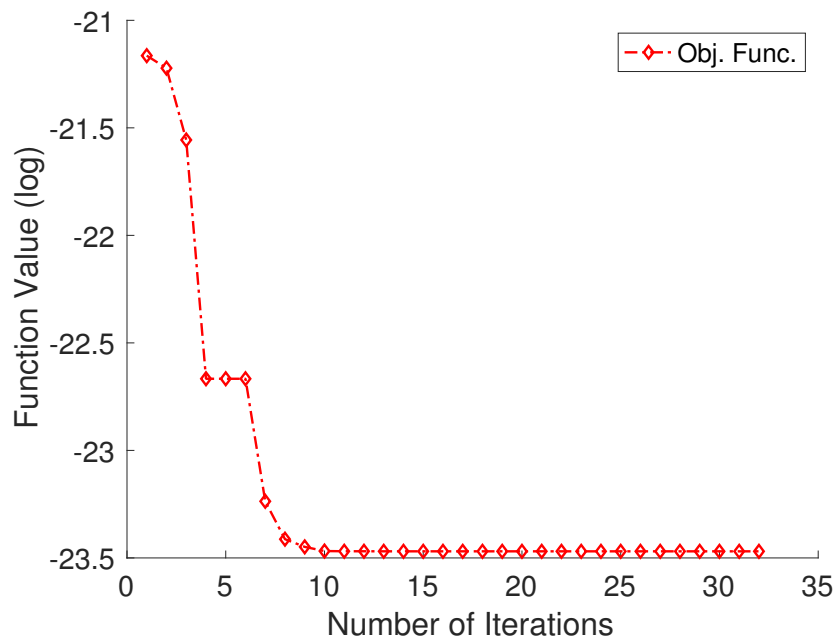


Figure 3.6: The plot of decrease of the objective function value.

that this stagnation occurs because of the regularization – without regularization, the relative errors may actually increase as the iterations proceed, which is a well-known behavior of ill-posed inverse problems, referred to as semi-convergence. We observe that when the relative errors stagnate, the current step approaches the first order optimality condition, which is an approximate solution to the KKT system.

To further validate convergence behavior of the proposed algorithm, we plot the curve of function value in Figure 3.6. From Figure 3.6, we can clearly identify that the function value drops fast in the beginning and then it stops for two iterations. After that, it starts to drop again and then stagnate. It cannot achieve lower value when reaching a specific level.

3.4.2 Limited Angle Reconstruction

In addition to the full CT case, it is important to also consider the case of limited angle reconstructions. Specifically, in the area of digital tomography, the limited angle reconstruction known as tomosynthesis has become an important diagnostic tool in breast imaging. The motivations for limited angle reconstruction are to reduce radiation dose to patients as well as to reduce the cost of this procedure. The limited angle reconstruction provides significantly more challenges to image reconstruction because the mathematical problems is much more ill-posed than the full CT case. This also means that the reconstruction quality is much more sensitive to the noise. In addition, the original objective function might have more stationary points and several of them are likely to satisfy the KKT condition. Under this situation, when the problem is effectively underdetermined, a poor initial guess may lead to an undesirable local minimum.

To test the limited angle reconstruction, the same test problem is used but with fewer projection angles. We shrink the projection angles from 180 degrees to 90 degrees, which ranges from 0 to 90 degrees. Furthermore, Poisson distribution is

used to generate the projection data. After implementing the previous algorithm, the reconstructed images are presented in Figure 3.7.

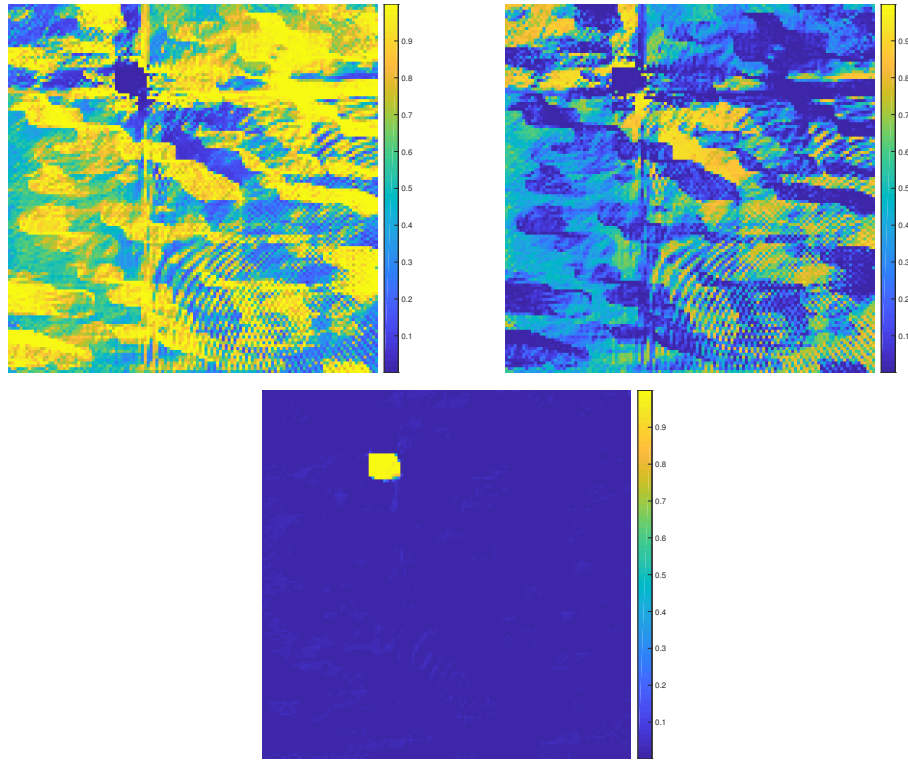


Figure 3.7: The reconstructed images for air (upper left), adipose (upper right) and calcium (middle) with 90 degrees projection.

From Figure 3.7, we can clearly see that the reconstructed images obtained from the limited angle case are more blurred than the images from the full CT case. As expected, with fewer projection angles, the images are of poorer quality. For the 90 degrees case, we can basically identify the distributions of materials roughly, while the details are more difficult to recognize. Only the material map of calcium is nearly fully separated from other materials. Moreover, we can see that the boundaries of different materials are not as clear as the boundaries in the full CT case. In several areas, the pixels are surrounded by shadows, which means that the materials are not completely separated. In the area of mixture, several pixels are colorful and the results depart slightly from the true solution. However, it is well known that due to the limited angle data, there are fundamental limitations when computing reconstructions [20].

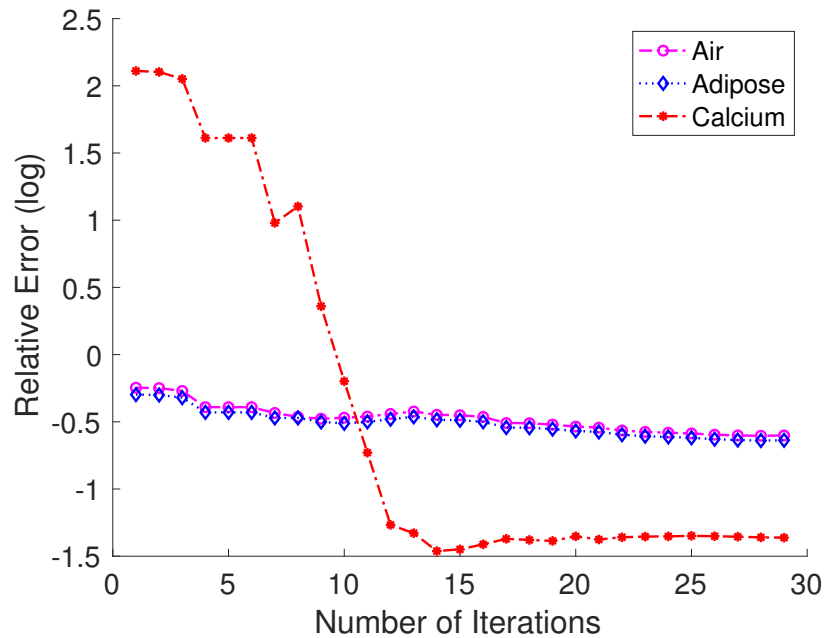


Figure 3.8: The plot of relative errors of air (magenta), adipose (blue) and calcium (red) for the 90 degrees limited angle simulation.

The plot of relative errors for 90 degrees case is presented in Figure 3.8. From Figure 3.8, we can find that the relative error curves corresponding to air and adipose decrease slowly and then stagnate. On the other hand, the relative error curve corresponding to calcium drops fast and converges to a lower level. This observation matches the phenomenon we conclude for the reconstructed images. Compared with the plot of full CT, the relative errors stagnate at higher levels. Moreover, we can find that the speed of convergence is not as rapid as the full CT case. However, we do observe that even if the regularization cannot completely compensate for the limited angle limitations, it does help to stabilize the solution.

3.5 Conclusions and Remarks

By taking multiple materials into consideration, the reconstructed images can reveal the weights of materials that compose the object, providing substantially more useful

information for the clinicians. Furthermore, the objective function and the gradient are uncomplicated to implement and the modified Hessian is a sufficient and stable estimate to the true Hessian. In addition, the merits of using a nonlinear interior point method are easy to identify. It is a globally convergent method with superlinear rate of convergence. It is also a stable and robust algorithm that can handle large-scale problems. Furthermore, there is substantial flexibility in choosing the initial guess because it does not need to satisfy the constraints.

Although this method has advantages such as faster convergence, robust computation and flexibility, it still has a few limitations. For example, implementation of nonlinear interior point method is not straightforward for large-scale problems. It requires solving a normal subproblem as well as a tangential subproblem. Furthermore, we need to decide the size of trust region in each iteration. Meanwhile, this method involves many parameters that we need to choose manually. So far, we have only tested 2D images rather than 3D images. For 3D images, the evaluation of each part might be more complicated, which is likely to increase the expense for solving this problem. For further research, we might consider the gradient-based methods such as the scaled gradient descent method or splitting methods such as the alternating direction method of multipliers (ADMM) [21].

Chapter 4

Nonlinear Optimization for Energy-windowed Spectral Computed Tomography

The development of new energy-windowed spectral computed tomography (CT) machines have received a great deal of interest in recent years; see, e.g. [2, 57]. These detectors assume that x-rays emitted by the x-ray source are composed of a spectrum of different energies, and in each energy window, the detector can detect a specific range of energy. Moreover, it assumes that the detector can perform photon counting and the data collected by the detector are nonnegative integers. Compared with traditional CT machines, we can avoid introducing beam-hardening artifacts [45] and improve quality of reconstructed images. To reconstruct the material maps of an object, we need to solve a nonlinear equation of the form

$$\mathbf{Y} = \exp(-\mathbf{A}\mathbf{W}\mathbf{C}^T) \mathbf{S} + \boldsymbol{\varepsilon}, \quad (4.1)$$

where \mathbf{Y} is a matrix that gathers the projected data of each energy window in the corresponding column and the exponential operator is applied element-wise (i.e., it is

not a matrix function). \mathbf{A} is a matrix that is related to the quantitative information of ray trace and \mathbf{C} is a matrix that contains linear attenuation coefficients for particular (known) materials at specified energies. \mathbf{S} is the matrix that accumulates the spectrum energies for each energy window in the corresponding column. Moreover, $\boldsymbol{\epsilon}$ represents the noise term and we assume that $y_{il} \sim \text{Poisson}([\exp(-\mathbf{A}\mathbf{W}\mathbf{C}^T)\mathbf{S}]_{il})$ for each component y_{il} in \mathbf{Y} and the corresponding entry in the exponential term. We also assume that these data are known and the target is to solve the unknown weight matrix \mathbf{W} . \mathbf{W} is of size N_v by N_m , where N_v is the number of voxels (pixels if 2D) for each material map and N_m is the number of materials. Since the weight matrix \mathbf{W} represents the material maps of different materials, it must be nonnegative and we need to add a lower bound constraint $\mathbf{W} \geq \mathbf{0}$.

To obtain the quantitative information about the material composition, we need to solve a nonlinear inverse problem. This nonlinear inverse problem is extremely ill-posed and direct implementation of regular methods is not effective. In this case, Barber et al. [2] suggests a preconditioner that is based on the eigenvalue decomposition of the matrix of linear attenuation coefficients, $\mathbf{C}^T\mathbf{C}$. With this preconditioner, they use the *Chambolle-Pock* (CP) primal-dual algorithm [10] to solve the inverse problem with Poisson log-likelihood loss function. However, they construct the preconditioner with only linear attenuation coefficients and do not include energy information. Moreover, the Poisson log-likelihood function contains multiple exponential terms and it is hard to handle. Because of these problems, we propose a new preconditioner that both involves linear attenuation coefficients and energy spectrum information based on a rank-1 approximation. To implement this preconditioner, we use a two-step method that includes finding an approximate Cauchy point in the first step and solving a quadratic programming problem in the second step.

This chapter is organized as follows. We review the energy-windowed spectral CT model in Section 4.1. Because of the challenges raised by this model, a new

preconditioning framework is introduced in Section 4.2. In order to implement this preconditioner, we suggest in Section 4.3 an optimization algorithm based on the projected line search and the trust region method. In Section 4.4 we illustrate the strength of the new preconditioner and the algorithm using numerical experiments. Finally, comments, limitations and future work are provided in Section 4.5.

4.1 The Energy-windowed Spectral CT Model

In computed tomography (CT), source x-ray beams are composed of a spectrum of different energies [7]. Recent technological developments have resulted in the design of new photon counting detectors that can discriminate the measured data into specific energy windows. Image reconstruction algorithms that exploit this information can avoid introducing beam-hardening artifacts, obtain material decomposition and improve the quality of reconstructed images. The mathematical model of energy-windowed spectral CT is expressed by

$$y_i^{(k)} = \int_E S^{(k)}(e) \exp\left(-\int_{t \in l} \mu(\vec{r}(t), e) dt\right) de + \eta_i^{(k)}, \quad \begin{cases} i = 1, 2, \dots, N_d \times N_p, \\ k = 1, 2, \dots, N_b, \end{cases} \quad (4.2)$$

where

- $y_i^{(k)}$ is the x-ray intensity of the i -th pixel in the k -th detector bin.
- E is the photon flux density. Figure 4.4 shows a curve of E versus photon energy with relative high potential (120 keV).
- N_d is the number of detector pixels. For a material map of the size n by n , we assume $N_d = n$ and the number of projection rays for each angle is also equivalent to N_d .

- N_p is the number of projections. For cone/fan beam CT, projections are distributed equally from 0 to 360 degrees.
- N_b is the number of detector bins (windows). For an energy-windowed CT machine, we usually assume that it has 5 to 6 energy bins.
- $S^{(k)}(e)$ represents the photon flux density for the k -th detector bin, which is the number of incident photons at the energy level e in the k -th energy window.
- $\mu(\vec{r}(t), e)$ denotes the linear attenuation coefficient that is related to the position function $\vec{r}(t)$ and energy level e .
- $\eta_i^{(k)}$ is the error term for the i -th element in k -th energy bin and it is assumed to be Gaussian for this model.

With the introduction of energy window, we need to use Beer's law k times to express the corresponding equation. Compared with Equation (3.2) in Chapter 3, this discrete model has multiple columns corresponding to the projected data of specified detector bins so it can be expressed in a form of matrix equation. However, we can still expand the unknown linear attenuation coefficient $\mu(\vec{r}(t), e)$ into a summation of multiplication of $u_{m,e}$, the linear attenuation coefficient for the m -th material at the energy level e , and $w_m(\vec{r})$, the unknown weight of the m -th material at the position \vec{r} :

$$\mu(\vec{r}(t), e) = \sum_{m=1}^{N_m} u_{m,e} w_m(\vec{r}), \quad (4.3)$$

where

- N_m is the number of materials that form the object.

Again, we can shrink the size of unknown variable with Expansion (4.3). The size of the unknown variable is made up of two dimensions and one dimension, the resolution of material map, remain the same but the other dimension, the number of

discrete energies, has been reduced to the number of materials. Usually, the number of discrete energies can be hundreds but the number of materials are 2 or 3 so the dimension of new solution space is significantly decreased. As we have seen in Chapter 3, the weight fraction, $w_m(\vec{r})$, can be represented as a summation of product of the weight $w_{j,m}$ and the basis function $\phi_j(\vec{r})$

$$w_m(\vec{r}) = \sum_{j=1}^{N_v} w_{j,m} \phi_j(\vec{r}), \quad (4.4)$$

where

- N_v is the number of voxels (pixels if 2D) of images that compose the object.
- $w_{j,m}$ is the weight fraction of the m -th material in the j -th voxel (pixels if 2D).
- $\phi_j(\vec{r})$ is the basis function of image representation. The line integral of the basis function, $a_{i,j}$, is the length of the x-ray beam through the j -th voxel (pixel if 2D), incident onto the i -th element of the product of detector pixels N_d and number of projections N_p :

$$a_{i,j} = \int_{t \in l} \phi_j(\vec{r}(t)) dt. \quad (4.5)$$

With (4.4) and (4.5), the unknown linear attenuation coefficients can be represented as

$$\int_{t \in l} \mu(\vec{r}(t), e) dt = \sum_{m=1}^{N_m} \sum_{j=1}^{N_v} u_{m,e} w_{j,m} \int_{t \in l} \phi_j(\vec{r}(t)) dt = \sum_{j=1}^{N_v} \sum_{m=1}^{N_m} a_{i,j} w_{j,m} u_{m,e}. \quad (4.6)$$

By ignoring quadrature errors and discretizing the integral with respect to the energy E , the discretization of the basic model (4.2) can be written as:

$$y_i^{(k)} = \sum_{e=1}^{N_e} s_e^{(k)} \exp \left(- \sum_{j=1}^{N_v} \sum_{m=1}^{N_m} a_{i,j} w_{j,m} u_{m,e} \right) + \eta_i^{(k)}, \quad (4.7)$$

where N_e is the number of discrete energies. If we collect $a_{i,j}$, $w_{i,j}$ and $u_{m,e}$ in matrix form and concatenate $y_i^{(k)}$, $s_e^{(k)}$, $\eta_i^{(k)}$ with respect to the specified energy windows, then the corresponding matrix equation can be represented as:

$$\mathbf{Y} = \exp(-\mathbf{A}\mathbf{W}\mathbf{C}^T) \mathbf{S} + \boldsymbol{\mathcal{E}}, \quad (4.8)$$

where

- \mathbf{Y} is a matrix of the size $(N_d \cdot N_p) \times N_b$ that gathers x-ray photons of each energy window in the corresponding column.
- \mathbf{A} is a matrix of the size $(N_d \cdot N_p) \times N_v$ that collects the fan-beam geometry and each element corresponds to $a_{i,j}$.
- \mathbf{C} is a matrix of the size $N_e \times N_m$ that accumulates linear attenuation coefficients and each entry corresponds to $u_{e,m}$, the linear attenuation coefficient of the m -th material at the energy level e . For similar materials, we expect their linear attenuation coefficients to be similar so it might introduce the collinearity of the matrix \mathbf{C} .
- \mathbf{S} is a matrix of the size $N_e \times N_b$ and each column collects the spectrum energy of a specific range.
- $\boldsymbol{\mathcal{E}}$ is the noise matrix that is of the size $(N_d \cdot N_p) \times N_b$. The assumption for noise is that $y_{il} \sim \text{Poisson}([\exp(-\mathbf{A}\mathbf{W}\mathbf{C}^T) \mathbf{S}]_{il})$ for y_{il} in \mathbf{Y} .

In Equation (4.8), the exponential function is point-wise rather than matrix function. In addition to Equation (4.8), we also require that weight fractions should be non-negative and this can be illustrated by the constraint $\mathbf{W} \geq \mathbf{0}$. Recall that in Chapter 3, we also require that the weight fractions in each voxel should add to 1. That is to say, $\sum_{m=1}^{N_m} w_{j,m} = 1$ for $j = 1, 2, \dots, N_v$. For the energy-windowed spectral CT model,

we drop this normalization requirement and allow unnormalized weights. So the only constraint is $\mathbf{W} \geq \mathbf{0}$.

Even if the constraints are simplified, we still do not want to handle the matrix form of this equation directly since the second order derivative of the objective function will introduce tensors. To avoid introducing tensors, we need to vectorize Equation (4.8) on both sides. By taking vectorization and using the properties of Kronecker product, we can transform Equation (4.8) into

$$\text{vec}(\mathbf{Y}) = (\mathbf{S}^T \otimes \mathbf{I}) \exp\{- (\mathbf{C} \otimes \mathbf{A}) \text{vec}(\mathbf{W})\} + \text{vec}(\boldsymbol{\mathcal{E}}), \quad (4.9)$$

where \mathbf{I} is the identity matrix of the size $N_d \cdot N_p$ by $N_d \cdot N_p$. To simplify the notations, we let $\mathbf{y} = \text{vec}(\mathbf{Y})$, $\mathbf{w} = \text{vec}(\mathbf{W})$ and $\boldsymbol{\eta} = \text{vec}(\boldsymbol{\mathcal{E}})$. Then Equation (4.9) can be rewritten as

$$\mathbf{y} = (\mathbf{S}^T \otimes \mathbf{I}) \exp\{- (\mathbf{C} \otimes \mathbf{A}) \mathbf{w}\} + \boldsymbol{\eta}. \quad (4.10)$$

In addition to Equation (4.10), we also require that the weight fractions should be bounded below by zero. So we introduce the constraint $\mathbf{w} \geq \mathbf{0}$.

4.2 Problem Set-up and Preconditioning

4.2.1 The Constrained Least Squares Problem

Based on Equation (4.10) and the nonnegative constraint $\mathbf{w} \geq \mathbf{0}$, we can formulate a constrained nonlinear least squares problem:

$$\begin{aligned} \min_{\mathbf{w}} \quad & \frac{1}{2} \|\mathbf{y} - (\mathbf{S}^T \otimes \mathbf{I}) \exp\{- (\mathbf{C} \otimes \mathbf{A}) \mathbf{w}\}\|_2^2 \\ \text{subject to} \quad & \mathbf{w} \geq \mathbf{0}. \end{aligned} \quad (4.11)$$

With the least squares loss function, we want to use the Gauss-Newton method to solve this problem. If we assume the residual is $r(\mathbf{w})$, then $r(\mathbf{w})$ can be represented as:

$$r(\mathbf{w}) = \mathbf{y} - (\mathbf{S}^T \otimes \mathbf{I}) \exp\{- (\mathbf{C} \otimes \mathbf{A}) \mathbf{w}\}. \quad (4.12)$$

So the Jacobian can be calculated as

$$\mathbf{J}(\mathbf{w}) = \nabla r(\mathbf{w})^T = (\mathbf{S}^T \otimes \mathbf{I}) \text{diag}\{\exp\{- (\mathbf{C} \otimes \mathbf{A}) \mathbf{w}\}\} (\mathbf{C} \otimes \mathbf{A}). \quad (4.13)$$

With the Jacobian and residual, we can represent the gradient in terms of these two terms:

$$\nabla f(\mathbf{w}) = \mathbf{J}(\mathbf{w})^T r(\mathbf{w}). \quad (4.14)$$

The Gauss-Newton approximation of Hessian can be expressed in terms of $\mathbf{J}(\mathbf{w})$:

$$\mathbf{H}(\mathbf{w}) = \mathbf{J}(\mathbf{w})^T \mathbf{J}(\mathbf{w}) = (\mathbf{C}^T \otimes \mathbf{A}^T) \mathbf{D} (\mathbf{S} \mathbf{S}^T \otimes \mathbf{I}) \mathbf{D} (\mathbf{C} \otimes \mathbf{A}), \quad (4.15)$$

where $\mathbf{D} = \text{diag}\{\exp\{- (\mathbf{C} \otimes \mathbf{A}) \mathbf{w}\}\}$. In the k -th iteration, we need to solve a Gauss-Newton system for the step \mathbf{p}_k :

$$\mathbf{H}(\mathbf{w}_k) \mathbf{p}_k = -\nabla f(\mathbf{w}_k). \quad (4.16)$$

The naive way to solve this problem is based on two steps. At first, we solve the Gauss-Newton system to obtain step \mathbf{p}_k in each iteration. Then we update the current step and project this step onto boundary $\mathbf{w} \geq \mathbf{0}$. However, this Gauss-Newton method with simple projections does not guarantee convergence [37]. Furthermore, the sparsity of the matrix \mathbf{A} and the possible collinearity of the matrix \mathbf{C} might contribute to the ill-conditioning of the matrix $\mathbf{H}(\mathbf{w})$. Therefore, direct implementation of Gauss-Newton method is ineffective and does not produce satisfactory results in

reality. Because of these problems, we think about using preconditioners and try to apply another optimization framework to solve this problem. For this new optimization framework, we will solve the same constrained least squares problem and use the same Gauss-Newton approximation.

