

Research Article

Blind Audio Source Separation with Sparse Nonnegative Matrix Factorization

Abd Majid Darsono, N.Z. Haron, Shakir Saat, M.M. Ibrahim and N.A. Manap

Faculty of Electronics and Computer Engineering, Universiti Teknikal Malaysia Melaka Hang Tuah Jaya, 76100 Durian Tunggal, Melaka, Malaysia

Abstract: In this study, a new technique in source separation using Two-Dimensional Nonnegative Matrix Factorization (NMF2D) with the Beta-divergence is proposed. The Time-Frequency (TF) profile of each source is modeled as two-dimensional convolution of the temporal code and the spectral basis. In addition, adaptive sparsity constraint was imposed to reduce the ambiguity and provide uniqueness to the solution. The proposed model used Beta-divergence as a cost function and updated by maximizing the joint probability of the mixing spectral basis and temporal codes using the multiplicative update rules. Experimental tests have been conducted in audio application to blindly separate the source in musical mixture. Results have shown the effectiveness of the algorithm in separating the audio sources from single channel mixture.

Keywords: Beta divergence, blind audio source separation, machine learning, nonnegative matrix factorization

INTRODUCTION

Nonnegative Matrix Factorization (NMF) (Lee and Seung, 1999) has become one of the promising and exciting techniques in signal processing. NMF has been successfully applied in various applications such as in automatic music transcription (FitzGerald, 2004) cryptography (Xie *et al.*, 2008), pattern recognition (Biciu *et al.*, 2007) and etc. One of the most useful property of NMF is that the nonnegative constraint by itself enforcing the sparse representation of the data. This representation makes the encoded data easy to be estimated because data was encoded by using only a few active components. In NMF, given the matrix, Y of a dimension of $F \times N$ with nonnegative elements, Nonnegative Matrix Factorization (NMF) is the problem of approximate the factorization:

$$Y \approx WH \quad (1)$$

where, $W \in \mathbb{R}^{F \times C}$ and $H \in \mathbb{R}^{C \times N}$ are a non-negative matrices. F represents the frequency bins while N represents the time slot in the TF domain. W contains the spectral basis vectors while H represents the amplitude of each basis vector at each time point. C is the numbers of component from data sources being used and it is determine such that $FC + CN \ll FN$ so that the data can be compressed to its integral component. This problem can be formulated as the minimization of an objective function:

$$D(Y | WH) = \sum_{f,n} d \left(Y_{f,n} \left| \sum_c W_{f,c} H_{c,n} \right. \right) \quad (2)$$

where, d is a scalar divergence. common way to measure how close Y and WH are to use a so-called Beta-divergence (Kompass, 2005; Fevotte *et al.*, 2009), defined by:

$$d_\beta(y|x) = \begin{cases} \frac{y^\beta}{\beta(\beta-1)} + \frac{x^\beta}{\beta} - \frac{yx^{\beta-1}}{\beta-1} & \beta \in \mathbb{R} \setminus \{0,1\} \\ y(\log y - \log x) + (x - y) & \beta = 1 \\ \frac{y}{x} - \log \frac{y}{x} - 1 & \beta = 0 \end{cases} \quad (3)$$

The limiting cases $\beta = 0, 1$ and 2 correspond to the Itakura-Saito (IS) divergence, the Kullback-Leibler (KL) divergence and Least Square (LS) distance, respectively. Particularly it underlies the multiplicative Gamma observation noise, Poisson noise and Gaussian additive observation noise, respectively.

