

A General Variational Bayesian Framework for Robust Feature Extraction in Multisource Recordings

Kamil Adiloglu, Emmanuel Vincent

► **To cite this version:**

Kamil Adiloglu, Emmanuel Vincent. A General Variational Bayesian Framework for Robust Feature Extraction in Multisource Recordings. IEEE International Conference on Acoustics, Speech and Signal Processing, Mar 2012, Kyoto, Japan. 2012. <hal-00656613>

HAL Id: hal-00656613

<https://hal.inria.fr/hal-00656613>

Submitted on 4 Jan 2012

HAL is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

A GENERAL VARIATIONAL BAYESIAN FRAMEWORK FOR ROBUST FEATURE EXTRACTION IN MULTISOURCE RECORDINGS

Kamil Adiloğlu Emmanuel Vincent

INRIA, Centre de Rennes - Bretagne Atlantique
Campus de Beaulieu, 35042, Rennes cedex, France
email: {kamil.adiloglu, emmanuel.vincent}@inria.fr

ABSTRACT

We consider the problem of extracting features from individual sources in a multisource audio recording using a general source separation algorithm. The main issue is to estimate and propagate the uncertainty over the separated source signals, so as to robustly estimate the features despite source separation errors. While state-of-the-art techniques estimate the uncertainty in a heuristic manner, we propose to integrate over the parameter space of the source separation algorithm. We apply variational Bayes to estimate the posterior probability of the sources and subsequently derive the expectation of the features by moment matching. Experiments over stereo mixtures of three or four sources show that the proposed method provides the best results in terms of the root mean square (RMS) error on the estimated features.

Index Terms— Bayesian source separation, robust feature extraction

1. INTRODUCTION

Many applications in the field of audio information retrieval are typically solved by extracting features describing the audio content and exploiting them for *e.g.* speech recognition, speaker recognition, cover version detection etc. However, most audio signals consist of a mixture of several sound sources, which have their own characteristics. Applying source separation prior to feature extraction can increase retrieval accuracy. For increased robustness, the uncertainty over the separated sources must be estimated in the complex-valued time-frequency domain and propagated to the features [1].

A heuristic approach is to assume that the uncertainty is proportional to the squared difference between the separated sources and the mixture [1, 2]. In [3], a more principled approach is taken whereby the separated sources are assumed to follow a Gaussian posterior distribution, whose mean and variance are those of the Wiener filter used for separation. Propagation to the features is then achieved either by moment matching [3] or unscented transform [1]. This approach remains mathematically inaccurate however, since the parameters of the Wiener filter are fixed to a certain value instead of being integrated over in a fully Bayesian approach.

In a preliminary study using a simple local Gaussian source model, we proposed a Gibbs sampling algorithm and a variational Bayes (VB) algorithm to address this integration and showed that the latter decreased the RMS error over the resulting Mel frequency cepstral coefficients (MFCC) [4]. In this paper, we extend this approach to the general modeling framework for source separation recently

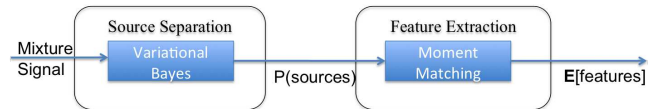


Fig. 1. Flow of the proposed Bayesian source separation and feature extraction approach.

introduced in [5]. This framework generalizes a wide class of existing source separation algorithms, including nonstationarity-based frequency-domain independent component analysis (FDICA) and single- or multi-channel nonnegative matrix factorization (NMF). We propose a VB algorithm to estimate the posterior distribution of the source time-frequency coefficients and subsequently derive the expectation of the features by moment matching. Figure 1 illustrates the workflow of the proposed approach.

This paper is organized as follows. Section 2 introduces the source separation framework. Section 3 presents the proposed VB inference algorithm for the estimation of the posterior distribution of the sources. Section 4 presents the uncertainty propagation method. In Section 5, we evaluate this framework over convolutive mixtures. We conclude in Section 6.