4.2.2 Preconditioning of the Hessian

Since the Hessian $\mathbf{H}(\mathbf{w})$ is extremely ill-conditioned, the corresponding constrained least squares problem is hard to solve. To overcome this difficulty, we want to add preconditioners to the Hessian matrix $\mathbf{H}(\mathbf{w})$. It is easy to see that the Hessian matrix $\mathbf{H}(\mathbf{w})$ is a product of several matrices, and among these matrices, it is hard to modify either $\mathbf{C} \otimes \mathbf{A}$ or $\mathbf{S}\mathbf{S}^T \otimes \mathbf{I}$. On the other hand, \mathbf{D} is a diagonal matrix so it might be convenient to construct preconditioners based on this matrix. If we can decompose the matrix \mathbf{D} into a Kronecker product of two matrices, then we can use the properties of Kronecker product to combine the left and right terms. With this idea, we try to decompose the matrix \mathbf{D} into a Kronecker product of two diagonal matrices, \mathbf{D}_1 and \mathbf{D}_2 , where \mathbf{D}_1 is of the size N_e by N_e and \mathbf{D}_2 is of the size $N_d \cdot N_p$ by $N_d \cdot N_p$:

$$\mathbf{D} \approx \mathbf{D}_1 \otimes \mathbf{D}_2. \quad (4.17)$$

Moreover, $\mathbf{D}_1 \otimes \mathbf{D}_2$ is chosen to minimize the distance to \mathbf{D} with respect to the Frobenius norm:

$$\min_{\mathbf{D}_1, \mathbf{D}_2} \|\mathbf{D} - \mathbf{D}_1 \otimes \mathbf{D}_2\|_F \quad (4.18)$$

This is a nearest Kronecker product (NKP) problem [55] and the solution has already been studied extensively. Since we require that \mathbf{D} , \mathbf{D}_1 and \mathbf{D}_2 are diagonal matrices, it is equivalent to minimizing the distance on behalf of their diagonal entries. In this case, we let $\mathbf{D} = \text{diag}\{\mathbf{d}\}$, $\mathbf{D}_1 = \text{diag}\{\mathbf{d}_1\}$ and $\mathbf{D}_2 = \text{diag}\{\mathbf{d}_2\}$, where \mathbf{d} , \mathbf{d}_1 and \mathbf{d}_2 are diagonal entries of the matrices \mathbf{D} , \mathbf{D}_1 and \mathbf{D}_2 , respectively. So Problem (4.18)

is equivalent to a least squares problem that is based on their diagonal elements with respect to the 2-norm:

$$\min_{\mathbf{d}_1, \mathbf{d}_2} \|\mathbf{d} - \mathbf{d}_1 \otimes \mathbf{d}_2\|_2 \quad (4.19)$$

It is easy to see that \mathbf{d} , \mathbf{d}_1 and \mathbf{d}_2 are vectors, which are also rank-1 matrices. In this case, the goal of this problem is to find two rank-1 matrices that will minimize the distance with respect to the 2-norm. The result is given by the largest singular value and its corresponding singular vectors. By using the singular value decomposition (SVD), one solution to this NKP problem is $\mathbf{d}_1 = \sqrt{\sigma_1} \mathbf{v}_1$ and $\mathbf{d}_2 = \sqrt{\sigma_1} \mathbf{u}_1$, where σ_1 is the largest singular value and \mathbf{u}_1 and \mathbf{v}_1 are the first left and right singular vectors of the SVD of the matrix $\tilde{\mathbf{D}} = \text{reshape}(\mathbf{d}, N_d \cdot N_p, N_e)$. In practice, we use an efficient MATLAB package, “PROPACK” [39], to compute the largest singular value and the corresponding singular vectors.

After we have obtained \mathbf{D}_1 and \mathbf{D}_2 , we can estimate $\mathbf{H}(\mathbf{w})$ using the properties of Kronecker product:

$$\begin{aligned} \mathbf{H}(\mathbf{w}) &\approx (\mathbf{C}^T \otimes \mathbf{A}^T) (\mathbf{D}_1 \otimes \mathbf{D}_2) (\mathbf{S}\mathbf{S}^T \otimes \mathbf{I}) (\mathbf{D}_1 \otimes \mathbf{D}_2) (\mathbf{C} \otimes \mathbf{A}) \\ &= (\mathbf{C}^T \mathbf{D}_1 \mathbf{S}\mathbf{S}^T \mathbf{D}_1 \mathbf{C}) \otimes (\mathbf{A}^T \mathbf{D}_2 \mathbf{D}_2 \mathbf{A}). \end{aligned} \quad (4.20)$$

The size of the matrix \mathbf{C} is N_e by N_m and N_m is the number of materials. In practice, we usually consider only 2 or 3 materials as the composition of the object. Therefore, the size of the matrix product, $\mathbf{C}^T \mathbf{D}_1 \mathbf{S}\mathbf{S}^T \mathbf{D}_1 \mathbf{C}$, is usually 2 by 2 or 3 by 3. Using this approximation, we can use the preconditioners to transform this term into identity so the new matrix should not depend on this part and thus be better-conditioned. To be specific, we let \mathbf{M}_1 be the preconditioner of the size $N_d \cdot N_p$ by $N_d \cdot N_p$ and \mathbf{M}_2 be the preconditioner of the size N_e by N_e . Then the preconditioned Hessian can be

represented as:

$$\begin{aligned}
\tilde{\mathbf{H}}(\mathbf{w}) &= (\mathbf{M}_2^T \otimes \mathbf{M}_1^T) \mathbf{H}(\mathbf{w}) (\mathbf{M}_2 \otimes \mathbf{M}_1) \\
&\approx (\mathbf{M}_2^T \otimes \mathbf{M}_1^T) ((\mathbf{C}^T \mathbf{D}_1 \mathbf{S} \mathbf{S}^T \mathbf{D}_1 \mathbf{C}) \otimes (\mathbf{A}^T \mathbf{D}_2 \mathbf{D}_2 \mathbf{A})) (\mathbf{M}_2 \otimes \mathbf{M}_1) \quad (4.21) \\
&= (\mathbf{M}_2^T \mathbf{C}^T \mathbf{D}_1 \mathbf{S} \mathbf{S}^T \mathbf{D}_1 \mathbf{C} \mathbf{M}_2) \otimes (\mathbf{M}_1^T \mathbf{A}^T \mathbf{D}_2 \mathbf{D}_2 \mathbf{A} \mathbf{M}_1).
\end{aligned}$$

Since we can choose \mathbf{D}_1 to guarantee $\mathbf{C}^T \mathbf{D}_1 \mathbf{S} \mathbf{S}^T \mathbf{D}_1 \mathbf{C}$ to be symmetric positive definite, then we can calculate the Cholesky decomposition to this matrix as

$$\mathbf{C}^T \mathbf{D}_1 \mathbf{S} \mathbf{S}^T \mathbf{D}_1 \mathbf{C} = \mathbf{G}^T \mathbf{G}, \quad (4.22)$$

where \mathbf{G} is an upper triangular matrix with positive diagonal entries. In this case, we can choose $\mathbf{M}_2 = \mathbf{G}^{-1}$ and the first part in $\tilde{\mathbf{H}}(\mathbf{w})$ has become identity. On the other hand, $\mathbf{A}^T \mathbf{D}_2 \mathbf{D}_2 \mathbf{A}$ is still large-scale because of the dimension of \mathbf{A} . It is challenging to find a preconditioner \mathbf{M}_1 that best fits all conditions. In this situation, we might choose \mathbf{M}_1 to be a diagonal matrix such that each column of the matrix $\mathbf{M}_1^T \mathbf{A}^T \mathbf{D}_2 \mathbf{D}_2 \mathbf{A} \mathbf{M}_1$ is scaled to 1 with respect to either 1-norm or 2-norm. Or we can even choose it to be identity so that we do not add any preconditioners to this part.

Using the SVD, we can analyze the condition number before and after preconditioning. Before preconditioning, the matrix $\mathbf{H}(\mathbf{w})$ depends on two parts, $\mathbf{C}^T \mathbf{D}_1 \mathbf{S} \mathbf{S}^T \mathbf{D}_1 \mathbf{C}$ and $\mathbf{A}^T \mathbf{D}_2 \mathbf{D}_2 \mathbf{A}$. If we assume that the singular value decompositions of these two matrices are $\mathbf{C}^T \mathbf{D}_1 \mathbf{S} \mathbf{S}^T \mathbf{D}_1 \mathbf{C} = \mathbf{U}_1 \mathbf{\Sigma}_1 \mathbf{V}_1^T$ and $\mathbf{A}^T \mathbf{D}_2 \mathbf{D}_2 \mathbf{A} = \mathbf{U}_2 \mathbf{\Sigma}_2 \mathbf{V}_2^T$, then the condition number of $\mathbf{H}(\mathbf{w})$ is closely related to the diagonal entries of $\mathbf{\Sigma}_1$ and $\mathbf{\Sigma}_2$. Let the largest and smallest diagonal elements of $\mathbf{\Sigma}_1$ be $\sigma_{1\max}$ and $\sigma_{1\min}$. In addition, let the largest and smallest diagonal elements of $\mathbf{\Sigma}_2$ be $\sigma_{2\max}$ and $\sigma_{2\min}$, then the condition number of $\mathbf{H}(\mathbf{w})$ can be estimated as

$$\kappa(\mathbf{H}(\mathbf{w})) \approx \frac{\sigma_{1\max} \sigma_{2\max}}{\sigma_{1\min} \sigma_{2\min}}, \quad (4.23)$$

where κ is the condition number. In contrast, the condition number of $\tilde{\mathbf{H}}(\mathbf{w})$ is mostly dependent on $\mathbf{A}^T \mathbf{D}_2 \mathbf{D}_2 \mathbf{A}$ after preconditioning, which is related to the largest and smallest diagonal entries of Σ_2 :

$$\kappa(\tilde{\mathbf{H}}(\mathbf{w})) \approx \frac{\sigma_{2\max}}{\sigma_{2\min}}. \quad (4.24)$$

In most cases, $\sigma_{1\max}/\sigma_{1\min}$ is not close to 1 and $\mathbf{M}_1 = \mathbf{I}$, then we can conclude that

$$\kappa(\tilde{\mathbf{H}}(\mathbf{w})) \ll \kappa(\mathbf{H}(\mathbf{w})). \quad (4.25)$$

In practice, the reduction of condition number can be at least two orders of magnitude. In addition to the condition number, we can check if the eigenvalues are more clustered after preconditioning. With more clustered eigenvalues, optimization methods such as conjugate gradient (CG) or generalized minimal residual method (GMRES) [51] will converge with faster speed. To check the eigenvalues before and after preconditioning, we construct an object of two materials and the material map corresponding to each material is of the size 16 by 16. Therefore, the resulting Hessian matrices, $\mathbf{H}(\mathbf{w})$ and $\tilde{\mathbf{H}}(\mathbf{w})$, are of the size 512 by 512. The plot of eigenvalues of these two matrices is presented in Figure 4.1. In Figure 4.1, we take the logarithm of eigenvalues to compare the clusters. From this plot, we can easily see that the eigenvalues of the original Hessian $\mathbf{H}(\mathbf{w})$ are scattered in a larger span, while the eigenvalues of the preconditioned Hessian $\tilde{\mathbf{H}}(\mathbf{w})$ are clustered within a smaller range. Usually, the clustered eigenvalues provide favorable convergence speed. Since the preconditioned Hessian, $\tilde{\mathbf{H}}(\mathbf{w})$, relies on \mathbf{w} , we need to compute the preconditioners, \mathbf{M}_1 and \mathbf{M}_2 , in each iteration. However, we only need the largest singular value and the corresponding singular vectors, which is very cheap.

Even if we have obtained the new preconditioners, it is still not straightforward about how to implement it for the optimization problem. If we follow the steps in

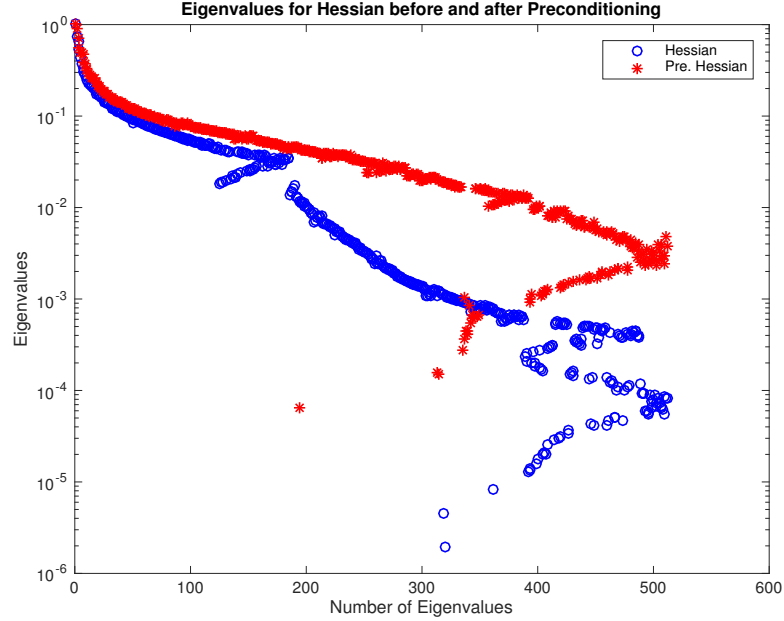


Figure 4.1: The comparison of eigenvalues before and after preconditioning (with scaling).

Section 4.2 directly, we need to solve a preconditioned system:

$$(\mathbf{M}_2^T \otimes \mathbf{M}_1^T) \mathbf{H}(\mathbf{w}_k) (\mathbf{M}_2 \otimes \mathbf{M}_1) (\mathbf{M}_2^{-1} \otimes \mathbf{M}_1^{-1}) \mathbf{p}_k = -(\mathbf{M}_2^T \otimes \mathbf{M}_1^T) \nabla f(\mathbf{w}_k). \quad (4.26)$$

If we let $\tilde{\mathbf{p}}_k = (\mathbf{M}_2^{-1} \otimes \mathbf{M}_1^{-1}) \mathbf{p}_k$ and $\nabla \tilde{f}(\mathbf{w}_k) = (\mathbf{M}_2^T \otimes \mathbf{M}_1^T) \nabla f(\mathbf{w}_k)$, then this system can be rewritten as

$$\tilde{\mathbf{H}}(\mathbf{w}_k) \tilde{\mathbf{p}}_k = -\nabla \tilde{f}(\mathbf{w}_k). \quad (4.27)$$

After solving this system, we need to retrieve the step \mathbf{p}_k from $\tilde{\mathbf{p}}_k$ and tailor the step to meet the nonnegative constraint $\mathbf{w} \geq \mathbf{0}$. In practice, if we update the current step and project the new step onto the boundary, it is not efficient and unable to guarantee convergence. Therefore, how to implement preconditioners efficiently under the constraint is worth considering.

4.3 Optimization and Regularization

4.3.1 Optimization with the Proposed Preconditioner

With the proposed preconditioners, we can implement them directly into Equation (4.16) and project the new step onto the boundary. However, numerical experiments are not in favor of this method. Without constraints, it is likely that the step obtained by this method can offer sufficient reduction. After we project the new step onto the boundary, it might not be a feasible choice anymore. Furthermore, we cannot guarantee convergence even if we add a line search scheme. To improve convergence speed and further stabilize solution, we apply the preconditioners into a two-step method. In the first step, we use the projected line search method to find the approximate Cauchy point. In the second step, we fix the elements of this Cauchy point that are active on the boundary and minimize the objective function with points that are inactive. In both steps, we include trust regions in terms of the infinite norm to further restrict the step size. The idea is that the step obtained using this method should be at least better than the approximate Cauchy point.

In the first step, we use a projected line search method to find the approximate Cauchy point. The basic idea is to find the point that will both satisfy the Wolfe condition [56] and the trust region constraint:

$$\begin{aligned} f(\mathbf{w}_k(\alpha)) &\leq f(\mathbf{w}_k) + c_1 \alpha \nabla f_k^T(\mathbf{w}_k(\alpha) - \mathbf{w}_k), \\ \|\mathbf{w}_k(\alpha) - \mathbf{w}_k\|_\infty &\leq c_2 \Delta_k, \end{aligned} \tag{4.28}$$

where $\mathbf{w}_k(\alpha) = P(\mathbf{w}_k - \alpha \nabla f_k)$ and P is the operator that is used to project \mathbf{w} onto the boundary. Δ_k is the size of trust region for the k -th step. c_1, c_2 are constants and $c_1, c_2 \in (0, 1)$. c_2 is chosen to be close to 1 to guarantee the trust region in the first step is smaller, but also close, to the trust region in the second step. Details are presented in Algorithm 1.

Algorithm 1 Projected Line Search

```

1: Initialization:
2:  $c_1 = 0.1$ ;  $c_2 = 0.8$ ;  $\alpha = \beta = 0.5$ ;
3: Set up the maximum number of iterations:  $\text{maxIter}$ .
4: while  $LS \leq \text{maxIter}$  do
5:    $\mathbf{w}_k(\alpha) = P(\mathbf{w}_k - \alpha \nabla f_k)$ ;
6:   if  $f(\mathbf{w}_k(\alpha)) \leq f(\mathbf{w}_k) + c_1 \alpha \nabla f_k^T(\mathbf{w}_k(\alpha) - \mathbf{w}_k)$  and  $\|\mathbf{w}_k(\alpha) - \mathbf{w}_k\|_\infty \leq c_2 \Delta_k$ 
     then
7:     Return  $\mathbf{w}_k(\alpha)$ ;
8:     Break;
9:    $\alpha = \alpha * \beta$ ;
10:   $LS = LS + 1$ ;
11: if  $LS > \text{maxIter}$  then
12:  Line search failed;
13:  Break;

```

If we assume that the approximate Cauchy point found in the first step is \mathbf{w}^c and the active set corresponding to this point is $\mathcal{A}(\mathbf{w}^c)$, then we can construct a quadratic programming problem that is based on the current step \mathbf{w}_k :

$$\begin{aligned}
\min_{\mathbf{w}} \quad & f(\mathbf{w}_k) + \nabla f(\mathbf{w}_k)^T (\mathbf{w} - \mathbf{w}_k) + \frac{1}{2} (\mathbf{w} - \mathbf{w}_k)^T \mathbf{H}(\mathbf{w}_k) (\mathbf{w} - \mathbf{w}_k) \\
\text{subject to} \quad & \mathbf{w} \geq \mathbf{0}, \quad i \notin \mathcal{A}(\mathbf{w}^c), \\
& \mathbf{w}_i = \mathbf{w}_i^c, \quad i \in \mathcal{A}(\mathbf{w}^c), \\
& \|\mathbf{w} - \mathbf{w}_k\|_\infty \leq \Delta_k.
\end{aligned} \tag{4.29}$$

Moreover, we should notice that $f(\mathbf{w}_k)$ is a constant and we can ignore this term. If we let $\mathbf{d} = \mathbf{w} - \mathbf{w}_k$, then we can simplify the objective function as $\nabla f(\mathbf{w}_k)^T \mathbf{d} + \frac{1}{2} \mathbf{d}^T \mathbf{H}(\mathbf{w}_k) \mathbf{d}$. The inequality constraint $\mathbf{w} \geq \mathbf{0}$ is equivalent to $\mathbf{w} - \mathbf{w}_k \geq -\mathbf{w}_k$, which is $\mathbf{d} \geq -\mathbf{w}_k$ using our notations. The third inequality constraint, $\|\mathbf{w} - \mathbf{w}_k\|_\infty \leq \Delta_k$, can be rewritten as $\|\mathbf{d}\|_\infty \leq \Delta_k$. This inequality constraint is equivalent to $-\Delta_k \leq \mathbf{d} \leq \Delta_k$. By combining the previous two constraints, we can obtain that $\max\{-\Delta_k, -\mathbf{w}_k\} \leq \mathbf{d} \leq \Delta_k$. For the elements in the active set, we have $\mathbf{w} - \mathbf{w}_k = \mathbf{w}^c - \mathbf{w}_k$, which is $\mathbf{d} = \mathbf{w}^c - \mathbf{w}_k$. So the previous optimization problem is equivalent

to

$$\begin{aligned}
& \min_{\mathbf{d}} \quad \nabla f(\mathbf{w}_k)^T \mathbf{d} + \frac{1}{2} \mathbf{d}^T \mathbf{H}(\mathbf{w}_k) \mathbf{d} \\
& \text{subject to} \quad \max\{-\Delta_k, -\mathbf{w}_k\} \leq \mathbf{d} \leq \Delta_k, \\
& \quad \quad \quad \mathbf{P}\mathbf{d} = \mathbf{w}^c - \mathbf{w}_k.
\end{aligned} \tag{4.30}$$

where \mathbf{P} is the projection matrix that keeps \mathbf{d}_i invariant for $i \in \mathcal{A}(\mathbf{w}^c)$. We can construct the matrix \mathbf{P} by assigning the i -th element in the i -th column to be 1 and otherwise 0, where the size of \mathbf{P} is the number of active elements by $N_v \cdot N_m$. It is obvious that this problem is a quadratic programming problem with both bound and equality constraints. However, for the large-scale cases, solving this problem exactly is likely to be as expensive as solving the original problem [48]. Since we have obtained the approximate Cauchy point in the first step, we only want the solution to be at least better than the approximate Cauchy point, which means that we do not require an exact solution for the second step. So we can either stop the iteration when the current step crosses the boundaries or ignore the bound constraints, project the step back to the boundaries after it meets the stopping criteria and compare it with the approximate Cauchy point. In this chapter, we will ignore the bound constraints at first and then project the step back to the boundaries because the step can cross the boundaries within only 1 or 2 iterations, and in this case, it does not give us a sufficient reduction.

So far, we have not added any preconditioners to this problem. As we can see, the preconditioners proposed in the last section are used to precondition the Hessian, so it is not necessary to apply the preconditioners in the first step. In the second step, we substitute the Hessian of Gauss-Newton approximation for the true Hessian so we can try to use the preconditioners in this step. If we ignore the bound constraints,

we can formulate a quadratic programming problem with the preconditioners:

$$\begin{aligned} \min_{\tilde{\mathbf{d}}} \quad & \nabla \tilde{f}(\mathbf{w}_k)^T \tilde{\mathbf{d}} + \frac{1}{2} \tilde{\mathbf{d}}^T \tilde{\mathbf{H}}(\mathbf{w}_k) \tilde{\mathbf{d}} \\ \text{subject to} \quad & \tilde{\mathbf{P}} \tilde{\mathbf{d}} = \mathbf{w}^c - \mathbf{w}_k, \end{aligned} \quad (4.31)$$

where $\nabla \tilde{f}(\mathbf{w}_k) = (\mathbf{M}_2^T \otimes \mathbf{M}_1^T) \nabla f(\mathbf{w}_k)$, $\tilde{\mathbf{d}} = (\mathbf{M}_2^{-1} \otimes \mathbf{M}_1^{-1}) \mathbf{d}$, $\tilde{\mathbf{P}} = \mathbf{P}(\mathbf{M}_2 \otimes \mathbf{M}_1)$ and $\tilde{\mathbf{H}}(\mathbf{w}) = (\mathbf{M}_2^T \otimes \mathbf{M}_1^T) \mathbf{H}(\mathbf{w})(\mathbf{M}_2 \otimes \mathbf{M}_1)$. We use the projected conjugated gradient (PCG) method to solve $\tilde{\mathbf{d}}$ directly and retrieve \mathbf{d} using $\mathbf{d} = (\mathbf{M}_2 \otimes \mathbf{M}_1) \tilde{\mathbf{d}}$. The details are shown in Algorithm 2. In Algorithm 2, we need to use the approxi-

Algorithm 2 Projected Conjugate Gradient

- 1: *Initialization:*
 - 2: $\mathbf{d} = \mathbf{w}^c - \mathbf{w}_k$;
 - 3: $\tilde{\mathbf{d}} = (\mathbf{M}_2^{-1} \otimes \mathbf{M}_1^{-1}) \mathbf{d}$;
 - 4: $\mathbf{r} = \tilde{\mathbf{H}}(\mathbf{w}_k) \tilde{\mathbf{d}} + \nabla \tilde{f}(\mathbf{w}_k)$;
 - 5: $\mathbf{g} = \mathbf{r} - \mathbf{P}^T (\mathbf{P}\mathbf{P}^T)^{-1} \mathbf{P}\mathbf{r}$;
 - 6: $\mathbf{p} = -\mathbf{g}$;
 - 7: Set up the stopping criterion: tol, and the maximum number of iterations: maxIter;
 - 8: **while** $i \leq \text{maxIter}$ **do**
 - 9: $\alpha = \mathbf{r}^T \mathbf{g} / \mathbf{p}^T \tilde{\mathbf{H}}(\mathbf{w}_k) \mathbf{p}$;
 - 10: $\tilde{\mathbf{d}} = \tilde{\mathbf{d}} + \alpha \mathbf{p}$;
 - 11: $\mathbf{r}^+ = \mathbf{r} + \alpha \tilde{\mathbf{H}}(\mathbf{w}_k) \mathbf{p}$;
 - 12: $\mathbf{g}^+ = \mathbf{r}^+ - \mathbf{P}^T (\mathbf{P}\mathbf{P}^T)^{-1} \mathbf{P}\mathbf{r}^+$;
 - 13: **if** $\mathbf{g}^{+T} \mathbf{r}^+ < \text{tol}$ **then**
 - 14: Return $\tilde{\mathbf{d}}$;
 - 15: Break;
 - 16: $\beta = \mathbf{r}^{+T} \mathbf{g}^+ / \mathbf{r}^T \mathbf{g}$;
 - 17: $\mathbf{p} = -\mathbf{g}^+ + \beta \mathbf{p}$;
 - 18: $\mathbf{g} = \mathbf{g}^+$;
 - 19: $\mathbf{r} = \mathbf{r}^+$;
 - 20: $i = i + 1$;
-

mate Cauchy point obtained in Algorithm 1 as the initial guess. Otherwise, we cannot guarantee the elements that are active will be fixed on the boundaries. Moreover, we can see that we need $(\mathbf{P}\mathbf{P}^T)^{-1}$ in each iteration, where $\mathbf{P}\mathbf{P}^T$ is large-scale, sparse and symmetric positive semi-definite. So we can use the Cholesky factorization to

find an upper triangular matrix \mathbf{R} such that $\mathbf{P}\mathbf{P}^T = \mathbf{R}^T\mathbf{R}$. Compared with inverting $\mathbf{P}\mathbf{P}^T$ directly, it is more efficient to invert $\mathbf{R}^T\mathbf{R}$ because \mathbf{R} is upper triangular and we only need to calculate the factorization once in Algorithm 2.

