

Joint Reconstruction and Low-Rank Decomposition for Dynamic Inverse Problems

Simon Arridge* Pascal Fernsel† Andreas Hauptmann‡

May 29, 2020

Abstract

A primary interest in dynamic inverse problems is to identify the underlying temporal behaviour of the system from outside measurements. In this work we consider the case, where the target can be represented by a decomposition of spatial and temporal basis functions and hence can be efficiently represented by a low-rank decomposition. We then propose a joint reconstruction and low-rank decomposition method based on the Nonnegative Matrix Factorisation to obtain the unknown from highly undersampled dynamic measurement data. The proposed framework allows for flexible incorporation of separate regularisers for spatial and temporal features. For the special case of a stationary operator, we can effectively use the decomposition to reduce the computational complexity and obtain a substantial speed-up. The proposed methods are evaluated for two simulated phantoms and we compare obtained results to a separate low-rank reconstruction and subsequent decomposition approach based on the widely used principal component analysis.

Keywords: Nonnegative matrix factorisation, dynamic inverse problems, low-rank decomposition, variational methods

AMS Subject Classification: 15A69, 15A23, 65K10

1 Introduction

Several inverse problems are concerned with reconstruction of solutions in multiple physical dimensions such as space, time and frequency. Generally such problems require very large datasets in order to satisfy conditions for accurate reconstruction, whereas in practice only subsets of such complete data can be

*Department of Computer Science, University College London, London, United Kingdom (S.Arridge@cs.ucl.ac.uk).

†Center for Industrial Mathematics, University of Bremen, Bremen, Germany (pfernsel@math.uni-bremen.de).

‡Department of Mathematical Sciences, University of Oulu, Oulu, Finland; Department of Computer Science, University College London, London, United Kingdom (andreas.hauptmann@oulu.fi).

measured. Furthermore, the information content of the solutions from such reduced data may be much less than suggested by the complete set. In these cases, regularisation in the reconstruction process is required to compensate for the reduced information content, for instance by correlating features between auxiliary physical dimensions.

For instance, dynamic inverse problems have gained considerable interest in recent years. This development is partly driven by the increase in computational resources and the possibility to handle large data size more efficiently, but also novel and more efficient imaging devices enabling wide areas of applications in medicine and industrial imaging. For instance in medical imaging, dynamic information is essential for accurate diagnosis of heart diseases or for applications in angiography to determine blood flow by injecting a contrast agent to the patient's blood stream. But also in nondestructive testing and chemical engineering, tomographic imaging has become increasingly popular to monitor dynamic processes. The underlying problem in these imaging scenarios is often, that a fine temporal sampling, i.e. in the discrete setting a large number of channels, is only possible under considerable restrictions to sampling density at each time instance. This limitation often renders time-discrete (static) reconstructions insufficient. Additionally, an underlying problem in many dynamic applications is given by the specific temporal dynamics of the process, which are often non-periodic and hence prevents temporal binning approaches. Thus, it is essential to include the dynamic nature of the imaging task in the reconstruction process.

With increasing computational resources, it has become more feasible to address the reconstruction task as a fully dynamic problem in a spatio-temporal setting. In these approaches it is essential to include the dynamic information in some form into the reconstruction task. This could be done for instance by including a regularisation on the temporal behaviour as penalty in a variational setting [36, 37]. Such approaches have been used in a wide variety of applications, such as magnetic resonance imaging [15, 31, 38], X-ray tomography [4, 32] and application to process monitoring with electrical resistance tomography [8].

More advanced approaches aim to include a physical motion model and estimate the motion of the target from the measurements itself. This can be done for instance by incorporating an image registration step into the reconstruction algorithm and reformulate the reconstruction problem as a joint motion-estimation and reconstruction task [5, 6, 14, 30]. Another possibility is the incorporation of an explicit motion model by methamorphosis as considered in [9, 20].

In this work we consider another possibility to incorporate regularisation, and in particular temporal regularity, to the reconstruction task by assuming a low-dimensional representation of the unknown. This leads naturally to a low-rank description of the underlying inverse problem and is especially suitable to reduce data size in cases where we have much fewer basis functions to represent the unknown than the temporal sampling. In a continuous setting, this yields the analysis of low-rank approximations in tensor product of Hilbert spaces, for which we refer the reader to [22, 42]. We rather focus on low-rank

approximation methods in a discretised framework, which leads to the use of specific matrix factorisation approaches and their optimisation techniques.

In particular, in this work we propose a joint reconstruction and decomposition in a variational framework using non-negative matrix factorisation, which naturally represents the physical assumption of nonnegativity of the dynamic target and allows for a variety of regularising terms on spatial and temporal basis functions. Following this framework we propose two algorithms, that either jointly recover the reconstruction and the low-rank decomposition, or alternatively recovers only the low-rank representation of the unknown without the need to construct the full spatio-temporal target in the reconstruction process. Here, the second approach effectively incorporates the dimension reduction and can lead under certain assumptions on the forward operator to a significant reduction in computational complexity. This can be particularly useful, if one is only interested in the dynamics of the system and not the full reconstruction.

This paper is organised as follows. In Section 3 we discuss our setting for dynamic inverse problems and continue to discuss low-rank decomposition approaches. Specifically, principal component analysis (PCA) and non-negative matrix factorisation (NMF), which is the focus in this study. As a baseline we first present a low-rank reconstruction method followed by either of the decomposition methods. We then continue to present the proposed framework of joint reconstruction and decomposition with the NMF. In particular, we prove that the proposed framework leads to a monotonic decrease of the cost functions. We then proceed in Section 3 to evaluate the algorithms under considerations with the use case of dynamic X-ray tomography and two simulated phantoms with different characteristics. We conclude the study in Section 4 with some thoughts on the extension of the proposed framework.

2 Reconstruction and low-rank decomposition methods

2.1 A setting for dynamic Inverse Problems

In this work, we consider a general multi-dimensional inverse problem, where the unknown $x(s, \tau)$ is defined on a spatial domain $\Omega_1 \subset \mathbb{R}^{d_1}$ with dependence on a secondary variable $t \in \mathbb{R}_{\geq 0}$ defined in a bounded interval $\mathcal{T} := [0, T]$. This setting admits some quite general applications where the secondary variable could have other physical interpretations, notably wavelength for hyper-spectral problems; however, to fix our ideas, we henceforth consider t to explicitly represent time, and our application to be that of *dynamic inverse problems*. Consequently, the underlying equation of the resulting inverse problem can be described in the following form

$$\mathcal{A}(x(s, t); t) = y(\sigma, t) \quad \text{for } t \in \mathcal{T}, \quad (1)$$

where $\mathcal{A} : \Omega_1 \times \mathcal{T} \rightarrow \Omega_2 \times \mathcal{T}$ is a linear bounded operator between suitable Hilbert spaces and $\Omega_2 \subset \mathbb{R}^{d_2}$. We will primarily consider the non-stationary

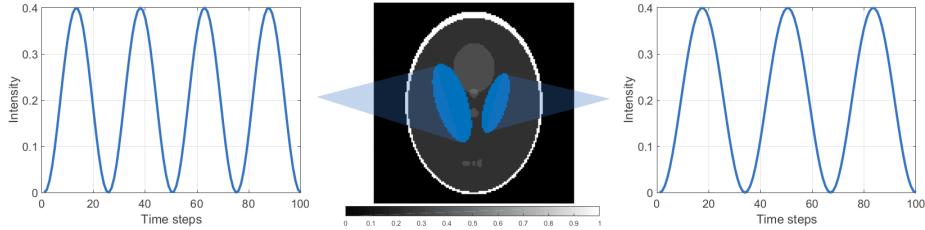


Figure 1: Illustration of a phantom that can be represented by the decomposition in (2). The phantom consists of $K = 3$ components, for the background and two dynamic components with periodically changing intensity (left and right plot). As such, this phantom can be efficiently represented by a low-rank decomposition considered in this study.

case here, where the forward operator \mathcal{A} is dependent on t . In the special case of a stationary operator $\mathcal{A}(\cdot; t) \equiv \mathcal{A}$ for all $t \in \mathcal{T}$, where for each t the operator follows the same sampling process, we can achieve possible computational improvements. The resulting implications will be discussed later in Section 2.5.

Furthermore, the underlying assumption in this work is that the unknown $x \in \Omega_1 \times \mathcal{T}$ can be decomposed into a set of spatial $b^k : \Omega_1 \rightarrow \mathbb{R}_{\geq 0}$ and channel basis functions $c^k(t) : \mathcal{T} \rightarrow \mathbb{R}_{\geq 0}$ for $1 \leq k \leq K$. Then the unknown can be represented by the decomposition

$$x(s, t) = \sum_{k=1}^K b^k(s)c^k(t). \quad (2)$$

This formulation naturally gives rise to the reconstruction and low-rank decomposition framework to extract the relevant features given by b^k and c^k . An illustration for a possible phantom represented by (2) is shown in Figure 1.

We intentionally keep the formulation general here to allow for applications different to dynamic inverse problems, such as multi-spectral imaging. Nevertheless, the derived reconstruction and feature extraction framework in this paper will be used in Section 3 for the specific application to dynamic computed tomography.

Furthermore, a suitable discretisation of the continuous formulation (1) is needed to introduce the feature extraction methods in the forthcoming sections. Let us first discretise the secondary variable, such that $t \in \mathbb{N}$ with $1 \leq t \leq T$. For the spatial domain, we assume a vectorised representation such that the resulting unknown can be represented as matrix $X \in \mathbb{R}^{N \times T}$, which leads to the matrix equation

$$A_t X_{\bullet, t} = Y_{\bullet, t} \quad \text{for } 1 \leq t \leq T, \quad (3)$$

where $A_t \in \mathbb{R}^{M \times N}$ is the discretised forward operator, $X_{\bullet, t}$ the t -th column of X and $Y_{\bullet, t}$ the t -th column of the data matrix $Y \in \mathbb{R}^{M \times T}$. Analogously, we will write $M_{n, \bullet}$ for the n -th row of an arbitrary matrix M .

Suitable restrictions to the matrices in Equation (3) will be made in the following sections to ensure the applicability of the considered frameworks and, if possible, to properly represent the decomposition (2).

2.2 Feature Extraction Methods

In this section, we introduce two feature extraction methods, namely the *Principle Component Analysis* (PCA) and the *Nonnegative Matrix Factorisation* (NMF). These approaches are used to compute the latent components of the reconstruction X . The NMF will be used in Section 2.4 to introduce a joint reconstruction and low-rank decomposition framework to tackle the problem stated in (3).

2.2.1 Principle Component Analysis

Large and high dimensional datasets demand modern data analysis approaches to reduce the dimensionality and increase the interpretability of the data while keeping the loss of information as low as possible. Many different techniques have been developed for this purpose, but PCA is one of the most widely used and goes back to [35].

For a given matrix $X \in \mathbb{R}^{N \times T}$ with N different observations of an experiment and T features, the PCA is a linear orthogonal transformation given by the weights $C_{\tilde{k}, \bullet} = (C_{\tilde{k}1}, \dots, C_{\tilde{k}T})$ with $C \in \mathbb{R}^{\tilde{K} \times T}$, which transforms each observation $X_{n, \bullet}$ to *principal component scores* given by $B_{n\tilde{k}} := \sum_t X_{nt} C_{\tilde{k}t}$ with $B = [B_{\bullet, 1}, \dots, B_{\bullet, \tilde{K}}] \in \mathbb{R}^{N \times \tilde{K}}$ and $\tilde{K} = \min(N - 1, T)$, such that

- the sample variance $\text{Var}(B_{\bullet, \tilde{k}})$ is maximised for all \tilde{k} ,
- each row $C_{\tilde{k}, \bullet}$ is constrained to be a unit vector
- and the sample covariance $\text{cov}(B_{\bullet, k}, B_{\bullet, \tilde{k}}) = 0$ for $k \neq \tilde{k}$.

Together with the usual assumption that the number of observations is higher than the underlying dimension, this leads to $\tilde{K} = T$ and the full transformation $B = XC^\top$, where C is an orthogonal matrix. The t -th column vector $(C_{t, \bullet})^\top$ defines the t -th *principal direction* and is an eigenvector of the covariance matrix $S = X^\top X / (N - 1)$. The corresponding t -th largest eigenvalue of S denotes the variance of the t -th principal component.

The above transformation is equivalent to the factorisation of the matrix X given by

$$X = BC, \tag{4}$$

which allows to decompose each observation into the principal components, such that $X_{n, \bullet} = \sum_{t=1}^T B_{nt} C_{t, \bullet}$. Therefore, we also have $X = \sum_{t=1}^T B_{\bullet, t} C_{t, \bullet}$ similarly to the continuous case in (2).

Furthermore, it is possible to obtain an approximation of the matrix X by truncating the sum at the first $K < T$ principle components for all n , which

yields a rank K matrix $X^{(K)}$ given by

$$X^{(K)} = \sum_{k=1}^K B_{\bullet,k} C_{k,\bullet}.$$

Based on the Eckart–Young–Mirsky theorem [19], $X^{(K)}$ is the best rank K approximation of X in the sense that it minimises the discrepancy $\|X - X^{(K)}\|$ for both the Frobenius and spectral norm.

One typical approach to compute the PCA is based on the Singular Value Decomposition (SVD) of the data matrix $X = U\Sigma V^\top$ and will be used in this work. Setting $B := U\Sigma$ and $C = V^\top$ gives already the desired factorisation in (4) based on the PCA.

2.2.2 Nonnegative Matrix Factorisation

Nonnegative Matrix Factorisation (NMF), originally introduced as positive matrix factorisation by Paatero and Tapper in 1994 [34], is an established tool to obtain low-rank approximations of nonnegative data matrices. It has been widely used in the machine learning and data mining community for compression, basis learning, clustering and feature extraction for high-dimensional classification problems with applications in music analysis [17], document clustering [13] and medical imaging problems such as tumor typing in matrix-assisted laser desorption/ionisation (MALDI) imaging in the field of bioinformatics [28].

Different from the PCA approach above, the NMF enforces nonnegativity constraints on the factor matrices without any orthogonality restrictions. This makes the NMF the method of choice for application fields, where the underlying physical model enforces the solution to be nonnegative assuming that each datapoint can be described as a superposition of some unknown characteristic features of the dataset. The NMF makes it possible to extract these features while constraining the matrix factors to have nonnegative entries, which simplifies their interpretation.

These data assumptions are true for many application fields including the ones mentioned above but also especially our considered problem of dynamic computed tomography, where the measurements consist naturally of the nonnegative absorption of photons.

