**Beyond NMF:** Time-Domain Audio Source Separation without Phase Reconstruction

Kazuyoshi Yoshii (AIST, Japan)
Ryota Tomioka (Univ. of Tokyo , Japan)
Daichi Mochihashi (ISM , Japan)
Masataka Goto (AIST , Japan)

2013/11/06 ISMIR 2013

#### **Summary**

# If you do not care about computational time, you could forget NMF and use our method PSDTF for audio source separation

We provide MATLAB source codes!

• Our goal: high-quality source separation of monaural audio signals

#### Time-domain separation is equivalent to

#### frequency-domain separation



- Our goal: high-quality source separation of monaural audio signals
  - Nonnegative matrix factorization (NMF) is very popular
    - Frequency-domain decomposition for power spectrogram



- Our goal: high-quality source separation of monaural audio signals
  - Nonnegative matrix factorization (NMF) is very popular
    - Frequency-domain decomposition for power spectrogram
    - Ignore phase → Cannot separate time-domain signals well



• Our goal: high-quality source separation of monaural audio signals

#### Our PSDTF can solve this problem by time-domain separation!



• Our goal: high-quality source separation of monaural audio signals

#### Time-domain PSDTF has an equivalent form of frequency-domain PSDTF considering the phase



- A nonnegative matrix (power spectrogram X ) is approximated by the product of two nonnegative matrices WH
  - Given the observation matrix X, NMF estimates W and H so that their product Y can be similar to X



 <u>Nonnegative vectors</u> are approximated by convex combinations of fewer <u>nonnegative vectors</u>



A series of nonnegative observed vectors along the time

Each vector corresponds to the power spectrogram of a locally observed mixture signal at each frame

 <u>Nonnegative vectors</u> are approximated by convex combinations of fewer <u>nonnegative vectors</u>



 <u>Nonnegative vectors</u> are approximated by convex combinations of fewer <u>nonnegative vectors</u>



The goal of NMF is to estimate basis and activation vectors so that the total reconstruction error can be minimized





The observed PSD matrices are a set of local covariance matrices calculated from a locally observed mixture signal at each frame

By calculating the outer product of the local signal, we can make the PSD matrix  $X_n$ 





The goal of PSDTF is to estimate basis matrices and activation vectors so that the total reconstruction error can be minimized



 <u>Nonnegative vectors</u> are approximated by convex combinations of fewer <u>nonnegative vectors</u>





- Wiener filtering based on NMF results  $s_{kmn} = y_{kmn} y_{mn}^{-1} s_{mn}$ 
  - <u>Element-wise</u> decomposition of <u>nonnegative</u> power spectrogram



#### As shown in this equation, Wiener filtering can be regarded as a simple power partition function

The power of each element  $s_{mn}$  is distributed into different sources  $s_{1mn}$ ,  $s_{2mn}$ ,  $s_{3mn}$ 

- Wiener filtering based on NMF results  $s_{kmn} = y_{kmn} y_{mn}^{-1} s_{mn}$ 
  - <u>Element-wise</u> decomposition of <u>nonnegative</u> power spectrogram
  - Ignore phases and treat frequency bins independently



- Wiener filtering based on NMF results  $s_{kmn} = y_{kmn} y_{mn}^{-1} s_{mn}$ 
  - <u>Element-wise</u> decomposition of <u>nonnegative</u> power spectrogram
  - Ignore phases and treat frequency bins independent



- Wiener filtering based on PSDTF results  $s_{kn} = Y_{kn}Y_n^{-1}s_n$ 
  - <u>Vector-wise</u> decomposition of <u>complex</u> spectrogram



- Wiener filtering based on PSDTF results  $m{s}_{kn} = m{Y}_{kn} m{Y}_n^{-1} m{s}_n$
- The complex spectrogram of the observed mixture signal is decomposed in a vector-wise/frame-wise manner unlike NMF



- Wiener filtering based on PSDTF results  $m{s}_{kn} = m{Y}_{kn} m{Y}_n^{-1} m{s}_n$ 
  - <u>Vector-wise</u> decomposition of <u>complex</u> spectrogram
  - Consider phases and correlations between frequency bins



• Wiener filtering based on PSDTF results  $s_{kn} = Y_{kn}Y_n^{-1}s_n$ PSDTF considers the correlations between frequency bins! Harmonic structure has strong correlations between harmonics



- Wiener filtering based on PSDTF results  $m{s}_{kn} = m{Y}_{kn}m{Y}_n^{-1}m{s}_n$ 
  - <u>Vector-wise</u> decomposition of <u>complex</u> spectrogram
  - Consider phases and correlations between frequency bins



- Wiener filtering based on NMF results  $s_{kmn} = y_{kmn} y_{mn}^{-1} s_{mn}$ 
  - <u>Element-wise</u> decomposition of <u>nonnegative</u> power spectrogram
  - Ignore phases and treat frequency bins independently



- Wiener filtering based on PSDTF results  $m{s}_{kn} = m{Y}_{kn} m{Y}_n^{-1} m{s}_n$ 
  - <u>Vector-wise</u> decomposition of <u>complex</u> spectrogram
  - Consider phases and correlations between frequency bins



#### **PSD Matrices**

## I intentionally skipped a very important explanation



#### Gaussian Process (GP)

- A stochastic process over <u>continuous time</u>
  - GP can be interpreted as a probability distribution p(f)over continuous functions  $f: t \rightarrow x$



#### Gaussian Process (GP)

- A stochastic process over <u>continuous time</u>
  - GP can be interpreted as a probability distribution p(f)over continuous functions  $f: t \rightarrow x$



