

CONVOLUTIVE NON-NEGATIVE MATRIX FACTORISATION WITH SPARSENESS CONSTRAINT

Paul D. O’Grady

Barak A. Pearlmutter

Hamilton Institute
National University of Ireland, Maynooth
Co. Kildare, Ireland.

ABSTRACT

Discovering a parsimonious representation that reflects the structure of audio is a requirement of many machine learning and signal processing methods. Such a representation can be constructed by Non-negative Matrix Factorisation (NMF), which is a method for finding parts-based representations of non-negative data. We present an extension to NMF that is convolutive and forces a sparseness constraint. Combined with spectral magnitude analysis of audio, this method discovers auditory objects and their associated sparse activation patterns.

1. INTRODUCTION

A preliminary step in many data analysis tasks is to find a suitable representation of the data. Typically, methods exploit the latent structure in the data. For example, ICA [1] reduces the redundancy of the data by projecting the data onto its independent components, which can be discovered by maximising a statistical measure such as independence [2] or non-Gaussianity [3].

Given a non-negative matrix \mathbf{V} , Non-negative Matrix Factorisation (NMF) approximately decomposes \mathbf{V} into a product of two non-negative matrices \mathbf{W} and \mathbf{H} [4, 5]. NMF is a parts-based approach that does not make a statistical assumption. Instead, it assumes that for the domain at hand, negative numbers would be physically meaningless. The lack of statistical assumptions makes it difficult to prove that NMF will give correct decompositions, although it has been shown geometrically that NMF provides a correct decomposition for some classes of images [6].

For data that contains negative components, for example audio, a non-negative representation must be found. In this case a spectrogram representation may be used. Spectrograms have been used in audio analysis for many years [7] and combined with NMF have been applied to variety of problems such as monaural speech separation [8], identification of au-

ditary objects with time-varying spectra [9], and automatic transcription of music [10].

In this paper we combine the previous extension of convolutive NMF [9] with a sparseness constraint [11] and apply it to the analysis of audio. The paper is structured as follows. In Section 2 we present NMF and discuss its performance using experiments on synthetic data. We then present convolutive NMF in Section 3 and discuss its advantages over conventional NMF. In Section 4 we add an additional sparseness constraint to the convolutive NMF objective and present an experiment on music.

2. NON-NEGATIVE MATRIX FACTORISATION

NMF is a linear non-negative approximate factorisation and is formulated as follows. Given a non-negative $M \times N$ matrix $\mathbf{V} \in \mathbb{R}^{\geq 0, M \times N}$ the goal is to approximate \mathbf{V} as a product of two non-negative matrices $\mathbf{W} \in \mathbb{R}^{\geq 0, M \times R}$ and $\mathbf{H} \in \mathbb{R}^{\geq 0, R \times N}$

$$\mathbf{V} \approx \mathbf{W} \cdot \mathbf{H} \quad (1)$$

where $R \leq M$, such that the reconstruction error is minimised. Two NMF algorithms were introduced by Lee and Seung [4, 12], each implementing a different cost function by which the quality of the approximation can be measured. The first cost function presented is the Euclidean distance between \mathbf{V} and $\mathbf{W}\mathbf{H}$, the second is a generalised version of Kullback-Leibler divergence. We will use the latter

$$D(\mathbf{V} \parallel \mathbf{W}, \mathbf{H}) = \left\| \mathbf{V} \otimes \log \frac{\mathbf{V}}{\mathbf{W} \cdot \mathbf{H}} - \mathbf{V} + \mathbf{W} \cdot \mathbf{H} \right\| \quad (2)$$

where \otimes denotes an element-wise (also known as Hadamard or Schur) product and division is also element-wise. NMF can now be written as an optimisation problem.

$$\min_{\mathbf{W}, \mathbf{H}} D(\mathbf{V} \parallel \mathbf{W}, \mathbf{H}) \quad \mathbf{W}, \mathbf{H} \geq 0$$

The above objective is convex in \mathbf{W} and \mathbf{H} individually but not together. Therefore algorithms usually alternate updates of \mathbf{W} and \mathbf{H} . These admit to multiplicative updates, which

