

# BAYESIAN SINGING-VOICE SEPARATION

Guan-Xiang Wang, Po-Kai Yang, Chung-Chien Hsu and Jen-Tzung Chien

Department of Electrical and Computer Engineering, National Chiao Tung University, Taiwan  
gxwang@chien.cm.nctu.edu.tw, {niceallen.cm01g, chien.cm97g, jtchien}@nctu.edu.tw

## ABSTRACT

This paper presents a Bayesian nonnegative matrix factorization (NMF) approach to extract singing voice from background music accompaniment. Using this approach, the reconstruction error based on NMF is represented by a Poisson distribution and the NMF parameters, consisting of basis and weight matrices, are characterized by the exponential priors. A variational Bayesian expectation-maximization algorithm is developed to learn variational parameters and model parameters for monaural source separation. A clustering algorithm is performed to establish two groups of bases: one is for singing voice and the other is for background music. Model complexity is controlled by adaptively selecting the number of bases for different mixed signals according to the variational lower bound.

## 1. INTRODUCTION

This paper presents a new model-based singing-voice separation. The novelties of this paper are twofold. The first one is to develop Bayesian approach to unsupervised singing-voice separation. Model uncertainty is compensated to improve the performance of source separation of vocal signal and background accompaniment signal. Number of bases is adaptively determined from the mixed signal according to the variational lower bound of the marginal likelihood over NMF basis and weight matrices. The second one is the theoretical contribution in Bayesian NMF. We construct a new Bayesian NMF where the modeling error in NMF is drawn from Poisson distribution and the model parameters are characterized by exponential distributions.

## 2. NEW BAYESIAN NONNEGATIVE MATRIX FACTORIZATION

This study aims to find an analytical solution to full Bayesian NMF by considering all dependencies of variational lower bound on regularization parameters. Regularization parameters are optimally estimated.

This document is licensed under the Creative Commons

Attribution-Noncommercial-Share Alike 3.0 License.

<http://creativecommons.org/licenses/by-nc-sa/3.0/>

© 2014 The Authors.

## 2.1 Bayesian Objectives

We adopt the Poisson distribution as likelihood function and the exponential distribution as *conjugate prior* for NMF parameters  $B_{mk}$  and  $W_{kn}$  with hyperparameters  $\lambda_{mk}^b$  and  $\lambda_{kn}^w$ , respectively. Maximum *a posteriori* (MAP) estimates of parameters  $\Theta = \{\mathbf{B}, \mathbf{W}\}$  are obtained by maximizing the posterior distribution or minimizing  $-\log p(\mathbf{B}, \mathbf{W}|\mathbf{X})$  which is arranged as a regularized KL divergence between  $\mathbf{X}$  and  $\mathbf{B}\mathbf{W}$

$$D_{\text{KL}}(\mathbf{X}||\mathbf{B}\mathbf{W}) + \sum_{m,k} \lambda_{mk}^b B_{mk} + \sum_{k,n} \lambda_{kn}^w W_{kn} \quad (1)$$

where the terms independent of  $B_{mk}$  and  $W_{kn}$  are treated as constants. Notably, the regularization terms (2nd and 3rd terms) in this objective are nonnegative and seen as the  $\ell_1$  regularizers [2] which are controlled by hyperparameters  $\{\lambda_{mk}^b, \lambda_{kn}^w\}$ . These regularizers impose sparseness in the estimated MAP parameters.

However, MAP estimates are seen as point estimates. The randomness of parameters is not considered in model construction. To conduct full Bayesian treatment, BNMF is developed by maximizing the marginal likelihood  $p(\mathbf{X}|\Theta)$  over latent variables  $\mathbf{Z}$  as well as NMF parameters  $\{\mathbf{B}, \mathbf{W}\}$

$$\int \sum_{\mathbf{Z}} p(\mathbf{X}|\mathbf{Z}, \mathbf{B}, \mathbf{W}) p(\mathbf{Z}|\mathbf{B}, \mathbf{W}) p(\mathbf{B}, \mathbf{W}|\Theta) d\mathbf{B}d\mathbf{W} \quad (2)$$

and estimating the sparsity-controlled hyperparameters or regularization parameters  $\Theta = \{\lambda_{mk}^b, \lambda_{kn}^w\}$ . The resulting evidence function is meaningful to act as an objective for model selection which balances the tradeoff between data fitness and model complexity [1]. In the singing-voice separation based on NMF, this objective is used to judge which number of bases  $K$  should be selected.

