

Efficient Algorithms for Convolutional Sparse Representations

Brendt Wohlberg

EDICS: SMR-REP

Abstract—When applying sparse representation techniques to images, the standard approach is to independently compute the representations for a set of overlapping image patches. This method performs very well in a variety of applications, but results in a representation that is multi-valued and not optimised with respect to the entire image. A recent development is the *convolutional sparse representation*, in which a sparse representation of an entire image is computed by replacing the linear combination of a set of dictionary vectors by the sum of a set of convolutions with dictionary filters. The resulting representation is both single-valued and jointly optimised over the entire image. Despite its advantages, however, this form of sparse representation has received little attention for image processing applications. This is, presumably, at least partially due to the considerable computational cost of initial algorithms for this form of sparse coding and dictionary learning. This paper presents new, efficient algorithms that substantially improve on the performance of other recent methods, contributing to the development of this type of representation into a practical tool for image processing inverse problems.

Index Terms—Sparse Representation, Sparse Coding, Dictionary Learning, Convolutional Sparse Representation, ADMM

I. INTRODUCTION

Sparse representation [1] is a widely used technique for a very broad range of signal and image processing applications. Given a signal s and a *dictionary* matrix D , *sparse coding* is the inverse problem of finding the sparse representation \mathbf{x} with only a few non-zero entries such that $D\mathbf{x} \approx s$. Most sparse coding algorithms optimize a functional consisting of a data fidelity term and a sparsity inducing penalty

$$\arg \min_{\mathbf{x}} \frac{1}{2} \|D\mathbf{x} - s\|_2^2 + \lambda R(\mathbf{x}), \quad (1)$$

or constrained forms such as

$$\arg \min_{\mathbf{x}} R(\mathbf{x}) \text{ such that } \|D\mathbf{x} - s\|_2 \leq \epsilon, \quad (2)$$

or

$$\arg \min_{\mathbf{x}} \|D\mathbf{x} - s\|_2^2 \text{ such that } R(\mathbf{x}) \leq \tau, \quad (3)$$

where $R(\cdot)$ denotes a sparsity inducing function such as the ℓ^1 norm or the ℓ^0 “norm”. The two leading families of sparse coding methods are a wide variety of convex optimization algorithms (e.g. Alternating Direction Method of Multipliers (ADMM) [2], [3]) for solving Eq. (1) when $R(\mathbf{x}) = \|\mathbf{x}\|_1$ and a family of greedy algorithms (e.g. Matching Pursuit (MP) [4] and Orthogonal Matching Pursuit (OMP) [5]) for approximate solution of Eq. (2) or Eq. (3) when $R(\mathbf{x}) = \|\mathbf{x}\|_0$.

If the dictionary D is analytically defined and corresponds to a linear operator with a fast transform (e.g. the Discrete Wavelet Transform), a representation for an entire signal or image can easily be computed. More recently, however, it has been realised that improved performance can be obtained by learning the dictionary from a set of training data relevant to a specific problem [6]; this inverse problem is known as *dictionary learning*. In this case computing a sparse representation for an entire signal is not feasible, the usual approach being to apply the decomposition independently to a set of overlapping blocks covering the signal, as illustrated in Fig. 1. This approach is easy to implement, but results in a representation that is multi-valued and suboptimal over the signal as a whole, and often necessitates somewhat ad hoc handling of the overlap between blocks¹.

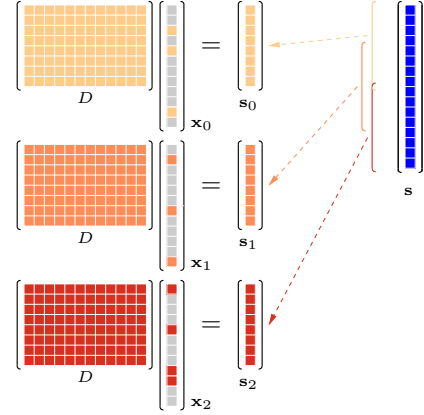


Fig. 1. Independent sparse representation of overlapping blocks.

An optimal representation for a given block structure would involve solving a single optimization problem for a block matrix dictionary \tilde{D} with blocks consisting of appropriately shifted versions of the single-block dictionary D , as illustrated in Fig. 2. If the block structure is modified so that the step between blocks is a single sample, as in Fig. 3, the representation on the resulting block dictionary \tilde{D} is both optimal for the entire signal and shift-invariant. Furthermore, in this case the linear representation on the block dictionary \tilde{D} can be expressed as a sum of convolutions of columns \mathbf{d}_m of D with a set of spatial coefficient maps corresponding to the subsets of \mathbf{x} associated with different shifts of each column of D . That is, by a simple change of indexing

¹In some applications, such as denoising of Gaussian white noise, the averaging of independent block estimates appears to provide a performance advantage [7].

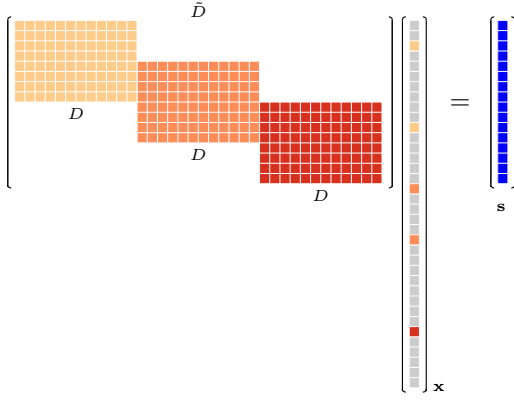


Fig. 2. Jointly optimized sparse representation of overlapping blocks.

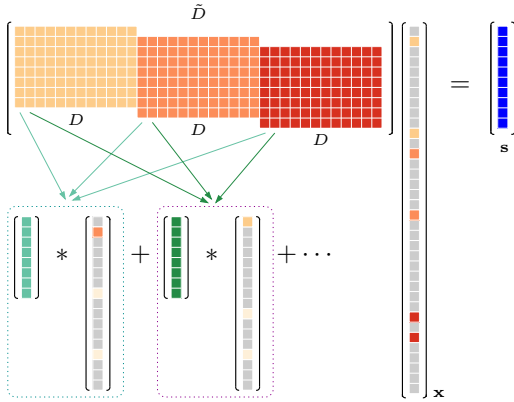


Fig. 3. Equivalence between jointly optimal sparse representation and convolutional sparse representation.

of coefficients, the representation $\tilde{D}\mathbf{x} \approx \mathbf{s}$ is equivalent to $\sum_m \mathbf{d}_m * \mathbf{x}_m \approx \mathbf{s}$. This form of sparse representation is referred to here as a *convolutional* sparse representation. While convolutional forms can be constructed for each of Eq. (1)–(3), the present paper focuses on the convolutional form of Eq. (1) with $R(\mathbf{x}) = \|\mathbf{x}\|_1$, i.e. the Basis Pursuit DeNoising (BPDN) [8] problem

$$\arg \min_{\mathbf{x}} \frac{1}{2} \|\mathbf{D}\mathbf{x} - \mathbf{s}\|_2^2 + \lambda \|\mathbf{x}\|_1. \quad (4)$$

Convolutional BPDN is defined as

$$\arg \min_{\{\mathbf{x}_m\}} \frac{1}{2} \left\| \sum_m \mathbf{d}_m * \mathbf{x}_m - \mathbf{s} \right\|_2^2 + \lambda \sum_m \|\mathbf{x}_m\|_1, \quad (5)$$

where $\{\mathbf{d}_m\}$ is a set of M dictionary *filters*, $*$ denotes convolution, and $\{\mathbf{x}_m\}$ is a set of coefficient maps. (For notational simplicity \mathbf{s} and each of the $\{\mathbf{x}_m\}$ are considered to be N dimensional vectors, where N is the number of pixels in an image.)

The present paper substantially extends an earlier work [9] to describe fast algorithms for both convolutional sparse coding and dictionary learning, including algorithm enhancements, penalty parameter selection methods, and performance comparisons with other methods. While these techniques are also applicable to 1-d signals, discussion here focuses on single-band images (the extension to color and other multi-band images is mathematically straightforward, at the price

of some additional notational clutter). In all computational experiments the 8 bit greyscale test images are divided by 255 so that pixel values are within the interval $[0,1]$, and convolutional sparse coding/dictionary learning algorithms are applied to test/training images after a highpass filtering preprocessing step, which is a common approach when using convolutional sparse representations (e.g. [10], [11]).

II. LITERATURE SURVEY

There are two independent streams of research, with almost no cross-citation, related to the convolutional form of sparse representations.

A. Convolutional Sparse Representations

The more recent of these streams is a product of the deep learning community; the convolutional form of the representation is primary, inspired by the desire to modify convolutional neural networks [12] to provide a *generative model*. The idea appears to have been proposed independently and almost simultaneously by Zeiler et al. [13] and Yang et al. [14], although the latter, while also originating within the machine learning community, does not seem to have influenced any later work in this direction. Subsequent work within this stream has also tended to consider the convolutional form of the representation as an integral part of a multi-layer neural network with application to image classification tasks [10], [11], [15], [16], although some more recent work has specifically addressed the concept from a more traditional sparse representations perspective [17]–[20].

Most of these works [10], [11], [13], [15], [19], [20] have posed the sparse coding and dictionary learning problems in the form of convolutional BPDN, the exceptions being probabilistic/Bayesian models [16], [18] and convolutional extensions [17] of MP and K-SVD [21].

The most frequent application area within this stream has been image classification [10], [11], [13], [14], [16] and detection [10], [15] tasks, but segmentation in biological imagery has also been considered [18], and image denoising has received very brief attention [13].

B. Translation-Invariant Sparse Representations

The other stream of research has origins in the neuroscience and signal processing communities, the unifying concept being the construction of *translation-invariant* representations. The resulting representation is equivalent to the convolutional representation, but the motivation is translation-invariance and the convolutional form is derived rather than given. The earliest work to construct representations directly equivalent to Eq. (5) is that of Lewicki and Sejnowski [22]. Work within this stream [22]–[43] has concentrated primarily on properties of the representation and signal processing applications. A related idea that is not equivalent to convolutional representations is the construction of representations that are invariant to continuous translations via Taylor series or other interpolation methods [44], [45]. Some of these works have considered constructing representations that are invariant to a more general

set of transformations [24], [39]–[41], most commonly being image translations and rotations.

A range of sparse coding and dictionary learning forms have been considered within this stream, including penalized form Eq. (1) with a log function based sparsity inducing regularization terms [23] or with ℓ^1 regularization [24], [34], [38]–[40]; probabilistic/Bayesian models [22], [25]–[29], [36], which generally correspond to various forms of Eq. (1) with either log [25], [28] or ℓ^1 [36] regularization; and constrained forms Eq. (2) [37], [41] or Eq. (3) [31], [33], [35], [42], [43].

Many of these works focus on properties of the representation rather than considering any specific application. The majority of applications that have been addressed within this stream have been to signal decomposition [26], [28], [29], [46], classification [34], and blind source separation [27], [29], [36] problems in audio [26]–[29], [34], [36], [46] or medical diagnostic signals [27], [31], [43]. Image denoising [37], image inpainting [38], [42], and video decomposition [25] have received some attention.

C. Sparse Coding Algorithms

Solutions for the convolutional constrained forms of sparse coding have employed convolutional extensions of MP [17], [18], [30], [35], [37] or variants of OMP [31], [33], [41]–[43]. In most of these cases the primary focus of the work is not on sparse coding algorithm development, and only [17], [18], [35], [37], [41] discuss efficient convolutional extensions of these methods in any detail.

A wide variety of algorithms have been proposed for convolutional BPDN and related unconstrained forms. A few different approaches combine heuristic dictionary subset selection and optimization within that subset; these include a greedy subset search via fast convolution followed by a conjugate gradient method to compute the optimal coefficients on the selected subset [22], a similar greedy subset selection followed by an IRLS algorithm to solve the BPDN problem on the subset [28], [29], [46], and the feature-sign search algorithm [34]. Other methods include block coordinate relaxation [38], [39], coordinate descent [10], an alternating minimization method (not ADMM) [13], ISTA [11] or equivalent to ISTA but not identified as such [36], [40], FISTA [15], [19], and ADMM [20].

Most of these methods solve the problem in the spatial domain, but [20], [36], [37], [42], [43] exploit the convolution theorem to solve at least part of the problem in the Discrete Fourier Transform (DFT) domain.

