

[ポスター講演] マルチチャンネル音源分離のための 低ランク音源モデルとスパース重畳過程に基づく ネスト型ベイズ混合・因子モデル

板倉 光佑[†] 坂東 宜昭[†] 中村 栄太[†] 糸山 克寿[†] 吉井 和佳[†]
河原 達也[†]

[†] 京都大学 大学院情報学研究科

E-mail: †{itakura,bando,enakamura,itoyama,yoshii,kawahara}@sap.ist.i.kyoto-u.ac.jp

あらまし 本稿では、低ランク音源モデルとスパース重畳過程に基づくネスト型ベイズ混合・因子モデルを用いたマルチチャンネル音源分離法について述べる。従来の音源分離では、音源モデルに対して低ランク性を仮定し因子モデルを用いたモデル化を行う非負値行列因子分解や重畳過程において音源のスパース性を仮定し混合モデルに基づいた時間周波数クラスタリングによる音源分離法などがあった。提案法ではこの音源モデルと重畳過程を統合したモデル化により音源分離を行う。また、因子モデルと混合モデルの関係性に着目し、音源モデルと重畳過程のそれぞれに対し因子モデルと混合モデルによるモデル化を行うことで複数の分離法を提案する。

キーワード マルチチャンネル音源分離, 時間周波数クラスタリング, 潜在的ディリクレ配分法, 非負値行列因子分解

1. はじめに

マイクロホンアレイを用いたマルチチャンネル音源分離において、これまでに多くの手法が提案されてきた。広く用いられている手法のうちの一つに、独立成分分析 (ICA) [1] がある。ICA は各音源の統計的な独立性を仮定することにより混合行列の逆行列である分離行列を推定し、分離を行う。この ICA をもとに独立ベクトル分析 (IVA) [2] や FastICA [3] などのさまざまな手法が提案されているが、これらの手法は共通してマイク数が音源数より少ない劣決定条件では分離できないという問題点がある。

これに対し、劣決定条件でも分離が可能な手法として時間・周波数クラスタリングに基づく音源分離法が着目されている [4-9]。このアプローチでは、各音源スペクトログラムが時間・周波数領域でスパースであると仮定することで、混合音スペクトログラムの各時間・周波数ビンにおける観測はそれぞれいずれか一つの音源成分が直接観測されたものであるとみなす。これにより、音の重畳過程が混合モデルにより表される。この混合モデルを推定するため、マイク間の音の位相差とパワー差を特徴量として用いた混合ワトソン分布のクラスタリング [5-7] や各マイクでの音の位相とパワーを特徴量として用いた混合ガウス分布のクラスタリング [9,10] による分離法が提案されている。これらのクラスタリングはそれぞれの時間周波数ビンでそれぞれ独立に行われるためパーミュテーション問題が発生するが、大塚ら [9] は潜在的ディリクレ配分法 (LDA) を用いて各時間周波数ビンを各音源に割り当てたあと、それらの音源をさらにいずれかの方向に割り当てることでこのパーミュテーション問題を解決した。

また、単チャンネル音源分離のための手法として非負値行列因子分解 (NMF) [11] が広く使われている。単チャンネル音源分離ではマイク間の位相差などの空間的な情報を用いることができないため、NMF では音源の低ランク性を仮定することで音源を構成する因子を推定し、音源分離を行う。具体的には、音源が基底とアクティベーションの二つの因子から成るとし、観測のパワースペクトログラムを基底行列とアクティベーション行列の積により近似する。このような因子分解によるモデルを因子モデルと呼ぶ。また、この NMF をマルチチャンネル音源分離のために拡張したマルチチャンネル NMF (MNMF) [12] も提案されている。MNMF では観測スペクトログラムを基底行列とアクティベーション行列と空間特徴量に分解する。MNMF は空間特徴量も因子分解により推定するため、音源モデル・重畳過程がともに因子モデルによりモデル化された手法である。

本稿では、音の重畳過程と音源モデルを統合したモデルにより音源分離を行う。また、従来法では重畳過程が混合モデルでも因子モデルでもモデル化されていたことから、その二つのモデル化の関係性に着目し、音源モデルにおいても混合モデルと因子モデルの二つのモデル化を行った。つまり、因子モデルでは音源パワースペクトログラムの各時間周波数ビンが複数の基底とアクティベーションの積の和で表されたが、混合モデルでは各時間周波数ビンをいずれか一つの基底とアクティベーションの積により表す。このように、本稿では重畳過程・音源モデルに対して、混合モデル・因子モデルのそれぞれのモデル化を行うことで、それらの組み合わせにより複数の音源分離法を実現した。具体的には、図 1 に示すように、因子モデルでは音源のパワースペクトログラムが基底とアクティベーションの積の和から成り、観測スペクトログラムは音源スペクトログラムと

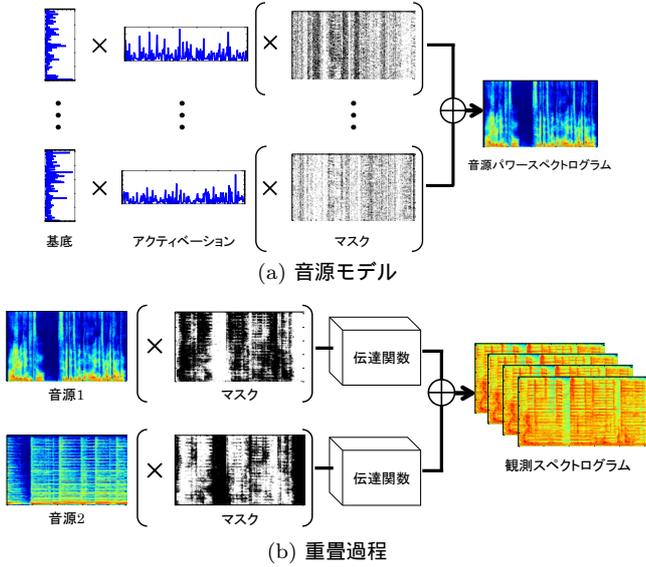


図 1 提案法の生成モデル．音源モデルでは各音源のパワースペクトログラムが基底・アクティベーション（・マスク）により構成される．重畳過程では混合音スペクトログラムが音源スペクトログラム・伝達関数（・マスク）により構成される．音源モデル・重畳過程ともに混合モデルではマスクによりどれか一つの基底，もしくは音源が選択されるが，因子モデルではマスクはなくすべての基底，もしくは音源が使用される．

