

多チャンネル階乗隠れマルコフモデルの スペクトル包絡事前学習によるセミブラインド音源分離*

◎樋口卓哉 (東大院情報理工), 亀岡弘和 (東大院情報理工, NTT)

1 はじめに

ブラインド音源分離 (Blind Source Separation; BSS) の問題とは、音源信号や音源からマイクまでの伝達特性が未知の場合に、複数の音源信号が混合された観測信号から音源信号を推定する問題である。一般的に BSS の問題を不良設定問題であり、解くためには、音源信号に対してなんらかの仮定を立て解を限定する必要がある。例えば、非負値行列因子分解 (Non-negative Matrix Factorization; NMF)[1, 2] では、観測信号のパワースペクトログラムが限られた数のパワースペクトルによって構成されていることを仮定することで、観測信号のパワースペクトログラムを 2 つの非負値行列の積に分解する。また筆者らは、音響イベント、残響、音源の到来方向などが BSS の問題を解くための手がかりと成りうることに着目し、多チャンネル階乗隠れマルコフモデルと呼ぶモデルを用いて、音響イベント、残響、音源の到来方向推定などを音源分離と統合的に行う手法を提案してきた [3, 4, 5]。この手法は、音源信号のパワースペクトルが限られた数のパワースペクトルの遷移によって表現できると仮定することで、BSS の問題において解を限定し、音源分離の手がかり (の一部) とする手法となっていた。しかし、NMF や多チャンネル階乗隠れマルコフモデルなどの、音源信号のスペクトルの数を制限するアプローチでは、音源信号をより詳細にモデル化しようとしてモデルにおけるパワースペクトルの数を増やせば増やすほど、BSS の問題において立てられていた仮定を弱め、解を限定することができなくなるというトレードオフが生じてしまう。そこで本稿では、[6] や [7] の手法のように、ソースフィルタモデルに基づいた新たな仮定を導入することを考える。ソースフィルタモデルとは、音声スペクトルを声帯駆動音源 (ソース) のスペクトルと声道 (フィルタ) スペクトルの積で表せるとしたモデルである。この仮定が成り立つならば、音声スペクトルは、複数のソーススペクトルと複数のフィルタスペクトルの組み合わせで表すことができる。これにより、 $Q^{(s)}$ 個のソーススペクトルと $Q^{(f)}$ 個のフィルタスペクトルを用いることで、 $Q^{(s)} + Q^{(f)}$ の自由度で $Q^{(s)} \times Q^{(f)}$ 個のスペクトルを表現することができる。本稿では、ソースフィルタモデルに基づいて音源信号をモデル化し、フィルタスペクトルを事前学習することで、音源信号をより詳細にモデル化しながらも、BSS の問題における解を十分に限定し、音源分離性能を向上させる手法を提案する。

2 多チャンネル階乗隠れマルコフモデル

2.1 畳み込み混合近似による観測信号の混合モデル

I 個の音源から到来する信号を M 個のマイクロフォンで観測する場合を考える。室内インパルス応答長が時間周波数展開の時間窓長に対して十分に短いと限らず、瞬時混合近似が成り立たない場合 (残響がある場合) を考え、時間周波数領域における畳み込み混合

近似によって時間周波数領域の観測信号を近似する。

$$\mathbf{y}(\omega_k, t_l) \approx \sum_{i=1}^I \sum_{\tau=0}^{\mathcal{T}} \mathbf{a}_i(\omega_k, t_\tau) s_i(\omega_k, t_l - t_\tau). \quad (1)$$

ただし、 $\mathbf{y}(\omega_k, t_l) = (y_1(\omega_k, t_l), \dots, y_M(\omega_k, t_l))^T \in \mathbb{C}^M$, $\mathbf{a}_i(\omega_k, t_\tau) = (a_{i,1}(\omega_k, t_\tau), \dots, a_{i,M}(\omega_k, t_\tau))^T \in \mathbb{C}^M$ である。ここで $y_m(\omega_k, t_l) \in \mathbb{C}$ は m 番目のマイクで観測された観測信号の周波数 ω_k , 時刻 t_l における時間周波数成分であり、 $s_i(\omega_k, t_l) \in \mathbb{C}$ は i 番目の音源信号の周波数 ω_k , 時刻 t_l における時間周波数成分である。また $\mathbf{a}_i(\omega_k, t_\tau)$ は i 番目の音源信号に対する周波数 ω_k における伝達周波数特性の時刻 t_τ の成分であり、 $0 \leq \tau \leq \mathcal{T}$ は伝達周波数特性の時間周波数領域における時間インデックスである。以下では ω_k, t_l をそれぞれ k, l の添え字で表す。