D. Dictionary Learning Algorithms

A number of different authors have considered the development of convolutional extensions of the K-SVD dictionary update [17], [18], [32], [35], [37], [42]. A block coordinate relaxation approach [43] shares the property of the K-SVD of updating a single dictionary element at a time. The MoTIF algorithm [30] is rather unusual in its strategy of iteratively growing the dictionary with a procedure somewhat reminiscent of Gram-Schmidt. These methods tend to be paired with

greedy MP/OMP sparse coding algorithms for the coefficient update step.

The most widely used dictionary updates are convolutional forms of the gradient descent [13], [25]–[28], [36], [38]–[40], [46] and MOD [11], [31], [33] or variants thereof [20], [34]. Other methods that have been considered include a quasi-Newton method [19], stochastic gradient descent [10], [15], stochastic Levenberg-Marquardt [41], and block coordinate relaxation [43].

Most dictionary learning methods are composed of alternation at the outer level between a sparse coding step, itself consisting of a number of iterations of a standard sparse coding algorithm, and a dictionary update step. A different approach is a more tightly coupled algorithm alternating between coefficient and dictionary update steps [20], [36], the more recent of which [20] is derived with an Augmented Lagrangian framework.

As in the case of sparse coding, most methods operate in the spatial domain, but [20], [34], [36], [42], [43] perform the dictionary update in the DFT domain.

III. ADMM ALGORITHM FOR CONVOLUTIONAL BPDN

Given the recent success of ADMM for the standard BPDN problem, it is a natural choice to consider for Convolutional BPDN. The general outline of the method derived here is very similar to the sparse coding component of the Augmented Lagrangian dictionary learning algorithm proposed by Bristow et al. [20] and described in more detail, with some errors corrected, in an independent technical report [47]. At a superficial level, the only difference is that here the ADMM algorithm is derived in the spatial domain, with one of the sub-problems being solved in the frequency domain, whereas they directly derive the algorithm using a mixture of spatial and frequency domain variables. A much more important difference, however, is the use of a far more efficient method (first introduced in [9]) for solving the linear system that represents the bulk of the computational cost of the algorithm.

The ADMM iterations for solving the optimization

$$\arg \min_{\mathbf{x}, \mathbf{y}} f(\mathbf{x}) + g(\mathbf{y}) \quad \text{such that} \quad \mathbf{Ax} + \mathbf{By} = \mathbf{c} \quad (6)$$

are, in *scaled form* [2, Sec. 3.1.1]

$$\mathbf{x}^{(j+1)} = \arg \min_{\mathbf{x}} f(\mathbf{x}) + \frac{\rho}{2} \left\| \mathbf{Ax} + \mathbf{By}^{(j)} - \mathbf{c} + \mathbf{u}^{(j)} \right\|_2^2 \quad (7)$$

$$\mathbf{y}^{(j+1)} = \arg \min_{\mathbf{y}} g(\mathbf{y}) + \frac{\rho}{2} \left\| \mathbf{Ax}^{(j+1)} + \mathbf{By} - \mathbf{c} + \mathbf{u}^{(j)} \right\|_2^2 \quad (8)$$

$$\mathbf{u}^{(j+1)} = \mathbf{u}^{(j)} + \mathbf{Ax}^{(j+1)} + \mathbf{By}^{(j+1)} - \mathbf{c} . \quad (9)$$

Rewriting Eq. (5) in a suitable form by introducing an auxiliary variable with a constraint, we have

$$\arg \min_{\{\mathbf{x}_m\}, \{\mathbf{y}_m\}} \frac{1}{2} \left\| \sum_m \mathbf{d}_m * \mathbf{x}_m - \mathbf{s} \right\|_2^2 + \lambda \sum_m \|\mathbf{y}_m\|_1 \quad \text{s.t.} \quad \mathbf{x}_m - \mathbf{y}_m = 0, \quad (10)$$

so the ADMM iterations for our problem are

$$\{\mathbf{x}_m\}^{(j+1)} = \arg \min_{\{\mathbf{x}_m\}} \frac{1}{2} \left\| \sum_m \mathbf{d}_m * \mathbf{x}_m - \mathbf{s} \right\|_2^2 + \frac{\rho}{2} \sum_m \left\| \mathbf{x}_m - \mathbf{y}_m^{(j)} + \mathbf{u}_m^{(j)} \right\|_2^2 \quad (11)$$

$$\{\mathbf{y}_m\}^{(j+1)} = \arg \min_{\{\mathbf{y}_m\}} \lambda \sum_m \|\mathbf{y}_m\|_1 + \frac{\rho}{2} \sum_m \left\| \mathbf{x}_m^{(j+1)} - \mathbf{y}_m + \mathbf{u}_m^{(j)} \right\|_2^2 \quad (12)$$

$$\mathbf{u}_m^{(j+1)} = \mathbf{u}_m^{(j)} + \mathbf{x}_m^{(j+1)} - \mathbf{y}_m^{(j+1)}. \quad (13)$$

The stopping criteria discussed in [2, Sec. 3.3] provide a more effective means of ensuring convergence to suitable accuracy than simply specifying a maximum allowed number of iterations.

Sub-Problem Eq. (12) is solved via shrinkage/soft thresholding as

$$\mathbf{y}_m^{(j+1)} = \mathcal{S}_{\lambda/\rho} \left(\mathbf{x}_m^{(j+1)} + \mathbf{u}_m^{(j)} \right), \quad (14)$$

where

$$\mathcal{S}_\gamma(\mathbf{u}) = \text{sign}(\mathbf{u}) \odot \max(0, |\mathbf{u}| - \gamma), \quad (15)$$

with $\text{sign}(\cdot)$ and $|\cdot|$ of a vector considered to be applied element-wise, and \odot denoting element-wise multiplication. The computational cost of this sub-problem is $\mathcal{O}(MN)$. The only computationally expensive step is solving Eq. (11), which is of the form

$$\arg \min_{\{\mathbf{x}_m\}} \frac{1}{2} \left\| \sum_m \mathbf{d}_m * \mathbf{x}_m - \mathbf{s} \right\|_2^2 + \frac{\rho}{2} \sum_m \|\mathbf{x}_m - \mathbf{z}_m\|_2^2. \quad (16)$$

A. DFT Domain Formulation

An obvious approach is to attempt to exploit the FFT for efficient implementation of the convolution via the DFT convolution theorem. Define linear operators D_m such that $D_m \mathbf{x}_m = \mathbf{d}_m * \mathbf{x}_m$, and denote the variables D_m , \mathbf{x}_m , \mathbf{s} , and \mathbf{z}_m in the DFT domain by \hat{D}_m , $\hat{\mathbf{x}}_m$, $\hat{\mathbf{s}}$, and $\hat{\mathbf{z}}_m$ respectively. It is easy to show via the DFT convolution theorem that Eq. (16) is equivalent to

$$\arg \min_{\{\hat{\mathbf{x}}_m\}} \frac{1}{2} \left\| \sum_m \hat{D}_m \hat{\mathbf{x}}_m - \hat{\mathbf{s}} \right\|_2^2 + \frac{\rho}{2} \sum_m \|\hat{\mathbf{x}}_m - \hat{\mathbf{z}}_m\|_2^2, \quad (17)$$

with the $\{\mathbf{x}_m\}$ minimizing Eq. (16) being given by the inverse DFT of the $\{\hat{\mathbf{x}}_m\}$ minimizing Eq. (17). Defining

$$\hat{D} = \begin{pmatrix} \hat{D}_0 & \hat{D}_1 & \dots \end{pmatrix} \quad \hat{\mathbf{x}} = \begin{pmatrix} \hat{\mathbf{x}}_0 \\ \hat{\mathbf{x}}_1 \\ \vdots \end{pmatrix} \quad \hat{\mathbf{z}} = \begin{pmatrix} \hat{\mathbf{z}}_0 \\ \hat{\mathbf{z}}_1 \\ \vdots \end{pmatrix}, \quad (18)$$

this problem can be expressed as

$$\arg \min_{\hat{\mathbf{x}}} \frac{1}{2} \|\hat{D}\hat{\mathbf{x}} - \hat{\mathbf{s}}\|_2^2 + \frac{\rho}{2} \|\hat{\mathbf{x}} - \hat{\mathbf{z}}\|_2^2, \quad (19)$$

the solution being given by the linear system

$$(\hat{D}^H \hat{D} + \rho I) \hat{\mathbf{x}} = \hat{D}^H \hat{\mathbf{s}} + \rho \hat{\mathbf{z}}. \quad (20)$$

B. Independent Linear Systems

Matrix \hat{D} has a block structure consisting of M concatenated $N \times N$ diagonal matrices, where M is the number of filters and N is the number of samples in \mathbf{s} . $\hat{D}^H \hat{D}$ is an $MN \times MN$ matrix, but due to the diagonal block (not block diagonal) structure of \hat{D} , a row of \hat{D}^H with its non-zero element at column n will only have a non-zero product with a column of \hat{D} that has its non-zero element at row n . As a result, there is no interaction between elements of \hat{D} corresponding to different frequency indices n , so that

(as pointed out in [20]) one need only solve N independent $M \times M$ linear systems to solve Eq. (20). The cost of the FFTs and these linear solves dominate the computational cost of the algorithm. Bristow et al. [20] do not specify how they solve these linear systems, but since they rate the computational cost of solving them as $\mathcal{O}(M^3)$, one must conclude that they apply a direct method such as Gaussian elimination. This can be a very effective solution when it is possible to precompute and cache an LU or Cholesky decomposition of the system matrix, but in this case that is, in general, impossible due to the $\mathcal{O}(M^2 N)$ storage requirements for these decompositions.

A careful analysis of the unique structure of this problem, however, reveals that there is an alternative, and vastly more effective solution. The independent systems have a left hand side consisting of a diagonal matrix plus a rank-one component, which can be solved very efficiently by re-arranging the independent systems in Eq. (20) into the form of Eq. (45), as described in Appendix A, and then solving using Eq. (49), derived by applying the Sherman-Morrison formula as described in Appendix B. The only vector operations in Eq. (49) are scalar multiplication, subtraction, and inner products, so that this method is $\mathcal{O}(M)$, instead of $\mathcal{O}(M^3)$ as in [20]. The cost of solving such a system at all M spatial indices is $\mathcal{O}(MN)$, and the cost of the FFTs is $\mathcal{O}(MN \log N)$; it is the cost of the FFTs that dominates the computational complexity, whereas in [20] the cost of the solutions of the linear systems in the DFT domain dominate the cost of the FFTs.

This approach can be implemented in an interpreted language such as Matlab in a form that avoids explicit loops over the N spatial indices by passing data for all N spatial indices as a single array to the relevant linear-algebraic routines (commonly referred to as *vectorization* in Matlab terminology). Further computational improvement is possible, at the cost of additional memory requirements, by precomputing components of Eq. (49), e.g. $\mathbf{a}_n^H / (\rho + \mathbf{a}_n^H \mathbf{a}_n)$.

C. Performance Comparison: Linear Solvers for ADMM

The execution times and solution accuracies of four different methods of solving the linear system are compared in Fig. 4. Computations were performed in both single and double precision arithmetic to allow evaluation of the corresponding execution time and solution accuracy trade-off. Dictionaries of different sizes were constructed by selecting subsets of a $12 \times 12 \times 512$ dictionary² learned using standard (non-convolutional) methods, and \mathbf{s} was the 512×512 greyscale Lena image. After an initial 20 iterations of the ADMM algorithm (to avoid any unique characteristics of the linear system directly after initialisation of the iterations), the four different methods were used to solve Eq. (20) for the following iteration, recording the time required by each. Reconstruction error was computed as the relative residual error³ with which solution $\hat{\mathbf{x}}$ satisfied the linear equation Eq. (20).

The following observations can be made from these results:

²Included in the demonstration software distributed by the authors of [48].

³The relative residual error of \mathbf{x} in the linear equation $\mathbf{A}\mathbf{x} = \mathbf{b}$ is $\|\mathbf{A}\mathbf{x} - \mathbf{b}\|_2 / \|\mathbf{b}\|_2$.

