# SEPARABLE DICTIONARY LEARNING FOR CONVOLUTIONAL SPARSE CODING VIA SPLIT UPDATES

*Jorge Quesada*[†], *Paul Rodriguez*[†], *and Brendt Wohlberg\**

[†]Department of Electrical Engineering, Pontificia Universidad Catolica del Peru, Lima, Peru
*T-5, Theoretical Division, Los Alamos National Laboratory, Los Alamos, NM, USA

## ABSTRACT

Existing methods for constructing separable 2D dictionary filter banks approximate a set of $K$ non-separable filters via a linear combination of $R \ll K$ separable filters. This approach involves the inefficiency of learning an initial set of non-separable filters, and places an upper bound on the quality of the separable filter banks. In this paper, we propose a method to directly learn a set of $K$ separable dictionary filters from a given image training set by drawing ideas from convolutional dictionary learning (CDL) methods. We show that the separable filters obtained by our method match the performance of an equivalent number of non-separable filters. Furthermore, the computational performance of our learning method is shown to be substantially faster than a state-of-the-art non-separable CDL method for large numbers of filters or large training sets.

*Index Terms*— Convolutional Sparse Representation, Dictionary Learning, Separable filters

## 1. INTRODUCTION

Sparse representations and dictionary learning are well-known techniques in the field of signal and image processing, yielding effective approaches in tasks such as denoising, object recognition, and machine learning applications [1]. In particular, convolutional formulations, which model an image as a sum over a set of convolutions between coefficient maps and dictionary filters, have received increasing attention for their ability to represent whole images, as opposed to their patch-based counterparts [2]. The most common form of Convolutional Sparse Coding (CSC) problem is Convolutional Basis Pursuit Denoising (CBPDN)

$$\min_{\{\mathbf{x}_k\}} \frac{1}{2} \Big\| \sum_{k=1}^{K} D_k * \mathbf{x}_k - \mathbf{b} \Big\|_2^2 + \lambda \sum_{k=1}^{K} \|\mathbf{x}_k\|_1 , \quad (1)$$

where $\mathbf{b}$ is the observed image, $\{\mathbf{x}_k\}$ is the coefficient map set, and $\{D_k\}$ are the non-separable dictionary filters. The corresponding Convolutional Dictionary Learning (CDL) problem is

$$\min_{\{D_k, \mathbf{x}_{k,s}\}} \frac{1}{2} \sum_{s=1}^{S} \Big\| \sum_{k=1}^{K} D_k * \mathbf{x}_{k,s} - \mathbf{b}_s \Big\|_2^2 + \lambda \sum_{s=1}^{S} \sum_{k=1}^{K} \|\mathbf{x}_{k,s}\|_1$$
$$\text{s.t.} \quad \|D_k\|_2 = 1 \; \forall k , \quad (2)$$

where $\{\mathbf{b}_s\}$ is the set of training images.

It has been shown [3, 4] that using separable filters as dictionaries in tasks such as CSC or Convolutional Neural Network (CNN) applications provides significant improvements in computational performance with respect to non-separable implementations, with little loss in accuracy or reconstruction quality. In general, most of these methods rely on learning the separable filter set as an approximation of a previously obtained set of non-separable filters, by using

the equivalence

$$D_k \approx \sum_{r=1}^{R} \alpha_{kr} G_r \quad k \in \{1, 2, \dots, K\} , \quad (3)$$

which represents each non-separable filter $\{D_k\}$ as a linear combination of a smaller number of separable filters $\{G_r\}$, where $R << K$ [5]. This approach, however, depends heavily on the quality of the originating non-separable filters to obtain a good separable approximation. Furthermore, it implies a two step procedure: learning first the whole set of standard filters, and only then approximating the separable ones.