伝達関数から生成されるとする．混合モデルでは音源のパワースペクトログラムがマスクによって選択された基底とアクティベーションから成り，観測スペクトログラムはマスクによって選択された音源スペクトログラムと伝達関数から生成されるとする．提案法では混合モデルには LDA，因子モデルには NMF の枠組みを用いてギブスサンプリングにより音源分離を行う．

2. モデル化

提案法では，音の重畳過程と各音源の音源モデルを考慮することで音源分離を行う．モデル化には NMF のような因子モデルと LDA のような混合モデルの二通りの方法がある．本章では重畳過程と音源モデルのそれぞれに対してこれらの二通りのモデル化を行う方法について述べる．それらの組み合わせにより，表 1 に示すように音源モデル・重畳過程を因子・混合モデル，混合・混合モデル，混合・因子モデル，因子・因子モデルでそれぞれモデル化した NMF-LDA，LDA-LDA，LDA-NMF，NMF-NMF について述べる．

2.1 問題設定

まず，単一の音源を観測した時の観測モデルについて述べる．提案法では時間領域の信号に対して短時間フーリエ変換 (STFT) を行うことにより得られる時間周波数領域の信号に対してモデル化を行う．まず K 個の音源を M 個のマイクを用いて録音するとし，時刻 t ，周波数 f での混合音の観測 x_{tf} と音源信号 y_{tf} を以下のように定義する．

$$\mathbf{x}_{tf} = [x_{tf1}, \dots, x_{tfM}]^T \in \mathbb{C}^M \quad (1)$$

$$\mathbf{y}_{tf} = [y_{tf1}, \dots, y_{tfK}]^T \in \mathbb{C}^K \quad (2)$$

表 1 モデル化方法の組み合わせ

音源モデル	重畳過程	
	混合モデル	因子モデル
混合モデル	LDA-LDA	LDA-NMF
因子モデル	NMF-LDA	NMF-NMF

このとき周波数領域での瞬時混合を仮定すると，音源 k のみを観測したときの観測スペクトル x_{tfk} は以下のように表される．

$$x_{tfk} = \mathbf{a}_{fk} \cdot y_{tfk} \quad (3)$$

ただし， \mathbf{a}_{fk} は周波数 f での音源 k の伝達関数である．ここで， y_{tfk} が次のような複素ガウス分布に従うとする．

$$y_{tfk} \sim \mathcal{N}_{\mathbb{C}}(0, \lambda_{tfk}) \quad (4)$$

$\lambda_{tfk} = \mathbb{E}[y_{tfk}^2]$ は時刻 t ，周波数 f での音源 k のパワーを表す．このとき音源 k のみを観測したときの観測 x_{tfk} は次のような複素ガウス分布に従う．

$$\mathbf{x}_{tfk} \sim \mathcal{N}_{\mathbb{C}}(\mathbf{0}, \lambda_{tfk} \mathbf{G}_{fk}^{-1}) \quad (5)$$

ここで， \mathbf{G}_{fk}^{-1} は周波数 f での音源 k の空間相関行列であり， $\mathbf{G}_{fk}^{-1} = \mathbf{a}_{fk} \mathbf{a}_{fk}^H$ である．ただし $*$ ^H はエルミート共役を示す．

2.2 重畳過程

本節では音の重畳過程に対して因子モデルと混合モデルを用いてモデル化を行う．

2.2.1 因子モデル

混合音の観測が各音源ごとの観測の和に等しいとすると式 (3) より，混合音の観測は次式で与えられる．

$$\mathbf{x}_{tf} = \sum_{k=1}^K \mathbf{x}_{tfk} = \sum_{k=1}^K \mathbf{a}_{fk} \cdot y_{tfk} \quad (6)$$

このとき，式 (5)，(6) とガウスの加法性により混合音の観測は次のような複素ガウス分布に従うとされる．

$$\mathbf{x}_{tf} | \lambda, \mathbf{S}, \mathbf{G} \sim \mathcal{N}_{\mathbb{C}}\left(\mathbf{0}, \sum_{k=1}^K \lambda_{tfk} \mathbf{G}_{fk}^{-1}\right) \quad (7)$$

ここで，空間相関行列は音源の種類ではなく音源の方向に依存するため，音源ごとの空間相関行列 \mathbf{G}_{fk} を方向ごとの空間相関行列 \mathbf{G}_{fd} の重み付き和を用いて表す．このとき，方向 d の音源 k への重みを r_{kd} とすると式 (7) は次のように表される．

$$\mathbf{x}_{tf} | \lambda, \mathbf{S}, \mathbf{G} \sim \mathcal{N}_{\mathbb{C}}\left(\mathbf{0}, \sum_{k=1}^K \sum_{d=1}^D \lambda_{tfk} r_{kd} \mathbf{G}_{fd}^{-1}\right), \quad (8)$$

