

## SPECTRAL COMPUTED TOMOGRAPHY WITH LINEARIZATION AND PRECONDITIONING\*

YUNYI HU<sup>†</sup>, MARTIN S. ANDERSEN<sup>‡</sup>, AND JAMES G. NAGY<sup>†</sup>

**Abstract.** In the area of image sciences, the emergence of spectral computed tomography (CT) detectors highlights the concept of *quantitative imaging*, in which not only are reconstructed images offered, but weights of different materials that compose the object are also provided. If a detector is made up of several energy windows and each energy window is assumed to detect a specific range of energy spectrum, then a nonlinear matrix equation is formulated to represent the discretized process of attenuation of x-ray intensity. In this paper, we present a linearization technique to transform this nonlinear equation into an optimization problem that is based on a weighted least squares term and a nonnegative bound constraint. To solve this optimization problem, we propose a new preconditioner that can significantly reduce the condition number, and with this preconditioner, we implement a highly efficient first order method, the 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  $\ell_1$  regularization to stabilize the solution. With the introduction of new preconditioning, a linear inequality constraint is introduced. In each iteration, we decompose this constraint into small-sized problems that can be solved with fast optimization solvers. Numerical experiments illustrate convergence, effectiveness, and significance of the proposed method.

**Key words.** preconditioning, image reconstruction, tomography, FISTA

**AMS subject classifications.** 65F22, 65F10, 49N45, 65K99

**DOI.** 10.1137/18M1194419

**1. Introduction.** The development of new energy-windowed spectral computed tomography (CT) machines has received a great deal of interest in recent years; see, e.g., [1, 24]. 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 that the data collected by the detector are nonnegative integers. Compared with traditional CT machines, we can avoid introducing beam-hardening artifacts [19] and improve quality of reconstructed images. To reconstruct images of an object, we need to solve a nonlinear equation

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

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 elementwise (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

\*Received by the editors June 14, 2018; accepted for publication (in revised form) June 18, 2019; published electronically October 29, 2019.

<https://doi.org/10.1137/18M1194419>

**Funding:** This work was supported by the U.S. National Science Foundation under grant DMS-1819042.

<sup>†</sup>Department of Mathematics and Computer Science, Emory University, Atlanta, GA 30322 (yunyi.hu@emory.edu, jnagy@emory.edu).

<sup>‡</sup>Department of Applied Mathematics and Computer Science, Technical University of Denmark, Kgs. Lyngby, 2800, Denmark (mskan@dtu.dk).

that  $\mathbf{S}$  is square and invertible. Moreover,  $\mathbf{\mathcal{E}}$  represents the noise term and we assume that  $E_{il} \sim \mathcal{N}(0, y_{il})$  for each component  $E_{il}$  in  $\mathbf{\mathcal{E}}$  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 the size  $N_v \times N_m$ , 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  $\mathbf{W} \geq \mathbf{0}$ .

To solve (1.1), we want to vectorize it first. Then we use the Taylor expansion to remove the pointwise exponential function and obtain an approximate linearized equation. Under the Gaussian assumption, as we show in section 2, we can transform this equation into a weighted least squares problem under bound constraints:

$$(1.2) \quad \begin{aligned} \min_{\mathbf{w}} \quad & \frac{1}{2} \|\mathcal{A}\mathbf{w} - \mathbf{b}\|_{\Sigma^{-1}}^2 \\ \text{subject to} \quad & \mathbf{w} \geq \mathbf{0}, \end{aligned}$$

where  $\mathcal{A} = \mathbf{C} \otimes \mathbf{A}$ ,  $\mathbf{b} = -\log(\mathbf{y})$ ,  $\mathbf{y} = \text{vec}(\mathbf{Y})$ , and  $\mathbf{w} = \text{vec}(\mathbf{W})$ .  $\Sigma^{-1}$ , which combines information from  $\mathbf{S}$  and  $\mathbf{y}$ , is the inverse covariance matrix generated by Gaussian noise and 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  $N_e \times N_m$ , where  $N_e$  is the number of energy bins 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, problem (1.2) is similar to a quadratic programming problem under bound constraints. However, direct implementation of an optimization solver 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.

Because of the ill-posedness, Barber et al. [1] proposed a preconditioner based on the eigenvalue decomposition of the matrix product of linear attenuation coefficients,  $\mathbf{C}^T \mathbf{C}$ , to orthogonalize columns of  $\mathbf{C}$ . They also suggest using a Poisson noise assumption and construct loss functions that are based on either the maximum likelihood estimator (MLE) or the nonlinear least squares term. Using these types of loss functions and the proposed preconditioner, a *Chambolle-Pock* primal-dual method [5] is implemented to solve the corresponding optimization problem. However, because the MLE for the Poisson model is nonlinear, it is not obvious how this preconditioner can reduce the condition number of the Hessian matrix. Moreover, because each iteration of a second order method for large 3D imaging problems is very costly (in terms of both the computations and storage requirements), in this paper we consider first order methods. With a first order method, it is not necessary to construct either the Hessian or Hessian-vector multiplication in each step.

To mitigate the ill-posedness, we propose a new preconditioner that is based on a rank-1 approximation of the matrix  $\mathbf{Y}$ . With this rank-1 approximation, we can estimate the Hessian of the objective function in (1.2) by 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 an 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 an identity,

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 future iterations. The preconditioner proposed in [1] includes only the data of  $\mathbf{C}$ , the matrix of linear attenuation coefficients of material and energy. Compared with this, the preconditioner proposed in this paper includes the information of linear attenuation coefficients, the energy spectrum, and photon counting data. It offers 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, the Fast Iterative Shrinkage-Thresholding Algorithm (FISTA) [2, 21, 20] 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 a form of  $f(\mathbf{x}) + g(\mathbf{x})$  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 (1.2), and  $g(\mathbf{x})$  can represent a nonsmooth regularization term such as  $\ell_1$  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 paper is organized as follows. In section 2, we review the continuous energy-windowed spectral CT model and the corresponding discretized nonlinear matrix equation. The linearization, vectorization, and setup of the optimization problem are also included in section 2. The key idea of this paper, preconditioning, is introduced in section 3. 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 4, we study FISTA and how we construct and solve the projection problems. Moreover, numerical experiments are presented in section 5, and concluding remarks are given in section 6.

**2. The energy-windowed spectral CT model.** In this section, we start with an introduction to the basic model. Then we show how to discretize this model to obtain a matrix equation. Since we do not want to solve this matrix equation directly, we therefore vectorize it and take the Taylor expansion to the first order term to remove the exponential function. In this case, we can obtain a linear system with transformed noise. With this transformed noise, we can build a weighted least squares optimization problem under bound constraints.