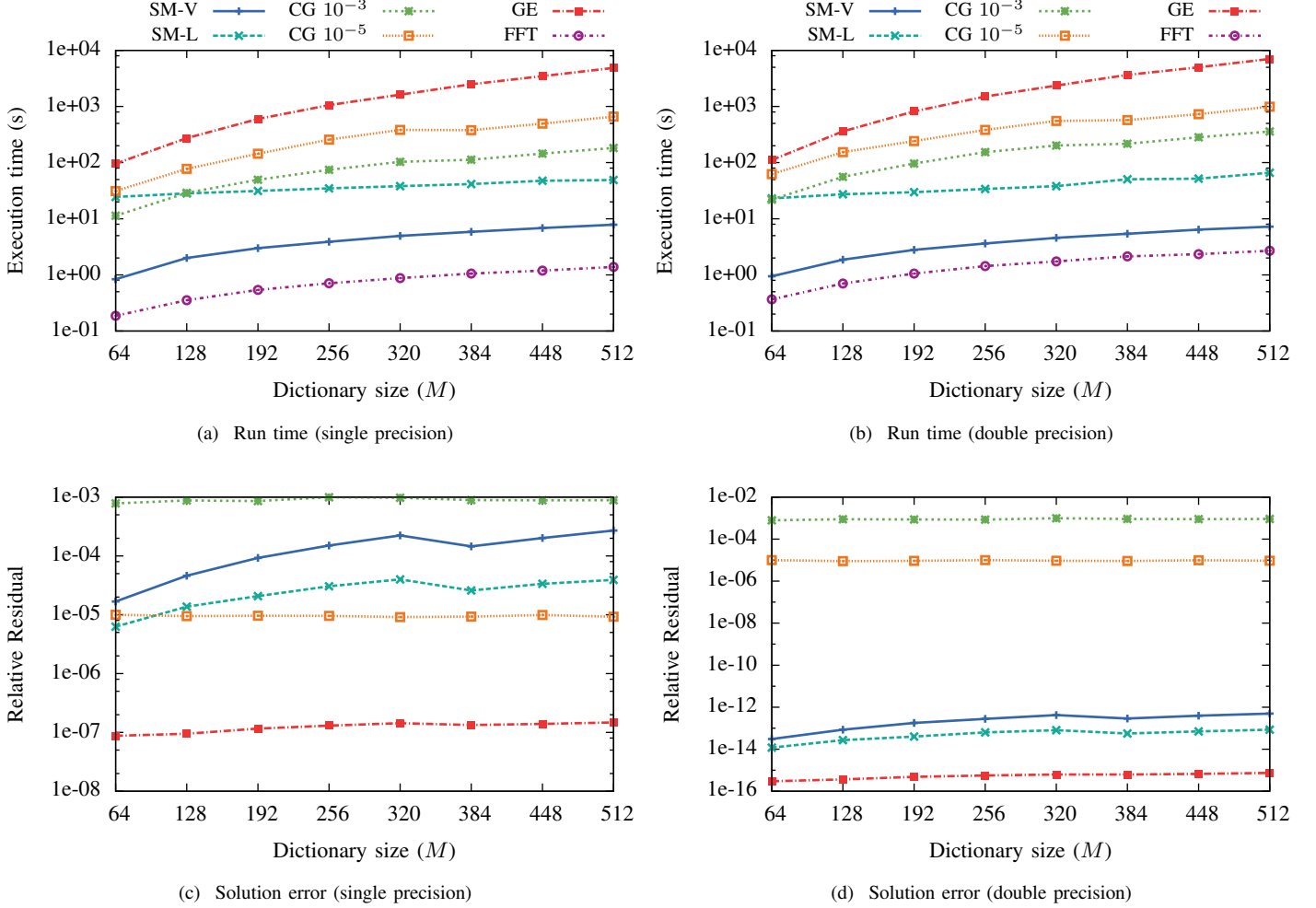


Fig. 4. A comparison of execution times and solution errors for solution of a linear system for the 512×512 grayscale Lena image with dictionaries consisting of 12×12 filters, with $\lambda = 0.01$. Solution methods are Gaussian elimination (GE), Conjugate Gradient (CG) with relative residual tolerances of 10^{-3} and 10^{-5} , and Sherman-Morrison implemented with a loop over each spatial location (SM-L) or jointly over all spatial locations (SM-V). The combined run time for a forward and inverse FFT is provided as a reference.

- The variation in execution times between single and double precision is far smaller than the variation between the different solution methods. SM-V run time has very small variation with precision, while CG can be almost twice as slow in double precision as in single, and GE run time dependence on precision increases with M .
- GE gives solutions accurate to floating point precision, but is the slowest of all the methods, being vastly slower than SM-V, and having much worse scaling with M .
- CG provides a better accuracy/time trade-off than GE, being substantially faster, while still providing solutions of acceptable accuracy.
- SM-V is by far the fastest method, giving acceptable solution accuracy in single precision, and very high accuracy in double precision. The negligible variation in run time with precision suggests that use of a double precision solution might be preferred, but empirical evidence (see below) indicates that, in practice, the single precision solution accuracy is sufficient, and gives approximately the same rate of convergence for the ADMM algorithm

as double precision solution.

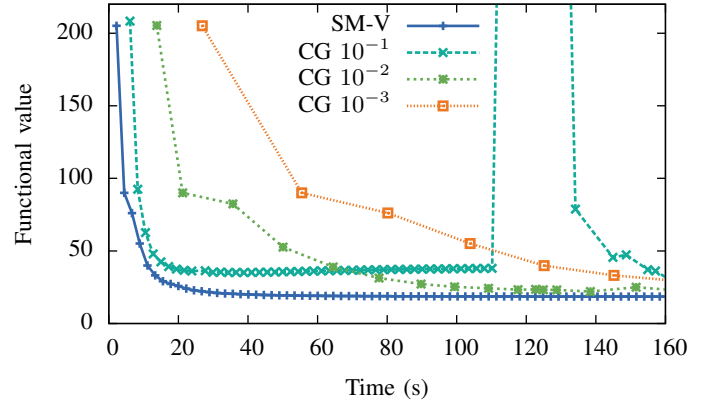


Fig. 5. A comparison of functional value evolution with time for the ADMM algorithm with Eq. (20) solved in single precision using Sherman-Morrison and CG with three different relative residual tolerances. The test image was the 512×512 grayscale Lena image with a learned $8 \times 8 \times 64$ dictionary and $\lambda = 0.01$.

Since these methods all exhibit a different trade-off between accuracy and speed, it is reasonable to ask what the effect is on algorithm performance. It is not worth expending a large amount of time to compute a very accurate solution since ADMM algorithms only require approximate solutions of each step [2, Sec. 3.4.4], but convergence will be slow or unreliable if the solutions are sufficiently inaccurate. The results presented in Fig. 5 show that SM-V gives a faster reduction in functional value⁴ with time than CG with any choice of tolerance⁵, even for the smallest M considered in Fig. 4. In this case CG with a relative residual tolerance of 10^{-1} is not much slower than SM-V, but the algorithm does not converge; larger tolerance values give more reliable decrease in functional value with each iteration (note, though, that on a longer timescale than displayed in Fig. 5, CG with tolerances 10^{-2} and 10^{-3} also begin to exhibit some oscillations), but are very much slower. If plotted on the same graph, the SM-V result for double precision is indistinguishable from the single precision SM-V result shown.

D. ADMM Parameter Selection

The selection of a suitable penalty parameter ρ is critical to obtaining a good convergence rate. In this section two different strategies for selecting ρ_0 are compared: the automatic method described in [2, Sec. 3.4.1], and the increasing parameter scheme advocated in [20]. Two variants of the automatic method are considered: one with a period of 1 iteration (i.e. parameter update allowed at every iteration) and another with a period of 10 iterations (i.e. update test applied only at every 10 iterations). An additional variation for all methods is comparing the results with and without the relaxation method described in [2, Sec. 3.4.3] (relaxation parameter $\alpha = 1.8$, which was found to give the best results, is used whenever this method is applied). All experiments used the 512×512 greyscale Lena test image with a learned $8 \times 8 \times 64$ dictionary.

The first set of experiments, reported in Fig. 6, compare the different methods across different values of λ using relative functional values, computed by dividing functional values for each λ by the smallest functional value attained for that λ , by any method, at the maximum number of iterations (500). These results show that:

- The multiplicative update scheme without over-relaxation is not at all competitive with other methods.
- Including over-relaxation with the multiplicative update scheme substantially improves performance, but it remains inferior to the automatic scheme for all but the smallest values of λ when performance is compared at 50 iterations, and is uniformly inferior when comparisons are made at larger iteration counts.

⁴In an ADMM algorithm such as this, different functional values will be obtained depending on whether it is computed using the primary ($\{\mathbf{x}_m\}$) or auxiliary ($\{\mathbf{y}_m\}$) variable. The auxiliary $\{\mathbf{y}_m\}$ is used for the experiments reported here since (i) $\{\mathbf{x}_m\}$ is merely an intermediate result and not part of the algorithm state [2, pg. 14], and (ii) it has been found empirically to give a curve that decays faster and is less subject to oscillations. Corresponding plots for primal and dual residuals for the ADMM algorithm are provided in the Supplementary Material.

⁵Smaller CG tolerances than in Fig. 4 are used in Fig. 5 since it is already clear from Fig. 4 that tolerances smaller than 10^{-3} are not competitive.

- Over-relaxation substantially improves performance for the auto-update scheme, and when over-relaxation is applied, the auto-update period makes very little difference to performance (period 1 is slightly better than period 10 for smaller λ , and slightly worse for larger λ).

The second set of experiments, reported in Fig. 7, evaluate the sensitivities of the different methods to the correct choice of ρ_0 by comparing the decrease in functional value with iteration number for the optimum choice of ρ_0 , as well as the optimum value multiplied by 0.1 and 10. These experiments show that:

- Auto-update with period 1 and 10 give similar results when using the optimal ρ_0 , but period 1 is less sensitive to a poor choice of ρ_0 .
- The multiplicative update scheme with ρ_0 chosen for the smallest functional value at 50 iterations converges substantially slower than the auto-update scheme. Initial convergence is much faster if ρ_0 is chosen to be 10 times larger, but then the functional value stagnates at a level substantially above the minimum. (If ρ_0 is instead chosen for best performance at 100 iterations then the stagnation effect is reduced, but at the expense of even slower decrease in the functional value, i.e. the initial plateau becomes larger.) The performance of the multiplicative update scheme with respect to the auto-update method, both in terms of decay rate and stagnation effect, appears to improve slightly for smaller values of λ .

By inspection of 2D arrays of relative functional values for different λ and ρ_0 values⁶, it was determined that $\rho_0 = 100\lambda + 0.5$ is a simple heuristic that provides good performance for the experimental conditions – type of data, data scaling, and preprocessing, etc. – considered here.

E. Performance Comparison: FISTA

While there is evidence that ADMM is significantly faster than FISTA for at least some problems [49], an empirical comparison with the recently proposed FISTA algorithm for Convolutional BPDN [19] is warranted. The experiments reported in Fig. 8 compare the ADMM algorithm proposed here with the FISTA algorithm of [19], as well as with a variant of that algorithm that computes the gradient in the frequency domain instead of in the spatial domain as in the original version. The comparisons were performed for a set of different dictionaries (learned on a separate set of images) and λ values, using the 512×512 greyscale Lena test image. The ADMM results were generated using over-relaxation with $\alpha = 1.8$, with $\rho_0 = 100\lambda + 0.5$, and with ρ autoscaling with period 1 iterations. FISTA parameters were selected by a grid search over a range of values around the fixed $L_0 = 1$ and $\eta = 5$ values of [19]; the best values of L_0 were 0.05, 0.5, or 1, and the best values of η were 2 or 5. (Note that ADMM parameters were all automatically set, or set via an effective heuristic, whereas it is not clear how to select the best FISTA parameters without a search over the parameter space.)

The results indicate that FISTA with spatial domain computation of the gradient is substantially slower than ADMM.

⁶See Supplementary Material.