Supported by Higher Education Authority of Ireland (An tÚdarás Um Ard-Oideachas) and Science Foundation Ireland grant 00/PL.1/C067.

can be interpreted as diagonally rescaled gradient descent [4]:

$$\mathbf{H} = \mathbf{H} \otimes \frac{\mathbf{W}^T \cdot \mathbf{V}}{\mathbf{W}^T \cdot \mathbf{1}}, \quad \mathbf{W} = \mathbf{W} \otimes \frac{\mathbf{V} \cdot \mathbf{H}^T}{\mathbf{1} \cdot \mathbf{H}^T} \quad (3)$$

where $\mathbf{1}$ is an $M \times N$ matrix with all its elements equal to unity and divisions are element-wise. As the algorithm iterates the factors converge to a local optimum of Eq. 2.

The parameter R , which defines the number of columns in \mathbf{W} and rows in \mathbf{H} , defines the rank of the approximation. If $R < M$ then \mathbf{W} is under-determined and NMF reveals low-rank features of the data. The columns of \mathbf{W} will contain the basis for the data while the rows of \mathbf{H} will contain activation patterns for each basis. The selection of an appropriate value of R is usually based on prior knowledge and is necessary for good approximation.

2.1. NMF applied on audio spectra

To illustrate the application of NMF on audio data consider the example shown in Figure 1. The signal under consideration is composed of two band-limited noise bursts with magnitude spectra constant over time. The first burst is centred around 2 kHz and occurs four times, while the second burst is centred around 6 kHz and occurs three times. The signal's spectrogram is a $M \times N$ matrix \mathbf{V} with magnitude information for M frequency bins at N time intervals. NMF is applied to \mathbf{V} with $R = 2$ and the resultant factors shown. In this example both the frequency spectra of the bursts (columns of \mathbf{W}) and their activations in time (rows of \mathbf{H}) have been identified.

This decomposition has successfully revealed the structure of the \mathbf{V} by correctly describing its constituent elements in both the frequency and time domains.

Now consider the example presented in Figure 2. Here, the signal under consideration is composed of two auditory objects that have differing frequency sweeps over time. The first object is centred around 2 kHz and the second object is centred around 6 kHz, each occurring four times. NMF is applied to the data with the same parameters as above and the factors are shown. It is evident from the columns of \mathbf{W} that the identified spectra contain frequency components that are centred around both 2 kHz and 6 kHz. Thus, NMF fails to identify the spectra of each object and instead discovers objects that are a combination of both. The reason for this is that the spectra of the auditory objects evolve over time and that NMF is not expressive enough to reveal this temporal structure. Therefore, in order to reveal a correct decomposition, the expressive properties of NMF need to be extended to consider the evolution of each object's spectrum.

3. CONVOLUTIVE NMF

Typically, the temporal relationship between multiple observations over nearby intervals of time are discovered using a convolutive generative model. Such a model has previously

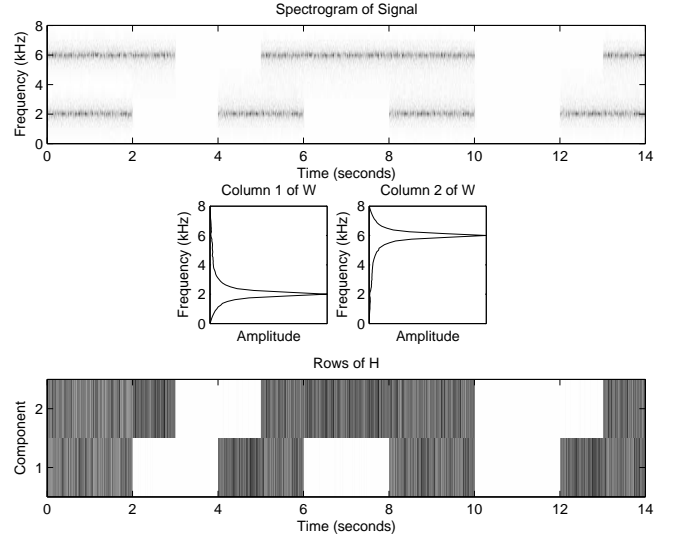


Fig. 1. Spectrogram of a signal composed of band-limited noise bursts, and its factors obtained by NMF.