After we have obtained the new direction \mathbf{d} , we can compute the new step as $\mathbf{w}_{k+1} = \mathbf{w}_k + \mathbf{d}$. To decide if we should accept the new step and the size of trust regions, we need to compute the actual reduction and the predicted reduction. The actual reduction is the difference between function values of the previous step and the new step:

$$ared = f(\mathbf{w}_k) - f(\mathbf{w}_{k+1}). \quad (4.32)$$

The predicted reduction is based on the quadratic expansion of the previous step:

$$pred = -\nabla f(\mathbf{w}_k)^T \mathbf{d} - \frac{1}{2} \mathbf{d}^T \mathbf{H}(\mathbf{w}_k) \mathbf{d}. \quad (4.33)$$

If $ared < 0$, then we accept the new step and update the trust region with respect to the ratio $ared/pred$. In this chapter, we use the standard way to update trust regions. After updating the current step, we should check if we should stop the iteration. If the current step meets the stopping criteria for local minimum, then we stop the iteration and report the results. Otherwise, we keep iterating until we find a local minimizer. By combining Algorithm 1 and Algorithm 2, we can express the main framework in Algorithm 3.

From the description, we can see that Algorithm 3 includes the strengths of both the projected line search method and the trust region method. With the projected line search method, we can guarantee a descent direction with proper reduction. With the trust region method, we are likely to obtain a better step and thus further reduction. However, this problem is a nonlinear least squares problem and the objective function might have multiple local minimizers. In addition, we truncate each step with respect to the lower bounds and it is hard to analyze the convergence properties. It is possible

Algorithm 3 Main Framework

- 1: *Initialization:*
 - 2: Set up the test problem, the initial guess \mathbf{w}_0 and the stopping criteria.
 - 3: **while** the stopping criteria do not satisfy **do**
 - 4: Use Algorithm 1 to compute the approximate Cauchy point \mathbf{w}^c ;
 - 5: Find the active set of \mathbf{w}^c ;
 - 6: Construct the project matrix \mathbf{P} based on the active set;
 - 7: Calculate the new boundaries based on the original boundaries and the current trust region;
 - 8: Generate the preconditioners based on the current step \mathbf{w}_i ;
 - 9: Use Algorithm 2 to compute the new step \mathbf{w}_{new} ;
 - 10: Project the new step \mathbf{w}_{new} onto the new boundaries to obtain \mathbf{w}_{i+1} ;
 - 11: Calculate $pred$ and $ared$ using \mathbf{w}_i and \mathbf{w}_{i+1} ;
 - 12: Decide if we should accept \mathbf{w}_{i+1} and update the trust region;
-

that the new step might not offer sufficient reduction but it is closer to the global minimizer. For example, we can think about an 1D problem of climbing the hill. On the other side of the hill, it might have another point that has even lower altitude. During the process of climbing, we might need to increase the altitude to a certain degree in order to achieve better results later. Therefore, it is not possible to decide specific criteria that fit all conditions. For example, we can even accept the new step every time and only decide the size of trust region with respect to the ratio $ared/pred$. In practice, this method still gives us a fast convergence and high-quality image reconstruction.

4.3.2 Regularization and Scaling

In this chapter, we assume that the measured data obtained from the detector follow a Poisson distribution. To remove the noise and stabilize the solution, we need to introduce regularization terms. As the weights of each material correspond to an image, we can construct the regularization term with respect to the specified material. If we assume \mathbf{w} is concatenated by $\mathbf{w}_1, \mathbf{w}_2, \dots, \mathbf{w}_m$ where \mathbf{w}_i is the vectorization

of the i -th material map, then we can build the regularization term as

$$R(\mathbf{w}) = \sum_{i=1}^{N_m} \alpha_i R(\mathbf{w}_i), \quad (4.34)$$

where $R(\mathbf{w}_i)$ represents the regularization term for the i -th material and α_i is the corresponding regularization parameter. Using this method, it gives us more freedom to choose the regularization term with respect to the material. In practice, we can usually find that one or several specified materials dominate the object and it is likely that the reconstructed material maps will contain many edges. So we select generalized Tikhonov regularization for these materials to smooth the edges. The forward difference operator and zero boundary conditions are used to build the discrete differential operator. On the other hand, other materials might only occupy a small area so we think about restrict the sum of weights for these materials. To realize this idea, we introduce ℓ_1 regularization to penalize the sum of weights. With the generalized Tikhonov regularization and the ℓ_1 regularization, we need to select the corresponding regularization parameters. It is not obvious how we can use regular methods, such as L-curve and generalized cross-validation, to find proper parameters. In this case, we generate a log space for each regularization and use the grid search method to find the “best” regularization parameters.

Another challenge associated with the regularization parameter is how to scale the problem. The point such that the objective function is zero might not be feasible and the residual corresponding to the global minimizer might still be a large number. In this case, choosing a proper regularization parameter is hard if the magnitude is large. In addition to selecting the regularization parameter, the infeasible step might result in the difficulty of meeting stopping criteria and thus increase the number of CG iterations. For this problem, we want to scale the objective function, gradient and Hessian with the spectrum radius of the Hessian based on the 2-norm. However,

it is not necessary to compute the largest singular value of the Hessian. It is easy to see that $\mathbf{H}(\mathbf{w})$ is a nonnegative and symmetric matrix, so we can use the Shur's test [40] to find an upper bound for this value with the initial step \mathbf{w}_0 :

$$\|\mathbf{H}(\mathbf{w}_0)\|_2 \leq \sqrt{\|\mathbf{H}(\mathbf{w}_0)\|_1 \|\mathbf{H}(\mathbf{w}_0)\|_\infty}. \quad (4.35)$$

With the nonnegativity of $\mathbf{H}(\mathbf{w}_0)$, we can calculate the right hand side of Inequality (4.35) Using matrix-vector multiplication rather than form the Hessian explicitly. Using the symmetry of the Hessian, we can obtain that $\|\mathbf{H}(\mathbf{w}_0)\|_1 = \|\mathbf{H}(\mathbf{w}_0)\|_\infty$. Then Inequality (4.35) can be simplified as

$$\|\mathbf{H}(\mathbf{w}_0)\|_2 \leq \|\mathbf{H}(\mathbf{w}_0)\|_\infty. \quad (4.36)$$

So we can use $\|\mathbf{H}(\mathbf{w}_0)\|_\infty$ as the scaling parameter for the objective function, the gradient and the Hessian. With this scaling parameter, we can choose regularization parameters with less effort.

4.4 Numerical Experiments

To test the preconditioners and the optimization method, we generate a 2D image of the size 128 by 128 as the object. We also assume that this object is made up of two materials, plexiglass and polyvinyl chloride (PVC). Thus we can obtain 2 material maps corresponding to the weights of these two materials. The original material maps are shown in Figure 4.2. In Figure 4.2, the yellow color represents that it has the corresponding material in this area, while the blue color shows that it does not have the corresponding material in this area. Therefore, we can see that the object is a circle and both materials are distributed inside this circle. Inside this circle, plexiglass dominates most areas except three dots occupied by PVC. We can also

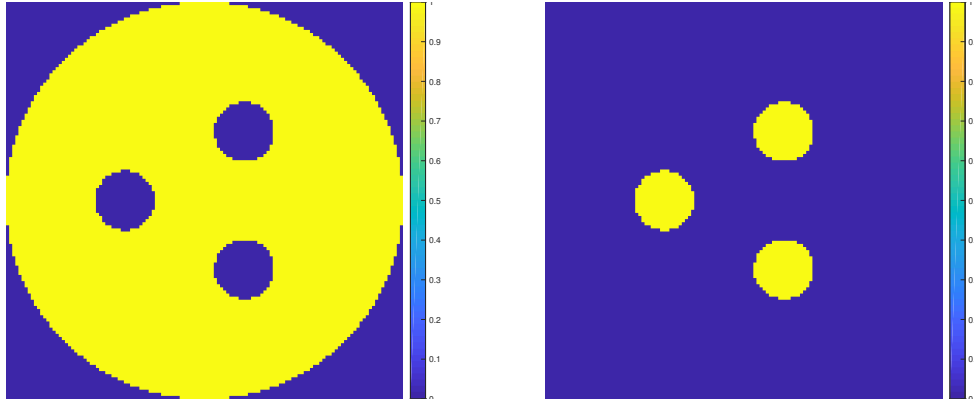


Figure 4.2: The original images for plexiglass (left) and PVC (right).

see that the images corresponding to these two materials compensate each other and they are completely separable. The goal of this numerical experiment is to reconstruct these two images such that different material maps present the corresponding material compositions.

Moreover, we present the plot of the linear attenuation coefficients for these two materials in Figure 4.3. In Figure 4.3, we can see that the slopes of these two curves

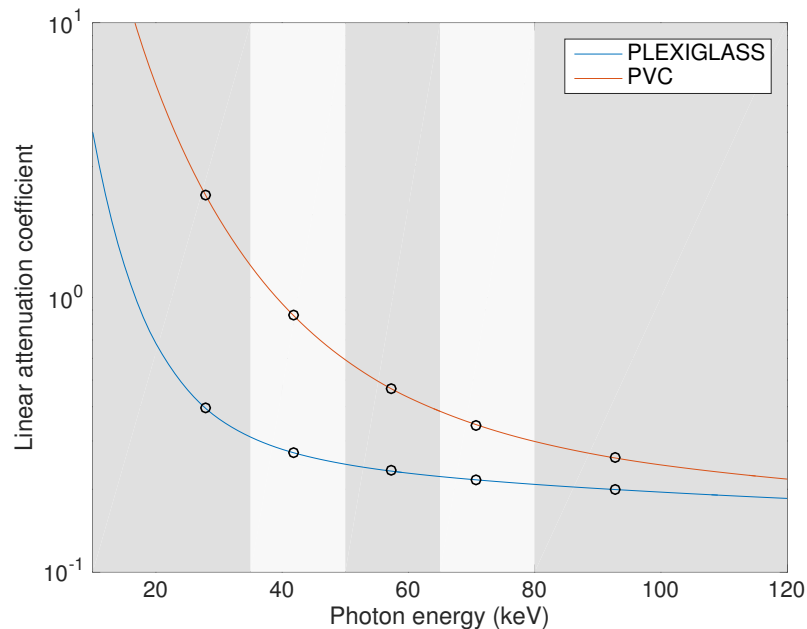


Figure 4.3: The linear attenuation coefficients for plexiglass and PVC.

are close to each other so it might indicate the collinearity of linear attenuation

coefficients and thus the linear correlation of columns of the matrix \mathbf{C} . It might also introduce small singular values and strengthen the ill-posedness.

In addition to the previous images, the geometry parameters of the CT machine are set as Table 4.1. With these parameters, we generate a distance matrix, \mathbf{A} , to represent the fan-beam geometry for flat detector using the MATLAB function `fanbeamtomolinar` [35]. Moreover, we choose 180 projections and they are dis-

Items	Parameters (cm)
Width of Domain	2.0
Distance from Source to Rotation Center	3.0
Distance from Source to Detector	5.0
Detector Width	4.0

Table 4.1: Geometry parameters of the CT machine.

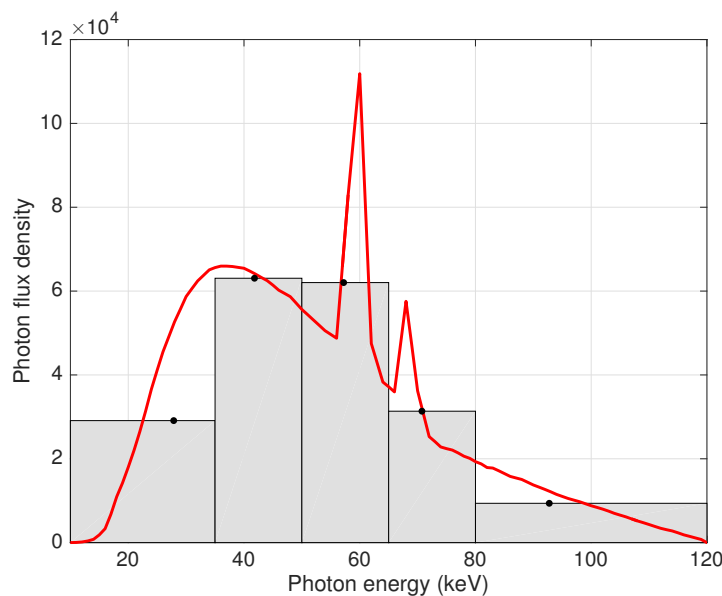


Figure 4.4: Detector bins and photon flux density.

tributed equally from 0 to 360 degrees. The voltage of the x-ray source is assumed to be 120 keV and the corresponding spectrum is generated using the MATLAB function `spektrSpectrum` [53]. We also assume that we have 5 detector bins and each of them can detect the specific range of photon energy: 10 to 34 keV, 35 to 49 keV, 50 to 64

keV, 65 to 79 keV and 80 to 120 keV. The detector is assumed to be photon counting detector and the data obtained from the detector are positive integers. Both the x-ray spectrum and the detector bins are presented in Figure 4.4. In Figure 4.4, the small black dots represent the values of mean energy in each bin. When we construct the forward problem, we use the full spectrum and all linear attenuation coefficients. However, we only use the mean energy and the corresponding linear attenuation coefficients to compute reconstruction. In this way, we can avoid inverse crime but the grid for reconstruction is coarser. With all these parameters, the goal is to reconstruct two material maps with the proposed methods.

With previous preparations, we set up random numbers between 0 and 1 as the initial guess and run the main algorithm. It only takes us around 116 seconds to converge even if we do not use the optimal implementation. The reconstructed images are presented in Figure 4.5. In Figure 4.5, we can see that we have successfully

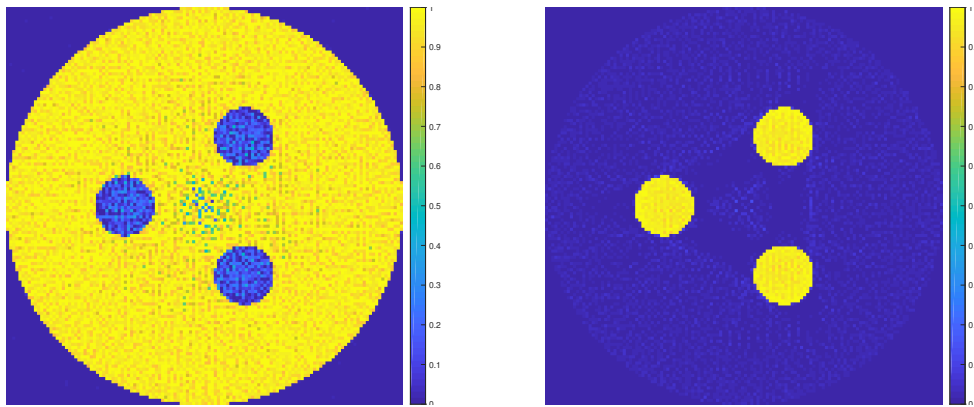


Figure 4.5: The reconstructed images for plexiglass (Left) and PVC (Right).

separated these two materials and it is hard to find many overlaps between these two material maps. Moreover, we can see that the shapes of these two material maps are exactly the same as the original images. In addition, the boundaries of large circle and three small dots are clear to identify, which shows the strength of results. By comparing Figure 4.5 with Figure 4.2, we can see that the reconstructed images have more shades than the original ones. For the first material map, most weights

are yellow but several of them are not exact 1. These light yellow weights represent the values that are not 1 but significantly close to 1. It might result from the noise contained in the data. In the center of this material map, we can see several dots of high frequency, which depart from the true solution to a certain degree. On the other hand, the reconstruction of the second material map is of higher quality and it is hard to tell the difference between the original image and the reconstructed one. It is possible that the weights corresponding to the second material map are much less than the weights of the first material map so that they contain less noise.

To check the convergence, we can also plot the relative errors for these two materials. The plot of relative errors is shown in Figure 4.6. In Figure 4.6, we can see that the curves of both materials drop and stagnate, which indicates the convergence. Moreover, the blue curve decreases faster than the red one. It shows that the convergence rate for the second material is faster than the first one. We can also see that

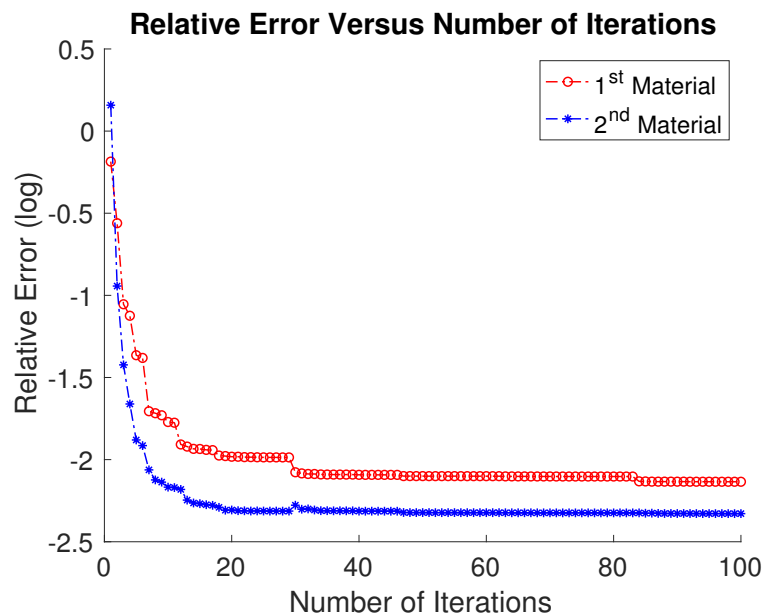


Figure 4.6: The Relative Errors for Plexiglass (red) and PVC (blue dash).

the blue curve reaches a lower level than the red one when they stagnate. It reflects the observation that the reconstruction of the second material map has better quality than the first one. Generally speaking, the ultimate relative errors for both materials

are around 10% and the overall reconstruction is satisfactory.

To further confirm the convergence, we can plot the decay of norm of the gradient. This plot is shown in Figure 4.7. In Figure 4.7, we can see that most of times, the

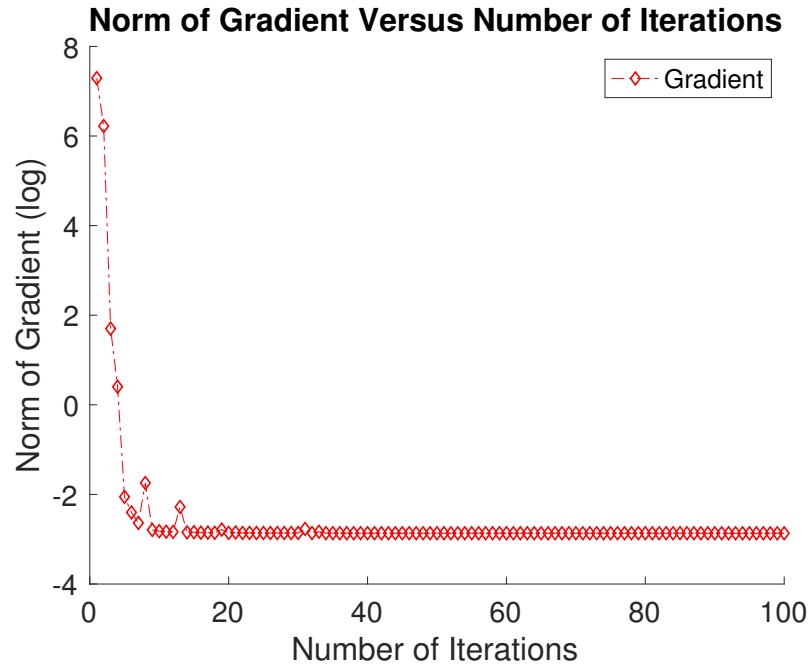


Figure 4.7: The decay of norm of the gradient.

curve drops as expected. However, there are two “spikes” that dissatisfied the trend of decrease. I think this might result from the projections onto the boundary or the irregular updating scheme of trust region. Furthermore, the best reconstructed images might not match the optimal solutions to the optimization problem.

To see the significance of the preconditioner, we can compare the number of CG iterations with and without preconditioners before the relative errors reach a specified level. We keep other parameters identical but shrink the size of images to 64 by 64 and the number of projections to 96. Moreover, we set the maximum number of CG iterations to 5000. After 5 Newton iterations, the methods with and without preconditioner both reach relative errors around 23%. The results are shown in Table 4.2. From Table 4.2, we can find that the number of CG iterations in each Newton step is significantly reduced after preconditioning. With the proposed method, the New-

#	No Precond.			Precond.		
	Rel. Err.	Num. CG	Time (s)	Rel. Err.	Num. CG	Time (s)
1	1.25	2625	87.8	0.84	9	0.3
2	1.67	318	11.5	0.46	6	0.3
3	1.26	111	3.8	0.30	13	0.5
4	0.27	4051	132.3	0.27	8	0.3
5	0.23	1347	44.2	0.22	12	0.5

Table 4.2: The comparison of CG iterations.

ton system is sufficiently better-conditioned or the eigenvalues are remarkably more clustered.

4.5 Conclusions and Remarks

With the energy-windowed spectral CT model, we set up a nonlinear least squares problem under bound constraints. To solve this optimization problem, we propose a new preconditioner and then implement it into a two-step method. The new preconditioner can transform the eigenvalues of the original Gauss-Newton system into more clustered ones, which will lead to faster convergence and higher accuracy. With the introduction of the two-step method, we can guarantee that the obtained step is at least better than the approximate Cauchy point. By solving the Gauss-Newton system in the second step, we expect further reduction from the solution. Therefore, the convergence rate should be better than linearity. Moreover, we further restrict the size of each step with the help of trust regions. In addition, we can remove parts of noise and speed up the convergence rate with the scaling parameter and the regularization terms.

On the other hand, it still has several limitations. Because of the nonlinearity of the objective function, it is hard to decide if we should accept the new step or not under certain circumstances. Furthermore, with multiple materials, we cannot use regular methods to choose regularization parameters. In each iteration, we need

to solve a NKP problem to obtain the preconditioner. Even if the computational cost is cheap, it might be better if we can find a preconditioner that is feasible for all steps. We might also think about how to implement the preconditioner into first order methods such as FISTA [3].

Chapter 5

Preconditioning and Optimization for Energy-windowed Spectral Computed Tomography

In Chapter 4, we discussed the energy-windowed spectral CT model and presented a preconditioning framework and a nonlinear optimization approach to compute the solution. In this chapter, we still focus on the energy-windowed spectral CT model

$$\mathbf{Y} = \exp(-\mathbf{A}\mathbf{W}\mathbf{C}^T) \mathbf{S} + \boldsymbol{\mathcal{E}}, \quad (5.1)$$

where \mathbf{Y} is a matrix that gathers the projected data of each energy window in the corresponding column and the exponential operator is applied element-wise (i.e., it is not a matrix function). \mathbf{A} is a matrix that is related to the quantitative information of ray trace and \mathbf{C} is a matrix that contains linear attenuation coefficients for particular (known) materials at specified energies. \mathbf{S} is the matrix that accumulates the spectrum energies for each energy window in the corresponding column. We assume that \mathbf{S} is square and invertible. Moreover, $\boldsymbol{\mathcal{E}}$ represents the noise term and we assume that $E_{il} \sim \mathcal{N}(0, y_{il})$ for each component E_{il} in $\boldsymbol{\mathcal{E}}$ and y_{il} in \mathbf{Y} . We assume that these

data are known and the target is to solve the unknown weight matrix \mathbf{W} . \mathbf{W} is of size N_v by N_m , where N_v is the number of voxels (pixels if 2D) for each material map and N_m is the number of materials. Since the weight matrix \mathbf{W} represents the material maps of different materials, then it must be nonnegative and we need to add a lower bound constraint $\mathbf{W} \geq \mathbf{0}$.

To solve Equation (5.1), we follow the instruction in Chapter 4 and take a vectorization at first. Rather than building a nonlinear optimization problem, we use the Taylor expansion to remove the point-wise exponential function and obtain an approximate linearized equation. Under the Gaussian assumption, as we show in Section 5.1, we can transform this equation into a weighted least squares problem under bound constraints:

$$\begin{aligned} \min_{\mathbf{w}} \quad & \frac{1}{2} \|\mathcal{A}\mathbf{w} - \mathbf{b}\|_{\Sigma^{-1}}^2 \\ \text{subject to} \quad & \mathbf{w} \geq \mathbf{0}, \end{aligned} \tag{5.2}$$

where $\mathcal{A} = \mathbf{C} \otimes \mathbf{A}$, $\mathbf{b} = -\log(\mathbf{y})$, $\mathbf{y} = \text{vec}(\mathbf{Y})$ and $\mathbf{w} = \text{vec}(\mathbf{W})$. Σ^{-1} , which combines information of \mathbf{S} and \mathbf{y} , is the inverse covariance matrix generated by the Gaussian noise and the log transformation. $\|\cdot\|_{\Sigma^{-1}}^2$ represents a weighted 2-norm and $\|\mathcal{A}\mathbf{w} - \mathbf{b}\|_{\Sigma^{-1}}^2 = (\mathcal{A}\mathbf{w} - \mathbf{b})^T \Sigma^{-1} (\mathcal{A}\mathbf{w} - \mathbf{b})$. \mathbf{C} is of the size N_e by N_m , where N_e is the number of energy and N_m is the number of materials. Since each column of \mathbf{C} collects the corresponding linear attenuation coefficients and two materials, such as adipose and glandular, might be similar to each other, the matrix \mathbf{C} is likely to be ill-conditioned. On the other hand, the problem (5.2) is similar to a quadratic programming problem under bound constraints. However, direct implementation of optimization solvers does not provide high-quality reconstruction because the ray trace matrix \mathbf{A} is large and ill-conditioned, and the columns of the linear attenuation coefficient matrix \mathbf{C} might be nearly collinear.