2. GENERAL SOURCE SEPARATION FRAMEWORK

Classically, we address the source separation problem in the time-frequency domain by means of the Short-Time Fourier Transform (STFT). For J sources and I channels, the mixing equation writes

$$\mathbf{x}_{fn} = \mathbf{A}_f \mathbf{s}_{fn} + \epsilon_{fn}, \quad (1)$$

where $\mathbf{x}_{fn} = [x_{1,fn}, \dots, x_{I,fn}]^T$ denotes the mixture STFT coefficients, $\mathbf{s}_{fn} = [s_{1,fn}, \dots, s_{J,fn}]^T$ the source STFT coefficients. $\mathbf{A}_f = [\mathbf{A}_{1,f}, \dots, \mathbf{A}_{J,f}]$ is the complex valued mixing matrix, where $\mathbf{A}_{j,f}$ is the vector of mixing coefficients for source j and ϵ_{fn} is the noise. In this formulation, f is the frequency index, n the time frame index, i the channel index and j the source index. Note that this framework also works for diffuse or reverberated sources by modeling each source as a subspace spanned by several point sources [5].

We adopt a local Gaussian model [5] for the source coefficients. We set a zero-mean complex Gaussian prior over the source coefficients $s_{j,fn}$ with variance $v_{j,fn}$:

$$s_{j,fn} \sim \mathcal{N}(0, v_{j,fn}). \quad (2)$$

The source variances $v_{j,fn}$, which encode the spectral power are decomposed via an excitation-filter model [5]:

$$v_{j,fn} = v_{j,fn}^{\text{ex}} v_{j,fn}^{\text{ft}}. \quad (3)$$

The excitation spectral power $v_{j,fn}^{\text{ex}}$ is decomposed into characteristic spectral patterns modulated by time activation coefficients. Finally, the characteristic spectral patterns are defined as the sum of narrowband spectral patterns $w_{j,fl}^{\text{ex}}$ with associated weights $u_{j,lk}^{\text{ex}}$. Similarly, the time activation coefficients are represented as a sum of time-localized patterns $h_{j,mn}^{\text{ex}}$ with their weights $g_{j,km}^{\text{ex}}$. The same decomposition applies to the filter spectral power $v_{j,fn}^{\text{ft}}$. This framework makes it possible to incorporate a wide range of constraints about the sources. For instance, harmonicity can be enforced by choosing $w_{j,fl}^{\text{ex}}$ as narrowband harmonic spectra and letting the spectral envelope and the active pitches be inferred from the data via the other parameters. For more details about how to constrain spectral and temporal structures, see [5]. As a result, the complete factorization scheme is as follows:

$$v_{j,fn}^{\text{ex}} = \sum_{k=1}^{K_j^{\text{ex}}} \sum_{m=1}^{M_j^{\text{ex}}} \sum_{l=1}^{L_j^{\text{ex}}} h_{j,mn}^{\text{ex}} g_{j,km}^{\text{ex}} u_{j,lk}^{\text{ex}} w_{j,fl}^{\text{ex}}, \quad (4)$$

$$v_{j,fn}^{\text{ft}} = \sum_{k'=1}^{K_j^{\text{ft}}} \sum_{m'=1}^{M_j^{\text{ft}}} \sum_{l'=1}^{L_j^{\text{ft}}} h_{j,m'n}^{\text{ft}} g_{j,k'm'}^{\text{ft}} u_{j,l'k'}^{\text{ft}} w_{j,fl'}^{\text{ft}}. \quad (5)$$

3. PROPOSED VARIATIONAL INFERENCE ALGORITHM

For the separation of the sources, we propose a VB algorithm [6].

3.1. The Bayesian Approach

Let us denote the set of all model parameters as

$$\theta = \{\mathbf{A}_f, h_{j,mn}^{\text{ex}}, g_{j,km}^{\text{ex}}, u_{j,lk}^{\text{ex}}, w_{j,fl}^{\text{ex}}, h_{j,m'n}^{\text{ft}}, g_{j,k'm'}^{\text{ft}}, u_{j,l'k'}^{\text{ft}}, w_{j,fl'}^{\text{ft}}\}. \quad (6)$$

For the sake of simplicity, we define all prior probabilities to be non-informative, even-though informative priors could be defined as well. Hence, \mathbf{A}_f follows a flat prior and the non-negative factors $\{h_{j,mn}^{\text{ex}}, g_{j,km}^{\text{ex}}, \dots, u_{j,l'k'}^{\text{ft}}, w_{j,fl'}^{\text{ft}}\}$ follow a Jeffreys prior.