In CT, source x-ray beams are composed of a spectrum of different energies [4]. 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 for image reconstruction uses Beer’s

law [12], which states that the change in x-ray intensity before and after illumination through the object is

$$(2.1) \quad y_i^{(k)} = \int_E S^{(k)}(e) \exp \left( - \int_{t \in l} \mu(\vec{r}(t), e) dt \right) de + \eta_i^{(k)}, \quad \begin{cases} i = 1, 2, \dots, N_d \times N_p, \\ k = 1, 2, \dots, N_b, \end{cases}$$

where the following hold:

- $y_i^{(k)}$  is the x-ray intensity of the  $i$ th pixel in the  $k$ th detector bin.
- $E$  is the photon flux density. Figure 5.2 shows a curve of  $E$  versus photon energy.
- $N_d$  is the number of detector pixels. For a material map of the size  $n \times n$ , we assume  $N_d = n$ .
- $N_p$  is the number of projections. For cone/fan beam CT, projections are uniformly distributed from 0 to 360 degrees.
- $N_b$  is the number of detector bins. For an energy-windowed CT machine, we usually assume that it has 5 to 6 energy bins.
- $S^{(k)}(e)$  represents photon flux density for the  $k$ th detector bin, which is the number of incident photons at the energy  $e$  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 the energy level  $e$ .
- $\eta_i^{(k)}$  is the error term for the  $i$ th element in the  $k$ th energy bin, and it is assumed to be Gaussian for this model.

In (2.1), the unknown linear attenuation coefficient  $\mu(\vec{r}(t), e)$  is dependent on the position function  $r(t)$  and the 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  $\mu(\vec{r}(t), e)$  [11]:

$$(2.2) \quad \mu(\vec{r}(t), e) = \sum_{m=1}^{N_m} u_{m,e} w_m(\vec{r}),$$

where the following hold:

- $N_m$  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 the weight fraction  $w_m(\vec{r})$ . If we also assume that  $w_m(\vec{r})$  can be represented as a sum of product of weights and basis functions  $\phi_j(\vec{r})$ , then another expansion can be expressed by

$$(2.3) \quad w_m(\vec{r}) = \sum_{j=1}^{N_v} w_{j,m} \phi_j(\vec{r}),$$

where the following hold:

- $N_v$  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 the number of projections  $N_p$ :

$$(2.4) \quad a_{i,j} = \int_{t \in l} \phi_j(\vec{r}(t)) \, dt.$$

Then the line integral in (2.1) can be simplified by expansion (2.3) and integral (2.4):

$$(2.5) \quad \int_{t \in l} \mu(\vec{r}(t), e) \, dt = \sum_{m=1}^{N_m} \sum_{j=1}^{N_v} u_{m,e} w_{j,m} \int_{t \in l} \phi_j(\vec{r}(t)) \, dt = \sum_{j=1}^{N_v} \sum_{m=1}^{N_m} a_{i,j} w_{j,m} u_{m,e}.$$

If we also discretize the integral over the energy  $E$  and ignore quadrature errors, then the discrete model of (2.1) can be written as

$$(2.6) \quad y_i^{(k)} = \sum_{e=1}^{N_e} s_e^{(k)} \exp \left( - \sum_{j=1}^{N_v} \sum_{m=1}^{N_m} a_{i,j} w_{j,m} u_{m,e} \right) + \eta_i^{(k)},$$

where  $N_e$  is the number of discrete energies. If we collect  $a_{i,j}$ ,  $w_{i,j}$ , and  $u_{m,e}$  in a matrix form and concatenate  $y_i^{(k)}$ ,  $s_e^{(k)}$ ,  $\eta_i^{(k)}$  with respect to their energy windows, then the corresponding matrix equation of (2.6) can be represented as

$$(2.7) \quad \mathbf{Y} = \exp(-\mathbf{A}\mathbf{W}\mathbf{C}^T) \mathbf{S} + \mathbf{\mathcal{E}},$$

where the following hold:

- $\mathbf{Y}$  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 energy  $e$ , and the  $m$ th material.
- $\mathbf{S}$  is a matrix of the size  $N_e \times N_b$  and each column collects the spectrum energy of a specific range. In the forward problem, we use the full spectrum, but when we solve the inverse problem, the average in each energy window is used to represent the corresponding spectral energy. Therefore,  $N_b = N_e$  for the inverse problem and  $\mathbf{S}$  is an invertible diagonal matrix because the means are placed in the diagonal. A detailed example is shown in Figure 5.2.
- $\mathbf{\mathcal{E}}$  is the noise matrix that is of the size  $(N_d \cdot N_p) \times N_b$ . The assumption for the noise is  $E_{il} \sim \mathcal{N}(0, y_{il})$  for each element  $E_{il}$  in  $\mathbf{\mathcal{E}}$  and  $y_{il}$  in  $\mathbf{Y}$ .

In (2.7), the exponential operator is applied elementwise (i.e., it is not a matrix function). In addition to (2.7), we also require that weight fractions should be nonnegative, and this can be illustrated by the constraint  $\mathbf{W} \geq \mathbf{0}$ .

In several cases, the composition of materials can be similar. For example, glandular and adipose have similar attenuation coefficients at the same energy level, which causes collinearity. After discretization, the columns of  $\mathbf{C}$  can be nearly dependent. Moreover,  $\mathbf{A}$  is large scale and sparse and is highly likely to have small singular values. As we will see later, the Hessian system involves the Kronecker product  $\mathbf{C} \otimes \mathbf{A}$ , and it can cause ill-posedness. Since it is challenging to solve this equation directly, it is important to consider approaches that facilitate the process. First, we can introduce a preconditioning matrix  $\mathbf{M}$  into (2.7):

$$(2.8) \quad \mathbf{Y} = \exp(-\mathbf{A}\mathbf{W}\mathbf{M}^{-T}\mathbf{M}^T\mathbf{C}^T) \mathbf{S} + \mathbf{\mathcal{E}}.$$

If we let  $\tilde{\mathbf{W}} = \mathbf{W}\mathbf{M}^{-T}$  and  $\tilde{\mathbf{C}} = \mathbf{C}\mathbf{M}$ , then (2.8) is equivalent to

$$(2.9) \quad \mathbf{Y} = \exp\left(-\mathbf{A}\tilde{\mathbf{W}}\tilde{\mathbf{C}}^T\right)\mathbf{S} + \boldsymbol{\varepsilon}.$$

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 3.1, 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 (2.9) directly because it might introduce a tensor when we compute second order derivatives. In this case, we want to vectorize (2.9) 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 then 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,  $N_b = N_e$  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 (2.9):

$$(2.10) \quad \mathbf{Y}\mathbf{S}^{-1} = \exp\left(-\mathbf{A}\tilde{\mathbf{W}}\tilde{\mathbf{C}}^T\right) + \boldsymbol{\varepsilon}\mathbf{S}^{-1}.$$

Vectorizing both sides of (2.10), and using properties of Kronecker products, we obtain

$$(2.11) \quad (\mathbf{S}^{-T} \otimes \mathbf{I})\mathbf{y} = \exp\left\{-\left(\tilde{\mathbf{C}} \otimes \mathbf{A}\right)\tilde{\mathbf{w}}\right\} + (\mathbf{S}^{-T} \otimes \mathbf{I})\mathbf{e},$$

