

# Sparse shift-invariant NMF

Vamsi K. Potluru  
MIND Research Network,  
Dept. of Computer Science,  
University of New Mexico  
ismav@cs.unm.edu

Sergey M. Plis  
Dept. of Computer Science,  
University of New Mexico  
pliz@cs.unm.edu

Vince D. Calhoun  
MIND Research Network,  
Dept. of Elec and Comp Engg,  
University of New Mexico  
vcalhoun@unm.edu

## Abstract

*Non-negative Matrix factorization (NMF) has increasingly been used for efficiently decomposing multivariate data into a signal dictionary and corresponding activations. In this paper, we propose an algorithm called sparse shift-invariant NMF (ssiNMF) for learning possibly overcomplete shift-invariant features. This is done by incorporating a circulant property on the features and sparsity constraints on the activations. The circulant property allows us to capture shifts in the features and enables efficient computation by the Fast Fourier Transform. The ssiNMF algorithm turns out to be matrix-free for we need to store only a small number of features. We demonstrate this on a dataset generated from an overcomplete set of bars.*

## 1 Introduction

Non-negative matrix factorization (NMF) is a tool to split the given data matrix into a product of two non-negative matrix factors. This process can be used to identify useful features in the dataset. The constraint of non-negativity results in a parts-based representation and is usually different from other factorization techniques which result in more holistic representations (e.g. principal component analysis (PCA) and vector quantization (VQ)). Another tool used commonly to find features is independent component analysis (ICA) [1]. ICA as-

sumes that the features thus found are statistically independent [2].

NMF intrinsically enforces certain amount of sparsity in its representations. However, in the case of overcomplete representations, we would like to explicitly enforce a sparsity constraint. NMF with a sparsity constraint on the activations was introduced in [4]. Convolutional NMF had previously been studied in [3] within application to audio data. It was extended with a sparsity constraint in [? ].

In this paper, we combine convolutional dictionary with sparse activations. Unlike the previous approaches, we constrain the features to be circulant to model arbitrary shifts in the data. This property is useful for training datasets that have misaligned instances, such as datasets of images. We demonstrate the utility of our approach using the dataset of [4] by learning a parsimonious dictionary to represent it.

## 2 NMF

Given a non-negative  $m \times n$  matrix  $\mathbf{X}$ , we want to represent it with a product of two non-negative matrices  $\mathbf{W}$ ,  $\mathbf{H}$  of sizes  $m \times r$  and  $r \times n$  respectively:

$$\mathbf{X} \approx \mathbf{W}\mathbf{H}. \quad (1)$$

The non-negativity constraint corresponds to the intuitive notion of features adding up to give the resulting data.

Lee and Seung [6] describe two simple multiplicative updates which work well in practice.

These correspond to two different cost functions representing the quality of approximation. Here, we use the Frobenius norm for the cost function. The cost function and the corresponding updates are:

$$E = \|X - WH\|_F \quad (2)$$

$$W = W \odot \frac{XH^T}{WHH^T} \quad (3)$$

$$H = H \odot \frac{W^T X}{W^T WH}, \quad (4)$$

where  $\|\cdot\|_F$  denotes the Frobenius norm and  $\|\cdot\|_1$  the  $L_1$  norm. The operator  $\odot$  represents element-wise multiplication and division is also element-wise. It should be noted that the cost function to be minimized is convex in either  $W$  or  $H$  but not in both. As the algorithm iterates using the updates given,  $W$  and  $H$  converge to a local minimum of the cost function. The value of  $r$  determines quality of approximation and is usually based on prior knowledge of the data being modelled.

### 3 Sparse NMF

NMF with a sparsity constraint was introduced in [4]. It was shown that explicitly controlling sparsity gives better decompositions. Using  $L_1$  norm for sparsity, sparse NMF is formulated as follows:

$$\min_{W, H} \frac{1}{2} \|X - WH\|_F + \lambda \|H\|_1 \quad (5)$$

The update equations for this objective are given by:

$$W = W - \eta[-XH^T + WHH^T] \quad (6)$$

$$H = H \odot \frac{W^T X}{W^T WH + \lambda \mathbf{1}} \quad (7)$$

The parameter  $\eta$  is the learning rate and has to be explicitly set. As has already been noted in [4], the objective function is not scale free. This can be seen by setting  $W \leftarrow \alpha W$  and  $H \leftarrow \frac{1}{\alpha} H$ , with  $\alpha > 1$ . To combat this problem, new multiplicative update rules were derived by Egart and Konner [?]. The update for matrix  $W$  is given by

$$W = W \odot \frac{XH^T + W \text{diag}(1(WHH^T \odot W))}{WHH^T + W \text{diag}(1(XH^T \odot W))} \quad (8)$$

