

2006-01-01

Shifted 2D Non-negative Tensor Factorisation

Derry Fitzgerald
Cork Institute of Technology

Matt Cranitch
Technological University Dublin, matt.cranitch@cit.ie

Eugene Coyle
Technological University Dublin, Eugene.Coyle@tudublin.ie

Follow this and additional works at: <https://arrow.tudublin.ie/argcon>

Recommended Citation

Fitzgerald, D., Cranitch, M. & Coyle, E. (2006) Shifted 2D non-negative tensor factorisation. Proceedings of the *Irish Signals and Systems Conference, Dublin, Ireland, 2006*.

This Conference Paper is brought to you for free and open access by the Audio Research Group at ARROW@TU Dublin. It has been accepted for inclusion in Conference papers by an authorized administrator of ARROW@TU Dublin. For more information, please contact arrow.admin@tudublin.ie, aisling.coyne@tudublin.ie, vera.kilshaw@tudublin.ie.

Shifted 2D Non-negative Tensor Factorisation

Derry FitzGerald, Matt Cranitch[♠] and Eugene Coyle*

[♠]*Dept. of Electronic Engineering,
Cork Institute of Technology,
IRELAND
E-mail: [♠]derry.fitzgerald@cit.ie*

**School of Control Systems and Electrical
Engineering,
Dublin Institute of Technology,
IRELAND
E-mail: *eugene.coyle@dit.ie*

Recently, Non-negative Matrix Factor 2D Deconvolution was developed as a means of separating harmonic instruments from single channel mixtures. This technique uses a model which is convolutive in both time and frequency, and so can capture instruments which have both time-varying spectra and time-varying fundamental frequencies simultaneously. However, in many cases two or more channels are available, in which case it would be advantageous to have a multi-channel version of the algorithm. To this end, a shifted 2D Non-negative Tensor Factorisation algorithm is derived, which extends Non-negative Matrix Factor 2D Deconvolution to the multi-channel case. The use of this algorithm for multi-channel sound source separation of pitched instruments is demonstrated.

Keywords – Non-negative tensor factorisation, sound source separation

I INTRODUCTION

In recent years, matrix factorisation techniques such as non-negative matrix factorisation (NMF) [1] have been used to attempt single channel sound source separation [2]. These techniques attempt to approximate a magnitude spectrogram \mathbf{X} as the product of low rank matrices \mathbf{A} and \mathbf{S} , i.e. $\mathbf{X} \approx \mathbf{AS}$. The columns of \mathbf{A} contain frequency basis functions, while the associated rows of \mathbf{S} contain corresponding amplitude envelopes for the frequency basis functions. Individual elements of \mathbf{A} and \mathbf{S} can then be used to attempt resynthesis of individual components or sources in the input data. Using a cost function which encourages sparseness in \mathbf{A} and \mathbf{S} results in a factorisation where the basis functions in \mathbf{A} and \mathbf{S} correspond to perceptually meaningful features, such as the frequency spectra of individual notes and their associated amplitude envelopes. Ensuring non-negativity is useful in obtaining the factorisation as a magnitude spectrogram is by definition non-negative, and it also reflects the intuition that sound sources add together.

However, these techniques are limited in that the decomposition is linear, and so each basis function pair typically corresponds to a single note played by a given pitched instrument. Therefore, to use these

techniques for sound source separation, the basis functions must be grouped by instrument or source. While grouping techniques have been developed, it has proved difficult to obtain good clustering in many situations [3]. To overcome this problem shifted non-negative matrix factorisation was developed, which models notes played by an instrument as translations of a single instrument basis function [4]. This necessitates the use of a time-frequency resolution with log-frequency resolution, such as the Constant Q transform [5]. If the centre frequencies are set so that $f_i = f_{i-1}2^{1/12}$ where f_i is the centre frequency of band i , then the spacing between centre frequencies matches that of the even-tempered tuning system. Therefore, translating a frequency basis function of a note up by one bin is equivalent to a rise in pitch of one semitone.

A further problem with matrix factorisation techniques is that they return fixed spectra for each note, whereas the spectra of real instruments evolve over time. To this end, convolutive forms of NMF have been developed which model sources as a sequence of successive spectra and a corresponding amplitude envelope which is translated across time to activate each successive spectrum [6].