In this paper, we present an algorithm to directly learn the separable filters from the image training set by solving

$$\min_{\{\mathbf{h}_r, \mathbf{v}_r, \mathbf{x}_{r,s}\}} \frac{1}{2} \sum_{s=1}^{S} \Big\| \sum_{r=1}^{R} \mathbf{v}_r * \mathbf{h}_r * \mathbf{x}_{r,s} - \mathbf{b}_s \Big\|_2^2 + \lambda \sum_{s=1}^{S} \sum_{r=1}^{R} \|\mathbf{x}_{r,s}\|_1$$
$$\text{s.t.} \quad \|\mathbf{h}_r\|_2 = \|\mathbf{v}_r\|_2 = 1 \; \forall r , \quad (4)$$

where $R = K$, and $\{\mathbf{h}_r\}$ and $\{\mathbf{v}_r\}$ are the horizontal and vertical components of each filter. The proposed method is derived as a natural extension of a well known CDL algorithm [2], and compared against both standard non-separable dictionaries and separable approximations learned via (3). The computational results in Section 4 show that the separable filter banks obtained by our method provide superior performance to the approximated filter banks of the same size (i.e. the same number of filters) when evaluated via a standard CBPDN problem. Furthermore, the computational runtime of our learning algorithm is shown to be faster than standard non-separable learning approaches for most configurations.

## 2. PREVIOUS RELATED WORK

### 2.1. Non-separable (standard) dictionary learning

CDL problem (2) is non-convex when dealing with both variables ($\{\mathbf{x}_{k,s}\}$ and $\{D_k\}$) simultaneously, but becomes convex when keeping either of them constant. Therefore, the most widely used minimization approach consists in alternating between the updates for the feature maps $\{\mathbf{x}_{k,s}\}$ (sparse coding) and the filters $\{D_k\}$ (dictionary learning). This section will address the main existing dictionary learning update methods[1], which require solving a constrained convolutional form of the Method of Optimal Directions (MOD) [8]

$$\min_{\{D_k\}} \frac{1}{2} \sum_{s=1}^{S} \Big\| \sum_{k=1}^{K} D_k * \mathbf{x}_{k,s} - \mathbf{b}_s \Big\|_2^2 \quad \text{s.t.} \quad \|D_k\|_2 = 1 \; \forall k , \quad (5)$$

---

[1]See [6, 7] for a thorough review and comparison of sparse coding and dictionary learning updates and their coupling mechanisms.

where $\{\mathbf{x}_{k,s}\}$ is a given coefficient map set.

Early methods solved this problem in the spatial domain, via variants of gradient descent [9] and MOD [10], among others [11, 12]. More recent implementations solve the most computationally demanding components of the problem in the frequency domain due to the associated speedup [7]. When performing the convolutions in the frequency domain, the filters must be zero-padded in order to have an adequate spatial support. This requirement can be denoted by a zero-padding projection operator $P$, and coupled with the normalization constraint into the constraint set

$$C_{PN} = \{x \in \mathbb{R}^N : (I - PP^T)x = 0, \|x\|_2 = 1\} , \qquad (6)$$

which allows to write the dictionary update in unconstrained form

$$\min_{\{D_k\}} \frac{1}{2} \sum_{s=1}^{S} \left\| \sum_{k=1}^{K} D_k * \mathbf{x}_{k,s} - \mathbf{b}_s \right\|_2^2 + \sum_{r=1}^{R} \iota_{C_{PN}}(D_k) , \qquad (7)$$

where $\iota_{C_{PN}}(\cdot)$ is the indicator function of the constraint set $C_{PN}$. Several algorithms have been proposed to solve (7), most of which are based on Augmented Lagrangian frameworks, differing primarily on the approach they take to solve the $\ell_2$ fidelity term sub-problem. [13] proposed an Alternating Direction Method of Multipliers (ADMM [14]) formulation, which [2] and [15] later improved by efficiently approaching the aforementioned sub-problem using Iterated Sherman Morrison and ADMM consensus solutions, respectively. Modified ADMM consensus and Fast Iterative Shrinkage-Thresholding Algorithm (FISTA) [16] based methods have recently been shown to significantly outperform earlier alternatives [7].

There are also variants of these methods that perform the dictionary update in an online fashion, such as [17, 18, 19], in order to achieve scalability to very large training sets.