The recent developed Two Dimensional NMF (NMF2D) model (Morup and Schmidt, 2006) extends the NMF model in order to provide decomposition that can capture the temporal dependency of the frequency patterns within the source efficiently. In NMF2D, the Time-Frequency (TF) profile of each source is modeled as two-dimensional convolution of the temporal code and the spectral basis. This significantly reduces the number of components per source needed in the decomposition. So far, for NMF2D, there is no research work has been done to apply the general framework of Beta-divergence. This study carried out to accommodate the Beta-divergence cost function in NMF2D model and investigate the effect of β in the

Corresponding Author: Abd Majid Darsono, Faculty of Electronics and Computer Engineering, Universiti Teknikal Malaysia Melaka Hang Tuah Jaya, 76100 Durian Tunggal, Melaka, Malaysia

This work is licensed under a Creative Commons Attribution 4.0 International License (URL: <http://creativecommons.org/licenses/by/4.0/>).

performance of the algorithm. To further improve the algorithm, this study proposed a sparseness constraint to be imposed in the cost function to reduce the ambiguity the ambiguity associated with the estimation of the spectral basis and temporal codes.

METHODOLOGY

Two-dimensional nonnegative matrix factorization: In derivation of nonnegative matrix factorization framework, firstly, we considered a source model of \mathbf{Y} which is defined as a follows:

$$\mathbf{Y} \approx \sum_{\tau=0}^{\tau_{\max}} \sum_{\phi=0}^{\phi_{\max}} \mathbf{W}^{\tau} \mathbf{H}^{\phi} \approx \sum_{\tau=0}^{\tau_{\max}} \sum_{\phi=0}^{\phi_{\max}} \left(\sum_{j=1}^J \mathbf{W}_j^{\tau} \mathbf{H}_j^{\phi} \right) \quad (4)$$

where, J is the number of sources. The matrix \mathbf{W}^{τ} represents the τ^{th} slice spectral basis and \mathbf{H}^{ϕ} represents the ϕ^{th} slice of temporal code for each spectral basis element. The vertical arrow in \mathbf{W}^{τ} denotes downward shift operator which moves each element in the matrix by ϕ row down. By the same token, the horizontal arrow in \mathbf{H}^{ϕ} denotes the right shift operator which moves each element in the matrix by τ column to the right.

The factorization for NMF2D source model in (4) is based on a model that represents temporal structure and pitch change. In audio processing, the model represents each instrument by a single time-frequency profile convolved in both time and frequency by a time-pitch weight matrix. This model thoroughly decreases the number components need to model various instruments and efficiently solves the monaural source separation problem. In the following, novel algorithm of sparse NMF2D with Beta-divergence is proposed to estimate the parameter of \mathbf{W}_j^{τ} and \mathbf{H}_j^{ϕ} from the mixture.

Proposed separation method: In this section, a new algorithm of two-dimensional sparse nonnegative matrix factorization using the sparse Beta-divergence NMF2D will be developed. The algorithm optimizes the parameters of the signal model based on the multiplicative update rule using gradient descent.

Sparse representation is a representation of data where most coefficients are zero. It is proving to be a particularly interesting and powerful tool especially for analysis and processing of audio signals. If each signal to be separated has a sparse representation, then there is a good chance that there will be little overlap between the small sets of coefficients used to represent the different source signals. Therefore by selecting the

coefficients used by each source signal, we can restore each of the original signals with most of the interference from the unwanted signals removed.

Now, we incorporated the Beta-divergence as defined in (3) with the sparsity constrained such that it will minimize the cost function as follow:

$$C(\mathbf{Y}|\hat{\mathbf{Y}}) = \sum_{f,n} \left(\frac{(\mathbf{Y}_{f,n})^{\beta}}{\beta(\beta-1)} + \frac{(\hat{\mathbf{Y}}_{f,n})^{\beta}}{\beta} - \frac{\mathbf{Y}_{f,n}(\hat{\mathbf{Y}}_{f,n})^{\beta-1}}{\beta-1} \right) + \lambda f(\mathbf{H}) \quad (5)$$