2.2.2 混合モデル

音源スペクトログラムがスパースである，すなわち，各時間周波数ビンにおいて観測される音は高々一つであるとし，そのときの混合音の観測モデルについて考える．まず各時間周波数ビンにおいて観測される音源を示すための変数を $z_{tf} = [z_{tf1}, \dots, z_{tfK}]^T$ とする．ただし， z_{tf} は 1 of K 表現のベクトルであり，音源 k が観測されるときは $z_{tfk} = 1$ となり，それ以外のときは 0 となる．このとき，観測 \mathbf{x}_{tf} は次式で与えられる．

$$\mathbf{x}_{tf} = \prod_{k=1}^K \mathbf{x}_{tfk}^{z_{tfk}} = \prod_{k=1}^K (\mathbf{a}_{fk} \cdot \mathbf{y}_{tfk})^{z_{tfk}} \quad (9)$$

式 (5), (9) より, 観測は次の分布にしたがって生成される.

$$\mathbf{x}_{tf} \sim \prod_{k=1}^K \mathcal{N}_{\mathbb{C}}(\mathbf{0}, \lambda_{tfk} \mathbf{G}_{fk}^{-1})^{z_{tfk}} \quad (10)$$

ここで, 空間相関行列は音源の種類ではなく音源の方向に依存するため, 空間相関行列 \mathbf{G}_{fk} を音源ごとに独立な変数ではなく, 方向ごとに独立な変数 \mathbf{G}_{fd} として考える. このとき, 式 (10) は音源 k の方向を示すベクトル $\mathbf{s}_k = [s_{k1}, \dots, s_{kD}]^T$ を用いると次のように表される.

$$\mathbf{x}_{tf} \sim \prod_{k=1}^K \prod_{d=1}^D \mathcal{N}_{\mathbb{C}}(\mathbf{0}, \lambda_{tfk} \mathbf{G}_{fd}^{-1})^{z_{tfk} s_{kd}} \quad (11)$$

ただし, s_k は 1 of D 表現のベクトルであり, 音源 k が方向 d にあるときは $s_{kd} = 1$, それ以外のときは 0 となる.

2.3 音源モデル

本節では音源 k のパワー λ_{tfk} に対して因子モデルと混合モデルを用いてモデル化を行う.

2.3.1 因子モデル

音源 k のパワースペクトログラム λ_{tfk} が低ランク性を持つと仮定すると, パワースペクトログラムは基底スペクトル w_{klf} とアクティベーション h_{klt} の積の和で表現される. この基底スペクトルとアクティベーションはそれぞれ各基底の周波数的特徴とその基底の各時刻ごとの音の大きさを示す. 基底の数を L とするとパワースペクトログラムは次式で表される.

$$\lambda_{tfk} = \sum_l w_{klf} h_{klt} \quad (12)$$

2.3.2 混合モデル

混合モデルでも因子モデルのように基底スペクトル w_{klf} とアクティベーション h_{klt} を用いてパワースペクトログラム λ_{tfk} を表現する. ただし, 混合モデルでは全ての基底の和でパワーを表現するのに対し, 提案法では各時間周波数ビンごとに一つの基底 l' を選択することによりパワー λ_{tfk} を

$$\lambda_{tfk} = w_{kl'f} h_{kl't} \quad (13)$$

と表す. その基底を選択するための変数を $\mathbf{u}_{tfk} = [u_{tfk1}, \dots, u_{tfkL}]^T$ とする. ただし, \mathbf{u}_{tfk} は 1 of L 表現のベクトルであり, 基底 l が用いられるときは $u_{tfkl} = 1$, それ以外のときは 0 となる.

2.4 各モデルの尤度関数

本章では, 音源モデル-重畳過程を因子-混合モデル, 混合-混合モデル, 因子-因子モデル, 混合-因子モデルでそれぞれモデル化した NMF-LDA, LDA-LDA, NMF-NMF, LDA-NMF の観測モデルについて述べる.

2.4.1 NMF-LDA

NMF-LDA は重畳過程が混合モデル, 音源モデルが因子モデルなので式 (8), (13) より尤度関数は次のようになる.

$$\mathbf{x}_{tf} | \mathbf{W}, \mathbf{H}, \mathbf{Z}, \mathbf{S}, \mathbf{G}$$

$$\sim \prod_{k,d} \mathcal{N}_{\mathbb{C}} \left(\mathbf{0}, \sum_l w_{klf} h_{klt} \mathbf{G}_{fd}^{-1} \right)^{z_{tfk} s_{kd}} \quad (14)$$

2.4.2 LDA-LDA

LDA-LDA は重畳過程が混合モデル, 音源モデルが混合モデルなので式 (11), (13) より尤度関数は次のようになる.

$$\mathbf{x}_{tf} | \mathbf{W}, \mathbf{H}, \mathbf{Z}, \mathbf{S}, \mathbf{U}, \mathbf{G} \\ \sim \prod_{k,d,l} \mathcal{N}_{\mathbb{C}} \left(\mathbf{0}, w_{klf} h_{klt} \mathbf{G}_{fd}^{-1} \right)^{z_{tfk} s_{kd} u_{tfkl}} \quad (15)$$

2.4.3 NMF-NMF

NMF-NMF は重畳過程が因子モデル, 音源モデルが因子モデルなので式 (8), (12) より尤度関数は次のようになる.

$$\mathbf{x}_{tf} | \mathbf{W}, \mathbf{H}, \mathbf{R}, \mathbf{G} \\ \sim \mathcal{N}_{\mathbb{C}} \left(\mathbf{0}, \sum_{k,d,l} w_{klf} h_{klt} r_{kd} \mathbf{G}_{fd}^{-1} \right) \quad (16)$$

2.4.4 LDA-NMF

LDA-NMF は重畳過程が混合モデル, 音源モデルが因子モデルなので式 (11), (12) より尤度関数は次のようになる.

$$\mathbf{x}_{tf} | \mathbf{W}, \mathbf{H}, \mathbf{R}, \mathbf{G} \\ \sim \prod_l \mathcal{N}_{\mathbb{C}} \left(\mathbf{0}, \sum_{k,d} w_{klf} h_{klt} r_{kd} \mathbf{G}_{fd}^{-1} \right)^{u_{tfkl}} \quad (17)$$

しかし, このモデルでは音源 k に関するパラメータはガウス分布の分散パラメータとして因子分解により推定されるべきであるが, 基底の選択を行う u_{tfkl} が混合ガウス分布における分布を選択するパラメータとして存在している. これにより, u_{tfkl} は因子分解では推定することができないためこのモデル化方法では LDA-NMF を実現することができない.