## 4 Matrix-free computations

If we allow features to be shifted within the data In the case of circulant matrices, matrix-vector product can be efficiently computed by using Fast Fourier Transform (FFT). In addition, a circulant matrix of size  $n \times n$  requires storage space of  $O(n)$ .

### 4.1 Circulant matrices

Let us introduce two operators :

- Circulant-shift operator  $S^i(\mathbf{v})$  : given a vector  $\mathbf{v}$  and a shift size  $i$ , returns the right circularly-shifted vector shifted by  $i$  positions.
- Flip operator  $FLIP(\mathbf{v})$  : returns a permuted vector with the  $i$ -th element replaced by the  $n - i + 1$ -th element of the given vector.

The circulant matrix with the first column equal to the vector  $\mathbf{c}$  is given by

$$C = \begin{bmatrix} S^0(\mathbf{c}) & S^1(\mathbf{c}) & \dots & S^{n-1}(\mathbf{c}) \end{bmatrix} \quad (9)$$

$$= \text{cm}(\mathbf{c})$$

We note that even though it has  $O(n^2)$  elements, it can be generated from  $\mathbf{c}$ , which has  $O(n)$  elements.

If  $\mathbf{f} = S^1(FLIP(\mathbf{c}))$  then

$$C^T = \begin{bmatrix} S^0(\mathbf{f}) & S^1(\mathbf{f}) & \dots & S^{n-1}(\mathbf{f}) \end{bmatrix} \quad (10)$$

$$= \text{cm}(\mathbf{f})$$

### 4.2 Circulant Matrix-vector product

Here, we outline an efficient routine to calculate the product of a circulant matrix  $C$  whose first column is  $\mathbf{c}$  with an appropriately sized vector  $\mathbf{r}$ . Let us denote by  $FFT$  and  $iFFT$  the routines of Fast Fourier Transform and inverse Fast Fourier Transform respectively. We then have :

$$C\mathbf{r} = iFFT(\text{diag}(FFT(\mathbf{c})) FFT(\mathbf{r})) \quad (11)$$

$$= iFFT(FFT(\mathbf{c}) \odot FFT(\mathbf{r})) \quad (12)$$

$$= \text{mvc}(\mathbf{c}, \mathbf{r}) \quad (13)$$

### 4.3 Composite Circulant products

Let us define the matrix  $A$  to be composite circulant if its elements are square circulant matrices. Matrix  $A$  is completely characterized by matrix  $B$  given the following relations:

$$\mathbf{B} = [\mathbf{b}_1 \ \mathbf{b}_2 \ \dots \ \mathbf{b}_r] \quad (14)$$

$$\mathbf{A} = [\text{cm}(\mathbf{b}_1) \ \text{cm}(\mathbf{b}_2) \ \dots \ \text{cm}(\mathbf{b}_r)] \quad (15)$$

The matrix-vector products given an appropriately sized vector  $\mathbf{y}$  with matrix  $\mathbf{A}$  are given by:

$$\begin{aligned} \mathbf{A}\mathbf{y} &= [\text{cm}(\mathbf{b}_1) \ \text{cm}(\mathbf{b}_2) \ \dots \ \text{cm}(\mathbf{b}_r)] \begin{bmatrix} \mathbf{y}_1 \\ \mathbf{y}_2 \\ \vdots \\ \mathbf{y}_r \end{bmatrix} \\ &= \sum_i \text{mvc}(\mathbf{b}_i, \mathbf{y}_i) \\ &= \text{fmvc}(\mathbf{B}, \mathbf{y}) \end{aligned} \quad (16)$$

$$\begin{aligned} \mathbf{A}^\top \mathbf{y} &= \begin{bmatrix} \text{cm}(\mathbf{b}_1)^\top \\ \text{cm}(\mathbf{g}_2)^\top \\ \vdots \\ \text{cm}(\mathbf{b}_r)^\top \end{bmatrix} \mathbf{y} \\ &= \begin{bmatrix} \text{mvc}(S^1(\text{FLIP}(\mathbf{b}_1)), \mathbf{y}) \\ \text{mvc}(S^1(\text{FLIP}(\mathbf{b}_2)), \mathbf{y}) \\ \vdots \\ \text{mvc}(S^1(\text{FLIP}(\mathbf{b}_r)), \mathbf{y}) \end{bmatrix} \\ &= \text{tfmvc}(\mathbf{B}, \mathbf{y}) \end{aligned} \quad (17)$$