## 2.2 Variational Bayesian Inference

The variational Bayesian expectation-maximization (VB-EM) algorithm is developed to implement Poisson-Exponential BNMF. VB-EM algorithm applies the Jensen's inequality and maximizes the lower bound of the logarithm of marginal likelihood

$$\log p(\mathbf{X}|\Theta) \geq \int \sum_{\mathbf{Z}} q(\mathbf{Z}, \mathbf{B}, \mathbf{W}) \log \frac{p(\mathbf{X}, \mathbf{Z}, \mathbf{B}, \mathbf{W}|\Theta)}{q(\mathbf{Z}, \mathbf{B}, \mathbf{W})} \times d\mathbf{B}d\mathbf{W} = \mathbb{E}_q[\log p(\mathbf{X}, \mathbf{Z}, \mathbf{B}, \mathbf{W}|\Theta)] + H[q(\mathbf{Z}, \mathbf{B}, \mathbf{W})] \quad (3)$$

where  $H[\cdot]$  is an entropy function. The factorized variational distribution  $q(\mathbf{Z}, \mathbf{B}, \mathbf{W}) = q(\mathbf{Z})q(\mathbf{B})q(\mathbf{W})$  is

assumed to approximate the true posterior distribution  $p(\mathbf{Z}, \mathbf{B}, \mathbf{W} | \mathbf{X}, \Theta)$ .

### 2.2.1 VB-E Step

In VB-E step, a general solution to variational distribution  $q_j$  of an individual latent variable  $j \in \{\mathbf{Z}, \mathbf{B}, \mathbf{W}\}$  is obtained by [1]

$$\log \hat{q}_j \propto \mathbb{E}_{q(i \neq j)} [\log p(\mathbf{X}, \mathbf{Z}, \mathbf{B}, \mathbf{W} | \Theta)]. \quad (4)$$

Given the variational distributions defined by

$$\begin{aligned} q(B_{mk}) &\propto \text{Gam}(B_{mk}; \alpha_{mk}^b, \beta_{mk}^b) \\ q(W_{kn}) &\propto \text{Gam}(W_{kn}; \alpha_{kn}^w, \beta_{kn}^w) \\ q(Z_{mkn}) &\propto \text{Mult}(Z_{mkn}; P_{mkn}) \end{aligned} \quad (5)$$

the variational parameters  $\{\alpha_{mk}^b, \beta_{mk}^b, \alpha_{kn}^w, \beta_{kn}^w, P_{mkn}\}$  in three distributions are estimated by

$$\begin{aligned} \hat{\alpha}_{mk}^b &= 1 + \sum_n \langle Z_{mkn} \rangle, \quad \hat{\beta}_{mk}^b = \left( \sum_n \langle W_{kn} \rangle + \lambda_{mk}^b \right)^{-1} \\ \hat{\alpha}_{kn}^w &= 1 + \sum_m \langle Z_{mkn} \rangle, \quad \hat{\beta}_{kn}^w = \left( \sum_k \langle B_{mk} \rangle + \lambda_{kn}^w \right)^{-1} \\ \hat{P}_{mkn} &= \frac{\exp(\langle \log B_{mk} \rangle + \langle \log W_{kn} \rangle)}{\sum_j \exp(\langle \log B_{mj} \rangle + \langle \log W_{jn} \rangle)} \end{aligned} \quad (6)$$

where the expectation function  $\mathbb{E}_q[\cdot]$  is replaced by  $\langle \cdot \rangle$  for simplicity. By substituting the variational distribution into Eq. (3), the variational lower bound is obtained by

$$\begin{aligned} \mathcal{B}_L &= - \sum_{m,n,k} \langle B_{mk} \rangle \langle W_{kn} \rangle \\ &+ \sum_{m,n} (-\log \Gamma(X_{mn} + 1) - \sum_k \langle Z_{mkn} \rangle \log \hat{P}_{mkn}) \\ &+ \sum_{m,k} \langle \log B_{mk} \rangle \sum_n \langle Z_{mkn} \rangle + \sum_{k,n} \langle \log W_{kn} \rangle \sum_m \langle Z_{mkn} \rangle \\ &+ \sum_{m,k} (\log \lambda_{mk}^b - \lambda_{mk}^b \langle B_{mk} \rangle) + \sum_{k,n} (\log \lambda_{kn}^w - \lambda_{kn}^w \langle W_{kn} \rangle) \\ &+ \sum_{m,k} (-\langle \hat{\alpha}_{mk}^b \rangle - 1) \Psi(\langle \hat{\alpha}_{mk}^b \rangle) + \log \langle \hat{\beta}_{mk}^b \rangle + \langle \hat{\alpha}_{mk}^b \rangle + \log \Gamma(\langle \hat{\alpha}_{mk}^b \rangle) \\ &+ \sum_{k,n} (-\langle \hat{\alpha}_{kn}^w \rangle - 1) \Psi(\langle \hat{\alpha}_{kn}^w \rangle) + \log \langle \hat{\beta}_{kn}^w \rangle + \langle \hat{\alpha}_{kn}^w \rangle + \log \Gamma(\langle \hat{\alpha}_{kn}^w \rangle) \end{aligned} \quad (7)$$

where  $\Psi(\cdot)$  is the derivative of the log gamma function, and is known as a digamma function.