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 (2.11) and obtain

$$(2.12) \quad \tilde{\mathbf{y}} - \tilde{\mathbf{e}} = \exp\left\{-\left(\tilde{\mathbf{C}} \otimes \mathbf{A}\right)\tilde{\mathbf{w}}\right\}.$$

By taking the logarithm on both sides of (2.12), we can obtain a linear equation

$$(2.13) \quad \log(\tilde{\mathbf{y}} - \tilde{\mathbf{e}}) = -\left(\tilde{\mathbf{C}} \otimes \mathbf{A}\right)\tilde{\mathbf{w}}.$$

However, the left-hand side of (2.13) 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 at  $\tilde{\mathbf{y}}$ :

$$(2.14) \quad \log(\tilde{\mathbf{y}} - \tilde{\mathbf{e}}) = \log(\tilde{\mathbf{y}}) - \text{diag}(\tilde{\mathbf{y}})^{-1}\tilde{\mathbf{e}} + \mathcal{O}(\|\tilde{\mathbf{e}}\|_2^2).$$

If we use the first two terms on the right-hand side of (2.14) to estimate the term  $\log(\tilde{\mathbf{y}} - \tilde{\mathbf{e}})$ , then (2.13) can be expressed by a linear equation with the error term  $\text{diag}(\tilde{\mathbf{y}})^{-1}\tilde{\mathbf{e}}$ . Letting  $\mathbf{b} = -\log(\tilde{\mathbf{y}})$ , then (2.13) is approximately equal to

$$(2.15) \quad \mathbf{b} \approx \left(\tilde{\mathbf{C}} \otimes \mathbf{A}\right)\tilde{\mathbf{w}} - \text{diag}(\tilde{\mathbf{y}})^{-1}\tilde{\mathbf{e}}.$$

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

$$(2.16) \quad \mathbf{b}|\tilde{\mathbf{w}} \sim \mathcal{N}\left(\left(\tilde{\mathbf{C}} \otimes \mathbf{A}\right)\tilde{\mathbf{w}}, \boldsymbol{\Sigma}\right),$$

where the noise covariance matrix  $\Sigma$  is expressed by

$$(2.17) \quad \Sigma = \text{diag}(\tilde{\mathbf{y}})^{-1} (\mathbf{S}^{-T} \otimes \mathbf{I}) \text{diag}(\mathbf{y}) (\mathbf{S}^{-1} \otimes \mathbf{I}) \text{diag}(\tilde{\mathbf{y}})^{-1},$$

and the inverse covariance matrix is given by

$$(2.18) \quad \Sigma^{-1} = \text{diag}(\tilde{\mathbf{y}}) (\mathbf{S} \otimes \mathbf{I}) \text{diag}(\mathbf{y})^{-1} (\mathbf{S}^T \otimes \mathbf{I}) \text{diag}(\tilde{\mathbf{y}}).$$

Since  $\mathbf{Y}$  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. As long as the noise does not dominate the projected data, we expect the entries of  $\tilde{\mathbf{y}}$  will be larger than zero. From expression (2.18), we can see that the structure of  $\Sigma^{-1}$  depends on the structure of the matrix  $\mathbf{S}$ . If  $\mathbf{S}$  is diagonal, then  $\Sigma$  is also diagonal. If we let  $\mathcal{A} = \tilde{\mathbf{C}} \otimes \mathbf{A}$ , then (see, e.g., [3]) the best unbiased linear estimator of  $\tilde{\mathbf{w}}$  for the Gaussian model (2.16) is the solution of

$$(2.19) \quad \min_{\tilde{\mathbf{w}}} \frac{1}{2} (\mathcal{A}\tilde{\mathbf{w}} - \mathbf{b})^T \Sigma^{-1} (\mathcal{A}\tilde{\mathbf{w}} - \mathbf{b}).$$

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

$$(2.20) \quad \begin{aligned} \min_{\tilde{\mathbf{w}}} \quad & \frac{1}{2} \|\mathcal{A}\tilde{\mathbf{w}} - \mathbf{b}\|_{\Sigma^{-1}}^2 \\ \text{subject to} \quad & (\mathbf{M} \otimes \mathbf{I}) \tilde{\mathbf{w}} \geq \mathbf{0}. \end{aligned}$$

In (2.20) the norm  $\|\cdot\|_{\Sigma^{-1}}^2$  corresponds to the weighted inner product given in (2.19). From this expression, we know that the objective function is convex. Moreover, the inverse covariance matrix  $\Sigma^{-1}$  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. First, we need to choose an appropriate preconditioning matrix to reduce the ill-conditioning of the Hessian. Second, we want to select suitable regularizations for the corresponding materials. Third, we have to find an efficient method for solving the constrained weighted least squares problem. These three challenges are related to each other, and an appropriate preconditioner with appropriate regularizations will be beneficial for the solver efficiency. Finally, we should handle linear inequality constraints in an efficient way. We will address these four challenges in the following sections.

### 3. Preconditioning and regularization.

**3.1. Preconditioning.** The choice of the preconditioning matrix  $\mathbf{M}$  is crucial for solving the optimization problem (2.20). If we do not have a preconditioner or we choose the preconditioner  $\mathbf{M}$  as the identity, the original Hessian for the weighted least squares problem (2.20) is expressed by

$$(3.1) \quad \mathbf{H} = (\mathbf{C}^T \otimes \mathbf{A}^T) \Sigma^{-1} (\mathbf{C} \otimes \mathbf{A}).$$

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

$$(3.2) \quad \tilde{\mathbf{H}} = \mathbf{A}^T \Sigma^{-1} \mathbf{A} = (\tilde{\mathbf{C}}^T \otimes \mathbf{A}^T) \Sigma^{-1} (\tilde{\mathbf{C}} \otimes \mathbf{A}),$$

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 separate the noise covariance matrix  $\Sigma^{-1}$  into a Kronecker product of two terms, we can merge several terms using properties of the Kronecker product and transform parts of the Hessian into an identity with the help of  $\mathbf{M}$ . To realize this idea, we review the expression of  $\Sigma^{-1}$  in (2.18), 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, which 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

$$(3.3) \quad \text{diag}\{\mathbf{y}\} \approx \text{diag}\{\text{vec}(\mathbf{u}\mathbf{v}^T)\} = \text{diag}\{\mathbf{v}\} \otimes \text{diag}\{\mathbf{u}\}.$$

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

$$(3.4) \quad \min_{\mathbf{u}, \mathbf{v}} \|\mathbf{Y} - \mathbf{u}\mathbf{v}^T\|_F.$$