## 5 Sparse shift-invariant NMF

In the ssiNMF framework, we model the dictionary  $\mathbf{W}$  to be a set of circularly shifted features. This is captured by matrix  $\mathbf{G}$  representing the features and matrix  $\mathbf{W}$  the set of all possible linear shifts. The relationship between the matrices is :

$$\mathbf{G} = [\mathbf{g}_1 \ \mathbf{g}_2 \ \dots \ \mathbf{g}_r] \quad (18)$$

$$\mathbf{W} = [\text{cm}(\mathbf{g}_1) \ \text{cm}(\mathbf{g}_2) \ \dots \ \text{cm}(\mathbf{g}_r)] \quad (19)$$

We need to store only the matrix  $\mathbf{G}$  to generate the full matrix  $\mathbf{W}$ . This makes the algorithm matrix-free and computationally efficient for using FFT's to compute matrix-vector products.

Given the data matrix  $\mathbf{X}$ , we apply Algorithm 1 denoted by ssiNMF to obtain the features and their corresponding activations. We note that the vectors with superscripts and subscripts denote the row and column vectors of the corresponding matrices respectively.

---

### Algorithm 1 *ssiNMF*

---

```

1: randomly initialize  $\mathbf{G}$  and  $\mathbf{H}$ 
2: normalize columns of  $\mathbf{G}$  to unit  $L_2$  norm
3: repeat
4:   update  $\mathbf{G}$ 
5:   for each column  $i$  in  $\mathbf{G}$  do
6:      $t \leftarrow 0$ 
7:     for each element  $j$  in  $\mathbf{g}_i$  do
8:        $t \leftarrow t + \text{fmvc}(\mathbf{G}, \mathbf{H}\mathbf{h}^{i*m+j}) - \mathbf{X}\mathbf{h}^{i*m+j}$ 
9:     end for
10:     $\mathbf{g}_i = \mathbf{g}_i - \eta t$ 
11:   end for
12:   update  $\mathbf{H}$ 
13:   for each column  $i$  in  $\mathbf{H}$  do
14:      $\mathbf{h}_i = \mathbf{h}_i \odot \text{tfmvc}(\mathbf{G}, \mathbf{x}_i) / (\text{tfmvc}(\mathbf{G}, \text{fmvc}(\mathbf{G}, \mathbf{h}_i)) + \lambda)$ 
15:   end for
16: until convergence

```

---

## 6 Experiments

To test our algorithm we generate the bars dataset as in [4]. As shown in Figure 1(a), its generating features are single and double bars aligned vertically and horizontally on a  $3 \times 3$  grid. Since all double bars can be expressed in terms of the single bars, this feature basis is overcomplete. These features form a generating feature matrix  $W_{gen}$ . Initializing  $H_{gen}$  to a sparse random matrix, we construct the dataset as  $X = W_{gen}H_{gen}$ . 12 random samples from the dataset are shown in Figure 1(b).

As previously demonstrated by Hoyer [4], the addition of sparsity assists in handling overcompleteness of the feature space. This is shown in the feature set learned by non-negative sparse coding in Figure 1(c). However, in the case of allowed translations the original overcomplete set can be represented by only 4 features: vertical, horizontal single bars and corresponding double bars.

We applied ssiNMF by setting the number of features to 4 and  $\lambda = 0.7$ . Each feature in  $\mathbf{W}$  was initialized by iid samples from the uniform distribution and normalized by its  $L_2$  norm. Activations  $\mathbf{H}$  were also randomly initialized from the uniform distribution.

Features identified by ssiNMF algorithm are shown in Figure 1(d). These features still represent an overcomplete basis since the double bar features can be represented in terms of the single bars. Shift-invariance leads to a smaller set of features while still enabling a sparse representation.