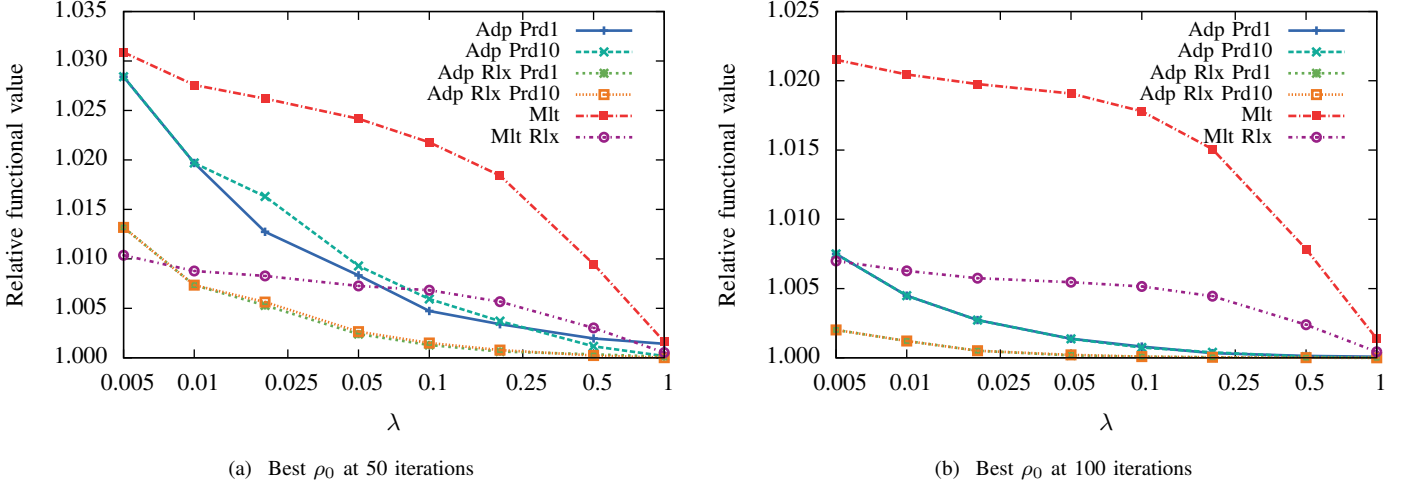


Fig. 6. Relative functional values for different ρ update strategies with ρ_0 selected to minimize the functional value at 50 and 100 iterations. The update strategies are denoted by “Adp” (adaptive) and “Mlt” (fixed multiplicative update), with “Rlx” indicating the use of over-relaxation with the parameter $\alpha = 1.8$, and “Prd1” and “Prd10” indicating respectively that the adaptive update is applied at every iteration or at every 10 iterations.

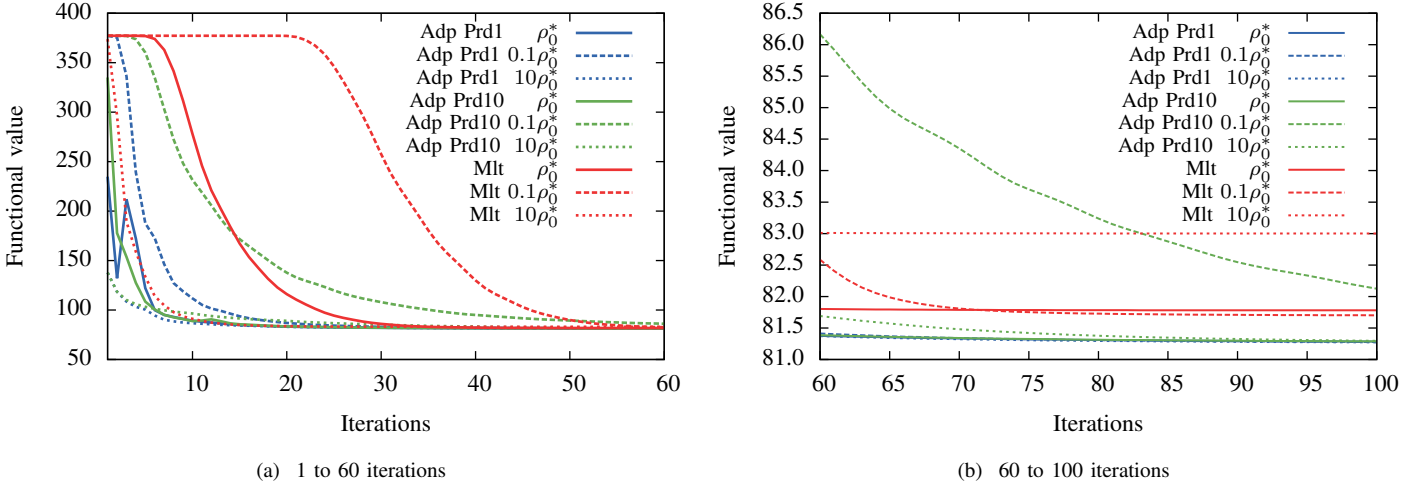


Fig. 7. Functional value behaviour for $\lambda = 0.05$ and different ρ update strategies with ρ_0 selected to minimize the functional value at 50 iterations. Both “Adp” (adaptive) and “Mlt” (fixed multiplicative update) strategies use over-relaxation parameter $\alpha = 1.8$. The adaptive update is shown for both updates at every iteration (“Prd1”) and every 10 iterations (“Prd10”). The evolution of the functional value for each method is shown for the optimum choice of ρ_0 for that method at 50 iterations, denoted by ρ_0^* as well as for ρ_0 set to 0.1 and 10 times ρ_0^* .

FISTA with frequency domain computation of the gradient is competitive with ADMM for larger values of λ (of less interest to image reconstruction problems) and smaller (M) dictionaries, but is also substantially slower than ADMM for larger λ and M .

IV. DICTIONARY LEARNING

In contrast to standard sparse representations, it is possible to define the convolutional dictionary learning problem in the Single Measurement Vector (SMV) case, but for full generality it should be constructed in the Multiple Measurement Vector (MMV) context. Including the usual constraint on the norm of dictionary elements, the problem can be expressed as

$$\arg \min_{\{\mathbf{d}_m\}, \{\mathbf{x}_{k,m}\}} \frac{1}{2} \sum_k \left\| \sum_m \mathbf{d}_m * \mathbf{x}_{k,m} - \mathbf{s}_k \right\|_2^2 + \lambda \sum_k \sum_m \|\mathbf{x}_{k,m}\|_1$$

$$\text{such that } \|\mathbf{d}_m\|_2 = 1 \quad \forall m. \quad (21)$$

The standard approach is to solve via alternating minimization with respect to $\{\mathbf{x}_{k,m}\}$ and $\{\mathbf{d}_m\}$. The minimization with respect to $\{\mathbf{x}_{k,m}\}$ involves solving the MMV extension of Convolutional BPDN, which is trivial since the problems for each k are decoupled from one another, but $\{\mathbf{d}_m\}$ is more challenging since the problems for different k are coupled.

Ignoring the constraint on the norm of \mathbf{d}_m , which is usually applied as a post-processing normalization step after the update, the minimization with respect to $\{\mathbf{d}_m\}$ can be expressed as

$$\arg \min_{\{\mathbf{d}_m\}} \frac{1}{2} \sum_k \left\| \sum_m \mathbf{d}_m * \mathbf{x}_{k,m} - \mathbf{s}_k \right\|_2^2, \quad (22)$$

which is a convolutional form of the Method of Optimal

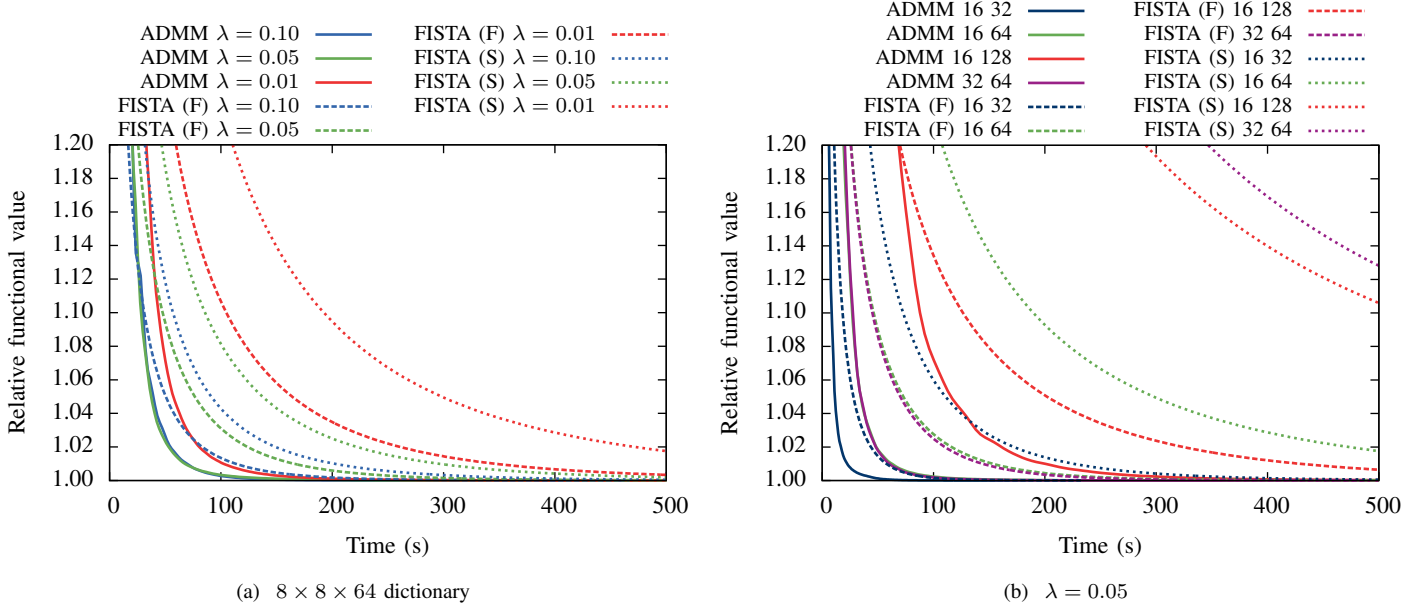


Fig. 8. Comparison of time evolution of relative functional values for ADMM, and FISTA with the gradient computed in the spatial domain “(S)” and frequency domain “(F)”. (a) provides a comparison for different values of λ with a dictionary consisting of 64 filters of size 8×8 , and (b) provides a comparison for different dictionary sizes (“method $n\ m$ ” indicates a dictionary of m filters of size $n \times n$) and $\lambda = 0.05$.

Directions (MOD) [50]. When computing the convolutions $\mathbf{d}_m * \mathbf{x}_{k,m}$ in the DFT domain, there is an implicit zero-padding of the filters \mathbf{d}_m to the size of the coefficient maps $\mathbf{x}_{k,m}$; this can be overlooked when minimizing with respect to the coefficient maps, but must be explicitly represented when minimizing with respect to the filters to ensure that the filters resulting from the optimization have an appropriately constrained support in the spatial domain. Defining zero-padding operator P , Eq. (22) can be expressed in the DFT domain as

$$\arg \min_{\{\mathbf{d}_m\}} \frac{1}{2} \sum_k \left\| \sum_m (\widehat{P\mathbf{d}_m}) \odot \hat{\mathbf{x}}_{k,m} - \hat{\mathbf{s}}_k \right\|_2^2. \quad (23)$$

Unfortunately, the spatial-domain operator P does not have a compact representation in the DFT domain, making an efficient direct DFT domain solution impossible. A variable splitting approach, however, makes it possible to solve this problem, including dictionary norm constraints, via an ADMM algorithm.

A. Constrained MOD Update

The desired filters can be obtained as $P^T \mathbf{d}_m$ after solving the constrained problem

$$\arg \min_{\{\mathbf{d}_m\}} \frac{1}{2} \sum_k \left\| \sum_m \mathbf{x}_{k,m} * \mathbf{d}_m - \mathbf{s}_k \right\|_2^2 \text{ s. t. } \mathbf{d}_m \in C_P \ \forall m, \quad (24)$$

where the \mathbf{d}_m have the same spatial support as the $\mathbf{x}_{k,m}$, and

$$C_P = \{\mathbf{x} \in \mathbb{R}^N : (I - PP^T)\mathbf{x} = 0\}. \quad (25)$$

Since we are setting up a constrained problem requiring an iterative solution, however, it is reasonable to also include the normalisation $\|\mathbf{d}_m\|_2 = 1$ or $\|\mathbf{d}_m\|_2 \leq 1$ of the dictionary

elements that is often, and suboptimally, performed as a postprocessing step after the dictionary update. Including the normalisation requirement $\|\mathbf{d}_m\|_2 = 1$, the constraint set becomes

$$C_{PN} = \{\mathbf{x} \in \mathbb{R}^N : (I - PP^T)\mathbf{x} = 0, \|\mathbf{x}\|_2 = 1\}. \quad (26)$$

Employing the indicator function⁷ $\iota_{C_{PN}}$ of the constraint set C_{PN} , the constrained problem can be written in unconstrained form [51]

$$\arg \min_{\{\mathbf{d}_m\}} \frac{1}{2} \sum_k \left\| \sum_m \mathbf{x}_{k,m} * \mathbf{d}_m - \mathbf{s}_k \right\|_2^2 + \sum_m \iota_{C_{PN}}(\mathbf{d}_m), \quad (27)$$

and rewriting with an auxiliary variable in a form suitable for ADMM gives

$$\arg \min_{\{\mathbf{d}_m\}} \frac{1}{2} \sum_k \left\| \sum_m \mathbf{x}_{k,m} * \mathbf{d}_m - \mathbf{s}_k \right\|_2^2 + \sum_m \iota_{C_{PN}}(\mathbf{g}_m) \text{ such that } \mathbf{d}_m - \mathbf{g}_m = 0 \ \forall m. \quad (28)$$