To handle these problems, we propose a new preconditioner that is based on rank-

1 approximation of the matrix \mathbf{Y} . With this rank-1 approximation, we can estimate the Hessian of the objective function in (5.2) using a Kronecker product of two parts. The first part of this Kronecker product is of the size $N_m \times N_m$, where N_m denotes the number of materials; usually this is quite small, e.g. $N_m = 2$ or 3 . This matrix product is also symmetric and positive definite so we can construct a preconditioner from its inverse Cholesky factorization, and thus transform it into identity in the preconditioned system. Because the conditioning of the Hessian is closely related to these two matrices and one of them has been transformed into the identity matrix, we have reduced the condition number significantly. Moreover, it is an economical preconditioner since we only need to compute the preconditioner once and can reuse it in the future iterations. In [2], Barber et al. propose an alternative preconditioner that is based on the eigenvalue decomposition of $\mathbf{C}^T \mathbf{C}$, where \mathbf{C} is the matrix of linear attenuation coefficients. Compared with this, the preconditioner proposed in this chapter involves not only \mathbf{C} , but also includes information of the energy spectrum \mathbf{S} and parts of the photon counting data, \mathbf{Y} , and thus provides a more physically meaningful approximation of the Hessian.

In addition, with the weighted least squares objective function, it is much easier to analyze the condition number before and after preconditioning. Since the performance of a first order method is closely related to the condition number of the Hessian, it is intuitive to implement a first order method if we can reduce the condition number significantly. Based on this idea, Fast Iterative Shrinkage-Thresholding Algorithm (FISTA) [3, 46, 47] comes into view. FISTA is a first order method that has an “optimal” function convergence rate, $\mathcal{O}(1/k^2)$, where k is the number of iterations. Furthermore, this method is suitable for solving problems that have the form $f(\mathbf{x}) + g(\mathbf{x})$, where both $f(\mathbf{x})$ and $g(\mathbf{x})$ are convex but $g(\mathbf{x})$ is possibly nonsmooth. This $f(\mathbf{x})$ can be the weighted least squares term in problem (5.2) and $g(\mathbf{x})$ can represent a nonsmooth regularization such as ℓ_1 regularization or nonnegative constraints. Even

if we can achieve fast convergence, the introduction of a preconditioner complicates the bound constraints. The previous bound constraints have become linear inequality constraints because of the preconditioner. However, we can construct a projection problem that can find the closest solutions to satisfy these constraints. Moreover, this projection problem is separable and we can apply highly efficient solvers to compute the solutions to these decomposed small-sized problems. Generally speaking, the implementations of our preconditioner, FISTA and projection problem complement each other and exhibit high-quality reconstructed images and fast convergence results.

This chapter is organized as follows. In Section 5.1, we review the continuous energy-windowed spectral CT model and the corresponding discretized nonlinear matrix equation. The key idea of this chapter, preconditioning, is introduced in Section 5.2. In this section, both the derivation of our preconditioner and an analysis of the reduction of the condition number are presented. The choice of regularization will be exhibited in this section as well. In Section 5.3, we study FISTA and how we construct and solve the projection problems. Moreover, numerical experiments are presented in Section 5.4 and concluding remarks are given in Section 5.5.

5.1 The Weighted Least Squares Problem

In Chapter 4, we have discussed energy-windowed spectral CT model and introduced a new preconditioner based on the corresponding nonlinear least squares problem and the Gauss-Newton approximation of the Hessian. A two-step optimization method that includes projected line search and trust region method is implemented to solve this problem. However, there are two main concerns related to the preconditioner and optimization. At first, the preconditioner requires the information of the current iteration and even if it is cheap to compute, we still need to repeat the computational process in each Newton iteration. Secondly, since the nonlinear optimization is based

on the Gauss-Newton approximation of the Hessian, it raises a question if we can come up with a new preconditioner and use a first order method to solve it. In this chapter, we still focus on the energy-windowed spectral CT model and we present a linearization technique to transform the nonlinear equation into an optimization problem that is based on a weighted least squares term and a bound constraint. Recall that the basic energy-windowed spectral CT model is expressed by

$$y_i^{(k)} = \int_E S^{(k)}(e) \exp\left(-\int_{t \in l} \mu(\vec{r}(t), e) dt\right) de + \eta_i^{(k)}. \quad \begin{cases} i = 1, 2, \dots, N_d \times N_p, \\ k = 1, 2, \dots, N_b, \end{cases} \quad (5.3)$$

Using the material decomposition, $\mu(\vec{r}(t), e) = \sum_{m=1}^{N_m} u_{m,e} w_m(\vec{r})$, discretizing Equation (5.3) over energy E , and concatenating corresponding variables with respect to energy windows, we can obtain a matrix equation

$$\mathbf{Y} = \exp(-\mathbf{A}\mathbf{W}\mathbf{C}^T) \mathbf{S} + \boldsymbol{\mathcal{E}}, \quad (5.4)$$

where

- \mathbf{Y} is a matrix of the size $(N_d \cdot N_p) \times N_b$ that gathers x-ray photons of each energy window in the corresponding column.
- \mathbf{A} is a matrix of the size $(N_d \cdot N_p) \times N_v$ that collects the fan-beam geometry and each element corresponds to $a_{i,j}$.
- \mathbf{C} is a matrix of the size $N_e \times N_m$ that accumulates linear attenuation coefficients and each entry corresponds to $u_{e,m}$, the linear attenuation coefficient of the m -th material at the energy level e . For similar materials such as adipose and glandular, the collinearity might cause the ill-conditioning of \mathbf{C} .
- \mathbf{S} is a matrix of the size $N_e \times N_b$ and each column collects the spectrum energy of a specific range. We assume that \mathbf{S} is square and invertible.

- $\boldsymbol{\mathcal{E}}$ is the noise matrix that is of size $(N_d \cdot N_p) \times N_b$. The assumption for noise is that $E_{il} \sim \mathcal{N}(0, y_{il})$ for each element E_{il} in $\boldsymbol{\mathcal{E}}$ and y_{il} in \mathbf{Y} .

In Figure 4.4, we present a plot of spectrum energy. The y-axis represents photon flux density while the x-axis indicates photon energy. The red curve in Figure 4.4 shows the relationship between photon flux density and energy under 120 keV voltage. The gray boxes represents the energy bins and as we can identify, there are five energy bins. The black dot in each energy bin illustrates the average photon density in the corresponding bin. If we use average energies rather than full spectrum when solving the inverse problem, the matrix \boldsymbol{S} is diagonal and it gives us more flexibility on manipulating the equation.

Moreover, the current assumption for the noise matrix $\boldsymbol{\mathcal{E}}$ is that each entry follows a normal distribution with mean 0 and variance y_{il} (the corresponding entry in \mathbf{Y}). That is to say, $E_{il} \sim \mathcal{N}(0, y_{il})$. In the last chapter, we assume that $y_{il} \sim \text{Poisson}([\exp(-\mathbf{A}\mathbf{W}\mathbf{C}^T)\boldsymbol{S}]_{il})$ and Poisson distribution is the most fundamental assumption because of randomness of x-ray photon motion and errors of photon counting. However, a Poisson distribution is accurately approximated by a Gaussian distribution if the mean of this Poisson distribution is large enough. For spectral CT problems, this mean corresponds to the number of photons and this number is often hundreds of thousands. With this large number as the mean, the Gaussian assumption is valid.

In several cases, the composition of materials can be similar. For example, glandular and adipose have similar attenuation coefficients at the same energy levels and this feature can cause the collinearity. After discretization, the columns of \boldsymbol{C} can be nearly dependent. Moreover, \mathbf{A} is large-scale and sparse and it is highly likely to have small singular values. As we will see later, the Hessian system involves the Kronecker product $\boldsymbol{C} \otimes \mathbf{A}$ and these small singular values can cause the ill-posedness. Since it is challenging to solve this equation directly, it is important to consider approaches

to facilitate the process. First, we can introduce a preconditioning matrix \mathbf{M} into Equation (5.4):

$$\mathbf{Y} = \exp(-\mathbf{A}\mathbf{W}\mathbf{M}^{-T}\mathbf{M}^T\mathbf{C}^T)\mathbf{S} + \boldsymbol{\varepsilon}. \quad (5.5)$$

If we let $\tilde{\mathbf{W}} = \mathbf{W}\mathbf{M}^{-T}$ and $\tilde{\mathbf{C}} = \mathbf{C}\mathbf{M}$, then Equation (5.5) is equivalent to

$$\mathbf{Y} = \exp(-\mathbf{A}\tilde{\mathbf{W}}\tilde{\mathbf{C}}^T)\mathbf{S} + \boldsymbol{\varepsilon}. \quad (5.6)$$

So far, we have not introduced how to choose the preconditioner \mathbf{M} . The choice of \mathbf{M} depends on linearization and approximation. In Section (5.2.1), we will state the process in detail, and in the new coordinate system defined by \mathbf{M} , the corresponding Hessian will be better-conditioned. With the help of the preconditioning matrix \mathbf{M} , we have transformed the original system of solving \mathbf{W} into the new system of solving $\tilde{\mathbf{W}}$. Since each entry of $\tilde{\mathbf{W}}$ is a linear combination of all entries in the corresponding row of \mathbf{W} , we can try to find a matrix \mathbf{M} such that the new system is better-conditioned than the original one.

On the other hand, we do not want to solve the nonlinear matrix equation (5.6) directly because it might introduce tensors when we compute second order derivatives. In this case, we want to vectorize Equation (5.6) on both sides and linearize it to construct a weighted least squares optimization problem. In the forward problem, we use the full spectrum and the matrix \mathbf{S} is usually rectangular. When we solve the inverse problem, we choose the average in each energy window to represent the corresponding energy spectrum. In this case, $N_e = N_b$ and the matrix \mathbf{S} in the inverse problem is a nonsingular diagonal matrix. So we can multiply \mathbf{S}^{-1} on both sides of Equation (5.6):

$$\mathbf{Y}\mathbf{S}^{-1} = \exp(-\mathbf{A}\tilde{\mathbf{W}}\tilde{\mathbf{C}}^T) + \boldsymbol{\varepsilon}\mathbf{S}^{-1}. \quad (5.7)$$

Vectorizing both sides of (5.7), and using properties of Kronecker products, we obtain

$$(\mathbf{S}^{-T} \otimes \mathbf{I}) \mathbf{y} = \exp \left\{ - \left(\tilde{\mathbf{C}} \otimes \mathbf{A} \right) \tilde{\mathbf{w}} \right\} + (\mathbf{S}^{-T} \otimes \mathbf{I}) \mathbf{e}, \quad (5.8)$$

where $\mathbf{y} = \text{vec}(\mathbf{Y})$, $\tilde{\mathbf{w}} = \text{vec}(\tilde{\mathbf{W}})$ and $\mathbf{e} = \text{vec}(\mathbf{E})$. If we let $\tilde{\mathbf{y}} = (\mathbf{S}^{-T} \otimes \mathbf{I}) \mathbf{y}$ and $\tilde{\mathbf{e}} = (\mathbf{S}^{-T} \otimes \mathbf{I}) \mathbf{e}$, then we can subtract $\tilde{\mathbf{e}}$ on both sides of (5.8) and obtain

$$\tilde{\mathbf{y}} - \tilde{\mathbf{e}} = \exp \left\{ - \left(\tilde{\mathbf{C}} \otimes \mathbf{A} \right) \tilde{\mathbf{w}} \right\}. \quad (5.9)$$

By taking the logarithm on both sides of Equation (5.9), we can obtain an equation:

$$\log(\tilde{\mathbf{y}} - \tilde{\mathbf{e}}) = - \left(\tilde{\mathbf{C}} \otimes \mathbf{A} \right) \tilde{\mathbf{w}}. \quad (5.10)$$

However, the left-hand side of Equation (5.10) contains the transformed error term $\tilde{\mathbf{e}}$ so we cannot solve this equation directly. In this case, we can separate the error term $\tilde{\mathbf{e}}$ from $\tilde{\mathbf{y}}$ using a first order Taylor expansion of the logarithm at $\tilde{\mathbf{y}}$:

$$\log(\tilde{\mathbf{y}} - \tilde{\mathbf{e}}) = \log(\tilde{\mathbf{y}}) - \text{diag}(\tilde{\mathbf{y}})^{-1} \tilde{\mathbf{e}} + \mathcal{O}(\|\tilde{\mathbf{e}}\|_2^2). \quad (5.11)$$

If we use the first two terms on the right-hand side of Equation (5.11) to estimate the term $\log(\tilde{\mathbf{y}} - \tilde{\mathbf{e}})$, then Equation (5.10) can be rewritten as a linear equation with the error term $\text{diag}(\tilde{\mathbf{y}})^{-1} \tilde{\mathbf{e}}$. Let $\mathbf{b} = -\log(\tilde{\mathbf{y}})$, then Equation (5.10) is equivalent to

$$\mathbf{b} \approx \left(\tilde{\mathbf{C}} \otimes \mathbf{A} \right) \tilde{\mathbf{w}} - \text{diag}(\tilde{\mathbf{y}})^{-1} \tilde{\mathbf{e}}. \quad (5.12)$$

With this equation and the Gaussian assumption of noise, $\mathbf{e} \sim \mathcal{N}(\mathbf{0}, \text{diag}(\mathbf{y}))$, we have

$$\mathbf{b} | \tilde{\mathbf{w}} \sim \mathcal{N} \left(\left(\tilde{\mathbf{C}} \otimes \mathbf{A} \right) \tilde{\mathbf{w}}, \Sigma \right), \quad (5.13)$$

where the noise covariance matrix Σ is expressed as

$$\Sigma = \text{diag}(\tilde{\mathbf{y}})^{-1} (\mathbf{S}^{-T} \otimes \mathbf{I}) \text{diag}(\mathbf{y}) (\mathbf{S}^{-1} \otimes \mathbf{I}) \text{diag}(\tilde{\mathbf{y}})^{-1}, \quad (5.14)$$

and the inverse noise covariance matrix is given by

$$\Sigma^{-1} = \text{diag}(\tilde{\mathbf{y}}) (\mathbf{S} \otimes \mathbf{I}) \text{diag}(\mathbf{y})^{-1} (\mathbf{S}^T \otimes \mathbf{I}) \text{diag}(\tilde{\mathbf{y}}). \quad (5.15)$$

Since $\mathbf{y} = \text{vec}(\mathbf{Y})$ and \mathbf{Y} is a matrix that collects the number of photons of each energy window in the corresponding column, each entry of \mathbf{Y} is a positive integer whose value can be on the order of hundreds of thousands. Moreover, $\tilde{\mathbf{y}} = (\mathbf{S}^{-T} \otimes \mathbf{I}) \mathbf{y}$ and as long as the noise does not dominate the projected data, we expect the entries of $\tilde{\mathbf{y}}$ to be larger than zero. From the expression (5.15), we can see that the structure of Σ^{-1} depends on the structure of matrix \mathbf{S} . If \mathbf{S} is diagonal, then Σ is also diagonal. Otherwise, Σ^{-1} is a block diagonal matrix. If we let $\mathcal{A} = \tilde{\mathbf{C}} \otimes \mathbf{A}$ and ignore constants, the corresponding probability density function is given by

$$f(\mathbf{b}; \tilde{\mathbf{w}}) = \exp \left\{ -\frac{1}{2} (\mathcal{A}\tilde{\mathbf{w}} - \mathbf{b})^T \Sigma^{-1} (\mathcal{A}\tilde{\mathbf{w}} - \mathbf{b}) \right\}. \quad (5.16)$$

So the log-likelihood function is given by

$$l(\tilde{\mathbf{w}}; \mathbf{b}) = -\frac{1}{2} (\mathcal{A}\tilde{\mathbf{w}} - \mathbf{b})^T \Sigma^{-1} (\mathcal{A}\tilde{\mathbf{w}} - \mathbf{b}). \quad (5.17)$$

We try to maximize the log-likelihood function $l(\tilde{\mathbf{w}}; \mathbf{b}, \Sigma)$, which is equivalent to minimizing the negative log-likelihood function $-l(\tilde{\mathbf{w}}; \mathbf{b}, \Sigma)$. In addition, we require that $\mathbf{W} \geq \mathbf{0}$, and with the preconditioner, these constraints are transformed into $(\mathbf{M} \otimes \mathbf{I}) \tilde{\mathbf{w}} \geq \mathbf{0}$. Therefore, we can formulate a weighted least squares problem

under bound constraints

$$\begin{aligned} \min_{\tilde{\mathbf{w}}} \quad & \frac{1}{2} \|\mathbf{A}\tilde{\mathbf{w}} - \mathbf{b}\|_{\Sigma^{-1}}^2 \\ \text{subject to} \quad & (\mathbf{M} \otimes \mathbf{I}) \tilde{\mathbf{w}} \geq \mathbf{0}. \end{aligned} \tag{5.18}$$

In Equation (5.18), the norm $\|\cdot\|_{\Sigma^{-1}}^2$ corresponds to the negative log-likelihood function $-l(\tilde{\mathbf{w}}; \mathbf{b}, \Sigma)$. From this expression, we know that the objective function is convex. Moreover, the inverse covariance matrix Σ^{-1} is diagonal as long as \mathbf{S} is diagonal and this optimization problem has linear inequality constraints. Based on these observations, we can identify four challenges involved in solving this optimization problem. At first, we need to choose an appropriate preconditioning matrix to reduce the ill-conditioning of the Hessian. Secondly, we want to choose suitable regularizations for the corresponding materials. Thirdly, we have to find an efficient method to solve the weighted least squares problem. These three challenges are related to each other and an appropriate preconditioner with feasible regularizations will be beneficial for the solver efficiency. Finally, we should handle linear inequality constraints in an efficient way. We address these four challenges in the following sections.

5.2 Preconditioning and Regularization

5.2.1 Preconditioning

The choice of the preconditioning matrix \mathbf{M} is crucial for solving the optimization problem (5.18). If we do not have a preconditioner or we choose the preconditioner \mathbf{M} as identity, the original Hessian for the weighted least squares problem is expressed as

$$\mathbf{H} = (\mathbf{C}^T \otimes \mathbf{A}^T) \Sigma^{-1} (\mathbf{C} \otimes \mathbf{A}). \tag{5.19}$$

An appropriate preconditioner can transform the original ill-posed system into a better-conditioned system and thus bring faster convergence speed as well as higher quality of reconstructed images. In general, the preconditioned Hessian $\tilde{\mathbf{H}}$ can be represented as

$$\tilde{\mathbf{H}} = \mathcal{A}^T \Sigma^{-1} \mathcal{A} = \left(\tilde{\mathbf{C}}^T \otimes \mathbf{A}^T \right) \Sigma^{-1} \left(\tilde{\mathbf{C}} \otimes \mathbf{A} \right), \quad (5.20)$$

where $\tilde{\mathbf{C}} = \mathbf{C}\mathbf{M}$. From this expression, it is still not obvious how to construct the preconditioner. However, if we can decompose the noise covariance matrix Σ^{-1} into a Kronecker product of two terms, then we can merge several terms using the properties of Kronecker product and transform parts of the products into identity. To realize this idea, we review the expression of Σ^{-1} in Equation (5.15), where we can see that it contains the Kronecker products $\mathbf{S} \otimes \mathbf{I}$ and $\mathbf{S}^T \otimes \mathbf{I}$ and it is not necessary to separate these two terms. So we focus on the other terms that include $\text{diag}\{\tilde{\mathbf{y}}\}$ and $\text{diag}\{\mathbf{y}\}^{-1}$. By definition, these two terms are related to each other by $\tilde{\mathbf{y}} = (\mathbf{S}^{-T} \otimes \mathbf{I}) \mathbf{y}$. In this case, if we can express $\text{diag}\{\mathbf{y}\}$ into a Kronecker product of two terms, then we will reach the goal.

Recall that $\mathbf{y} = \text{vec}(\mathbf{Y})$. Therefore, if we can find two rank-1 matrices, \mathbf{u} and \mathbf{v} , such that $\mathbf{Y} \approx \mathbf{u}\mathbf{v}^T$, then

$$\text{diag}\{\mathbf{y}\} \approx \text{diag}\{\text{vec}(\mathbf{u}\mathbf{v}^T)\} = \text{diag}\{\mathbf{v}\} \otimes \text{diag}\{\mathbf{u}\}. \quad (5.21)$$

These two rank-1 matrices can be obtained by solving a nearest Kronecker product (NKP) problem, which is equivalent to the rank-1 approximation of \mathbf{Y} in terms of the Frobenius norm.

$$\min_{\mathbf{u}, \mathbf{v}} \|\mathbf{Y} - \mathbf{u}\mathbf{v}^T\|_F. \quad (5.22)$$

The solution to this problem is similar to the one in Chapter 4, but it has several variations. Using the singular value decomposition (SVD), one solution to Problem (5.22) can be expressed by $\mathbf{u} = \sqrt{\sigma_1} \mathbf{u}_1$ and $\mathbf{v} = \sqrt{\sigma_1} \mathbf{v}_1$, where \mathbf{u}_1 and \mathbf{v}_1 are the

first left and right singular vectors and σ_1 is the corresponding largest singular value. Since we only need these terms rather than a full SVD, we can use MATLAB's `svds` function, or other efficient approaches, such as "PROPACK" [39], to compute only σ_1 , \mathbf{u}_1 and \mathbf{v}_1 .

After we have obtained \mathbf{u} and \mathbf{v} , we can estimate the matrix $\text{diag}\{\mathbf{y}\}$ as a Kronecker product of two terms as Equation (5.21). In addition, the term $\text{diag}\{\tilde{\mathbf{y}}\}$ can be represented as

$$\begin{aligned}
\text{diag}\{\tilde{\mathbf{y}}\} &= \text{diag}\{(\mathbf{S}^{-T} \otimes \mathbf{I}) \text{vec}(\mathbf{Y})\} \\
&\approx \text{diag}\{(\mathbf{S}^{-T} \otimes \mathbf{I}) \text{vec}(\mathbf{u}\mathbf{v}^T)\} \\
&= \text{diag}\{\text{vec}(\mathbf{u}\mathbf{v}^T \mathbf{S}^{-1})\} \\
&= \text{diag}\{\mathbf{S}^{-T} \mathbf{v}\} \otimes \text{diag}\{\mathbf{u}\}.
\end{aligned} \tag{5.23}$$

If we substitute the terms in (5.21) and (5.23) for the same terms in (5.15), then we can obtain that

$$\boldsymbol{\Sigma}^{-1} \approx (\text{diag}\{\mathbf{S}^{-T} \mathbf{v}\} \mathbf{S} \text{diag}\{\mathbf{v}\}^{-1} \mathbf{S}^T \text{diag}\{\mathbf{S}^{-T} \mathbf{v}\}) \otimes \text{diag}\{\mathbf{u}\}. \tag{5.24}$$

So the preconditioned Hessian matrix is given by

$$\begin{aligned}
\tilde{\mathbf{H}} &= (\tilde{\mathbf{C}}^T \otimes \mathbf{A}^T) \boldsymbol{\Sigma}^{-1} (\tilde{\mathbf{C}} \otimes \mathbf{A}) \\
&\approx (\tilde{\mathbf{C}}^T \otimes \mathbf{A}^T) (\text{diag}\{\mathbf{S}^{-T} \mathbf{v}\} \mathbf{S} \text{diag}\{\mathbf{v}\}^{-1} \mathbf{S}^T \text{diag}\{\mathbf{S}^{-T} \mathbf{v}\} \otimes \text{diag}\{\mathbf{u}\}) (\tilde{\mathbf{C}} \otimes \mathbf{A}) \\
&= (\tilde{\mathbf{C}}^T \text{diag}\{\mathbf{S}^{-T} \mathbf{v}\} \mathbf{S} \text{diag}\{\mathbf{v}\}^{-1} \mathbf{S}^T \text{diag}\{\mathbf{S}^{-T} \mathbf{v}\} \tilde{\mathbf{C}}) \otimes (\mathbf{A}^T \text{diag}\{\mathbf{u}\} \mathbf{A}).
\end{aligned} \tag{5.25}$$

Since the size of $\tilde{\mathbf{C}}$ is $N_e \times N_m$, then the first part of the Kronecker product in (5.25) is a square matrix of the size N_m . In other words, this part only depends on the number of materials that compose the object. Usually, we only consider 2 or 3 materials to form

the object so that the size of this part is usually either 2×2 or 3×3 . Moreover, the matrix \mathbf{Y} gathers the number of photons of each energy window in the corresponding column so all of its entries are positive integers. In this case, we can choose \mathbf{u} and \mathbf{v} to be positive such that the matrix product, $\mathbf{C}^T \text{diag} \{ \mathbf{S}^{-T} \mathbf{v} \} \mathbf{S} \text{diag} \{ \mathbf{v} \}^{-1} \mathbf{S}^T \text{diag} \{ \mathbf{S}^{-T} \mathbf{v} \} \mathbf{C}$, is a symmetric positive definite (SPD) matrix. Therefore, we can calculate \mathbf{M} using the Cholesky decomposition:

$$\mathbf{C}^T \text{diag} \{ \mathbf{S}^{-T} \mathbf{v} \} \mathbf{S} \text{diag} \{ \mathbf{v} \}^{-1} \mathbf{S}^T \text{diag} \{ \mathbf{S}^{-T} \mathbf{v} \} \mathbf{C} = \mathbf{G}^T \mathbf{G}, \quad (5.26)$$

where \mathbf{G} is an upper triangular matrix with positive diagonal entries. Since $\tilde{\mathbf{C}} = \mathbf{C} \mathbf{M}$, we can choose $\mathbf{M} = \mathbf{G}^{-1}$ to transform this part into identity. From Expression (5.25), we see that the preconditioned Hessian matrix $\tilde{\mathbf{H}}$ is dependent on a Kronecker product of two parts and the first part has been transformed into an identity matrix. In particular, since the condition number of this part is typically significantly greater than 1, the condition number of the preconditioned Hessian $\tilde{\mathbf{H}}$ is significantly smaller than the original Hessian \mathbf{H} .