## 2.2. Separable from non-separable approximation

A straightforward approach to estimate $G_r$ (as defined in Equation (3)) from a given set of standard filters $\{D_k\}$ was proposed in [5, 3], involving placing a penalty on high-rank filters, namely

$$\min_{\{G_r, \alpha_{rk}\}} \frac{1}{2} \sum_{k=1}^{K} \left\| D_k - \sum_{r=1}^{R} \alpha_{rk} \cdot G_r \right\|_F^2 + \lambda \sum_{r=1}^{R} \|G_r\|_* , \quad (8)$$

where $\|\cdot\|_*$ is the nuclear norm. [5, 3] highlighted that the choice of $\lambda$ is a challenging task, and that convergence was slow when estimating high-rank filters. They also proposed a second approach based on the Canonical Polyadic Decomposition [20] that provides faster performance

$$\min_{\{\alpha_{rk}, x_r, y_r\}} \frac{1}{2} \sum_{k=1}^{K} \left\| D_k - \sum_{r=1}^{R} \alpha_{rk} \cdot x_r \circ y_r \right\|_F^2 , \qquad (9)$$

where $x_r$ and $y_r$ are rank-1 tensors and $\circ$ represents tensor outer product. A reformulation of this problem as a special case of the low-rank basis problem was proposed in [21], but the authors reported that the tensor approach was significantly faster and attained the same accuracy.

An auxiliary variable formulation of (8) given by

$$\min_{\{G_r, \alpha_{rk}, F_r\}} \frac{1}{2} \sum_{k=1}^{K} \left\| D_k - \sum_{r=1}^{R} \alpha_{rm} G_r \right\|_F^2 + \frac{\lambda}{2} \sum_{r=1}^{R} \|G_r - F_r\|_F^2$$
$$\text{s.t. rank}(F_r) = 1 \quad \forall r \qquad (10)$$

was proposed in [22] along with an efficient SVD-based generation of the initial solution. The method was shown to be faster than the

tensor decomposition approach for small $R$ ($< 40$) values while attaining comparable accuracy.

## 3. PROPOSED METHOD

Writing the dictionary update for (4) and coupling the norm constraint with the zero-padding restriction described in Section 2.1 gives the unconstrained problem

$$\min_{\{\mathbf{h}_r, \mathbf{v}_r\}} \frac{1}{2} \sum_{s=1}^{S} \left\| \sum_{r=1}^{R} \mathbf{v}_r * \mathbf{h}_r * \mathbf{x}_{r,s} - \mathbf{b}_s \right\|_2^2 +$$
$$\sum_{r=1}^{R} \left( \iota_{C_{\text{PhN}}}(\mathbf{h}_r) + \iota_{C_{\text{PvN}}}(\mathbf{v}_r) \right) , \qquad (11)$$

where $\iota_{C_{\text{PhN}}}(\cdot)$ and $\iota_{C_{\text{PvN}}}(\cdot)$ are the indicator functions of the constraint sets $C_{\text{PhN}}$ and $C_{\text{PvN}}$ (analogous to (6)), with zero-padding operators $P_h$ (applied along the horizontal dimension) and $P_v$ (applied along the vertical dimension), respectively.

We approach the solution of (11) by alternating between updating the horizontal filters $\mathbf{h}_r$ and the vertical ones $\mathbf{v}_r$. Considering only the solution for the vertical filters $\mathbf{v}_r$ (assuming fixed horizontal filters), and reformulating the problem in ADMM-compatible form in a fashion reminiscent of [2] leads to

$$\min_{\{\mathbf{v}_r, g_r\}} \frac{1}{2} \sum_{s=1}^{S} \left\| \sum_{r=1}^{R} \mathbf{v}_r * \mathbf{x}'_{r,s} - \mathbf{b}_s \right\|_2^2 + \sum_{r=1}^{R} \iota_{C_{\text{PvN}}}(g_r)$$
$$\text{s.t.} \quad \mathbf{v}_r - g_r = 0 \ \forall r , \qquad (12)$$

where $\mathbf{x}'_{r,s}$ is the result of convolving the horizontal filters $\mathbf{h}_r$ with the feature maps $\mathbf{x}_{r,s}$. The associated subproblems are then given by

$$\mathbf{v}_r^{(i+1)} = \arg\min_{\mathbf{v}_r} \frac{1}{2} \sum_{s=1}^{S} \left\| \sum_{r=1}^{R} \mathbf{v}_r * \mathbf{x}'_{r,s} - \mathbf{b}_s \right\|_2^2$$
$$+ \frac{\rho}{2} \sum_{r=1}^{R} \left\| \mathbf{v}_r - g_r^{(i)} + f_r^{(i)} \right\|_2^2 \qquad (13)$$

