



A source/filter model with adaptive constraints for NMF-based speech separation

Damien Bouvier¹, Nicolas Obin¹, Axel Roebel¹, Marco Liuni¹

¹IRCAM, UPMC - Sorbonne Universités, CNRS

International Conference on Acoustics, Speech and Signal Processing
24 March 2016

State of the art

- Speech separation using NMF
- Semi-supervised NMF
- Source/filter model

Proposed method

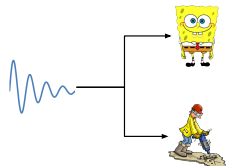
- Semi-supervised constrained NMF
- Contribution 1: speech-specific source/filter coherence constraint
- Contribution 2: adaptive weight method

Experimental evaluation

- Experiment description
- Effect of weight's adaptation
- Algorithm comparison

Speech separation using NMF

Signal has only 2 sources:
speech and background sound



Supervised algorithms

- [Mysore and Smaragdis, 2012]: language model
- [Virtanen et al., 2013]: new updates using Newton algorithm
- [Sun and Mysore, 2013]: Universal Speech Model (USM)

Semi-supervised algorithms

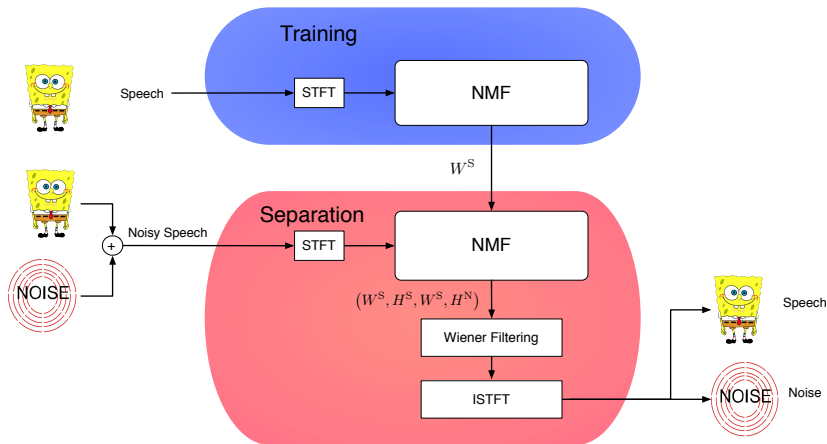
- [Germain and Mysore, 2015]: USM & online noise adaptation

Unsupervised algorithms (but informed)

- [Le Magoarou et al., 2014]: use of textual information
- [Durrieu et al., 2009]: source/filter model for NMF

Semi-supervised NMF

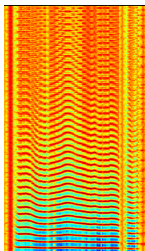
$$\underset{W^S, W^N, H^N \geq 0}{\operatorname{argmin}} \mathcal{C}(V|\tilde{V}) \text{ with } \begin{cases} \tilde{V} = W^S H^S + W^N H^N \\ W^S \text{ learned} \end{cases} \quad (1)$$



Source/filter model [Durrieu et al., 2009]

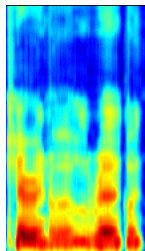


Source/filter model [Durrieu et al., 2009]


 V^{ex}

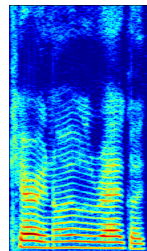
.*

.*

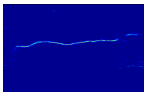
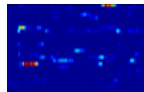
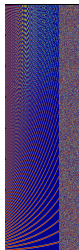
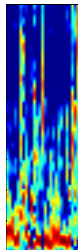