The solution to this problem has been studied extensively [23]. Using the singular value decomposition (SVD), one solution to (3.4) 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  $\sigma_1$  is the corresponding largest singular value of  $\mathbf{Y}$ . Since we only need these terms rather than a full SVD, we can use the MATLAB function `svds`, or other efficient approaches, such as PROPACK [14], to calculate only  $\sigma_1$ ,  $\mathbf{u}_1$ , 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 (3.3). In addition, the term  $\text{diag}\{\tilde{\mathbf{y}}\}$  can be represented by

$$(3.5) \quad \begin{aligned} \text{diag}\{\tilde{\mathbf{y}}\} &= \text{diag}\{(\mathbf{S}^{-T} \otimes \mathbf{I}) \text{vec}(\mathbf{Y})\} \approx \text{diag}\{(\mathbf{S}^{-T} \otimes \mathbf{I}) \text{vec}(\mathbf{u}\mathbf{v}^T)\} \\ &= \text{diag}\{\text{vec}(\mathbf{u}\mathbf{v}^T \mathbf{S}^{-1})\} = \text{diag}\{\mathbf{S}^{-T} \mathbf{v}\} \otimes \text{diag}\{\mathbf{u}\}. \end{aligned}$$

If we substitute the terms in (3.3) and (3.5) for the same terms in (2.18), we can obtain that

$$(3.6) \quad \Sigma^{-1} \approx \left( \text{diag}\{\mathbf{S}^{-T} \mathbf{v}\} \mathbf{S} \text{diag}\{\mathbf{v}\}^{-1} \mathbf{S}^T \text{diag}\{\mathbf{S}^{-T} \mathbf{v}\} \right) \otimes \text{diag}\{\mathbf{u}\}.$$

So the preconditioned Hessian matrix is given by

$$(3.7) \quad \begin{aligned} \tilde{\mathbf{H}} &= (\tilde{\mathbf{C}}^T \otimes \mathbf{A}^T) \Sigma^{-1} (\tilde{\mathbf{C}} \otimes \mathbf{A}) \\ &\approx (\tilde{\mathbf{C}}^T \otimes \mathbf{A}^T) \left[ \text{diag}\{\mathbf{S}^{-T} \mathbf{v}\} \mathbf{S} \text{diag}\{\mathbf{v}\}^{-1} \mathbf{S}^T \text{diag}\{\mathbf{S}^{-T} \mathbf{v}\} \otimes \text{diag}\{\mathbf{u}\} \right] (\tilde{\mathbf{C}} \otimes \mathbf{A}) \\ &= (\tilde{\mathbf{C}}^T \text{diag}\{\mathbf{S}^{-T} \mathbf{v}\} \mathbf{S} \text{diag}\{\mathbf{v}\}^{-1} \mathbf{S}^T \text{diag}\{\mathbf{S}^{-T} \mathbf{v}\} \tilde{\mathbf{C}}) \otimes (\mathbf{A}^T \text{diag}\{\mathbf{u}\} \mathbf{A}). \end{aligned}$$

Since the size of  $\tilde{\mathbf{C}}$  is  $N_e \times N_m$ , then the first part of the Kronecker product in (3.7) is a square matrix of the size  $N_m \times N_m$ . In other words, this part only depends on the number of materials that compose the object. Usually, we only consider 2 or 3 materials for the object so that the size of the matrix products for this



part is usually either  $2 \times 2$  or  $3 \times 3$ . 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  $\mathbf{C}^T \text{diag}\{\mathbf{S}^{-T}\mathbf{v}\} \mathbf{S} \text{diag}\{\mathbf{v}\}^{-1} \mathbf{S}^T \text{diag}\{\mathbf{S}^{-T}\mathbf{v}\} \mathbf{C}$  is a symmetric positive definite matrix. Therefore, we can calculate  $\mathbf{M}$  with the Cholesky decomposition:

$$(3.8) \quad \mathbf{C}^T \text{diag}\{\mathbf{S}^{-T}\mathbf{v}\} \mathbf{S} \text{diag}\{\mathbf{v}\}^{-1} \mathbf{S}^T \text{diag}\{\mathbf{S}^{-T}\mathbf{v}\} \mathbf{C} = \mathbf{G}^T \mathbf{G},$$

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 (3.7), we see that the preconditioned Hessian,  $\tilde{\mathbf{H}}$ , is dependent on a Kronecker product of two parts, and the first part has been transformed into an identity. 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 the SVD 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  $\sigma_{1\max}, \sigma_{1\min}$  and  $\sigma_{2\max}, \sigma_{2\min}$ , respectively; then the condition number of the original Hessian,  $\kappa(\mathbf{H})$ , can be estimated as

$$(3.9) \quad \kappa(\mathbf{H}) = \frac{\sigma_{1\max} \sigma_{2\max}}{\sigma_{1\min} \sigma_{2\min}}.$$

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

$$(3.10) \quad \kappa(\tilde{\mathbf{H}}) = \frac{\sigma_{2\max}}{\sigma_{2\min}}.$$

Since the fraction  $\sigma_{1\max}/\sigma_{1\min}$  is most likely to be significantly greater than 1, the condition number of  $\tilde{\mathbf{H}}$  is likely to be much smaller than  $\mathbf{H}$ . To validate this phenomenon, we can build a numerical example to compare the condition numbers. For an object that is composed of two materials, with each material map of the size  $16 \times 16$ , 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 3.1.

TABLE 3.1  
Comparison of condition numbers.

Matrix types	Condition numbers
Original Hessian	2.00 e+06
Preconditioned Hessian	2.59 e+04

From Table 3.1, we can see that the difference between  $\kappa(\mathbf{H})$  and  $\kappa(\tilde{\mathbf{H}})$  is around 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. With a better-conditioned system, we can compute the solution 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.

**3.2. Regularization.** With the help of our preconditioner, we can speed up an optimization algorithm and achieve higher accuracy. To further alleviate the 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 should 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 the generalized Tikhonov regularization to smooth the edges. For other materials, we choose the  $\ell_1$  regularization to penalize the sum of weights. Based on this idea, we can represent the regularization term as a sum of  $m$  parts:

$$(3.11) \quad R(\mathbf{w}) = \sum_{i=1}^m \frac{\alpha_i}{2} R_i(\mathbf{w}_i),$$

where  $\mathbf{w}_i$  is the vectorization form of the  $i$ th weight matrix,  $R_i(\mathbf{w}_i)$  is the corresponding regularization term, and  $\alpha_i$  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. In breast imaging, if the object is dominated by glandular and adipose tissue, 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. 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 regularizations for the dominating materials (e.g., glandular and adipose tissue) and an  $\ell_1$  regularization for the sparse material. We illustrate this with two materials, one that dominates, and one that is sparse:

$$(3.12) \quad R(\mathbf{w}) = \frac{\alpha_1}{2} \|\mathbf{L}\mathbf{w}_1\|_2^2 + \frac{\alpha_2}{2} \|\mathbf{w}_2\|_1.$$

If we add these regularization terms to the objective function in (2.20), we can rewrite it as an augmented system:

$$(3.13) \quad \min_{\tilde{\mathbf{w}}} \left\| \begin{bmatrix} \frac{\sqrt{2}}{2} \Sigma^{-\frac{1}{2}} (\tilde{\mathbf{C}} \otimes \mathbf{A}) \\ \sqrt{\frac{\alpha_1}{2}} \tilde{\mathbf{L}} \end{bmatrix} \tilde{\mathbf{w}} - \begin{bmatrix} \Sigma^{-\frac{1}{2}} \mathbf{b} \\ \mathbf{0} \end{bmatrix} \right\|_2^2 + \frac{\alpha_2}{2} [\mathbf{0} \quad \mathbf{1}] (\mathbf{M} \otimes \mathbf{I}) \tilde{\mathbf{w}} \\ \text{subject to} \quad (\mathbf{M} \otimes \mathbf{I}) \tilde{\mathbf{w}} \geq \mathbf{0},$$

where  $\tilde{\mathbf{L}} = [\mathbf{L} \quad \mathbf{0}] (\mathbf{M} \otimes \mathbf{I})$ . 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 non-smooth. Because of these properties, we can think about using FISTA [2] 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.