2.2 ソースフィルタモデルに基づく音源信号の生成モデル

では、音源信号の生成プロセスを確率的に記述する。まず音源信号が区分的に定常であることを仮定し、各時間周波数点で $s_{i,k,l}$ が平均 0, 分散 $\sigma_{i,k,l}^2$ の複素正規分布に従うとすると、音源信号の生成プロセスは、

$$s_{i,k,l} | \sigma_{i,k,l} \sim \mathcal{N}_{\mathbb{C}}(s_{i,k,l}; 0, \sigma_{i,k,l}^2), \quad (2)$$

と書き表せる。ここで $\sigma_{i,k,l}^2$ は周波数 k , 時刻 l における i 番目の音源のパワースペクトル密度を表す。上記のモデルに NMF の仮定を組み込むと、

$$\sigma_{i,k,l}^2 = w_{i,k} h_{i,l}, \quad (3)$$

となる。上記の式による多チャンネル観測信号の生成モデルは、多チャンネル NMF [8] と呼ばれている。ここで、ソースフィルタモデルに基づき、 i 番目の音源の周波数 k における基底パワースペクトルが、ソーススペクトル $w_{i,k}^{(s)}$ と、フィルタスペクトル $w_{i,k}^{(f)}$ の積によって構成されると仮定すると、

$$w_{i,k} = w_{i,k}^{(s)} w_{i,k}^{(f)}, \quad (4)$$

となる。しかしながら、ソーススペクトルやフィルタスペクトルは、発話内容によって時間変化すると考えられるので、ソーススペクトルの状態を表す隠れ変数 $z_{i,l}^{(s)} \in 1, \dots, Q^{(s)}$, フィルタスペクトルの状態を表す隠れ変数 $z_{i,l}^{(f)} \in 1, \dots, Q^{(f)}$ を導入し、それぞれの状態の時系列がマルコフ連鎖に従うと仮定する。

$$z_{i,l}^{(s)} | z_{i,l-1}^{(s)} \sim \text{Categorical}(z_{i,l}^{(s)}; \boldsymbol{\rho}_{i,z_{i,l-1}^{(s)}}^{(s)}), \quad (5)$$

$$z_{i,l}^{(f)} | z_{i,l-1}^{(f)} \sim \text{Categorical}(z_{i,l}^{(f)}; \boldsymbol{\rho}_{i,z_{i,l-1}^{(f)}}^{(f)}). \quad (6)$$

ここで $\text{Categorical}(x; \mathbf{y}) = y_x$ であり、 $\boldsymbol{\rho}_{i,q^{(s)}}^{(s)} = (\rho_{i,q^{(s)},1}^{(s)}, \dots, \rho_{i,q^{(s)},Q^{(s)}}^{(s)})$ は i 番目の音源におけるソ-

*Semi-blind source separation with multichannel factorial hidden Markov model using pretrained spectral envelope templates. by HIGUCHI, Takuya (The University of Tokyo), KAMEOKA Hirokazu (The University of Tokyo, NTT)

スの状態 $q^{(s)}$ から各状態 $1, \dots, Q^{(s)}$ への遷移確率を表し、 $\rho_{i,q^{(f)}}^{(f)} = (\rho_{i,q^{(f)},1}^{(f)}, \dots, \rho_{i,q^{(f)},Q^{(f)}}^{(f)})$ は i 番目の音源におけるフィルタの状態 $q^{(f)}$ から各状態 $1, \dots, Q^{(f)}$ への遷移確率を表す。状態 $q^{(s)}$ である i 番目の音源のソーススペクトルを $w_{i,k,q^{(s)}}^{(s)}$ 、状態 $q^{(f)}$ である i 番目の音源のフィルタスペクトルを $w_{i,k,q^{(f)}}^{(f)}$ と表すとすると、時刻 l における i 番目の音源信号のソーススペクトルとフィルタスペクトルは $z_{i,l}^{(s)}$ と $z_{i,l}^{(f)}$ にそれぞれ依存し、 $s_{i,k,l}$ の生成モデルは以下のように書き直せる。

$$\mathbf{s}_{i,k,l} | w_{i,k,1:Q^{(s)}}^{(s)}, w_{i,k,1:Q^{(f)}}^{(f)}, h_{i,l}, z_{i,l}^{(s)}, z_{i,l}^{(f)} \sim \mathcal{N}_{\mathbb{C}}(\mathbf{s}_{i,k,l}; 0, w_{i,k,z_{i,l}^{(s)}}^{(s)} w_{i,k,z_{i,l}^{(f)}}^{(f)} h_{i,l}). \quad (7)$$