After we have obtained the matrix \mathbf{M} , we can analyze the effect of preconditioning using the SVD. Without preconditioning, the Hessian matrix \mathbf{H} depends on two parts, $\mathbf{C}^T \text{diag} \{ \mathbf{S}^{-T} \mathbf{v} \} \mathbf{S} \text{diag} \{ \mathbf{v} \}^{-1} \mathbf{S}^T \text{diag} \{ \mathbf{S}^{-T} \mathbf{v} \} \mathbf{C}$ and $\mathbf{A}^T \text{diag} \{ \mathbf{u} \} \mathbf{A}$. If we assume that the singular value decomposition for these two matrices are $\mathbf{U}_1 \mathbf{\Sigma}_1 \mathbf{V}_1^T$ and $\mathbf{U}_2 \mathbf{\Sigma}_2 \mathbf{V}_2^T$, then the condition number of the original Hessian \mathbf{H} is closely related to $\mathbf{\Sigma}_1$ and $\mathbf{\Sigma}_2$. Let the largest and smallest singular values of $\mathbf{\Sigma}_1$ and $\mathbf{\Sigma}_2$ be σ_{1max} , σ_{1min} , σ_{2max} and σ_{2min} , respectively, then the condition number of the original Hessian, $\kappa(\mathbf{H})$, can be estimated as

$$\kappa(\mathbf{H}) \approx \frac{\sigma_{1max} \sigma_{2max}}{\sigma_{1min} \sigma_{2min}}. \quad (5.27)$$

On the other hand, the condition number of the preconditioned Hessian can be approximated as

$$\kappa(\tilde{\mathbf{H}}) \approx \frac{\sigma_{2max}}{\sigma_{2min}}. \quad (5.28)$$

Since the fraction $\sigma_{1max}/\sigma_{1min}$ is most likely to be significantly greater than 1, then the condition number of $\tilde{\mathbf{H}}$ is likely to be much smaller than \mathbf{H} . Furthermore, we can build a numerical example to validate this phenomenon. For an object that is composed of two materials and each material map is of the size 16×16 , we can construct the original Hessian \mathbf{H} and the preconditioned Hessian $\tilde{\mathbf{H}}$ explicitly and compute the estimations of condition numbers for these two Hessian matrices. The result is presented in Table 5.1. From Table 5.1, we can see that the difference of condition number for \mathbf{H} and $\tilde{\mathbf{H}}$ is around

Matrix Types	Condition Numbers
Original Hessian	2.0042e+06
Preconditioned Hessian	2.5874e+04

Table 5.1: The comparison of condition numbers.

two orders of magnitude, which indicates the significance of this preconditioner. For a linear system that involves the preconditioned Hessian $\tilde{\mathbf{H}}$, the convergence rate is highly dependent on the condition number. So we can solve the preconditioned system in a more efficient way. Moreover, we will validate the strength of this preconditioner by solving the preconditioned system versus the original system. More details are presented in section 5.4.

5.2.2 Regularization

With the help of our preconditioner, we can speed up an optimization algorithm and achieve higher accuracy. To further alleviate noise amplification, it is important to add regularization terms to the objective function. In total, we have m materials and the weights of these m materials are not equal. Rather than adding a single regularization to all weights, we instead add a specific regularization to each material. In addition, for different materials, we can choose distinct regularizations to match their properties. For the dominant material, we select generalized Tikhonov regularization to smooth the edges. For other materials, we choose ℓ_1 regularizations to penalize the sum of weights. Based on this idea, we can

represent the regularization term as a sum of m parts:

$$R(\mathbf{w}) = \sum_{i=1}^m \frac{\alpha_i}{2} R_i(\mathbf{w}_i), \quad (5.29)$$

where \mathbf{w}_i is the i -th column of the weight matrix \mathbf{W} , $R_i(\mathbf{w}_i)$ is the corresponding regularization term and α_i is the regularization parameter.

The choice of what type of regularization to use is problem-specific, and *a priori* knowledge of the object being imaged could inform this decision. For example, if it is known that the object contains two material maps with relatively equal distributions, we might select two generalized Tikhonov regularizations. For example, in breast imaging the object is dominated by glandular and adipose tissue, and so if the aim is to determine the weights of these two materials, it might make sense to use a generalized Tikhonov regularization for each of them. On the other hand, it could be the case that the object is dominated by one material (or one set of materials), with a relatively sparse distribution of another material. For example, in the breast imaging situation, the object may contain small microcalcifications or areas highlighted by an iodine tracer. In this case, one can use generalized Tikhonov regularization for the dominating materials (e.g., glandular and adipose tissue) and a ℓ_1 regularization for the sparse material. We illustrate this with two materials, one that dominates, and one that is sparse:

$$R(\mathbf{w}) = \frac{\alpha_1}{2} \|\mathbf{L}\mathbf{w}_1\|_2^2 + \frac{\alpha_2}{2} \|\mathbf{w}_2\|_1, \quad (5.30)$$

where \mathbf{L} is the specified discrete differential operator. If we add these regularization terms to the objective function in Equation (5.18), we can rewrite it as an augmented system:

$$\min_{\tilde{\mathbf{w}}} \left\| \begin{bmatrix} \frac{\sqrt{2}}{2} \boldsymbol{\Sigma}^{-\frac{1}{2}} (\tilde{\mathbf{C}} \otimes \mathbf{A}) \\ \sqrt{\frac{\alpha_1}{2}} \tilde{\mathbf{L}} \end{bmatrix} \tilde{\mathbf{w}} - \begin{bmatrix} \boldsymbol{\Sigma}^{-\frac{1}{2}} \mathbf{b} \\ \mathbf{0} \end{bmatrix} \right\|_2^2 + \frac{\alpha_2}{2} \begin{bmatrix} \mathbf{0} & \mathbf{1} \end{bmatrix} (\mathbf{M} \otimes \mathbf{I}) \tilde{\mathbf{w}} \quad (5.31)$$

subject to $(\mathbf{M} \otimes \mathbf{I}) \tilde{\mathbf{w}} \geq \mathbf{0}$,

where $\tilde{\mathbf{L}} = \begin{bmatrix} \mathbf{L} & \mathbf{0} \end{bmatrix} (\mathbf{M} \otimes \mathbf{I})$. As we can see, the objective function in this problem consists

of two parts: one is smooth and convex and the other one is possibly nonsmooth. Because of these properties, we can think about using FISTA [3] to solve this problem. It not only fits the features of the objective function but also provides an optimal convergence rate. In addition, we are concerned about the linear inequality constraints, and in each step, we can maintain these constraints by solving a projection problem that is based on the 2-norm.

5.3 FISTA and Projections

In this section, we first present the main algorithm FISTA briefly. To implement FISTA to solve the target optimization problem, we need to decide the step size and handle the non-negative constraints. For the step size, we introduce how to compute the Lipschitz constant numerically then we choose a constant step size based on the calculated Lipschitz constant. For the nonnegative constraints, we build another quadratic programming problem and solve it with delicate decomposition and efficient algorithms.

5.3.1 FISTA

Fast Iterative Shrinkage-Thresholding Algorithm (FISTA) is a first order method that belongs to the family of Iterative Shrinkage-Thresholding Algorithm (ISTA) [13]. This method is proposed by Beck et al., and compared with the $\mathcal{O}(1/k)$ rate of convergence of ISTA, it has a best function value convergence rate $\mathcal{O}(1/k^2)$, where k is the number of iterations. In most situations, FISTA is considered to be the “optimal” first order method with respect to convergence speed. Moreover, it best fits the problems in imaging science because it is usually used to solve the nonsmooth convex problem

$$\min_{\mathbf{x}} f(\mathbf{x}) + g(\mathbf{x}), \quad (5.32)$$

where $f(\mathbf{x})$ and $g(\mathbf{x})$ are both convex functions and $g(\mathbf{x})$ might not be smooth. In imaging sciences, $f(\mathbf{x})$ is likely to be a least squares loss function to test the goodness of fit and $g(\mathbf{x})$ can be a regularization term such as ℓ_1 penalty or total variation regularization.

For Problem (5.31), we construct an augmented loss function that merges the generalized Tikhonov regularization term, which corresponds to $f(\mathbf{x})$ in (5.32). For the regularization term, ℓ_1 regularization is nonsmooth but convex and this matches $g(\mathbf{x})$ in (5.32).

The details of this algorithm are shown in Algorithm 4. For the main algorithm, we need to calculate the smallest Lipschitz constant K at first. Then we can update the current step using FISTA. Because of the linear inequality constraints, we need to project the new step onto boundaries to keep the solution feasible. We would like to implement FISTA with a constant step size to solve the optimization problem (5.31). To implement this method, we need several preparations.

Algorithm 4 FISTA and Projections [3]

- 1: *Initialization:*
 - 2: Calculate the smallest Lipschitz constant K in (5.34) by the power method.
 - 3: Set up initial guess $\tilde{\mathbf{W}}_0$; Let $\mathbf{y}_0 = \text{vec}(\tilde{\mathbf{W}}_0)$, $\mathbf{x}_{old} = \mathbf{y}_0$ and $t_1 = 1$;
 - 4: **for** $k = 1, 2, \dots$ **do**
 - 5: Calculate the gradients, $\nabla f(\mathbf{y}_k)$ and $\nabla g(\mathbf{y}_k)$, of $f(\mathbf{y}_k)$ and $g(\mathbf{y}_k)$ in (5.33);
 - 6: $\mathbf{x}_k = \mathbf{y}_k - \frac{1}{L(f)} [\nabla f(\mathbf{y}_k) + \nabla g(\mathbf{y}_k)]$;
 - 7: Reshape \mathbf{x}_k into a matrix and use CVXGEN to solve the projection problems to obtain \mathbf{x}_{new} as (5.37);
 - 8: $t_{k+1} = \frac{1 + \sqrt{1 + 4t_k^2}}{2}$;
 - 9: $\mathbf{y}_{k+1} = \mathbf{x}_{new} + \left(\frac{t_k - 1}{t_{k+1}}\right) (\mathbf{x}_{new} - \mathbf{x}_{old})$;
 - 10: $\mathbf{x}_{old} = \mathbf{x}_{new}$.
-

5.3.2 Lipschitz Constant

The first step is to calculate the smallest Lipschitz constant. If we let

$$f(\tilde{\mathbf{w}}) = \left\| \begin{bmatrix} \frac{\sqrt{2}}{2} \Sigma^{-\frac{1}{2}} (\tilde{\mathbf{C}} \otimes \mathbf{A}) \\ \sqrt{\frac{\alpha_1}{2}} \tilde{\mathbf{L}} \end{bmatrix} \tilde{\mathbf{w}} - \begin{bmatrix} \Sigma^{-\frac{1}{2}} \mathbf{b} \\ \mathbf{0} \end{bmatrix} \right\|_2^2, \quad (5.33)$$

$$g(\tilde{\mathbf{w}}) = \frac{\alpha_2}{2} \begin{bmatrix} \mathbf{0} & \mathbf{1} \end{bmatrix} (\mathbf{M} \otimes \mathbf{I}) \tilde{\mathbf{w}},$$

then we need the smallest Lipschitz constant K for $\nabla f(\tilde{\mathbf{w}})$, which is the largest eigenvalue for $\nabla^2 f(\tilde{\mathbf{w}})$. That is to say,

$$K = \lambda_{\max} \left[\left(\tilde{\mathbf{C}}^T \otimes \mathbf{A}^T \right) \boldsymbol{\Sigma}^{-1} \left(\tilde{\mathbf{C}} \otimes \mathbf{A} \right) + \alpha_1 \tilde{\mathbf{L}}^T \tilde{\mathbf{L}} \right]. \quad (5.34)$$

Since we only need the largest eigenvalue, it is not necessary for us to construct these matrices explicitly; instead we can use an iterative method, such as the power method [24]. Note that we only need to compute K once for all FISTA iterations. The details of the power method are shown in Algorithm 5.

Algorithm 5 Power Method [24]

- 1: *Initialization:*
 - 2: Generate a random vector \mathbf{q}_0 and normalize \mathbf{q}_0 ;
 - 3: **for** $i = 1, 2, \dots$ **do**
 - 4: $\mathbf{z}_i = \left[\left(\tilde{\mathbf{C}}^T \otimes \mathbf{A}^T \right) \boldsymbol{\Sigma}^{-1} \left(\tilde{\mathbf{C}} \otimes \mathbf{A} \right) + \alpha_1 \tilde{\mathbf{L}}^T \tilde{\mathbf{L}} \right] \mathbf{q}_{i-1}$;
 - 5: $\mathbf{q}_i = \mathbf{z}_i / \|\mathbf{z}_i\|_2$;
 - 6: $\lambda_i = \mathbf{q}_i^T \left[\left(\tilde{\mathbf{C}}^T \otimes \mathbf{A}^T \right) \boldsymbol{\Sigma}^{-1} \left(\tilde{\mathbf{C}} \otimes \mathbf{A} \right) + \alpha_1 \tilde{\mathbf{L}}^T \tilde{\mathbf{L}} \right] \mathbf{q}_i$;
-

5.3.3 Projections

In addition to the Lipschitz constant, we also need to handle the linear inequality constraints $(\mathbf{M} \otimes \mathbf{I}) \tilde{\mathbf{w}} \geq \mathbf{0}$. Generally speaking, we can regard Problem (5.31) as a quadratic programming problem under these specific constraints. To impose the linear inequality constraints, we can construct another quadratic programming problem that can find the nearest solution to satisfy these constraints. If we assume that we have obtained $\tilde{\mathbf{w}}_k$ in the k -th step, then we can build a projection problem:

$$\begin{aligned} \min_{\tilde{\mathbf{w}}_{new}} \quad & \|\tilde{\mathbf{w}}_{new} - \tilde{\mathbf{w}}_k\|_2^2 \\ \text{subject to} \quad & (\mathbf{M} \otimes \mathbf{I}) \tilde{\mathbf{w}}_{new} \geq \mathbf{0}. \end{aligned} \quad (5.35)$$

For small and medium size problems, we can solve it efficiently by direct implementation of standard optimization algorithms. For example, we can use CVX [26, 27] to solve Prob-

lem (5.35), which turns to be low-cost both in storage and computational consumptions. However, there are challenges for large-scale problems. For example, saving long vectors or constructing sparse matrices might require large storage space. Therefore, we should find a method to decompose Problem (5.35) into small pieces and try to solve each small-sized problem accurately and efficiently.

Suppose we reshape vectors into matrices, for example using MATLAB's "reshape" function, $\tilde{\mathbf{W}}_{new} = \text{reshape}(\tilde{\mathbf{w}}_{new}, N_v, N_m)$ and $\tilde{\mathbf{W}}_k = \text{reshape}(\tilde{\mathbf{w}}_k, N_v, N_m)$, then by Kronecker product properties and the connection between the 2-norm and the Frobenius norm, Problem (5.35) is equivalent to

$$\begin{aligned} \min_{\tilde{\mathbf{W}}_{new}} \quad & \left\| \tilde{\mathbf{W}}_{new} - \tilde{\mathbf{W}}_k \right\|_F^2 \\ \text{subject to} \quad & \tilde{\mathbf{W}}_{new} \mathbf{M}^T \geq \mathbf{0}. \end{aligned} \quad (5.36)$$

If we focus on each row of $\tilde{\mathbf{W}}_k$, $\tilde{\mathbf{W}}_k(i, :)$, then Problem (5.36) can be rewritten as

$$\begin{aligned} \min_{\tilde{\mathbf{W}}_{new}} \quad & \sum_{i=1}^{N_v} \left\| \tilde{\mathbf{W}}_{new}(i, :) - \tilde{\mathbf{W}}_k(i, :) \right\|_2^2 \\ \text{subject to} \quad & \tilde{\mathbf{W}}_{new}(i, :) \mathbf{M}^T \geq \mathbf{0}, \quad i = 1, 2, \dots, N_v, \end{aligned} \quad (5.37)$$

where $\tilde{\mathbf{W}}_{new}(i, :)$ is the i -th row of $\tilde{\mathbf{W}}_{new}$. It is obvious that this problem is separable, and the original problem (5.36) can be separated into small-sized problems that only involve each row of $\tilde{\mathbf{W}}_{new}$ and $\tilde{\mathbf{W}}_k$. Since each row only depends on the number of materials N_m , then the size of problem is usually 1×2 or 1×3 . In this case, we can solve each small-sized problem efficiently and concatenate the solutions into a large matrix later. To realize this idea, we can find a highly efficient solver for small-sized problems and loop around the number of voxels (pixels if 2D) N_v . In this chapter, we choose CVXGEN [41–44] to generate a customized solver for small quadratic programming problems. It is a problem-specific, fast and accurate code generator which can achieve advance performance in particular for small-sized quadratic programming problems. In addition, if computer clusters are available, we can write parallel programming codes, such as MPI or OpenMP, and compute the solution

to this projection problem in parallel. The speedup in this case relies on the number of available compute nodes, but clearly there is potential for significant speedup with such an approach.

In conclusion, we can see that this algorithm incorporates the advantages of the power method, FISTA and the fast solver, CVXGEN, for small problems. With the power method, we only need to save the Hessian-vector multiplication rather than the full Hessian, and it is very cheap to compute. Moreover, we can achieve a rapid convergence using FISTA in the main loop. Finally, the projection problem is decomposed into many small pieces and each can be solved by CVXGEN efficiently.

5.4 Numerical Experiments

To test the performance of our preconditioner and the main algorithm, we set up a test problem that is composed of two materials, plexiglass and polyvinyl chloride (PVC). The size of each material map is 128×128 . The first material map is a circular mask that dominates the object, while the second material map consists of small “spikes” that are scattered randomly inside the circle. The number of “spikes” is chosen to be 50. Outside of the circle, we assume that there exist no weights of the object. These two images are shown in Figure 5.1.

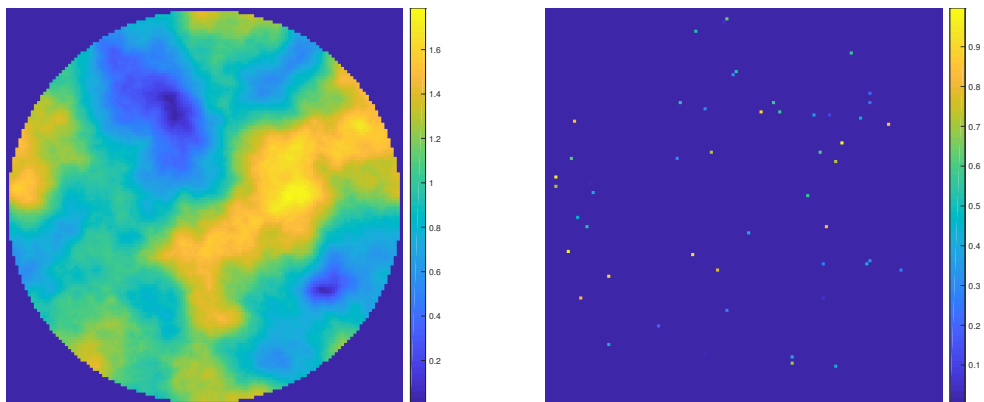


Figure 5.1: The original material maps for plexiglass (Left) and PVC (Right).

Inside the mask, the darker blue areas for the first material map are mainly located

in the upper left and lower right corners, which corresponds to blank points. Other areas inside the circle are represented by heavily weighted yellow and green color. In the second material map, the weights are scattered around the image and only occupy a small part of the area in total. This test problem can be regarded as a simplification of a real life application. For example, in medical imaging for cancer detection, the first material map is similar to a small area of human body or tissue, while the second material map can represent the calcium located inside this area.

In addition to the test images, we also need other parameters in Equation (5.1). To generate the ray trace matrix \mathbf{A} , we use the MATLAB function `fanbeamtomolinear` from AIR Tools [29,32,35] to simulate a fan-beam geometry with a flat detector. The parameters that we need to choose are presented in Table 4.1 in Chapter 4. In addition, we use 180 projections in total which are equally distributed from 0 to 360 degrees. The spectral energy of the x-ray source is generated by the MATLAB function `spektrSpectrum` [53] with 120 keV voltage as input. The detector is assumed to be photon-counting with 5 energy windows. From the first energy window to the fifth energy window, we assume that they can detect the range of photon energies 10 to 34 keV, 35 to 49 keV, 50 to 64 keV, 65 to 79 keV and 80 to 120 keV, respectively.

The plot of photon flux density versus photon energy is presented in Figure 4.4. In Figure 4.4, the red curve represents photon intensity of x-ray source and the gray boxes indicate energy windows of the detector. Moreover, the black dots are the values of mean photon energy in each energy window. When we build the test problem, the full energy spectrum and all the corresponding linear attenuation coefficients are used, while only the mean photon energies and the corresponding linear attenuation coefficients are applied for reconstruction. As it is well-known, this strategy of generating data on a finer grid and solve it on a coarser grid is a standard approach to avoiding what is called the inverse crime.

The curves of linear attenuation coefficients versus photon energies are presented in Figure 4.3 in Chapter 4. From Figure 4.3, we can see that the slopes of these two curves are close to each other, which are likely to introduce the collinearity between coefficients. Moreover, we assume that the entries of the matrix \mathbf{Y} follow a Poisson distribution, and for large

scale problems, from the Central Limit Theorem, the Poisson distribution is approximated well by a Gaussian distribution. So the assumption of Gaussian model is valid.

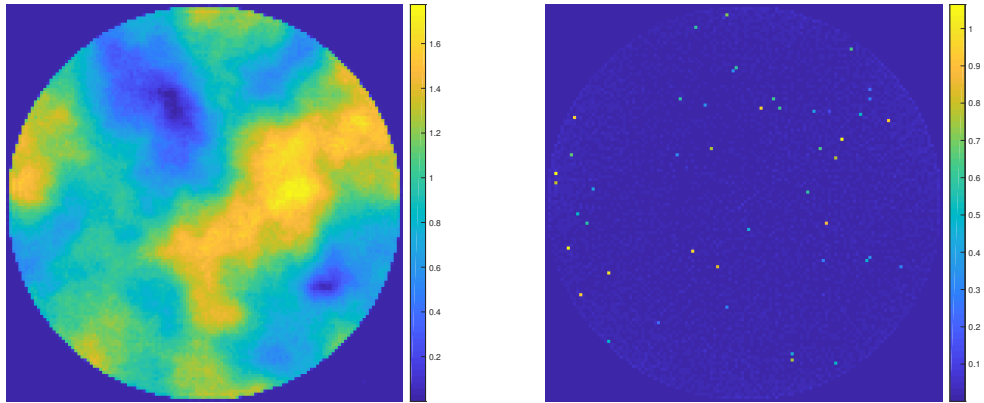


Figure 5.2: The reconstructed images for plexiglass (Left) and PVC (Right).

The reconstructed images are shown in Figure 5.2. From Figure 5.2, we can see that we achieve almost perfect separation for these two materials. Moreover, the reconstructed images have excellent quality in terms of visibility. Both two material maps are relatively close to the true images. In the first material map, the distribution of weights is clear to identify. The low intensity pixels are located in the upper left and lower right areas of the circle, while other places are occupied by the yellow and green colors. Moreover, we can easily recognize the edges of the circle that indicate the boundary of the object, which is a plus. As we can see, the reconstruction of small “spikes” are of great difficulty because of the randomness of weights and spots. However, we can see that the small “spikes” are scattered in the same positions as the true image, while they are masked by the shade of a circle. These results present the significance of the methods proposed in this chapter.