For $f=1, \dots, F$, $n=1, \dots, N$ where $\hat{\mathbf{Y}} = \sum_{j,\tau,\phi} \tilde{\mathbf{W}}_j^{\tau} \mathbf{H}_j^{\phi}$ with $\tilde{\mathbf{W}}_{f,j}^{\tau} = \mathbf{W}_{f,j}^{\tau} / \sqrt{\sum_{\tau,j} (\mathbf{W}_{f,j}^{\tau})^2}$. Parameter λ is the sparsity constraint and $f(\mathbf{H})$ can be any function with positive derivative such as $L_{\zeta} = \text{norm}(\zeta > 0)$ given by $f(\mathbf{H}) = \|\mathbf{H}\|_{\zeta} = \left(\sum_{\phi,j,n} \mathbf{H}_{j,n}^{\phi} \right)^{1/\zeta}$. This will resolve the ambiguity between the factors by imposing sparseness on \mathbf{H}^{ϕ} and forcing the structure onto \mathbf{W}^{τ} .

In this study, we employed the multiplicative gradient descent approach which consists in updating each parameter by multiplying its value at the previous iteration by a certain coefficient. The derivatives of (5) corresponding to \mathbf{W}^{τ} and \mathbf{H}^{ϕ} are given by:

$$\frac{\partial C_{\beta}}{\partial \mathbf{W}_{f,j}^{\tau}} = \frac{\partial}{\partial \mathbf{W}_{f,j}^{\tau}} \left(\sum_{f,n} \left(\frac{(\mathbf{Y}_{f,n})^{\beta}}{\beta(\beta-1)} + \frac{(\hat{\mathbf{Y}}_{f,n})^{\beta}}{\beta} - \frac{\mathbf{Y}_{f,n}(\hat{\mathbf{Y}}_{f,n})^{\beta-1}}{\beta-1} \right) + \lambda f(\mathbf{H}) \right) = \sum_{\phi,n} \left((\hat{\mathbf{Y}}_{f'+\phi,n}^{\beta-1}) - \mathbf{Y}_{f'+\phi,n} (\hat{\mathbf{Y}}_{f'+\phi,n}^{\beta-2}) \right) \mathbf{H}_{j',n-\tau}^{\phi} \quad (6)$$

$$\frac{\partial C_{\beta}}{\partial \mathbf{H}_{j',n'}^{\phi}} = \frac{\partial}{\partial \mathbf{H}_{j',n'}^{\phi}} \left(\sum_{f,n} \left(\frac{(\mathbf{Y}_{f,n})^{\beta}}{\beta(\beta-1)} + \frac{(\hat{\mathbf{Y}}_{f,n})^{\beta}}{\beta} - \frac{\mathbf{Y}_{f,n}(\hat{\mathbf{Y}}_{f,n})^{\beta-1}}{\beta-1} \right) + \lambda f(\mathbf{H}) \right) = \sum_{\tau,f} \tilde{\mathbf{W}}_{f-\phi,j'}^{\tau} \left((\hat{\mathbf{Y}}_{f,n'+\tau}^{\beta-1}) - \mathbf{Y}_{f,n'+\tau} (\hat{\mathbf{Y}}_{f,n'+\tau}^{\beta-2}) \right) + \lambda \frac{\partial f(\mathbf{H})}{\partial \mathbf{H}_{j',n'}^{\phi}} \quad (7)$$

Thus, by applying the standard gradient descent approach:

$$\mathbf{W}_{f,j'}^{\tau'} \leftarrow \tilde{\mathbf{W}}_{f,j'}^{\tau'} - \eta_w \frac{\partial C_{\beta}}{\partial \mathbf{W}_{f,j'}^{\tau'}} \quad \mathbf{H}_{j',n'}^{\phi'} \leftarrow \mathbf{H}_{j',n'}^{\phi'} - \eta_H \frac{\partial C_{\beta}}{\partial \mathbf{H}_{j',n'}^{\phi'}} \quad (8)$$