 V^{Φ}

=

=

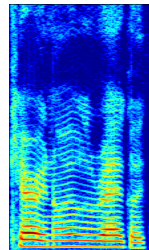
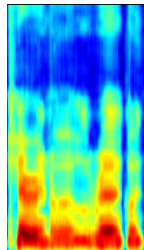
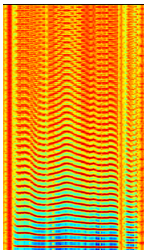

 V

Source/filter model [Durrieu et al., 2009]

 H^{ex}  H^{Φ}  W^{ex}  \hat{W}^{Φ} 

.*

=

 $(W^{\text{ex}}$ $H^{\text{ex}})$

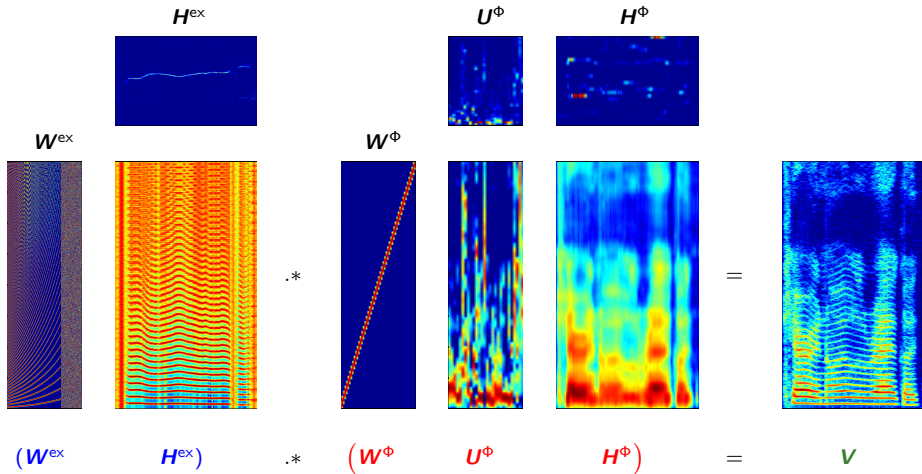
.*

 $(\hat{W}^{\Phi}$ $H^{\Phi})$

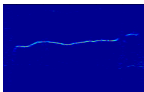
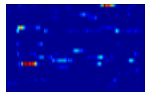
=

 V

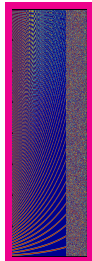
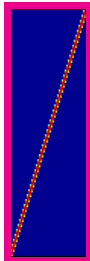
Source/filter model [Durrieu et al., 2009]



Source/filter model [Durrieu et al., 2009]

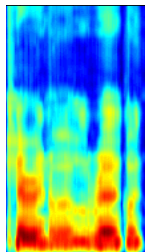
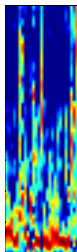
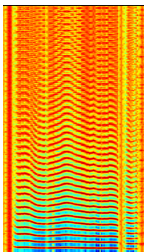
 H^{ex}  U^{Φ}  H^{Φ} 

Fixed

 W^{ex}  W^{Φ} 

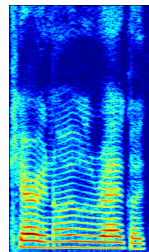
.*

.*

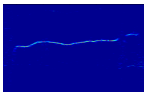
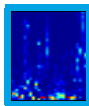
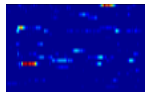
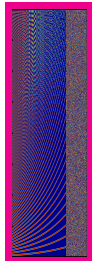
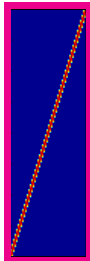


=

=

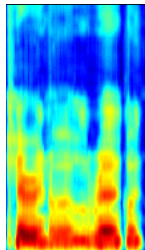
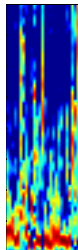
 $(W^{\text{ex}}$ $H^{\text{ex}})$ $(W^{\Phi}$ U^{Φ} $H^{\Phi})$ V