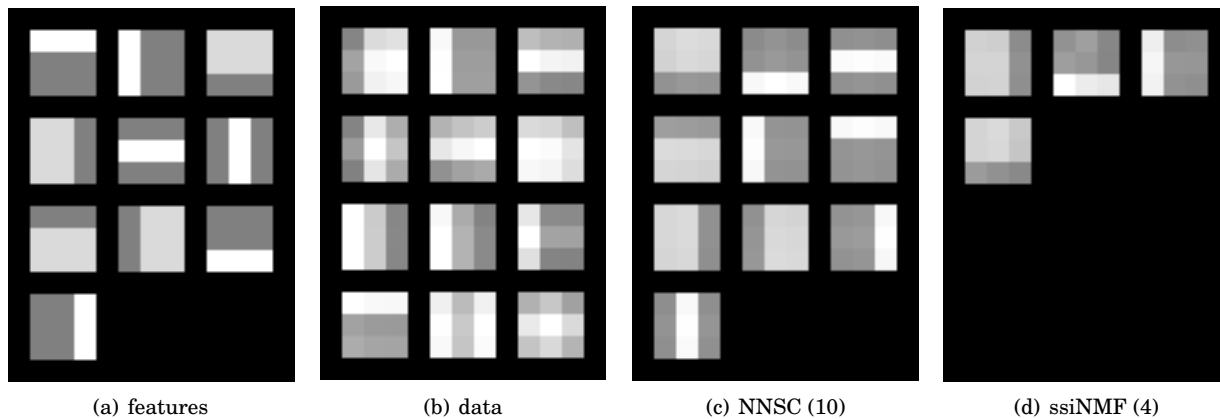


Figure 1. Experimental results on bars dataset. (a) The features used to generate training data. (b) A random sample from the dataset. (c) 10 features as learned by non-negative sparse coding of [4]. (d) 4 features that can represent the data in circulant case as learned by shift-invariant sparse NMF.

## 7 Discussion and Future Work

Circulant constraints make the computation of matrix-vector products fast and reduce storage space in case of dictionaries with shift-invariant features. The gradient descent rule for updating the dictionary matrix  $W$  is additive and it would be interesting to come up with a suitable multiplicative rule.

The algorithm is also potentially useful for datasets which are misaligned. For example, ssiNMF could be applied to a dataset of fMRI images where the head is not stabilized. Shift-invariance in this case would compensate for the motion typically observed in fMRI experiments or for coregistration differences between subjects.

Our approach can also be extended to non-negative tensor factorization (NTF) [5] which is a rich framework with which to model additional factors.

## 8 Acknowledgements

The first author would like to acknowledge the support from NIBIB grants 1 R01 EB 000840 and 1 R01 EB 005846. The second author was supported by NIMH grant 1 R01 MH076282-01. The latter two grants were funded as part of the NSF/NIH Collaborative Research in Computational Neuroscience Program. The authors would like to thank Barak Pearlmutter for the initial idea.

## References

- [1] P. Comon, "Independent component analysis: A new concept," *Signal Processing*, vol. 36, pp. 287–314, 1994.
- [2] Anthony J. Bell and Terrence J. Sejnowski, "An information-maximization approach to blind separation and blind deconvolution," *Neu. Comp.*, vol. 7, no. 6, pp. 1129–59, 1995.
- [3] Paris Smaragdis, "Non-negative matrix factor deconvolution; extraction of multiple sound sources from monophonic inputs," in *Fifth International Conference on Independent Component Analysis*, Granada, Spain, Sept. 22–24 2004, LNCS 3195, pp. 494–9, Springer-Verlag.
- [4] Patrik O. Hoyer, "Non-negative sparse coding," in *IEEE Workshop on Neural Networks for Signal Processing*, 2002.
- [5] A.Cichocki, R.Zdunek, S.Choi, R.Plemmons, and S.Amari, "Non-negative tensor factorization using alpha and beta divergences," *International Conference on Acoustics, Speech and Signal Processing (ICASSP 07)*, Honolulu, Hawaii, USA, 2007.
- [6] Daniel D. Lee and H. Sebastian Seung, "Learning the parts of objects with nonnegative matrix factorization," *Nature*, vol. 401, pp. 788–91, 1999.
- [7] David Donoho and Victoria Stodden, "When does non-negative matrix factorization give a correct decomposition into parts?," in *Adv. in Neu. Info. Proc. Sys. 16*. 2004, MIT Press.