次に、音源の音量に着目すると、無音状態と有音状態では当然音量の大きな値の取りやすさが異なると考えられるので、音量もまた音源の状態（音響イベント）に依存して異なる振る舞いをするといえる。そこで、まず音量の状態を表す隠れ変数 $z_{i,l}^{(j)} \in 1, \dots, J$ を導入し、状態の時系列 $z_{i,1}^{(j)}, \dots, z_{i,L}^{(j)}$ がマルコフ連鎖に従うと仮定する。

$$z_{i,l}^{(j)} | z_{i,l-1}^{(j)} \sim \text{Categorical}(z_{i,l}^{(j)}; \boldsymbol{\rho}_{i,z_{i,l-1}^{(j)}}^{(j)}). \quad (8)$$

このとき、 $h_{i,l}$ が音量の状態 $z_{i,l}^{(j)} \in 1, \dots, J$ によって異なるハイパーパラメータを持つガンマ分布に従うと仮定すると、

$$h_{i,l} | z_{i,l}^{(j)} \sim \text{Gamma}(h_{i,l}; \alpha_{z_{i,l}^{(j)}} \beta_{z_{i,l}^{(j)}}), \quad (9)$$

となる。ここで $\alpha_{1:J}$ と $\beta_{1:J}$ はそれぞれガンマ分布の形状パラメータとスケールパラメータであり、 $\text{Gamma}(x; \alpha, \beta) = \frac{x^{\alpha-1} e^{-x/\beta}}{\Gamma(\alpha) \beta^{\alpha}}$ 、ただし $\Gamma(\cdot)$ はガンマ関数である。 $z_{i,l}^{(j)}$ が無音状態に対応するときは $h_{i,l}$ は小さな値をとってほしいので、小さな値をとる確率が高くなるようにガンマ分布のハイパーパラメータを設定し、 $z_{i,l}^{(j)}$ が有音状態に対応するときは一様分布に近くなるようにガンマ分布のハイパーパラメータを設定すればよい。

2.3 混合 DOA モデルによる空間相関行列の生成モデル

次に空間相関行列の生成プロセスを確率的にモデル化する。点音源と平面波到来を仮定すると、空間相関行列は音源の到来方向に応じてある特定の構造を持つ。例えば、マイクロフォンの数 $M = 2$ の場合では、方向 $\theta (0 \leq \theta \leq \pi)$ にある音源の空間相関行列は、以下のように陽に記述できる。

$$\mathbf{J}(\theta, \omega) = \begin{bmatrix} 1 \\ e^{j\omega B \cos \theta / C} \end{bmatrix} \begin{bmatrix} 1 & e^{j\omega B \cos \theta / C} \end{bmatrix}^* \quad (10)$$

ここで j は虚数単位、 B [m] はマイクロフォン間の距離、 C [m/s] は音速である。 i 番目の音源の到来方向 θ_i が既知の場合では、直接波に対応する空間相関行列は $\mathbf{J}(\theta_i, \omega_k)$ と等しくなることが期待される。しかしながら実際には、音源の到来方向は音響信号から直接観測できないばかりでなく、フレーム内の残響やノイズなどによって、空間相関行列は理想的な構造から逸脱することがある。

そこでまず、[9] の手法と同様に、 $\theta_i \in \{\vartheta_1, \dots, \vartheta_O\}$ を i 番目の音源の到来方向とし、各音源の DOA がこの DOA 候補値の中から決定されると仮定することで、音源 i の到来方向 θ_i が生成されるプロセスを以下のように記述する。

$$z_i^{(o)} | \boldsymbol{\rho}_i^{(o)} \sim \text{Categorical}(z_i^{(o)}; \boldsymbol{\rho}_i^{(o)}), \quad (11)$$

$$\theta_i = \vartheta_{z_i^{(o)}}. \quad (12)$$

$z_i^{(o)} \in \{1, \dots, O\}$ は i 番目の音源にどの DOA 候補値が割り当てられるかを表すインジケータ変数であり、上式はこれが離散分布（各確率値が $\rho_1^{(o)}, \dots, \rho_O^{(o)}$ から生成されることを意味している。そして、直接波に対応する空間相関行列 $\mathbf{C}_{i,k,0}$ が、 $z_i^{(o)}$ が既知の条件下で、以下のようなウィッシュヤート分布に従うと仮定する。

$$\mathbf{C}_{i,k,0} | z_i^{(o)} \sim \mathcal{W}_{\mathbb{C}}(\mathbf{C}_{i,k,0}; \gamma, \mathbf{J}_{\vartheta_{z_i^{(o)}}} + \epsilon \mathbf{I}). \quad (13)$$

