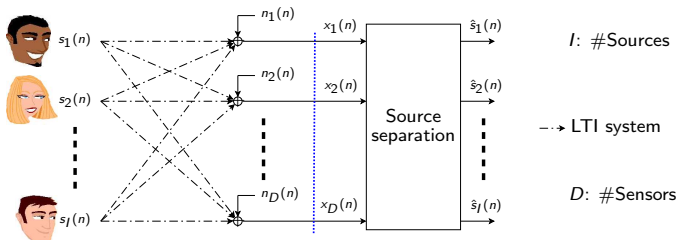


Blind Speech Separation Exploiting Temporal and Spectral Correlations Using 2D-HMMS

Dang Hai Tran Vu and Reinhold Haeb-Umbach

EUSIPCO 2013

Frequency domain blind source separation



Short-time Fourier-transform (STFT) domain mixing model

$$\mathbf{X}(m, k) \approx \sum_{i=1}^I \mathbf{H}_i(k) S_i(m, k) + \mathbf{N}(m, k)$$

where m : frame index, k : frequency bin

Sparseness-based BSS

Generative model in STFT domain

- At most one source dominant in each time-frequency slot (m, k)
- Generative model:

$$\mathbf{X}(m, k) = \begin{cases} \mathbf{N}(m, k) & \text{if } Z(m, k) = 0 \\ \mathbf{H}_i(k)S_i(m, k) + \mathbf{N}(m, k) & \text{if } Z(m, k) \in \{i; i = 1, \dots, I\} \end{cases}$$

- $Z(m, k)$ hidden random variable indicating which source is active

Goal: Extract $S_i(m, k)$ solely from $\mathbf{X}(m, k)$

- By computation of source activity probability

$$P(Z(m, k) | \mathbf{X}(m, k))$$

Sparseness-based BSS

Generative model in STFT domain

- At most one source dominant in each time-frequency slot (m, k)
- Generative model:

$$\mathbf{X}(m, k) = \begin{cases} \mathbf{N}(m, k) & \text{if } Z(m, k) = 0 \\ \mathbf{H}_i(k)S_i(m, k) + \mathbf{N}(m, k) & \text{if } Z(m, k) \in \{i; i = 1, \dots, I\} \end{cases}$$

- $Z(m, k)$ hidden random variable indicating which source is active
- Temporal and spectral correlations in $\mathbf{X}(1..M, 1..K)$ present

Goal: Extract $S_i(m)$ solely from $\mathbf{X}(1..M, 1..K)$

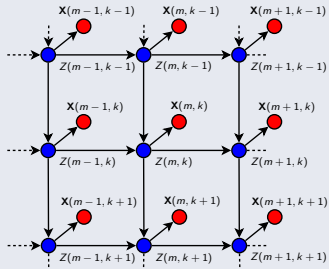
- By computation of source activity probability

$$P(Z(m, k) | \mathbf{X}(1..M, 1..K))$$

using temporal and spectral correlations!

Capturing temporal and spectral correlations

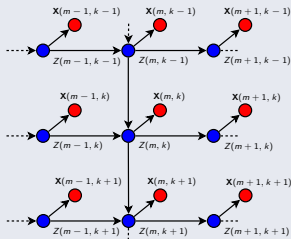
2D Hidden Markov Model



- Modeling correlations among adjacent time-frequency slots
- **Problem:** Exact inference is computational intractable!
 - ▶ **Solution:** Alternate inference in time and frequency direction