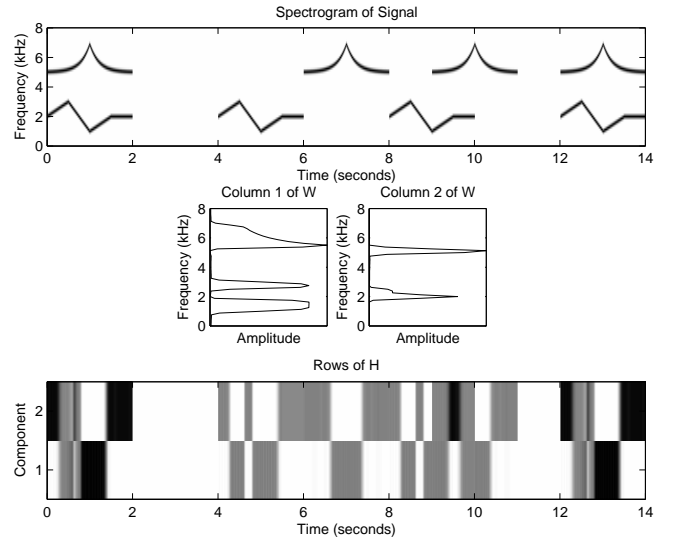


Fig. 2. Spectrogram of a signal composed of auditory objects with time-varying spectra, and its factors obtained by NMF.

been used to extend ICA [13] and NMF [9]. For conventional NMF, each object is described by its spectrum and corresponding activation in time, while for convolutive NMF each object has a sequence of successive spectra and corresponding activation pattern across time. The generative model of Eq. 1 is extended to the convolutive case

$$\mathbf{V} \approx \sum_{t=0}^{T-1} \mathbf{W}_t \cdot \mathbf{H}^{t \rightarrow}$$

where $\mathbf{V} \in \mathbb{R}^{\geq 0, M \times N}$ is the input to be decomposed, $\mathbf{W}_t \in \mathbb{R}^{\geq 0, M \times R}$ and $\mathbf{H} \in \mathbb{R}^{\geq 0, R \times N}$ are its two factors, and T is

the length of each spectrum sequence. The i -th column of \mathbf{W}_t describes the spectrum of the i -th object t time steps after the object has begun. The $\overset{i \rightarrow}{(\cdot)}$ denotes a column shift operator that moves its argument i places to the right, as each column is shifted off to the right the leftmost columns are zero filled. Conversely, the $\overset{\leftarrow i}{(\cdot)}$ operator shifts columns off to the left, with zero filling on the right.

Using the previously presented framework for NMF, the new cost function for the convolutive generative model is

$$D(\mathbf{V}||\mathbf{\Lambda}) = \left\| \mathbf{V} \otimes \log \frac{\mathbf{V}}{\mathbf{\Lambda}} - \mathbf{V} + \mathbf{\Lambda} \right\| \quad (4)$$

where $\mathbf{\Lambda}$ is the approximation to \mathbf{V} and is defined as

$$\mathbf{\Lambda} = \sum_{t=0}^{T-1} \mathbf{W}_t \cdot \overset{t \rightarrow}{\mathbf{H}}$$

This new cost function can be viewed as a set of T conventional NMF operations that are summed to produce the final result. Consequently, as opposed to updating two matrices (\mathbf{W} and \mathbf{H}) as in conventional NMF, $T + 1$ matrices require an update, including all \mathbf{W}_t and \mathbf{H} . The resultant convolutive NMF update equations are

$$\mathbf{H} = \mathbf{H} \otimes \frac{\mathbf{W}_t^T \cdot \overset{\leftarrow t}{[\mathbf{V}]} }{\mathbf{W}_t^T \cdot \mathbf{1}}, \quad \mathbf{W}_t = \mathbf{W}_t \otimes \frac{\mathbf{V} \cdot \overset{t \rightarrow}{\mathbf{H}} }{\mathbf{1} \cdot \mathbf{H}} \quad (5)$$

At each iteration \mathbf{H} and all \mathbf{W}_t are updated, where \mathbf{H} is updated to the average result of its updates for all \mathbf{W}_t [9]. It can easily be seen that for $T = 1$ this reduces to conventional NMF (Eq. 3).