Source/filter model [Durrieu et al., 2009]

 H^{ex}  U^Φ  H^Φ Fixed
Trained W^{ex}  W^Φ 

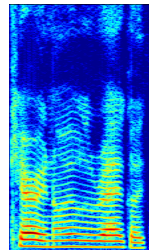
.*

.*

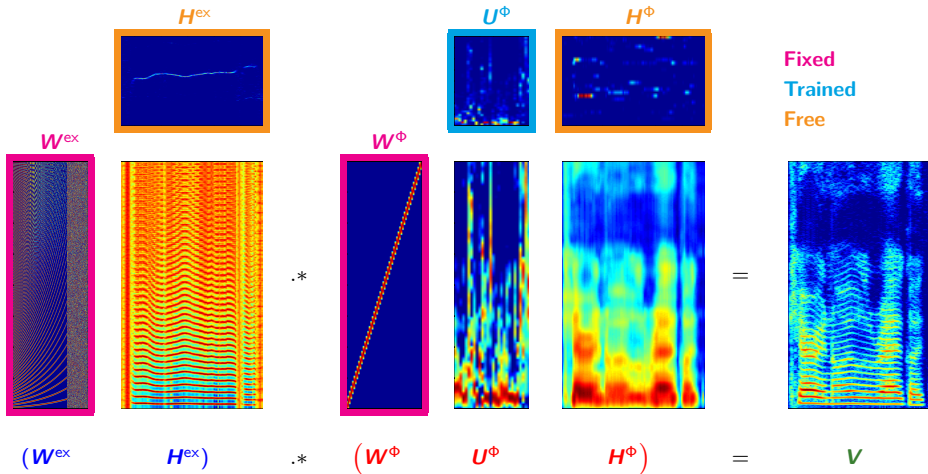


=

=

 $(W^{\text{ex}}$ $H^{\text{ex}})$ $(W^\Phi$ U^Φ $H^\Phi)$ V

Source/filter model [Durrieu et al., 2009]



State of the art

- Speech separation using NMF
- Semi-supervised NMF
- Source/filter model

Proposed method

- Semi-supervised constrained NMF
- Contribution 1: speech-specific source/filter coherence constraint
- Contribution 2: adaptive weight method

Experimental evaluation

- Experiment description
- Effect of weight's adaptation
- Algorithm comparison

Semi-supervised constrained NMF

Physically-informed model

$$\underset{H^{\text{ex}}, H^{\Phi}, W^N, H^N \geq 0}{\text{argmin}} \quad \mathcal{C}(V | \tilde{V}) \text{ with } \begin{cases} \tilde{V} = W^{\text{ex}} H^{\text{ex}} \otimes W^{\Phi} U^{\Phi} H^{\Phi} + W^N H^N \\ W^{\text{ex}} \text{ and } W^{\Phi} \text{ fixed} \\ U^{\Phi} \text{ learned} \end{cases} \quad (2)$$

But still no physically-coherent behavior.

Semi-supervised constrained NMF

Physically-informed model

$$\underset{H^{\text{ex}}, H^{\Phi}, W^N, H^N \geq 0}{\text{argmin}} \mathcal{C}(V|\tilde{V}) \text{ with } \begin{cases} \tilde{V} = W^{\text{ex}} H^{\text{ex}} \otimes W^{\Phi} U^{\Phi} H^{\Phi} + W^N H^N \\ W^{\text{ex}} \text{ and } W^{\Phi} \text{ fixed} \\ U^{\Phi} \text{ learned} \end{cases} \quad (2)$$

But still no physically-coherent behavior.