Table 1: Algorithm of beta-divergence sparse NMF2D

1. Initialize \mathbf{W}^τ and \mathbf{H}^ϕ with nonnegative random values
2. Normalize $\tilde{\mathbf{W}}_{f,j}^\tau = \mathbf{W}_{f,j}^\tau / \sqrt{\sum_{\tau,j} (\mathbf{W}_{f,j}^\tau)^2}$
3. Compute $\hat{\mathbf{Y}} = \sum_{j,\tau,\phi} \tilde{\mathbf{W}}_j^\tau \mathbf{H}_j^\phi$
4. Update $\mathbf{H}^\phi \leftarrow \mathbf{H}^\phi \cdot \frac{\sum_{\tau} \tilde{\mathbf{W}}^\tau \left(\left(\hat{\mathbf{Y}} \right)^{(\beta-2)} \cdot \hat{\mathbf{Y}} \right)}{\sum_{\tau} \tilde{\mathbf{W}}^\tau \left(\left(\hat{\mathbf{Y}} \right)^{(\beta-1)} + \lambda \frac{\partial f(\mathbf{H})}{\partial \mathbf{H}^\phi} \right)}$
5. Compute $\hat{\mathbf{Y}} = \sum_{j,\tau,\phi} \tilde{\mathbf{W}}_j^\tau \mathbf{H}_j^\phi$
6. Update $\mathbf{W}^\tau \leftarrow \tilde{\mathbf{W}}^\tau \cdot \frac{\sum_{\phi} \left(\left(\hat{\mathbf{Y}} \right)^{(\beta-2)} \cdot \hat{\mathbf{Y}} \right) \mathbf{H}^\phi}{\sum_{\phi} \left(\left(\hat{\mathbf{Y}} \right)^{(\beta-1)} + \lambda \frac{\partial f(\mathbf{H})}{\partial \mathbf{H}^\phi} \right)}$

Repeat from steps 2 to 6 until convergence

where, η_W and η_H are positive learning rates which can be obtained by following (Lee and Seung, 1999), namely:

$$\eta_W = \frac{\tilde{\mathbf{W}}_{f',j'}^{\tau'}}{\sum_{\phi,n} \left(\hat{\mathbf{Y}}_{f'+\phi,n} \right)^{\beta-1} \mathbf{H}_{j',n-\tau}^\phi}$$

$$\eta_H = \frac{\mathbf{H}_{j',n'}^{\phi'}}{\sum_{\tau,j} \tilde{\mathbf{W}}_{f-\phi',j'}^\tau \left(\hat{\mathbf{Y}}_{f',n'+\tau} \right)^{\beta-1} + \lambda \frac{\partial f(\mathbf{H})}{\partial \mathbf{H}_{j',n'}^{\phi'}}}$$
(9)

Thus, the multiplicative learning rules become:

$$\mathbf{H}^\phi \leftarrow \mathbf{H}^\phi \cdot \frac{\sum_{\tau} \tilde{\mathbf{W}}^\tau \left(\left(\hat{\mathbf{Y}} \right)^{(\beta-2)} \cdot \hat{\mathbf{Y}} \right)}{\sum_{\tau} \tilde{\mathbf{W}}^\tau \left(\left(\hat{\mathbf{Y}} \right)^{(\beta-1)} + \lambda \frac{\partial f(\mathbf{H})}{\partial \mathbf{H}^\phi} \right)}$$
(10)

$$\mathbf{W}^\tau \leftarrow \tilde{\mathbf{W}}^\tau \cdot \frac{\sum_{\phi} \left(\left(\hat{\mathbf{Y}} \right)^{(\beta-2)} \cdot \hat{\mathbf{Y}} \right) \mathbf{H}^\phi}{\sum_{\phi} \left(\left(\hat{\mathbf{Y}} \right)^{(\beta-1)} + \lambda \frac{\partial f(\mathbf{H})}{\partial \mathbf{H}^\phi} \right)}$$
(11)

For (10) and (11), $A.B$ denotes element wise multiplication and $\frac{A}{B}$ denotes the element wise division. Table 1 presents the main steps of the proposed algorithm for blind separation using sparse NMF2D with Beta-divergence.