3.1. Convolutive NMF applied on audio spectra

We have shown that conventional NMF reveals a correct decomposition for auditory objects with constant spectra but fails for objects that exhibit time-varying spectra. Now let us consider convolutive NMF applied to this example. The performance of the algorithm now depends on two parameters R and T , where T must be larger than the time each object exists. Convolutive NMF is applied to the data with $R = 2$ and $T = 2$ seconds, and the resultant factors presented in Figure 3. It is evident from spectra sequences obtained (i -th column of \mathbf{W}_t , for $t = 0, 1, \dots, T - 1$) that the time-varying spectra of each object has been revealed and that the rows of \mathbf{H} identify the start of each object. This decomposition has successfully revealed the structure of \mathbf{V} by correctly describing the spectral evolution of each object and its position in time.

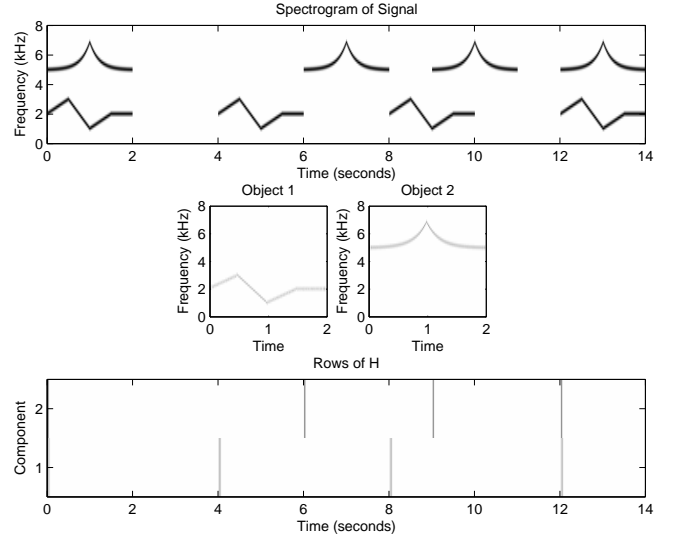


Fig. 3. Spectrogram of a signal composed of auditory objects with time-varying spectra, and its factors obtained by convolutive NMF.

4. CONVOLUTIVE NMF WITH ADDITIONAL CONSTRAINTS

For some tasks it may be advantageous to perform NMF with additional constraints placed on either \mathbf{W} or \mathbf{H} . One increasingly popular and powerful constraint is that the rows of \mathbf{H} have a parsimonious activation pattern for each basis contained in the columns of \mathbf{W} . This is the so called *Sparse-ness Constraint* [14, 15]. A signal is said to be sparse when it is zero or nearly zero more than might be expected from its variance. Such a signal has a probability density function or distribution of values with a sharper peak at zero and fatter tails than a Gaussian. A standard sparse distribution is the Laplacian distribution ($p(c) \propto \exp -|c|$). The advantage of a sparse signal representation is that the probability of two or more activation patterns being active simultaneously is low. Thus, sparse representations lend themselves to good separability [16]. Although convolutive NMF produces activation patterns that tend to be sparse, the addition of the sparseness constraint on \mathbf{H} provides a means of trading off the sparseness of the representation against accurate reconstruction.