**4. FISTA and projections.** In this section, we first briefly present the main algorithm FISTA. To implement FISTA to solve the target optimization problem, we need to determine the step size and handle the nonnegative constraints. For the step size, we introduce how to compute the Lipschitz constant numerically and then 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 a delicate decomposition and efficient algorithms.

**4.1. FISTA.** The Fast Iterative Shrinkage-Thresholding Algorithm (FISTA) is a first order method that belongs to the family of Iterative Shrinkage-Thresholding Algorithms (ISTAs). This method was 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. Moreover, it is very appropriate for problems in imaging science because it is usually used to solve the nonsmooth convex problem

$$(4.1) \quad \min_{\mathbf{x}} \quad f(\mathbf{x}) + g(\mathbf{x}),$$

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 an  $\ell_1$  penalty or a total variation regularization. For problem (3.13), we construct an augmented loss function that merges the generalized Tikhonov regularization term, which corresponds to  $f(\mathbf{x})$  in (4.1). For the regularization term, the  $\ell_1$  regularization is nonsmooth but convex, and this matches  $g(\mathbf{x})$  in (4.1).

The details of this algorithm are shown in Algorithm 4.1. For the main algorithm, we need to first compute the smallest Lipschitz constant  $K$ . Then we can update the current step using FISTA. Because of the linear inequality constraints, we need to project the new step onto these constraints to keep the solution feasible. We would like to implement FISTA with a constant step size to solve the optimization problem (3.13). To implement this method, we need several preparations, which we will discuss in the following sections.

---

**Algorithm 4.1** FISTA and projections [2].

---

- 1: *Initialization:*
  - 2: Calculate the smallest Lipschitz constant  $K$  in (4.3) by the Power Method.
  - 3: Set up the 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 (4.2);
  - 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 (4.6);
  - 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:    $\mathbf{x}_{old} = \mathbf{x}_{new}$ .
-

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

$$(4.2) \quad f(\tilde{\mathbf{w}}) = \left\| \begin{bmatrix} \frac{\sqrt{2}}{2} \Sigma^{-\frac{1}{2}} (\tilde{\mathbf{C}} \otimes \mathbf{A}) \\ \sqrt{\frac{\alpha_1}{2}} \tilde{\mathbf{L}} \end{bmatrix} \tilde{\mathbf{w}} - \begin{bmatrix} \Sigma^{-\frac{1}{2}} \mathbf{b} \\ \mathbf{0} \end{bmatrix} \right\|_2^2,$$

$$g(\tilde{\mathbf{w}}) = \frac{\alpha_2}{2} [\mathbf{0} \quad \mathbf{1}] (\mathbf{M} \otimes \mathbf{I}) \tilde{\mathbf{w}},$$

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,

$$(4.3) \quad K = \lambda_{\max} \left[ \left( \tilde{\mathbf{C}}^T \otimes \mathbf{A}^T \right) \Sigma^{-1} (\tilde{\mathbf{C}} \otimes \mathbf{A}) + \alpha_1 \tilde{\mathbf{L}}^T \tilde{\mathbf{L}} \right].$$

Since we only need the largest eigenvalue, it is not necessary for us to construct these matrices explicitly; instead we can use an iterative approach, such as the Power Method [6]. Note that we only need to calculate  $K$  once for all FISTA iterations. The details are shown in Algorithm 4.2.

---

**Algorithm 4.2** Power Method [6].

---

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

**4.3. Projections.** In addition to the largest eigenvalue, 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 (3.13) as a quadratic programming problem under these specific constraints. To impose the linear inequality constraints, we can construct another quadratic programming problem that offers a nearest solution to satisfy these constraints. If we assume that we have obtained  $\tilde{\mathbf{w}}_k$  in the  $k$ th step, then we build a projection problem of the form

$$(4.4) \quad \begin{aligned} \min_{\tilde{\mathbf{w}}_{new}} \quad & \|\tilde{\mathbf{w}}_{new} - \tilde{\mathbf{w}}_k\|_2^2 \\ \text{subject to} \quad & (\mathbf{M} \otimes \mathbf{I}) \tilde{\mathbf{w}}_{new} \geq \mathbf{0}. \end{aligned}$$

For small- and medium-sized problems, we can solve it efficiently by direct implementation of standard optimization algorithms. For example, we can use CVX [7, 8] to solve problem (4.4), which turns out to be low-cost both in storage and calculation 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 for decomposing problem (4.4) into small pieces and try to solve each small problem accurately and efficiently.

Suppose we reshape vectors into matrices, for example, using the MATLAB `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 (4.4) is equivalent to

$$(4.5) \quad \begin{aligned} & \min_{\tilde{\mathbf{W}}_{new}} \quad \left\| \tilde{\mathbf{W}}_{new} - \tilde{\mathbf{W}}_k \right\|_F^2 \\ & \text{subject to} \quad \tilde{\mathbf{W}}_{new} \mathbf{M}^T \geq \mathbf{0}. \end{aligned}$$

If we focus on each row of  $\tilde{\mathbf{W}}_k$ ,  $\tilde{\mathbf{W}}_k(i, :)$ , then problem (4.5) can be rewritten as

$$(4.6) \quad \begin{aligned} & \min_{\tilde{\mathbf{W}}_{new}} \quad \sum_{k=1}^{N_v} \left\| \tilde{\mathbf{W}}_{new}(i, :) - \tilde{\mathbf{W}}_k(i, :) \right\|_2^2 \\ & \text{subject to} \quad \tilde{\mathbf{W}}_{new}(i, :) \mathbf{M}^T \geq \mathbf{0}, \end{aligned}$$

where  $\tilde{\mathbf{W}}_{new}(i, :)$  is the corresponding  $i$ th row in  $\tilde{\mathbf{W}}_{new}$ . It is obvious that this problem is separable, and the original problem (4.5) 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 each problem is usually  $2 \times 1$  or  $3 \times 1$ . In this case, we can solve each small-sized problem efficiently and concatenate the solutions into a large matrix. 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 paper, we choose CVXGEN [15, 16, 17, 18] to generate a customized solver for small quadratic programming problems. It is a problem-specific, fast, and accurate code generator which can achieve advanced 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-sized 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 by FISTA in the main loop. Finally, the projection problem is decomposed into many small pieces, and each can be solved by CVXGEN efficiently.