$$g_r^{(i+1)} = \arg\min_{g_r} \sum_{r=1}^{R} \iota_{C_{\text{PvN}}}(g_r) + \frac{\rho}{2} \sum_{r=1}^{R} \left\| \mathbf{v}_r^{(i+1)} - g_r + f_r^{(i)} \right\|_2^2 \qquad (14)$$

$$f_r^{(i+1)} = f_r^{(i)} + \mathbf{v}_r^{(i+1)} - g_r^{(i+1)} . \qquad (15)$$

Since (14) is of the form

$$\arg\min_{\mathbf{x}} \frac{1}{2} \|\mathbf{x} - \mathbf{y}\|_2^2 + \iota_{C_{\text{PvN}}}(\mathbf{x}) = \text{prox}_{\iota_{C_{\text{PvN}}}}(\mathbf{y}) , \qquad (16)$$

its minimizer is given by

$$\text{prox}_{\iota_{C_{\text{PvN}}}}(\mathbf{y}) = P_v P_v^T \mathbf{y} \ / \ \|P_v P_v^T \mathbf{y}\|_2 . \qquad (17)$$

For notational simplicity we rewrite (13) as

$$\mathbf{v}_r^{(i+1)} = \arg\min_{\mathbf{v}_r} \frac{1}{2} \sum_{s=1}^{S} \left\| \sum_{r=1}^{R} \mathbf{v}_r * \mathbf{x}'_{r,s} - \mathbf{b}_s \right\|_2^2$$
$$+ \frac{\rho}{2} \sum_{r=1}^{R} \|\mathbf{v}_r - z_r\|_2^2 , \qquad (18)$$

where $z_r = g_r^{(i)} - f_r^{(i)}$.

When performing standard CDL [2], the non-separable equivalent of (18) is solved by switching to the Fourier domain and solving

the associated linear system. In the separable case, however, it is worth noting that since the filters $\{\mathbf{v}_r\}$ are 1-D, whereas the coefficient maps $\{\mathbf{x}'_{r,s}\}$ are 2-D, moving directly onto the frequency domain would require the DFT solution ($\hat{v}_r$) to be a 2-D matrix composed of replicating columns. This would mean including an additional constraint and further increasing the complexity of the problem. Instead we choose to rewrite the fidelity $\ell_2$-norm term as a sum over columns

$$\left\| \sum_{r=1}^{R} \mathbf{v}_r * \mathbf{x}'_{r,s} - \mathbf{b}_s \right\|_2^2 = \sum_{i=1}^{I} \left\| \sum_{r=1}^{R} \mathbf{v}_r * \mathbf{x}'_{r,s}[i] - \mathbf{b}_s[i] \right\|_2^2 , \quad (19)$$

where $\mathbf{x}'_{r,s}[i]$ and $\mathbf{b}_s[i]$ are the i-th columns of the corresponding feature map and training image respectively. Replacing the equality in (18) gives

$$\arg\min_{\mathbf{v}_r} \frac{1}{2} \sum_{s=1}^{S} \sum_{i=1}^{I} \left\| \sum_{r=1}^{R} \mathbf{v}_r * \mathbf{x}'_{r,s}[i] - \mathbf{b}_s[i] \right\|_2^2 + \frac{\rho}{2} \sum_{r=1}^{R} \left\| \mathbf{v}_r - z_r \right\|_2^2 .$$

Switching to the DFT domain, and defining $\hat{\mathbf{X}}'_{r,s}[i] = \mathrm{diag}(\hat{\mathbf{x}}'_{r,s}[i])$ gives

$$\arg\min_{\mathbf{v}_r} \frac{1}{2} \sum_{s=1}^{S} \sum_{i=1}^{I} \left\| \sum_{r=1}^{R} \hat{\mathbf{X}}'_{r,s}[i] \hat{v}_r - \hat{\mathbf{b}}_s[i] \right\|_2^2 + \frac{\rho}{2} \sum_{r=1}^{R} \| \hat{v}_r - \hat{z}_r \|_2^2$$