2.5 事前分布の設計

提案法では式 (14)–(16) のパラメータに対し, それぞれ適切な事前分布を与えることで推論を行う. まず, z_{tf} , \mathbf{s}_k , \mathbf{u}_{tfk} はクラスタリングによる推論を行うためにカテゴリカル分布から生成されるとする.

$$z_{tf} | \boldsymbol{\pi}_t \sim \text{Categorical}(\boldsymbol{\pi}_t) \quad (18)$$

$$\mathbf{s}_k | \boldsymbol{\phi} \sim \text{Categorical}(\boldsymbol{\phi}) \quad (19)$$

$$\mathbf{u}_{tfk} | \boldsymbol{\psi}_{tk} \sim \text{Categorical}(\boldsymbol{\psi}_{tk}) \quad (20)$$

ここで, ハイパーパラメータ $\boldsymbol{\pi}_t$, $\boldsymbol{\phi}$, $\boldsymbol{\psi}_{tk}$ は観測に依存して変動するため推論を必要とする. したがって $\boldsymbol{\pi}_t$, $\boldsymbol{\phi}$, $\boldsymbol{\psi}_{tk}$ はカテゴリカル分布と共役なディリクレ分布から生成されるとする.

$$\boldsymbol{\pi}_t \sim \text{Dirichlet}(a_0^\pi \mathbf{1}_K) \quad (21)$$

$$\boldsymbol{\phi} \sim \text{Dirichlet}(a_0^\phi \mathbf{1}_D) \quad (22)$$

$$\boldsymbol{\psi}_{tk} \sim \text{Dirichlet}(a_0^\psi \mathbf{1}_L) \quad (23)$$

ここで, $\mathbf{1}_N$ は要素が全て 1 の N 次元のベクトルとし, a_0^* はハイパーパラメータとする. また, 空間相関行列 \mathbf{G}_{fd} , 基底 w_{klf} , アクティベーション h_{klt} , 重み r_{kd} は式 (14)–(16) と共

役になるように事前分布として次のような分布を与える．

$$\mathbf{G}_{fd} \sim \mathcal{W}_C(\nu, \mathbf{G}_{fd}^0) \quad (24)$$

$$w_{klf} \sim \text{Gamma}(a_0^w, b_0^w) \quad (25)$$

$$h_{klt} \sim \text{Gamma}(a_0^h, b_0^h) \quad (26)$$

$$r_{kd} \sim \text{Gamma}(a_0^r, b_0^r) \quad (27)$$

ここで， ν ， a_0^* ， b_0^* はハイパーパラメータであり， \mathcal{W}_C は複素ウィシャート分布 (付録参照) とする．

2.6 分離音の生成

各音源信号は，重畳過程が因子モデルのときはマルチチャネル Wiener フィルタ，混合モデルのときは時間周波数マスキングにより復元することができる．マルチチャネル Wiener フィルタでは音源 k の信号 \mathbf{x}_{tf}^k は以下の式により復元される．

$$\mathbf{x}_{tf}^k = \mathbf{Y}_{tfk} \mathbf{Y}_{tf}^{-1} \mathbf{x}_{tf} \quad (28)$$

ここで，モデルの推論の結果を $\mathbf{Y}_{tfk} = \sum_{d,l} w_{klf} h_{klt} r_{kd} \mathbf{G}_{fd}^{-1}$ ， $\mathbf{Y}_{tf} = \sum_k \mathbf{Y}_{tfk}$ と定義した．

時間周波数マスキングでは，時間周波数領域で音源方向ごとにマスクを推定することにより分離音を生成する．ギブスサンプリングで i 回目の試行により得られるサンプルを $\mathbf{W}^{(i)}$ ， $\mathbf{H}^{(i)}$ ， $\mathbf{G}^{(i)}$ ， $\mathbf{Z}^{(i)}$ ， $\mathbf{S}^{(i)}$ ， $(\mathbf{U}^{(i)})$ とすると，時刻 t 周波数 f での音源方向 d に対するマスク M_{tf}^d は次のようになる．

$$M_{tf}^d = \frac{1}{I} \sum_{i=1}^I \sum_{k=1}^K z_{tfk}^{(i)} s_{kd}^{(i)} \quad (29)$$

このマスクを用いると方向 d の音源の信号 \mathbf{x}_{tf}^d は

$$\mathbf{x}_{tf}^d = M_{tf}^d \mathbf{x}_{tf} \quad (30)$$

により復元される．

3. 推 論

提案法では，観測データ集合 \mathbf{X} に対するすべてのパラメータの事後分布 $p(\mathbf{G}, \mathbf{W}, \mathbf{H}, (\mathbf{Z}, \mathbf{S}, \mathbf{R}, \mathbf{U}, \pi, \phi, \psi) | \mathbf{X})$ を最大とするパラメータを求めることを目標とする．ただし，これらのパラメータを解析的に求めることは困難なので，提案法ではこれらのパラメータをギブスサンプリングにより求めることとする．ただし，提案法では π ， ϕ ， ψ は積分消去を行い，残りのパラメータ $\Theta = \{\mathbf{G}, \mathbf{W}, \mathbf{H}, (\mathbf{Z}, \mathbf{S}, \mathbf{R}, \mathbf{U})\}$ を求めることとする．ギブスサンプリングではそれぞれのパラメータの事後分布を求め，それらの事後分布からサンプリングを繰り返すことにより推定を行う．