Constraints for controlling its behavior

$$\underbrace{\mathcal{C}(V|\tilde{V})}_{\text{Total cost}} = \underbrace{D(V|\tilde{V})}_{\text{Reconstruction cost}} + \underbrace{\lambda}_{\text{Weight parameter}} \underbrace{\mathcal{P}(\Theta)}_{\text{Constraint penalty}} \quad (3)$$

Semi-supervised constrained NMF

Physically-informed model

$$\underset{\mathbf{H}^{\text{ex}}, \mathbf{H}^\Phi, \mathbf{W}^N, \mathbf{H}^N \geq 0}{\text{argmin}} \quad \mathcal{C}(\mathbf{V} | \tilde{\mathbf{V}}) \quad \text{with} \quad \begin{cases} \tilde{\mathbf{V}} = \mathbf{W}^{\text{ex}} \mathbf{H}^{\text{ex}} \otimes \mathbf{W}^\Phi \mathbf{U}^\Phi \mathbf{H}^\Phi + \mathbf{W}^N \mathbf{H}^N \\ \mathbf{W}^{\text{ex}} \text{ and } \mathbf{W}^\Phi \text{ fixed} \\ \mathbf{U}^\Phi \text{ learned} \end{cases} \quad (2)$$

But still no physically-coherent behavior.

Constraints for controlling its behavior

$$\underbrace{\mathcal{C}(\mathbf{V} | \tilde{\mathbf{V}})}_{\text{Total cost}} = \underbrace{D(\mathbf{V} | \tilde{\mathbf{V}})}_{\text{Reconstruction cost}} + \underbrace{\lambda}_{\text{Weight parameter}} \underbrace{\mathcal{P}(\Theta)}_{\text{Constraint penalty}} \quad (3)$$

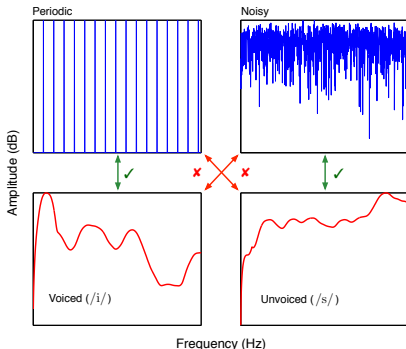
New multiplicative update rules :

$$\Theta^{(i+1)} \leftarrow \Theta^{(i)} \otimes \frac{\nabla_{\Theta}^- D + \lambda \nabla_{\Theta}^- \mathcal{P}}{\nabla_{\Theta}^+ D + \lambda \nabla_{\Theta}^+ \mathcal{P}} \quad \forall \Theta \in \{\mathbf{H}^{\text{ex}}, \mathbf{H}^\Phi, \mathbf{W}^N, \mathbf{H}^N\} \quad (4)$$

Constraints from literature [Bertin, 2009]: Sparsity, Decorrelation, Smoothness.

Contribution 1: speech-specific source/filter coherence constraint

Problem: unrealistic source/filter combinations possible



We only want to allow:

- periodic excitation with adequate filter (e.g. vowels, voice consonants)
- noisy excitation with adequate filter (e.g., unvoiced consonants)

Contribution 1: speech-specific source/filter coherence constraint

New constraint (that requires phoneme-labelled spectral filter basis)

$$\mathcal{P}_\phi(\mathbf{H}^{\text{ex}}, \mathbf{H}^\Phi) \tag{5}$$

Contribution 1: speech-specific source/filter coherence constraint

New constraint (that requires phoneme-labelled spectral filter basis)

$$\mathcal{P}_\phi(\mathbf{H}^{\text{ex}}, \mathbf{H}^\Phi) = \sum_{\substack{k \in \text{periodics} \\ l \in \text{unvoiced}}} \frac{[\mathbf{H}^{\text{ex}} \mathbf{H}^\Phi]^T]_{kl}}{\quad} \quad (5)$$