Reconstruction of the separated sources: From mixture \mathbf{Y} , we seek the two estimated sources which are:

$$\hat{\mathbf{X}}_1 = \sum_{\tau=0}^{\tau_{\max}} \sum_{\phi=0}^{\phi_{\max}} \tilde{\mathbf{W}}_1^\tau \mathbf{H}_1^\phi$$

and

$$\hat{\mathbf{X}}_2 = \sum_{\tau=0}^{\tau_{\max}} \sum_{\phi=0}^{\phi_{\max}} \tilde{\mathbf{W}}_2^\tau \mathbf{H}_2^\phi.$$

Then, by using binary masking technique (Wang, 2005) we obtained mask, \mathbf{M}_j as follows:

$$\mathbf{M}_j = \begin{cases} 1, & \text{if } \hat{\mathbf{X}}_j > \hat{\mathbf{X}}_k \\ 0, & \text{Otherwise} \end{cases} \quad (12)$$

Then, the time domain estimated signal \hat{x}_j is obtained by resynthesizing \mathbf{M}_j with the mixture \mathbf{Y} i.e., $\hat{x}_j = \text{resynthesize}(\mathbf{M}_j, \mathbf{Y})$. Here, 'resynthesize' signifies the inverse mapping of log-frequency axis to the original frequency axis and then followed by inverse Short-Time Fourier Transform (STFT) back to the time domain.

RESULTS AND DISCUSSION

The proposed algorithm is tested on audio signals containing synthetic piano sound and trumpet sound. The mixture is approximately 5 sec long and sampled at 16 kHz. In this experiment, STFT using 2048-point Hamming window with 50% overlap was used and this gives 175 frequency bins in the log-frequency spectrogram within the range of 50 Hz to 8 kHz with 24 bins/octave. This corresponds to twice the resolution of the equal source signal scale. Figure 1 shows the original TF domains of the source of piano and trumpet as well as its mixture. For audio separation, after conducting the Monte-Carlo experiments over 50 independent realizations of the mixture, the parameters of the convolutive factors of τ and ϕ shifts are set to be $\tau_{\max} = 8$ and $\phi_{\max} = 32$. This is the best attainable parameter setting to represent the temporal code and spectral basis in the factorization for most of music signals. To evaluate the proposed algorithm, the performance will be measured using the Signal-to-Distortion Ratio (SDR), Source-to-Artifacts Ratio (SAR) and Source-to-Interference Ratio (SIR) which measures an overall sound quality of the source separation. The MATLAB implementation of these measures can be found in Fevotte *et al.* (2005).

Beta performance analysis: Now, we investigate the effect of β in terms of performance of the proposed algorithm. Figure 2 shows the average SDR values obtained from various values of β using multiplicative update NMF2D algorithm. The value of β tested was