Finally, we assume Gaussian noise ϵ_{fn} with fixed variance $\sigma_b^2 \mathbf{I}$, where \mathbf{I} is the identity matrix.

Under these assumptions, we aim to estimate the posterior probability of the source coefficients, which is given by

$$p(\mathbf{s}|\mathbf{x}) \propto \int p(\mathbf{x}|\mathbf{s}, \theta) p(\mathbf{s}|\theta) p(\theta) d\theta. \quad (7)$$

This integral is intractable. By contrast with previous approaches [1], which consist of estimating the maximum likelihood (ML) value $\hat{\theta}$ for the parameters and considering $p(\mathbf{x}|\mathbf{s}, \hat{\theta})$, we resort to more accurate VB inference [6].

3.2. Algorithm

VB minimizes the Kullback-Leibler (KL) divergence between the true posterior distribution $p(\mathbf{s}, \theta|\mathbf{X})$ and some approximation $q(\mathbf{s}, \theta)$, which is typically specified by assuming some factorization. Unfortunately, direct application of this rule does not yield a closed form solution. Therefore we pursue a decomposition of the sources as in [7] for the above model.

A source $s_{j,fn}$ is sub-divided into $\Lambda_j = \Lambda_j^{\text{ex}} \Lambda_j^{\text{ft}}$ sub-components such that $\Lambda_j^{\text{ex}} = K_j^{\text{ex}} M_j^{\text{ex}} L_j^{\text{ex}}$ and $\Lambda_j^{\text{ft}} = K_j^{\text{ft}} M_j^{\text{ft}} L_j^{\text{ft}}$. Let us define λ to be a joint index of $\{k, m, l\}$ and λ' is a joint index of $\{k', m', l'\}$. We then decompose the sources as

$$s_{j,fn} = \sum_{\lambda} \sum_{\lambda'} c_{j,fn,\lambda,\lambda'}, \quad (8)$$

where each sub-component $c_{j,fn,\lambda,\lambda'}$ follows a complex Gaussian distribution

$$c_{j,fn,\lambda,\lambda'} \sim \mathcal{N}(0, v_{j,fn,\lambda,\lambda'}), \quad (9)$$

with

$$v_{j,fn,\lambda,\lambda'} = h_{j,mn}^{\text{ex}} g_{j,km}^{\text{ex}} u_{j,lk}^{\text{ex}} w_{j,fl}^{\text{ex}} h_{j,m'n}^{\text{ft}} g_{j,k'm'}^{\text{ft}} u_{j,l'k'}^{\text{ft}} w_{j,fl'}^{\text{ft}}. \quad (10)$$

The mixing equation for the source sub-components can be rewritten as

$$\mathbf{x}_{fn} = \mathbf{A}_f^+ \mathbf{c}_{fn} + \epsilon_f, \quad (11)$$

where \mathbf{A}_f^+ indicates the extended mixing matrix where the elements $a_{i,j,f}$ of \mathbf{A}_f are repeated Λ_j times.

Assuming the following factorization of the variational approximation of the posterior

$$q(\mathbf{c}, \theta) = \left(\prod_{f,n} q(\mathbf{c}_{fn}) \right) \left(\prod_f q(\mathbf{A}_f) \right) \left(\prod_{j,m,n} q(h_{j,mn}^{\text{ex}}) \cdots \prod_{j,f,l'} q(w_{j,fl'}^{\text{ft}}) \right), \quad (12)$$

we obtain the optimal factors of \mathbf{c}_{fn} , \mathbf{A}_f and $w_{j,fl}^{\text{ex}}$ as in the following. Due to space limitations, we give the update equations of $w_{j,fl}^{\text{ex}}$ among all of the NMF components. However the rest of the factors follow a similar shape:

$$q^*(\mathbf{c}_{fn}) \sim \mathcal{N}(\boldsymbol{\mu}_{\mathbf{c},fn}, \mathbf{R}_{\mathbf{cc},fn}), \quad (13)$$