3.1 NMF-LDA の更新

NMF-LDA では \mathbf{G}_{fd} ， z_{tf} ， s_k は尤度と事前分布の積により得られる以下の事後分布からサンプリングすることにより更新される．

$$\mathbf{G}_{fd} | \mathbf{X}, \Theta_{-\mathbf{G}_{fd}} \sim \mathcal{W}_C(\nu'_{fd}, \mathbf{G}'_{fd}) \quad (31)$$

$$z_{tf} | \mathbf{X}, \Theta_{-z_{tf}} \sim \text{Categorical}(\boldsymbol{\pi}'_{tf}) \quad (32)$$

$$s_k | \mathbf{X}, \Theta_{-s_k} \sim \text{Categorical}(\boldsymbol{\phi}'_k) \quad (33)$$

ここで， Θ_{-*} は Θ から $*$ の要素のみを除いた集合とする．また，ハイパーパラメータ ν'_{fd} ， \mathbf{G}'_{fd} ， $\boldsymbol{\pi}'_{tf}$ ， $\boldsymbol{\phi}'_k$ は次のようになる．

$$\nu'_{fd} = \nu + \sum_{t,k} z_{tfk} s_{kd} \quad (34)$$

$$\mathbf{G}'_{fd} = (\mathbf{G}_{fd}^0)^{-1} + \sum_{t,k} \frac{\mathbf{x}_{tf} \mathbf{x}_{tf}^H}{\sum_l w_{klf} h_{klt}} z_{tfk} s_{kd} \quad (35)$$

$$\boldsymbol{\pi}'_{tf} = \prod_d \left\{ \left| \frac{\mathbf{G}_{fd}}{\sum_l w_{klf} h_{klt}} \right| \exp \left(-\frac{\mathbf{x}_{tf}^H \mathbf{G}_{fd} \mathbf{x}_{tf}}{\sum_l w_{klf} h_{klt}} \right) \right\}^{s_{kd}} \times (a_0^\pi + n_{tk}^{-tf}) \quad (36)$$

$$\boldsymbol{\phi}'_k = \prod_{t,f} \left\{ \left| \frac{\mathbf{G}_{fd}}{\sum_l w_{klf} h_{klt}} \right| \exp \left(-\frac{\mathbf{x}_{tf}^H \mathbf{G}_{fd} \mathbf{x}_{tf}}{\sum_l w_{klf} h_{klt}} \right) \right\}^{z_{tfk}} \times (a_0^\phi + c_d^{-k}) \quad (37)$$

ここで， n_{tk}^{-tf} は時刻 t 周波数 f でのサンプルを除いて，時刻 t において音源 k に割り当てられた時間周波数ピンの数を表し， c_d^{-k} は音源 k を除いて方向 d に割り当てられた音源の数を表す．また， n_{tkl} は時刻 t において音源 k ，基底 l に割り当てられた時間周波数ピンの数を示し， n_{tkl}^{-tfk} は n_{tkl} から時刻 t 周波数 f で音源 k に割り当てられた要素を除いたものである．

w_{klf} と h_{klt} は事後分布から直接サンプルを得ることは難しいため，補助関数法を用いて得られる下限からサンプリングを行う．ただし，この下限は補助変数について最大化したとき，もとの対数尤度と一致する．まず，対数尤度は次式で与えられる．

$$\log p(\mathbf{X} | \mathbf{W}, \mathbf{H}, \mathbf{G}, \mathbf{Z}, \mathbf{S}) \stackrel{c}{=} \sum_{t,f,k,d} z_{tfk} s_{kd} \left(-M \log \left| \sum_l w_{klf} h_{klt} \right| - \frac{\mathbf{x}_{tf}^H \mathbf{G}_{fd} \mathbf{x}_{tf}}{\sum_l w_{klf} h_{klt}} \right) \quad (38)$$

この対数尤度に対しテイラー展開と Jensen の不等式を用いて下限を計算する [13]．まず一次のテイラー展開により，凸関数 $f(y) = -\log y$ の下限は以下のように計算される．

$$-\log y \geq -\log \alpha + 1 - \frac{y}{\alpha} \quad (39)$$

α は補助変数であり，等号は $\alpha = y$ のとき成立する．次に Jensen の不等式により，凹関数 $g(y) = -\frac{1}{y}$ に対し次の不等式が成立する．

$$-\frac{1}{\sum_k y_k} \geq -\sum_k \beta_k^2 \frac{1}{y_k} \quad (40)$$

β_k は $\beta_k \geq 0$ かつ $\sum_k \beta_k = 1$ を満たす補助変数であり，等号は $\beta_k = y_k (\sum_k y_k)^{-1}$ のとき成立する．

これらの二つの不等式により対数尤度の下限は次式により求められる．

$$\log p(\mathbf{X} | \mathbf{W}, \mathbf{H}, \mathbf{G}, \mathbf{Z}, \mathbf{S}) \stackrel{c}{\geq} \sum_{t,f,k,d} z_{tfk} s_{kd} \left(-M \left(\log |\alpha_{tf}| - 1 + \frac{\sum_l w_{klf} h_{klt}}{\alpha_{tf}} \right) - \mathbf{x}_{tf}^H \mathbf{G}_{fd} \mathbf{x}_{tf} \sum_l \frac{\beta_{tfl}}{w_{klf} h_{klt}} \right) \quad (41)$$

ここで, α_{tf} と β_{tfl} は補助変数であり, 次の式を満たすとき等号が成立する .

$$\alpha_{tf} = \sum_l w_{klf} h_{klt} \quad (42)$$

$$\beta_{tfl} = w_{klf} h_{klt} \left(\sum_l w_{klf} h_{klt} \right)^{-1} \quad (43)$$

w_{klf} と h_{klt} は式 (25), (26), (41) の積により得られる次の分布からサンプルされる .