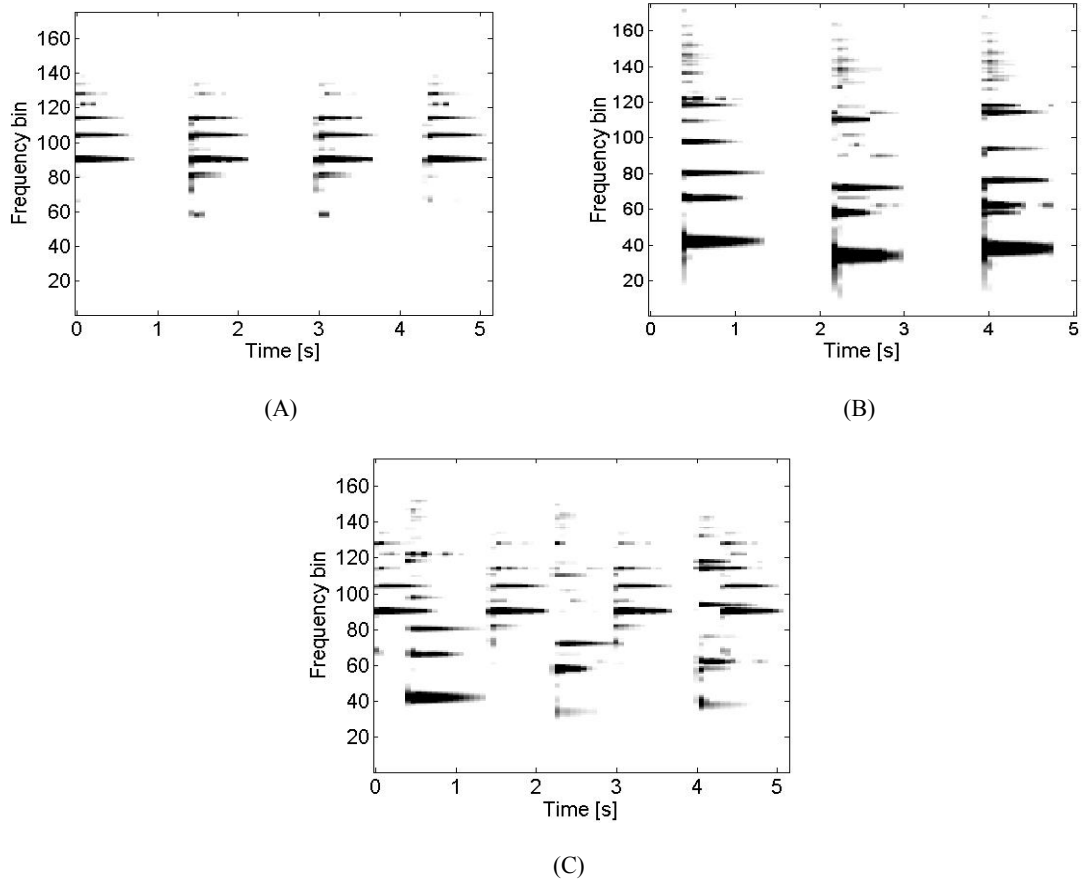


Fig. 1: Log-frequency spectrogram of (A) trumpet, (B) piano, (C) convolutive mixed signal

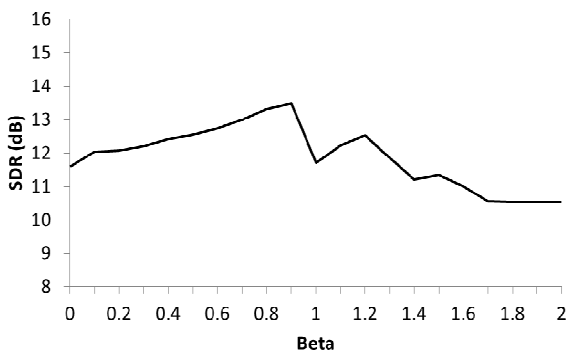


Fig. 2: Separation results for various values of β using beta-divergence NMF2D

varied from 0 to 2 in steps of 0.1. It ought to cover Least Square (LS) distant, the Kullback-Leibler (KL) divergence and the Itakura-Saito (IS) divergence of NMF2D. The average separation performance was obtained from the estimated SDR value for each source in a trumpet-piano mixture, thereby providing a measure of overall separation for each signal. From Fig. 2, as we increase the value of β , the performance also increase and it reach its peak value when $\beta = 0.9$ with average

SDR value of 13.5 dB is obtained for each source. A tail-off in performance occurs as the value of β increases from 0.9 goes up to 2. From this experiment, it suggests that β around 0.9 is an optimal value for audio separation which will be used in our experiment in the next sub-section.