To further validate the results, we plot the relative errors of these two materials versus the number of FISTA iterations. The decrease of relative errors of corresponding materials is shown in Figure 5.3. From this figure, we can see that the relative error of the first material drops sharply as the number of iterations increase. It then stagnates after around 150 iterations. However, the relative error of the second material only decreases fast in the beginning, and after several iterations, the rate of change slows down and the relative error cannot reduce further. We can also identify the same phenomenon by comparing

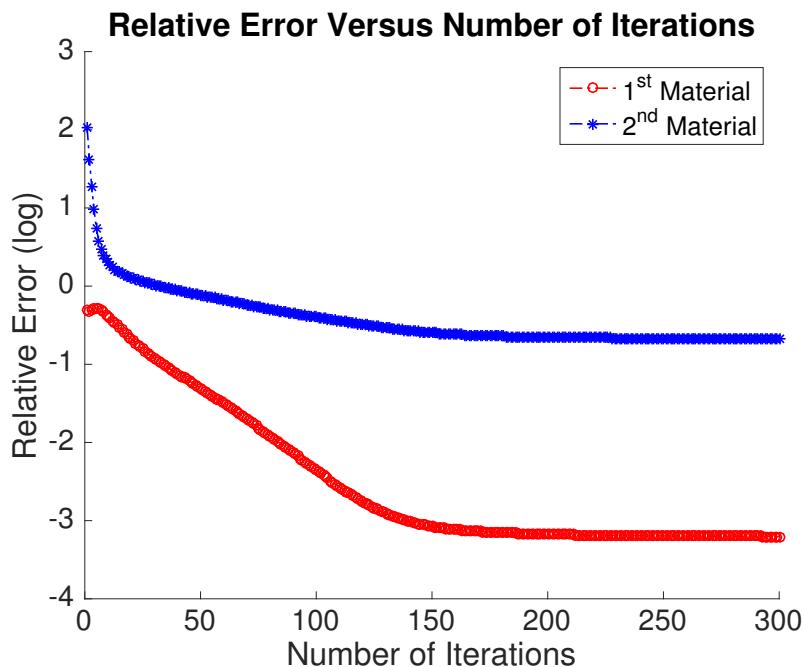


Figure 5.3: The related errors of each iteration (with preconditioner) for plexiglass and PVC.

the true and reconstructed images of the second material map. Even if the spots of these “spikes” are approximately correct, the numerical weights of these dots might not be the same. Moreover, there are a large number of small values in the background of reconstructed image, causing somewhat large relative errors, even though visually the result looks quite good.

Other accuracy measures illustrate this phenomenon. In Figure 5.4 we plot the mean squared error (MSE) at each iteration, in Figure 5.5 we plot the peak signal to noise ratio (PSNR), and in Figure 5.6 we plot the structural similarity index (SSIM). Not surprisingly the MSE produces information very similar to the relative errors, but it also shows a clear diminution for the second material from Figure 5.4. PSNR provides a similar measure to MSE, with an inverted interpretation (higher values correspond to better solutions). The SSIM is a metric for image quality, and as with PSNR, large values correspond to better solutions. From Figure 5.6, it can be found that the quality of the reconstructed first material map improves slowly in the early iterations but it achieves a higher quality measure in the end compared with the second material map. In summary, all of these errors

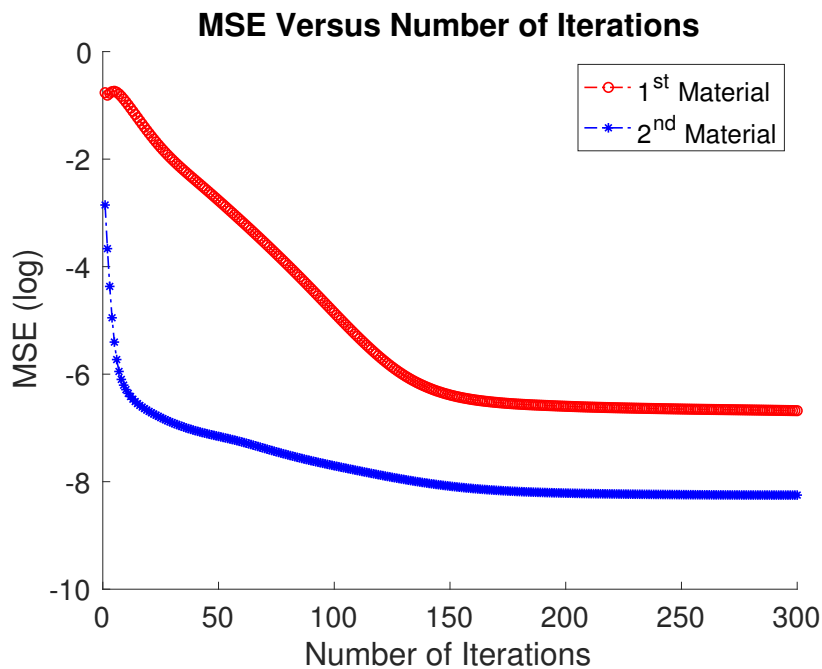


Figure 5.4: The MSE of each iteration (with preconditioner) for plexiglass and PVC.

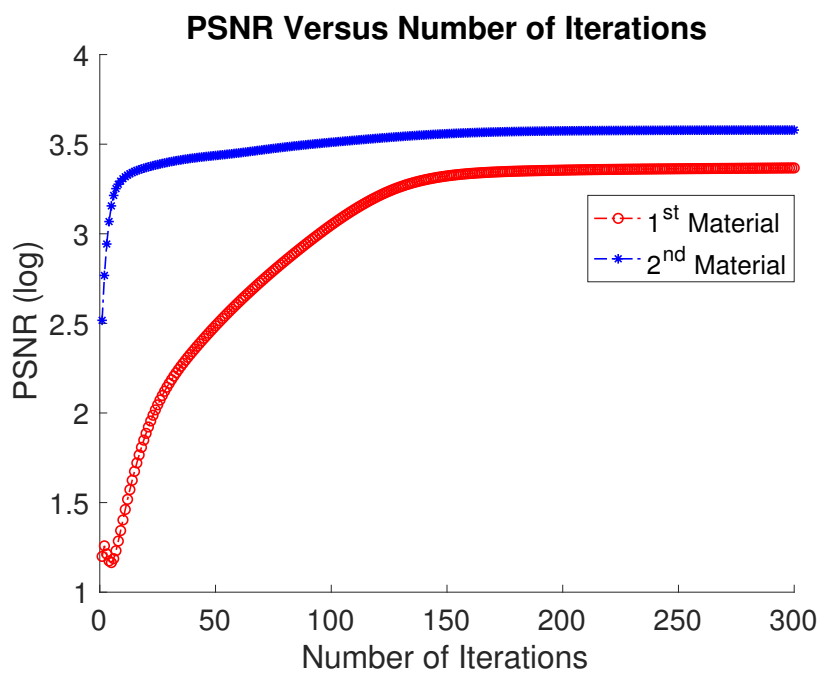


Figure 5.5: The PSNR of each iteration (with preconditioner) for plexiglass and PVC.

and quality measures illustrate fast convergence to high quality reconstructions.

It may also be of interest to observe the decay of norm of the gradient at each iteration, which is shown in Figure 5.7. From this figure we can see that the norm of the gradient

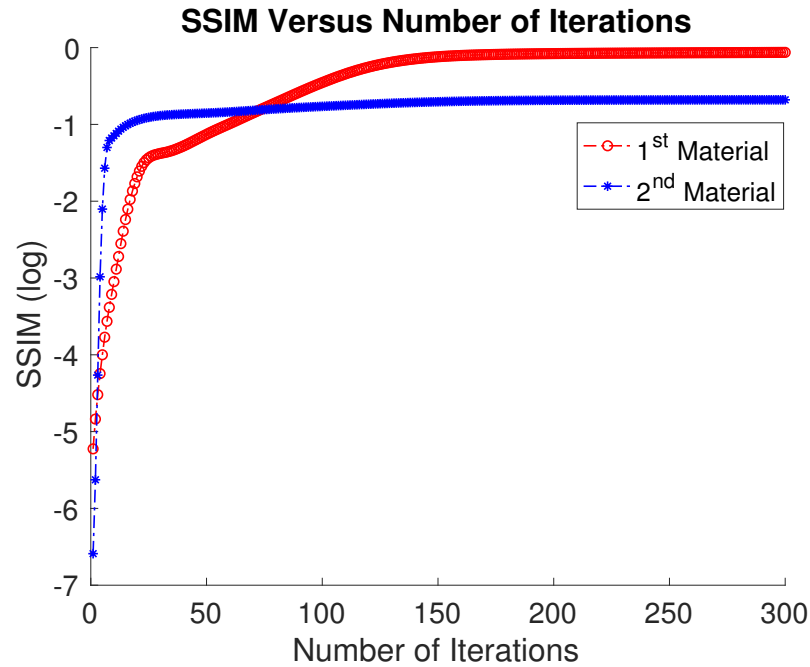


Figure 5.6: The SSIM of each iteration (with preconditioner) for plexiglass and PVC.

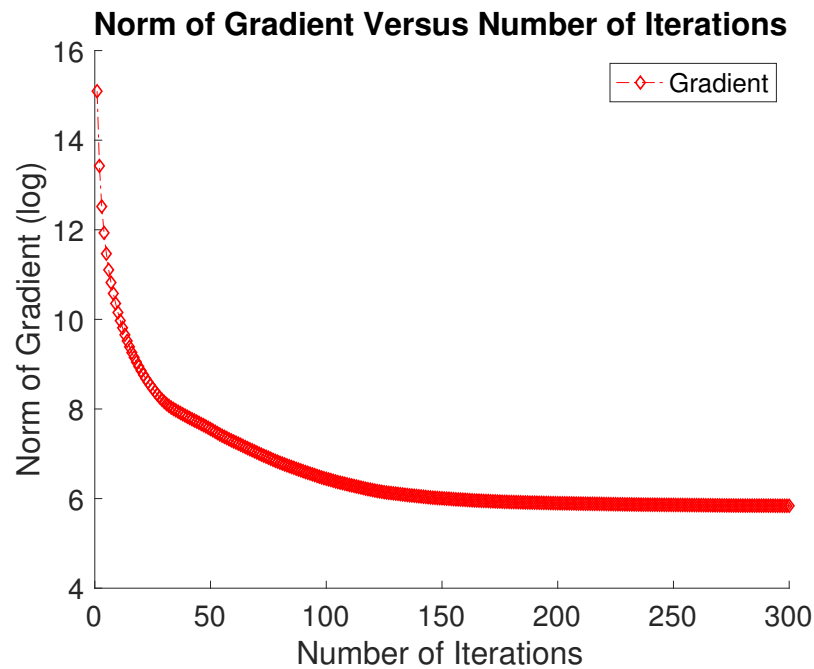


Figure 5.7: The decay of norm of the gradient for overall materials.

decreases significantly in the beginning and levels off after a sufficient number of iterations, indicating the convergence to a minimizer.

To further validate the strength of our proposed preconditioner, we compare the perfor-

mance with a preconditioner proposed by Barber [2], and the performance without using any preconditioners. As previously mentioned, the approach proposed in [2] is based on the eigenvalue decomposition of $\mathbf{C}^T \mathbf{C}$. The results are shown in Figure 5.8, where we plot the decay of relative errors for these three cases. To reduce clutter in this plot, we only show results for the first material; the behavior for the second material is the same. From

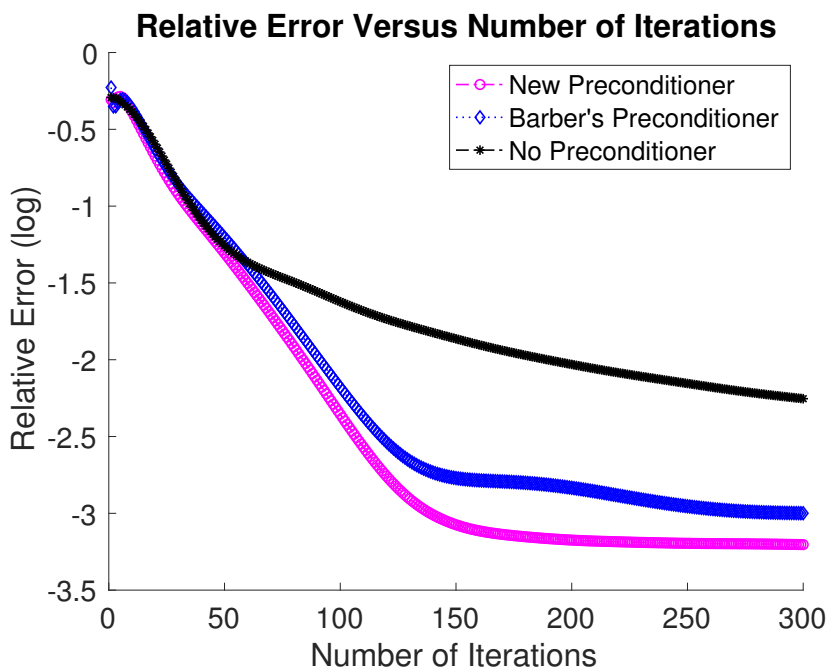


Figure 5.8: The decay of related errors with new preconditioner, Barber's [2] preconditioner, and with no preconditioner.

this figure, we can easily observe that both preconditioners are effective at accelerating convergence, with our approach producing the fastest convergence and the lowest relative errors.

5.5 Conclusions and Remarks

In this chapter, we use the Gaussian assumption of noise to construct a weighted least squares problem under bound constraints for energy discriminating x-ray detectors in computed tomography. Based on this problem, we propose a new preconditioner that includes not only the information of the linear attenuation coefficient matrix \mathbf{C} but also the pro-

jected data matrix \mathbf{Y} and the energy spectrum matrix \mathbf{S} . With this new preconditioner, the condition number of the Hessian can be reduced significantly. To implement this new preconditioner within an optimization framework, we suggest to use a first order method, FISTA, that can generate fast convergence speed. Because of the introduction of the new preconditioner, we recommend to construct a projection problem and compute the nearest step that will satisfy the linear inequality constraints for each iteration. Finally, numerical experiments also specify the advantages of the method mentioned in this chapter. For future work, it would be interesting to consider other regularization schemes to emphasize the edges of the object, such as the total variation.

Chapter 6

A Two-Step Method for Energy-windowed Spectral Computed Tomography

In this chapter, we still focus on the energy-windowed spectral computed tomography model

$$\mathbf{Y} = \exp(-\mathbf{A}\mathbf{W}\mathbf{C}^T) \mathbf{S} + \boldsymbol{\varepsilon}, \quad (6.1)$$

where \mathbf{Y} is a matrix that gathers the projected data of each energy window in the corresponding column and the exponential operator is applied element-wise (i.e., it is not a matrix function). \mathbf{A} is a matrix that is related to the quantitative information of ray trace and \mathbf{C} is a matrix that contains linear attenuation coefficients for particular (known) materials at specified energies. \mathbf{S} is the matrix that accumulates the spectrum energies for each energy window in the corresponding column. We assume that \mathbf{S} is square and invertible. Moreover, $\boldsymbol{\varepsilon}$ represents the noise term and we assume that $E_{il} \sim \mathcal{N}(0, y_{il})$ for each component E_{il} in $\boldsymbol{\varepsilon}$ and y_{il} in \mathbf{Y} . We assume that these data are known and the target is to solve the unknown weight matrix \mathbf{W} . \mathbf{W} is of size N_v by N_m , where N_v is the number of voxels (pixels if 2D) for each material map and N_m is the number of materials. Since the weight matrix \mathbf{W} represents the material maps of different materials, then it must be nonnegative

and we need to add a lower bound constraint $\mathbf{W} \geq \mathbf{0}$.

In the previous chapters, we try to compute the solution directly. However, we can also introduce an auxiliary variable and solve it using two steps. If we let $\mathbf{X} = \mathbf{A}\mathbf{W}$, then a two-step model is expressed as

$$\begin{aligned} \mathbf{Y} &= \exp(-\mathbf{X}\mathbf{C}^T) \mathbf{S} + \boldsymbol{\varepsilon}, \\ \mathbf{X} &= \mathbf{A}\mathbf{W}, \quad \mathbf{W} \geq \mathbf{0}. \end{aligned} \tag{6.2}$$

Based on the Gaussian assumption of noise, we can repeat the process of Chapter 5 to develop a weighted least squares framework. However, we do not want to solve the entire model using an iterative method. Instead, we want to construct a series of small-sized problems and solve each of them efficiently and accurately. Recall that when we implement the Algebraic Reconstruction Technique (ART) to solve a least squares problem, the iteration format is based on each row. Inspired by this idea, we want to come up with a similar method that takes advantage of each row of a weighted least squares problem. In the first step, we can build a row-wise optimization problem of small-sized and solve each one with cheap cost. In the second step, we can sum up the results obtained from the first step and solve a linear system under bound constraints. A challenge comes from the noise propagation, where we have to quantify the noise in the linear system of the second step. This quantity can be obtained using the properties of linear transform of the covariance matrix in the first step and we can use it to build another weighted least squares problem.

We should also mention that splitting the original problem into two steps and solving it in sequence might put the solution away from the truth. Instead, we can solve (6.2) directly by constructing an optimization problem of two coupled terms,

$$\begin{aligned} \min_{\bar{\mathbf{x}}, \bar{\mathbf{w}}} \quad & \frac{1}{2} \|(\mathbf{I} \otimes \mathbf{C}) \bar{\mathbf{x}} - \hat{\mathbf{b}}\|_{\boldsymbol{\Sigma}_{\hat{\mathbf{b}}}^{-1}}^2 + \frac{1}{2} \|(\mathbf{A} \otimes \mathbf{I}) \bar{\mathbf{w}} - \bar{\mathbf{x}}\|_{\boldsymbol{\Sigma}_{\bar{\mathbf{x}}_{\text{ml}}}^{-1}}^2, \\ \text{subject to} \quad & \bar{\mathbf{w}} \geq \mathbf{0}, \end{aligned} \tag{6.3}$$

where $\bar{\mathbf{x}}$ and $\bar{\mathbf{w}}$ are vectorizations of \mathbf{X}^T and \mathbf{W}^T , respectively, and $\boldsymbol{\Sigma}_{\hat{\mathbf{b}}}^{-1}$ and $\boldsymbol{\Sigma}_{\bar{\mathbf{x}}_{\text{ml}}}^{-1}$ are noise covariance matrices. Compared with the two-step method, the solution to (6.3) can

be obtained using the traditional iterative scheme and the solution is likely to be closer to the true solution of (6.2). On the other hand, it is likely that the solution to problem (6.3) might be better in theory, but it can be less attractive in terms of image quality. In this case, we should conduct numerical experiments to demonstrate strengths and weaknesses of these methods.

This chapter is organized as follows. In section (6.1), we review the limitations of previous methods and set up the two-step model. How to obtain the solution in two steps is also explained in this section. In section (6.2), we discuss another method, the coupled method, to solve the two-step model. This method combines the two steps together and uses an iterative framework to compute the solution. Numerical experiments are presented in (6.3). Concluding remarks and comparisons with the previous method are discussed in (6.4).

6.1 The Two-step Method

6.1.1 The Framework of Two-step Model

Compared with traditional CT models, the energy-windowed spectral CT model expands the room for developing mathematical theories and computing numerical solutions. In Chapter 4, the basic equation related to this model is reformulated as a nonlinear least squares problem and solved by nonlinear optimization. In Chapter 5, we have transformed the basic equation into a weighted least squares problem, developed a new preconditioner and used FISTA to compute the solution. Even if these two methods are based on different objective functions, they both include a preconditioner to mitigate the influence of ill-posedness. So it raises a question if we can solve this equation efficiently to a high accuracy without using preconditioners. In this section, we will focus on a two-step method which can calculate the solution to the weighted least squares model accurately and economically.

Recall that the discretized energy-windowed spectral CT model is expressed as

$$\mathbf{Y} = \exp(-\mathbf{A}\mathbf{W}\mathbf{C}^T) \mathbf{S} + \boldsymbol{\varepsilon}, \quad (6.4)$$

where

- \mathbf{Y} is a matrix of size $(N_d \cdot N_p) \times N_b$ that gathers x-ray photons of each energy window in the corresponding column.
- \mathbf{A} is a matrix of size $(N_d \cdot N_p) \times N_v$ that collects the fan-beam geometry and each element corresponds to $a_{i,j}$.
- \mathbf{C} is a matrix of size $N_e \times N_m$ that accumulates linear attenuation coefficients and each entry corresponds to $u_{e,m}$, the linear attenuation coefficient of e -th energy and m -th material.
- \mathbf{W} is a matrix of size $N_v \times N_m$ and each column corresponds the unknown material map we want to reconstruct.
- \mathbf{S} is a matrix of size $N_e \times N_b$ and each column collects the spectrum energy of a specific range. We assume that \mathbf{S} is square and invertible.
- \mathcal{E} is the noise matrix that is of size $(N_d \cdot N_p) \times N_b$. The assumption for noise is that $E_{il} \sim \mathcal{N}(0, y_{il})$ for each element E_{il} in \mathcal{E} and y_{il} in \mathbf{Y} .

As we have shown in Chapter 5, we can vectorize (6.4) and build a weighted least squares problem using the Gaussian assumption of noise

$$\begin{aligned} \min_{\mathbf{w}} \quad & \frac{1}{2} \|\mathcal{A}\mathbf{w} - \mathbf{b}\|_{\Sigma^{-1}}^2 \\ \text{subject to} \quad & \mathbf{w} \geq \mathbf{0}, \end{aligned} \tag{6.5}$$

where $\mathcal{A} = \mathbf{C} \otimes \mathbf{A}$, $\mathbf{b} = -\log(\mathbf{y})$, $\mathbf{y} = \text{vec}(\mathbf{Y})$ and $\mathbf{w} = \text{vec}(\mathbf{W})$. Σ^{-1} , which combines information from \mathbf{S} and \mathbf{y} , is the inverse covariance matrix generated by the Gaussian noise and logarithmic transformation. $\|\cdot\|_{\Sigma^{-1}}^2$ represents a weighted 2-norm and $\|\mathcal{A}\mathbf{w} - \mathbf{b}\|_{\Sigma^{-1}}^2 = (\mathcal{A}\mathbf{w} - \mathbf{b})^T \Sigma^{-1} (\mathcal{A}\mathbf{w} - \mathbf{b})$. Because of the collinearity of \mathbf{C} and the sparsity of \mathbf{A} , solutions to Equation (6.5) might not be satisfactory. In Chapter 5, we introduce a preconditioner

\mathbf{M} and the optimization problem (6.5) has been transformed into

$$\begin{aligned} \min_{\tilde{\mathbf{w}}} \quad & \frac{1}{2} \|\tilde{\mathbf{A}}\tilde{\mathbf{w}} - \mathbf{b}\|_{\Sigma^{-1}}^2 \\ \text{subject to} \quad & (\mathbf{M} \otimes \mathbf{I}) \tilde{\mathbf{w}} \geq \mathbf{0}. \end{aligned} \quad (6.6)$$

where $\tilde{\mathbf{C}} = \mathbf{C}\mathbf{M}$, $\tilde{\mathbf{A}} = \tilde{\mathbf{C}} \otimes \mathbf{A}$, $\tilde{\mathbf{W}} = \mathbf{W}\mathbf{M}^{-T}$ and $\tilde{\mathbf{w}} = \text{vec}(\tilde{\mathbf{W}})$. To solve problem (6.6), we implement FISTA to solve the main problem and build a projection subproblem (6.7) to maintain nonnegativity

$$\begin{aligned} \min_{\tilde{\mathbf{w}}_{new}} \quad & \|\tilde{\mathbf{w}}_{new} - \tilde{\mathbf{w}}_k\|_2^2 \\ \text{subject to} \quad & (\mathbf{M} \otimes \mathbf{I}) \tilde{\mathbf{w}}_{new} \geq \mathbf{0}, \end{aligned} \quad (6.7)$$

where $\tilde{\mathbf{w}}_k$ is the current iteration of $\tilde{\mathbf{w}}$. To compute the solution to subproblem (6.7), we reshape all variables back to a matrix form and split this form into several pieces. Each piece is a tiny size problem with a row in the original matrix and the corresponding constraints. For each problem, we solve it using a high performance solver such as CVXGEN. Even if this method is efficient for large-scale problems, solving the projection problem might still be a very expensive cost. So we want to figure out a method that can both solve the problem (6.4) accurately and avoid large expense on keeping the nonnegativity of the solution.

Rather than computing the solution directly, we propose a two-step method with an auxiliary step and bridge it to the desired solution. Recall that our goal is to solve \mathbf{W} from the equation

$$\mathbf{Y} = \exp(-\mathbf{A}\mathbf{W}\mathbf{C}^T) \mathbf{S} + \boldsymbol{\varepsilon}. \quad (6.8)$$

We first let $\mathbf{X} = \mathbf{A}\mathbf{W}$ and solve \mathbf{X} from Equation (6.9)

$$\mathbf{Y} = \exp(-\mathbf{X}\mathbf{C}^T) \mathbf{S} + \boldsymbol{\varepsilon}. \quad (6.9)$$

After we have obtained \mathbf{X} , we solve \mathbf{W} from the equation

$$\mathbf{X} = \mathbf{A}\mathbf{W}, \quad \mathbf{W} \geq \mathbf{0}. \quad (6.10)$$

Recall that the size of \mathbf{W} is $N_v \times N_m$ and the auxiliary step is bound to \mathbf{X} , which is of size $(N_d \cdot N_p) \times N_m$. Usually, $N_d \cdot N_p$ is much larger than N_v so we seem to introduce more unknown variables and this behavior is undesirable. Another problem we need to consider is how to quantify the noise propagation from Equation (6.9) to Equation (6.10). In Equation (6.9), we solve \mathbf{X} with the noise \mathcal{E} so the solution must include probability properties of \mathcal{E} . When we solve \mathbf{W} from (6.10), the variable \mathbf{X} is assumed to be known so we should also identify the noise.

6.1.2 A Solution to the Two-step Model

To solve the first problem, we try to use the statistical properties of noise to decompose the entire problem into small pieces and solve each of these problems accurately. Recall that \mathcal{E} is the noise matrix and $E_{il} \sim \mathcal{N}(0, y_{il})$ for each element E_{il} in \mathcal{E} and y_{il} in \mathbf{Y} , so the noise is independent for each entry of projected data. With this independence, we can decompose Equation (6.9) in a row-wise or column-wise way. Since \mathbf{Y} is of size $(N_d \cdot N_p) \times N_b$ and the number of energy windows, N_b , is normally much fewer than $N_d \cdot N_p$, the splitting over rows is preferable.