Defining

$$\hat{\mathbf{X}}'_s[i] = (\hat{\mathbf{X}}'_{0,s}[i] \quad \hat{\mathbf{X}}'_{1,s}[i] \quad \dots ) \quad \hat{v} = \begin{bmatrix} \hat{v}_1 \\ \hat{v}_2 \\ \vdots \\ \hat{v}_R \end{bmatrix} \quad \hat{z} = \begin{bmatrix} \hat{z}_1 \\ \hat{z}_2 \\ \vdots \\ \hat{z}_R \end{bmatrix} \quad (20)$$

the problem can be expressed as

$$\arg\min_{\mathbf{v}_r} \frac{1}{2} \sum_{s=1}^{S} \sum_{i=1}^{I} \left\| \hat{\mathbf{X}}'_s[i] \hat{v} - \hat{\mathbf{b}}_s[i] \right\|_2^2 + \frac{\rho}{2} \| \hat{v} - \hat{z} \|_2^2 . \quad (21)$$

Finally, to further simplify (21) we define

$$\hat{\mathbf{X}}'_s = \begin{bmatrix} \hat{\mathbf{X}}'_s[1] \\ \hat{\mathbf{X}}'_s[2] \\ \vdots \\ \hat{\mathbf{X}}'_s[I] \end{bmatrix} \quad (22)$$

Substituting $\hat{\mathbf{X}}'_s$ and recovering the full vectorized DFT training images $\hat{\mathbf{b}}_s$ leads to the problem being expressed as

$$\arg\min_{\mathbf{v}_r} \frac{1}{2} \sum_{s=1}^{S} \left\| \hat{\mathbf{X}}'_s \hat{v} - \hat{\mathbf{b}}_s \right\|_2^2 + \frac{\rho}{2} \| \hat{v} - \hat{z} \|_2^2 , \quad (23)$$

with solution

$$\left( \sum_s \hat{\mathbf{X}}'^H_s \hat{\mathbf{X}}'_s + \rho I \right) \hat{v} = \sum_s \hat{\mathbf{X}}'^H_s \hat{\mathbf{b}}_s + \rho \hat{z} . \quad (24)$$

Due to the commutativity property of the convolution operation, the update for the horizontal filters $\mathbf{h}_r$ can be easily derived by fixing the vertical filters, defining $\mathbf{x}'_{r,s} = \mathbf{v}_r * \mathbf{x}_{r,s}$, and following an analogous chain of derivations as the one described in this section.

### 3.1. Implementation remarks

The linear system given by Eq. (24) is solved by applying Conjugate Gradient (CG)[2]. Furthermore, in order to minimize the number of

---

[2]While another widely used method to deal with (24) is the Iterative Sherman Morrison (ISM) approach from [2], the column indexing we introduce here entails an additional computational overhead that renders ISM imprac-

inner CG iterations, we use the solution for each previous update as the initial value, as suggested in [7].

The full dictionary learning algorithm is implemented by combining the proposed update method for $\{\mathbf{v}_r\}$ and $\{\mathbf{h}_r\}$ with the ADMM-based sparse coding update proposed in [2]. Based on standard non-separable implementations, and the results provided by [6], we interleave a single iteration of each update per outer loop, and transfer the auxiliary variables of each ADMM framework across the other update steps, which has been shown to provide the most stable convergence ratio among the other possible choices [6].

## 4. RESULTS

In this section we assess the performance of the proposed separable dictionary learning method in terms of reconstruction performance for a CSC-based denoising problem, along with convergence and computational runtime for the learning process.

### 4.1. Experimental framework

For the denoising comparisons, we used a set of 5 well-known images corrupted with AWGN ($\sigma = 0.2$), to perform CBPDN using the following labeled set of filters of different sizes (see Table 1):

- **Nat-sep**: 36 Natively learned separable filters (our proposed method)
- **Apr-sep**: 36 Separable filters approximated from 36 non-separable ones via [22]
- **Non-sep**: 36 Standard non-separable filters learned via [2]