ただし $\mathcal{W}_{\mathbb{C}}(\mathbf{X}; \gamma, \mathbf{Y}) \propto |\mathbf{X}|^{(\gamma-M)/2} \exp(-\frac{1}{2} \text{tr}(\mathbf{X}\mathbf{Y}^{-1}))$ であり、 γ は $\gamma \geq M + 1$ を満たすハイパーパラメータである。ここでは、逆行列演算を可能にするために、近似的に $\mathbf{J}_{\vartheta_{z_i^{(o)}}}$ に微小値 ϵ を用いて単位行列 \mathbf{I} を足

してある。ここで隠れ変数 $z_i^{(o)}$ を推定することが、音源の到来方向推定を行うことに相当している。

観測信号の最終的な生成モデルは $\mathbf{a}_{1:I,k,0:T}$ 、 $w_{1:I,k,1:Q^{(s)}}^{(s)}$ 、 $w_{1:I,k,1:Q^{(f)}}^{(f)}$ 、 $h_{1:I,l-T:l}$ 、 $z_{1:I,l-T:l}^{(s)}$ 、 $z_{1:I,l-T:l}^{(f)}$ が既知の条件下で、式 (5)、式 (6)、式 (8)、式 (9)、式 (13) と合わせて以下のように書き直せる。

$$\mathbf{y}_{k,l} | \mathbf{a}_{1:I,k,0:T}, w_{1:I,k,1:Q^{(s)}}^{(s)}, w_{1:I,k,1:Q^{(f)}}^{(f)}, h_{1:I,l-T:l}, z_{1:I,l-T:l}^{(s)}, z_{1:I,l-T:l}^{(f)} \sim \mathcal{N}_{\mathbb{C}}(\mathbf{y}_{k,l}; 0, \sum_{i,\tau} \mathbf{C}_{i,k,\tau} w_{i,k,z_{i,l}^{(s)}}^{(s)} w_{i,k,z_{i,l}^{(f)}}^{(f)} h_{i,l-\tau}). \quad (14)$$

この生成モデルに基づいて最適なパラメータを求めることは、ソースフィルタモデルによって音源信号の自由度を限定しながら、音源の到来方向推定・残響除去・音響イベント検出・音源分離の問題を統合的に解くことに相当している。

3 補助関数法に基づくパラメータ推論アルゴリズム

目的関数である対数事後確率では各パラメータが複雑に関係し合っており、一般的に最適化が困難である。そこで効率的に最適なパラメータを求める反復アルゴリズムを、補助関数法の原理に基づき導出する。モデルにおける推定したい変数は $\mathbf{W}^{(s)} = w_{1:I,1:K,1:Q^{(s)}}^{(s)}$ 、 $\mathbf{W}^{(f)} = w_{1:I,1:K,1:Q^{(f)}}^{(f)}$ 、 $\mathbf{H} = h_{1:I,1:L}$ 、 $\mathbf{C} = \mathbf{C}_{1:I,1:K,0:T}$ である。上記の変数の集合を Θ で表す。 $\mathbf{Z}^{(j)} = z_{1:I,1:L}^{(j)}$ 、 $\mathbf{Z}^{(o)} = z_{1:I}^{(o)}$ 、 $\mathbf{Z}^{(s)} = z_{1:I,1:L}^{(s)}$ 、 $\mathbf{Z}^{(f)} = z_{1:I,1:L}^{(f)}$ は隠れ変数とする。以下では $\boldsymbol{\rho}^{(s)}$ 、 $\boldsymbol{\rho}^{(f)}$ 、 $\boldsymbol{\rho}^{(j)}$ 、 $\boldsymbol{\rho}^{(o)}$ は実験的に定められた定数とする。補助変数の集合を Λ とすると、目的関数 $L(\Theta) = \log p(\Theta | \mathbf{Y})$ に対する補助関数 $L^+(\Theta, \Lambda)$ は以下のように設計できる。