Blind audio source separation results: Here, we compare the performance of audio source separation of proposed algorithms of Beta-divergence sparse NMF2D with the one without the sparsity constraint. We set $\beta = 0.9$ for both algorithm. The best value of sparsity parameter λ was identified as 0.5 after conducting the Monte-Carlo experiments over 50 independent realizations. L_1 -norm regularization is used to resolve the ambiguity by forcing all structure in H onto W . Figure 3 shows the separation result in log-frequency spectrogram for both algorithms. Compared with original sources in Fig. 1, it is visually clear that separation of Beta-divergence NMF2D without the sparseness in Fig. 3A and B led to poor result since the factorization still contains the mixed signal (indicated by the box marked area). This is because without the sparsity constraint, it leads to component ambiguity, i.e.,

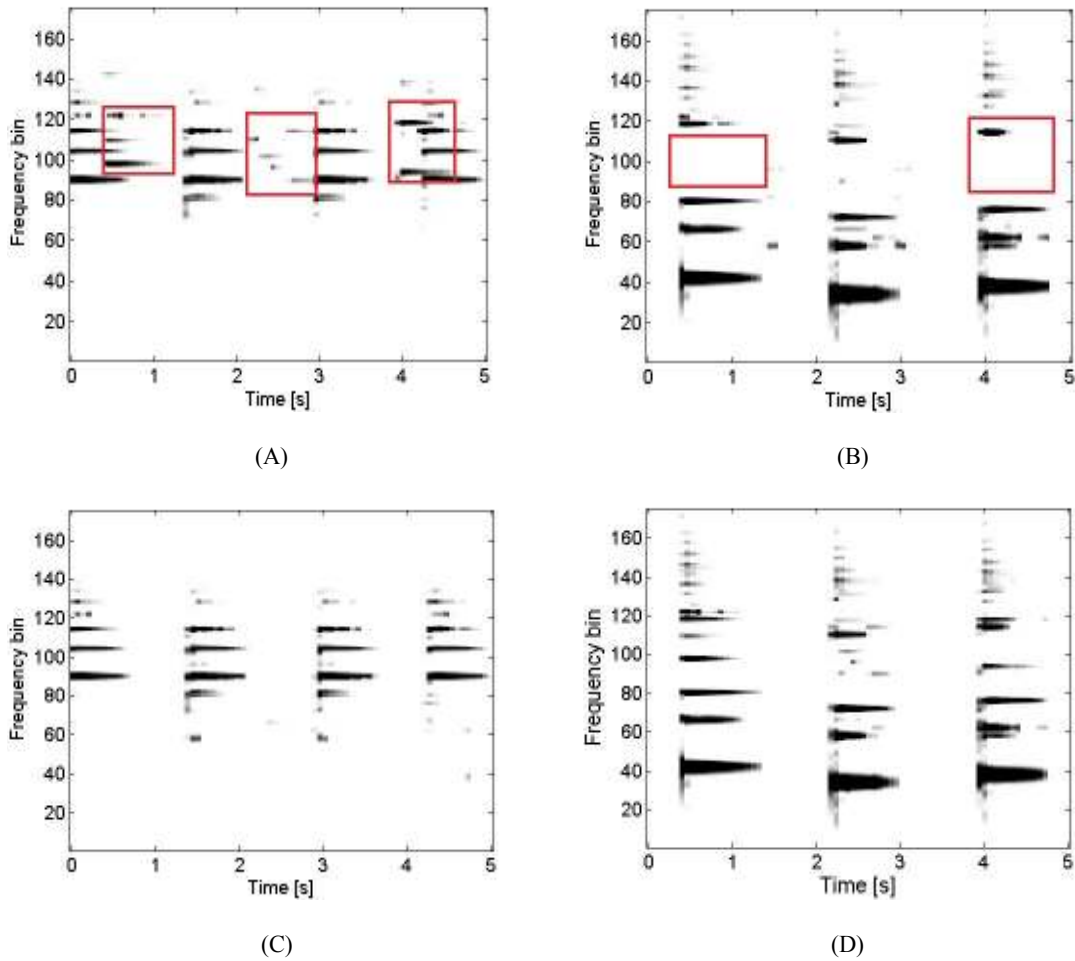


Fig. 3: Separated sound in log-frequency spectrogram for beta-divergence NMF2D (A)-(B) trumpet and piano sound without sparsity (C)-(D) trumpet and piano sound with sparsity