Since the CBPDN problem Eq. (1) has a tunable parameter $\lambda$, we ensure a fair evaluation by solving for a grid of $\lambda$ values and comparing only the optimal performance for each of the evaluated filter sets. An example of the entire simulation results for a single image is given in Figure 3.
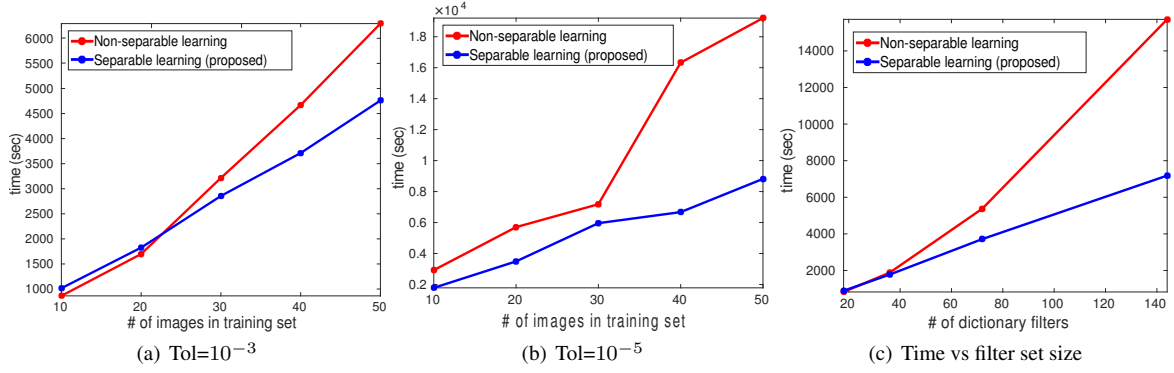
For the separable dictionaries (nat-sep and apr-sep), we use the $\ell_1$ version of the FISTA-based CBPDN solver proposed in [23] that exploits filter separability by computing the convolutions in the spatial domain. For the non-separable dictionaries (non-sep), we use the ADMM-based solver from [2], which is considered to be state-of-the-art for this problem.

For the computational performance simulations, we evaluate the learning time on the full CDL task for different training set sizes ($S$) and filter set sizes ($R = K$) against a state-of-the-art ADMMM-based non-separable CDL method [2]. These simulations were performed on an Intel Xeon E5-2640 CPU (2,50 GHz , 128Gb RAM, 2x NVidia Tesla K40m GPU). Our Matlab code [24] can be used to reproduce our experimental results.

### 4.2. Experiments

In Table 1 we illustrate the results of the denoising comparisons between the 3 evaluated filter sets in terms of the SSIM metric for different dictionary sizes, and report the average runtime for each method across the grid of $\lambda$ values. It can be observed that the natively separable filters consistently outperform the approximated (separable) ones, and show equivalent performance to the non-separable filters. The runtime results also show that performing CSC with separable filters is almost two times faster than doing it with non separable ones, which is consistent with the results reported in [23]. We also show in Figure 3 the entire set of denoising simulations across the $\lambda$ grid for a single image.
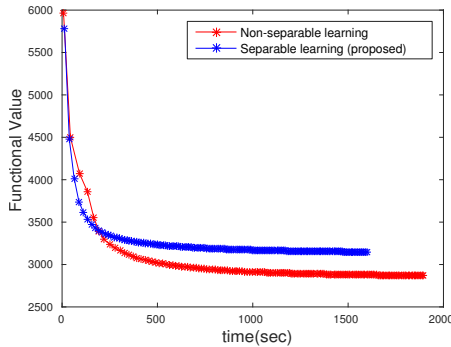
---

tical for this task.

(a) Tol=$10^{-3}$     (b) Tol=$10^{-5}$     (c) Time vs filter set size

**Fig. 1**: Computational performance (separable vs. non-separable) results for CDL simulations, with termination at 200 iterations.