Mathematically, the basic NMF problem can be formulated as follows: For a given nonnegative matrix $X \in \mathbb{R}_{\geq 0}^{N \times T}$, find nonnegative matrices $B \in \mathbb{R}_{\geq 0}^{N \times K}$ and $C \in \mathbb{R}_{\geq 0}^{K \times T}$ with $K \ll \min\{N, T\}$, such that

$$X \approx BC.$$

The factorisation allows to approximate the rows $X_{n,\bullet}$ and columns $X_{\bullet,t}$ of the data matrix as a superposition of the K rows $B_{k,\bullet}$ of B and columns $C_{\bullet,k}$ of C respectively, such that $X_{n,\bullet} \approx \sum_{k=1}^K B_{nk} C_{k,\bullet}$ and $X_{\bullet,t} \approx \sum_{k=1}^K C_{kt} B_{\bullet,k}$.

Similarly, it holds that

$$X \approx BC = \sum_{k=1}^K B_{\bullet,k} C_{k,\bullet},$$

where the K terms of the sum are rank-one matrices. Hence, the sets $\{B_{\bullet,k}\}_k$ and $\{C_{k,\bullet}\}_k$ can be interpreted as a low-dimensional basis to approximate the data matrix, i.e. the NMF performs the task of basis learning with additional nonnegativity constraints.

The usual approach to compute the factorisation is to define a suitable discrepancy term \mathcal{D}_{NMF} , which has to be chosen according to the noise assumption of the underlying problem, and to reformulate the NMF as a minimisation problem. Typical discrepancies include the default case of the Frobenius Norm on which we will focus on, the Kullback–Leibler divergence, the Itakura–Saito distance or other generalized divergences [10].

Furthermore, NMF problems are usually ill-posed due to the non-uniqueness of the solution [21] and require the application of suitable regularisation techniques. One common method is to include penalty terms in the minimisation problem to tackle the ill-posedness of the problem but also to enforce desirable properties of the factorisation matrices. Typical examples range from ℓ_1, ℓ_2 and total variation regularisation terms [25] to more problem specific terms, which enforce additional orthogonality of the matrices or even allow supervised classification workflows if the NMF is used as a prior feature extraction method [16, 28].

Hence, the general regularised NMF problem can be written as

$$\min_{B,C \geq 0} \mathcal{D}_{\text{NMF}}(X, BC) + \sum_{\ell=1}^L \gamma_\ell \mathcal{P}_\ell(B, C) =: \min_{B,C \geq 0} \mathcal{F}(B, C), \quad (5)$$

where \mathcal{P}_ℓ denote the penalty terms, $\gamma_\ell \geq 0$ the corresponding regularisation parameters and \mathcal{F} the cost function of the NMF.

The considered optimisation approach in this work is based on the so-called Majorise-Minimisation principle and gives rise to multiplicative update rules of the matrices in (5), which automatically preserve the nonnegativity of the iterates provided that they are initialised nonnegative. For more details on this optimisation technique, we refer the reader to Appendix A.

The idea of the feature extraction procedure based on the NMF can be well illustrated by considering the example from Figure 1 that satisfies the decomposition assumption from (2). Here, the highlighted spatial regions change their intensities according to the given dynamics. The NMF allows a natural interpretation of the factorisation matrices B and C as the spatial and temporal basis functions for this case, as illustrated in Figure 2. The column $X_{\bullet,t}$ of X denotes the reconstruction of the t -th time step of the inverse problem in (3). The NMF allows to decompose the spatial and temporal features of the data: The matrix B contains the spatial features in its columns with the corresponding temporal features in the rows of C .

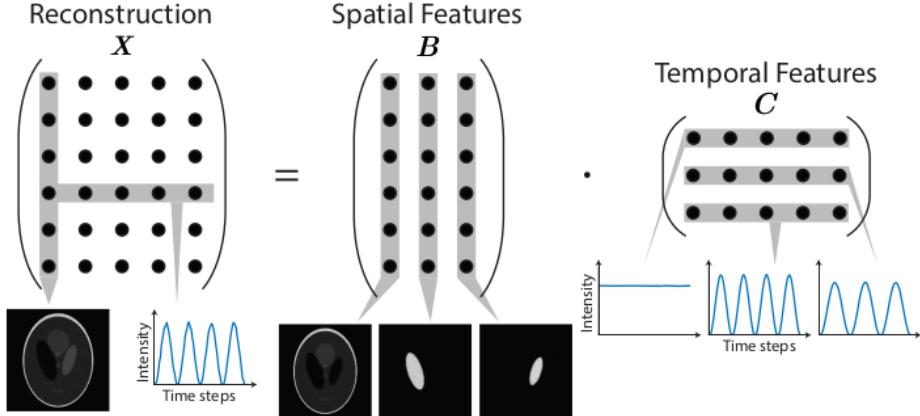


Figure 2: Structure of the NMF in the context of the dynamic Shepp-Logan phantom as shown in Figure 1. Here, the nonnegative spatial and temporal basis functions can be naturally represented by the matrices B and C .

2.3 Separated Reconstruction and Low-rank Decomposition

Let us first discuss a separated reconstruction and feature extraction approach to solve the inverse problem in (3), that means we compute first a reconstruction and perform then subsequently the feature extraction with one of the previously discussed methods. We consider this method as baseline for our comparison.

The considered reconstruction method for this separated framework involves a basic gradient descent approach together with a regularisation step and a subsequent total variation denoising, which will be henceforth referred to as **gradTV**. The details on the algorithm are provided in Algorithm 1. In particular, we aim to compute solutions to the least squares problem and incorporate the low-rank assumptions as additional penalty of the nuclear norm of $X_{\bullet,t}$, that is

$$\min_{X_{\bullet,t} \geq 0} \|Y_{\bullet,t} - A_t X_{\bullet,t}\|_2^2 + \alpha \|X_{\bullet,t}\|_*$$

for all t ; see e.g. [29, 41]. This can then be efficiently solved by a proximal gradient descent scheme with a soft-thresholding on the singular values and hence enforcing the low-rank structure. Ideally, one would like to include the total variation regularisation as penalty term, but as this tends to be computationally expensive for the fine temporal sampling, we included this as a subsequent denoiser.

In practice, after a suitable initialisation of the reconstruction matrix, the gradient descent step is computed with an, *a priori* defined, fixed stepsize ρ_{grad} . For the proximal step, the truncated SVD of X is computed and a soft thresholding of the singular values is performed with a fixed threshold ρ_{thr} . Afterwards, we enforce the nonnegativity with a projection step on the reconstruction X . When

the stopping criterion is satisfied, a TV denoising algorithm¹ based on [18, 40] with the corresponding parameter ρ_{TV} is applied.

Algorithm 1 gradTV

```

1: Initialise:  $X$ 
2: Input:  $\rho_{\text{grad}}, \rho_{\text{thr}}, \rho_{\text{TV}} > 0$ 
3: repeat
4:    $X_{\bullet,t} \leftarrow X_{\bullet,t} - \rho_{\text{grad}}(A_t^\top A_t X_{\bullet,t} - A_t^\top Y_{\bullet,t})$  for all  $t$ 
5:    $(U, \Sigma, V) \leftarrow \text{SVD}(X)$ 
6:    $\Sigma \leftarrow \text{SOFTTHRESH}_{\rho_{\text{thr}}}(\Sigma)$ 
7:    $X \leftarrow U\Sigma V^\top$ 
8:    $X \leftarrow \max(X, 0)$ 
9: until STOPPINGCRITERION satisfied
10:  $X \leftarrow \text{TVDENOSER}_{\rho_{\text{TV}}}(X)$ 
11: return  $X$ 
```

After the reconstruction procedure given by Algorithm 1, we perform the feature extraction of the reconstruction X via both the PCA and the NMF and call the approach **gradTV_PCA** and **gradTV_NMF** respectively.

For **gradTV_PCA**, we simply compute the PCA of X based on its SVD. Concerning the method **gradTV_NMF**, we consider the standard NMF model

$$\min_{B, C \geq 0} \|X - BC\|_F^2 + \frac{\tilde{\mu}_C}{2} \|C\|_F^2 \quad (6)$$

with the parameter $\tilde{\mu}_C$. The ℓ_2 regularisation penalty term on C is motivated by our application in Section 3. The corresponding multiplicative algorithms to solve (6) are well-known [11, 16] and a special case of the derived update rules in the next Section.

2.4 Joint Reconstruction and Low-rank Decomposition

Instead of the previously discussed separated reconstruction we now aim to include the feature extraction into the reconstruction procedure. This gives rise to consider a joint reconstruction and low-rank decomposition approach based on the NMF, rather than one based on a low-rank plus sparsity approach based on PCA [7, 39, 44]. The basic idea of the method is to incorporate the reconstruction procedure of the inverse problem in (3) into the NMF workflow. To do this, we have to additionally assume that $A_t \in \mathbb{R}_{\geq 0}^{M \times N}$, $Y \in \mathbb{R}_{\geq 0}^{M \times T}$ and $X \in \mathbb{R}_{\geq 0}^{N \times T}$ to ensure the desired nonnegativity of the factorisation matrices B and C , which corresponds to the assumptions of the decomposition in (2). The main motivation is that this joint approach allows the reconstruction process to exploit the underlying latent NMF features of the dataset, which can therefore

¹<https://www.mathworks.com/matlabcentral/fileexchange/36278-split-bregman-method-for-total-variation-denoising>

enhance the quality of the reconstructions by enabling regularisation of temporal and spatial features separately.

This can be achieved by including a discrepancy term $\mathcal{D}_{\text{IP}}(Y_{\bullet,t}, A_t X_{\bullet,t})$ of the inverse problem into the NMF cost function in (5). This leads together with some possible penalty terms for the reconstruction X to the model

$$\min_{B,C,X \geq 0} \mathcal{D}_{\text{IP}}(Y_{\bullet,t}, A_t X_{\bullet,t}) + \alpha \mathcal{D}_{\text{NMF}}(X, BC) + \sum_{\ell=1}^L \gamma_\ell \mathcal{P}_\ell(B, C, X), \quad (7)$$

with $\alpha \geq 0$ for the joint reconstruction and low-rank decomposition problem, which we will call **BC-X**. Furthermore, we can enforce $X := BC$ as a hard constraint, such that the reconstruction matrix will have at most rank K . In this case, the discrepancy \mathcal{D}_{NMF} vanishes and we end up with the model **BC**:

$$\min_{B,C \geq 0} \mathcal{D}_{\text{IP}}(Y_{\bullet,t}, A_t(BC)_{\bullet,t}) + \sum_{\ell=1}^L \gamma_\ell \mathcal{P}_\ell(B, C). \quad (8)$$

2.4.1 Considered NMF Models

For both models (7) and (8), we use the standard Frobenius norm for both the discrepancy terms \mathcal{D}_{NMF} and \mathcal{D}_{IP} . Furthermore, the optimisation method discussed in Section 2.4.2 allows to include a variety of penalty terms into the cost function. This makes it possible to construct suitable regularised NMF models and to enforce additional properties to the matrices depending on the specific application. In this work, we will consider standard ℓ_1 and ℓ_2 regularisation terms on each matrix and an isotropic total variation penalty on the matrix B . The latter is motivated by our considered application, which denoises the spatial features and thus also the reconstruction matrix. Hence, we will focus on the following NMF models in the remainder of this work:

$$\begin{aligned} & \min_{B,C,X \geq 0} \left\{ \sum_{t=1}^T \frac{1}{2} \|A_t X_{\bullet,t} - Y_{\bullet,t}\|_2^2 + \frac{\alpha}{2} \|BC - X\|_F^2 + \lambda_B \|B\|_1 + \frac{\mu_B}{2} \|B\|_F^2 \right. \\ & \quad \left. + \lambda_C \|C\|_1 + \frac{\mu_C}{2} \|C\|_F^2 + \lambda_X \|X\|_1 + \frac{\mu_X}{2} \|X\|_F^2 + \frac{\tau}{2} \text{TV}(B) \right\}, \quad \text{BC-X} \\ & \min_{B,C \geq 0} \left\{ \sum_{t=1}^T \frac{1}{2} \|A_t(BC)_{\bullet,t} - Y_{\bullet,t}\|_2^2 + \lambda_C \|C\|_1 + \frac{\mu_C}{2} \|C\|_F^2 + \lambda_B \|B\|_1 \right. \\ & \quad \left. + \frac{\mu_B}{2} \|B\|_F^2 + \frac{\tau}{2} \text{TV}(B) \right\}. \quad \text{BC} \end{aligned}$$

The regularisation parameters $\alpha, \lambda_C, \mu_C, \lambda_B, \mu_B, \lambda_X, \mu_X, \tau \geq 0$, are chosen *a priori*. Furthermore, $\|\cdot\|_F$ denotes the Frobenius norm, $\|M\|_1 := \sum_{ij} |M_{ij}|$ the 1-norm for matrices M and $\text{TV}(\cdot)$ is the following smoothed isotropic total variation [12, 16, 25].

Definition 2.1. The total variation of a matrix $B \in \mathbb{R}^{N \times K}$ is defined as

$$\text{TV}(B) := \sum_{k=1}^K \sum_{n=1}^N |\nabla_{nk} B| := \sum_{k=1}^K \sum_{n=1}^N \sqrt{\varepsilon_{\text{TV}}^2 + \sum_{\ell \in N_n} (B_{nk} - B_{\ell k})^2},$$

where $\varepsilon_{\text{TV}} > 0$ is a small positive constant and N_n are index sets referring to spatially neighboring pixels.

A typical example for the neighbourhood of the pixel $(0, 0)$ in two dimensions is $N_{(0,0)} = \{(1, 0), (0, 1)\}$ to get an estimate of the gradient components in both directions of the axes. The parameter ε_{TV} ensures the differentiability of the TV penalty term.

In the following section, we will present the multiplicative update rules for the NMF models BC-X and BC and derive the algorithms in Appendix B based on the majorise minimisation principle.