Recently, these methods have been combined in an attempt to model a source or pitched instrument as translations of successive spectra in both frequency and time, thereby allowing time-varying spectra and fundamental frequencies. This leads to a more realistic model of the sources present. This technique, called Non-negative Matrix Factor 2D Deconvolution (NM2D) has been used to separate mixtures of single channel instruments [7].

All of the above techniques work on single channel mixtures, however, most recordings of popular music from the past 40 years are stereo or two channel recordings. These two channel recordings are typically created by linear mixing of single channel recordings of individual instruments, with the only difference between each channel for a given instrument lying in the gain of the instrument in each channel. Therefore, the same model or set of basis functions can be used to describe a given instrument in either channel. Further, this gain difference is a source of extra information which can be used to aid the separation process. It can be seen then that extending NMF-based techniques to the multi-channel case would be advantageous, firstly due to the widespread use of 2 channel recordings, and secondly, as it provides extra information to aid the sound source separation process. To this end, techniques such as non-negative tensor factorisation (NTF) and shifted NTF have been proposed to deal with the multi-channel case [8],[9]. For sound source separation of multi-channel recordings, the NTF model can be written as:

$$\mathcal{X} \approx \hat{\mathcal{X}} = \sum_{k=1}^K \mathbf{G}_{:k} \circ \mathbf{A}_{:k} \circ \mathbf{S}_{:k} \quad (1)$$

where \mathcal{X} is a $r \times n \times m$ tensor containing the spectrograms of the r channels, containing n frequency bins and m time frames. $\hat{\mathcal{X}}$ is an approximation to \mathcal{X} , \mathbf{G} is a $r \times K$ matrix containing the gains of each factor in each channel, \mathbf{A} is an $n \times K$ matrix containing the frequency basis functions, and \mathbf{S} is an $m \times K$ matrix containing the amplitude basis functions, \circ denotes outer product multiplication, and $:k$ denotes the k th column of a given matrix. The tensor factorisation is obtained by minimising the generalised Kullback-Liebler divergence between \mathcal{X} and $\hat{\mathcal{X}}$. This is defined as:

$$D(\mathcal{X} \parallel \hat{\mathcal{X}}) = \sum_{i,j,l} \mathcal{X} \log \frac{\mathcal{X}}{\hat{\mathcal{X}}} - \mathcal{X} + \hat{\mathcal{X}} \quad (2)$$

where l , i , and j index over channel, frequency bin and time frame respectively.

For the rest of the paper, the following conventions, in line with those adapted by Bader and Kolda in [10] are used. Tensors are denoted by upper case letters such as \mathcal{X} , and contracted tensor product multiplication is defined as follows. If \mathcal{W} is a tensor

of size $I_1 \times \dots \times I_N \times J_1 \times \dots \times J_M$ and \mathcal{Y} is a tensor of size $I_1 \times \dots \times I_N \times K_1 \times \dots \times K_P$ then contracted tensor multiplication along the first N modes is given as:

$$\begin{aligned} & \langle \mathcal{W} \mathcal{Y} \rangle_{\{1:N, 1:N\}}(j_1, \dots, j_M, k_1, \dots, k_P) \\ &= \sum_{i_1=1}^{I_1} \dots \sum_{i_N=1}^{I_N} \mathcal{W}(i_1, \dots, i_N, j_1, \dots, j_M) \\ & \quad \mathcal{Y}(i_1, \dots, i_N, k_1, \dots, k_P) \end{aligned} \quad (3)$$

where element indexing occurs within $()$ brackets, and where the modes to be multiplied are specified in the subscripts within the angle brackets.

II SHIFTED 2D NON-NEGATIVE TENSOR FACTORISATION

This research aims to extend NM2D to the multi-channel case. To do so requires the extension of the basic NTF model to allow for translations of the underlying sources in both frequency, and time. Where NM2D carries out translations by means of a shift operator, in line with our previous work, we carry out shifting by means of translation tensors. For an $n \times 1$ vector, an n by n translation matrix is required. This translation matrix can be easily obtained by permuting the columns of the identity matrix. For example, in order to shift a vector by one position, the required matrix is $\mathbf{I}(:, [n, 1:n-1])$, where \mathbf{I} is the identity matrix and the ordering of the columns is defined in the square brackets. For z translations, each of the z translation matrices can then be grouped into a translation tensor, \mathcal{Q} of size $n \times z \times n$.