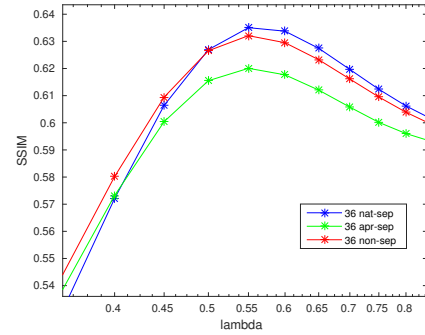|  | Dict. Size | barbara | mandrill | parrots | boats | goldhill | Time |
|---|---|---|---|---|---|---|---|
| | 8x8 | 0.6175 | 0.5188 | 0.7188 | 0.6438 | 0.6709 | 40,77 |
| nat-sep | 12x12 | 0.6370 | 0.5248 | 0.7219 | 0.6532 | 0.6730 | 60,57 |
| | 16x16 | 0.6285 | 0.5223 | 0.7197 | 0.6554 | 0.6728 | 70,11 |
| | 8x8 | 0.6189 | 0.5218 | 0.7207 | 0.6449 | 0.6732 | 70,3 |
| non-sep | 12x12 | 0.6330 | 0.5300 | 0.7225 | 0.6507 | 0.6763 | 104,7 |
| | 16x16 | 0.6283 | 0.5257 | 0.7218 | 0.6536 | 0.6741 | 112,6 |
| | 8x8 | 0.6147 | 0.5015 | 0.7118 | 0.6335 | 0.6659 | 40,82 |
| apr-sep | 12x12 | 0.6122 | 0.5186 | 0.7132 | 0.6396 | 0.6640 | 60,64 |
| | 16x16 | 0.6207 | 0.5157 | 0.7151 | 0.6502 | 0.6686 | 70,52 |

**Table 1**: Denoising performance (SSIM) for different filter sizes



**Fig. 2**: Functional value behaviour comparison for CDL task

We report in Figure 1 (a) and (b) the computational performance comparisons in the learning process for 2 different CG tolerance values, in terms of runtime (seconds) vs image training set size. We consider a fixed number of 36 separable and non-separable filters for this simulation, and measure the runtime for both training methods for a fixed number of iterations (200). As can be observed in the graph, when the CG tolerance is $10^{-3}$ the proposed separable method is slightly slower than its non-separable counterpart [2] for small values of $S$, and outperforms it when $S$ increases. When the tolerance value is $10^{-5}$, the proposed method significantly outperforms [2] as $S$ increases. Figure 1 (c) depicts a similar runtime comparison where the training set size is fixed ($S = 20$) and the dictionary size (number of filters) is varied (the tolerance value used

is $10^{-3}$). In tmakehis case it is also clear that the proposed method is substantially faster than the non-separable method as the number of filters increases. An example of the functional value behaviour



**Fig. 3**: Denoising results on $\lambda$ grid for 'barbara' image, where *apr-sep, nat-sep* and *non-sep* are the labels defined in Section 4.1

for a training set size of $S = 20$ is shown in Figure 2 for 200 iterations. It can be seen from the graph that the proposed separable method converges to a slightly higher functional value than the non-separable method. However, this difference does not seem to have a significant impact on the performance quality of the learned separable filters, as can be seen on Table 1.

## 5. CONCLUSIONS

We have proposed an efficient method to learn separable dictionary filters directly from an image training set, without the need to previously compute a set of non-separable filters. Our results show that the separable filters learned through this method, when evaluated through a CSC denoising task, consistently outperform approximated separable filters, and attain the same reconstruction quality as obtained from standard non-separable filters. Furthermore, the proposed separable learning method is substantially faster than its non-separable counterpart when either the training set or the number of filters to estimate is large.

## 6. REFERENCES

[1] J. Mairal, F. Bach, and J. Ponce, "Sparse modeling for image and vision processing," *Foundations and Trends in Computer Graphics and Vision*, vol. 8, no. 2-3, pp. 85–283, 2014. doi:10.1561/0600000058

[2] B. Wohlberg, "Efficient algorithms for convolutional sparse representations," *IEEE Transactions on Image Processing*, vol. 25, no. 1, pp. 301–315, Jan. 2016. doi:10.1109/TIP.2015.2495260

[3] A. Sironi, B. Tekin, R. Rigamonti, V. Lepetit, and P. Fua, "Learning separable filters," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 37, no. 1, pp. 94–106, Jan. 2015. doi:10.1109/TPAMI.2014.2343229

[4] M. Jaderberg, A. Vedaldi, and A. Zisserman, "Speeding up convolutional neural networks with low rank expansions," in *Proceedings of the British Machine Vision Conference*. BMVA Press, 2014

[5] R. Rigamonti, A. Sironi, V. Lepetit, and P. Fua, "Learning separable filters," in *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2013, pp. 2754–2761. doi:10.1109/CVPR.2013.355

[6] C. Garcia-Cardona and B. Wohlberg, "Subproblem coupling in convolutional dictionary learning," in *Proceedings of IEEE International Conference on Image Processing (ICIP)*, Beijing, China, Sep. 2017