### 2.2.2 VB-M Step

In VB-M step, the optimal regularization parameters  $\Theta = \{\lambda_{mk}^b, \lambda_{kn}^w\}$  are derived by maximizing Eq. (7) with respect to  $\Theta$  and yielding

$$\begin{aligned} \frac{\partial \mathcal{B}_L}{\partial \lambda_{mk}^b} &= \frac{1}{\lambda_{mk}^b} - \langle B_{mk} \rangle + \frac{\partial \log \beta_{mk}^b}{\partial \lambda_{mk}^b} = 0 \\ \frac{\partial \mathcal{B}_L}{\partial \lambda_{kn}^w} &= \frac{1}{\lambda_{kn}^w} - \langle W_{kn} \rangle + \frac{\partial \log \beta_{kn}^w}{\partial \lambda_{kn}^w} = 0. \end{aligned} \quad (8)$$

Accordingly, the solution to BNMF hyperparameters is derived by solving a quadratic equation where nonnegative

constraint is considered to find positive values of hyperparameters by

$$\begin{aligned} \hat{\lambda}_{mk}^b &= \frac{1}{2} \left( - \sum_n \langle W_{kn} \rangle + \sqrt{\left( \sum_n \langle W_{kn} \rangle \right)^2 + 4 \frac{\sum_n \langle W_{kn} \rangle}{\langle B_{mk} \rangle}} \right) \\ \hat{\lambda}_{kn}^w &= \frac{1}{2} \left( - \sum_m \langle B_{mk} \rangle + \sqrt{\left( \sum_m \langle B_{mk} \rangle \right)^2 + 4 \frac{\sum_m \langle B_{mk} \rangle}{\langle W_{kn} \rangle}} \right) \end{aligned} \quad (9)$$

where  $\langle B_{mk} \rangle = \alpha_{mk}^b \beta_{mk}^b$  and  $\langle W_{kn} \rangle = \alpha_{kn}^w \beta_{kn}^w$  are obtained as the means of gamma distributions.

## 2.3 Poisson-Exponential Bayesian NMF

In this study, total number of basis vectors  $K$  is adaptively selected for individual mixed signal according to the variational lower bound in Eq. (7) with the converged variational parameters  $\{\hat{\alpha}_{mk}^b, \hat{\beta}_{mk}^b, \hat{\alpha}_{kn}^w, \hat{\beta}_{kn}^w, \hat{P}_{mkn}\}$  and model parameters  $\{\hat{\lambda}_{mk}^b, \hat{\lambda}_{kn}^w\}$ .

Considering the pairs of likelihood function and prior distribution in NMF, the proposed method is also called the Poisson-Exponential BNMF.

## 2.4 Unsupervised Singing-Voice Separation

We implemented the unsupervised singing-voice separation where total number of bases ( $K$ ) and the grouping of these bases into vocal source and music source were both learned from test data in an unsupervised way. We conduct NMF-based clustering for the proposed BNMF method. To do so, we transformed the basis vectors  $\mathbf{B}$  into Mel-scaled spectrum to form the Mel-scaled basis matrix. ML-NMF was applied to factorize this Mel-scaled basis matrix into two matrices  $\tilde{\mathbf{B}}$  of size  $N$ -by-2 and  $\tilde{\mathbf{W}}$  of size 2-by- $K$ . The soft mask scheme based on Wiener gain was applied to smooth the separation of  $\mathbf{B}$  into basis vectors for vocal signal and music signal. This same soft mask was performed for the separation of mixed signal  $X$  into vocal signal and music signal based on the K-means clustering and NMF clustering. Finally, the separated singing voice and music accompaniment signals were obtained by the overlap-and-add method using the original phase.

## 3. CONCLUSIONS

We proposed a new unsupervised Bayesian nonnegative matrix factorization approach to extract the singing voice from background music accompaniment and illustrated the novelty on an analytical and true optimum solution to the Poisson-Exponential BNMF. Through the VB-EM inference procedure, the proposed method automatically selected different number of bases to fit various experimental conditions.

## 4. REFERENCES

- [1] C. M. Bishop. *Pattern Recognition and Machine Learning*. Springer Science, 2006.
- [2] R. Tibshirani. Regression shrinkage and selection via the lasso. *Journal of the Royal Statistical Society. Series B*, 58(1):267–288, 1996.