Separate translation tensors are defined to deal with shifts in frequency and across time respectively, resulting in the following model for shifted 2D Non-negative Tensor Factorisation (SNTF):

$$\hat{\mathcal{X}} = \sum_{k=1}^K \left\langle \mathcal{G}_{:k} \left\langle \langle \mathcal{T} \mathcal{A}_{:k} \rangle_{\{3,1\}} \langle \mathcal{S}_{:k} \mathcal{P} \rangle_{\{3,1\}} \right\rangle_{\{2,4,1,3\}} \right\rangle_{\{2,2\}} \quad (4)$$

where $\hat{\mathcal{X}}$ is a tensor of size $r \times n \times m$ which is an approximation to \mathcal{X} , and \mathcal{G} is a tensor of size $r \times K$, containing the gains of each of the k sources in each of the r channels. \mathcal{T} is an $n \times z \times n$ translation tensor, which translates the frequency basis functions in \mathcal{A} up or down in frequency, thereby approximating different notes played by a given source. \mathcal{A} is a tensor of size $n \times K \times p$, where p is the number of translations across time. \mathcal{S} is a tensor of size $z \times K \times m$ and \mathcal{P} is a translation tensor of size $m \times p \times m$, which translates the amplitude envelopes contained in \mathcal{S} across time, thereby allowing time-varying source spectra. For simplicity of notation, we adopt the convention that $:k$ denotes the tensor slice associated with the k^{th} source, with the singleton dimension included in the size of the slice.

$\mathcal{A}(\cdot, d)_k$ can then be interpreted as the frequency spectrum associated with the d^{th} translation in time of the k^{th} source. Similarly $\mathcal{S}(e, \cdot)_k$ can be interpreted as the amplitude envelope or time basis function associated with the e^{th} translation in frequency of the k^{th} source.

The required tensor factorisation is obtained again using Eqn. (2) as a cost function. By eliminating terms in \mathcal{X} which are constant, this reduces to:

$$D(\mathcal{X} \parallel \hat{\mathcal{X}}) = \sum_{i,j,l} -\mathcal{X} \log \hat{\mathcal{X}} + \hat{\mathcal{X}} \quad (5)$$

Substituting for Eqn. (4) and taking the gradient with respect to \mathcal{G}_k yields the following update equation:

$$\mathcal{G}_k = \mathcal{G}_k + \lambda^* \left[\begin{array}{c} \left\langle \left\langle \left\langle \left\langle \mathcal{DT} \right\rangle_{\{2,1\}} \mathcal{A}_k \right\rangle_{\{4,1\}} \mathcal{S}_k \right\rangle_{\{3,4,1,2\}} \mathcal{P} \right\rangle_{\{2,4,3,1\}} \\ \left\langle \left\langle \left\langle \left\langle \mathcal{OT} \right\rangle_{\{2,1\}} \mathcal{A}_k \right\rangle_{\{4,1\}} \mathcal{S}_k \right\rangle_{\{3,4,1,2\}} \mathcal{P} \right\rangle_{\{2,4,3,1\}} \end{array} \right] \quad (6)$$

where $\mathcal{D} = \mathcal{X} ./ \hat{\mathcal{X}}$, \mathcal{O} is an all-ones tensor of size equal to \mathcal{X} , and $./$ denotes elementwise division. This can be converted to a multiplicative update rule by setting λ equal to:

$$\lambda = \mathcal{G}_k ./ \left\langle \left\langle \left\langle \left\langle \mathcal{OT} \right\rangle_{\{2,1\}} \mathcal{A}_k \right\rangle_{\{4,1\}} \mathcal{S}_k \right\rangle_{\{3,4,1,2\}} \mathcal{P} \right\rangle_{\{2,4,3,1\}}$$