[7] ——, "Convolutional Dictionary Learning," Sep. 2017. arXiv:1709.02893

[8] K. Engan, S. O. Aase, and J. H. Husoy, "Method of optimal directions for frame design," in *Proceedings of IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, vol. 5, 1999, pp. 2443–2446 vol.5. doi:10.1109/ICASSP.1999.760624

[9] M. Mrup and M. N. Schmidt, "Transformation invariant sparse coding," in *IEEE International Workshop on Machine Learning for Signal Processing*, Sep. 2011, pp. 1–6. doi:10.1109/MLSP.2011.6064547

[10] M. D. Zeiler, G. W. Taylor, and R. Fergus, "Adaptive deconvolutional networks for mid and high level feature learning," in *International Conference on Computer Vision*, Nov. 2011, pp. 2018–2025. doi:10.1109/ICCV.2011.6126474

[11] Q. Barthelemy, A. Larue, A. Mayoue, D. Mercier, and J. I. Mars, "Shift & 2d rotation invariant sparse coding for multivariate signals," *IEEE Transactions on Signal Processing*, vol. 60, no. 4, pp. 1597–1611, Apr. 2012. doi:10.1109/TSP.2012.2183129

[12] R. Chalasani, J. C. Principe, and N. Ramakrishnan, "A fast proximal method for convolutional sparse coding," in *International Joint Conference on Neural Networks (IJCNN)*, Aug. 2013, pp. 1–5. doi:10.1109/IJCNN.2013.6706854

[13] H. Bristow, A. Eriksson, and S. Lucey, "Fast convolutional sparse coding," in *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, Jun. 2013, pp. 391–398. doi:10.1109/CVPR.2013.57

[14] S. Boyd, N. Parikh, E. Chu, B. Peleato, and J. Eckstein, "Distributed optimization and statistical learning via the alternating direction method of multipliers," *Found. Trends Mach. Learn.*, vol. 3, no. 1, pp. 1–122, Jan. 2011. doi:10.1561/2200000016

[15] M. Šorel and F. Šroubek, "Fast convolutional sparse coding using matrix inversion lemma," *Digital Signal Processing*, vol. 55, pp. 44–51, 2016. doi:10.1016/j.dsp.2016.04.012

[16] A. Beck and M. Teboulle, "A fast iterative shrinkage-thresholding algorithm for linear inverse problems," *SIAM Journal on Imaging Sciences*, vol. 2, no. 1, pp. 183–202, 2009. doi:10.1137/080716542

[17] K. Degraux, U. S. Kamilov, P. T. Boufounos, and D. Liu, "Online Convolutional Dictionary Learning for Multimodal Imaging," Jun. 2017. arXiv:1706.04256

[18] J. Liu, C. Garcia-Cardona, B. Wohlberg, and W. Yin, "Online Convolutional Dictionary Learning," Jun. 2017. arXiv:1706.09563

[19] J. Liu, C. Garcia-Cardona, B. Wohlberg, and W. Yin, "First and second order methods for online convolutional dictionary learning," *SIAM Journal on Imaging Sciences*, 2018, accepted for publication. arXiv:1709.00106

[20] T. G. Kolda and B. W. Bader, "Tensor decompositions and applications," *SIAM Review*, vol. 51, no. 3, pp. 455–500, 2009. doi:10.1137/07070111X

[21] Y. Nakatsukasa, T. Soma, and A. Uschmajew, "Finding a low-rank basis in a matrix subspace," *Mathematical Programming*, vol. 162, no. 1, pp. 325–361, Mar 2017. doi:10.1007/s10107-016-1042-2

[22] P. Rodríguez, "Alternating optimization low-rank expansion algorithm to estimate a linear combination of separable filters to approximate 2d filter banks," in *50th Asilomar Conference on Signals, Systems and Computers*, Nov. 2016, pp. 954–958. doi:10.1109/ACSSC.2016.7869190

[23] G. Silva, J. Quesada, P. Rodríguez, and B. Wohlberg, "Fast convolutional sparse coding with separable filters," in *IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, Mar. 2017, pp. 6035–6039. doi:10.1109/ICASSP.2017.7953315

[24] J. Quesada and P. Rodriguez, "Separable filter learning," available at https://sites.google.com/pucp.pe/jquesada