2.4.2 Algorithms

In this section, we present in Theorem 2.1 and 2.2 the multiplicative algorithms for the NMF problems in BC-X and BC. As mentioned in Section 2.2.2, the multiplicative structure of the iteration scheme ensures automatically the nonnegativity of the matrices B and C as long as they are initialised nonnegative. The derivation of such algorithms in this work are based on the so-called Majorise-Minimisation principle. The main idea of this approach is to replace the considered NMF cost function \mathcal{F} with a suitable auxiliary function $\mathcal{Q}_{\mathcal{F}}$, whose minimisation is much easier to handle and leads to a monotone decrease of \mathcal{F} . Furthermore, specific construction techniques of these *surrogate functions* lead to the desired multiplicative update rules which fulfill the nonnegativity constraint. We provide a short description of the main principles in Appendix A. A more detailed discussion of different construction methods for various kinds of discrepancy and penalty terms of \mathcal{F} can be found in the survey paper [16]. For better readability we present only the main results here and a detailed construction of the surrogate functions as well as derivation of the algorithms for both cost functions BC-X and BC can be found in Appendix B. Consequently, we will only state the main results in Theorem 2.1 and 2.2 here. Nevertheless, due to the construction of a suitable surrogate function for the TV penalty term (see Appendix B and [16] for more details), we first introduce the following matrices $P(B), Z(B) \in \mathbb{R}_{\geq 0}^{N \times K}$ as

$$P(B)_{nk} := \frac{1}{|\nabla_{nk} B|} \sum_{\ell \in N_n} 1 + \sum_{\ell \in \bar{N}_n} \frac{1}{|\nabla_{\ell k} B|}, \quad (9)$$

$$Z(B)_{nk} := \frac{1}{P(B)_{nk}} \left(\frac{1}{|\nabla_{nk} B|} \sum_{\ell \in N_n} \frac{B_{nk} + B_{\ell k}}{2} + \sum_{\ell \in \bar{N}_n} \frac{B_{nk} + B_{\ell k}}{2|\nabla_{\ell k} B|} \right), \quad (10)$$

where \bar{N}_n is the set of the so-called *adjoint* neighbourhood pixels, which is given by the relation

$$\ell \in \bar{N}_n \Leftrightarrow n \in N_\ell.$$

Furthermore, we write $\mathbf{1}_{M \times N}$ for an $M \times N$ matrix with ones in every entry. We then obtain the two algorithms for both models under consideration. First for the BC-X model that jointly obtains the reconstruction X and the decomposition:

Theorem 2.1 (Algorithm for BC-X). *For $A_t \in \mathbb{R}_{\geq 0}^{M \times N}, Y \in \mathbb{R}_{\geq 0}^{M \times T}$ and initialisations $X^{[0]} \in \mathbb{R}_{>0}^{N \times T}, B^{[0]} \in \mathbb{R}_{>0}^{N \times K}, C^{[0]} \in \mathbb{R}_{>0}^{K \times T}$, the alternating update rules*

$$\begin{aligned} X_{\bullet,t}^{[d+1]} &= X_{\bullet,t}^{[d]} \circ \frac{A_t^\top Y_{\bullet,t} + \alpha B^{[d]} C_{\bullet,t}^{[d]}}{A_t^\top A_t X_{\bullet,t}^{[d]} + (\mu_X + \alpha) X_{\bullet,t}^{[d]} + \lambda_X \mathbf{1}_{N \times 1}} \\ B^{[d+1]} &= B^{[d]} \circ \frac{\alpha X^{[d+1]} C^{[d]\top} + \tau P(B^{[d]}) \circ Z(B^{[d]})}{\alpha B^{[d]} C^{[d]\top} + \mu_B B^{[d]} + \lambda_B \mathbf{1}_{N \times K} + \tau B^{[d]} \circ P(B^{[d]})} \\ C^{[d+1]} &= C^{[d]} \circ \frac{\alpha B^{[d+1]\top} X^{[d+1]}}{\alpha B^{[d+1]\top} B^{[d+1]} C^{[d]} + \mu_C C^{[d]} + \lambda_C \mathbf{1}_{K \times T}} \end{aligned}$$

lead to a monotonic decrease of the cost function in BC-X.

Similarly, for the BC model we obtain the updates rules without constructing the matrix X during the reconstruction process:

Theorem 2.2 (Algorithm for BC). *For $A_t \in \mathbb{R}_{\geq 0}^{M \times N}, Y \in \mathbb{R}_{\geq 0}^{M \times T}$ and initialisations $B^{[0]} \in \mathbb{R}_{>0}^{N \times K}, C^{[0]} \in \mathbb{R}_{>0}^{K \times T}$, the alternating update rules*

$$\begin{aligned} B^{[d+1]} &= B^{[d]} \circ \frac{\sum_{t=1}^T A_t^\top Y_{\bullet,t} \cdot (C^{[d]\top})_{t,\bullet} + \tau P(B^{[d]}) \circ Z(B^{[d]})}{\sum_{t=1}^T A_t^\top A_t (B^{[d]} C^{[d]})_{\bullet,t} \cdot (C^{[d]\top})_{t,\bullet} + \mu_B B^{[d]} + \lambda_B \mathbf{1}_{N \times K} + \tau B^{[d]} \circ P(B^{[d]})} \\ C_{\bullet,t}^{[d+1]} &= C_{\bullet,t}^{[d]} \circ \frac{B^{[d+1]\top} A_t^\top Y_{\bullet,t}}{B^{[d+1]\top} A_t^\top A_t (B^{[d+1]} C^{[d]})_{\bullet,t} + \mu_C C_{\bullet,t}^{[d]} + \lambda_C \mathbf{1}_{K \times 1}} \end{aligned}$$

lead to a monotonic decrease of the cost function in BC.

We remind that the derivation is described in Appendix B with lead to the update rules in Theorems above. Due to the multiplicative structure of the algorithms, zero entries in the matrices stay zero during the iteration scheme and can cause divisions by zero. This issue is handled via the strict positive initialisation in both Theorems. Furthermore, very small or high numbers can cause numerical instabilities and lead to undesirable results. As a standard procedure, this problem is handled by suitable projection steps after every iteration step [10].

2.5 Complexity Reduction for Stationary Operator

Let us now consider the case for a stationary operator, i.e. $\mathcal{A}(\cdot; t)$ in equation (1) does not change with t . Then we simply write \mathcal{A} or A for the matrix representation in (3). If further the number of channels T is large, the application of the forward operator represented a major computational burden per channel. In particular, we make use here of the assumption $T \gg K$, i.e. the number of channels is much larger than the basis functions for the decomposition. In this case, we can effectively reduce the computational cost by shifting the application of the forward operator to the spatial basis functions contained in B . That means, we make essential use of the decomposition $X \approx BC$ in the reconstruction task and as such avoid to construct the approximation to X . Consequently, we will only consider the case of BC here. Since A is independent from t , the NMF model BC becomes

$$\begin{aligned} \min_{B,C \geq 0} & \left\{ \frac{1}{2} \|ABC - Y\|_F^2 + \lambda_C \|C\|_1 + \frac{\mu_C}{2} \|C\|_F^2 + \lambda_B \|B\|_1 \right. \\ & \left. + \frac{\mu_B}{2} \|B\|_F^2 + \frac{\tau}{2} \text{TV}(B) \right\}. \end{aligned} \quad \text{sBC}$$

To illustrate this, let us consider the update equation in Theorem 2.2 for B , where we can simplify the first term in the denominator as follows:

$$\sum_{t=1}^T A^\top A (B^{[d]} C^{[d]})_{\bullet,t} \cdot (C^{[d]\top})_{t,\bullet} = A^\top A \sum_{t=1}^T (B^{[d]} C^{[d]})_{\bullet,t} \cdot (C^{[d]\top})_{t,\bullet} = A^\top AB^{[d]} C^{[d]} C^{[d]\top}$$

The other terms in the update rules can be simplified similarly, such that we obtain the following reduced update equations:

Corollary 2.1 (Algorithm for sBC). *For $A \in \mathbb{R}_{\geq 0}^{M \times N}$, $Y \in \mathbb{R}_{\geq 0}^{M \times T}$ and initialisations $B^{[0]} \in \mathbb{R}_{>0}^{N \times K}$, $C^{[0]} \in \mathbb{R}_{>0}^{K \times T}$, the alternating update rules*

$$\begin{aligned} B^{[d+1]} &= B^{[d]} \circ \frac{A^\top Y C^{[d]\top} + \tau P(B^{[d]}) \circ Z(B^{[d]})}{A^\top AB^{[d]} C^{[d]\top} + \mu_B B^{[d]} + \lambda_B \mathbf{1}_{N \times K} + \tau B^{[d]} \circ P(B^{[d]})} \\ C^{[d+1]} &= C^{[d]} \circ \frac{B^{[d+1]\top} A^\top Y}{B^{[d+1]\top} A^\top AB^{[d+1]} C^{[d]} + \mu_C C^{[d]} + \lambda_C \mathbf{1}_{K \times T}}. \end{aligned}$$

lead to a monotonic decrease of the cost function in sBC.

Finally, the order of application is essential here to obtain the complexity reduction. In particular, in the following we implemented the algorithm such that A acts on the basis functions in B . That means, we compute first the product $A^\top AB$ followed by multiplication with C . That means, we can expect a reduction of computational complexity by a factor T/K with the sBC model and hence is especially useful for dimension reduction under fine temporal sampling.

3 Application to Dynamic CT

In the following we will apply the presented methods to the use case of dynamic computerised tomography (CT). Here the quantity of interest is given as the

attenuation coefficient $x(s, t)$ at time $t \in [0, T]$ on a bounded domain in two dimensions $s \in \Omega_1 \subset \mathbb{R}^2$. Following the formulation in (1), the time-dependent forward operator is given by the Radon transform

$$y(\theta, \sigma, t) := (\mathcal{R}_{\mathcal{I}(t)}x(s, t))(\theta, \sigma) = \int_{s \cdot \theta = \sigma} x(s, t) \, ds \quad (11)$$

Here, the measurement $y(\theta, \sigma, t)$ consist of line integrals over the domain Ω_1 for each time point $t \in \mathcal{T}$, and is referred to as the sinogram. This measurement depends on two parameters, the angle $\theta \in S^1$ on the unit circle and a signed distance to the origin $\sigma \in \mathbb{R}$. Consequently, the measurements depend on a set of angles at each time step $\mathcal{I}(t)$, such that $(\theta, \sigma) \in \mathcal{I}(t)$ at time t , we will refer to this as the sampling patterns. In a slight abuse of notation, we will use $|\mathcal{I}(t)|$ for the number of angles, i.e. directions for the line integrals, at each time point.

In the following we consider two scenarios for the choice of angles in $\mathcal{I}(t)$ and by that defining the nature of the forward operator, as discussed in Section 2.1. In the general case of a nonstationary forward operator, that means the sampling patterns are time-dependent, we assume that the angles change but the amount of angles is constant over time $|\mathcal{I}(t)| \equiv c$. Additionally, we will consider the case for stationary operators, which in our setting means that the set of angles does not change over time, we can write for instance $\mathcal{I}(t) \equiv \mathcal{I}(t=0)$, and hence this leads to a stationary measurement operator of the dynamic process in (11). We note that even though the measurement process is stationary, the obtained measurement $y(\theta, \sigma, t)$ itself is still time dependent.

For the computations, we discretise (11) to obtain a matrix vector representation as in (3). In the following we will write R_t for the discrete Radon transform with respect to the sampling pattern $\mathcal{I}(t)$ at time point t , which gives rise to the discrete reconstruction problem for dynamic CT

$$R_t X_{\bullet, t} = Y_{\bullet, t} \quad \text{for } 1 \leq t \leq T. \quad (12)$$

We note, that due to the definition of the Radon transform by line integrals, the matrix $R_t \in \mathbb{R}_{\geq 0}^{M \times N}$ has only nonnegative entries and hence satisfies the assumption for Theorem 2.1 and 2.2. Furthermore, N denotes here the number of pixels in the original image and M is given by the product $M := |\mathcal{I}(t)|n_S$, where n_S is the number of detection points.

3.1 Results and Discussion

For a qualitative evaluation of the proposed NMF approaches, we consider in the following sections two simulated datasets. Due to the known ground truth in both cases, we are able to measure the performance of each method via computing the mean of the Peak Signal to Noise Ratio (PSNR) and the mean of the Structural Similarity Index Measure (SSIM) index [3] over all time steps for every experiment.

For each dataset, the parameters of all methods are chosen empirically to provide good reconstructions. For the NMF models of the joint reconstruction and

low-rank decomposition approach, we restrict ourselves to the total variation penalty term on B to provide some denoising effect on the spatial features and the ℓ^2 penalty on C for the time features, since we expect and enforce smooth changes in time. We consider the standard case for the TV term with the default pixel neighbourhood and choose the smoothing parameter $\varepsilon_{\text{TV}} = 10^{-5}$ relatively small.

Furthermore, for both datasets we measure different angles at each time step based on a tiny golden angle sampling [43] using consecutive projections with increasing angle of $\varphi = 32.039\dots$, such that projection angles are not repeated. Nevertheless, we remind that we keep the total number of observed angles constant for each time step.

For all considered approaches we use the unfiltered backprojection, given by the adjoint of the Radon transform, applied to the noisy data matrix Y as the initialisation for the reconstruction matrix X . In case of the NMF approaches, the matrices B and C are initialised via SVD of X based on [2]. After the initialisation and at every iteration of the NMF algorithm, a suitable projection step for small values is performed to prevent numerical instabilities and zero entries during the multiplicative algorithm [10].

The algorithms were implemented with MATLAB® R2019b and run on an Intel® Core™ i7-7700K quad core CPU @4.20 GHz with 32 GB of RAM.

In Table 1, a list of all considered algorithms is provided.

To this end we present a summary and short explanation of all considered algorithms in this experimental section in Table 1.

| Algorithm | Description |
|------------|---|
| BC | Joint reconstruction and feature extraction with the NMF model BC without constructing X , see algorithm in Theorem 2.2 |
| BC-X | Joint reconstruction and feature extraction with NMF model BC-X and explicit construction of X , see algorithm in Theorem 2.1 |
| sBC | Joint reconstruction and feature extraction method with NMF model sBC for stationary operator, see algorithm in Corollary 2.1 |
| gradTV | Low-rank based reconstruction method for X , see Algorithm 1 |
| gradTV_PCA | Separated reconstruction and feature extraction with Algorithm 1 and subsequent PCA computation |
| gradTV_NMF | Separated reconstruction and feature extraction with Algorithm 1 and subsequent NMF computation (6) |

Table 1: Summary and short explanation of considered algorithms in the experimental section.

3.1.1 Shepp–Logan Phantom

This synthetic dataset consists of a dynamic two-dimensional Shepp–Logan phantom with $T = 100$ and spatial size 128×128 , see Figure 1 for the ground-truth. During the whole time, two of the inner ellipsoids change their intensities sinusoidally with different frequencies while the rest of the phantom remains constant.