If we let

$$\begin{aligned}\mathbf{Y}^T &= [\mathbf{y}_1, \mathbf{y}_2, \dots, \mathbf{y}_{N_d \cdot N_p}], \\ \mathbf{X}^T &= [\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_{N_d \cdot N_p}], \\ \mathcal{E}^T &= [\mathbf{e}_1, \mathbf{e}_2, \dots, \mathbf{e}_{N_d \cdot N_p}],\end{aligned}\tag{6.11}$$

then Equation (6.9) can be rewritten with respect to each row

$$\mathbf{y}_i^T = \exp(-\mathbf{x}_i^T \mathbf{C}^T) \mathbf{S} + \mathbf{e}_i^T,\tag{6.12}$$

for $i = 1, 2, \dots, N_d \cdot N_p$. Meanwhile, we can notice that each row of \mathbf{Y} only depends on the corresponding row of \mathbf{X} and \mathcal{E} as well as matrices \mathbf{C} and \mathbf{S} . If we take the transpose

for Equation (6.12), then we have

$$\mathbf{y}_i = \mathbf{S}^T \exp(-\mathbf{C}\mathbf{x}_i) + \mathbf{e}_i. \quad (6.13)$$

By assumption, \mathbf{S} is square and invertible, so we can multiply \mathbf{S}^{-T} on both sides

$$\mathbf{S}^{-T} \mathbf{y}_i = \exp(-\mathbf{C}\mathbf{x}_i) + \mathbf{S}^{-T} \mathbf{e}_i. \quad (6.14)$$

If we let $\hat{\mathbf{y}}_i = \mathbf{S}^{-T} \mathbf{y}_i$ and $\hat{\mathbf{e}}_i = \mathbf{S}^{-T} \mathbf{e}_i$, then Equation (6.14) is equivalent to

$$\hat{\mathbf{y}}_i = \exp(-\mathbf{C}\mathbf{x}_i) + \hat{\mathbf{e}}_i. \quad (6.15)$$

By subtracting $\hat{\mathbf{e}}_i$ on both sides and taking the logarithm, we have

$$\log(\hat{\mathbf{y}}_i - \hat{\mathbf{e}}_i) = -\mathbf{C}\mathbf{x}_i. \quad (6.16)$$

Using the Taylor expansion, we can expand the left hand side of Equation (6.16) as

$$\log(\hat{\mathbf{y}}_i - \hat{\mathbf{e}}_i) = \log(\hat{\mathbf{y}}_i) - \text{diag}(\hat{\mathbf{y}}_i)^{-1} \hat{\mathbf{e}}_i + \mathcal{O}(\|\hat{\mathbf{e}}_i\|_2^2). \quad (6.17)$$

Therefore, Equation (6.16) can be estimated as

$$\hat{\mathbf{b}}_i \approx \mathbf{C}\mathbf{x}_i - \text{diag}(\hat{\mathbf{y}}_i)^{-1} \hat{\mathbf{e}}_i. \quad (6.18)$$

where $\hat{\mathbf{b}}_i = -\log(\hat{\mathbf{y}}_i)$. By assumption, we have $\mathbf{e}_i \sim \mathcal{N}(\mathbf{0}, \text{diag}(\mathbf{y}_i))$. Then we can obtain that

$$\hat{\mathbf{b}}_i | \mathbf{x}_i \sim \mathcal{N}\left(\mathbf{C}\mathbf{x}_i, \boldsymbol{\Sigma}_{\hat{\mathbf{b}}_i}\right), \quad (6.19)$$

where $\boldsymbol{\Sigma}_{\hat{\mathbf{b}}_i} = \text{diag}(\hat{\mathbf{y}}_i)^{-1} \mathbf{S}^{-T} \text{diag}(\mathbf{y}_i) \mathbf{S}^{-1} \text{diag}(\hat{\mathbf{y}}_i)^{-1}$. By ignoring constants, the corresponding probability density function can be expressed as

$$f(\hat{\mathbf{b}}_i; \mathbf{x}_i) = \exp\left\{-\frac{1}{2} (\mathbf{C}\mathbf{x}_i - \hat{\mathbf{b}}_i)^T \boldsymbol{\Sigma}_{\hat{\mathbf{b}}_i}^{-1} (\mathbf{C}\mathbf{x}_i - \hat{\mathbf{b}}_i)\right\}. \quad (6.20)$$

So the log-likelihood function is given by

$$l(\mathbf{x}_i; \hat{\mathbf{b}}_i) = -\frac{1}{2} (\mathbf{C}\mathbf{x}_i - \hat{\mathbf{b}}_i)^T \boldsymbol{\Sigma}_{\hat{\mathbf{b}}_i}^{-1} (\mathbf{C}\mathbf{x}_i - \hat{\mathbf{b}}_i). \quad (6.21)$$

Our goal is to maximize the log-likelihood function $l(\mathbf{x}_i; \hat{\mathbf{b}}_i)$ and it is equivalent to minimizing the negative log-likelihood function $-l(\mathbf{x}_i; \hat{\mathbf{b}}_i)$. Based on Equation (6.21), the maximum likelihood estimator for \mathbf{x}_i can be represented as

$$(\mathbf{x}_i)_{\text{ml}} = \underset{\mathbf{x}_i}{\operatorname{argmin}} \left\{ \left\| \mathbf{C}\mathbf{x}_i - \hat{\mathbf{b}}_i \right\|_{\boldsymbol{\Sigma}_{\hat{\mathbf{b}}_i}^{-1}}^2 \right\}. \quad (6.22)$$

This equation can be solved analytically and the corresponding solution is

$$(\mathbf{x}_i)_{\text{ml}} = \left(\mathbf{C}^T \boldsymbol{\Sigma}_{\hat{\mathbf{b}}_i}^{-1} \mathbf{C} \right)^{-1} \mathbf{C}^T \boldsymbol{\Sigma}_{\hat{\mathbf{b}}_i}^{-1} \hat{\mathbf{b}}_i. \quad (6.23)$$

Since the matrix \mathbf{C} is of size N_e by N_m , then solving each $(\mathbf{x}_i)_{\text{ml}}$ is of low cost. So for the first step, we can loop around all $i = 1, 2, \dots, N_d \cdot N_p$ and solve each linear system to a high accuracy. After we have obtained all $(\mathbf{x}_i)_{\text{ml}}$, we can build the maximum likelihood solution for the matrix \mathbf{X} by concatenation

$$\mathbf{X}_{ml}^T = [(\mathbf{x}_1)_{\text{ml}}, (\mathbf{x}_2)_{\text{ml}}, \dots, (\mathbf{x}_{N_d \cdot N_p})_{\text{ml}}]. \quad (6.24)$$

So far we have obtained the maximum likelihood solution, \mathbf{X}_{ml}^T , for Equation (6.9). However, we cannot substitute it for the same variable in Equation (6.10) because of the unknown noise.

To quantify the noise propagation, we repeat a similar process as we have shown in Chapter 5. To be consistent with the row-wise model, we first take a transpose of Equation (6.9):

$$\mathbf{Y}^T = \mathbf{S}^T \exp(-\mathbf{C}\mathbf{X}^T) + \boldsymbol{\varepsilon}^T. \quad (6.25)$$

Multiplying \mathbf{S}^{-T} on both sides, we can obtain that

$$\mathbf{S}^{-T} \mathbf{Y}^T = \exp(-\mathbf{C} \mathbf{X}^T) + \mathbf{S}^{-T} \boldsymbol{\varepsilon}^T. \quad (6.26)$$

By vectorizing the variables on both sides, we have

$$(\mathbf{I} \otimes \mathbf{S}^{-T}) \bar{\mathbf{y}} = \exp(-(\mathbf{I} \otimes \mathbf{C}) \bar{\mathbf{x}}) + (\mathbf{I} \otimes \mathbf{S}^{-T}) \bar{\mathbf{e}}, \quad (6.27)$$

where $\bar{\mathbf{y}} = \text{vec}(\mathbf{Y}^T)$, $\bar{\mathbf{x}} = \text{vec}(\mathbf{X}^T)$ and $\bar{\mathbf{e}} = \text{vec}(\boldsymbol{\varepsilon}^T)$. Using the former notations, we define $\hat{\mathbf{y}} = (\mathbf{I} \otimes \mathbf{S}^{-T}) \bar{\mathbf{y}}$, $\hat{\mathbf{e}} = (\mathbf{I} \otimes \mathbf{S}^{-T}) \bar{\mathbf{e}}$ and $\hat{\mathbf{b}} = -\log(\hat{\mathbf{y}})$. Then Equation (6.27) can be written as

$$\hat{\mathbf{y}} = \exp(-(\mathbf{I} \otimes \mathbf{C}) \bar{\mathbf{x}}) + \hat{\mathbf{e}}. \quad (6.28)$$

We take the logarithm as in (6.16) and expand the logarithmic term as Equation (6.17), then Equation (6.28) can be approximated as

$$\hat{\mathbf{b}} \approx (\mathbf{I} \otimes \mathbf{C}) \bar{\mathbf{x}} - \text{diag}(\hat{\mathbf{y}})^{-1} (\mathbf{I} \otimes \mathbf{S}^{-T}) \bar{\mathbf{e}}, \quad (6.29)$$

By assumption, we assume that $\bar{\mathbf{e}} \sim \mathcal{N}(\mathbf{0}, \text{diag}(\bar{\mathbf{y}}))$. Then we have

$$\hat{\mathbf{b}} | \bar{\mathbf{x}} \sim \mathcal{N}((\mathbf{I} \otimes \mathbf{C}) \bar{\mathbf{x}}, \boldsymbol{\Sigma}_{\hat{\mathbf{b}}}), \quad (6.30)$$

where $\boldsymbol{\Sigma}_{\hat{\mathbf{b}}} = \text{diag}(\hat{\mathbf{y}})^{-1} (\mathbf{I} \otimes \mathbf{S}^{-T}) \text{diag}(\bar{\mathbf{y}}) (\mathbf{I} \otimes \mathbf{S}^{-1}) \text{diag}(\hat{\mathbf{y}})^{-1}$ is the block-diagonal covariance matrix. If we build the maximum likelihood function with $\boldsymbol{\Sigma}_{\hat{\mathbf{b}}}$, then the maximum likelihood estimator of $\bar{\mathbf{x}}$ can be expressed as

$$\begin{aligned} \bar{\mathbf{x}}_{\text{ml}} &= \underset{\mathbf{x}}{\text{argmin}} \left\{ \frac{1}{2} \|(\mathbf{I} \otimes \mathbf{C}) \mathbf{x} - \hat{\mathbf{b}}\|_{\boldsymbol{\Sigma}_{\hat{\mathbf{b}}}^{-1}}^2 \right\} \\ &= \left((\mathbf{I} \otimes \mathbf{C}^T) \boldsymbol{\Sigma}_{\hat{\mathbf{b}}}^{-1} (\mathbf{I} \otimes \mathbf{C}) \right)^{-1} (\mathbf{I} \otimes \mathbf{C}^T) \boldsymbol{\Sigma}_{\hat{\mathbf{b}}}^{-1} \hat{\mathbf{b}}. \end{aligned} \quad (6.31)$$

By using assumption (6.30), we know that $\bar{\mathbf{x}}_{\text{ml}}$ is itself a random variable with mean

$(\mathbf{I} \otimes \mathbf{C}) \bar{\mathbf{x}}$ and covariance

$$\boldsymbol{\Sigma}_{\bar{\mathbf{x}}_{\text{ml}}} = \left((\mathbf{I} \otimes \mathbf{C}^T) \boldsymbol{\Sigma}_{\hat{\mathbf{b}}}^{-1} (\mathbf{I} \otimes \mathbf{C}) \right)^{-1}. \quad (6.32)$$

So far we have obtained the noise covariance matrix for the second step. With this covariance matrix, we can build another optimization problem under bound constraints. For the second step, we transpose the equation first in order to be consistent with $\bar{\mathbf{x}}$:

$$\mathbf{X}_{\text{ml}}^T = \mathbf{W}^T \mathbf{A}^T. \quad (6.33)$$

Taking vectorization on both sides, we can obtain that

$$\bar{\mathbf{x}}_{\text{ml}} = (\mathbf{A} \otimes \mathbf{I}) \bar{\mathbf{w}}. \quad (6.34)$$

Using the noise covariance matrix (6.32), we can formulate a weighted least squares problem under bound constraints to solve $\bar{\mathbf{w}}$:

$$\begin{aligned} \underset{\bar{\mathbf{w}}}{\operatorname{argmin}} \quad & \frac{1}{2} \| (\mathbf{A} \otimes \mathbf{I}) \bar{\mathbf{w}} - \bar{\mathbf{x}}_{\text{ml}} \|_{\boldsymbol{\Sigma}_{\bar{\mathbf{x}}_{\text{ml}}}^{-1}}^2 \\ \text{subject to} \quad & \bar{\mathbf{w}} \geq \mathbf{0}. \end{aligned} \quad (6.35)$$

For the regularization term, we take ℓ_1 regularization to keep sparsity. So the optimization problem we need to solve is

$$\begin{aligned} \underset{\bar{\mathbf{w}}}{\operatorname{argmin}} \quad & \frac{1}{2} \| (\mathbf{A} \otimes \mathbf{I}) \bar{\mathbf{w}} - \bar{\mathbf{x}}_{\text{ml}} \|_{\boldsymbol{\Sigma}_{\bar{\mathbf{x}}_{\text{ml}}}^{-1}}^2 + \beta \| \bar{\mathbf{w}} \|_1 \\ \text{subject to} \quad & \bar{\mathbf{w}} \geq \mathbf{0}, \end{aligned} \quad (6.36)$$

where β is the regularization parameter. For problem (6.36), we solve it using FISTA with a projection step onto the boundary. The iteration involved in FISTA only requires the current gradient as well as the current and previous steps. The current and previous steps are generated iteratively while the current gradient can be computed using the properties of the Kronecker product. In this case, we only need a matrix-vector multiplication in each

iteration and thus can avoid forming the Hessian explicitly.

To conclude, this two-step method can damp the influence of the matrix \mathbf{C} and it has a similar effect as the preconditioner in the last chapter. However, rather than building a preconditioned system, we solve several row-wise, small-size problems to high accuracy. Compared with the preconditioning method, it is not necessary to solve a projection problem in each iteration so the convergence speed is fast. Moreover, the solution of each auxiliary system does not depend on each other, so we can further implement the solver for the first step in parallel. On the other hand, the possible drawbacks come from the cost of solving the linear system in the first step, which might be increased significantly when the size of image is large. These limitations can also be mitigated by parallel computation.

6.2 The Coupled Method

In the last section, we have discussed how to solve Equation (6.4) using a two step method. We use an auxiliary variable \mathbf{X} such that $\mathbf{X} = \mathbf{A}\mathbf{W}$. Then we solve \mathbf{X} in the first step and \mathbf{W} in the second step. We derive the optimization problem using a log-likelihood function and use the Gaussian noise assumption to obtain the propagation of noise in the second step. Then we solve these two problems in sequence and obtain the final result. However, in each of these two steps, we might either overestimate or underestimate the solutions and the combined solution is likely to be distorted. To compare and evaluate the two-step method, we can come up with an integrated framework that merges these separated steps.

Instead of solving these two equations individually, we can try to combine previous two steps and calculate the solution alternatively. In each iteration, we solve the first equation and evaluate the result to the corresponding part in the second equation. Then we solve the second equation under the current circumstances. In this case, we have coupled the results generated from two separated steps and they might be influential to the following iterations. Moreover, each equation is still separable and we can use optimization techniques, such as the coordinate descent method, to update each step. In the previous inferences, we try to

solve the following two problems

$$\begin{aligned} \mathbf{Y} &= \exp(-\mathbf{X}\mathbf{C}^T) \mathbf{S} + \boldsymbol{\varepsilon}, \\ \mathbf{X} &= \mathbf{A}\mathbf{W}, \quad \mathbf{W} \geq \mathbf{0}. \end{aligned} \tag{6.37}$$

The noise covariance matrix is $\boldsymbol{\Sigma}_{\hat{\mathbf{b}}-1}$ for the first equation and for the second equation, the noise covariance matrix can be represented as $\boldsymbol{\Sigma}_{\bar{\mathbf{x}}_{\text{ml}}-1}$. Recall that the noise covariance matrix for the second step is based on the maximum likelihood estimate of the first step so if we couple these two problems, it is only based on the current maximum likelihood estimator of the first step. With these two covariance matrices, we can build two optimization problems

$$\begin{aligned} \min_{\bar{\mathbf{x}}} \quad & \frac{1}{2} \|(\mathbf{I} \otimes \mathbf{C}) \bar{\mathbf{x}} - \hat{\mathbf{b}}\|_{\boldsymbol{\Sigma}_{\hat{\mathbf{b}}-1}}^2, \\ \min_{\bar{\mathbf{w}}} \quad & \frac{1}{2} \|(\mathbf{A} \otimes \mathbf{I}) \bar{\mathbf{w}} - \bar{\mathbf{x}}_{\text{ml}}\|_{\boldsymbol{\Sigma}_{\bar{\mathbf{x}}_{\text{ml}}-1}}^2, \\ & \text{subject to } \bar{\mathbf{w}} \geq \mathbf{0}. \end{aligned} \tag{6.38}$$

where

$$\begin{aligned} \boldsymbol{\Sigma}_{\hat{\mathbf{b}}} &= \text{diag}(\hat{\mathbf{y}})^{-1} (\mathbf{I} \otimes \mathbf{S}^{-T}) \text{diag}(\bar{\mathbf{y}}) (\mathbf{I} \otimes \mathbf{S}^{-1}) \text{diag}(\hat{\mathbf{y}})^{-1}, \\ \boldsymbol{\Sigma}_{\bar{\mathbf{x}}_{\text{ml}}} &= \left((\mathbf{I} \otimes \mathbf{C}^T) \boldsymbol{\Sigma}_{\hat{\mathbf{b}}-1}^{-1} (\mathbf{I} \otimes \mathbf{C}) \right)^{-1}. \end{aligned} \tag{6.39}$$

These two problems can be merged into one optimization problem as

$$\begin{aligned} \min_{\bar{\mathbf{x}}, \bar{\mathbf{w}}} \quad & \frac{1}{2} \|(\mathbf{I} \otimes \mathbf{C}) \bar{\mathbf{x}} - \hat{\mathbf{b}}\|_{\boldsymbol{\Sigma}_{\hat{\mathbf{b}}-1}}^2 + \frac{1}{2} \|(\mathbf{A} \otimes \mathbf{I}) \bar{\mathbf{w}} - \bar{\mathbf{x}}\|_{\boldsymbol{\Sigma}_{\bar{\mathbf{x}}_{\text{ml}}-1}}^2, \\ & \text{subject to } \bar{\mathbf{w}} \geq \mathbf{0}. \end{aligned} \tag{6.40}$$

The nonnegative constraint $\bar{\mathbf{w}} \geq \mathbf{0}$ is equivalent to

$$\chi(w_i) = \begin{cases} 0, & w_i \geq 0, \\ \infty, & w_i < 0, \end{cases} \tag{6.41}$$

for $i = 1, 2, \dots, N_v \times N_m$. For the regularization, we want to penalize the total sum of weights. So we add ℓ_1 regularization and use another variable, $\bar{\mathbf{z}}$, to facilitate the process.

The integrated optimization is expressed as

$$\min_{\bar{\mathbf{x}}, \bar{\mathbf{w}}, \bar{\mathbf{z}}} \frac{1}{2} \|(\mathbf{I} \otimes \mathbf{C}) \bar{\mathbf{x}} - \hat{\mathbf{b}}\|_{\Sigma_{\hat{\mathbf{b}}}^{-1}}^2 + \frac{1}{2} \|(\mathbf{A} \otimes \mathbf{I}) \bar{\mathbf{z}} - \bar{\mathbf{x}}\|_{\Sigma_{\bar{\mathbf{x}}_{ml}}^{-1}}^2 + \frac{1}{2} \|\bar{\mathbf{z}} - \bar{\mathbf{w}}\|_2^2 + \beta \|\bar{\mathbf{w}}\|_1 + \chi(\bar{\mathbf{w}}). \quad (6.42)$$

As we can see, the first two terms in problem (6.42) are weighted least squares. Furthermore, the last three terms can be regarded as a soft shrinkage function with a mutation. Given $\bar{\mathbf{z}}$, the optimization problem can be described as

$$\min_{\bar{\mathbf{w}}} \frac{1}{2} \|\bar{\mathbf{w}} - \bar{\mathbf{z}}\|_2^2 + \beta \|\bar{\mathbf{w}}\|_1 + \chi(\bar{\mathbf{w}}), \quad (6.43)$$

Let the i -th element in $\bar{\mathbf{w}}$ be \bar{w}_i . If $\bar{w}_i \leq 0$, we have $\bar{w}_i = 0$. If $\bar{w}_i > 0$, we differentiate the equation with respect to \bar{w}_i and let the derivative equal 0 for $i = 1, 2, \dots, N_v \times N_m$. Then we can obtain that

$$\frac{1}{2} (2\bar{w}_i - 2\bar{z}_i) + \beta = 0. \quad (6.44)$$

So we have

$$\bar{w}_i = \bar{z}_i - \beta. \quad (6.45)$$

With the bound constraint, each \bar{w}_i should satisfy the condition

$$\bar{w}_i = \max(\bar{z}_i - \beta, 0). \quad (6.46)$$

This function is only defined on the half domain of a regular soft shrinkage function. The plot of this function is shown in Figure 6.1. Therefore, problem (6.43) has an analytical solution

$$\bar{\mathbf{w}} = \max(\bar{\mathbf{z}} - \beta, \mathbf{0}). \quad (6.47)$$

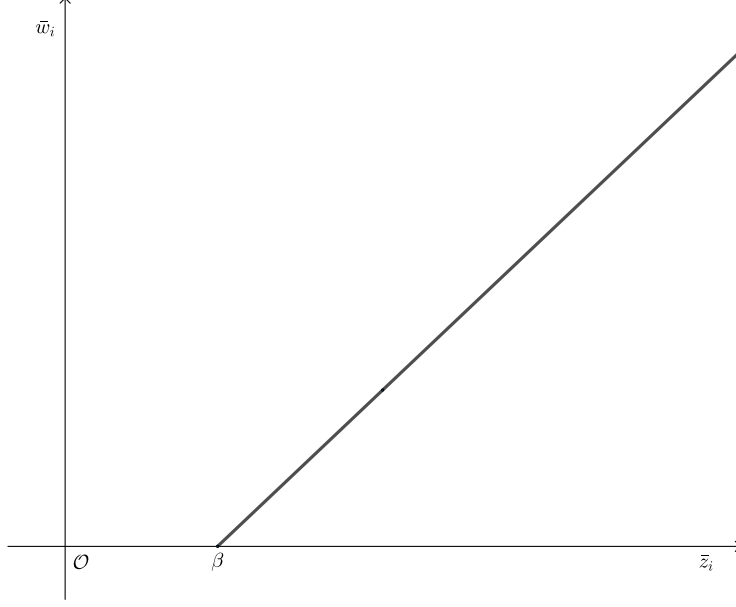


Figure 6.1: The soft shrinkage function with a bound constraint.

The updating framework to solve problem (6.42) is presented as

$$\begin{aligned}
 \bar{\mathbf{x}}_{k+1} &= \bar{\mathbf{x}}_k - \alpha_1 \left[(\mathbf{I} \otimes \mathbf{C}^T) \boldsymbol{\Sigma}_{\hat{\mathbf{b}}}^{-1} (\mathbf{I} \otimes \mathbf{C}) \bar{\mathbf{x}}_k - (\mathbf{I} \otimes \mathbf{C}^T) \boldsymbol{\Sigma}_{\hat{\mathbf{b}}}^{-1} \hat{\mathbf{b}} \right], \\
 \bar{\mathbf{z}}_{k+1} &= \bar{\mathbf{z}}_k - \alpha_2 \left[(\mathbf{A}^T \otimes \mathbf{I}) \boldsymbol{\Sigma}_{\hat{\mathbf{x}}_{\text{ml}}}^{-1} (\mathbf{A} \otimes \mathbf{I}) \bar{\mathbf{z}}_k - (\mathbf{A}^T \otimes \mathbf{I}) \boldsymbol{\Sigma}_{\hat{\mathbf{x}}_{\text{ml}}}^{-1} \bar{\mathbf{x}}_{k+1} \right], \\
 \bar{\mathbf{w}}_{k+1} &= \max(\bar{\mathbf{z}}_{k+1} - \beta, \mathbf{0}),
 \end{aligned} \tag{6.48}$$

where α_1 and α_2 are selected by line search. This iterative framework is similar to the coordinate descent method, while the result in the first step is reused in the second step and the result obtained from the second step is indispensable for the third step. This framework is simple and easy to implement and it is also clear that we can use the properties of the Kronecker product to avoid saving large matrices.

To conclude, this coupled method might generate a better solution to the optimization problem (6.2). It has the similarity to the coordinate descent method, in which we only need to consider one direction in each update of variables. Furthermore, this method is easy to implement because the iteration in each direction is straightforward. However, we should consider the quality of reconstructed images as well. Similar to the semi-convergence properties [16], the solution obtained with extensive iteration might give us unsatisfactory

results. In this case, we should conduct numerical experiments to compare this method with the two-step method proposed in the previous section.