Table 2: Separation results for NMF2D with beta-divergence

Algorithms	Separated trumpet			Separated piano		
	SDR	SIR	SAR	SDR	SIR	SAR
Beta-divergence NMF2D	12.7	16.2	13.2	14.2	21.2	14.7
Sparse beta-divergence NMF2D	15.5	23.1	16.6	16.8	23.1	17.9

lack of uniqueness in decomposition. In contrary, by employing the sparseness, it has yielded the better performance when the decomposition of spectral bases and temporal codes is performed with the sparsity constraint.

From Table 2, in general both algorithms of Beta-divergence NMF2D provide decent results with the performance of SDR, SIR and SAR that can be considered good. Over 10 dB of SDR measurement have been recorded for both methods. However, performance of Beta-divergence NMF2D with sparsity constraint is superior with the average SDR improvement of 2.7 dB per source compare with the one without imposing the sparseness. In percentage, this translates to an average improvement of 20%. From this result, it can be inferred

that the sparsity constraint have significant effects on the separation performance.

CONCLUSION

The use of the Beta-divergence for audio source separation using NMF2D model has been investigated. The value of Beta-divergence with $\beta = 0.5$ was found to produce an optimal result. Furthermore, the proposed sparse Beta-divergence NMF2D is developed under the probabilistic framework which enables sparseness to be incorporated in the solution. This will significantly resolve the ambiguities problem in the factorization. We confirmed through an experiment that the proposed algorithm performs exceptionally well in separation of

an audio mixture with the high value of SDR has been achieved.

ACKNOWLEDGMENT

The authors would like to thank Universiti Teknikal Malaysia Melaka for the research grant funding PJP/2012/FKEKK (15D)/S01127.

REFERENCES

- Biciu, I., N. Nikolaidis and I. Pitas, 2007. Nonnegative matrix factorization in polynomial feature space. *IEEE T. Neural Networ.*, 19: 1090-1100.
- Fevotte, C., R. Gribonval and E. Vincent, 2005. BSS EVAL toolbox user guide. IRISA Technical Report 1706, Rennes, France.
- Fevotte, C., N. Bertin and J.L. Durrieu, 2009. Nonnegative matrix factorization with the Itakura-Saito divergence with application to music analysis. *Neural Comput.*, 21: 793-830.
- FitzGerald, D., 2004. Automatic drum transcription and source separation. Ph.D. Thesis, Dublin Institute of Technology, Dublin, Ireland.
- Kompass, R., 2005. Generalized divergence measure for non-negative matrix factorization. *Proceeding of the Neuroinformatics Workshop*. Torun, Poland.
- Lee, C. and H. Seung, 1999. Learning the parts of objects by nonnegative matrix factorisation. *Nature*, 401(6755): 788-791.
- Morup, M. and M.N. Schmidt, 2006. Sparse nonnegative matrix factor 2-D deconvolution. Technical Report, Technical University of Denmark, Copenhagen, Denmark.
- Wang, D.L., 2005. On Ideal Binary Mask as the Computational Goal of Auditory Scene Analysis. In: Diventi, P. (Ed.), *Speech Separation by Humans and Machines*. Kluwer, Norwell, MA, pp: 181-197.
- Xie, S., Z. Yang and Y. Fu, 2008. Nonnegative matrix factorization applied to nonlinear speech and image cryptosystems. *IEEE T. Circuits-I*, 55(8): 2356-2367.