This problem can be solved via an ADMM algorithm

$$\{\mathbf{d}_m\}^{(j+1)} = \arg \min_{\{\mathbf{d}_m\}} \frac{1}{2} \sum_k \left\| \sum_m \mathbf{x}_{k,m} * \mathbf{d}_m - \mathbf{s}_k \right\|_2^2 + \frac{\sigma}{2} \sum_m \left\| \mathbf{d}_m - \mathbf{g}_m^{(j)} + \mathbf{h}_m^{(j)} \right\|_2^2 \quad (29)$$

$$\{\mathbf{g}_m\}^{(j+1)} = \arg \min_{\{\mathbf{g}_m\}} \sum_m \iota_{C_{PN}}(\mathbf{g}_m) + \frac{\sigma}{2} \sum_m \left\| \mathbf{d}_m^{(j+1)} - \mathbf{g}_m + \mathbf{h}_m^{(j)} \right\|_2^2 \quad (30)$$

⁷The indicator function of a set S is defined as

$$\iota_S(X) = \begin{cases} 0 & \text{if } X \in S \\ \infty & \text{if } X \notin S \end{cases}.$$

$$\mathbf{h}_m^{(j+1)} = \mathbf{h}_m^{(j)} + \mathbf{d}_m^{(j+1)} - \mathbf{g}_m^{(j+1)}. \quad (31)$$

The $\{\mathbf{g}_m\}$ update is of the form

$$\arg \min_{\mathbf{x}} \frac{1}{2} \|\mathbf{x} - \mathbf{y}\|_2^2 + \iota_{C_{\text{PN}}}(\mathbf{x}) = \text{prox}_{\iota_{C_{\text{PN}}}}(\mathbf{y}). \quad (32)$$

It is immediately clear from the geometry of the problem that

$$\text{prox}_{\iota_{C_{\text{PN}}}}(\mathbf{y}) = \frac{PP^T \mathbf{y}}{\|PP^T \mathbf{y}\|_2}, \quad (33)$$

or, if the normalisation $\|\mathbf{d}_m\|_2 \leq 1$ is desired instead,

$$\text{prox}_{\iota_{C_{\text{PN}}}}(\mathbf{y}) = \begin{cases} PP^T \mathbf{y} & \text{if } \|PP^T \mathbf{y}\|_2 \leq 1 \\ \frac{PP^T \mathbf{y}}{\|PP^T \mathbf{y}\|_2} & \text{if } \|PP^T \mathbf{y}\|_2 > 1 \end{cases}. \quad (34)$$

The problem for the $\{\mathbf{d}_m\}$ update is of form

$$\arg \min_{\{\mathbf{d}_m\}} \frac{1}{2} \sum_k \left\| \sum_m \mathbf{x}_{k,m} * \mathbf{d}_m - \mathbf{s}_k \right\|_2^2 + \frac{\sigma}{2} \sum_m \|\mathbf{d}_m - \mathbf{z}_m\|_2^2. \quad (35)$$

In the DFT domain, where $\hat{X}_{k,m} = \text{diag}(\hat{\mathbf{x}}_{k,m})$, this becomes

$$\arg \min_{\{\mathbf{d}_m\}} \frac{1}{2} \sum_k \left\| \sum_m \hat{X}_{k,m} \hat{\mathbf{d}}_m - \hat{\mathbf{s}}_k \right\|_2^2 + \frac{\sigma}{2} \sum_m \|\hat{\mathbf{d}}_m - \hat{\mathbf{z}}_m\|_2^2. \quad (36)$$

Defining

$$\hat{X}_k = (\hat{X}_{k,0} \quad \hat{X}_{k,1} \quad \dots) \quad \hat{\mathbf{d}} = \begin{pmatrix} \hat{\mathbf{d}}_0 \\ \hat{\mathbf{d}}_1 \\ \vdots \end{pmatrix} \quad \hat{\mathbf{z}} = \begin{pmatrix} \hat{\mathbf{z}}_0 \\ \hat{\mathbf{z}}_1 \\ \vdots \end{pmatrix} \quad (37)$$

this problem can be expressed as

$$\arg \min_{\hat{\mathbf{d}}} \frac{1}{2} \sum_k \|\hat{X}_k \hat{\mathbf{d}} - \hat{\mathbf{s}}_k\|_2^2 + \frac{\sigma}{2} \|\hat{\mathbf{d}} - \hat{\mathbf{z}}\|_2^2 \quad (38)$$

with solution

$$\left(\sum_k \hat{X}_k^H \hat{X}_k + \sigma I \right) \hat{\mathbf{d}} = \sum_k \hat{X}_k^H \hat{\mathbf{s}}_k + \sigma \hat{\mathbf{z}}. \quad (39)$$

This linear system can be solved by iterated application of the Sherman-Morrison formula, as described in Appendices C and D.

B. Interleaved $\mathbf{x}_{k,m}$ and \mathbf{d}_m Updates

Given iterative algorithms for the $\mathbf{x}_{k,m}$ and \mathbf{d}_m updates (i.e. the ADMM algorithms for Convolutional BPDN and Constrained MOD respectively), the immediate question is how these should be combined into a full dictionary learning algorithm. The standard approach is to alternate between sparse coding and dictionary update, solving the sparse coding problem with reasonable accuracy before moving on to the dictionary update. Since each sub-problem is an iterative algorithm, this would entail performing multiple iterations of each sub-problem before switching to the other. A reasonable strategy is to maintain the auxiliary variables for each sub-problem across the top-level iterations since performing a cold start at each switch of sub-problems would entail substantially reduced efficiency. The necessity of retaining these variables, in turn, suggests interleaving the ADMM iterations for each sub-problem into a single set of iterations, rather than combining the $\mathbf{x}_{k,m}$ and \mathbf{d}_m updates in a way that treats each as

a single functional unit⁸. A single iteration of the resulting algorithm consists of updates Eq. (11), (12), (13), (29), (30), and (31) in sequence.

The most obvious way of combining these updates is to transfer the primary variables across to the other update steps, i.e. \mathbf{d}_m represent the dictionary in the sparse coding steps, and $\mathbf{x}_{k,m}$ represent the sparse code in the dictionary update steps. Such a combination turns out to be quite unstable in practice, with convergence being very sensitive to suitable choices of the ρ and σ parameters for each update. A far more stable algorithm is obtained if the updates are interleaved on their auxiliary variables, i.e. \mathbf{g}_m represent the dictionary in the sparse coding steps, and $\mathbf{y}_{k,m}$ represent the sparse code in the dictionary update steps. It is worth noting that such a combination is not derivable from single Augmented Lagrangian functional, as is the otherwise-similar dictionary learning approach of Bristow et al. [20]. Additional differences in derivation and algorithm are:

- they construct the Augmented Lagrangian using a mixture of spatial and frequency domain variables, while the derivation presented here presents the problem in the spatial domain, switching into the frequency domain where appropriate for efficient solution of relevant sub-problems, and
- they derive the ADMM algorithm in unscaled rather than scaled form [2, Sec. 3.1.1].

These choices appear to lead to a slightly more complicated path to deriving solutions to at least one of the sub-problems (see section “Sub-Problem s” in [20], describing the sub-problem corresponding to the $\{\mathbf{g}_m\}$ update here).

C. Performance Comparison: Linear Solvers for ADMM

The execution times and solution accuracies of three different methods of solving the linear system are compared in Figs. 9 and 10. Dictionaries, D , of different sizes, M , were constructed by initialising with Gaussian random dictionaries of the appropriate sizes, and the different training image sets, S , of size, K , were selected from a set of 256×256 training images derived from the MIRFlickr dataset [53]. After an initial 20 iterations of the ADMM algorithm, the different methods were used to solve Eq. (39) for the following iteration, recording the time required by each. Reconstruction error was computed as the relative residual error with which solution $\hat{\mathbf{d}}$ satisfies the linear equation Eq. (39).

The following observations can be made from these results:

- The CG 10^{-1} method is by far the fastest, but the solution accuracy is inadequate (see below).
- The CG 10^{-5} method has very similar solution accuracy to SM, but is substantially slower for small K values.
- In the sparse coding problem, the SM solution greatly outperforms the alternatives, but that is not always the case here, with other methods, including GE, becoming competitive for larger K values. This should not be

⁸This general strategy, of interleaving the algorithms for sparse code and dictionary updates, has previously been proposed, but with substantially different algorithms for each update [52].

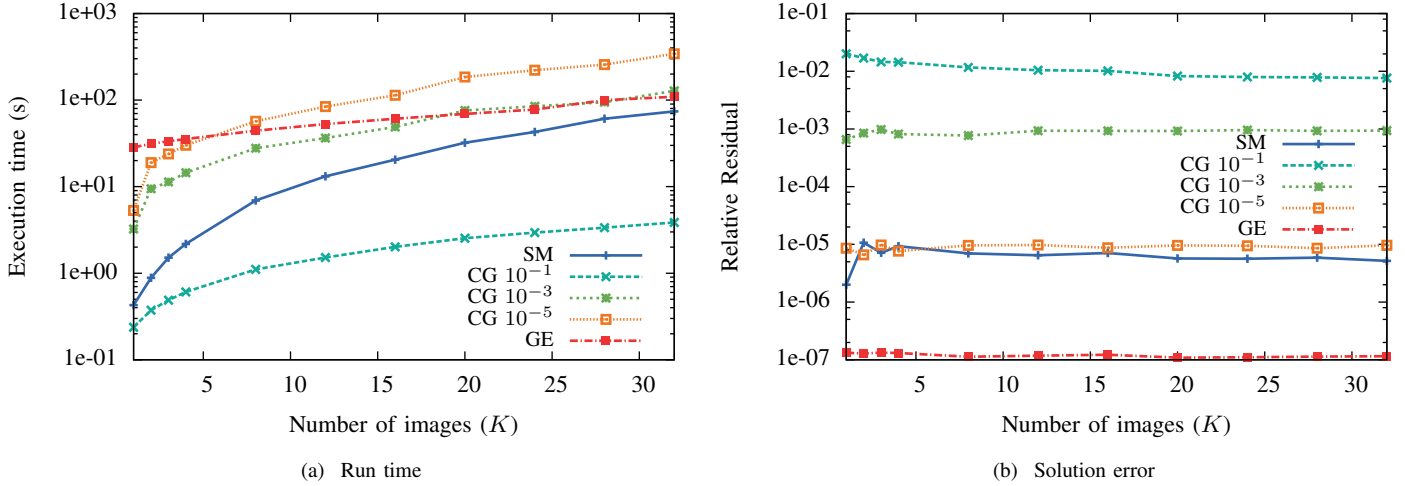


Fig. 9. A comparison of execution times and solution errors for single-precision solution of a linear system corresponding to an $8 \times 8 \times 64$ dictionary with $\lambda = 0.1$. Solution methods are Gaussian elimination (GE), Conjugate Gradient (CG), and Sherman-Morrison (SM).

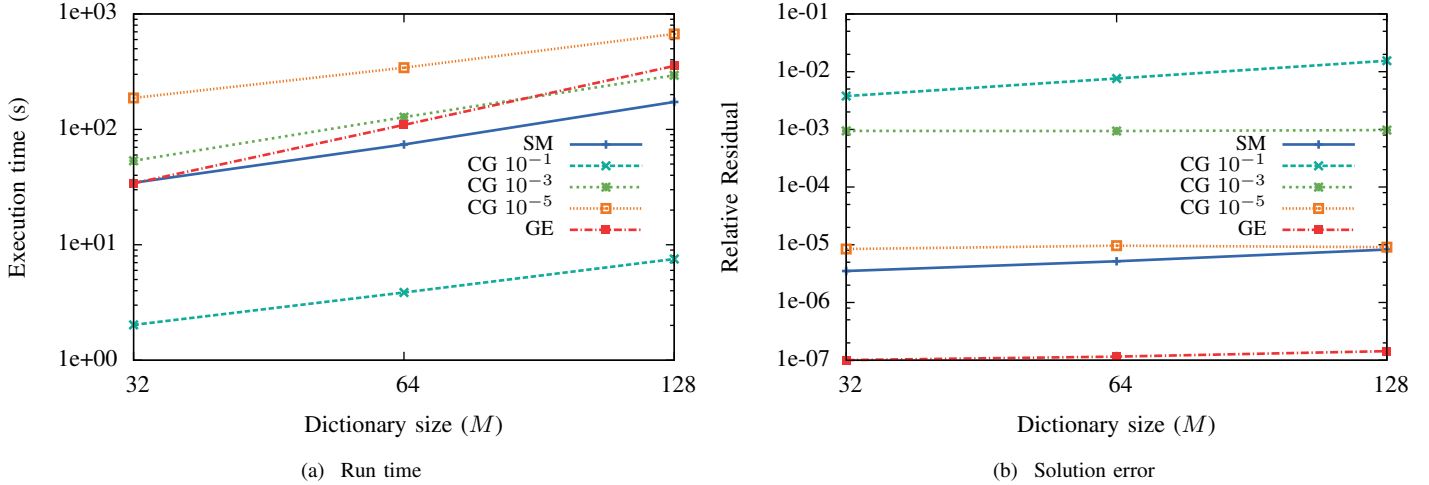


Fig. 10. A comparison of execution times for single-precision solution of linear system corresponding to an $8 \times 8 \times M$ dictionary with $K = 32$ and $\lambda = 0.1$. Solution methods are Gaussian elimination (GE), Conjugate Gradient (CG), and Sherman-Morrison (SM).

surprising since GE has $\mathcal{O}(KM^3N)$ cost, while SM has $\mathcal{O}(K^2MN)$ cost (see Alg. 1 in Appendix D).