$$q^*(\mathbf{A}_f) \sim \mathcal{N}(\boldsymbol{\mu}_{\mathbf{A},f}, \mathbf{R}_{\mathbf{AA},f}), \quad (14)$$

$$q^*(w_{j,fl}^{\text{ex}}) \sim \text{Inv-Gamma}(\alpha_{w,j,fl}^{\text{ex}}, \beta_{w,j,fl}^{\text{ex}}), \quad (15)$$

where \mathbf{A}_f denotes the column vector obtained by concatenating the rows of the mixing matrix \mathbf{A}_f and

$$\mathbf{R}_{\mathbf{cc},fn} = \left(\mathbf{B}_{\mathbf{c},fn} + (\sigma_b^2 \mathbf{I})^{-1} \sum_i ([\boldsymbol{\mu}_{\mathbf{A},f} \boldsymbol{\mu}_{\mathbf{A},f}^H + \mathbf{R}_{\mathbf{AA},f}]_{ii})^T \right)^{-1}, \quad (16)$$

$$\boldsymbol{\mu}_{\mathbf{c},fn} = \mathbf{R}_{\mathbf{cc},fn} (\boldsymbol{\mu}_{\mathbf{A},f}^+)^H (\sigma_b^2 \mathbf{I})^{-1} \mathbf{x}_{fn}, \quad (17)$$

$$\mathbf{R}_{\mathbf{AA},f} = \left(\frac{1}{\sigma_b^2} \sum_n \text{diag} \left(\underbrace{\mathbf{R}_{\mathbf{s},fn}^T, \dots, \mathbf{R}_{\mathbf{s},fn}^T}_{\text{I times}} \right) \right)^{-1}, \quad (18)$$

$$\boldsymbol{\mu}_{\mathbf{A},f} = \boldsymbol{\Sigma}_{\mathbf{A},f} \left(\frac{1}{\sigma_b^2} \sum_n \mathbf{R}_{\mathbf{xs},fn} \right), \quad (19)$$

$$\alpha_{w,j,fl}^{\text{ex}} = N K_j^{\text{ex}} M_j^{\text{ex}} \Lambda_j^{\text{ft}}, \quad (20)$$

$$\beta_{w,j,fl}^{\text{ex}} = N^5 F^6 (K_j^{\text{ex}})^5 (M_j^{\text{ex}})^5 (L_j^{\text{ex}})^6 (\Lambda_j^{\text{ft}})^5 \sum_{n,km,\lambda'} \frac{(\boldsymbol{\mu}_{\mathbf{c},fn})_{(j,km,\lambda')}^2 + (\mathbf{R}_{\mathbf{cc},fn})_{(j,km,\lambda'),(j,km,\lambda')}}{\beta_{h,j,mn}^{\text{ex}} \beta_{g,j,km}^{\text{ex}} \beta_{u,j,lk}^{\text{ex}} \beta_{h,j,m'n}^{\text{ft}} \beta_{g,j,k'm'}^{\text{ft}} \beta_{u,j,l'k'}^{\text{ft}} \beta_{w,j,fl'}^{\text{ft}}}. \quad (21)$$

In (16), $[\cdot]_{ii}$ denotes the diagonal $J \times J$ block corresponding to channel i . In the same equation, the matrix \mathbf{B}_{c,f_n} writes

$$\mathbf{B}_{c,f_n} = \text{diag} \left(\frac{N^6 F^6 (\Lambda_j^{\text{ex}})^6}{\beta_{w,j,fl}^{\text{ex}} \beta_{u,j,lk}^{\text{ex}} \beta_{g,j,km}^{\text{ex}} \beta_{h,j,mn}^{\text{ex}}} \frac{(\Lambda_j^{\text{ft}})^6}{\beta_{w,j,fl}^{\text{ft}} \beta_{u,j,lk'}^{\text{ft}} \beta_{g,j,k'm'}^{\text{ft}} \beta_{h,j,m'n}^{\text{ft}}} \right)_{j,klm, k'l'm'}. \quad (22)$$