The most widely used method for multi-objective optimization is the weighted sum method. This method creates an aggregate objective function by multiplying each constituent cost function by a weighting factor and summing the weighted costs. Combining our reconstruction cost function (Eq. 4) with a sparseness constraint on \mathbf{H} results in the following objective function

$$G(\mathbf{V}||\mathbf{\Lambda}) = D(\mathbf{V}||\mathbf{\Lambda}) + \lambda \sum_{ij} \mathbf{H}_{ij} \quad (6)$$

The left term of the objective function corresponds to NMF,

while the right term is an additional constraint on \mathbf{H} that enforces sparsity by minimising the L_1 -norm of its columns [17]. The parameter λ controls the trade off between sparseness and accurate reconstruction.

This objective creates a new problem: the right term is a strictly increasing function of the absolute value of its argument, so it is possible that the objective can be decreased by scaling up \mathbf{W} and scaling down \mathbf{H} ($\mathbf{W} \mapsto \alpha\mathbf{W}$ and $\mathbf{H} \mapsto (1/\alpha)\mathbf{H}$, with $\alpha > 1$). This situation does not alter the left term in the objective function, but will cause the right term to decrease, resulting in the elements of \mathbf{W} growing without bound and \mathbf{H} tending toward zero. Consequently, the solution arrived at by the optimisation algorithm is not influenced by the right term of the objective function and the resultant \mathbf{H} matrix is not sparse. Therefore another constraint needs to be introduced in order to make the cost function well-defined. This is achieved by fixing the norm of the i -th object of \mathbf{W} (over all $t = [0, 1, \dots, T - 1]$) to unity which constrains the scale of the elements in \mathbf{W} and \mathbf{H} .

4.1. New Update rules

The classic NMF update rules [4] implement gradient descent and our new updates will also follow this approach. First we consider the update for \mathbf{H} , where the gradient descent update is

$$\mathbf{H} = \mathbf{H} + \eta_{\mathbf{H}} \nabla_{\mathbf{H}} G(\mathbf{V} \|\mathbf{\Lambda})$$

Taking the gradient of Eq. 6 with respect to \mathbf{H} gives

$$\nabla_{\mathbf{H}} G(\mathbf{V} \|\mathbf{\Lambda}) = \mathbf{W}_t^T \cdot \left[\frac{\mathbf{V}}{\mathbf{\Lambda}} \right] - \mathbf{W}_t^T \cdot \mathbf{1} + \lambda \cdot \mathbf{1}$$

Diagonally rescaling the variables [4, 11] and setting the learning rate to

$$\eta_{\mathbf{H}} = \frac{\mathbf{H}}{\mathbf{W}_t^T \cdot \mathbf{1} + \lambda \cdot \mathbf{1}}$$

gives the new update rule for \mathbf{H}

$$\mathbf{H} = \mathbf{H} \otimes \frac{\mathbf{W}_t^T \cdot \left[\frac{\mathbf{V}}{\mathbf{\Lambda}} \right]}{\mathbf{W}_t^T \cdot \mathbf{1} + \lambda \cdot \mathbf{1}} \quad (7)$$

Similarly, we define the update for \mathbf{W}

$$\mathbf{W} = \mathbf{W} + \eta_{\mathbf{W}} \nabla_{\mathbf{W}} G(\mathbf{V} \|\mathbf{\Lambda})$$

where the gradient of Eq. 6 with respect to \mathbf{W} is

$$\nabla_{\mathbf{W}} G(\mathbf{V} \|\mathbf{\Lambda}) = \frac{\mathbf{V}}{\mathbf{\Lambda}} \cdot \mathbf{H}^{\top} - \mathbf{1} \cdot \mathbf{H}^{\top}$$

The additional unit norm constraint on \mathbf{W} complicates the update rule and impedes the discovery of a suitable form for $\eta_{\mathbf{W}}$ that would result in a multiplicative update [11], thus resulting in the following update

$$\mathbf{W} = \mathbf{W} + \eta_{\mathbf{W}} \left[\frac{\mathbf{V}}{\mathbf{\Lambda}} \cdot \mathbf{H}^{\top} - \mathbf{1} \cdot \mathbf{H}^{\top} \right] \quad (8)$$

