



Separating Time-Frequency Sources from Time-Domain Convolutive Mixtures Using Non-negative Matrix Factorization

Simon Leglaive, Roland Badeau, Gaël Richard LTCI, Télécom ParisTech, Université Paris Saclay

IEEE Workshop on Applications of Signal Processing to Audio and Acoustics (WASPAA) New Paltz, NY

October 17, 2017

Probabilistic model

Inference

Experiments

Conclusion O

Multichannel audio source separation

Objective: Recover source signals from the observation of several mixtures.

Context: Under-determined.



Probabilistic model

Inference

Experiments

Conclusion

Time-frequency source representation

$\label{eq:time-frequency} \ensuremath{\mathsf{TF}}\xspace) \ensuremath{\mathsf{TF}}\x$



Probabilistic model

Inference

Experiments

Conclusion

Reverberant mixtures (1)

Convolutive mixing process in the time domain: $x_i(t) = \sum_{j=1}^{J} [a_{ij} \star s_j](t)$



Probabilistic model

Inference

Experiments

Conclusion O

Reverberant mixtures (2)

Convolutive mixing process in the STFT domain: $x_{i,fn} \approx \sum_{j=1}^{J} a_{ij,f} s_{j,fn}$



Introduction 0000●0 Probabilistic model

Inference

Experiments

Conclusion O

Proposed approach

- Time-domain mixture representation: $x_i(t) = \sum_{j=1}^{J} [a_{ij} \star s_j](t)$
- Time-frequency source representation: $s_j(t) = \mathcal{T}^{-1}(\{s_{j,fn}\}_{f,n})$



Introduction	
000000	

Choosing the time-frequency transform

- Modified Discrete Cosine Transform (MDCT)
 - Real valued;
 - Critically sampled;
 - No shift-invariance (phase information is contained in the amplitude of the MDCT coefficients).
- Odd-Frequency Short-Time Fourier Transform (OFSTFT)
 - Complex valued;
 - Redundant;
 - Shift-invariance.
- General TF synthesis equation:

$$s_j(t) = \frac{2}{\phi} \Re \left(\sum_{f=0}^{F-1} \sum_{n=0}^{N-1} s_{j,fn} \psi_{fn}(t) \right), \text{ with } \phi = \begin{cases} 2 & \text{ if MDCT } (s_{j,fn} \in \mathbb{R}) \\ 1 & \text{ if OFSTFT } (s_{j,fn} \in \mathbb{C}) \end{cases}$$

Probabilistic model

Inference

Experiments

Conclusion O

Outline

Probabilistic model

Inference

Experiments

Conclusion

7/23

ntroduction	Probabilistic model	Inference	Experiments	Conclusion
00000	● ○ ○	000	000000	0

Probabilistic modeling with latent variables

- ▶ Latent TF source random variables: $\mathbf{s} = \{s_{j,fn} \in \mathbb{R} \text{ or } \mathbb{C}\}_{j,f,n}$
- Observed time-domain random variables: $\mathbf{x} = \{x_i(t) \in \mathbb{R}\}_{i,t}$



How are the data generated from the latent unobserved variables?

Probabilistic model ○●○ Inference

Conclusion O

Prior distribution of the latent variables

Gaussian source model based on Non-negative Matrix Factorization [1]:

 $s_{j,fn} \sim \mathcal{N}(0, [\mathbf{W}_{j}\mathbf{H}_{j}]_{fn})$



[1] C. Févotte, N. Bertin, J.-L. Durrieu. "Nonnegative matrix factorization with the Itakura-Saito divergence: With application to music analysis". *Neural computation*, 2009.

Introduction	Probabilistic model	Inference	Experiments	Conclusion	
000000	00●	000	0000000	0	
			- ·		

Conditional distribution of x given s

Gaussian modeling error

$$x_i(t) = \sum_{j=1}^J [a_{ij} \star s_j](t) + b_i(t),$$

with
$$b_i(t) \stackrel{i.i.d}{\sim} \mathcal{N}(0,\sigma_i^2)$$
 and $s_j(t) = \frac{2}{\phi} \Re\left(\sum_{f=0}^{F-1} \sum_{n=0}^{N-1} s_{j,fn} \psi_{fn}(t)\right).$

Conditional distribution

$$|\mathbf{x}_i(t)|\mathbf{s}; \boldsymbol{\theta} \sim \mathcal{N}\left(rac{2}{\phi} \Re\left(\sum_{j=1}^J \sum_{f=0}^{F-1} \sum_{n=0}^{N-1} s_{j,fn}[a_{ij} \star \psi_{fn}](t)
ight), \sigma_i^2
ight)$$