Decoding of a 2D-HMM

Modified forward-backward algorithm (FBA) in frequency direction



- Ignoring vertical dependencies in all other columns
- Modified forward-backward algorithm (FBA):

$$\gamma\alpha(m, k) = 1 \propto \gamma\mathbf{T}^T (\gamma\alpha(m, k-1) \circ \mathbf{o}(m, k-1) \circ \gamma\mathbf{u}(m, k-1))$$

$$\gamma\beta(m, k) = \pi \propto \gamma\mathbf{T} (\gamma\beta(m, k+1) \circ \mathbf{o}(m, k+1) \circ \gamma\mathbf{u}(m, k+1))$$

$$\gamma\gamma(m, k) = 1 \propto \mathbf{o}(m, k) \circ \gamma\mathbf{u}(m, k) \circ \gamma\alpha(m, k) \circ \gamma\beta(m, k)$$

vertical forward prediction variable: $\gamma\alpha(m, k)$,

prior probabilities: π ,

observation likelihood: $\mathbf{o}(m, k)$,

vertical backward variable: $\gamma\beta(m, k)$,

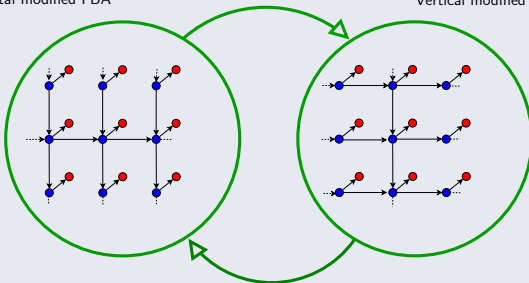
vertical transition matrix: $\gamma\mathbf{T}$,

vertical junction variable: $\gamma\mathbf{u}(m, k)$

Decoding of a 2D-HMM

Iterative Turbo Decoding Scheme

Horizontal modified FBA $\gamma \mathbf{u}(m, k) = (\gamma \alpha(m, k) \oslash \pi) \circ \gamma \beta(m, k)$ Vertical modified FBA



$$\gamma \mathbf{u}(m, k) = (\gamma \alpha(m, k) \oslash \pi) \circ \gamma \beta(m, k)$$

- Based on modified forward-backward algorithm (FBA) along time and frequency direction
- Extrinsic information exchange between both directions

Polar model of $p(\mathbf{X}(m, k)|Z(m, k))$

Average *a-posteriori* SNR

$$\varphi(m, k) := \frac{1}{D} \mathbf{X}^H(m, k) \Phi_{\mathbf{NN}}^{-1}(k) \mathbf{X}(m, k)$$

where $\Phi_{\mathbf{NN}} = \mathbb{E}[\mathbf{NN}^H]$

- Modeled by scaled chi-squared distribution

Frequency and unit-norm normalize observation vector

$$\tilde{Y}_j(m, k) := |X_j(m, k)| \exp \left\{ i \frac{\arg[X_j(m, k) X_1^*(m, k)]}{2(k-1)f_s d_{\max}(K_{C_V})^{-1}} \right\}$$

$$\mathbf{Y}(m, k) := \tilde{\mathbf{Y}}(m, k) / \|\tilde{\mathbf{Y}}(m, k)\|$$

- Modeled by complex Watson distribution

Assumption

$$p(\mathbf{X}(m, k)|Z(m, k)) = p(\mathbf{Y}(m, k)|Z(m, k)) \cdot p(\varphi(m, k)|Z(m, k))$$

Expectation Maximization algorithm for BSS

Parameters

- Set of unknown parameters $\Theta = \{\mathbf{W}_1, \dots, \mathbf{W}_I\}$
- Parameters of 2D-HMM are pretrained with speech mixtures

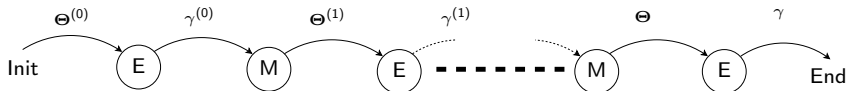
E-step: Estimate posterior speech activity probabilities

- Use proposed turbo decoding scheme to estimate posterior speech activity probabilities $\gamma_i^{(\nu)}(m, k) := P(Z(m, k) | \mathbf{X}(1..M, 1..K); \Theta^{(\nu)})$