**5. 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  $128 \times 128$ . 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 50. Outside of the circle, we assume that there exist no weights of the object. These two images are shown in Figure 5.1.

Inside the mask, the darker blue areas for the first material map are mainly located in the upper left and lower right corners, which correspond to blank points. Other areas inside the circle are represented by heavily weighted yellows and greens. 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.

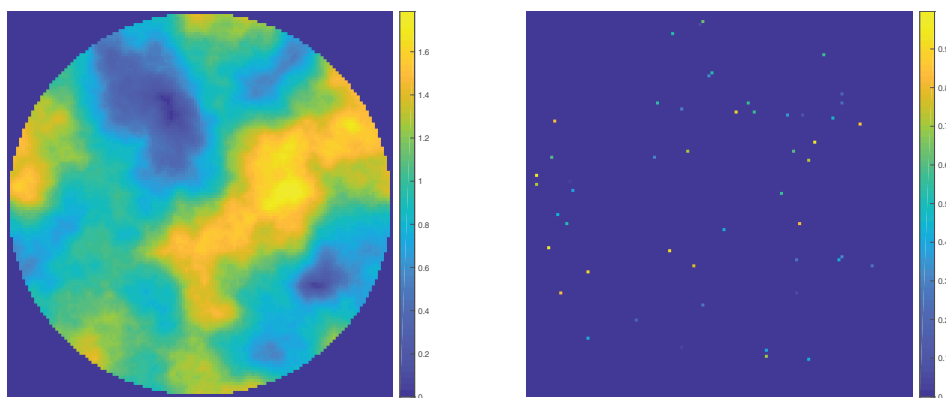


FIG. 5.1. The original material maps for plexiglass (left) and PVC (right).

In addition to the test images, we also need other parameters in (1.1). To generate the ray trace matrix  $\mathbf{A}$ , we use the MATLAB function `fanbeamtomolinear` from AIR Tools [13, 10, 9] to simulate a fan-beam geometry with a flat detector. Other parameters that we need to choose in this function are presented in Table 5.1.

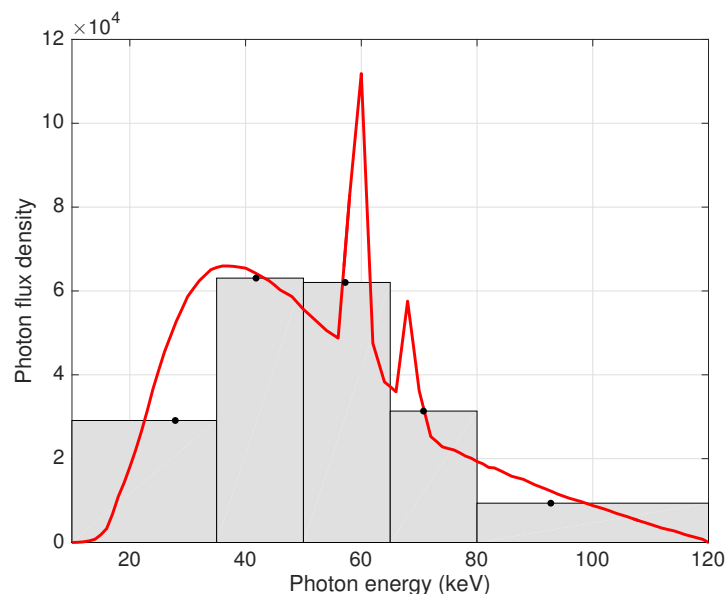
TABLE 5.1  
Geometry parameters of CT machine.

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

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` [22] 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 5.2. In Figure 5.2, the curve represents the photon intensity of the x-ray source, and the gray boxes indicate energy windows of the detector. Moreover, the black dots are the values of the 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 is well known, this strategy of generating data on a finer grid and solving it on a coarser grid is a standard approach to avoiding what is called the inverse crime.

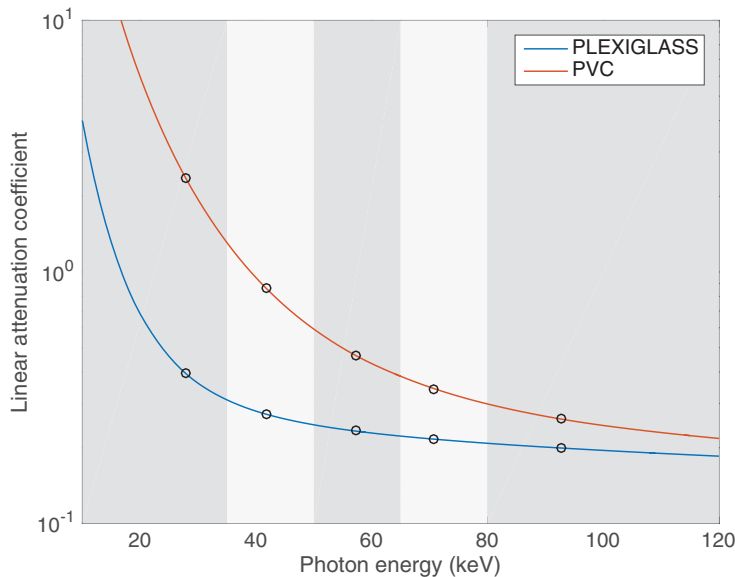
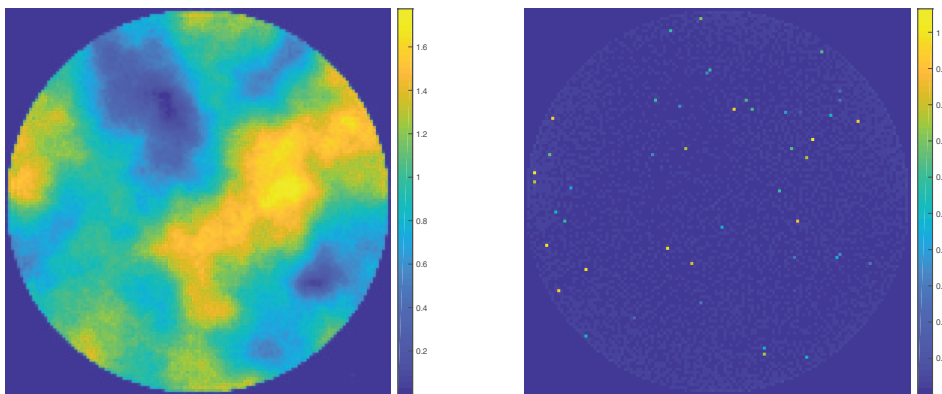
We also plot the curves of linear attenuation coefficients with respect to photon energy in Figure 5.3. From the figure, we can see that the slopes of these two curves are close to each other, which is likely to introduce collinearity between coefficients. Moreover, we assume that the entries of the 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 the Gaussian model is valid.