$$w_{kfl} | \mathbf{X}, \Theta_{-w_{kfl}} \sim \text{GIG}(\gamma_{klf}^w, \rho_{kfl}^w, \tau_{kfl}^w), \quad (44)$$

$$h_{klt} | \mathbf{X}, \Theta_{-h_{klt}} \sim \text{GIG}(\gamma_{klt}^h, \rho_{klt}^h, \tau_{klt}^h), \quad (45)$$

ここで, GIG は一般化逆ガウス分布 (付録参照) を示す . また, ハイパーパラメータ γ^* , ρ^* , τ^* は次のようになる .

$$\gamma_{klf}^w = a_0^w, \quad (46)$$

$$\rho_{klf}^w = b_0^w + \sum_{t,d} \frac{M h_{klt}}{\alpha_{tfk}}, \quad (47)$$

$$\tau_{klf}^w = \sum_{t,d} \text{tr}(\mathbf{X}_{tf} \mathbf{G}_{fd}) \frac{\beta_{t f k l}^2}{h_{klt}}, \quad (48)$$

$$\gamma_{klt}^h = a_0^h, \quad (49)$$

$$\rho_{klt}^h = b_0^h + \sum_{f,d} \frac{M w_{klf}}{\alpha_{tfk}}, \quad (50)$$

$$\tau_{klt}^h = \sum_{f,d} \text{tr}(\mathbf{X}_{tf} \mathbf{G}_{fd}) \frac{\beta_{t f k l}^2}{w_{klf}}. \quad (51)$$

3.2 LDA-LDA の更新

LDA-LDA ではすべてのパラメータが, 尤度関数と事前分布の積により得られる事後分布からサンプルすることが可能である . 事後分布は次のように定められる .

$$\mathbf{G}_{fd} | \mathbf{X}, \Theta_{-\mathbf{G}_{fd}} \sim \mathcal{W}_C(\nu'_{fd}, \mathbf{G}'_{fd}) \quad (52)$$

$$\mathbf{z}_{tf} | \mathbf{X}, \Theta_{-\mathbf{z}_{tf}} \sim \text{Categorical}(\boldsymbol{\pi}'_{tf}) \quad (53)$$

$$s_k | \mathbf{X}, \Theta_{-s_k} \sim \text{Categorical}(\phi'_k) \quad (54)$$

$$\mathbf{u}_{t f k l} | \mathbf{X}, \Theta_{-\mathbf{u}_{t f k l}} \sim \text{Categorical}(\boldsymbol{\psi}'_{t f k l}) \quad (55)$$

$$w_{klf} | \mathbf{X}, \Theta_{-w_{klf}} \sim \text{GIG}(\gamma_{klf}^w, \rho_{klf}^w, \tau_{klf}^w) \quad (56)$$

$$h_{klt} | \mathbf{X}, \Theta_{-h_{klt}} \sim \text{GIG}(\gamma_{klt}^h, \rho_{klt}^h, \tau_{klt}^h) \quad (57)$$

ここで, ハイパーパラメータ ν'_{fd} , \mathbf{G}'_{fd} , $\boldsymbol{\pi}'_{tf}$, ϕ'_k , $\boldsymbol{\psi}$, γ^* , ρ^* , τ^* は次のようになる .

$$\nu'_{fd} = \nu + \sum_{t,k} z_{t f k} s_{kd} \quad (58)$$

$$\mathbf{G}'_{fd} = (\mathbf{G}_{fd}^0)^{-1} + \sum_{t,k,l} \frac{\mathbf{x}_{tf} \mathbf{x}_{tf}^H}{w_{klf} h_{klt}} z_{t f k} s_{kd} \mathbf{u}_{t f k l} \quad (59)$$

$$\begin{aligned} \boldsymbol{\pi}'_{t f k} &= \prod_{d,l} \left\{ \left| \frac{\mathbf{G}_{fd}}{w_{klf} h_{klt}} \right| \exp \left(- \frac{\mathbf{x}_{t f}^H \mathbf{G}_{fd} \mathbf{x}_{t f}}{w_{klf} h_{klt}} \right) \right\}^{s_{kd} \mathbf{u}_{t f k l}} \\ &\times (a_0^\pi + n_{t k}^{-t f}) \end{aligned} \quad (60)$$

$$\phi'_{kd} = \prod_{t,f,l} \left\{ \left| \frac{\mathbf{G}_{fd}}{w_{klf} h_{klt}} \right| \exp \left(- \frac{\mathbf{x}_{t f}^H \mathbf{G}_{fd} \mathbf{x}_{t f}}{w_{klf} h_{klt}} \right) \right\}^{z_{t f k} \mathbf{u}_{t f k l}}$$

$$\times (a_0^\phi + c_d^{-k}) \quad (61)$$

$$\begin{aligned} \boldsymbol{\psi}'_{t f k l} &= \prod_d \left\{ \left| \frac{\mathbf{G}_{fd}}{w_{klf} h_{klt}} \right| \exp \left(- \frac{\mathbf{x}_{t f}^H \mathbf{G}_{fd} \mathbf{x}_{t f}}{w_{klf} h_{klt}} \right) \right\}^{z_{t f k} s_{kd}} \\ &\times (a_0^\psi + n_{t k l}^{-t f k}) \end{aligned} \quad (62)$$

$$\gamma_{klf}^w = a_0^w - M n_{f k l} \quad (63)$$

$$\rho_{klf}^w = b_0^w \quad (64)$$

$$\tau_{klf}^w = \sum_{t,d} \frac{\mathbf{x}_{t f} \mathbf{G}_{fd} \mathbf{x}_{t f}}{h_{klt}} z_{t f k} s_{kd} \mathbf{u}_{t f k l} \quad (65)$$

$$\gamma_{klt}^h = a_0^h - M n_{t k l} \quad (66)$$

$$\rho_{klt}^h = b_0^h \quad (67)$$

$$\tau_{klt}^h = \sum_{f,d} \frac{\mathbf{x}_{t f} \mathbf{G}_{fd} \mathbf{x}_{t f}}{w_{klf}} z_{t f k} s_{kd} \mathbf{u}_{t f k l} \quad (68)$$

ここで, $n_{f k l}$ は周波数 f において音源 k , 基底 l に割り当てられた時間周波数ピンの数を示す .