The effects of the different trade-offs between time and accuracy represented by the different solution methods can only be evaluated by comparing performance within a number of iterations of the full dictionary learning algorithm. Such a comparison was performed using the same dictionary size, λ value, and training data (with $K = 32$) as for the experiments above, the results being presented in Fig. 11.

The double-precision SM result is not shown since it would be indistinguishable from the single-precision version on this graph; for this experiment the maximum fractional difference between the two (occurring within the first few iterations) is 3.3×10^{-5} , and the median is 3.8×10^{-6} . Thus double-precision SM gives effectively the same solution, but at the cost of doubling the memory requirements (and is approximately 30% slower than single-precision). Since CG with a sufficiently large tolerance gives approximately the same

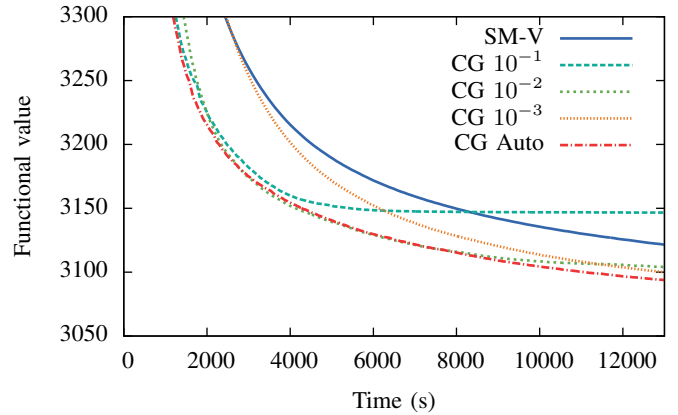


Fig. 11. A comparison of functional value evolution with time for the ADMM algorithm with Eq. (39) solved in single precision using Sherman-Morrison and CG with three different relative residual tolerances as well as an automatic tolerance selection method. The dictionary was $8 \times 8 \times 64$ and $\lambda = 0.1$.

solution accuracy independent of precision, at an even greater computation time penalty, double-precision CG is also not a useful option.

In addition to three different CG tolerances, these experiments also include an automatic CG tolerance method that starts with a tolerance of 10^{-1} , at each iteration setting it to $1/20$ times the primal residual for the \mathbf{d}_m update if that quantity is smaller than the current tolerance. In this experiment that auto-tolerance method is quite competitive, giving the overall best performance⁹ in terms of functional value reduction with time.

D. Multi-scale Dictionaries

The vast majority of dictionary learning research has focused on the learning of single-scale dictionaries (i.e. all dictionary elements are of the same size) within a patch-based context. Given the obvious benefits of a multi-scale representation (see e.g. [54, Sec. 3] for a discussion), it is not surprising, however, that the construction of multi-scale dictionaries has received at least some attention. The approaches that have been considered follow a few basic schemes:

Analytically defined The simplest approach is to analytically define a multi-scale family of basis functions, without any learning or adaptation of the dictionary from data. The design of a fast sparse coding algorithm using a multi-scale Gaussian chirp dictionary was described in [55].

Learned wavelet filters A wavelet basis can be optimised for sparse representation of images by a learning procedure applied to the associated filter bank [56], [57].

Transform domain dictionaries Standard dictionary learning applied in the wavelet transform domain [58], [59] or in a Laplacian pyramid [60] provides the dictionary with a multi-scale structure in the spatial domain.

Quadtree spatial decomposition A set of dictionaries can be learned to give an image representation at the multiple patch sizes occurring within a quadtree decomposition [54], [61].

These methods all have structural constraints (either being imposed by the properties of the transform domain within which the sparse representation is computed, or from the quadtree spatial structure imposed on the dictionary), ultimately resulting from the difficulty of applying a multi-scale dictionary in a natural way within a patch-based framework.

In the convolutional sparse representation framework, in contrast, there is absolutely no reason why the dictionary filters should be of the same size, and multi-scale dictionaries can be defined in a very natural way, without any structural constraints on their form (an example is displayed in Fig. 12). Learning of such dictionaries is no more difficult than learning a single-scale dictionary, simply by replacing P by P_m in Eq. (23), (25), (26), (33), and (34). It is surprising, therefore, that the learning and use of multi-scale convolutional dictionaries has not previously been considered in imaging applications, and has only received very limited attention in

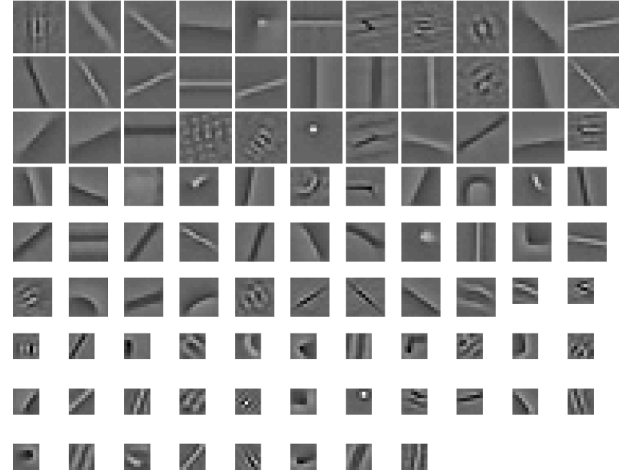


Fig. 12. Example multiscale dictionary with 8×8 , 12×12 , and 16×16 filters learned with $\lambda = 0.1$ from a set of 32 512×512 training images derived from a selection of Flickr Creative Commons images.

a signal processing context (the possibility of using filters of different lengths is pointed out in [41], but not discussed in any detail).

The potential advantage of a multi-scale dictionary is demonstrated by some simple experiments. A set of five different dictionaries, all with 72 filters, was learned using the same set of 16 images, and then used to solve the Convolutional BPDN problem using a separate set of 16 image as a test set. A corresponding experiment was also performed using a set of four dictionaries, all with 96 filters, trained on the same set of 32 images and tested on the same set of 23 images. The results in Table I show that a multi-scale dictionary can provide a lower cost representation than a single-scale dictionary with the same number of filters of the same size as one of the filter sizes in the multi-scale dictionary. A far more thorough investigation than this preliminary demonstration is warranted, but is beyond the scope of the present paper.

V. CONCLUSION

At a high level, the derivation of ADMM algorithms for convolutional sparse coding and the dictionary update for dictionary learning is straightforward, but correct choice of methods for solving the linear sub-problems are critical to the design of efficient algorithms, with vastly inferior computational performance being possible if these methods are not properly chosen.

The results presented here show that the proposed Sherman-Morrison solution of the main linear system is clearly the best choice for the Convolutional BPDN sparse coding problem, giving much better asymptotic performance, $\mathcal{O}(M)$ instead of $\mathcal{O}(M^3)$, than the previously proposed direct method [20], and is also much faster in practical application. Furthermore, over-relaxation methods are shown to improve performance, and effective heuristic techniques for selection of the penalty parameter ρ are presented. The resulting ADMM algorithm is shown to be much faster than a previously proposed FISTA approach [19], and faster to varying degrees depending on regularization parameter λ when the gradient required by FISTA is computed in the DFT domain.

⁹In contrast, the best reduction with respect to number of iterations is provided by SM and CG 10^{-3} (see Supplementary Material).

Dict.	$8 \times 8 \times 72$	$12 \times 12 \times 72$	$8 \times 8 \times 24, 12 \times 12 \times 48$	$16 \times 16 \times 72$	$8 \times 8 \times 24, 16 \times 16 \times 48$
Func.	4.96×10^2	4.81×10^2	4.79×10^2	4.79×10^2	4.76×10^2
Dict.	$8 \times 8 \times 96$	$12 \times 12 \times 96$	$16 \times 16 \times 96$	$8 \times 8 \times 16, 12 \times 12 \times 32, 16 \times 16 \times 48$	
Func.	9.02×10^2	8.86×10^2	8.88×10^2	8.79×10^2	

TABLE I

A COMPARISON OF FUNCTIONAL VALUES FOR THE CONVOLUTIONAL BPDN PROBLEM COMPUTED WITH $\lambda = 0.01$ FOR THE SAME SET OF IMAGES AND WITH DIFFERENT DICTIONARIES. THE COMPARISON IN THE UPPER ROWS WAS COMPUTED ON THE SAME SET OF 16 IMAGES, AND THE COMPARISON IN THE LOWER ROWS WAS COMPUTED ON THE SAME SET OF 23 IMAGES.

In the case of Convolutional BPDN dictionary learning, an iterated Sherman-Morrison solution of the main linear system in the dictionary update stage is the best choice for a small number of training images K , but other methods become competitive for larger K . Overall, a CG solution with a moderate accuracy tolerance, or with an adaptive tolerance appears to be the best choice when K is large. It is also demonstrated that the proposed dictionary learning algorithm can be used to construct convolutional multi-scale dictionaries that do not have any of the structural constraints or complications inherent in previous approaches to constructing multi-scale dictionaries within the standard patch-based framework.

In the interests of reproducible research, software implementations of the main algorithms proposed here are made publicly available [62].

ACKNOWLEDGMENT

The author is grateful to R. Chalasani for kindly providing a copy of his implementation of the method described in [19].

APPENDIX A

DIAGONAL BLOCK LINEAR SYSTEMS

Given sets of vectors $\mathbf{a}_m \in \mathbb{C}^N$, $\mathbf{b}_m \in \mathbb{C}^N$, and $\mathbf{c}_m \in \mathbb{C}^N$, and unknown vectors $\mathbf{x}_m \in \mathbb{C}^N$, define $A_m = \text{diag}(\mathbf{a}_m)$, $B_m = \text{diag}(\mathbf{b}_m)$, $A = \begin{pmatrix} A_0 & A_1 & \dots & A_{M-1} \end{pmatrix}$, and

$$B = \text{diag} \begin{pmatrix} \mathbf{b}_0 \\ \mathbf{b}_1 \\ \vdots \\ \mathbf{b}_{M-1} \end{pmatrix} \quad \mathbf{x} = \begin{pmatrix} \mathbf{x}_0 \\ \mathbf{x}_1 \\ \vdots \\ \mathbf{x}_{M-1} \end{pmatrix} \quad \mathbf{c} = \begin{pmatrix} \mathbf{c}_0 \\ \mathbf{c}_1 \\ \vdots \\ \mathbf{c}_{M-1} \end{pmatrix}. \quad (40)$$

The A_m are $N \times N$ matrices and A contains M blocks so that A is $N \times MN$. The goal is to solve the linear system $(A^H A + B)\mathbf{x} = \mathbf{c}$, which can be expanded as

$$\begin{pmatrix} A_0^H A_0 + B_0 & A_0^H A_1 & A_0^H A_2 & \dots \\ A_1^H A_0 & A_1^H A_1 + B_1 & A_1^H A_2 & \dots \\ A_2^H A_0 & A_2^H A_1 & A_2^H A_2 + B_2 & \dots \\ \vdots & \vdots & \vdots & \ddots \end{pmatrix} \begin{pmatrix} \mathbf{x}_0 \\ \mathbf{x}_1 \\ \mathbf{x}_2 \\ \vdots \end{pmatrix} = \begin{pmatrix} \mathbf{c}_0 \\ \mathbf{c}_1 \\ \mathbf{c}_2 \\ \vdots \end{pmatrix}. \quad (41)$$

Since the A_m are diagonal, so are the blocks $A_m^H A_{m'}$, and it is not difficult to confirm that each row of $A^H A$ has non-zero entries corresponding to a single index in each of the \mathbf{x}_m . Denoting entry n of vector \mathbf{x}_m by $\mathbf{x}_m(n)$, the rows of Eq. (41) can be written as a set of equations

$$\sum_m (\mathbf{a}_{m'}(n)^* \mathbf{a}_m(n) + \mathbf{b}_{m'}(n)) \mathbf{x}_m(n) = \mathbf{c}_{m'}(n), \quad (42)$$

which can be simplified by swapping the vector and entry indexing by defining

$$\begin{aligned} \tilde{\mathbf{a}}_n(m) &= \mathbf{a}_m(n)^* & \tilde{\mathbf{b}}_n(m) &= \mathbf{b}_m(n) \\ \tilde{\mathbf{x}}_n(m) &= \mathbf{x}_m(n) & \tilde{\mathbf{c}}_n(m) &= \mathbf{c}_m(n), \end{aligned} \quad (43)$$

and re-ordering the equations to give

$$\sum_m (\tilde{\mathbf{a}}_n(m') \tilde{\mathbf{a}}_n(m)^* + \tilde{\mathbf{b}}_n(m')) \tilde{\mathbf{x}}_n(m) = \tilde{\mathbf{c}}_n(m'), \quad (44)$$

which can now be written using vector products as

$$(\tilde{\mathbf{a}}_n \tilde{\mathbf{a}}_n^H + \text{diag}(\tilde{\mathbf{b}}_n)) \tilde{\mathbf{x}}_n = \tilde{\mathbf{c}}_n. \quad (45)$$

The $MN \times MN$ system $(A^H A + B)\mathbf{x} = \mathbf{c}$ has been replaced by N independent linear systems of size $M \times M$, each of which consists of a rank one component plus a diagonal component. Systems of this form can be efficiently solved by application of the Sherman-Morrison formula, as shown in Appendix B.