FIG. 5.2. *Detector bins and photon flux density.*

The reconstructed images are shown in Figure 5.4, which shows that we achieve almost perfect separation for these two materials. Moreover, the reconstructed images have excellent quality in terms of visually. Both material maps are relatively close to the true images. In the first material map, the distribution of weights is easy 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 yellows and greens. 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” is extremely difficult 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 methods proposed in this paper.

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 5.5. From this figure, we can see that the relative error of the first material drops sharply as the number of iterations increases. 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 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 the reconstructed image, causing somewhat large relative errors, even though visually the result looks quite good.

Other accuracy measures illustrate this phenomenon. In Figure 5.6, we plot the mean squared error (MSE) at each iteration. In Figure 5.7, the structural similarity index (SSIM) is presented.

FIG. 5.3. *Linear attenuation coefficients and photon flux density.*FIG. 5.4. *The reconstructed images for plexiglass (left) and PVC (right).*

Not surprisingly MSE produces information very similar to the relative errors, but it also shows a clear diminution for the second material in Figure 5.6. The SSIM is a metric for image quality, and large values correspond to better solutions. From Figure 5.7, it can be seen 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 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 5.8. From this figure, we can see that the norm of the gradient 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



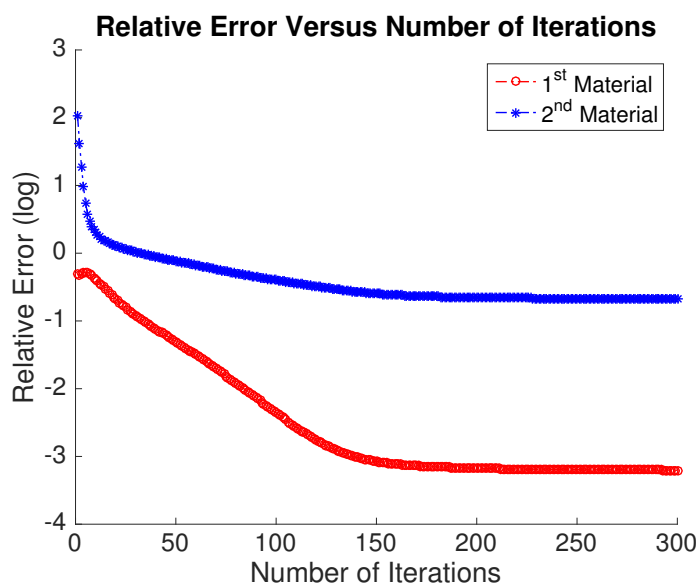


FIG. 5.5. The related errors for each iteration (with preconditioner) for plexiglass and PVC.

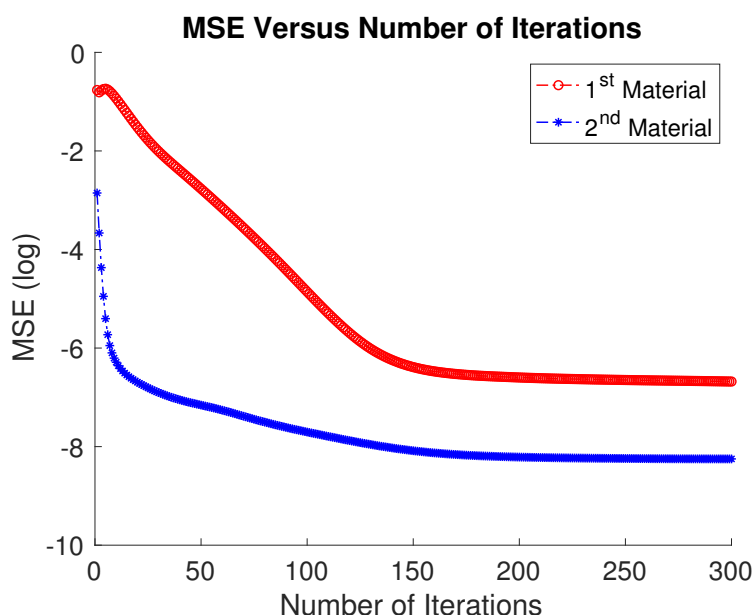


FIG. 5.6. MSE for each iteration (with preconditioner) for plexiglass and PVC.

performance with a preconditioner proposed by Barber et al. [1], and the performance without using any preconditioners. As previously mentioned, the approach proposed in [1] is based on the eigenvalue decomposition of  $\mathbf{C}^T \mathbf{C}$ . The results are shown in Figure 5.9, 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.

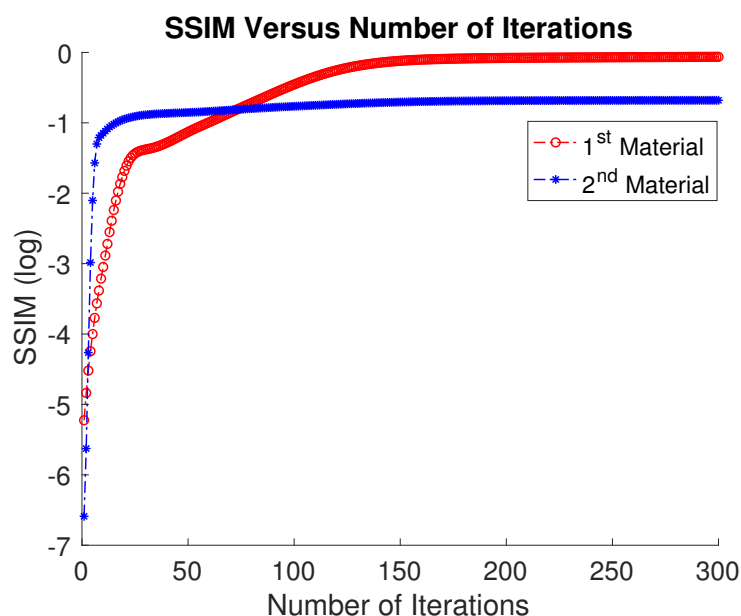


FIG. 5.7. *SSIM for each iteration (with preconditioner) for plexiglass and PVC.*

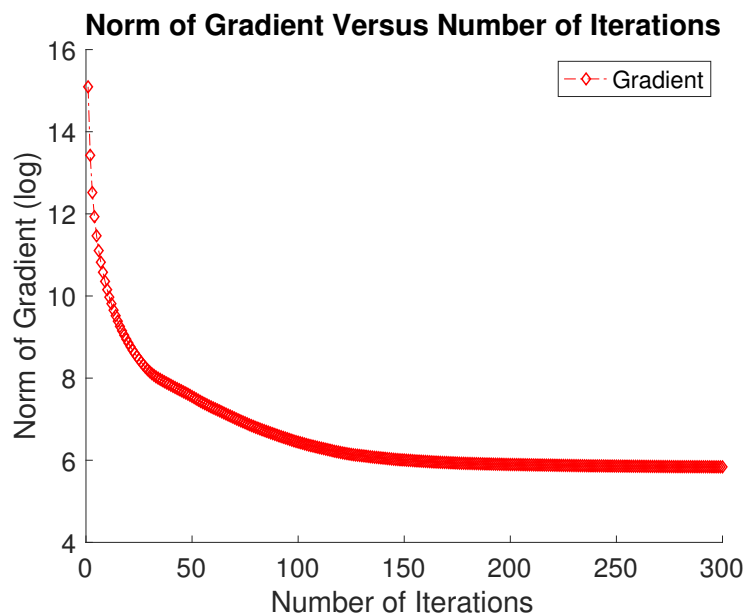


FIG. 5.8. *The norm of the gradient for overall materials, normalized by the 2-norm of the image.*

From 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.