3.3 NMF-NMF の更新

NMF-NMF ではどのパラメータも事後分布から直接サンプルを得ることは難しいため, 補助関数法を用いて得られる下限からサンプリングを行う . ただし, この補助関数は補助変数について最大化したとき, もとの対数尤度と一致する . まず, 対数尤度は次式で与えられる .

$$\begin{aligned} \log p(\mathbf{X} | \mathbf{W}, \mathbf{H}, \mathbf{R}, \mathbf{G}) &\stackrel{c}{=} \\ &\sum_{t f} (-\log |\mathbf{Y}_{t f}| - \text{tr}(\mathbf{X}_{t f} \mathbf{Y}_{t f}^{-1})) \end{aligned} \quad (69)$$

ここで, $\mathbf{X}_{t f} = \mathbf{x}_{t f}^H \mathbf{x}_{t f}$, $\mathbf{Y}_{t f k l d} = w_{klf} h_{klt} r_{kd} \mathbf{G}_{fd}^{-1}$, $\mathbf{Y}_{t f} = \sum_{k l d} \mathbf{Y}_{t f k l d}$ とした . 式 (69) の下限を求めするために二つの不等式を用いる [14] . まず, 多変量の凸関数 $f(\mathbf{Y}) = -\log |\mathbf{Y}|$ ($\mathbf{Y} \succeq \mathbf{0} \in \mathbb{C}^{M \times M}$) に対して一次の多変量のテイラー展開を用いると次の不等式が得られる .

$$-\log |\mathbf{Y}| \geq -\log |\boldsymbol{\Omega}| - \text{tr}(\boldsymbol{\Omega}^{-1} \mathbf{Y}) + M, \quad (70)$$

等号は $\boldsymbol{\Omega} = \mathbf{Z}$ のとき成立する . 次に, 凹関数 $g(\mathbf{Y}) = -\text{tr}(\mathbf{Y}^{-1} \mathbf{A})$ に対して次の不等式が成り立つ .

$$-\text{tr} \left(\left(\sum_{k=1}^K \mathbf{Y}_k \right)^{-1} \mathbf{A} \right) \geq -\sum_{k=1}^K \text{tr}(\mathbf{Y}_k^{-1} \boldsymbol{\Phi}_k \mathbf{A} \boldsymbol{\Phi}_k^H) \quad (71)$$

ここで $\{\boldsymbol{\Phi}_k\}_{k=1}^K$ は補助変数であり, $\sum_k \boldsymbol{\Phi}_k = \mathbf{I}$ を満たす . また, 等号は $\boldsymbol{\Phi}_k = \mathbf{Z}_k (\sum_{k'} \mathbf{Z}_{k'})^{-1}$ のとき成立する .

不等式 (70), (71) を用いると対数尤度 (式 (69)) の下限は次のように求められる .

$$\begin{aligned} \log p(\mathbf{X} | \mathbf{W}, \mathbf{H}, \mathbf{S}, \mathbf{G}) &\stackrel{c}{\geq} \sum_{t f} (-\text{tr}(\mathbf{Y}_{t f} \boldsymbol{\Omega}_{t f}^{-1}) - \log |\boldsymbol{\Omega}_{t f}| + M) \\ &- \sum_{t f k l d} \text{tr}(\mathbf{Y}_{t f k l d}^{-1} \boldsymbol{\Phi}_{t f k l d} \mathbf{X}_{t f} \boldsymbol{\Phi}_{t f k l d}) \end{aligned} \quad (72)$$

ここで, 補助変数 $\boldsymbol{\Omega}_{t f}$ と $\boldsymbol{\Phi}_{t f k l d}$ が次式を満たすとき等号が成立する .

$$\mathbf{\Omega}_{tf} = \mathbf{Y}_{tf} \quad (73)$$

$$\mathbf{\Phi}_{tfkld} = \mathbf{Y}_{tfkld} \mathbf{Y}_{tf}^{-1} \quad (74)$$

w_{klf} , h_{klt} , s_{kd} , \mathbf{G}_{fd} は式 (24)–(27), (72) の積により得られる次の分布からサンプルされる。

$$w_{klf} | \mathbf{X}, \Theta_{-w_{klf}} \sim \text{GIG}(\gamma_{klf}^w, \rho_{klf}^w, \tau_{klf}^w) \quad (75)$$

$$h_{klt} | \mathbf{X}, \Theta_{-h_{klt}} \sim \text{GIG}(\gamma_{klt}^h, \rho_{klt}^h, \tau_{klt}^h) \quad (76)$$

$$r_{kd} | \mathbf{X}, \Theta_{-r_{kd}} \sim \text{GIG}(\gamma_{kd}^r, \rho_{kd}^r, \tau_{kd}^r) \quad (77)$$

$$\mathbf{G}_{fd} | \mathbf{X}, \Theta_{-\mathbf{G}_{fd}} \sim \text{MGIG}_{\mathbb{C}}(\nu_0, \mathbf{R}_{fd}, \mathbf{U}_{fd}) \quad (78)$$