This yields:

$$\mathcal{G}_k = \mathcal{G}_k \cdot^* \left[\begin{array}{c} \left\langle \left\langle \left\langle \left\langle \mathcal{DT} \right\rangle_{\{2,1\}} \mathcal{A}_k \right\rangle_{\{4,1\}} \mathcal{S}_k \right\rangle_{\{3,4,1,2\}} \mathcal{P} \right\rangle_{\{2,4,3,1\}} \\ \left\langle \left\langle \left\langle \left\langle \mathcal{OT} \right\rangle_{\{2,1\}} \mathcal{A}_k \right\rangle_{\{4,1\}} \mathcal{S}_k \right\rangle_{\{3,4,1,2\}} \mathcal{P} \right\rangle_{\{2,4,3,1\}} \end{array} \right] ./$$

Similarly, update rules can be derived for \mathcal{A}_k and \mathcal{S}_k . These are given by:

$$\mathcal{A}_k = \mathcal{A}_k \cdot^* \left[\begin{array}{c} \left\langle \left\langle \left\langle \left\langle \mathcal{G}_k \circ \mathcal{T} \right\rangle_{\{1,3,1,2\}} \mathcal{D} \right\rangle_{\{3,1\}} \mathcal{S}_k \mathcal{P} \right\rangle_{\{1,2,4\}, \{2,1,4\}} \right\rangle_{\{1,2,4\}, \{2,1,4\}} \\ \left\langle \left\langle \left\langle \left\langle \mathcal{G}_k \circ \mathcal{T} \right\rangle_{\{1,3,1,2\}} \mathcal{O} \right\rangle_{\{3,1\}} \mathcal{S}_k \mathcal{P} \right\rangle_{\{1,2,4\}, \{2,1,4\}} \right\rangle_{\{1,2,4\}, \{2,1,4\}} \end{array} \right] ./$$

$$\mathcal{S}_k = \mathcal{S}_k \cdot^* \left[\begin{array}{c} \left\langle \left\langle \left\langle \left\langle \mathcal{G}_k \circ \mathcal{T} \right\rangle_{\{2,5\}, \{2,1\}} \mathcal{A}_k \right\rangle_{\{1,2,1,2\}} \mathcal{D} \right\rangle_{\{2,3\}, \{2,1\}} \right\rangle_{\{2,3\}, \{2,1\}} \\ \left\langle \left\langle \left\langle \left\langle \mathcal{G}_k \circ \mathcal{T} \right\rangle_{\{2,5\}, \{2,1\}} \mathcal{O} \right\rangle_{\{1,2,1,2\}} \mathcal{A}_k \right\rangle_{\{2,3\}, \{2,1\}} \right\rangle_{\{2,3\}, \{2,1\}} \end{array} \right] ./$$

Once \mathcal{G}_k , \mathcal{A}_k and \mathcal{S}_k are randomly initialised to positive values, non-negativity of the factorisation is guaranteed through the use of multiplicative updates. The algorithm was implemented in Matlab, using the tensor classes for Matlab available at [11].

III SOUND SOURCE SEPARATION USING SNTF

To demonstrate the use of SNTF for the purposes of sound source separation, a 2 channel mixture of flute, piano and trumpet was made using sampled instruments. A separate sample was used for each note of each instrument to make the test as realistic as possible. The flute was panned to mid-left, the piano to the centre and the trumpet to mid-right. Figures 1 to 3 show the Constant Q spectrograms of the flute, piano and trumpet signals respectively, while Figures 4 and 5 show the left channel and right channel mixtures respectively.