Probabilistic model

Inference

Experiments

Conclusion O



Outline

Probabilistic model

Inference

Experiments

Conclusion

11/23

Probabilistic model

Inference ●○○ Experiments

Conclusion O

Inference

Posterior distribution

We are interested in the posterior distribution of the latent variables:

$$p(\mathbf{s}|\mathbf{x}; \boldsymbol{\theta}^{\star})$$
 with $\boldsymbol{\theta}^{\star} = \arg \max_{\boldsymbol{\theta}} p(\mathbf{x}; \boldsymbol{\theta})$

• Model parameters:
$$\boldsymbol{ heta} = \left\{ \{ \mathbf{W}_j, \mathbf{H}_j \}_j, \{ a_{ij}(t) \}_{i,j,t}, \{ \sigma_i^2 \}_i
ight\}$$

Semi-blind setting: the mixing filters are assumed to be known.

The posterior distribution is Gaussian but with a high-dimensional full covariance matrix \rightarrow variational inference to reduce the computational cost.

Probabilistic model

Inference ○●○ Experiments

Conclusion O

Variational inference

- We want to find $q \in \mathcal{F}$ which approximates $p(\mathbf{s}|\mathbf{x}; \boldsymbol{\theta})$.
- ► Taking the KL divergence as a measure of fit, we can show that:

$$\mathcal{K}L(q(\mathbf{s}) || p(\mathbf{s} | \mathbf{x}; \boldsymbol{\theta})) = \underbrace{\ln p(\mathbf{x}; \boldsymbol{\theta})}_{\text{Log-likelihood}} - \underbrace{\mathcal{L}(q; \boldsymbol{\theta})}_{\text{Variational Free Energy}},$$
(1)
where $\mathcal{L}(q; \boldsymbol{\theta}) = \left\langle \ln \left(\frac{p(\mathbf{x}, \mathbf{s}; \boldsymbol{\theta})}{q(\mathbf{s})} \right) \right\rangle_{q}$ and $\langle f(\mathbf{z}) \rangle_{q} = \int f(\mathbf{z}) q(\mathbf{z}) d\mathbf{z}.$

Variational Expectation-Maximization algorithm:

► **E-step**:
$$q^* = \underset{q \in \mathcal{F}}{\operatorname{arg\,min}} KL(q(\mathbf{s}) || p(\mathbf{s} | \mathbf{x}; \boldsymbol{\theta}^*)) = \underset{q \in \mathcal{F}}{\operatorname{arg\,max}} \mathcal{L}(q; \boldsymbol{\theta}^*)$$

• **M-step**:
$$\theta^* = \arg \max_{\theta} \mathcal{L}(q^*; \theta)$$

Probabilistic model

Inference

Experiments

Conclusion O

Mean-field approximation





 ${\cal F}$ is the set of pdfs that factorize as:

$$q(\mathbf{s}) = \prod_{j=1}^{J} \prod_{f=0}^{F-1} \prod_{n=0}^{N-1} q_{jfn}(s_{j,fn}).$$

Under the mean-field approximation we can show that:

$$q_{jfn}^{\star}(s_{j,fn}) = \begin{cases} N_{\mathbb{R}}(\hat{s}_{j,fn}^{r}, \gamma_{j,fn}^{r}) & \text{if MDCT} \\ N_{\mathbb{C}}(\rho_{j,fn}, \hat{s}_{j,fn}^{r}, \hat{s}_{j,fn}^{i}, \gamma_{j,fn}^{r}, \gamma_{j,fn}^{i}) & \text{if OFSTFT} \end{cases}$$
(2)

In the OFSTFT case, $\Re(s_{j,fn})$ and $\Im(s_{j,fn})$ are correlated a posteriori.

Probabilistic model

Inference

Experiments

Conclusion O

Outline

Probabilistic model

Inference

Experiments

Conclusion

15/23

Probabilistic model

Inference

Experiments •000000 Conclusion O

Experiments

- Dataset:
 - ► 8 stereo mixtures created with measured room impulse responses from the RWCP database [2].
 - Reverberation time: 470 ms.
 - Number of sources per mixture: 3 to 5.
 - Mixture length: 12 to 28 seconds.
- Semi-blind setting: Mixing filters are known while all other parameters are blindly estimated.

^[2] S. Nakamura et al. "Acoustical sound database in real environments for sound scene understanding and hands-free speech recognition". *Proc. of LREC*, 2000.

Probabilistic model

Inference

Experiments ⊙●○○○○○ Conclusion O

MDCT vs. OFSTFT

Comparison of the source separation results using:

- ▶ the MDCT;
- ▶ the OFSTFT with several overlap ratios: 25%, 50% and 75%.