► : measure of correlation

Contribution 1: speech-specific source/filter coherence constraint

New constraint (that requires phoneme-labelled spectral filter basis)

$$\mathcal{P}_\phi(\mathbf{H}^{\text{ex}}, \mathbf{H}^\Phi) = \sum_{\substack{k \in \text{periodics} \\ l \in \text{unvoiced}}} \frac{\left[\mathbf{H}^{\text{ex}} \mathbf{H}^\Phi \right]_{kl}}{\|\mathbf{H}_k^{\text{ex}}\|_{\ell_2} \|\mathbf{H}_l^\Phi\|_{\ell_2}} \quad (5)$$

- : measure of correlation
- : normalized

Contribution 1: speech-specific source/filter coherence constraint

New constraint (that requires phoneme-labelled spectral filter basis)

$$\mathcal{P}_\phi(\mathbf{H}^{\text{ex}}, \mathbf{H}^\Phi) = \sum_{\substack{k \in \text{periodics} \\ l \in \text{unvoiced}}} \frac{[\mathbf{H}^{\text{ex}} \mathbf{H}^\Phi \mathbf{T}]_{kl}}{\|\mathbf{H}_k^{\text{ex}}\|_{\ell_2} \|\mathbf{H}_l^\Phi\|_{\ell_2}} + \sum_{\substack{k \in \text{noisy} \\ l \in \text{voiced}}} \frac{[\mathbf{H}^{\text{ex}} \mathbf{H}^\Phi \mathbf{T}]_{kl}}{\|\mathbf{H}_k^{\text{ex}}\|_{\ell_2} \|\mathbf{H}_l^\Phi\|_{\ell_2}} \quad (5)$$

- : measure of correlation
- : normalized
- : for both type of unwanted combination

Contribution 2: adaptive weight method

Main issue with constrained NMF

Adjusting the weight parameter λ :

- if too small, no effect is visible;
- if too big, convergence becomes extremely sensitive to initialization (which is typically random).

Contribution 2: adaptive weight method

Main issue with constrained NMF

Adjusting the weight parameter λ :

- if too small, no effect is visible;
- if too big, convergence becomes extremely sensitive to initialization (which is typically random).

Idea

Adjust the constraint weight at each iteration of the NMF:

- constraint relaxed during strong evolution of the reconstruction cost;
- constraint enforced when the reconstruction is more stable;

Contribution 2: adaptive weight method

Main issue with constrained NMF

Adjusting the weight parameter λ :

- if too small, no effect is visible;
- if too big, convergence becomes extremely sensitive to initialization (which is typically random).

Idea

Adjust the constraint weight at each iteration of the NMF:

- constraint relaxed during strong evolution of the reconstruction cost;
- constraint enforced when the reconstruction is more stable;

Adaptive method

$$\lambda^{(i)} = \lambda_{Max} \frac{D(\mathbf{V}|\tilde{\mathbf{V}}^{(i-1)})}{D(\mathbf{V}|\tilde{\mathbf{V}}^{(i-2)})} \quad (6)$$

$$\left(\begin{array}{l} D(\mathbf{V}|\tilde{\mathbf{V}}) \text{ the reconstruction cost} \\ \lambda \in [0 \lambda_{Max}] \end{array} \right)$$

State of the art

- Speech separation using NMF
- Semi-supervised NMF
- Source/filter model

Proposed method

- Semi-supervised constrained NMF
- Contribution 1: speech-specific source/filter coherence constraint
- Contribution 2: adaptive weight method

Experimental evaluation

- Experiment description
- Effect of weight's adaptation
- Algorithm comparison

Experiment description

Database

TIMIT [Zue et al., 1990] : 20 speakers, 2 learnings sentences and 8 test sentences

QUT-NOISE [Dean et al., 2010] : 5 types of noise

Mixed at 3 Signal-to-Noise Ratio (-6dB , $+0\text{dB}$ and $+6\text{dB}$)