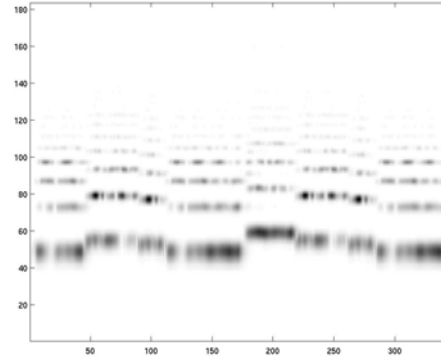


Figure 1: Spectrogram of flute signal

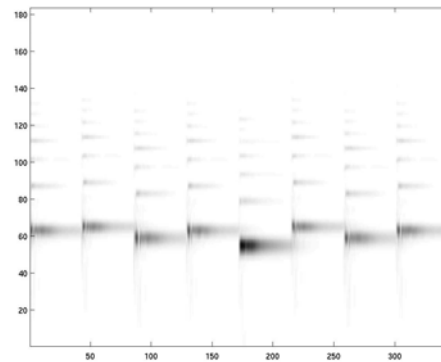


Figure 2: Spectrogram of piano signal

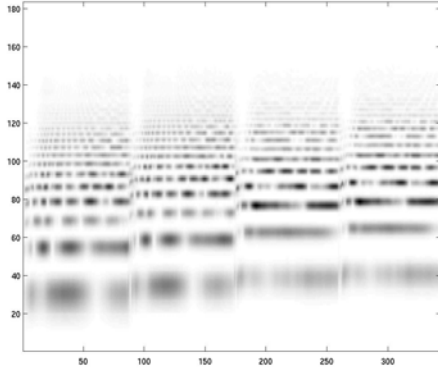


Figure 3: Spectrogram of trumpet signal

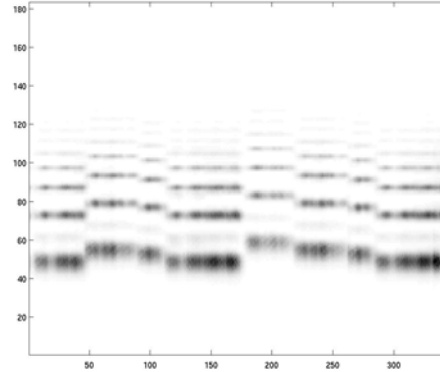


Figure 6: Separated spectrogram of flute.

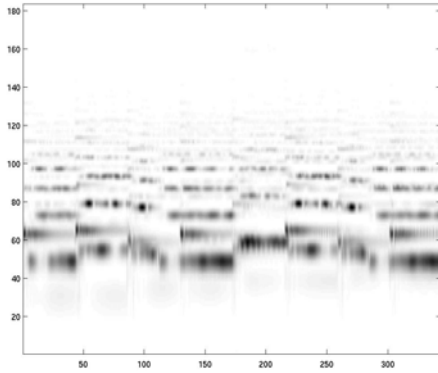


Figure 4: Left channel mixture spectrogram of flute, piano and trumpet.

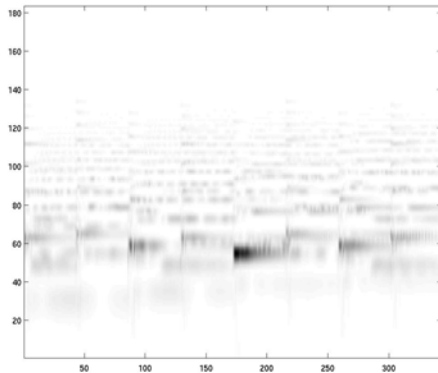


Figure 5: Right channel mixture spectrogram of flute, piano and trumpet.

Figures 6, 7 and 8 show the separated spectrograms of flute, piano and trumpet obtained by applying SNTF to the mixture spectrograms. The range of frequency translations was set to 7 frequency bins, corresponding to a pitch range of 7 semitones and the range of time translations set to 15 time frames, which corresponds to a maximum allowable shift in time of 0.18 seconds when using a hopsize of 512 samples at a sample rate of 44.1 kHz. The number of sources set to 3. Using these parameters, the SNTF algorithm succeeded in separating the instruments. This can be seen in the spectrograms of the separated instruments, with only the notes of each instrument visible in the separated spectrograms. The harmonic structure of the instruments has been successfully recovered, though some of the finer detail has been lost in the temporal evolution of the sources. This demonstrates that the SNTF algorithm is capable of separating mixtures of pitched instruments.

On listening to the separated waveforms, the separated instruments clearly predominate, with little or no traces of the other instruments to be heard in the separated signals. Indeed, the principal artefacts in the resynthesis is as a result of the approximate nature of mapping the log-frequency spectrograms back to linear-frequency spectrograms to allow for resynthesis. Methods of overcoming this problem are currently being researched [12]. A further shortcoming, inherent in all matrix factorisation and tensor factorisation algorithms, is that the algorithm is sensitive to the choices of shift ranges for both frequency and time, though we have observed that this is less of a problem for time shifts than for frequency shifts.