Audio examples available online (website address in the paper)

Probabilistic model

Inference

Experiments

Conclusion O

MDCT vs. OFSTFT

Comparison of the source separation results using:

- ▶ the MDCT;
- ▶ the OFSTFT with several overlap ratios: 25%, 50% and 75%.



Audio examples available online (website address in the paper)

ntroducti	ion Probabilisti 000	c model Inference	Experiments 000●000	Conclusion ○
	Baseline	methods		
=		Time-frequency source model (STFT domain)	Convolutive mixture representation	
-	Ozerov et al. [3]	Gaussian NMF-based	approximate (STFT)	
_	Kowalski et al. [4]	sparse $(\ell_1 {\sf norm})$	exact (time)	

Length of the TF analysis/synthesis window: 128 ms.

^[3] A. Ozerov, C. Févotte, "Multichannel nonnegative matrix factorization in convolutive mixtures for audio source separation", IEEE Trans. Audio, Speech, Language Process., 2010.

^[4] M. Kowalski, E. Vincent, R. Gribonval, "Beyond the narrowband approximation: Wideband convex methods for under-determined reverberant audio source separation", *IEEE Trans. Audio, Speech, Language Process.*, 2010.

Probabilistic model

Inference

Experiments 0000●00 Conclusion O

Source separation results



Probabilistic model

Inference

Experiments

Conclusion O

Computational time



Separation of a 12 second-long mixture of 3 sources at 16 kHz: \sim 2 hours with the MDCT-based method.

20/23

Probabilistic model

Inference

Experiments

Conclusion O





	Guitar 1	Guitar 2	Voice	Drums	Bass
Original source (stereo)	0	0	0	0	0
Ozerov et al.	0	0	0	0	0
Kowalski et al.	0	0	0		0
Proposed (MDCT)	0	0	0	0	0

Excerpt from "Ana" by Vieux Farka Toure.

Probabilistic model

Inference

Experiments

Conclusion



Probabilistic model

Inference

Experiments

Conclusion

22/23

Probabilistic model

Inference

Experiments

Conclusion •

Conclusion

Conclusion:

 Working with the MDCT is computationally cheaper than with the OFSTFT and leads to similar perceived results.

Further work:

- Source-specific time-frequency resolution;
- Probabilistic priors on the mixing filters in the time domain.

Journal preprint:

- S. L., R. Badeau, G. Richard, "Student's t source and mixing models for multichannel audio source separation", submitted, 2017;
- Available online: https://hal.archives-ouvertes.fr/hal-01584755;
- Poster at the SANE Workshop.



Thank you

More audio examples and Matlab code available at: https://perso.telecom-paristech.fr/leglaive/

Probabilistic model

Inference

Experiments

Conclusion O

Modified discrete cosine transform

MDCT synthesis equation:

$$s_j(t) = \sum_{f=0}^{F-1} \sum_{n=0}^{N-1} s_{j,fn} \psi_{fn}(t),$$

•
$$\psi_{fn}(t) = \sqrt{\frac{2}{F}}w(t-nH)\cos\left(\frac{2\pi}{L_w}\left(t-nH+\frac{1}{2}+\frac{L_w}{4}\right)\left(f+\frac{1}{2}\right)\right);$$

•
$$w(t)$$
: synthesis window of length L_w ;

$$\blacktriangleright F = H = L_w/2.$$

Probabilistic model

Inference

Experiments

Conclusion O

Short-time Fourier transform

STFT synthesis equation:

$$s_j(t) = \sum_{f=0}^{F-1} \sum_{n=0}^{N-1} s_{j,fn} \psi_{fn}(t),$$

$$\psi_{fn}(t) = \sqrt{\frac{1}{L_w}} w(t - nH) \exp\left(i\frac{2\pi}{L_w}f(t - nH)\right);$$

$$F = L_w.$$

- Hermitian symmetry: deterministic relation between TF coefficients.
- Using the Hermitian symmetry property:

$$s_{j}(t) = \sum_{n=0}^{N-1} \left[\underbrace{\underbrace{s_{j,0n} \psi_{0n}(t)}_{\text{Zero frequency}} + \underbrace{s_{j,\frac{F}{2}n} \psi_{\frac{F}{2}n}(t)}_{\text{Nyquist frequency}} + 2\Re \left(\sum_{f=1}^{F/2-1} s_{j,fn} \psi_{fn}(t) \right) \right]$$

 Introduction
 Probabilistic model
 Inference
 Experiments
 Conclusion

 000000
 000
 000
 0000000
 0
 0

Odd-frequency short-time Fourier transform

OFSTFT synthesis equation:

$$s_j(t) = 2\Re\left(\sum_{f=0}^{F-1}\sum_{n=0}^{N-1}s_{j,fn}\psi_{fn}(t)
ight),$$

All TF coefficients are complex valued.