$$L^+(\Theta, \Lambda)$$

$$\begin{aligned}
&= -\frac{1}{2} \sum_{k,l} \left[\sum_{i,q^{(s)},q^{(f)},\tau} \lambda_{q^{(s)},i,l-\tau}^{(s)} \lambda_{q^{(f)},i,l-\tau}^{(f)} \left(\text{tr} \left(\mathbf{y}_{k,l} \mathbf{y}_{k,l}^H \right. \right. \right. \\
&\quad \left. \left. \mathbf{R}_{i,k,l,\tau,q^{(s)},q^{(f)}} \mathbf{C}_{i,k,\tau}^{-1} \right. \right. \\
&\quad \left. \left. \mathbf{R}_{i,k,l,\tau,q^{(s)},q^{(f)}} \right) / w_{i,k,q^{(s)}}^{(s)} w_{i,k,q^{(f)}}^{(f)} h_{i,l-\tau} \right. \\
&\quad \left. + \text{tr} \left(\mathbf{U}_{k,l}^{-1} \mathbf{C}_{i,k,\tau} \right) w_{i,k,q^{(s)}}^{(s)} w_{i,k,q^{(f)}}^{(f)} h_{i,l-\tau} \right) + \log |\mathbf{U}_{k,l}| \Big] \\
&\quad + \sum_{j,i,l} \lambda_{j,i,l}^{(j)} \left((\alpha_j - 1) \log h_{i,l} - h_{i,l} / \beta_j \right. \\
&\quad \left. - \alpha_j \log \beta_j - \log \Gamma(\alpha_j) \right) \\
&\quad + \sum_{i,k,o} d_{i,o} \left[-\frac{1}{2} \text{tr} \left(\mathbf{C}_{i,k,0} (\mathbf{J}_{k,o} + \epsilon \mathbf{I})^{-1} \right) \right. \\
&\quad \left. + \frac{\gamma - M}{2} \left(\text{tr} \left(\tilde{\mathbf{U}}_{i,k}^{-1} \mathbf{C}_{i,k,0} \right) + \log |\tilde{\mathbf{U}}_{i,k}| \right) \right] \\
&\quad + \sum_{q^{(s)},i,l} \lambda_{q^{(s)},i,l}^{(s)} \log p(\mathbf{Z}^{(s)}) + \sum_{q^{(f)},i,l} \lambda_{q^{(f)},i,l}^{(f)} \log p(\mathbf{Z}^{(f)}) \\
&\quad + \sum_{j,i,l} \lambda_{j,i,l}^{(j)} \log p(\mathbf{Z}^{(j)}) + C_{L+}, \tag{15}
\end{aligned}$$

ここで $\mathbf{R}_{i,k,l,\tau,q^{(s)},q^{(f)}}$, $\mathbf{U}_{k,l}$, $\tilde{\mathbf{U}}_{i,k}$ は $\sum_{i,\tau} \mathbf{R}_{i,k,l,\tau,q^{(s)},q^{(f)}} = \mathbf{I}$ を満たすエルミート正定値行列であり, $\lambda_{q^{(s)},i,l}^{(s)}$, $\lambda_{q^{(f)},i,l}^{(f)}$, $\lambda_{j,i,l}^{(j)}$, $d_{i,o}$ は $\sum_{q^{(s)}} \lambda_{q^{(s)},i,l}^{(s)} = 1$, $\sum_{q^{(f)}} \lambda_{q^{(f)},i,l}^{(f)} = 1$, $\sum_j \lambda_{j,i,l}^{(j)} = 1$, $\sum_o d_{i,o} = 1$ をそれぞれ満たす非負値のスカラー値である. C_{L+} は定数項をまとめたものである. 式 (15) の等号成立条件は, $\mathbf{R}_{i,k,l,\tau,q^{(s)},q^{(f)}}$ と $\mathbf{U}_{k,l}$ に関しては

$$\mathbf{R}_{i,k,l,\tau,q^{(s)},q^{(f)}} = \mathbf{C}_{i,k,\tau} w_{i,k,q^{(s)}}^{(s)} w_{i,k,q^{(f)}}^{(f)} h_{i,l-\tau} \hat{\mathbf{X}}_{k,l}^{-1}, \tag{16}$$

$$\mathbf{U}_{k,l} = \hat{\mathbf{X}}_{k,l}, \tag{17}$$

である. ただし $\hat{\mathbf{X}}_{k,l} = \sum_{i,\tau} \mathbf{C}_{i,k,\tau} w_{i,k,z_{i,l-\tau}^{(s)}}^{(s)} w_{i,k,z_{i,l-\tau}^{(f)}}^{(f)} h_{i,l-\tau}$ である. $\lambda_{q^{(s)},i,l}^{(s)}$ の等号成立条件に関しては Forward-Backward アルゴリズムを用いて,

$$\lambda_{q^{(s)},i,l}^{(s)} = F_{q^{(s)},i,l}^{(s)} B_{q^{(s)},i,l}^{(s)} / \sum_{q^{(s)}} F_{q^{(s)},i,l}^{(s)} B_{q^{(s)},i,l}^{(s)}, \tag{18}$$

$$F_{q^{(s)},i,l}^{(s)} = p(\Theta | z_{i,l}^{(s)} = q^{(s)}) \sum_{q^{(s)'}} F_{q^{(s)',i,l-1}^{(s)}} \rho_{i,q^{(s)',q^{(s)}}}^{(s)}, \tag{19}$$

$$B_{q^{(s)},i,l}^{(s)} = \sum_{q^{(s)'}} B_{q^{(s)',i,l+1}^{(s)}} p(\Theta | z_{i,l+1}^{(s)} = q^{(s)'}) \rho_{i,q^{(s)',q^{(s)}}}^{(s)}, \tag{20}$$