Benchmark

SoA	<ul style="list-style-type: none">● ASNA [Virtanen et al., 2013]● IMM [Durrieu et al., 2009]	supervised unsupervised
Proposed	<ul style="list-style-type: none">● S-IMM: without constraints● SC-IMM1: state-of-the art constraints● SC-IMM2: source/filter coherence constraint● SC-IMM3: all constraints	semi-supervised

Measures

- SDR : Signal to Distortion Ratio (in dB)
- PESQ : Perceptual Evaluation of Speech Quality (from 1 (bad) to 5 (excellent))

Effect of weight's adaptation

SNR	Measure	With adaptation			Without adaptation		
		SC-IMM1	SC-IMM2	SC-IMM3	SC-IMM1	SC-IMM2	SC-IMM3
-6dB	SDR (dB)	4.1	5.2	5.4	4.1	5.0	5.2
	PESQ	1.91	2.01	2.01	1.91	1.94	1.92
+0dB	SDR (dB)	9.2	9.8	9.8	9.2	9.0	8.9
	PESQ	2.30	2.34	2.35	2.30	2.24	2.23
+6dB	SDR (dB)	13.0	12.8	12.9	12.8	11.1	10.9
	PESQ	2.62	2.59	2.62	2.61	2.46	2.44
Mean	SDR (dB)	8.7	9.3	9.4	8.7	8.4	8.3
	PESQ	2.28	2.31	2.33	2.27	2.21	2.20

⇒ **adaptation gives best results**

Algorithm comparison

SNR	Measure	Algorithms					
		ASNA	IMM	S-IMM	SC-IMM1	SC-IMM2	SC-IMM3
-6dB	SDR (dB)	5.8	4.4	4.0	4.1	5.2	5.4
	PESQ	2.00	1.22	1.91	1.91	2.01	2.01
+0dB	SDR (dB)	10.7	7.8	9.1	9.2	9.8	9.8
	PESQ	2.44	1.54	2.30	2.30	2.34	2.35
+6dB	SDR (dB)	15.0	9.7	13.0	13.0	12.8	12.9
	PESQ	2.85	1.82	2.62	2.62	2.59	2.62
Mean	SDR (dB)	10.5	7.3	8.7	8.7	9.3	9.4
	PESQ	2.43	1.52	2.28	2.28	2.31	2.33

Algorithm comparison

SNR	Measure	Algorithms					
		ASNA	IMM	S-IMM	SC-IMM1	SC-IMM2	SC-IMM3
-6dB	SDR (dB)	5.8	4.4	4.0	4.1	5.2	5.4
	PESQ	2.00	1.22	1.91	1.91	2.01	2.01
+0dB	SDR (dB)	10.7	7.8	9.1	9.2	9.8	9.8
	PESQ	2.44	1.54	2.30	2.30	2.34	2.35
+6dB	SDR (dB)	15.0	9.7	13.0	13.0	12.8	12.9
	PESQ	2.85	1.82	2.62	2.62	2.59	2.62
Mean	SDR (dB)	10.5	7.3	8.7	8.7	9.3	9.4
	PESQ	2.43	1.52	2.28	2.28	2.31	2.33

⇒ supervision helps separation

Algorithm comparison

SNR	Measure	Algorithms					
		ASNA	IMM	S-IMM	SC-IMM1	SC-IMM2	SC-IMM3
-6dB	SDR (dB)	5.8	4.4	4.0	4.1	5.2	5.4
	PESQ	2.00	1.22	1.91	1.91	2.01	2.01
+0dB	SDR (dB)	10.7	7.8	9.1	9.2	9.8	9.8
	PESQ	2.44	1.54	2.30	2.30	2.34	2.35
+6dB	SDR (dB)	15.0	9.7	13.0	13.0	12.8	12.9
	PESQ	2.85	1.82	2.62	2.62	2.59	2.62
Mean	SDR (dB)	10.5	7.3	8.7	8.7	9.3	9.4
	PESQ	2.43	1.52	2.28	2.28	2.31	2.33

⇒ best proposed algorithm: SC-IMM3 (with all constraints)