In (19), $\mathbf{R}_{\mathbf{x}s,f_n}$ is given by

$$\mathbf{R}_{\mathbf{x}s,f_n} = [x_{1,f_n} \boldsymbol{\mu}_{s,f_n}^H, \dots, x_{I,f_n} \boldsymbol{\mu}_{s,f_n}^H]^T. \quad (23)$$

Finally, the posterior distribution of the source coefficients is calculated by summing up the corresponding elements of the mean and the covariance of the source sub-components as follows:

$$\mu_{s,j,f_n} = \sum_{\lambda} \sum_{\lambda'} \mu_{c,j,f_n,\lambda,\lambda'}, \quad (24)$$

$$(\mathbf{R}_{\mathbf{ss},f_n})_{j,j'} = \sum_{\kappa} \sum_{\lambda'} (\mathbf{R}_{\mathbf{cc},f_n})_{(j,\lambda,\lambda'),(j',\lambda,\lambda')}. \quad (25)$$

and the second raw moment of the source coefficients \mathbf{R}_{s,f_n} is given by

$$\mathbf{R}_{s,f_n} = \boldsymbol{\mu}_{s,nf} \boldsymbol{\mu}_{s,nf}^H + \mathbf{R}_{\mathbf{ss},f_n}. \quad (26)$$

The VB update equations of the scale parameters of the NMF components $\{\beta_{w,j,fl}^{\text{ex}}, \dots, \beta_{h,j,m'n}^{\text{ft}}\}$ turn out to be identical to the state-of-the-art expectation maximization (EM) updates of the NMF components $\{\hat{w}_{j,fl}^{\text{ex}}, \dots, \hat{h}_{j,m'n}^{\text{ft}}\}$ [7] up to a variable change. For $\beta_{w,j,fl}^{\text{ex}}$ this variable change writes

$$\hat{w}_{j,fl}^{\text{ex}} = \frac{\beta_{w,j,fl}^{\text{ex}}}{N K_j^{\text{ex}} M_j^{\text{ex}} K_j^{\text{ft}} M_j^{\text{ft}} L_j^{\text{ft}}}. \quad (27)$$

Therefore, we replaced these EM updates by the multiplicative updates in [5], which converge much faster.

Equations (16) to (27) depend on each other. After proper initialization, we cycle through these equations by replacing the dependent values with their new estimates.

4. UNCERTAINTY PROPAGATION

We now present a moment matching method for propagating the uncertainty over the source components to the source images and in a following step to the MFCC features.

4.1. Uncertainty Propagation for the Source Images

Due to phase and scale indeterminacies in the source estimates \mathbf{s}_{j,f_n} , we use the spatial source images $\mathbf{y}_{j,f_n} = \mathbf{A}_{j,f} \mathbf{s}_{j,f_n}$ instead for our experiments, which do not suffer from such indeterminacies [5].

Once the posterior distribution of the source coefficients \mathbf{s}_{f_n} has been computed, the posterior distribution of the source images is calculated by propagating the first two moments of the sources to the source images as follows:

$$\boldsymbol{\mu}_{\mathbf{y},j,f_n} = \boldsymbol{\mu}_{\mathbf{A},f}^{\text{post}} \boldsymbol{\mu}_{\mathbf{s},f_n}^{\text{post}} \quad (28)$$

$$\begin{aligned} \mathbf{R}_{\mathbf{y}\mathbf{y},j,f_n} &= \left(\sum_{r,r'} ((\mathbf{R}_{\mathbf{AA},f})_{(ir,i'r')} + (\boldsymbol{\mu}_{\mathbf{A},f})_{(ir)} (\boldsymbol{\mu}_{\mathbf{A},f})_{(i'r')}^H) \right)_{ii'} \\ &+ ((\mathbf{R}_{\mathbf{ss},f_n})_{(r,r')} + (\boldsymbol{\mu}_{\mathbf{s},f_n})_{(r)} (\boldsymbol{\mu}_{\mathbf{s},f_n})_{(r')}^H)_{ii'} \\ &- \boldsymbol{\mu}_{\mathbf{y},j,f_n} \boldsymbol{\mu}_{\mathbf{y},j,f_n}^H. \end{aligned} \quad (29)$$