APPENDIX B

SHERMAN-MORRISON SOLUTION

Consider the solution of linear systems of the form

$$(J + \mathbf{a}\mathbf{a}^H) \mathbf{x} = \mathbf{b} \quad (46)$$

where J is a diagonal matrix. Applying the Sherman-Morrison formula [63]

$$(A + \mathbf{u}\mathbf{v}^H)^{-1} = A^{-1} - \frac{A^{-1}\mathbf{u}\mathbf{v}^H A^{-1}}{1 + \mathbf{u}^H A^{-1}\mathbf{v}} \quad (47)$$

to derive $(J + \mathbf{a}\mathbf{a}^H)^{-1}$, and re-arranging to avoid matrix-vector products and vector outer products (noting that $\mathbf{a}^H J^{-1} \mathbf{b}$ is a scalar so that $\mathbf{a}\mathbf{a}^H J^{-1} \mathbf{b} = (\mathbf{a}^H J^{-1} \mathbf{b})\mathbf{a}$), gives

$$\mathbf{x} = J^{-1} \left(\mathbf{b} - \frac{\mathbf{a}^H J^{-1} \mathbf{b}}{1 + \mathbf{a}^H J^{-1} \mathbf{a}} \mathbf{a} \right). \quad (48)$$

If $J = \rho I$, i.e. a scaled identity matrix, then

$$\mathbf{x} = \rho^{-1} \left(\mathbf{b} - \frac{\mathbf{a}^H \mathbf{b}}{\rho + \mathbf{a}^H \mathbf{a}} \mathbf{a} \right). \quad (49)$$

APPENDIX C

MULTIPLE DIAGONAL BLOCK LINEAR SYSTEMS

Consider a generalization of the situation in Appendix A such that vectors $\mathbf{a}_{k,m} \in \mathbb{C}^N$ take an additional index, with $A_{k,m} = \text{diag}(\mathbf{a}_{k,m})$ and $A_k = \begin{pmatrix} A_{k,0} & A_{k,1} & \dots & A_{k,M-1} \end{pmatrix}$, with other variables as before, and the goal now being to solve the linear system $(\sum_k A_k^H A_k + B)\mathbf{x} = \mathbf{c}$, which can be expanded as

$$\begin{pmatrix} A_{0,0}^H A_{0,0} + A_{0,0}^H A_{1,0} + \dots + B_0 & A_{0,0}^H A_{0,1} + A_{0,0}^H A_{1,1} + \dots & \dots \\ A_{0,1}^H A_{0,0} + A_{0,1}^H A_{1,0} + \dots & A_{0,1}^H A_{0,1} + A_{0,1}^H A_{1,1} + \dots + B_1 & \dots \\ A_{0,2}^H A_{0,0} + A_{0,2}^H A_{1,0} + \dots & A_{0,2}^H A_{0,1} + A_{0,2}^H A_{1,1} + \dots & \dots \\ \vdots & \vdots & \ddots \end{pmatrix} \begin{pmatrix} \mathbf{x}_0 \\ \mathbf{x}_1 \\ \mathbf{x}_2 \\ \vdots \end{pmatrix} = \begin{pmatrix} \mathbf{c}_0 \\ \mathbf{c}_1 \\ \mathbf{c}_2 \\ \vdots \end{pmatrix}.$$

By the same re-indexing and reordering as in Appendix A, but now with $\tilde{\mathbf{a}}_{k,n}(m) = \mathbf{a}_{k,m}(n)^*$, this system can be represented by the equations

$$\sum_m \left(\sum_k \tilde{\mathbf{a}}_{k,n}(m') \tilde{\mathbf{a}}_{k,n}(m)^* + \tilde{\mathbf{b}}_n(m') \right) \tilde{\mathbf{x}}_n(m) = \tilde{\mathbf{c}}_n(m'), \quad (50)$$

which can be written using vectors products as

$$\left(\sum_k \tilde{\mathbf{a}}_{k,n} \tilde{\mathbf{a}}_{k,n}^H + \text{diag}(\tilde{\mathbf{b}}_n) \right) \tilde{\mathbf{x}}_n = \tilde{\mathbf{c}}_n. \quad (51)$$

Systems of this form can be efficiently solved by iterated application of the Sherman-Morrison formula, as shown in Appendix D.

APPENDIX D

ITERATED SHERMAN-MORRISON SOLUTION

Consider the solution of problems of the form

$$(J + \mathbf{a}_0 \mathbf{a}_0^H + \mathbf{a}_1 \mathbf{a}_1^H + \dots + \mathbf{a}_K \mathbf{a}_K^H) \mathbf{x} = \mathbf{b}. \quad (52)$$

Define $A_0 = J$ and $A_{k+1} = A_k + \mathbf{a}_k \mathbf{a}_k^H$. From the Sherman-Morrison formula we have

$$A_{k+1}^{-1} = (A_k + \mathbf{a}_k \mathbf{a}_k^H)^{-1} = A_k^{-1} - \frac{A_k^{-1} \mathbf{a}_k \mathbf{a}_k^H A_k^{-1}}{1 + \mathbf{a}_k^H A_k^{-1} \mathbf{a}_k}. \quad (53)$$

Now define $\alpha_{l,k} = A_l^{-1} \mathbf{a}_k$ and $\beta_k = A_k^{-1} \mathbf{b}$ so that $\alpha_{0,k} = J^{-1} \mathbf{a}_k$ and $\beta_0 = J^{-1} \mathbf{b}$, so that

$$\beta_{k+1} = A_{k+1}^{-1} \mathbf{b} \quad (54)$$

$$= A_k^{-1} \mathbf{b} - \frac{A_k^{-1} \mathbf{a}_k \mathbf{a}_k^H A_k^{-1} \mathbf{b}}{1 + \mathbf{a}_k^H A_k^{-1} \mathbf{a}_k} \quad (55)$$

$$= \beta_k - \frac{\alpha_{k,k} \mathbf{a}_k^H \beta_k}{1 + \mathbf{a}_k^H \alpha_{k,k}}, \quad (56)$$

and

$$\alpha_{l+1,k} = A_{l+1}^{-1} \mathbf{a}_k \quad (57)$$

$$= A_l^{-1} \mathbf{a}_k - \frac{A_l^{-1} \mathbf{a}_l \mathbf{a}_l^H A_l^{-1} \mathbf{a}_k}{1 + \mathbf{a}_l^H A_l^{-1} \mathbf{a}_l} \quad (58)$$

$$= \alpha_{l,k} - \frac{\alpha_{l,l} \mathbf{a}_l^H \alpha_{l,k}}{1 + \mathbf{a}_l^H \alpha_{l,l}}. \quad (59)$$

An iterative algorithm to compute the solution for the system Eq. (52), given by β_K , is easily derived from these equations. (This algorithm is equivalent to one previously proposed by Egidi and Maioni [64].) An algorithm for solving Eq. (52) when $J = \rho I$ is presented in Alg. 1.

REFERENCES

- [1] A. M. Bruckstein, D. L. Donoho, and M. Elad, "From sparse solutions of systems of equations to sparse modeling of signals and images," *SIAM Review*, vol. 51, no. 1, pp. 34–81, 2009.
- [2] S. Boyd, N. Parikh, E. Chu, B. Peleato, and J. Eckstein, "Distributed optimization and statistical learning via the alternating direction method of multipliers," *Foundations and Trends in Machine Learning*, vol. 3, no. 1, pp. 1–122, 2010.
- [3] J. Eckstein, "Augmented Lagrangian and alternating direction methods for convex optimization: A tutorial and some illustrative computational results," Rutgers University, Tech. Rep. RRR 32-2012, Dec. 2012.
- [4] S. Mallat and Z. Zhang, "Matching pursuits with time-frequency dictionaries," *IEEE Trans. Signal Process.*, vol. 41, no. 12, pp. 3397–3415, Dec. 1993.

Input: vectors $\{\mathbf{a}_k\}$, parameter ρ
Initialize: $\alpha = \rho^{-1} \mathbf{a}_0$, $\beta = \rho^{-1} \mathbf{b}$

```

for  $k \in \{1 \dots K\}$  do
     $\alpha$ 
     $\gamma_{k-1} = \frac{\alpha}{1 + \mathbf{a}_{k-1}^H \alpha}$ 
     $\beta = \beta - \gamma_{k-1} \mathbf{a}_{k-1}^H \beta$ 
    if  $k \leq K - 1$  then
         $\alpha = \rho^{-1} \mathbf{a}_k$ 
        for  $l \in \{1 \dots k\}$  do
             $\alpha = \alpha - \gamma_{l-1} \mathbf{a}_{l-1}^H \alpha$ 
        end
    end
end

```

Output: linear equation solution β

Algorithm 1: Iterated Sherman-Morrison

- [5] Y. C. Pati, R. Rezaifar, and P. S. Krishnaprasad, "Orthogonal matching pursuit: recursive function approximation with applications to wavelet decomposition," in *Conference Record of The Twenty-Seventh Asilomar Conference on Signals, Systems and Computers*, Nov. 1993, pp. 40–44.
- [6] R. Rubinstein, A. M. Bruckstein, and M. Elad, "Dictionaries for sparse representation modeling," *Proc. IEEE*, vol. 98, no. 6, pp. 1045–1057, Jun. 2010.
- [7] M. Elad and M. Aharon, "Image denoising via sparse and redundant representations over learned dictionaries," *IEEE Trans. Image Process.*, vol. 15, no. 12, pp. 3736–3745, Dec. 2006.
- [8] S. S. Chen, D. L. Donoho, and M. A. Saunders, "Atomic decomposition by basis pursuit," *SIAM J. Sci. Comput.*, vol. 20, no. 1, pp. 33–61, 1998.
- [9] B. Wohlberg, "Efficient convolutional sparse coding," in *Proc. IEEE Int. Conf. Acoust. Speech Signal Process. (ICASSP)*, May 2014, pp. 7173–7177.
- [10] K. Kavukcuoglu, P. Sermanet, Y.-L. Boureau, K. Gregor, M. Mathieu, and Y. A. LeCun, "Learning convolutional feature hierarchies for visual recognition," in *Adv. Neural Inf. Process. Syst.*, vol. 23, 2010, pp. 1090–1098.
- [11] M. D. Zeiler, G. W. Taylor, and R. Fergus, "Adaptive deconvolutional networks for mid and high level feature learning," in *Proc. IEEE Int. Conf. Comp. Vis. (ICCV)*, 2011, pp. 2018–2025.
- [12] Y. A. LeCun, L. Bottou, Y. Bengio, and P. Haffner, "Gradient-based learning applied to document recognition," *Proc. IEEE*, vol. 86, no. 11, pp. 2278–2324, Nov. 1998.
- [13] M. D. Zeiler, D. Krishnan, G. W. Taylor, and R. Fergus, "Deconvolutional networks," in *Proc. IEEE Conf. Comp. Vis. Pat. Recog. (CVPR)*, Jun. 2010, pp. 2528–2535.
- [14] J. Yang, K. Yu, and T. S. Huang, "Supervised translation-invariant sparse coding," *Proc. IEEE Conf. Comp. Vis. Pat. Recog. (CVPR)*, pp. 3517–3524, 2010.
- [15] P. Sermanet, K. Kavukcuoglu, S. Chintala, and Y. A. LeCun, "Pedestrian detection with unsupervised multi-stage feature learning," in *Proc. IEEE Conf. Comp. Vis. Pat. Recog. (CVPR)*, Jun. 2013, pp. 3626–3633.
- [16] B. Chen, G. Polatkan, G. Sapiro, D. Blei, D. Dunson, and L. Carin, "Deep learning with hierarchical convolutional factor analysis," *IEEE Trans. Pattern Anal. Mach. Intell.*, vol. 35, no. 8, pp. 1887–1901, Aug. 2013.
- [17] A. D. Szlam, K. Kavukcuoglu, and Y. A. LeCun, "Convolutional matching pursuit and dictionary training," *CoRR*, vol. arXiv:1010.0422, 2010.
- [18] M. Pachitariu, A. M. Packer, N. Pettit, H. Dalgleish, M. Hausser, and M. Sahani, "Extracting regions of interest from biological images with convolutional sparse block coding," in *Adv. Neural Inf. Process. Syst.*, 2013, vol. 26, pp. 1745–1753.
- [19] R. Chalasani, J. C. Principe, and N. Ramakrishnan, "A fast proximal method for convolutional sparse coding," in *Proc. Int. Joint Conf. Neural Net. (IJCNN)*, Aug. 2013.
- [20] H. Bristow, A. Eriksson, and S. Lucey, "Fast convolutional sparse coding," in *Proc. IEEE Conf. Comp. Vis. Pat. Recog. (CVPR)*, Jun. 2013, pp. 391–398.
- [21] M. Aharon, M. Elad, and A. M. Bruckstein, "K-SVD: An algorithm for designing overcomplete dictionaries for sparse representation," *IEEE Trans. Signal Process.*, vol. 54, no. 11, pp. 4311–4322, 2006.