Algorithm comparison

SNR	Measure	Algorithms					
		ASNA	IMM	S-IMM	SC-IMM1	SC-IMM2	SC-IMM3
-6dB	SDR (dB)	5.8	4.4	4.0	4.1	5.2	5.4
	PESQ	2.00	1.22	1.91	1.91	2.01	2.01
+0dB	SDR (dB)	10.7	7.8	9.1	9.2	9.8	9.8
	PESQ	2.44	1.54	2.30	2.30	2.34	2.35
+6dB	SDR (dB)	15.0	9.7	13.0	13.0	12.8	12.9
	PESQ	2.85	1.82	2.62	2.62	2.59	2.62
Mean	SDR (dB)	10.5	7.3	8.7	8.7	9.3	9.4
	PESQ	2.43	1.52	2.28	2.28	2.31	2.33

⇒ better than Durrieu & close of Virtanen in low SNRs

Speech    

Noise   

Text : "Computers are being used to keep branch inventories at more workable levels."

Conclusion

Summary

- Semi-supervised speech separation
- Source/filter model

Contributions

- Weight adaptation method for constraints
- Source/filter coherence constraint for speech
- Good results close to literature in supervised separation

Further research

- Speaker-independant model [Sun and Mysore, 2013]
- Integration of a language model [Mysore and Smaragdis, 2012]
- Integration of a noise adaptation method [Roebel et al., 2015]

Bibliography I

- Gautham J. Mysore and Paris Smaragdis. A non-negative approach to language informed speech separation. In *Latent Variable Analysis and Signal Separation*, pages 356–363. Springer, 2012.
- Tuomas Virtanen, Jort F Gemmeke, and Bhiksha Raj. Active-set Newton algorithm for overcomplete non-negative representations of audio. *IEEE Transactions on Audio, Speech, and Language Processing*, 21(11):2277–2289, 2013.
- Dennis L. Sun and Gautham J. Mysore. Universal speech models for speaker independent single channel source separation. In *IEEE International Conference on Acoustics, Speech, and Signal Processing (ICASSP)*, pages 141–145, 2013.
- François G. Germain and Gautham J. Mysore. Speaker and noise independent online single-channel speech enhancement. In *IEEE International Conference on Acoustics, Speech, and Signal Processing (ICASSP)*, pages 71–75, 2015.
- Luc Le Magoarou, Alexey Ozerov, and Ngoc QK Duong. Text-informed audio source separation. Example-based approach using non-negative matrix partial co-factorization. *Journal of Signal Processing Systems*, pages 1–15, 2014.
- Jean-Louis Durrieu, Alexey Ozerov, Cédric Févotte, Gaël Richard, and Bertrand David. Main instrument separation from stereophonic audio signals using a source/filter model. In *European Signal Processing Conference (EUSIPCO)*, pages 15–19, 2009. URL <http://www.durrieu.ch/phd/eusipco09/>.

Bibliography II

- Nancy Bertin. *Les factorisations en matrices non-négatives. Approches contraintes et probabilistes, application à la transcription automatique de musique polyphonique*. PhD thesis, Télécom ParisTech, 2009.
- Victor Zue, Stephanie Seneff, and James Glass. Speech database development at MIT: TIMIT and beyond. *Speech Communication*, 9(4):351–356, 1990.
- David B. Dean, Sridha Sridharan, Robert J. Vogt, and Michael W. Mason. The QUT-NOISE-TIMIT corpus for the evaluation of voice activity detection algorithms. *Proceedings of Interspeech 2010*, pages 3110–3113, 2010.
- Axel Roebel, Jordi Pons, Marco Liuni, and Mathieu Lagrange. On automatic drum transcription using non-negative matrix deconvolution and Itakura-Saito divergence. In *IEEE International Conference on Acoustics, Speech, and Signal Processing (ICASSP)*, pages 414–418, 2015.