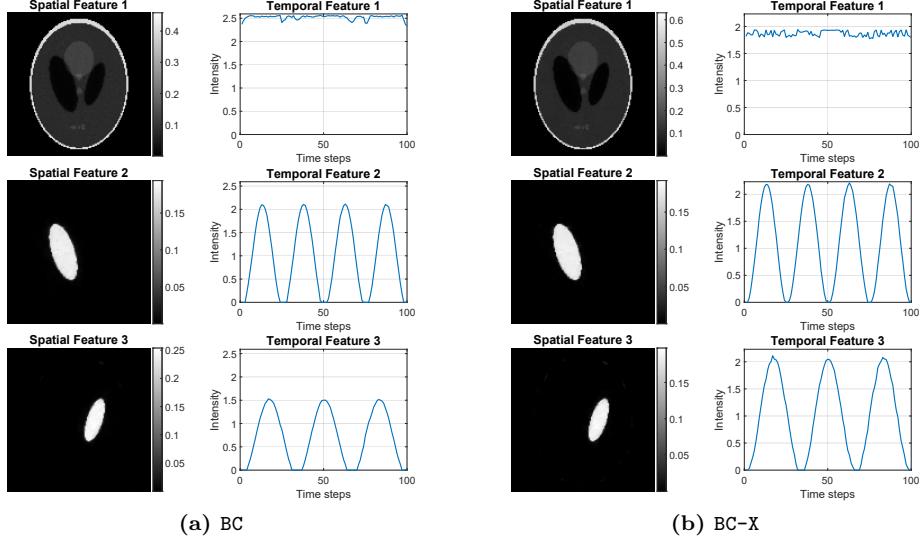


Figure 3: Leading extracted features of the dynamic Shepp–Logan phantom based on the models BC and BC–X for $|\mathcal{I}_t| = 6$ and 1% Gaussian noise.

In the following, we perform a variety of experiments for $|\mathcal{I}_t| \in \{2, \dots, 12\}$ with 1% and 3% Gaussian noise respectively. For all cases, we choose $K = 5$ for the number of the NMF features. In such a way, the NMF is also capable to approximate minor characteristics such as noise or other artefacts of the reconstruction matrix besides the three main features.

The parameters of all methods were determined empirically and are displayed in Appendix C for both noise levels. The stopping criterion for all methods is met, if 1200 iteration steps are reached or if the relative change of all matrices B, C and X goes below $5 \cdot 10^{-5}$.

We show first some results for the case with $|\mathcal{I}_t| = 6$ and 1% Gaussian noise are shown in Figure 3 for the joint NMF methods and Figure 4 for the separate reconstruction and extraction. The order of shown features is based on the singular values of B for `gradTV_PCA` and on the ℓ_2 -norm of the spatial features for NMF approaches.

In this case, all considered approaches are able to successfully identify the constant and dynamic parts of the dataset and extract meaningful spatial and temporal features. The extracted spatial features of BC, BC–X and `gradTV_NMF` show very clearly the dynamic and non–dynamic parts of the Shepp–Logan phantom. However, the spatial features of `gradTV_NMF` are slightly more blurred and affected by minor artefacts especially in both dynamic features. This underlines the positive effect of the separate TV regularisation on the spatial feature matrix B in the joint methods. In contrast, `gradTV_PCA` is able to identify the

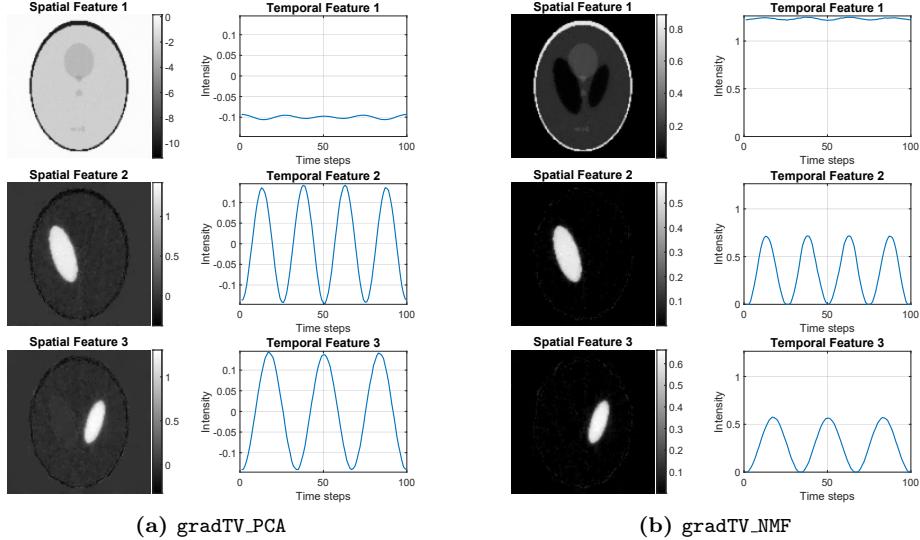


Figure 4: Leading extracted features of the dynamic Shepp–Logan phantom based on the models `gradTV_PCA` and `gradTV_NMF` for $|\mathcal{I}_t| = 6$ and 1% Gaussian noise.

main components of the dataset correctly, but there is a clear corruption of the dynamic features with other parts from the phantom. Furthermore, all spatial features contain negative parts due to the non-existent nonnegativity constraint of the `gradTV_PCA` approach which makes their interpretation more challenging. Hence, the additional nonnegativity constraint of the NMF methods improve significantly the quality and interpretability of the extracted components in comparison with the PCA based extraction method.

The temporal features of all methods are clearly extracted and are consistent with the underlying ground truth of the dataset. However, we note that BC and BC-X have a slight difficulty to resolve the lower intensity part close to 0, which is probably due to the multiplicative structure of the algorithms.

Similar observations can be made for the case $|\mathcal{I}_t| = 6$ and 3% Gaussian noise. We present the reconstructed features in in Figure 5 for BC and `gradTV_PCA` only. The higher amount of noise can be observed especially in the spatial features of `gradTV_PCA`, whereas it only has a slight effect in the BC model.

Finally, we present the reconstructed features with BC and BC-X in Figure 6 for $|\mathcal{I}_t| = 3$, i.e. only three three angles per time step, with noise level of 1%. The major difference to the previous cases can be seen in the results of the BC model. Here, the method splits up the dynamics of the right ellipse into two different temporal features, such that the true dynamics are not retained. However, the BC-X approach perform remarkably well with respect to the feature extraction despite the rather low number of projection angles. This might indicate, that

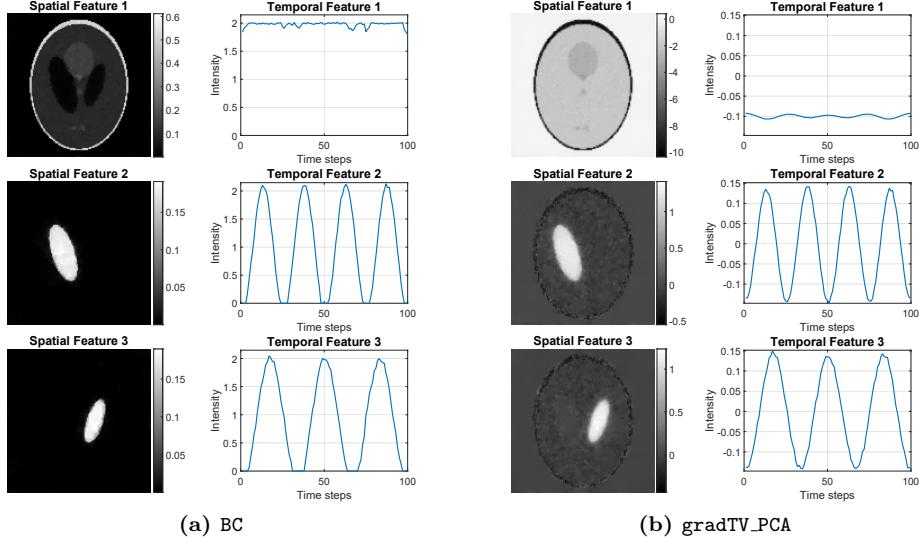


Figure 5: Leading extracted features of the dynamic Shepp–Logan phantom based on the models BC and gradTV_PCA for $|\mathcal{I}_t| = 6$ and 3% Gaussian noise.

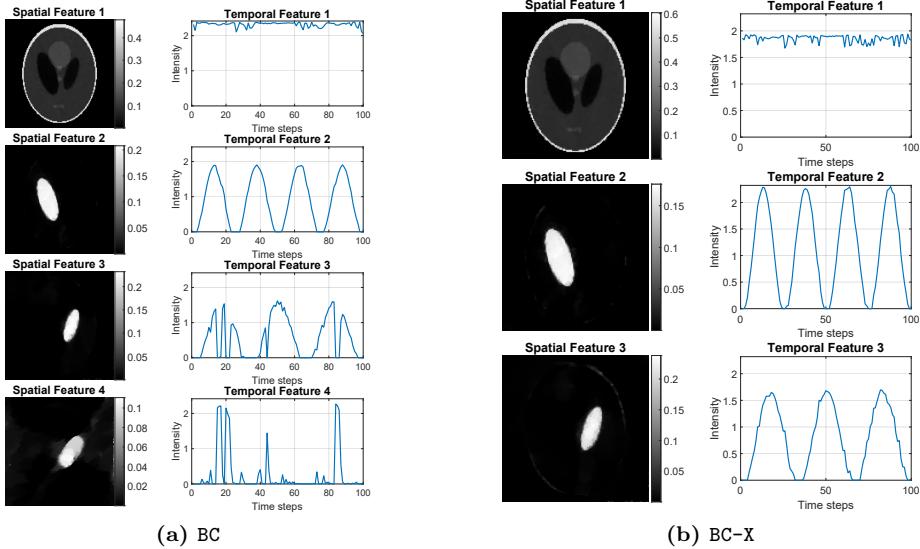


Figure 6: Leading extracted features of the dynamic Shepp–Logan phantom based on the models BC and BC-X for $|\mathcal{I}_t| = 3$ and 1% Gaussian noise.

enforcing the reconstruction X to have small data error helps in the BC-X model to stabilise the reconstruction in highly sparse data settings.

Let us shortly discuss other considered values of $|\mathcal{I}_t|$, that are not shown here. First of all, the performance of `gradTV_PCA` and `gradTV_NMF` with respect to the feature extraction behaves very similar for both noise cases. Besides the above mentioned drawbacks, both approaches give remarkably consistent results especially for low number of angles and do not tend as much to split up features like in BC and BC-X. The latter occurs in different degrees for several numbers of angles. For 1% noise, it occurs for $|\mathcal{I}_t| \in \{3, 7, 8, 10\}$ in BC and for $|\mathcal{I}_t| = 10$ in BC-X. In the case of a noise level of 3%, the split up effect only occurs for $|\mathcal{I}_t| = 10$ in BC. However, for $|\mathcal{I}_t| = 10$, it is possible to partially recover the correct temporal feature by simply adding up both features. Nevertheless, both approaches provide better reconstruction quality of X than `gradTV` as we will discuss in the following.

Quantitative evaluation Let us now discuss the quantitative reconstruction quality for all methods. In Figure 7 and 8 we show the mean PSNR and SSIM of the reconstruction for 1% and 3% noise over all time steps for all considered numbers of projection angles. Note that for the NMF model BC-X, we compute the quality measures for X . The same goes for `gradTV`, where we only compute the quality measures of X after the reconstruction procedure independently of the subsequent feature extraction method. In the case of BC, the reconstruction is computed as $X = BC$.

As expected, the reconstruction quality tends to get better if more angles per time step are considered. More importantly, we see that it is possible to obtain reasonable reconstructions with just a few projections per time step especially in the case of the joint reconstruction and feature extraction method via the NMF approach. In particular, we reach a stable reconstruction quality already with 5 or more angles for both joint methods and 1% noise.

The BC model clearly performs best with respect to the reconstruction quality. For almost every number of angles, the mean PSNR and SSIM values outperform the ones of the BC-X and `gradTV` method for both noise levels. In the case of 3% noise (see Figure 8) we can see that `gradTV` performs slightly better than BC-X in most of the cases in terms of their SSIM values. Still, the mean PSNR values of `gradTV` are significantly lower than the ones in BC-X for all numbers of angles. A selection of reconstructions for the experiments in Figure 7 and 8 are provided as videos in the Supplementary files.

Note that for BC-X, it is also possible to compute the reconstruction based on the decomposition $B \cdot C$ instead of the joint reconstruction X in the algorithm. Interestingly, our experiments showed that the reconstruction quality of $B \cdot C$ is in almost all cases better than the one of the matrix X itself and also mostly outperforms the `gradTV` approach. We believe, that this is due to the stronger regularising effect on the components B and C , which especially influences SSIM.

The computation times for the reconstruction and feature extraction with 1%

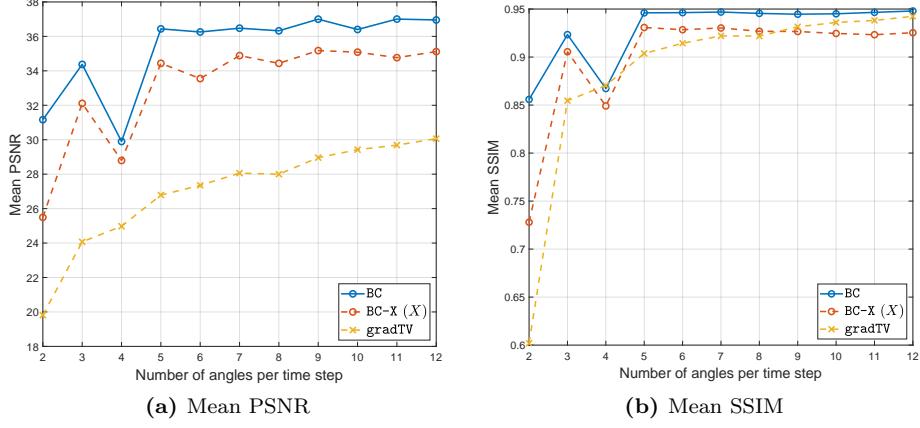


Figure 7: Mean PSNR and SSIM values of the reconstructions of the dynamic Shepp–Logan phantom with 1% Gaussian noise for different numbers of projection angles.

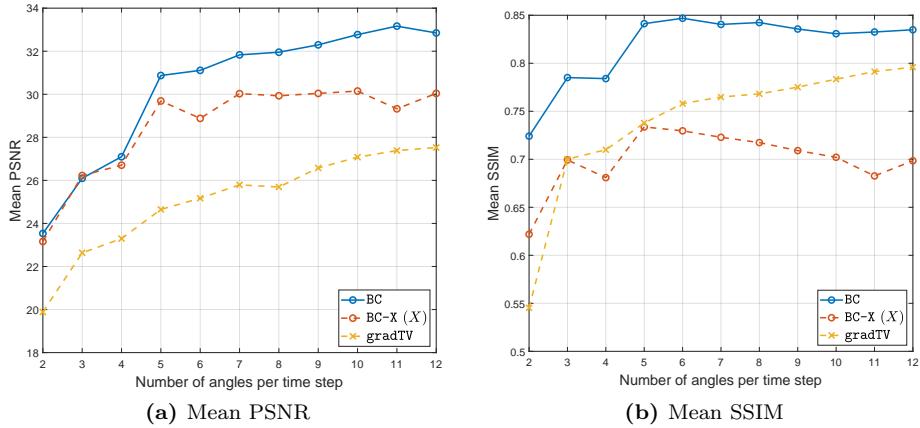


Figure 8: Mean PSNR and SSIM values of the reconstructions of the dynamic Shepp–Logan phantom with 3% Gaussian noise for different numbers of projection angles.