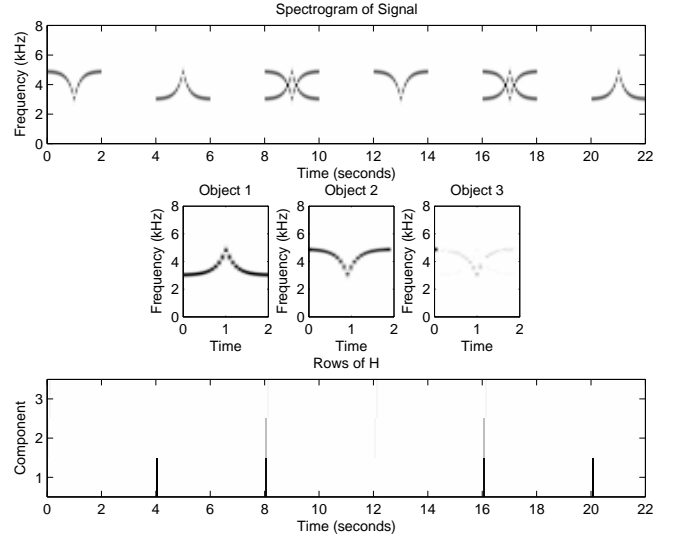


Fig. 4. Spectrogram of a signal composed of an over-complete basis, and its factors obtained by convolutive NMF.

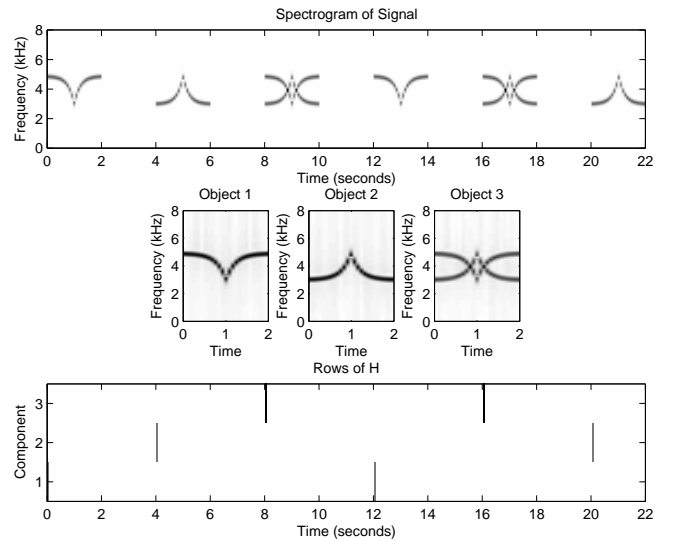


Fig. 5. Spectrogram of a signal composed of an over-complete basis, and its factors obtained by sparse convolutive NMF.

As long as $\eta_{\mathbf{W}}$ is sufficiently small the update should reduce Eq. 6. Subsequent to this update any negative values in \mathbf{W} are set to zero (non-negativity constraint) and each object contained in \mathbf{W} is rescaled to unit norm.

4.2. Sparse Convolutive NMF applied on audio spectra

An interesting property of the sparseness constraint is that it enables the discovery of an over-complete basis *i.e.* a basis that contains more basis functions than are necessary to span the