4.2. Uncertainty Propagation for Feature Extraction

We calculate the expectation of the MFCCs for each source as

$$\mu_{j_n}^{\text{MFCC}} = \int \text{MFCC}(\mathbf{y}_{j1n}) P(\mathbf{y}_{j1n}) d\mathbf{y}_{j1n} \quad (30)$$

where $\mathbf{y}_{j1n} = [y_{j,1fn}]_{f=1\dots F}$ are the STFT coefficients of the first channel of source image j in time frame n . Deterministic calculation without the use of the uncertainty model simply yields $\text{MFCC}(\mathbf{y}_{j1n}) = 20 \mathbf{D} \log_{10}(\mathbf{M}|\mathbf{y}_{j1n}|)$. In this formulation, \mathbf{D} is the DCT matrix and \mathbf{M} is the matrix containing the mel filter coefficients. Note that we chose the scaling so that the MFCCs are expressed in decibels (dB).

In order to obtain the mean estimate of the MFCCs $\mu_{j_n}^{\text{MFCC}}$, we first compute the mean and variance of $|\mathbf{y}_{j1n}|$ using formulae for the Rice distribution in [1] and then propagate them through the logarithm using the moment matching formulae in [8]. Note that this also makes it possible to estimate the variance of the MFCCs and exploit it for *e.g.* classification tasks [1]. In the following, due to space limitations, we simply consider the mean.

5. EXPERIMENTAL EVALUATION

5.1. Data and Algorithmic Settings

We considered the development dataset of the 2008 Signal Separation Evaluation Campaign (SiSEC)¹. This dataset contains synthetic and live recorded convolutive, under-determined, stereo mixtures. There are 32 mixtures of 3 sources and 24 mixtures of 4 sources. Each mixture has a duration of 10 s.

We performed experiments with eight different sets of constraints over the parameters as considered in the experiments section of [5]. These scenarios consist of all combinations of the following possibilities:

- Rank: Each source is either a single point source (1) or a subspace spanned by two point sources (2).
- Spectral Structure: The narrowband spectral patterns $w_{j,fl}^{\text{ex}}$ are either unconstrained (un) or fixed (co) to harmonic and noise-like patterns.
- Temporal Structure: The time-localized patterns $h_{j,mn}^{\text{ex}}$ are either unconstrained (un) or fixed (co) to decreasing exponential patterns.

All other parameters are free. We initialize the mixing matrix \mathbf{A}_j using the direction of arrival (DOA) estimation algorithm proposed in [9]. The noise variance σ_b^2 is initialized to 10^{-6} . Finally, we performed 200 iterations for convergence.

5.2. Evaluation Criteria

In order to assess the impact of source separation on feature extraction, we evaluate the proposed algorithm according to both tasks. Source separation quality is evaluated in terms of the Signal-to-Distortion Ratio (SDR) in [10] between the mean of the estimated source images $\boldsymbol{\mu}_{\mathbf{y},j,f_n}$ and the true source images.

Feature extraction accuracy is evaluated in terms of the RMS error [4] between the estimated $\mu_{j_n}^{\text{MFCC}}$ and the true MFCCs. We ignore the first MFCC coefficient and consider the MFCCs 2 to 20 only.

¹<http://sisecc2008.wiki.irisa.fr/tiki-index.php?page=Under-determined+speech+and+music+mixtures>

	1-un-un	2-un-un	1-co-un	2-co-un	1-un-co	2-un-co	1-co-co	2-co-co
ML	1.58	1.68	1.87	2.07	1.77	1.75	2.28	2.25
VB	1.61	1.72	1.93	2.06	1.84	1.80	2.49	2.29

Table 1. SDR in dB achieved by ML or VB source separation over all mixtures.

	1-un-un		2-un-un		1-co-un		2-co-un		1-un-co		2-un-co		1-co-co		2-co-co	
	ML	VB	ML	VB	ML	VB	ML	VB	ML	VB	ML	VB	ML	VB	ML	VB
det	7.55	7.45	8.62	8.55	7.35	7.28	8.01	7.96	7.43	7.32	8.41	8.34	7.67	7.56	9.28	8.74
mm	6.66	6.63	6.85	6.84	6.72	6.69	6.82	6.80	6.59	6.54	6.76	6.76	6.61	6.57	7.23	6.92

Table 2. Total RMS error in dB for the MFCCs 2-20 obtained by ML- or VB-based source separation followed by deterministic (det) or moment matching (mm) feature extraction over all mixtures.