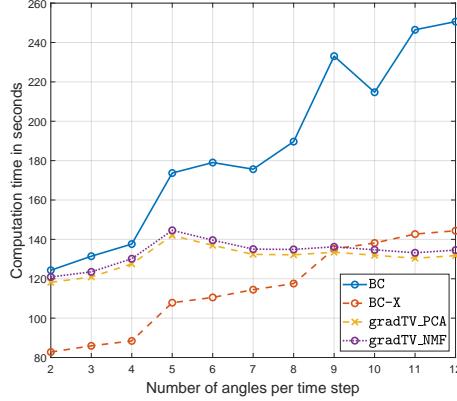


Figure 9: Needed time in seconds for the reconstruction and feature extraction of the dynamic Shepp–Logan phantom with 1% Gaussian noise.

noise for all algorithms until the stopping criterion is fulfilled are shown Figure 9. As expected, the computation time tends to increase with the number of projection angles and, considering all methods, ranges approximately from 1 to 5 minutes. For $|\mathcal{I}_t| \leq 8$, the BC–X method is the fastest while it is outperformed by `gradTV_PCA` for $|\mathcal{I}_t| \geq 9$. `gradTV_NMF` and BC needs more time in all experiments compared to `gradTV_PCA`. The significant temporal difference between BC–X and BC is due to its higher computational complexity: Owing to the model formulation of BC with the discrepancy term $\|R_t(BC)_{\bullet,t} - Y_{\bullet,t}\|_2^2$, the update rules in Theorem 2.2 for both matrices B and C contain the discretised Radon transform R_t . This is in contrast to the BC–X algorithm, where R_t only appears in the update rule of X .

Based on the presented results for the dynamic Shepp–Logan phantom, we can conclude that the joint approaches BC and BC–X outperform both other methods with respect to the reconstruction quality and for most cases of the extracted features. Nevertheless, the models `gradTV_PCA` and `gradTV_NMF` give remarkably consistent and stable results of the extracted components throughout all numbers of angles. Furthermore, the nonnegativity constraint of the NMF improves significantly the interpretability and quality of the extracted spatial features.

Stationary Operator As we have seen, the computational complexity of the BC model with the non-stationary operator is clearly higher than for all other cases. Thus, let us now consider the possibility to speed up the reconstructions with a stationary operator, which leads us to the complexity reduced formulation presented in Corollary 2.1 as the `sBC` model. Here we present similarly to the case above experiments with the dynamic Shepp–Logan phantom for $|\mathcal{I}_t| \in \{2, \dots, 30\}$ and 1% Gaussian noise, as we primarily aim to illustrate the reduction of the computational cost. Furthermore, the same hyperparameters and stopping criteria are used as before.

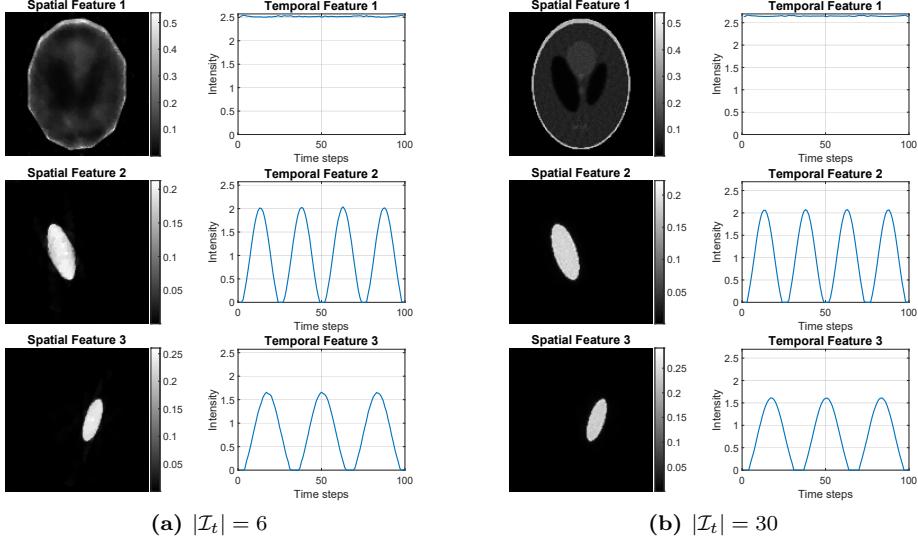


Figure 10: Leading extracted features of the dynamic Shepp–Logan phantom based on the model sBC for $|\mathcal{I}_t| \in \{6, 30\}$ and 1% Gaussian noise.

The reconstructed features for the cases $|\mathcal{I}_t| = 6$ and $|\mathcal{I}_t| = 30$ are shown in Figure 10. In particular, comparing the results in Figure 10a to the corresponding results of BC in Figure 3a, one can immediately see a significant difference between the extracted spatial features. This is clearly due to the fact that the same projection angles are used at every time step and the individual projection directions are clearly visible for stationary model sBC. Consequently, the details in the Shepp-Logan phantom are not well recovered, such that the extracted constant feature is significantly inferior to the one of BC. As one would expect, more projection angles per time step are needed to reconstruct finer details, this effect can be clearly seen for 30 angles in Figure 10b.

However, all temporal basis functions with sBC for $|\mathcal{I}_t| = 6$ are remarkably well reconstructed despite the low number of projection angles. This is also true for the other considered values of $|\mathcal{I}_t|$. Moreover, we observe that sBC is able to extract the correct three main features for every $|\mathcal{I}_t| \in \{2, \dots, 30\}$. Even for $|\mathcal{I}_t| = 2$ and the quality of the dynamic time features are similar to the ones in Figure 10.

This behaviour is different from the dynamic case discussed above. The reason for this is probably based on the different projection directions at every time step in the dynamic case, which results in directional dependencies of the occurring reconstruction artefacts in contrast to the stationary case. This can make it difficult for the NMF to distinguish the main features in the non-stationary case and thus leads to a more stable feature extraction in the here presented stationary case.

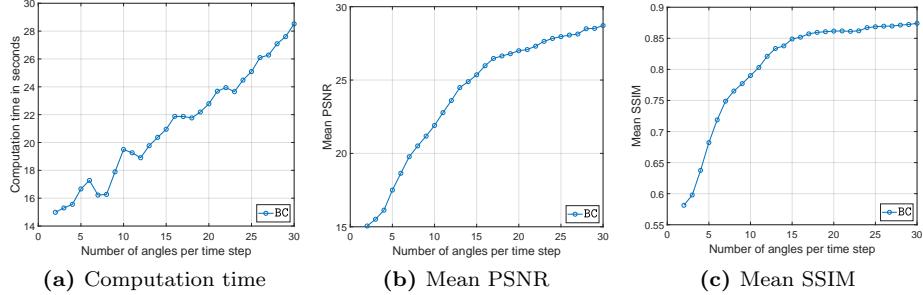


Figure 11: Needed time in seconds, mean PSNR and mean SSIM values of the reconstructions of the dynamic Shepp–Logan phantom with 1% Gaussian noise for the stationary case **sBC** and different numbers of projection angles.

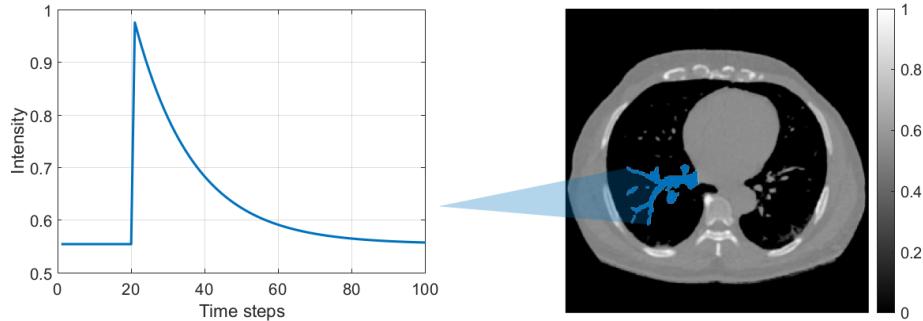


Figure 12: Illustration of Vessel Phantom dataset consisting of $T = 100$ phantoms of dimension 264×264 , where the intensity of the blue highlighted area changes according to blue curve on the left.

The quantitative measures are shown in Figure 11 for all experiments. Comparing the computation time of BC with the one of **sBC**, we obtain a clear speed-up by a factor of 10–20 with the stationary model. However, as expected, comparing Figure 11b and 11c with the quality measures of BC in Figure 7, one can observe that significantly more projection angles per time step are needed in the stationary case to provide a sufficient reconstruction quality. In conclusion, we can say that the **sBC** model is especially recommended if one is primarily interested in the dynamics of the system under consideration, as we could extract the temporal basis functions stably for all considered angles with $|\mathcal{I}_t| \geq 2$.

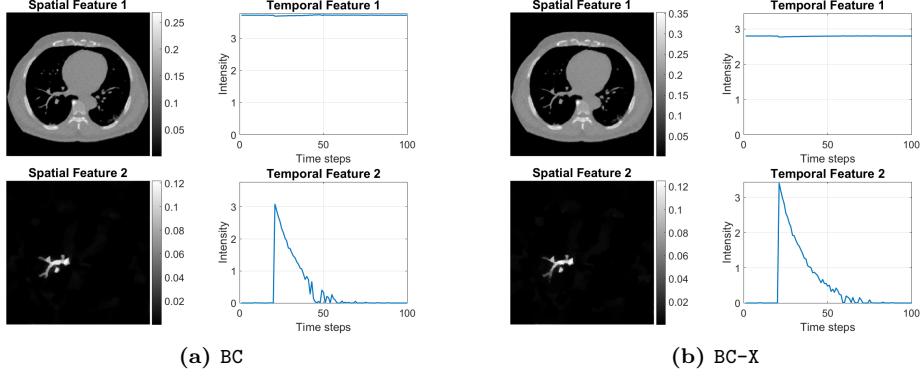


Figure 13: Leading extracted features of the vessel phantom based on the models BC and BC-X for $|\mathcal{I}_t| = 12$ and 1% Gaussian noise.

3.1.2 Vessel Phantom

The second test case is based on a CT scan of a human lung², see Figure 12. Here, the decomposition is given by the constant background and a segmented vessel that exhibits a sudden increase in attenuation followed by an exponential decay. This could for instance represent the injection of a tracer to the blood stream.

In contrast to the previous dataset, we perform only selected experiments for specific choices of noise levels and numbers of projection angles. More precisely, we present results for 1% Gaussian noise together with $|\mathcal{I}_t| \in \{7, 12\}$ and 3% Gaussian noise with $|\mathcal{I}_t| = 12$. In all cases, we choose $K = 4$ NMF features. Furthermore, the stopping criterion from the experiments with the dynamic Shepp-Logan phantom is changed for this dataset in such a way, that the maximum number of iterations is raised to 1400 to ensure sufficient convergence. The regularisation parameters of all methods are chosen empirically and are displayed in Appendix C.

Figure 13 and 14 show the feature extraction results for the noise level of 1% and $|\mathcal{I}_t| = 12$, where all approaches are able to extract both the main constant and dynamic component of the underlying ground truth. The order of the features here is based on a manual sorting.

Similar to the results of the Shepp-Logan phantom in Section 3.1.1, the joint methods BC and BC-X have difficulties to recover the lower intensities in the temporal features, whereas gradTV_PCA produce slight artefacts in the dynamic spatial feature due to the missing nonnegativity constraint. In addition, gradTV_NMF is able to recover more details in the vessel compared to the joint approaches. This is due to the relatively high choice of the total variation regularisation parameter τ in BC and BC-X to ensure a sufficient denoising effect on the matrix

²The phantom is based on the CT scans in the *ELCAP Public Lung Image database*: <http://www.via.cornell.edu/lungdb.html>

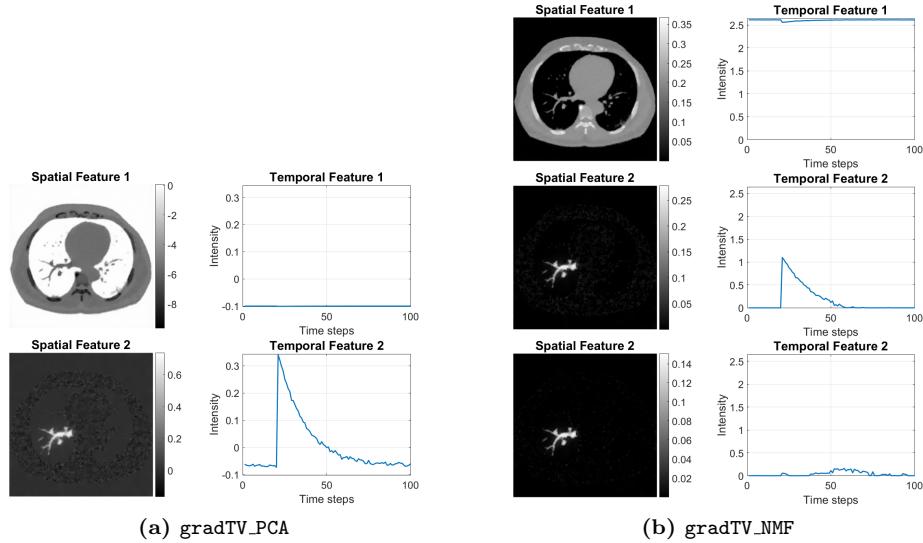


Figure 14: Leading extracted features of the vessel phantom based on the models gradTV_PCA and gradTV_NMF for $|I_t| = 12$ and 1% Gaussian noise.