- [22] M. S. Lewicki and T. J. Sejnowski, "Coding time-varying signals using sparse, shift-invariant representations," in *Adv. Neural Inf. Process. Syst.*, vol. 11, 1999, pp. 730–736.
- [23] W. Hashimoto and K. Kurata, "Properties of basis functions generated by shift invariant sparse representations of natural images," *Biological Cybernetics*, vol. 83, no. 2, pp. 111–118, 2000.
- [24] H. Wersing, J. Eggert, and E. Körner, "Sparse coding with invariance constraints," in *Artificial Neural Networks and Neural Information Processing – ICANN/ICONIP 2003*, ser. Lecture Notes in Computer Science, 2003, vol. 2714, pp. 385–392.
- [25] B. A. Olshausen, "Learning sparse, overcomplete representations of time-varying natural images," in *Proc. IEEE Int. Conf. Image Process. (ICIP)*, vol. 1, Sep. 2003, pp. 1–41.
- [26] T. Blumensath and M. E. Davies, "On shift-invariant sparse coding," in *Independent Component Analysis and Blind Signal Separation*, ser. Lecture Notes in Computer Science, 2004, vol. 3195, pp. 1205–1212.
- [27] —, "Shift-invariant sparse coding for single channel blind source separation," in *Signal Processing with Adaptive Sparse Structured Representations (SPARS)*, Nov. 2005.
- [28] M. D. Plumbley, S. A. Abdallah, T. Blumensath, and M. E. Davies, "Sparse representations of polyphonic music," *Signal Processing*, vol. 86, no. 3, pp. 417–431, 2006.
- [29] T. Blumensath and M. E. Davies, "Sparse and shift-invariant representations of music," *IEEE Trans. Audio, Speech, Language Process.*, vol. 14, no. 1, pp. 50–57, Jan. 2006.
- [30] P. Jost, P. Vandergheynst, S. Lesage, and R. Gribonval, "MoTIF: An efficient algorithm for learning translation invariant dictionaries," in *Proc. IEEE Int. Conf. Acoust. Speech Signal Process. (ICASSP)*, vol. 5, May 2006.
- [31] K. Skretting, J. H. Husøy, and S. O. Aase, "General design algorithm for sparse frame expansions," *Signal Process.*, vol. 86, no. 1, pp. 117–126, Jan. 2006.
- [32] M. Aharon, "Overcomplete dictionaries for sparse representation of signals," Ph.D. dissertation, Israel Institute of Technology, Nov. 2006.
- [33] K. Engan, K. Skretting, and J. H. Husøy, "Family of iterative LS-based dictionary learning algorithms, ILS-DLA, for sparse signal representation," *Digital Signal Processing*, vol. 17, no. 1, pp. 32–49, 2007.
- [34] R. Grosse, R. Raina, H. Kwong, and A. Y. Ng, "Shift-invariant sparse coding for audio classification," in *Proc. Twenty-Third Conf. on Uncertainty in Artificial Intell. (UAI)*, Jul. 2007, pp. 149–158.
- [35] B. Mailhé, S. Lesage, R. Gribonval, F. Bimbot, and P. Vandergheynst, "Shift-invariant dictionary learning for sparse representations: extending K-SVD," in *European Signal Processing Conference (EUSIPCO)*, 2008.
- [36] M. Mørup, M. N. Schmidt, and L. K. Hansen, "Shift invariant sparse coding of image and music data," Technical University of Denmark, Tech. Rep. IMM2008-04659, 2008.
- [37] J. J. Thiagarajan, K. N. Ramamurthy, and A. Spanias, "Shift-invariant sparse representation of images using learned dictionaries," in *IEEE Int. Workshop Mach. Learn. Signal Process. (MLSP)*, Oct. 2008, pp. 145–150.
- [38] M. Nakashizuka, H. Nishiura, and Y. Iiguni, "Sparse image representations with shift-invariant tree-structured dictionaries," in *Proc. IEEE Int. Conf. Image Process. (ICIP)*, Nov. 2009, pp. 2145–2148.
- [39] Y.-k. Tomokusa, M. Nakashizuka, and Y. Iiguni, "Sparse image representations with shift and rotation invariance constraints," in *Int. Symp. Intell. Signal Process. Commun. Syst. (ISPACS)*, Jan. 2009, pp. 256–259.
- [40] M. Mørup and M. N. Schmidt, "Transformation invariant sparse coding," in *IEEE Int. Workshop Mach. Learn. Signal Process. (MLSP)*, 2011, pp. 1–6.
- [41] Q. Barthélemy, A. Larue, A. Mayoue, D. Mercier, and J. I. Mars, "Shift & 2d rotation invariant sparse coding for multivariate signals," *IEEE Trans. Signal Process.*, vol. 60, no. 4, pp. 1597–1611, Apr. 2012.
- [42] G. Pope, C. Aubel, and C. Studer, "Learning phase-invariant dictionaries," in *Proc. IEEE Int. Conf. Acoust. Speech Signal Process. (ICASSP)*, May 2013, pp. 5979–5983.
- [43] C. Rusu, B. Dumitrescu, and S. A. Tsaftaris, "Explicit shift-invariant dictionary learning," *IEEE Signal Process. Lett.*, vol. 21, no. 1, pp. 6–9, Jan. 2014.
- [44] R. P. N. Rao and D. H. Ballard, "Development of localized oriented receptive fields by learning a translation-invariant code for natural images," *Network: Computation in Neural Systems*, vol. 9, no. 2, pp. 219–234, 1998.
- [45] C. Ekanadham, D. Tranchina, and E. P. Simoncelli, "Recovery of sparse translation-invariant signals with continuous basis pursuit," *IEEE Trans. Signal Process.*, vol. 59, no. 10, pp. 4735–4744, Oct. 2011.
- [46] T. Blumensath and M. E. Davies, "Unsupervised learning of sparse and shift-invariant decompositions of polyphonic music," in *Proc. IEEE Int. Conf. Acoust. Speech Signal Process. (ICASSP)*, vol. 5, May 2004, pp. 497–500.
- [47] B. Kong and C. C. Fowlkes, "Fast convolutional sparse coding (FCSC)," University of California, Irvine, Tech. Rep., May 2014.
- [48] J. Mairal, F. Bach, J. Ponce, G. Sapiro, and A. Zisserman, "Non-local sparse models for image restoration," in *Proc. IEEE Int. Conf. Comp. Vis. (ICCV)*, 2009, pp. 2272–2279.
- [49] M. V. Afonso, J. M. Bioucas-Dias, and M. A. T. Figueiredo, "Fast image recovery using variable splitting and constrained optimization," *IEEE Trans. Image Process.*, vol. 19, no. 9, pp. 2345–2356, 2010.
- [50] K. Engan, S. O. Aase, and J. H. Husøy, "Method of optimal directions for frame design," in *Proc. IEEE Int. Conf. Acoust. Speech Signal Process. (ICASSP)*, vol. 5, 1999, pp. 2443–2446.
- [51] M. V. Afonso, J. M. Bioucas-Dias, and M. A. T. Figueiredo, "An Augmented Lagrangian approach to the constrained optimization formulation of imaging inverse problems," *IEEE Trans. Image Process.*, vol. 20, no. 3, pp. 681–695, Mar. 2011.
- [52] Q. Liu, J. Luo, S. Wang, M. Xiao, and M. Ye, "An augmented Lagrangian multi-scale dictionary learning algorithm," *EURASIP J. Adv. Signal Process.*, vol. 2011, pp. 1–16, 2011.
- [53] M. J. Huiskes and M. S. Lew, "The MIR flicker retrieval evaluation," in *Proc. ACM Int. Conf. Multimedia Information Retrieval (MIR)*, 2008.
- [54] J. Mairal, G. Sapiro, and M. Elad, "Learning multiscale sparse representations for image and video restoration," *Multiscale Model. Simul.*, vol. 7, no. 1, pp. 214–241, 2008.
- [55] R. Gribonval, "Fast matching pursuit with a multiscale dictionary of Gaussian chirps," *IEEE Trans. Signal Process.*, no. 5, pp. 994–1001, 2001.
- [56] B. A. Olshausen, P. Sallee, and M. S. Lewicki, "Learning sparse image codes using a wavelet pyramid architecture," in *Adv. Neural Inf. Process. Syst.*, 2000, vol. 13, pp. 887–893.
- [57] P. Sallee and B. A. Olshausen, "Learning sparse multiscale image representations," in *Adv. Neural Inf. Process. Syst.*, 2002, vol. 15, pp. 1351–1358.
- [58] B. Ophir, M. Lustig, and M. Elad, "Multi-scale dictionary learning using wavelets," *IEEE J. Sel. Topics Signal Process.*, vol. 5, no. 5, pp. 1014–1024, Sep. 2011.
- [59] K. Skretting and K. Engan, "Image compression using learned dictionaries by RLS-DLA and compared with K-SVD," in *Proc. IEEE Int. Conf. Acoust. Speech Signal Process. (ICASSP)*, May 2011, pp. 1517–1520.
- [60] J. M. Hughes, D. N. Rockmore, and Y. Wang, "Bayesian learning of sparse multiscale image representations," *IEEE Trans. Image Process.*, vol. 22, no. 12, pp. 4972–4983, Dec. 2013.
- [61] J. Mairal, G. Sapiro, and M. Elad, "Multiscale sparse image representation with learned dictionaries," in *Proc. IEEE Int. Conf. Image Process. (ICIP)*, vol. 3, Sep. 2007, pp. 105–108.
- [62] B. Wohlberg, "SParse Optimization Research COde (SPORCO)," Matlab library available from <http://math.lanl.gov/~brendt/Software/SPORCO/>, 2014, version 0.01 **This library is being prepared for public release, and will be made available before publication of this manuscript.**
- [63] W. W. Hager, "Updating the inverse of a matrix," *SIAM Review*, vol. 31, no. 2, pp. 221–239, Jun. 1989.
- [64] N. Egidi and P. Maponi, "A Sherman-Morrison approach to the solution of linear systems," *Journal of Computational and Applied Mathematics*, vol. 189, no. 1-2, pp. 703–718, May 2006.

PLACE
PHOTO
HERE

Brendt Wohlberg received the BSc(Hons) degree in applied mathematics, and the MSc(Applied Science) and PhD degrees in electrical engineering from the University of Cape Town, South Africa, in 1990, 1993 and 1996 respectively. He is currently a staff scientist in Theoretical Division at Los Alamos National Laboratory, Los Alamos, NM. His research interests include sparse representations, exemplar-based methods for image restoration, signal and image processing inverse problems, and machine learning.