ただし

$$\begin{aligned}
&p(\Theta | z_{i,l}^{(s)} = q^{(s)}) \\
&\propto \exp \left[-\frac{1}{2} \sum_{k,\tau,q^{(f)}} \lambda_{q^{(f)},i,l}^{(f)} \left(\text{tr} \left(\mathbf{y}_{k,l+\tau} \mathbf{y}_{k,l+\tau}^H \right. \right. \right. \\
&\quad \left. \left. \mathbf{R}_{i,k,l+\tau,\tau,q^{(s)},q^{(f)}} \mathbf{C}_{i,k,\tau}^{-1} \right. \right. \\
&\quad \left. \left. \mathbf{R}_{i,k,l+\tau,\tau,q^{(s)},q^{(f)}} \right) / w_{i,k,q^{(s)}}^{(s)} w_{i,k,q^{(f)}}^{(f)} h_{i,l} \right. \\
&\quad \left. + \text{tr} \left(\mathbf{U}_{k,l+\tau}^{-1} \mathbf{C}_{i,k,\tau} \right) w_{i,k,q^{(s)}}^{(s)} w_{i,k,q^{(f)}}^{(f)} h_{i,l} \right) \Big], \tag{25}
\end{aligned}$$

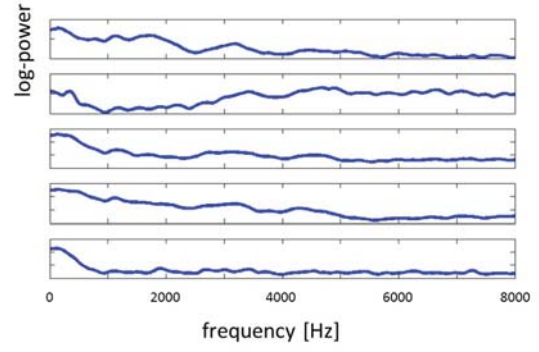


Fig. 1 クリーンな音声から STRAIGHT[12] によって抽出し実験に用いたフィルタスペクトルの例.

$$\begin{aligned}
&\mathbf{R}_{i,k,l+\tau,\tau,q^{(s)},q^{(f)}} / w_{i,k,q^{(s)}}^{(s)} w_{i,k,q^{(f)}}^{(f)} h_{i,l} \\
&+ \text{tr} \left(\mathbf{U}_{k,l+\tau}^{-1} \mathbf{C}_{i,k,\tau} \right) w_{i,k,q^{(s)}}^{(s)} w_{i,k,q^{(f)}}^{(f)} h_{i,l} \Big], \tag{21}
\end{aligned}$$

である. さらに $\lambda_{q^{(f)},i,l}^{(f)}$ の等号成立条件についても Forward-Backward アルゴリズムを用いて,

$$\lambda_{q^{(f)},i,l}^{(f)} = F_{q^{(f)},i,l}^{(f)} B_{q^{(f)},i,l}^{(f)} / \sum_{q^{(f)}} F_{q^{(f)},i,l}^{(f)} B_{q^{(f)},i,l}^{(f)}, \tag{22}$$

$$F_{q^{(f)},i,l}^{(f)} = p(\Theta | z_{i,l}^{(f)} = q^{(f)}) \sum_{q^{(f)'}} F_{q^{(f)',i,l-1}^{(f)}} \rho_{i,q^{(f)',q^{(f)}}}^{(f)}, \tag{23}$$

$$B_{q^{(f)},i,l}^{(f)} = \sum_{q^{(f)'}} B_{q^{(f)',i,l+1}^{(f)}} p(\Theta | z_{i,l+1}^{(f)} = q^{(f)'}) \rho_{i,q^{(f)',q^{(f)}}}^{(f)}, \tag{24}$$

ただし

$$\begin{aligned}
&p(\Theta | z_{i,l}^{(f)} = q^{(f)}) \\
&\propto \exp \left[-\frac{1}{2} \sum_{k,\tau,q^{(s)}} \lambda_{q^{(s)},i,l}^{(s)} \left(\text{tr} \left(\mathbf{y}_{k,l+\tau} \mathbf{y}_{k,l+\tau}^H \right. \right. \right. \\
&\quad \left. \left. \mathbf{R}_{i,k,l+\tau,\tau,q^{(s)},q^{(f)}} \mathbf{C}_{i,k,\tau}^{-1} \right. \right. \\
&\quad \left. \left. \mathbf{R}_{i,k,l+\tau,\tau,q^{(s)},q^{(f)}} \right) / w_{i,k,q^{(s)}}^{(s)} w_{i,k,q^{(f)}}^{(f)} h_{i,l} \right. \\
&\quad \left. + \text{tr} \left(\mathbf{U}_{k,l+\tau}^{-1} \mathbf{C}_{i,k,\tau} \right) w_{i,k,q^{(s)}}^{(s)} w_{i,k,q^{(f)}}^{(f)} h_{i,l} \right) \Big], \tag{25}
\end{aligned}$$