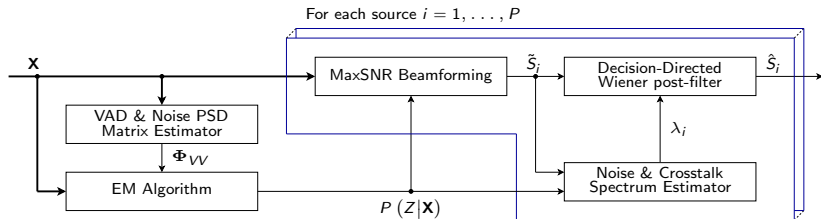
M-step: Update unknown parameter

- Principal eigenvector of source dependent matrix

$$\Phi_{\mathbf{Y}\mathbf{Y},i}^{(\nu)} := \frac{\sum_{m=1}^M \sum_{k=1}^K \gamma_i^{(\nu)}(m, k) \mathbf{Y}(m, k) \mathbf{Y}^H(m, k)}{\sum_{m=1}^M \sum_{k=1}^K \gamma_i^{(\nu)}(m, k)}$$



System overview



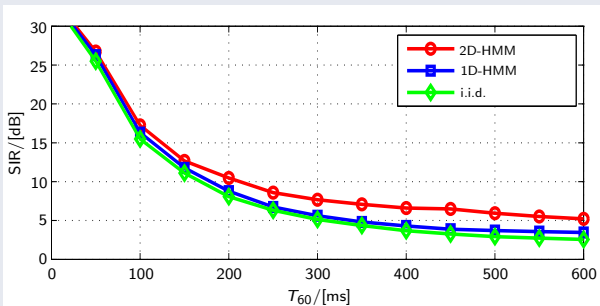
Properties

- Spatial filtering with generalized eigenvector (MaxSNR)-beamforming
- Spectral filtering with Wiener post-filter
- Adaptation control with source activity probability $P(Z|\mathbf{X})$

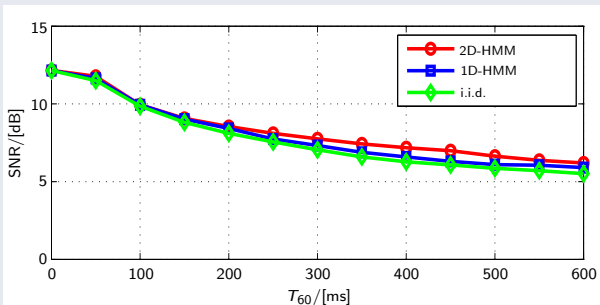
Setup

- Four sensor array arranged at the vertices of regular tetrahedron with lateral length of 2cm (No spatial aliasing!)
- 3 sources randomly positioned around the microphone
- Noise recordings of the fan noise of a video projector at -10 dB
- White noise at -20 dB
- Image method for reverberant room simulation

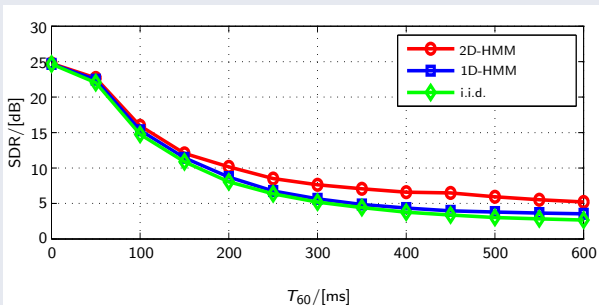
Gain in signal-to-interference-ratio (SIR)



Gain in signal-to-noise-ratio (SNR)



Gain in signal-to-distortion-ratio (SDR)



Summary and outlook

Summary

- Exploiting correlations of adjacent TF-slots for noisy BSS
 - ▶ 2D-HMM to capture temporal and spectral correlations
 - ▶ Iterative decoding scheme with modified FBA algorithm using extrinsic information exchange
- Improved performance in all cases and w.r.t. all measures
 - ▶ Advantages are evident especially in highly reverberant recording conditions

Outlook

- Exploitation of harmonic structures of speech
- Low latency block online implementation



Thank you for your attention!

Questions ?

Reinhold Haeb-Umbach

University of Paderborn
Department of Communications
Engineering

haeb@nt.uni-paderborn.de
nt.uni-paderborn.de