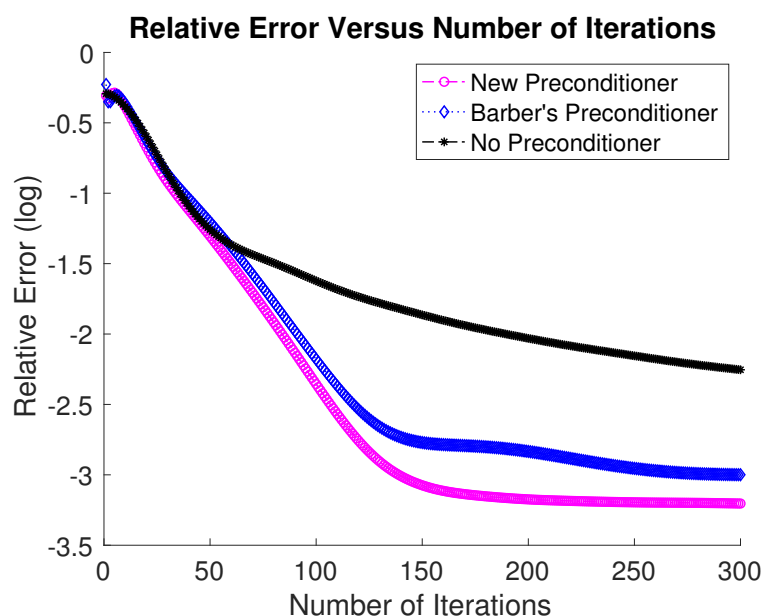


FIG. 5.9. The decay of related errors with new preconditioner, with Barber's [1] preconditioner, and with no preconditioner.

**6. Conclusions and remarks.** In this paper, 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 projected 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 using a first order method, FISTA, that can generate fast convergence speed. Because of the introduction of the new preconditioner, we recommend constructing a projection problem and computing 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 paper. For future work, it would be interesting to consider other regularization schemes to emphasize the edges of the object, such as total variation.

#### REFERENCES

- [1] R. F. BARBER, E. Y. SIDKY, T. G. SCHMIDT, AND X. PAN, *An algorithm for constrained one-step inversion of spectral CT data*, Phys. Med. Biol., 61 (2016), pp. 3784–3818.
- [2] A. BECK AND M. TEOULLE, *A fast iterative shrinkage-thresholding algorithm for linear inverse problems*, SIAM J. Imaging Sci., 2 (2009), pp. 183–202, <https://doi.org/10.1137/080716542>.
- [3] Å. BJÖRCK, *Numerical Methods for Least Squares Problems*, SIAM, Philadelphia, 1996, <https://doi.org/10.1137/1.9781611971484>.
- [4] V. M. BUSTAMANTE, J. G. NAGY, S. S. J. FENG, AND I. SECHOPOULOS, *Iterative breast tomosynthesis image reconstruction*, SIAM J. Sci. Comput., 35 (2013), pp. S192–S208, <https://doi.org/10.1137/120881440>.
- [5] A. CHAMBOLE AND T. POCK, *A first-order primal-dual algorithm for convex problems with applications to imaging*, J. Math. Imaging Vis., 40 (2011), pp. 120–145.

- [6] G. H. GOLUB AND C. F. VAN LOAN, *Matrix Computations*, Vol. 3, Johns Hopkins University Press, Baltimore, MD, 2012.
- [7] M. GRANT, S. BOYD, AND Y. YE, *CVX: MATLAB Software for Disciplined Convex Programming*, 2008, <http://cvxr.com/cvx/>.
- [8] 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.
- [9] P. C. HANSEN AND J. S. JØRGENSEN, *AIR Tools II: Algebraic iterative reconstruction methods, improved implementation*, Numer. Algorithms, 79 (2018), pp. 107–137.
- [10] P. C. HANSEN AND M. SAXILD-HANSEN, *AIR Tools: A MATLAB package of algebraic iterative reconstruction methods*, J. Comput. Appl. Math., 236 (2012), pp. 2167–2178.
- [11] B. J. HEISMAN, B. T. SCHMIDT, AND T. FLOHR, *Spectral Computed Tomography*, SPIE, Bellingham, WA, 2012.
- [12] J. D. INGLE, JR., AND S. R. CROUCH, *Spectrochemical Analysis*, Prentice-Hall, 1988.
- [13] A. C. KAK AND M. SLANEY, *Principles of Computerized Tomographic Imaging*, SIAM, Philadelphia, 2001, <https://doi.org/10.1137/1.9780898719277>.
- [14] R. M. LARSEN, *Lanczos bidiagonalization with partial reorthogonalization*, DAIMI Report Series, 27 (1998).
- [15] J. MATTINGLEY AND S. BOYD, *Automatic code generation for real-time convex optimization*, in Convex Optimization in Signal Processing and Communications, Cambridge University Press, 2009, pp. 1–41.
- [16] J. MATTINGLEY AND S. BOYD, *Real-time convex optimization in signal processing*, IEEE Signal Process. Mag., 27 (2010), pp. 50–61.
- [17] J. MATTINGLEY AND S. BOYD, *CVXGEN: A code generator for embedded convex optimization*, Optim. Engrg., 13 (2012), pp. 1–27.
- [18] J. MATTINGLEY, Y. WANG, AND S. BOYD, *Code generation for receding horizon control*, in Proceedings of the 2010 IEEE International Symposium on Computer-Aided Control System Design (CACSD), IEEE, 2010, pp. 985–992.
- [19] J. L. MUELLER AND S. SILTANEN, *Linear and Nonlinear Inverse Problems with Practical Applications*, SIAM, Philadelphia, 2012, <https://doi.org/10.1137/1.9781611972344>.
- [20] A. S. NEMIROVSKY AND D. B. YUDIN, *Problem Complexity and Method Efficiency in Optimization*, Wiley, Chichester, UK, 1983.
- [21] Y. E. NESTEROV, *A method for solving the convex programming problem with convergence rate  $O(1/k^2)$* , Dokl. Akad. Nauk SSSR, 269 (1983), pp. 543–547.
- [22] 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*, Med. Phys., 31 (2004), pp. 3057–3067.
- [23] C. F. VAN LOAN, *The ubiquitous Kronecker product*, J. Comput. Appl. Math., 123 (2000), pp. 85–100.
- [24] 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.