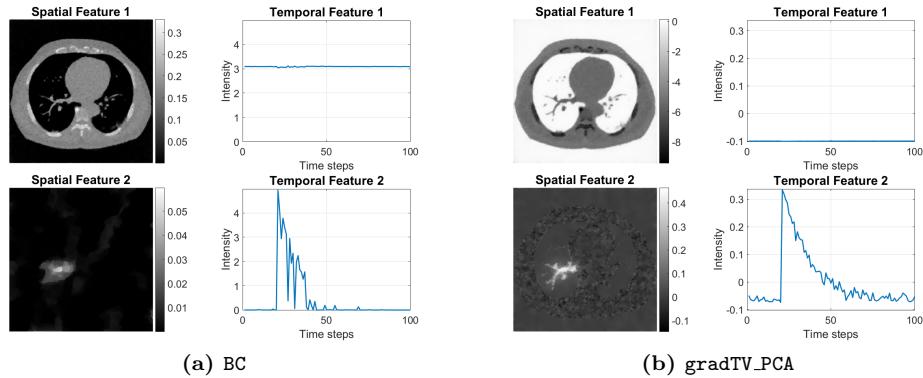


Figure 15: Leading extracted features of the vessel phantom based on the models BC and gradTV_PCA for $|I_t| = 12$ and 3% Gaussian noise.

| Noise | $ \mathcal{I}_t $ | BC | | BC-X | | gradTV | |
|-------|-------------------|----------|----------|--------|--------|--------|--------|
| | | PSNR | SSIM | PSNR | SSIM | PSNR | SSIM |
| 1% | 7 | (34.130) | (0.9016) | 32.969 | 0.8382 | 31.463 | 0.8414 |
| 1% | 12 | 35.050 | 0.9068 | 33.919 | 0.8496 | 34.309 | 0.8839 |
| 3% | 12 | 30.148 | 0.7484 | 28.119 | 0.5708 | 29.375 | 0.6698 |

Table 2: Mean PSNR and SSIM values of the reconstruction results of the vessel phantom for different noise levels and numbers of projection angles. Values in brackets indicate that the dynamic part of the dataset in the corresponding experiment could not be reconstructed sufficiently well.

B. The low peak in the second temporal feature of `gradTV_NMF` is likely caused by the choice of the ℓ_2 regularisation parameter $\hat{\mu}_C$.

Further experiments show that the quality of the extracted components of BC-X decreases steadily for lower angles until the main features cannot be identified anymore for $|\mathcal{I}_t| \leq 8$. BC produces inferior results and cannot extract reasonable components anymore for $|\mathcal{I}_t| \leq 10$.

In comparison, both separated approaches `gradTV_PCA` and `gradTV_NMF` are still able to extract decent features for $|\mathcal{I}_t| = 7$. For $|\mathcal{I}_t| \leq 6$, the performance of both methods decreases significantly.

Similar results for `gradTV_PCA` could be obtained for 3% noise and $|\mathcal{I}_t| = 12$, which are shown in Figure 15b. Its constant feature is inferior to the one of BC in Figure 15a due to the additional nonnegativity constraint of the NMF model. However, the details of the vessel in the dynamic spatial feature of BC are lost due to the choice of large regularisation parameter τ and the temporal features are affected by several disturbances. Further tests with the noise level of 3% showed that both joint methods are not able to recover the underlying features for $|\mathcal{I}_t| \leq 10$, while the separated approaches gives still acceptable results for $|\mathcal{I}_t| = 6$.

The reconstruction quality of the experiments are shown in Table 2. Similar to the Shepp-Logan phantom, the joint approach BC produces the best results compared to all other methods in terms of the mean PSNR and SSIM values. Further experiments confirm this observation for $4 \leq |\mathcal{I}_t| \leq 11$.

However, these observations have to be treated with caution. BC is not able to recover the dynamics for $|\mathcal{I}_t| \leq 10$ and 1% noise. In the case of BC-X, the dynamics can be reconstructed to some degree within the angle range $9 \leq |\mathcal{I}_t| \leq 11$, but are not recognizable anymore for $|\mathcal{I}_t| \leq 8$. In the case of 3% Gaussian noise, `gradTV` is still able to give acceptable reconstruction results for $|\mathcal{I}_t| = 10$. For less angles, the reconstructed dynamics of `gradTV` get constantly worse until they are not apparent anymore for $|\mathcal{I}_t| \leq 6$.

The computation times of the experiments in Table 2 range approximately from 7 to 15 minutes. The corresponding reconstructions can be found as video files in the Supplementary information.

4 Conclusion

In this work we consider dynamic inverse problems with the assumption that the target of interest has a low-rank structure and can be efficiently represented by spatial and temporal basis functions. This assumption leads naturally to a reconstruction and low-rank decomposition framework. In particular, we concentrate here on the Nonnegative Matrix Factorisation as decomposition because it exhibits three main advantages:

- i.) It naturally incorporates the physical assumption of nonnegativity
- ii.) Basis functions are not restricted to being orthogonal and therefore correspond more naturally to actual components
- iii.) It allows the flexibility to incorporate separate regularisation on each of the factorisation matrices

In particular, the last point is of importance, as it allows to consider different regularisers for spatial and temporal basis functions, and as such can be tailored to different applications.

We then proposed two approaches to obtain a joint reconstruction and low-rank decomposition based on the NMF, termed BC-X and BC. Both methods performed better than a baseline method, that computes a reconstruction with low-rank constraint followed by a subsequent decomposition. In particular, the second BC model has shown to have a stronger regularising effect on the reconstructed features as well as the reconstruction, which can be simply obtained as $X = BC$. We believe this is due to the fact, that only the decomposition is recovered during the reconstruction without the need to build the reconstruction X explicitly and hence the resulting features at the end exhibit a higher regularity. More importantly, if one considers a stationary operator in the complexity reduced sBC model we can obtain a considerable computational speed-up. Even though, due to constant projection angles the spatial basis functions are not as well recovered as in the non-stationary case, but the temporal features can be nicely extracted even for as low as 2 angles. This might be especially of interest in applications, where one is primarily interested in the underlying dynamics of the imaged target.

The primary limitation of the presented approach is the assumption on the decomposition of the target into spatial and temporal basis functions, as this does not allow for spatial movements in the target. However it opens up the possibility of combination with other methods, that do in fact allow for movements but assume a brightness consistency in the target, such as the optical flow constraint in CT [5]. Furthermore, the presented low-rank decomposition may be combined with a morphological motion model [20] to allow for a flexible and general model for dynamic inverse problems.

Acknowledgments

This project was supported by the Deutsche Forschungsgemeinschaft (DFG) within the framework of GRK 2224/1 “ π^3 : Parameter Identification - Analysis, Algorithms, Applications”. This work was partially supported by the Academy of Finland Project 312123 (Finnish Centre of Excellence in Inverse Modelling and Imaging, 2018–2025) EPSRC grant EDCLIRS (EP/N022750/1) as well as CMIC-EPSRC platform grant (EP/M020533/1).

References

- [1] D. BÖHNING AND B. G. LINDSAY, *Monotonicity of quadratic approximation algorithms*, Annals of the Institute of Statistical Mathematics, 40 (1988), pp. 641–663.
- [2] C. BOUTSIDIS AND E. GALLOPOULOS, *Svd based initialization: A head start for nonnegative matrix factorization*, Pattern Recognition, 41 (2008), pp. 1350–1362.
- [3] D. BRUNET, E. R. VRSCAY, AND Z. WANG, *On the mathematical properties of the structural similarity index*, IEEE Transactions on Image Processing, 21 (2012), pp. 1488–1499.
- [4] T. A. BUBBA, M. MÄRZ, Z. PURISHA, M. LASSAS, AND S. SILTANEN, *Shearlet-based regularization in sparse dynamic tomography*, in Wavelets and Sparsity XVII, vol. 10394, International Society for Optics and Photonics, 2017, p. 103940Y.
- [5] M. BURGER, H. DIRKS, L. FRERKING, A. HAUPTMANN, T. HELIN, AND S. SILTANEN, *A variational reconstruction method for undersampled dynamic x-ray tomography based on physical motion models*, Inverse Problems, 33 (2017), p. 124008.
- [6] M. BURGER, H. DIRKS, AND C.-B. SCHONLIEB, *A variational model for joint motion estimation and image reconstruction*, SIAM Journal on Imaging Sciences, 11 (2018), pp. 94–128.
- [7] E. J. CANDÈS, X. LI, Y. MA, AND J. WRIGHT, *Robust principal component analysis?*, Journal of the ACM, 58 (2011).
- [8] B. CHEN, J. ABASCAL, AND M. SOLEIMANI, *Extended joint sparsity reconstruction for spatial and temporal ERT imaging*, Sensors, 18 (2018), p. 4014.
- [9] C. CHEN AND O. ÖKTEM, *Indirect Image Registration with Large Diffeomorphic Deformations*, SIAM Journal on Imaging Sciences, 11 (2018), pp. 575–617.

- [10] A. CICHOCKI, R. ZDUNEK, A. H. PHAN, AND S. AMARI, *Nonnegative Matrix and Tensor Factorizations: Applications to Exploratory Multi-Way Data Analysis and Blind Source Separation*, Wiley Publishing, 2009.
- [11] C. DE MOL, *Blind deconvolution and nonnegative matrix factorization*, Oberwolfach Reports 51/2012, Mathematisches Forschungsinstitut Oberwolfach, 2012.
- [12] M. DEFRISE, C. VANHOVE, AND X. LIU, *An algorithm for total variation regularization in high-dimensional linear problems*, Inverse Problems, 27 (2011), p. 065002.
- [13] C. DING, X. HE, AND H. SIMON, *On the equivalence of nonnegative matrix factorization and spectral clustering*, in SIAM International Conference on Data Mining, 2005.
- [14] N. DJURABEKOVA, A. GOLDBERG, A. HAUPTMANN, D. HAWKES, G. LONG, F. LUCKA, AND M. BETCKE, *Application of proximal alternating linearized minimization (PALM) and inertial PALM to dynamic 3D CT*, in 15th International Meeting on Fully Three-Dimensional Image Reconstruction in Radiology and Nuclear Medicine, vol. 11072, International Society for Optics and Photonics, 2019, p. 1107208.
- [15] L. FENG, R. GRIMM, K. T. BLOCK, H. CHANDARANA, S. KIM, J. XU, L. AXEL, D. K. SODICKSON, AND R. OTAZO, *Golden-angle radial sparse parallel mri: combination of compressed sensing, parallel imaging, and golden-angle radial sampling for fast and flexible dynamic volumetric mri*, Magnetic resonance in medicine, 72 (2014), pp. 707–717.
- [16] P. FERNSEL AND P. MAASS, *A survey on surrogate approaches to non-negative matrix factorization*, Vietnam Journal of Mathematics, 46 (2018), pp. 987–1021.
- [17] C. FÉVOTTE, N. BERTIN, AND J.-L. DURRIEU, *Nonnegative matrix factorization with the itakura-saito-divergence: With application to music analysis*, Neural Computation, 21 (2009), pp. 793–830.
- [18] T. GOLDSTEIN AND S. OSHER, *The split Bregman method for l_1 -regularized problems*, SIAM Journal on Imaging Sciences, 2 (2009), pp. 323–343.
- [19] G. H. GOLUB AND C. F. VAN LOAN, *Matrix Computations*, The Johns Hopkins University Press, 4 ed., 2013.
- [20] B. GRIS, C. CHEN, AND O. ÖKTEM, *Image reconstruction through metamorphosis*, Inverse Problems, 36 (2020), p. 025001 (27pp).
- [21] B. KLINGENBERG, J. CURRY, AND A. DOUGHERTY, *Non-negative matrix factorization: Ill-posedness and a geometric algorithm*, Pattern Recognition, 42 (2009), pp. 918–928.

- [22] D. KRESSNER AND A. USCHMAJEW, *On low-rank approximability of solutions to high-dimensional operator equations and eigenvalue problems*, Linear Algebra and its Applications, 493 (2016), pp. 556–572.
- [23] K. LANGE, *Optimization*, vol. 95 of Springer Texts in Statistics, Springer-Verlag New York, 2 ed., 2013.
- [24] K. LANGE, D. R. HUNTER, AND I. YANG, *Optimization transfer using surrogate objective functions*, Journal of Computational and Graphical Statistics, 9 (2000), pp. 1–20.
- [25] L. LECHARLIER AND C. DE MOL, *Regularized blind deconvolution with poisson data*, Journal of Physics: Conference Series, 464 (2013), p. 012003.
- [26] D. D. LEE AND H. S. SEUNG, *Learning the parts of objects by non-negative matrix factorization*, Nature, 401 (1999), pp. 788–791.
- [27] D. D. LEE AND H. S. SEUNG, *Algorithms for non-negative matrix factorization*, in Proceedings of the 13th International Conference on Neural Information Processing Systems, 2000, pp. 535–541.
- [28] J. LEUSCHNER, M. SCHMIDT, P. FERNSEL, D. LACHMUND, T. BOSKAMP, AND P. MAASS, *Supervised non-negative matrix factorization methods for MALDI imaging applications*, Bioinformatics, 35 (2018), pp. 1940–1947.
- [29] S. G. LINGALA, Y. HU, E. V. R. DiBELLA, AND M. JACOB, *Accelerated dynamic MRI exploiting sparsity and low-rank structure: k-t SLR*, IEEE Transactions on Medical Imaging, 30 (2011), pp. 1042–1054.
- [30] F. LUCKA, N. HUYNH, M. BETCKE, E. ZHANG, P. BEARD, B. COX, AND S. ARRIDGE, *Enhancing compressed sensing 4D photoacoustic tomography by simultaneous motion estimation*, SIAM Journal on Imaging Sciences, 11 (2018), pp. 2224–2253.
- [31] M. LUSTIG, J. M. SANTOS, D. L. DONOHO, AND J. M. PAULY, *k_t SPARSE: High frame rate dynamic MRI exploiting spatio-temporal sparsity*, in Proceedings of the 13th annual meeting of ISMRM, Seattle, vol. 2420, 2006.
- [32] E. NIEMI, M. LASSAS, A. KALLONEN, L. HARHANEN, K. HÄMÄLÄINEN, AND S. SILTANEN, *Dynamic multi-source x-ray tomography using a space-time level set method*, Journal of Computational Physics, 291 (2015), pp. 218–237.
- [33] J. P. OLIVEIRA, J. M. BIOUCAS-DIAS, AND M. A. T. FIGUEIREDO, *Review: Adaptive total variation image deblurring: A majorization-minimization approach*, Signal Processing, 89 (2009), pp. 1683–1693.
- [34] P. PAATERO AND U. TAPPER, *Positive matrix factorization: A non-negative factor model with optimal utilization of error estimates of data values*, Environmetrics, 5 (1994), pp. 111–126.