である. $\tilde{\mathbf{U}}_{i,k}$, $d_{i,o}$, $\lambda_{j,i,l}^{(j)}$ の等号成立条件は, [5] と同様である. 紙面の都合上詳細は省略するが, その他のパラメータの更新則は $L+$ の偏微分を 0 と置くことで導ける.

4 評価実験

4.1 音源分離性能の評価実験

提案法の音源分離性能の評価のために, 残響下での半教師あり音源分離実験を行った. ATR 音声データベース [10] の中の 3 人の発話者 (女性 2 人, 男性 1 人) による 15 種類の発話に対して, RWCP データベース [11] のインパルス応答 (残響時間 380 ms, マイク間距離 11.48 cm, マイクの数 $M = 2$) を畳み込み, 人工

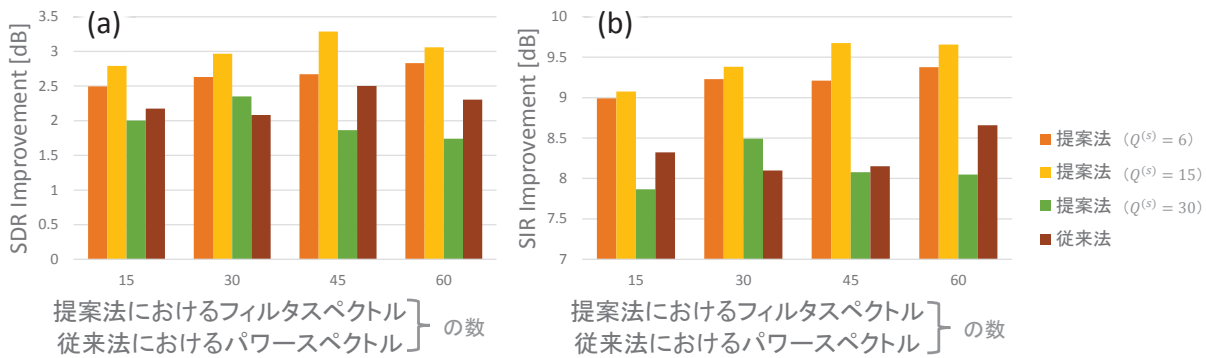


Fig. 2 提案法におけるフィルタスペクトルの数や従来法 [5] における基底パワースペクトルの数を変えて実験を行ったときに得られた分離音の (a)SDR の改善量と (b)SIR の改善量。