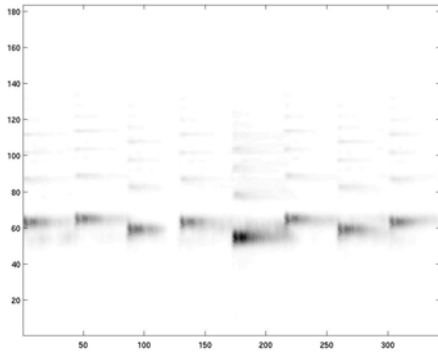


Figure 7: Separated spectrogram of piano.

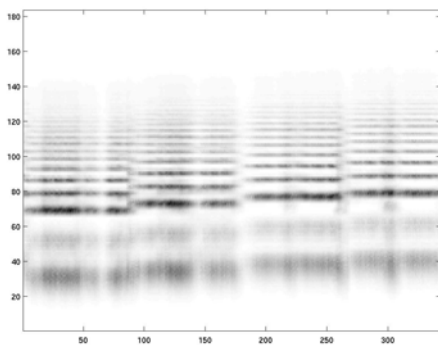


Figure 8: Separated spectrogram of trumpet.

IV CONCLUSIONS

A shifted 2D non-negative tensor factorisation algorithm has been derived, which extends non-negative matrix factor 2D deconvolution to the multi-channel case. The resultant algorithm is capable of modelling shifts in both time and frequency, thereby modelling sources as a group of successive spectra which can be shifted in frequency to approximate different notes played by the sources. The algorithm takes advantage of the spatial information available when dealing with multi-channel mixtures to improve source separation. The algorithm is demonstrated to be capable of separating simple mixtures of pitched instruments. Future work will concentrate on improving the separation capabilities of the algorithm by incorporating models of the instruments to be separated.

V ACKNOWLEDGEMENTS

This research was funded by the Irish Research Council for Science, Engineering and Technology.

VI REFERENCES

- [1] D. Lee, and H. Seung, "Algorithms for non-negative matrix factorization." *Adv. Neural Info. Proc. Syst.* 13, 556-562 (2001).
- [2] P. Smaragdis, J.C. Brown, "Non-negative Matrix Factorization for Polyphonic Music Transcription", *IEEE Workshop on Applications of Signal Processing to Audio and Acoustics (WASPAA)*, pp. 177-180, October 2003.
- [3] D. FitzGerald, "Automatic Drum Transcription and Source Separation", PhD. Thesis, Dublin Institute of Technology, 2004.
- [4] D. FitzGerald, M. Cranitch, and E. Coyle, "Shifted Non-negative Matrix Factorisation for Sound Source Separation", *Proceedings of the IEEE Conference on Statistics in Signal Processing, Bordeaux, France, 2005*.
- [5] J. Brown, "Calculation of a Constant Q Spectral Transform", *J. Acoust. Soc. Am.* 89 425-434, 1991
- [6] T. Virtanen, "Separation of sound sources by convolutive sparse coding", *Proceedings of the ISCA Tutorial and Research Workshop on Statistical and Perceptual Audio Processing, 2004*
- [7] M. Schmidt and M. Morup, "Nonnegative Matrix Factor 2D Deconvolution for Blind Single Channel Source Separation", *Proceedings of the International Conference on Independent Component Analysis 2006*.
- [8] D. FitzGerald, M. Cranitch and E. Coyle, "Non-negative Tensor Factorisation for Sound Source Separation", *Proceedings of the Irish Signals and Systems Conference, Dublin 2005*.
- [9] D. FitzGerald, M. Cranitch and E. Coyle, "Sound Source Separation using Shifted Non-negative Tensor Factorisation", *Proceedings of the IEE Conference on Audio and Speech Signal Processing (ICASSP), Toulouse, 2006*.
- [10] B. Bader and T. Kolda, "MATLAB Tensor Classes for Fast Algorithm Prototyping", *Technical Report SAND2004-5187, Sandia National Laboratories, Livermore, California, 2004*.
- [11] Tensor Classes for Matlab, available at: <http://csmr.ca.sandia.gov/tgkolda/>
- [12] D. FitzGerald, M. Cranitch, and M. Cychowski, "Towards an Inverse Constant Q Transform", *120th AES Convention, Paris, 2006*.