- [35] K. PEARSON, *On lines and planes of closest fit to systems of points in space*, The London, Edinburgh, and Dublin Philosophical Magazine and Journal of Science, 2 (1901), pp. 559–572.
- [36] U. SCHMITT AND A. K. LOUIS, *Efficient algorithms for the regularization of dynamic inverse problems: I. Theory*, Inverse Problems, 18 (2002), p. 645.
- [37] U. SCHMITT, A. K. LOUIS, C. WOLTERS, AND M. VAUHKONEN, *Efficient algorithms for the regularization of dynamic inverse problems: II. Applications*, Inverse Problems, 18 (2002), p. 659.
- [38] J. A. STEEDEN, G. T. KOWALIK, O. TANN, M. HUGHES, K. H. MORTENSEN, AND V. MUTHURANGU, *Real-time assessment of right and left ventricular volumes and function in children using high spatiotemporal resolution spiral bssfp with compressed sensing*, Journal of Cardiovascular Magnetic Resonance, 20 (2018), p. 79.
- [39] M. TAO AND X. YUAN, *Recovering low-rank and sparse components of matrices from incomplete and noisy observations*, SIAM Journal on Optimization, 21 (2011), pp. 57–81.
- [40] B. R. TRÉMOULHÉAC, *Low-rank and sparse reconstruction in dynamic magnetic resonance imaging via proximal splitting methods*, PhD thesis, University College London, 2014.
- [41] B. R. TRÉMOULHÉAC, N. DIKAIOS, D. ATKINSON, AND S. ARRIDGE, *Dynamic MR image reconstruction-separation from undersampled (k,t) -space via low-rank plus sparse prior*, IEEE Transactions on Medical Imaging, 33 (2014), pp. 1689–1701.
- [42] A. USCHMAJEW, *On low-rank approximation in tensor product Hilbert spaces*, doctoral thesis, Technical University of Berlin, 2013.
- [43] S. WUNDRAK, J. PAUL, J. ULRICI, E. HELL, AND V. RASCHE, *A small surrogate for the golden angle in time-resolved radial mri based on generalized fibonacci sequences*, IEEE transactions on medical imaging, 34 (2014), pp. 1262–1269.
- [44] X. M. YUAN AND J. F. YANG, *Sparse and low-rank matrix decomposition via alternating direction methods*, Pacific Journal of Optimization, 9 (2013), pp. 167–180.

A Optimisation Techniques for NMF Problems

The majority of optimisation techniques for NMF problems are based on alternating minimisation schemes. This is due to the fact that the corresponding cost function in (5) is usually convex in B for fixed C and C for fixed B and non-convex in (B, C) together, which yields algorithms of the form

$$\begin{aligned} B^{[d+1]} &:= \arg \min_{B \geq 0} \mathcal{F}(B, C^{[d]}), \\ C^{[d+1]} &:= \arg \min_{C \geq 0} \mathcal{F}(B^{[d+1]}, C). \end{aligned}$$

Typical minimisation approaches are based on alternating least squares methods, multiplicative algorithms as well as projected gradient descent and quasi-newton methods [10]. In this work, we focus on the derivation of multiplicative update rules based on the so-called *majorise minimisation* (MM) principle [23]. This approach allows the derivation of multiplicative update rules for non-standard NMF cost functions and gives therefore the flexibility to adjust the discrepancy and penalty terms according to the NMF model motivated by the corresponding application [16]. What is more, the update rules consist only of multiplications and summations of matrices, which allow very simple implementations of the algorithms and ensure automatically the nonnegativity of the iterates B and C without the need of any inversion process, provided they are initialised nonnegative.

A.1 Multiplicative Algorithms

The works of Lee and Seung [26, 27] brought much attention to NMF methods in general and, in particular, the multiplicative algorithms, which they derived based on the MM principle for the standard case with the Frobenius norm and the Kullback–Leibler divergence as discrepancy terms.

The main idea of the MM approach is to replace the original cost function \mathcal{F} by a majorizing so-called *surrogate function* $\mathcal{Q}_{\mathcal{F}}$, which is easier to minimise and leads to the desired multiplicative algorithms due to its tailored construction.

Definition A.1 (Surrogate Function). Let $\Omega \subset \mathbb{R}^n$ be an open subset and $\mathcal{F} : \Omega \rightarrow \mathbb{R}$ a function. Then $\mathcal{Q}_{\mathcal{F}} : \Omega \times \Omega \rightarrow \mathbb{R}$ is called a **surrogate function** or **surrogate** of \mathcal{F} , if it fulfills the following properties:

- i) $\mathcal{Q}_{\mathcal{F}}(x, \tilde{x}) \geq \mathcal{F}(x)$ for all $x, \tilde{x} \in \Omega$
- ii) $\mathcal{Q}_{\mathcal{F}}(x, x) = \mathcal{F}(x)$ for all $x \in \Omega$

The minimisation step of the MM approach is then defined by the update rule

$$x^{[d+1]} := \arg \min_{x \in \Omega} \mathcal{Q}_{\mathcal{F}}(x, x^{[d]}), \quad (13)$$

assuming that the $\arg \min_{x \in \Omega} \mathcal{Q}_{\mathcal{F}}(x, \tilde{x})$ exists for all $\tilde{x} \in \Omega$. Due to the defining properties of a surrogate function in Definition A.1, the monotonic decrease of

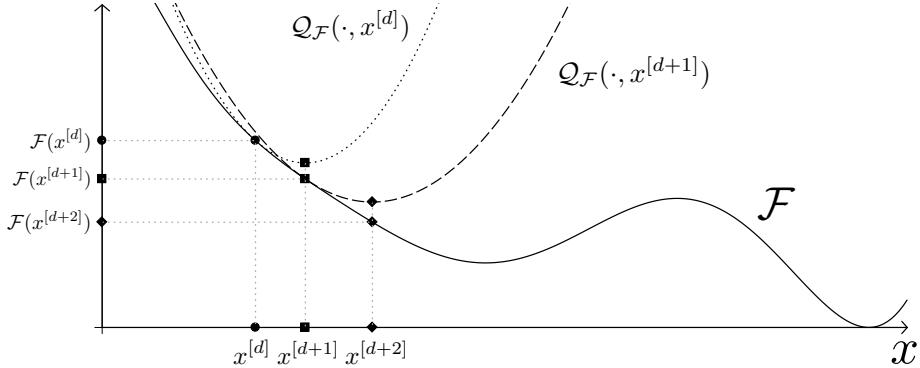


Figure 16: Illustration of two iteration steps of the MM principle for a cost function \mathcal{F} with bounded curvature and a surrogate function $\mathcal{Q}_{\mathcal{F}}$, which is strictly convex in the first argument.

the cost function \mathcal{F} is easily shown:

$$\mathcal{F}(x^{[d+1]}) \leq \mathcal{Q}_{\mathcal{F}}(x^{[d+1]}, x^{[d]}) \leq \mathcal{Q}_{\mathcal{F}}(x^{[d]}, x^{[d]}) = \mathcal{F}(x^{[d]}). \quad (14)$$

This principle is also illustrated in Figure 16. Typical construction techniques lead to surrogate functions, which are strictly convex in the first component to ensure the unique existence of the corresponding minimiser. Furthermore, the surrogates must be constructed in such a way, that the minimisation in Equation (13) yields multiplicative updates to ensure the nonnegativity of the matrix iterates. Finally, another useful property is the separability of $\mathcal{Q}_{\mathcal{F}}$ with respect to the first variable. This ensures, that $\mathcal{Q}_{\mathcal{F}}(x, \tilde{x})$ can be written as a sum, where each component just depends on one entry of x and allows the derivation of the multiplicative algorithm via the zero gradient condition $\nabla_x \mathcal{Q}_{\mathcal{F}} = 0$. One typical construction method is the so-called *quadratic upper bound principle* (QUBP) [1, 23], which forms one of the main approaches to construct suitable surrogate functions for NMF problems. Nice overviews of other construction principles, which will not be used in this work, can be found in [23, 24]. The QUBP is described in the following Lemma.

Lemma A.1. *Let $\Omega \subset \mathbb{R}^n$ be an open and convex subset, $\mathcal{F} : \Omega \rightarrow \mathbb{R}$ twice continuously differentiable with bounded curvature, i.e. there exists a matrix $\Lambda \in \mathbb{R}^{n \times n}$, such that $\Lambda - \nabla^2 \mathcal{F}(x)$ is positive semi-definite for all $x \in \Omega$. We then have*

$$\begin{aligned} \mathcal{F}(x) &\leq \mathcal{F}(\tilde{x}) + \nabla \mathcal{F}(\tilde{x})^\top (x - \tilde{x}) + \frac{1}{2}(x - \tilde{x})^\top \Lambda(x - \tilde{x}) \quad \forall x, \tilde{x} \in \Omega \\ &=: \mathcal{Q}_{\mathcal{F}}(x, \tilde{x}), \end{aligned}$$

where $\mathcal{Q}_{\mathcal{F}}$ is a surrogate function of \mathcal{F} .

This is a classical result based on the second-order Taylor polynomial and will not be proven here.

If the matrix Λ is additionally symmetric and positive definite, it can be shown [16] that the update rule for x according to (13) via the zero gradient condition $\nabla_x \mathcal{Q}_F(x, \tilde{x}) = 0$ gives the unique minimiser

$$x_{\tilde{x}}^* = \tilde{x} - \Lambda^{-1} \nabla \mathcal{F}(\tilde{x}). \quad (15)$$

In this work, we will only apply the QUBP for quadratic cost functions F , whose Hessian is automatically a constant matrix. For these functions, typical choices of Λ are diagonal matrices of the form

$$\Lambda(\tilde{x})_{ii} := \frac{(\nabla^2 f(\tilde{x}) \tilde{x})_i + \kappa_i}{\tilde{x}_i}, \quad (16)$$

which are dependent on the second argument of the corresponding surrogate $\mathcal{Q}_F(x, \tilde{x})$. The parameters $\kappa_i \geq 0$, are constants and have to be chosen depending on the considered penalty terms of the NMF cost function.

The diagonal structure of $\Lambda(\tilde{x})$ ensures its simple invertibility, the separability of the corresponding surrogate and the desired multiplicative algorithms based on (13). Hence, the update rule in (15) can be viewed as a gradient descent approach with a suitable stepsize defined by the diagonal matrix Λ .

B Derivation of the Algorithms

In this section, we derive the multiplicative update rules for the NMF minimisation problems in BC-X and BC.

B.1 Model BC-X

B.1.1 Algorithm for X

We start first of all with the NMF model BC-X and the minimisation with respect to X . The cost function of the NMF problem in BC-X for the minimisation with respect to X reduces to

$$\mathcal{F}(X) := \underbrace{\sum_{t=1}^T \frac{1}{2} \|A_t X_{\bullet,t} - Y_{\bullet,t}\|_2^2 + \frac{\mu_X}{2} \|X\|_F^2 + \lambda_X \|X\|_1}_{=: \mathcal{F}_1(X)} + \underbrace{\frac{\alpha}{2} \|X - BC\|_F^2}_{=: \mathcal{F}_2(X)} \quad (17)$$

by neglecting the constant terms. To apply the QUBP and to avoid fourth-order tensors during the computation of the Hessians, we use the separability of \mathcal{F}_1 with respect to the columns of X , i.e. it can be written as sum, where each term depends only on the respective column $X_{\bullet,t}$. Hence, we write

$$\mathcal{F}_1(X) = \sum_{t=1}^T \left[\frac{1}{2} \|A_t X_{\bullet,t} - Y_{\bullet,t}\|_2^2 + \frac{\mu_X}{2} \|X_{\bullet,t}\|_2^2 + \lambda_X \|X_{\bullet,t}\|_1 \right] =: \sum_{t=1}^T f_t(X_{\bullet,t}).$$

We can assume that X contains only strictly positive entries due to the strict positive initialisations of the multiplicative algorithms. Hence, the functions f_t are twice continuously differentiable despite the occurring ℓ^1 regularisation term. The computations of the gradient and the Hessian of f_t are straightforward and we obtain

$$\begin{aligned}\nabla f_t(X_{\bullet,t}) &= A_t^\top A_t X_{\bullet,t} - A_t^\top Y_{\bullet,t} + \mu_X X_{\bullet,t} + \lambda_X \mathbf{1}_{N \times 1}, \\ \nabla^2 f_t(X_{\bullet,t}) &= A_t^\top A_t + \mu_X I_{N \times N},\end{aligned}$$

where $I_{N \times N}$ is the $N \times N$ identity matrix. Choosing $\kappa_n = \lambda_X$ for all n in (16), we define the surrogate Q_{f_t} according to Lemma A.1. It is then easy to see, that

$$Q_{\mathcal{F}_1}(X, \tilde{X}) := \sum_{t=1}^T Q_{f_t}(X_{\bullet,t}, \tilde{X}_{\bullet,t})$$

defines a separable and convex surrogate function for \mathcal{F}_1 . For \mathcal{F}_2 , we set simply $Q_{\mathcal{F}_2}(X, \tilde{X}) := \alpha/2 \|X - BC\|_F^2$, such that we end up with

$$Q_{\mathcal{F}}(X, A) := Q_{\mathcal{F}_1}(X, A) + Q_{\mathcal{F}_2}(X, A)$$

as a suitable surrogate for F . Based on the update rule in (13), we consider the zero gradient condition $\nabla_X Q_{\mathcal{F}}(X, \tilde{X}) = 0$ and compute

$$\begin{aligned}\frac{\partial Q_F}{\partial X_{nt}}(X, \tilde{X}) &= \frac{\partial f_t}{\partial X_{nt}}(\tilde{X}_{\bullet,t}) + \left(\Lambda(\tilde{X}_{\bullet,t})(X_{\bullet,t} - \tilde{X}_{\bullet,t}) \right)_n + \frac{\alpha}{2} \frac{\partial}{\partial X_{nt}} \|X - BC\|_F^2 \\ &= \left(A_t^\top A_t \tilde{X}_{\bullet,t} \right)_n - (A_t^\top Y_{\bullet,t})_n + \mu_X \tilde{X}_{nt} + \lambda_X \\ &\quad + \frac{\left((A_t^\top A_t + \mu_X I_{N \times N}) \tilde{X}_{\bullet,t} \right)_n + \lambda_X}{\tilde{X}_{nt}} (X_{nt} - \tilde{X}_{nt}) + \alpha (X_{nt} - (BC)_{nt}) \\ &= - (A_t^\top Y_{\bullet,t})_n + X_{nt} \frac{\left(A_t^\top A_t \tilde{X}_{\bullet,t} \right)_n + \mu_X \tilde{X}_{nt} + \lambda_X}{\tilde{X}_{nt}} + \alpha (X_{nt} - (BC)_{nt}) \\ &= 0.\end{aligned}$$

Rearranging the equation leads to

$$X_{nt} = \frac{(A_t^\top Y_{\bullet,t})_n + \alpha (BC)_{nt}}{\frac{\left(A_t^\top A_t \tilde{X}_{\bullet,t} \right)_n + \mu_X \tilde{X}_{nt} + \lambda_X}{\tilde{X}_{nt}} + \alpha}.$$

We therefore have

$$X_{\bullet,t} = \tilde{X}_{\bullet,t} \circ \frac{A_t^\top Y_{\bullet,t} + \alpha BC_{\bullet,t}}{A_t^\top A_t \tilde{X}_{\bullet,t} + (\mu_X + \alpha) \tilde{X}_{\bullet,t} + \lambda_X \mathbf{1}_{N \times 1}},$$

which yields the multiplicative update rule

$$X_{\bullet,t} \leftarrow X_{\bullet,t} \circ \frac{A_t^\top Y_{\bullet,t} + \alpha BC_{\bullet,t}}{A_t^\top A_t X_{\bullet,t} + (\mu_X + \alpha) X_{\bullet,t} + \lambda_X \mathbf{1}_{N \times 1}}$$

based on (13). Note that the correct choice of the matrix Λ together with the κ_i is crucial to ensure the multiplicative structure of the algorithm.