#### Three functions are stochastically sampled from the same GP

#### **Gaussian Process (GP)**

- A stochastic process over <u>continuous time</u>
  - GP can be interpreted as a probability distribution p(f)over continuous functions  $f: t \to x$



#### **How to Understand GP?**

Gaussian process is characterized by its marginal distribution

- A marginal distribution over any M points can be Gaussian



#### **How to Understand GP?**

Gaussian process is characterized by its marginal distribution

– A marginal distribution over any M points can be Gaussian



#### **Example: Identity Kernel**

- Stationary Gaussian white noise can be generated
  - Marginal:  $\boldsymbol{x} \sim \mathcal{N}(\boldsymbol{0}, \boldsymbol{K})$   $\boldsymbol{f} \subset \boldsymbol{x} = [t_1, t_2, \cdots, t_M]^T$  $\boldsymbol{x} = [x_1, x_2, \cdots, x_M]^T$



#### **Example: Squared Exponential Kernel (RBF)**

• Smooth functions can be generated - Marginal:  $\boldsymbol{x} \sim \mathcal{N}(\boldsymbol{0}, \boldsymbol{K})$   $f \searrow \boldsymbol{x} = [x_1, x_2, \cdots, x_M]^T$ 



#### **Example: Periodic Kernel (Stripe Pattern)**

• Periodic functions can be generated - Marginal:  $\boldsymbol{x} \sim \mathcal{N}(\boldsymbol{0}, \boldsymbol{K})$   $f \longrightarrow \boldsymbol{x} = [x_1, x_2, \cdots, x_M]^T$ 





#### NMF vs. PSDTF

- Nonnegative Matrix Factorization (NMF)
  - Nonnegative vectors  $\rightarrow$  sums of nonnegative vectors
  - Bregman divergence
    - Itakura-Saito (IS) divergence  $\mathcal{D}_{IS}(\boldsymbol{x}_n | \boldsymbol{y}_n) = \sum_m \left( -\log x_{mn} y_{mn}^{-1} + x_{mn} y_{mn}^{-1} - 1 \right)$
    - Kullback-Leibler (KL) divergence  $\mathcal{D}_{\text{KL}}(\boldsymbol{x}_n|\boldsymbol{y}_n) = \sum_m \left(x_{mn}\log x_{mn}y_{mn}^{-1} x_{mn} + y_{mn}\right)$
- Positive Semidefinite Tensor Factorization (PSDTF)
  - PSD matrices  $\rightarrow$  sums of PSD matrices
- $oldsymbol{X}_n pprox oldsymbol{Y}_n = \sum_{k=1}^K h_{kn} oldsymbol{V}_k$

 $oldsymbol{x}_n pprox oldsymbol{y}_n = \sum h_{kn} oldsymbol{w}_k$ 

k=1

- Bregman "matrix" divergence
  - Log-Determinant (LD) divergence  $\mathcal{D}_{\text{LD}}(\boldsymbol{X}_n|\boldsymbol{Y}_n) = -\log \left|\boldsymbol{X}_n\boldsymbol{Y}_n^{-1}\right| + \operatorname{tr}\left(\boldsymbol{X}_n\boldsymbol{Y}_n^{-1}\right) - M$
  - von Neumann (vN) divergence  $\mathcal{D}_{vN}(\boldsymbol{X}_n|\boldsymbol{Y}_n) = \operatorname{tr}\left(\boldsymbol{X}_n\log \boldsymbol{X}_n - \boldsymbol{X}_n\log \boldsymbol{Y}_n - \boldsymbol{X}_n + \boldsymbol{Y}_n\right)$

### **Frequency-domain PSDTF**

• PSDTF is a natural extension of NMF



#### **Comparative Evaluation**

- Source separation performance on synthetic sounds
  - 3 bases (K=3)
  - 3 instruments (from RWC Music Database)
    - Piano / Electric guitar (decaying) / Clarinet (sustaining)
  - BSS Eval Toolbox [Vincent2006] was used



#### **Experimental Results**

• LD-PSDTF outperformed KL-NMF and IS-NMF



### **Conclusion and Future Work**

- Positive Semidefinite Tensor Factorization (PSDTF)
  - A natural extension of Nonnegative Matrix Factorization (NMF)
  - Nonparametric Bayesian extension [Yoshii et al. ICML 2013]
    - Automatically optimize the number of bases

If you do not care about computational time, you could forget NMF and use our method PSDTF for audio source separation

*We provide MATLAB source codes! (2-Clause BSD License)* http://staff.aist.go.jp/k.yoshii/psdtf/ or the last tweet of @MasatakaGoto **Beyond NMF:** Time-Domain Audio Source Separation without Phase Reconstruction

Kazuyoshi Yoshii (AIST, Japan)
Ryota Tomioka (Univ. of Tokyo , Japan)
Daichi Mochihashi (ISM , Japan)
Masataka Goto (AIST , Japan)

2013/11/06 ISMIR 2013

### We Cannot Ignore Computational Time

- Current issue: Big computational cost!
  - Comparison of computational cost
    - #samples: M=512 (frame size: 32ms)
    - #frames: N=1400
    - **#bases:** K=3 (Mixtures of 3 sounds)

(duration: 14s)

- MATLAB / Intel Xeon 3.4GHz
- LD-PSDTF: 600s/iteration  $O(M^3NK)$ - 16.7 hours (60000s) / 100iterations
- IS-NMF: 0.1s/iteration - 10s / 100iterations

O(MNK)