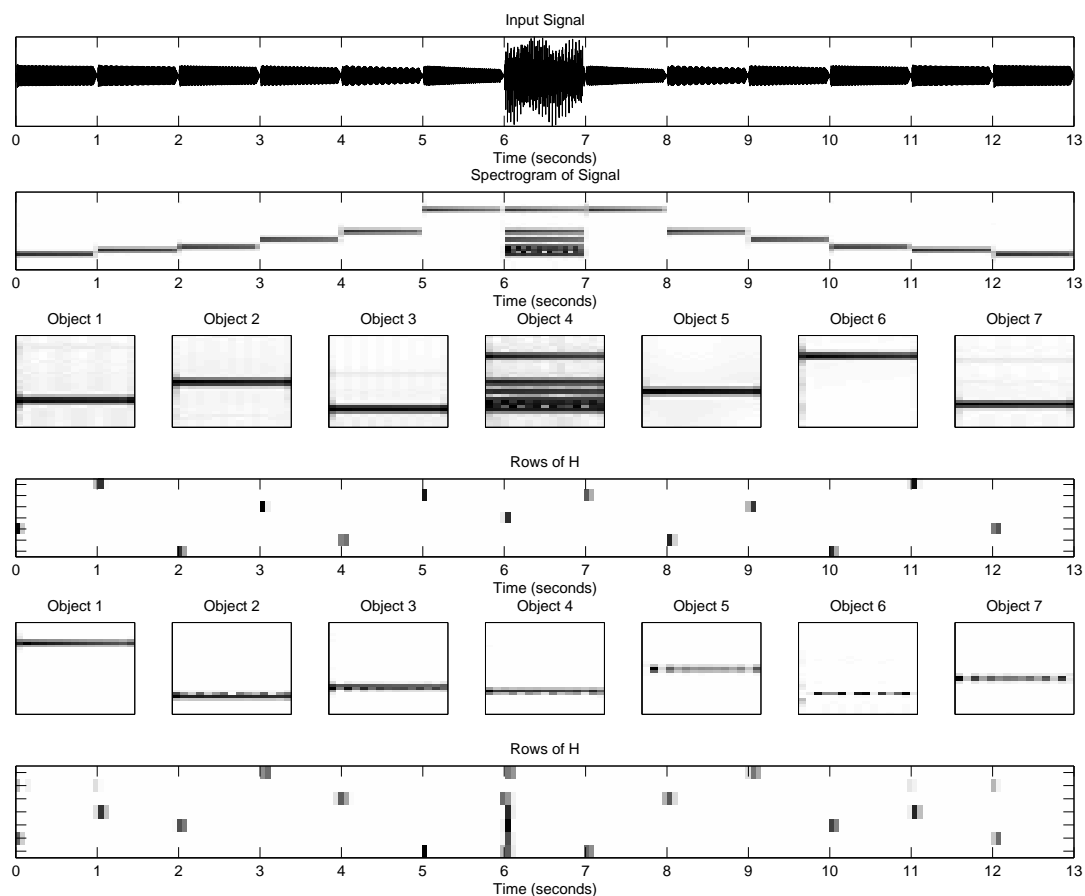


Fig. 6. Music waveform and its associated spectrogram along with its factors obtained by sparse convolutive NMF (rows 3 & 4) and conventional convolutive NMF (rows 5 & 6).

projection space.

To illustrate the performance of convolutive NMF on data generated from an over-complete basis consider the example presented in Figure 4. The signal under consideration is composed of three auditory objects each occurring twice, where the first object is an exponentially decreasing then increasing frequency sweep centred around 4 kHz, the second object is the reverse of the first, and the third object is a combination of the first two. Convolutive NMF is applied to the data with $R = 3$ and $T = 2$ seconds, and the resultant factors presented. It is evident from the results that only the first two auditory objects are identified. This is because the third object can be expressed in term of the first two and the signal can be described by using just the first two objects. Thus, convolutive NMF achieves its optimum with just the first two linearly independent objects without the need for an over-complete representation.

When the sparseness constraint is introduced to the objective the existence of an over-complete representation helps minimise the objective and allows for a sparser description of the signal. Sparse convolutive NMF is applied to the same

signal (Figure 5). Here, all three objects and their associated activation patterns are identified. Therefore, this decomposition has successfully revealed the over-complete basis used to generate the signal.

4.3. Sparse Convolutive NMF applied on music

To see the performance of sparse convolutive NMF in a real-world context, we apply it to a simple music example. The data consists of synthesised rudimentary guitar sounds, where each string produces only its fundamental frequency. The arrangement is simple, composed of three sections: the six notes of a G chord are played individually in descending order; all six notes of the chord are played simultaneously; and each note is played in reverse order. Each note is played for one second, and the frequencies of the notes are 98.00 Hz (G), 123.47 Hz (B), 146.83 Hz (D), 196.00 Hz (G), 246.94 Hz (B) and 392.00 Hz (G).