B.1.2 Algorithm for \mathbf{B}

The minimisation with respect to \mathbf{B} reduces the cost function in BC-X to

$$\mathcal{F}(\mathbf{B}) := \underbrace{\frac{\alpha}{2}\|\mathbf{BC} - \mathbf{X}\|_F^2 + \frac{\mu_B}{2}\|\mathbf{B}\|_F^2 + \lambda_B\|\mathbf{B}\|_1}_{=: \mathcal{F}_1(\mathbf{B})} + \underbrace{\frac{\tau}{2}\text{TV}(\mathbf{B})}_{=: \mathcal{F}_2(\mathbf{B})} \quad (18)$$

and involves the TV regularisation on \mathbf{B} of the NMF model. Analogously to the previous section, we use the separability of \mathcal{F}_1 and write

$$\mathcal{F}_1(\mathbf{B}) = \sum_{n=1}^N \left[\frac{\alpha}{2}\|X_{n,\bullet} - B_{n,\bullet}C\|_F^2 + \frac{\mu_B}{2}\|B_{n,\bullet}\|_2^2 + \lambda_B\|B_{n,\bullet}\|_1 \right] =: \sum_{n=1}^N f_n(B_{n,\bullet}).$$

By computing the gradients

$$\begin{aligned} \nabla f_n(B_{n,\bullet}) &= \alpha(B_{n,\bullet}C - X_{n,\bullet})C^\top + \mu_B B_{n,\bullet} + \lambda_B \mathbf{1}_{1 \times K} \\ \nabla^2 f_n(B_{n,\bullet}) &= \alpha CC^\top + \mu_B I_{K \times K} \end{aligned}$$

and choosing $\kappa_k = \lambda_B$ in (16), we define analogously the surrogates Q_{f_n} , which leads to the convex surrogate

$$Q_{\mathcal{F}_1}(B, \tilde{B}) := \sum_{n=1}^N Q_{f_n}(B_{n,\bullet}, \tilde{B}_{n,\bullet})$$

for \mathcal{F}_1 . The derivation of a suitable surrogate for the TV regularisation term \mathcal{F}_2 is based on an approach different from the quadratic upper bound principle and shall not be discussed in detail. We just state the result and refer the reader for details to [33, 12, 16]. A convex and separable surrogate function for \mathcal{F}_2 is given by

$$Q_{\mathcal{F}_2}(B, \tilde{B}) = \frac{\tau}{2} \sum_{k=1}^K \sum_{n=1}^N \left[P(\tilde{B})_{nk}(B_{nk} - Z(\tilde{B})_{nk})^2 \right] + G(\tilde{B}), \quad (19)$$

with the matrices $P(\tilde{B}), Z(\tilde{B}) \in \mathbb{R}_{\geq 0}^{N \times K}$ defined in (9) and (10) and a function G depending only on the matrix \tilde{B} . Hence, we finally end up with $Q_{\mathcal{F}}(B, \tilde{B}) := Q_{\mathcal{F}_1}(B, \tilde{B}) + Q_{\mathcal{F}_2}(B, \tilde{B})$ as a suitable surrogate for \mathcal{F} .

Similar to the computations in the previous paragraph, the zero gradient condition yields then

$$\frac{\partial Q_{\mathcal{F}}}{\partial B_{nk}}(B, \tilde{B}) = -\alpha(XC^\top)_{nk} + B_{nk} \frac{\alpha(\tilde{B}CC^\top)_{nk} + \mu_B \tilde{B}_{nk} + \lambda_B}{\tilde{B}_{nk}} + \tau P(\tilde{B})_{nk}(B_{nk} - Z(\tilde{B})_{nk}) = 0$$

and therefore

$$B_{nk} = \tilde{B}_{nk} \cdot \frac{\alpha(XC^\top)_{nk} + \tau P(\tilde{B})_{nk}Z(\tilde{B})_{nk}}{\alpha(\tilde{B}CC^\top)_{nk} + \mu_B \tilde{B}_{nk} + \lambda_B + \tau P(\tilde{B})_{nk}\tilde{B}_{nk}}.$$

Hence, we have the update rule

$$B \leftarrow B \circ \frac{\alpha X C^\top + \tau P(B) \circ Z(B)}{\alpha B C C^\top + \mu_B B + \lambda_B \mathbf{1}_{N \times K} + \tau P(B) \circ B}.$$

B.1.3 Algorithm for C

The optimisation with respect to the matrix C can be tackled analogously with the quadratic upper bound principle and will not be described in detail. In this case, the cost function can be reduced to well-known regularised NMF problems [11], which leads to the update rule

$$C \leftarrow C \circ \frac{\alpha B^\top X}{\alpha B^\top B C + \mu_C C + \lambda_C \mathbf{1}_{K \times T}}.$$

B.2 Model BC

In this section, we discuss the computation of the optimisation algorithms for the NMF model BC.

B.2.1 Algorithm for B

In this case, the cost function reduces to

$$\mathcal{F}(B) := \underbrace{\sum_{t=1}^T \frac{1}{2} \|A_t(BC)_{\bullet,t} - Y_{\bullet,t}\|_2^2 + \frac{\mu_B}{2} \|B\|_F^2 + \lambda_B \|B\|_1}_{=: \mathcal{F}_1(B)} + \underbrace{\frac{\tau}{2} \text{TV}(B)}_{=: \mathcal{F}_2(B)}.$$

Analogously to the previous cases, we analyze the functions \mathcal{F}_1 and \mathcal{F}_2 separately. The difference is here, that \mathcal{F}_1 is not separable with respect to the rows of B due to the discrepancy term and therefore, it is necessary to compute the gradient and the Hessian of the whole function \mathcal{F}_1 . Hence, the gradient $\nabla \mathcal{F}_1(B)$ is an $N \times K$ matrix and the Hessian $\nabla^2 \mathcal{F}_1(B)$ a fourth-order tensor, which are given by their entries

$$\begin{aligned} \nabla \mathcal{F}_1(B)_{nk} &= \sum_{t=1}^T C_{kt} (A_t^\top A_t(BC)_{\bullet,t})_n - \sum_{t=1}^T C_{kt} (A_t^\top Y_{\bullet,t})_n + \mu_B B_{nk} + \lambda_B, \\ \nabla^2 \mathcal{F}_1(B)_{(n,k),(\tilde{n},\tilde{k})} &= \sum_{t=1}^T C_{\tilde{k}t} C_{kt} (A_t^\top A_t)_{n\tilde{n}} + \mu_B \delta_{(n,k),(\tilde{n},\tilde{k})}, \end{aligned}$$

where $\delta_{(n,k),(\tilde{n},\tilde{k})} = 1$ if and only if $(n, k) = (\tilde{n}, \tilde{k})$. The natural expansion of the quadratic upper bound principle given in Lemma A.1 is the ansatz function

$$\begin{aligned} Q_{\mathcal{F}_1}(B, \tilde{B}) &:= \mathcal{F}_1(\tilde{B}) + \langle B - \tilde{B}, \nabla \mathcal{F}_1(\tilde{B}) \rangle_F \\ &\quad + \frac{1}{2} \sum_{(n,k)} \sum_{(\tilde{n},\tilde{k})} (B - \tilde{B})_{nk} \Lambda(\tilde{B})_{(n,k),(\tilde{n},\tilde{k})} (B - \tilde{B})_{\tilde{n}\tilde{k}} \end{aligned}$$

with the fourth order tensor

$$\Lambda(\tilde{B})_{(n,k),(\tilde{n},\tilde{k})} := \begin{cases} \frac{\sum_{(i,j)} \nabla^2 \mathcal{F}_1(\tilde{B})_{(n,k),(i,j)} \tilde{B}_{ij} + \lambda_B}{\tilde{B}_{nk}} & \text{for } (n,k) = (\tilde{n},\tilde{k}), \\ 0 & \text{for } (n,k) \neq (\tilde{n},\tilde{k}), \end{cases}$$

where $\langle \cdot, \cdot \rangle_F$ denotes the Frobenius inner product.

Taking the same surrogate $Q_{\mathcal{F}_2}$ for the TV penalty term as in (19), we end up with the surrogate function

$$Q_{\mathcal{F}}(B, \tilde{B}) := Q_{\mathcal{F}_1}(B, \tilde{B}) + Q_{\mathcal{F}_2}(B, \tilde{B})$$

for F . Its partial derivative with respect to B_{nk} is given by

$$\begin{aligned} \frac{\partial Q_{\mathcal{F}}}{\partial B_{nk}}(B) &= - \sum_{t=1}^T C_{kt} (A_t^\top Y_{\bullet,t})_n + B_{nk} \frac{\sum_{t=1}^T C_{kt} \left(A_t^\top A_t (\tilde{B}C)_{\bullet,t} \right)_n + \mu_B \tilde{B}_{nk} + \lambda_B}{\tilde{B}_{nk}} \\ &\quad + \tau P(\tilde{B})_{nk} (B_{nk} - Z(\tilde{B})_{nk}). \end{aligned}$$

The zero-gradient condition gives then the equation

$$B_{nk} = A_{nk} \left(\frac{\sum_{t=1}^T C_{kt} (A_t^\top Y_{\bullet,t})_n + \tau P(\tilde{B})_{nk} Z(\tilde{B})_{nk}}{\sum_{t=1}^T C_{kt} \left(A_t^\top A_t (\tilde{B}C)_{\bullet,t} \right)_n + \mu_B \tilde{B}_{nk} + \lambda_B + \tilde{B}_{nk} \tau P(\tilde{B})_{nk}} \right),$$

which can be extended to the whole matrix B . Therefore, based on (13), we have the update rule

$$B \leftarrow B \circ \left(\frac{\sum_{t=1}^T A_t^\top Y_{\bullet,t} (C^\top)_{t,\bullet} + \tau P(B) \circ Z(B)}{\sum_{t=1}^T A_t^\top A_t (BC)_{\bullet,t} \cdot (C^\top)_{t,\bullet} + \mu_B B + \lambda_B \mathbf{1}_{N \times K} + \tau B \circ P(B)} \right).$$

B.2.2 Algorithm for C

In this case, the cost function is separable with respect to the columns of C , such that

$$\mathcal{F}(C) := \sum_{t=1}^T \frac{1}{2} \|A_t BC_{\bullet,t} - Y_{\bullet,t}\|_2^2 + \frac{\mu_C}{2} \|C_{\bullet,t}\|_2^2 + \lambda_C \|C_{\bullet,t}\|_1 =: \sum_{t=1}^T f_t(C_{\bullet,t}).$$

Hence, we can split the minimisation into the columns of C to use the standard quadratic upper bound principle without considering higher order tensors. We compute

$$\begin{aligned} \nabla f_t(C_{\bullet,t}) &= B^\top A_t^\top A_t (BC)_{\bullet,t} - B^\top A_t^\top Y_{\bullet,t} + \mu_C C_{\bullet,t} + \lambda_C \mathbf{1}_{K \times 1}, \\ \nabla^2 f_t(C_{\bullet,t}) &= B^\top A_t^\top A_t B + \mu_C I_{K \times K}. \end{aligned}$$

By choosing $\kappa_k = \lambda_C$ for all k in (16), we define $Q_{f_t}(C_{\bullet,t}, \tilde{C}_{\bullet,t})$ as a surrogate function for f_t according to Lemma A.1. The update rule in (15) gives then

$$C_{\bullet,t} = \tilde{C}_{\bullet,t} - \Lambda^{-1}(\tilde{C}_{\bullet,t}) \nabla f_t(\tilde{C}_{\bullet,t}),$$

which leads to

$$C_{\bullet,t} \leftarrow C_{\bullet,t} \circ \frac{B^\top A_t^\top Y_{\bullet,t}}{B^\top A_t^\top A_t B C_{\bullet,t} + \mu_C C_{\bullet,t} + \lambda_C \mathbf{1}_{K \times 1}}.$$

C Parameter Choice

| Parameter | BC | | BC-X | | gradTV | |
|----------------------|----------|----------|----------|----------|-------------------|---------------------|
| | 1% noise | 3% noise | 1% noise | 3% noise | 1% noise | 3% noise |
| α | — | — | 70 | 70 | — | — |
| μ_C | 0.1 | 0.1 | 0.1 | 0.1 | — | — |
| τ | 10 | 50 | 6 | 20 | — | — |
| ρ_{grad} | — | — | — | — | $1 \cdot 10^{-3}$ | $8 \cdot 10^{-4}$ |
| ρ_{thr} | — | — | — | — | $7 \cdot 10^{-4}$ | $1 \cdot 10^{-3}$ |
| ρ_{TV} | — | — | — | — | $1 \cdot 10^{-2}$ | $2.5 \cdot 10^{-2}$ |
| $\tilde{\mu}_C$ | — | — | — | — | 0.1 | 0.1 |

Table 3: Parameter choice of the considered methods for the dynamic Shepp–Logan Phantom for 1% and 3% Gaussian noise.

| Parameter | BC | | BC-X | | gradTV | |
|----------------------|------------------|------------------|----------------|----------------|-------------------|---------------------|
| | 1% noise | 3% noise | 1% noise | 3% noise | 1% noise | 3% noise |
| α | — | — | $3 \cdot 10^2$ | $3 \cdot 10^2$ | — | — |
| μ_C | 1 | 1 | 1 | 1 | — | — |
| τ | $1.3 \cdot 10^2$ | $4.3 \cdot 10^2$ | 90 | $3 \cdot 10^2$ | — | — |
| ρ_{grad} | — | — | — | — | $2 \cdot 10^{-4}$ | $8 \cdot 10^{-5}$ |
| ρ_{thr} | — | — | — | — | $2 \cdot 10^{-4}$ | $2.5 \cdot 10^{-4}$ |
| ρ_{TV} | — | — | — | — | $2 \cdot 10^{-2}$ | $4 \cdot 10^{-2}$ |
| $\tilde{\mu}_C$ | — | — | — | — | 0.1 | 0.1 |

Table 4: Parameter choice of the considered methods for the vessel phantom for 1% and 3% Gaussian noise.