6.3 Numerical Experiments

For the numerical experiments, we use the same test problem as Chapter 5. The object is composed of two materials, plexiglass and polyvinyl chloride (PVC). The first material map is made up of a circular mask in the center of the image while the second material map includes only small “spikes”. In real-life applications, the first material map can be represented by main tissues and these small spikes can indicate calcium in bones. The original images are shown in Figure 6.2. In addition, we use the same parameters as

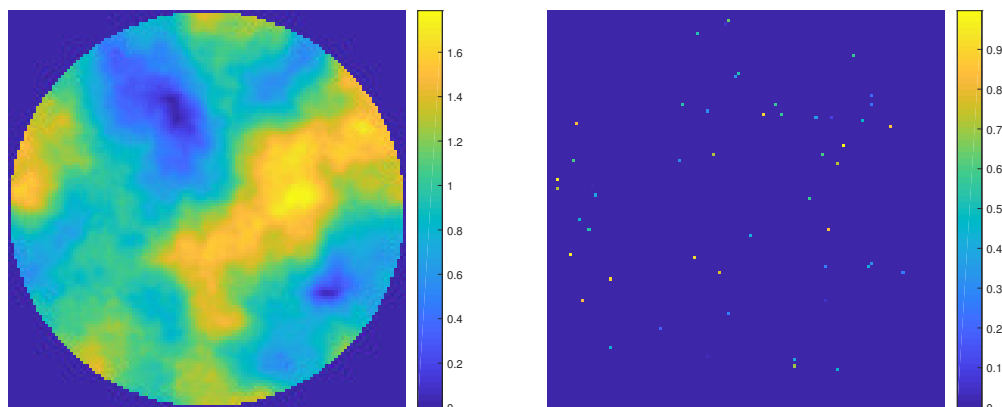


Figure 6.2: The original material maps for plexiglass (Left) and PVC (Right).

Table 4.1 to generate the matrix \mathbf{A} and we take 180 projections that are equally distributed from 0 to 360 degrees. The spectrum is built with 120 kV voltage and the detector is assumed to be photon-counting with five energy windows. These energy windows can detect photon energies of 10 to 34 KeV, 35 to 49 KeV, 50 to 64 KeV, 65 to 79 KeV and 80 to 120 KeV, respectively.

The reconstructed images obtained using the two-step method are shown in Figure 6.3. As we can identify, the images are of high quality in general. The distributions of weights are located in the same positions as the origins. For the first figure, the main profile is similar to the origin but the areas of yellow colors are shrank while the areas of blue color

are expanded. For the second figure, it still has shades but dots are located in the same places as the origin. Compared with the reconstructed images in Chapter 5, the first

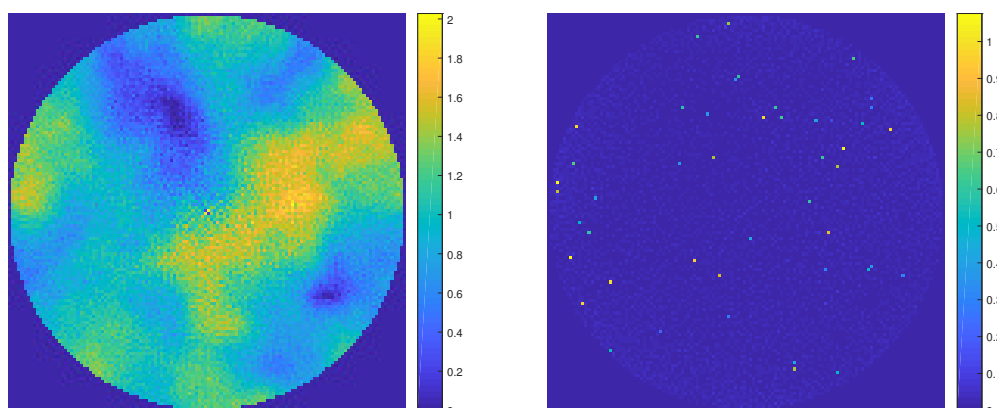


Figure 6.3: The reconstructed material maps for plexiglass (Left) and PVC (Right) using the two-step method.

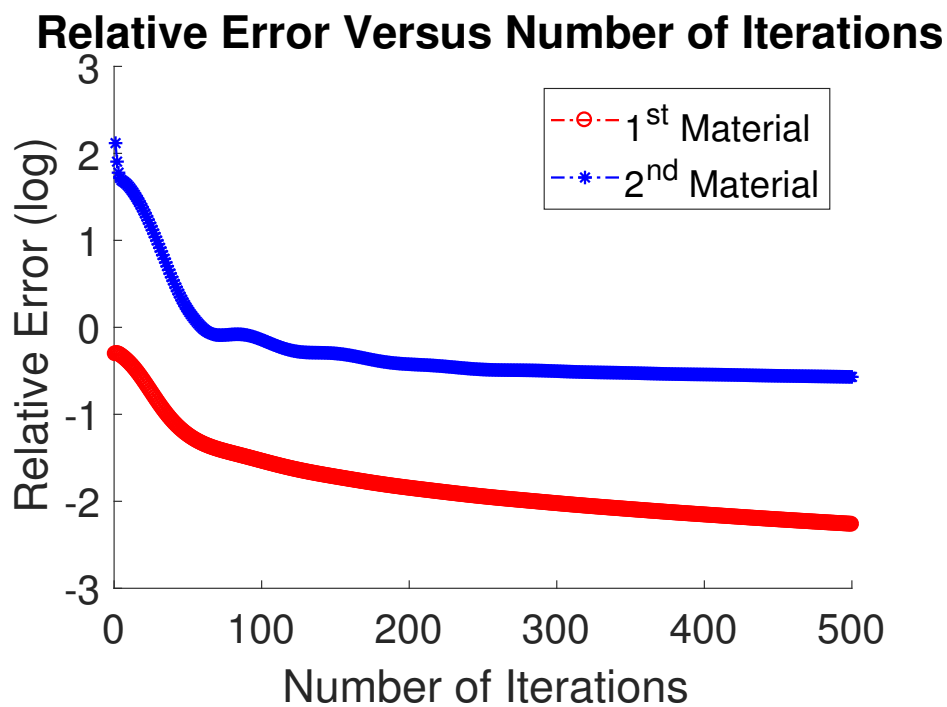


Figure 6.4: The relative errors of two material maps solved by the two-step method.

image has artifacts that are located in the center of the images. Moreover, the edges in the first image are not as smooth as the reconstructed image in Chapter 5. This might result from the ℓ_1 regularization used in this method or early termination of iterations. The

second images displayed in this chapter and Chapter 5 present similar shades and thus the reconstruction qualities are close.

We can also present the convergence properties by showing the plot of relative errors. The relative error figure is shown in Figure 6.4. From Figure 6.4, we can see that the relative error corresponding to the second material drops fast in the beginning but eventually stagnates. The relative error corresponding to the first material drops in a similar way but it arrives to a lower level. It also shows that the first material has a better convergence property than the second material. Moreover, these two curves both indicate convergences. The red curve converges to a lower relative error and the blue curve converges to a higher relative error. The slopes of these two curves in the last several iterations confirm the convergence phenomenon.

To compare with the two-step method, we also show the results obtained using the coupled method. The reconstructed images are shown in Figure 6.5. Compared with Fig-

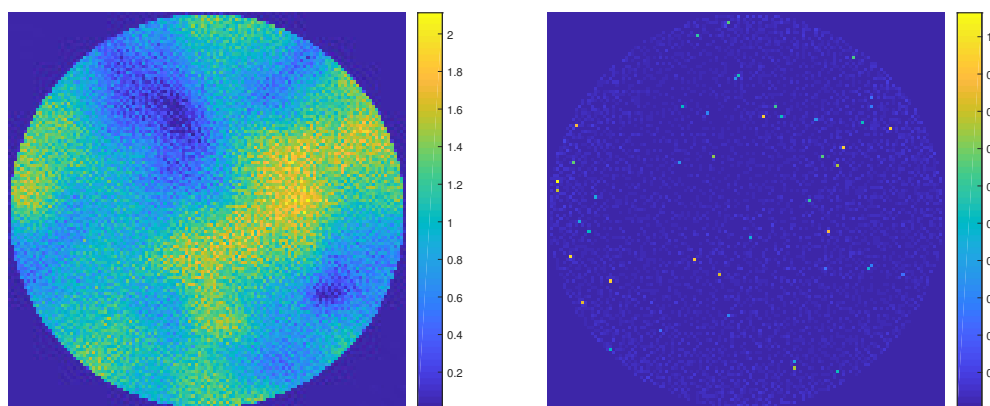


Figure 6.5: The reconstructed material maps for Plexiglass (Left) and PVC (Right) using the coupled method.

ure 6.3, these two images are more blurred and have more shades. For the first image, there are several blue dots scattering around the yellow areas and the yellow areas in the upper right corner and lower bottoms are not clear to identify. For the second material map, the locations of spikes are clear but the shades are more obvious compared with the reconstructed image obtained from the two-step method.

The relative errors corresponding to the iterative process of the coupled method is shown in Figure 6.6. From the slopes of curves, we can see that this method converges slower than

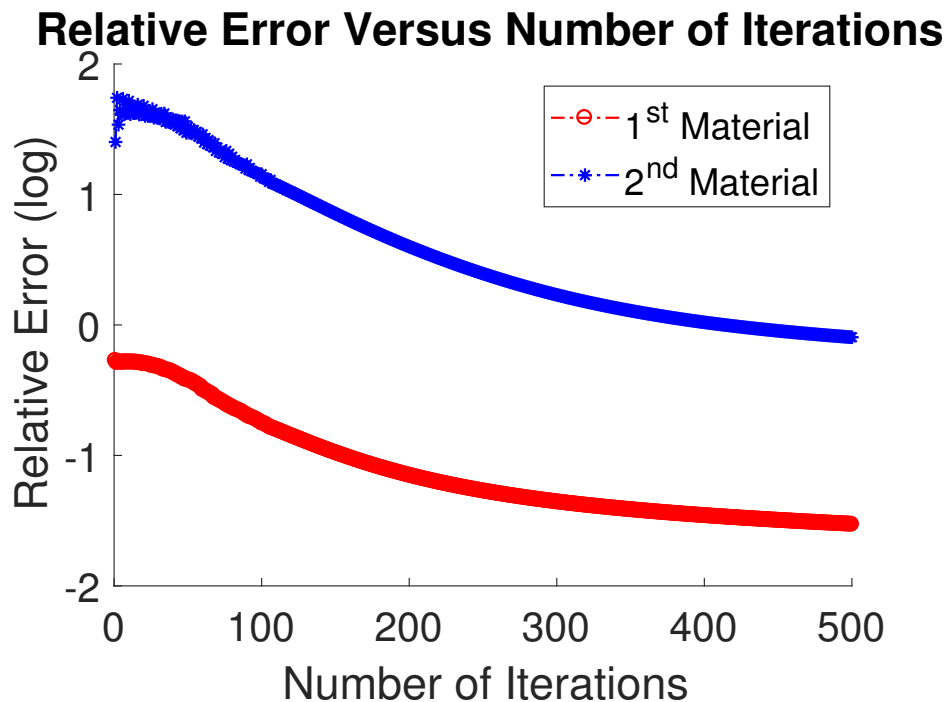


Figure 6.6: The relative errors of two material maps solved by the coupled method.

the previous two-step method and both materials stagnate at higher levels of relative errors. This phenomenon coincides with the observations we have found from the reconstructed images. This coupled method might generate a better solution to the optimization problem (6.2), but it might be a worse solution in terms of image quality. On the other hand, we can update the step alternatingly with this coupled method. In contrast, the updating framework of the two-step method does not have this property.

6.4 Conclusions and Remarks

In this chapter, we propose a two-step method that can solve the spectral CT model in sequence. In the first step, a row-wise model that is based on a weighted least squares term is used as a bridge to the second step. Since each row does not depend on each other, each tiny problem can also be solved in parallel. After we have obtained the results from the first step, we need to quantify the noise propagation and set up another least squares problem under bound constraints to compute the final solution. We implement FISTA with

projections to solve this problem to a high accuracy. Instead of using a two-step method, we can build an optimization problem that consists of these two terms and solve it directly. This coupled method is likely to provide us a better solution to the optimization problem. However, based on the numerical results, the two-step method beats the coupled method in terms of image quality, convergence speed and relative errors.

Compared with the method proposed in Chapter 5, this two-step method is more efficient when the product of three parameters, number of projections, range of angles and number of energy bins, is limited. However, for high-resolution images or 3D reconstruction, extensive projections are necessary to guarantee the image quality. Even if we can use parallel processing to reduce this influence, we should also choose different methods under proper circumstances.

Chapter 7

Conclusions and Future Works

Spectral computed tomography problems involve nonlinearity and require more efforts to obtain quantitative information. On the other hand, it can also offer material composition as well as images of higher quality with alleviation of artifacts. Based on different assumptions, spectral computed tomography problems display distinct forms and various math tools are used to compute solutions. For the simple energy discriminating model, we build a nonlinear optimization problem based on a Poisson likelihood estimator, and with this model, a nonlinear interior-point trust region method is introduced to obtain robust reconstruction. For the energy-windowed spectral CT model, we contribute to the optimization frameworks and preconditioners. First, we build a nonlinear least squares problem and solve it with a two-step method, projected line search plus the trust region method. We also propose an adaptive preconditioner to further mitigate the influence of small singular values. With the Gaussian assumption of noise, we transform the energy-windowed spectral CT model into a weighted least squares problem under bound constraints. An efficient preconditioner derived from the corresponding Hessian is presented to reduce the ill-posedness significantly. This preconditioner is inspired from the interlacing of the Kronecker products and diagonal matrices and is built by using a rank-1 approximation. Moreover, we introduce a subproblem to simplify the projection step that is used to keep nonnegativity. Even if we can calculate the solution of the spectral CT model using optimization and preconditioners, we can also consider further separation of the model. In this case, we propose a two-step method using

an auxiliary variable. The first step can be further decomposed into row slices, which can be solved in parallel. After quantifying the noise propagation, we solve the problem corresponding the second step with FISTA.

Even if we only focus on computing the solutions to spectral CT problems, other research areas deserve considerations. For example, the model we use is only based on a single energy source. We can also consider dual-source, or even three or more sources, CT models. The problem of reconstructing images from limited data is even more challenging, and it is important to think about how to find a robust reconstruction under this situation. Furthermore, we can try to combine the state-of-the-art machine learning techniques with spectral CT. Machine learning can be used to conduct data-preprocessing, predict similar features and validate the strength of proposed methods. To conclude, spectral CT is still an active and promising research area and it can be expected that new inspiring work will continue to be done for years to come.

Bibliography

- [1] R. E. ALVAREZ AND A. MACOVSKI, *Energy-selective reconstructions in x-ray computerised tomography*, *Physics in Medicine and Biology*, 21 (1976), p. 733.
- [2] R. F. BARBER, E. Y. SIDKY, T. G. SCHMIDT, AND X. PAN, *An algorithm for constrained one-step inversion of spectral CT data*, *Physics in Medicine and Biology*, 61 (2016), p. 3784.
- [3] A. BECK AND M. TEBoulLE, *A fast iterative shrinkage-thresholding algorithm for linear inverse problems*, *SIAM Journal on Imaging Sciences*, 2 (2009), pp. 183–202.
- [4] E. C. BECKMANN, *CT scanning the early days*, *The British Journal of Radiology*, (2014).
- [5] P. BOUGUER, *Essai d’optique sur la gradation de la lumière*, chez Claude Jombert, rue S. Jacques, au coin de la rue des Mathurins, à l , 1729.
- [6] V. M. BUSTAMANTE, *Iterative Polyenergetic Digital Tomosynthesis Reconstructions for Breast Cancer Screening*, PhD thesis, Emory University, 2013.
https://etd.library.emory.edu/file/view/pid/emory:d6wn4/mejia%20bustamante_dissertation.pdf.
- [7] V. M. BUSTAMANTE, J. G. NAGY, S. S. FENG, AND I. SECHOPOULOS, *Iterative breast tomosynthesis image reconstruction*, *SIAM Journal on Scientific Computing*, 35 (2013), pp. S192–S208.

- [8] R. BYRD, J. NOCEDAL, AND R. WALTZ, *KNITRO: An integrated package for nonlinear optimization*, Large-scale Nonlinear Optimization, (2006), pp. 35–59.
- [9] R. H. BYRD, M. E. HRIBAR, AND J. NOCEDAL, *An interior point algorithm for large-scale nonlinear programming*, SIAM Journal on Optimization, 9 (1999), pp. 877–900.
- [10] A. CHAMBOLLE AND T. POCK, *A first-order primal-dual algorithm for convex problems with applications to imaging*, Journal of Mathematical Imaging and Vision, 40 (2011), pp. 120–145.
- [11] J. CHUNG, J. G. NAGY, AND I. SECHOPOULOS, *Numerical algorithms for polyenergetic digital breast tomosynthesis reconstruction*, SIAM Journal on Imaging Sciences, 3 (2010), pp. 133–152.
- [12] R. E. CROCHIERE AND L. R. RABINER, *Multirate Digital Signal Processing*, vol. 18, Prentice-Hall Englewood Cliffs, NJ, 1983.
- [13] I. DAUBECHIES, M. DEFRISE, AND C. DE MOL, *An iterative thresholding algorithm for linear inverse problems with a sparsity constraint*, Communications on Pure and Applied Mathematics: A Journal Issued by the Courant Institute of Mathematical Sciences, 57 (2004), pp. 1413–1457.
- [14] S. R. DEANS, *The Radon Transform and Some of Its Applications*, Courier Corporation, MA, 2007.
- [15] I. A. ELBAKRI AND J. A. FESSLER, *Statistical image reconstruction for polyenergetic x-ray computed tomography*, IEEE Transactions on Medical Imaging, 21 (2002), pp. 89–99.
- [16] T. ELFVING, P. C. HANSEN, AND T. NIKAZAD, *Semi-convergence properties of kaczmarz's method*, Inverse Problems, 30 (2014), p. 055007.
- [17] H. ENGL, *Discrepancy principles for tikhonov regularization of ill-posed problems leading to optimal convergence rates*, Journal of Optimization Theory and Applications, 52 (1987), pp. 209–215.

- [18] C. L. EPSTEIN, *Introduction to the Mathematics of Medical Imaging*, SIAM, PA, 2007.
- [19] T. G. FLOHR, C. H. MCCOLLOUGH, H. BRUDER, M. PETERSILKA, K. GRUBER, C. SÜß, M. GRASRUCK, K. STIERSTORFER, B. KRAUSS, R. RAUPACH, ET AL., *First performance evaluation of a dual-source CT (DSCT) system*, *European Radiology*, 16 (2006), pp. 256–268.
- [20] J. FRIKEL AND E. T. QUINTO, *Characterization and reduction of artifacts in limited angle tomography*, *Inverse Problems*, 29 (2013), p. 125007.
- [21] D. GABAY AND B. MERCIER, *A Dual Algorithm for the Solution of Nonlinear Variational Problems via Finite Element Approximation*, Institut de recherche d’informatique et d’automatique, 1975.
- [22] S. GAZZOLA, P. C. HANSEN, AND J. G. NAGY, *IR Tools: a matlab package of iterative regularization methods and large-scale test problems*, *Numerical Algorithms*, (2018), pp. 1–39.
- [23] G. H. GOLUB, M. HEATH, AND G. WAHBA, *Generalized cross-validation as a method for choosing a good ridge parameter*, *Technometrics*, 21 (1979), pp. 215–223.
- [24] G. H. GOLUB AND C. F. VAN LOAN, *Matrix Computations*, vol. 3, JHU Press, MD, 2012.
- [25] R. GORDON, R. BENDER, AND G. T. HERMAN, *Algebraic reconstruction techniques (art) for three-dimensional electron microscopy and x-ray photography*, *Journal of Theoretical Biology*, 29 (1970), pp. 471–481.
- [26] M. GRANT, S. BOYD, AND Y. YE, *CVX: Matlab software for disciplined convex programming*, 2008.
- [27] M. C. GRANT AND S. P. BOYD, *Graph implementations for nonsmooth convex programs*, in *Recent Advances in Learning and Control*, Springer, 2008, pp. 95–110.
- [28] J. HADAMARD, *Sur les problèmes aux dérivées partielles et leur signification physique*, *Princeton University Bulletin*, (1902), pp. 49–52.

- [29] P. C. HANSEN AND J. S. JØRGENSEN, *AIR Tools II: algebraic iterative reconstruction methods, improved implementation*, Numerical Algorithms, (2017), pp. 1–31.
- [30] P. C. HANSEN, J. G. NAGY, AND D. P. O’LEARY, *Deblurring Images: Matrices, Spectra, and Filtering*, vol. 3, SIAM, PA, 2006.
- [31] P. C. HANSEN AND D. P. OLEARY, *The use of the l-curve in the regularization of discrete ill-posed problems*, SIAM Journal on Scientific Computing, 14 (1993), pp. 1487–1503.
- [32] P. C. HANSEN AND M. SAXILD-HANSEN, *AIR Tools: A MATLAB package of algebraic iterative reconstruction methods*, Journal of Computational and Applied Mathematics, 236 (2012), pp. 2167–2178.
- [33] B. J. HEISMANN, B. T. SCHMIDT, AND T. FLOHR, *Spectral Computed Tomography*, SPIE Bellingham, WA, 2012.
- [34] J. D. INGLE JR AND S. R. CROUCH, *Spectrochemical Analysis*, Prentice Hall College Book Division, NJ, 1988.
- [35] A. C. KAK AND M. SLANEY, *Principles of Computerized Tomographic Imaging*, SIAM, PA, 2001.
- [36] S. KARZMARZ, *Angenaherte auflösung von systemen linearer gleichungen*, Bulletin International de l’Académie Polonaise des Sciences et des Lettres, (1937), pp. 355–357.
- [37] C. T. KELLEY, *Iterative Methods for Optimization*, vol. 18, SIAM, PA, 1999.
- [38] J. H. LAMBERT, *Photometria sive de mensura et gradibus luminis, colorum et umbrae*, Klett, 1760.
- [39] R. M. LARSEN, *Lanczos bidiagonalization with partial reorthogonalization*, DAIMI Report Series, 27 (1998).
- [40] R. MATHIAS, *The spectral norm of a nonnegative matrix*, Linear Algebra and its Applications, 139 (1990), pp. 269–284.

- [41] J. MATTINGLEY AND S. BOYD, *Automatic code generation for real-time convex optimization*, Convex Optimization in Signal Processing and Communications, (2009), pp. 1–41.
- [42] J. MATTINGLEY AND S. BOYD, *Real-time convex optimization in signal processing*, IEEE Signal Processing Magazine, 27 (2010), pp. 50–61.
- [43] J. MATTINGLEY AND S. BOYD, *CVXGEN: A code generator for embedded convex optimization*, Optimization and Engineering, 13 (2012), pp. 1–27.
- [44] J. MATTINGLEY, Y. WANG, AND S. BOYD, *Code generation for receding horizon control*, in Computer-Aided Control System Design (CACSD), 2010 IEEE International Symposium on, IEEE, 2010, pp. 985–992.
- [45] J. L. MUELLER AND S. SILTANEN, *Linear and Nonlinear Inverse Problems with Practical Applications*, SIAM, PA, 2012.
- [46] A. S. NEMIROVSKY AND D. B. YUDIN, *Problem Complexity and Method Efficiency in Optimization*, Chichester: Wiley, 1983.
- [47] Y. E. NESTEROV, *A method for solving the convex programming problem with convergence rate $o(1/k^2)$* , in Doklady Akademii Nauk SSSR, vol. 269, 1983, pp. 543–547.
- [48] J. NOCEDAL AND S. WRIGHT, *Numerical Optimization*, Springer Science & Business Media, NY, 2006.
- [49] X. PAN, E. Y. SIDKY, AND M. VANNIER, *Why do commercial CT scanners still employ traditional, filtered back-projection for image reconstruction?*, Inverse problems, 25 (2009), p. 123009.
- [50] J. RADON, *1.1 über die bestimmung von funktionen durch ihre integralwerte längs gewisser mannigfaltigkeiten*, Classic Papers in Modern Diagnostic Radiology, 5 (2005), p. 21.

- [51] Y. SAAD AND M. H. SCHULTZ, *Gmres: A generalized minimal residual algorithm for solving nonsymmetric linear systems*, SIAM Journal on Scientific and Statistical Computing, 7 (1986), pp. 856–869.
- [52] J. SCHLOMKA, E. ROESSL, R. DORSCHIED, S. DILL, G. MARTENS, T. ISEL, C. BÄUMER, C. HERRMANN, R. STEADMAN, G. ZEITLER, ET AL., *Experimental feasibility of multi-energy photon-counting k-edge imaging in pre-clinical computed tomography*, Physics in Medicine & Biology, 53 (2008), p. 4031.
- [53] J. H. SIEWERDSEN, A. M. WAESE, D. J. MOSELEY, S. RICHARD, AND D. A. JAFFRAY, *Spektr: A computational tool for x-ray spectral analysis and imaging system optimization*, Medical Physics, 31 (2004), pp. 3057–3067.
- [54] T. STEihaug, *The conjugate gradient method and trust regions in large scale optimization*, SIAM Journal on Numerical Analysis, 20 (1983), pp. 626–637.
- [55] C. F. VAN LOAN, *The ubiquitous kronecker product*, Journal of Computational and Applied Mathematics, 123 (2000), pp. 85–100.
- [56] P. WOLFE, *Convergence conditions for ascent methods*, SIAM Review, 11 (1969), pp. 226–235.
- [57] V. S. K. YOKHANA, B. D. ARHATARI, T. E. GUREYEV, AND B. ABBEY, *Soft-tissue differentiation and bone densitometry via energy-discriminating x-ray microct*, Optics Express, 25 (2017), pp. 29328–29341.