Both sparse convolutive NMF and convolutive NMF are applied to the music and the resultant factors are presented in Figure 6. It is evident from the spectrogram that the music can

be represented by an over-complete representation consisting of each individual note and the chord. Convolutional NMF is applied with $R = 7$, $T = 1$ second and the resultant factors are presented in rows 5 & 6. As can be seen from the activation pattern, the algorithm has failed to represent the chord as an individual auditory object and instead represents it as a combination of notes. Sparse convolutional NMF is applied with the same parameters above and with λ selected on an *ad hoc* basis. The resultant factors are presented in rows 3 & 4. Here, it is evident that an over-complete representation is discovered and that the chord is represented as an individual auditory object.

5. CONCLUSIONS

In this paper we have presented a sparse convolutional version of NMF that effectively discovers a sparse parts based representation for non-negative data. This method extends the convolutional NMF objective by including a sparseness constraint on the activation patterns, enabling the discovery of over-complete representations. We have shown how the expressive properties of NMF can be improved by reformulation of the problem in a convolutional framework and how the addition of a sparseness constraint can lead to the discovery of over-complete representations in music.

References

- [1] P. Comon. Independent component analysis: A new concept. *Signal Processing*, 36:287–314, 1994.
- [2] A. J. Bell and T. J. Sejnowski. An information-maximization approach to blind separation and blind deconvolution. *Neu. Comp.*, 7(6):1129–1159, 1995.
- [3] A. Hyvärinen and E. Oja. A fast fixed-point algorithm for independent component analysis. *Neu. Comp.*, 9(7):1483–1492, Oct. 1997.
- [4] D. D. Lee and H. S. Seung. Algorithms for non-negative matrix factorization. In *Adv. in Neu. Info. Proc. Sys. 13*, pages 556–562. MIT Press, 2001.
- [5] P. Paatero and U. Tapper. Positive matrix factorization: A nonnegative factor model with optimal utilization of error estimates of data values. *Environmetrics*, 5:111–126, 1994.
- [6] D. Donoho and V. Stodden. When does non-negative matrix factorization give a correct decomposition into parts? In *Adv. in Neu. Info. Proc. Sys. 16*. MIT Press, 2004.
- [7] R. K. Potter, G. A. Kopp, and H. C. Green. *Visible Speech*. D. Van Nostrand Company, 1947.
- [8] T. Virtanen. Sound source separation using sparse coding with temporal continuity objective. In *Proceedings of the International Computer Music Conference (ICMC 2003)*, 2003.
- [9] P. Smaragdis. Non-negative matrix factor deconvolution; extraction of multiple sound sources from monophonic inputs. In *Fifth International Conference on Independent Component Analysis*, LNCS 3195, pages 494–499, Granada, Spain, Sep. 22–24 2004. Springer-Verlag.
- [10] S. A. Abdallah and M. D. Plumbley. Polyphonic transcription by non-negative sparse coding of power spectra. In *Proceedings of the 5th International Conference on Music Information Retrieval (ISMIR 2004)*, pages 318–325, 2004.
- [11] P. O. Hoyer. Non-negative sparse coding. In *IEEE Workshop on Neural Networks for Signal Processing*, 2002.
- [12] D. D. Lee and H. S. Seung. Learning the parts of objects with nonnegative matrix factorization. *Nature*, 401:788–791, 1999.
- [13] R. H. Lambert. *Multichannel Blind Deconvolution: FIR Matrix Algebra and Separation of Multipath Mixtures*. PhD thesis, Univ. of Southern California, 1996.
- [14] B. A. Olshausen and D. J. Field. Sparse coding of sensory inputs. *Curr Opin Neurobiol.*, 14(4):481–487, 2004.
- [15] D. J. Field. What is the goal of sensory coding? *Neural Computation*, 6:559–601, 1994.
- [16] M. Zibulevsky and B. A. Pearlmutter. Blind source separation by sparse decomposition in a signal dictionary. *Neu. Comp.*, 13(4):863–882, 2001.
- [17] S. S. Chen, D. L. Donoho, and M. A. Saunders. Atomic decomposition by basis pursuit. *SIAM Journal on Scientific Computing*, 20(1):33–61, 1998.