5.3. Results

Table 1 shows the source separation performance of VB compared to that of the state-of-the-art ML method in [5]. ML and VB perform similarly in almost all of the eight configurations. However VB is 0.2 dB better than the state-of-the-art ML method and yields the best performance for the 1 point source, spectrally and temporally constrained model (1-co-co) with 2.5 dB SDR. The baseline binary masking method [9] yields 0.95 dB SDR.

Table 2 shows the total RMS error in dB over the MFCCs obtained either by deterministic computation or by moment matching for both VB and ML algorithms. As one can see, VB based MFCC estimation performs 0.05 dB better than the ML based estimation. Besides, the moment matching method outperforms deterministic MFCC estimation in all configurations with around 0.9 dB. Again, VB performs best for the RMS error by 6.54 dB for the 1 point source, spectrally unconstrained, temporally constrained model (1-un-co). The baseline binary masking method [9] performs significantly worse than both VB and ML algorithms and yields 17.25 dB RMS error.

6. CONCLUSION

In this paper, we presented a general, fully Bayesian source separation algorithm and an uncertainty propagation algorithm for the computation of the expectation of the MFCCs of individual sources in multisource recordings.

This algorithm provides a fundamental breakthrough towards mathematically rigorous estimation of uncertainty for robust feature extraction. The resulting MFCC coefficients are slightly more accurate than those obtained via the standard ML method, proving that the ML method provides a reasonable approximation, but that it is possible to obtain more accurate estimates.

In the future, we will seek to improve the tightness of the variational bound (and thereby the accuracy of the resulting features) and exploit the variance of the estimated features for uncertainty decoding.

7. REFERENCES

- [1] D. Kolossa, R. F. Astudillo, E. Hoffmann, and R. Orglmeister, "Independent component analysis and time-frequency masking for speech recognition in multitalker conditions," *EURASIP Journal on Audio, Speech, and Music Processing*, vol. 2010, Article ID 651420, 2010.
- [2] M. Delcroix, T. Nakatani, and S. Watanabe, "Static and dynamic variance compensation for recognition of reverberant

speech with dereverberation preprocessing," *IEEE Transactions on Audio, Speech and Language Processing*, vol. 17, no. 2, pp. 324–334, 2009.

- [3] R. F. Astudillo and R. Orglmeister, "A MMSE estimator in mel-cepstral domain for robust large vocabulary automatic speech recognition using uncertainty propagation," in *Proceedings of Interspeech*, 2010, pp. 713–716.
- [4] K. Adiloglu and E. Vincent, "An uncertainty estimation approach for the extraction of source features in multisource recordings," in *Proceedings of 19th European Signal Processing Conference (EUSIPCO)*, 2011, pp. 1663–1667.
- [5] A. Ozerov, E. Vincent, and F. Bimbot, "A general flexible framework for the handling of prior information in audio source separation," *IEEE Transactions on Audio, Speech and Language Processing*, to appear, <http://hal.inria.fr/inria-00536917/PDF/RR-7453.pdf>.
- [6] C. M. Bishop, *Pattern Recognition and Machine Learning*, Springer, 2006.
- [7] A. Ozerov and C. Févotte, "Multichannel nonnegative matrix factorization in convolutive mixtures for audio source separation," *IEEE Transactions on Audio, Speech and Language Processing*, vol. 18, no. 3, pp. 550–563, 2010.
- [8] M. J. F. Gales, *Model-Based Techniques for Noise Robust Speech Recognition*, Ph.D. thesis, Gonville and Caius College, University of Cambridge, 1995.
- [9] C. Blandin, E. Vincent, and A. Ozerov, "Multi-source TDOA estimation using SNR-based angular spectra," in *Proceedings of IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, 2011, pp. 2616–2619.
- [10] E. Vincent, R. Gribonval, and C. Févotte, "Performance measures in blind audio source separation," *IEEE Transactions on Audio, Speech and Language Processing*, vol. 14, no. 4, pp. 1462–1469, 2006.