的に多チャンネルの混合信号を 5 つ作成した。3 つの音源はインパルス応答を畳み込まれることでマイクから 30° , 90° , 130° の位置に人工的にそれぞれ配置された。サンプリング周波数は 16 kHz とした。フレーム長 64 ms, フレームシフト長 16 ms で STFT を行い, 時間周波数展開を行った。ウィッシュャート分布のハイパーパラメータ γ は 3 とした。ソーススペクトルの状態数 $Q^{(s)}$ を 6, 15, 30, フィルタスペクトルの状態数 $Q^{(f)}$ を 15, 30, 45, 60 とそれぞれ変えて実験した。 $\epsilon = 10^{-3}$, 到来方向の分割数 $O = 30$ とした。音量の状態数 $J = 2$ とし, α_1 と β_1 を 1, 10^{-3} とそれぞれ設定し, α_2 と β_2 を 1 と 10^{10} と設定することで, $j = 1$ を無音状態とみなした。 $\rho^{(j)}$ は $i = 1, \dots, 3$ において $\rho_{i,1}^{(j)} = (0.9, 0.1)$, $\rho_{i,2}^{(j)} = (0.1, 0.9)$ とし, 自己遷移の確率を高め設定した。 $\mathcal{T} = 3$ とした。 \mathbf{C} の初期値は, $\tau = 0$ の成分に関しては $1/\sqrt{M} \times \mathbf{I}$ とし, $\tau = 1 : \mathcal{T}$ では $10^{-2} \times \mathbf{I}$ とした。 \mathbf{H} の初期値はすべての要素を 1 とし与えた。 $\lambda^{(q)}$ の初期値は $1/Q$, $\lambda^{(j)}$ の初期値は $1/J$ とした。 $\mathbf{W}^{(f)}$ は, STRAIGHT[12] を用いて混合音声とは別の発話から抽出したスペクトル包絡に対して HMM 学習することで 3 人の話者の発話から $Q^{(f)}/3$ つずつ学習し, 固定した。 Fig. 1 に学習したフィルタスペクトルの一例を示す。 $\mathbf{W}^{(s)}$ は, STRAIGHT を用いて抽出したスペクトル包絡系列でスペクトログラムを割ったものをソース成分と見なし HMM 学習をしたものを初期値とし, 反復アルゴリズムにより更新した。 $\rho^{(s)}$, $\rho^{(f)}$ は一様とした。 d の初期値は到来方向を三等分するように設定した。また, θ_i の取りうる範囲を $\theta_{1+10(i-1)}$ から θ_{10i} に限定し, 各音源ごとに推定する到来方向が重ならないように $\rho^{(o)}$ を設定した。パラメータ推論アルゴリズムは 50 回反復した。前章と同様に望ましくない局所解を避けるため, 最初の 25 回の反復は $\mathcal{T} = 0$ とし, その後徐々に \mathcal{T} を 3 となるまで増やしながら反復した。また, 最初の 25 回の反復は \mathbf{W} を更新し, その後 \mathcal{T} を増やし $\mathcal{T} > 0$ となつてからは \mathbf{W} を更新せずに固定した。比較対象には [5] の手法を用いた。比較手法では, 分離に用いる混合音とは別の発話から学習した基底パワースペクトルを \mathbf{W} の初期値として用いた。また基底パワースペクトルの状態数 Q を 15, 30, 45, 60 とそれぞれ変えて実験した。分離音は多チャンネルウィナーフィルタによって得た。音源分離, 残響除去の客観評価基準として, 残響除去済み分離音の振幅スペクトログラムに対して, 室内インパルス応答を畳み込む前の音源信号の振幅スペクトログラムを参照することで得られた SDR, SIR[13] を用いた。分離処理前の SDR, SIR はそれぞれ -4.77, -0.96 [dB] であった。

Fig. 2 に, 本実験によって得られた分離音の SDR, SIR の改善量をそれぞれ示す。提案法におけるソーススペクトルやフィルタスペクトルの数, 従来法にお

ける基底パワースペクトルの数によって変動してはいるが, 提案法によって得られた分離音の客観評価値は, 従来法によって得られた分離音の客観評価値を上回る傾向が見られた。特に $Q^{(s)} = 6, 15$ のときの提案法を用いて得られた分離音は, 従来法によって得られた分離音より総じて高い SDR と SIR を示した。SDR は $Q^{(s)} = 15, Q^{(f)} = 45$ のとき最も高い値を示した。これにより, 音源分離における提案法の有効性が示された。

5 おわりに

本稿では, ソースフィルタモデルに基づいて音源信号をモデル化し, 多チャンネル階乗隠れマルコフモデルに組み込むことによって, 従来法を上回る性能を実現する音源分離手法を提案した。

謝辞 本研究は JSPS 科研費 26730100 の助成を受けたものです。

参考文献

- [1] D. D. Lee, and H. S. Seung, *Nature*, vol. 401, pp.788–791, 1999.
- [2] P. Smaragdis, and J. C. Brown, *WASPAA 2003*, pp. 177–180, Oct. 2003.
- [3] T. Higuchi, *et al.*, *Interspeech 2014*, pp. 850–854, 2014.
- [4] T. Higuchi and H. Kameoka, “Joint audio source separation and dereverberation based on multi-channel factorial hidden Markov model,” *MLSP 2014*.
- [5] 樋口 他, “多チャンネル階乗隠れマルコフモデルと混合 DOA モデルによる音源分離・到来方向推定・音響イベント検出・残響除去の統合的アプローチ,” 電子情報通信学会技術研究報告, 電気音響, 2015 [to appear].
- [6] H. Kameoka and K. Kashino, *ISCAS2009*, pp. 2477–2480, May 2009.
- [7] 吉井和佳, 後藤真孝, 情報処理学会研究報告, 2012–MUS–96–8, Aug. 2012.
- [8] H. Sawada *et al.*, *ICASSP 2012*, pp. 261–264, 2012.
- [9] 亀岡 他, 音講論 (春), 1–1–19, pp.713–716, Mar. 2012.
- [10] A. Kurematsu *et al.*, *Speech Communication*, pp. 357–363, 1990.
- [11] S. Nakamura *et al.*, *LREC 2000*, pp. 965–968, 2000.
- [12] H. Kawahara *et al.*, *Speech Communication*, vol. 27, no. 3–4, pp. 187–207, 1999.
- [13] E. Vincent *et al.*, *IEEE Transactions on Audio, Speech, and Language Processing*, pp. 1462–1469, 2006.