ここで, $\text{MGIG}_{\mathbb{C}}$ は complex matrix generalized inverse Gaussian (付録参照) を示す。ここで, ハイパーパラメータ γ^* , ρ^* , τ^* , \mathbf{R}_{fd} , \mathbf{U}_{fd} は次のように定められる。

$$\gamma_{klf}^w = a_0^w \quad (79)$$

$$\rho_{klf}^w = b_0^w + \sum_{td} h_{klt} r_{kd} \text{tr}(\mathbf{G}_{fd}^{-1} \mathbf{\Omega}_{tf}^{-1}) \quad (80)$$

$$\tau_{klf}^w = \sum_{td} h_{klt}^{-1} r_{kd}^{-1} \text{tr}(\mathbf{G}_{fd} \mathbf{\Phi}_{tfkld} \mathbf{X}_{tf} \mathbf{\Phi}_{tfkld}) \quad (81)$$

$$\gamma_{klt}^h = a_0^h \quad (82)$$

$$\rho_{klt}^h = b_0^h + \sum_{fd} w_{klf} r_{kd} \text{tr}(\mathbf{G}_{fd}^{-1} \mathbf{\Omega}_{tf}^{-1}) \quad (83)$$

$$\tau_{klt}^h = \sum_{fd} w_{klf}^{-1} r_{kd}^{-1} \text{tr}(\mathbf{G}_{fd} \mathbf{\Phi}_{tfkld} \mathbf{X}_{tf} \mathbf{\Phi}_{tfkld}) \quad (84)$$

$$\gamma_{klf}^r = a_0^r \quad (85)$$

$$\rho_{kd}^r = b_0^r + \sum_{tfl} w_{klf} h_{klt} \text{tr}(\mathbf{G}_{fd}^{-1} \mathbf{\Omega}_{tf}^{-1}) \quad (86)$$

$$\tau_{kd}^r = \sum_{tfl} w_{klf}^{-1} h_{klt}^{-1} \text{tr}(\mathbf{G}_{fd} \mathbf{\Phi}_{tfkld} \mathbf{X}_{tf} \mathbf{\Phi}_{tfkld}) \quad (87)$$

$$\mathbf{R}_{fd} = (\mathbf{G}_{fd}^0)^{-1} + \sum_{tkl} w_{klf}^{-1} h_{klt}^{-1} r_{kd}^{-1} \mathbf{\Phi}_{tfkld} \mathbf{X}_{tf} \mathbf{\Phi}_{tfkld} \quad (88)$$

$$\mathbf{U}_{fd} = \sum_{tkl} w_{klf} h_{klt} r_{kd} \mathbf{\Omega}_{tf}^{-1} \quad (89)$$

4. 評価実験

提案法の分離性能を評価するため, シミュレーションにより混合した音を用いた実験を行った。比較手法として, IVA [15], マルチチャンネル NMF (MNMF) [12], 音源モデルをスパースとし空間モデルに LDA を用いた分離法 (LDA) [9] を用いた。また, LDA では音源数の同時推定も行うが, 条件を対等にするため音源数は既知とした。

4.1 実験条件

図 2 に音源とマイクの配置を示す。残響時間 400 ms のインパルス応答を用いて 3 音源を混合した音声を用いた。マイク数は 4 とした。混合音には, 音声のみの混合音と音楽のみの混合音, 音声と音楽の混合音をそれぞれ 10 個ずつ使用した。用いる音楽と音声は SISEC [16] と JNAS の音素バランス文 [17] から選択した。サンプリング周波数は 16 kHz とし, STFT では窓幅 512 のハミング窓をシフト幅 256 で使用した。基底数 $L = 20$ とし, ハイパーパラメータは, $\nu = M + 1$, $a_0^\pi = a_0^\phi = 10$,

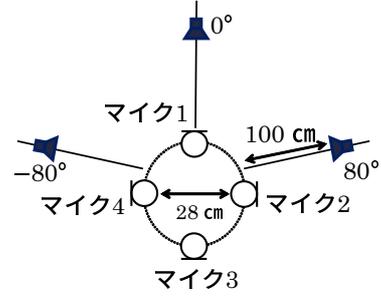


図 2 マイク配置

表 2 音楽による評価

	SDR	SIR	SAR
IVA	0.3 dB	4.9 dB	5.7 dB
MNMF	1.0 dB	6.2 dB	6.7 dB
LDA	0.7 dB	7.4 dB	4.1 dB
NMF-LDA	0.5 dB	8.7 dB	3.2 dB
LDA-LDA	0.7 dB	8.7 dB	3.3 dB
NMF-NMF	3.2 dB	8.1 dB	7.5 dB

表 3 音声による評価

	SDR	SIR	SAR
IVA	3.4 dB	7.5 dB	7.1 dB
MNMF	4.8 dB	10.0 dB	7.7 dB
LDA	5.5 dB	15.1 dB	6.3 dB
NMF-LDA	4.2 dB	14.0 dB	5.2 dB
LDA-LDA	5.8 dB	17.0 dB	6.3 dB
NMF-NMF	6.0 dB	12.6 dB	7.5 dB

表 4 音楽 + 音声による評価

	SDR	SIR	SAR
IVA	0.1 dB	5.3 dB	5.3 dB
MNMF	1.8 dB	8.6 dB	6.1 dB
LDA	2.4 dB	11.5 dB	4.5 dB
NMF-LDA	1.1 dB	9.8 dB	4.1 dB
LDA-LDA	2.8 dB	14.2 dB	3.9 dB
NMF-NMF	4.9 dB	13.0 dB	6.6 dB

$a_0^\psi = a_0^w = a_0^h = a_0^r = b_0^w = b_0^h = b_0^r = 1$ とした。また, $\mathbf{G}_{fd}^0 = (\mathbf{a}_{fd} \mathbf{a}_{fd}^H + 0.01 \times \mathbf{I})^{-1}$ とし, \mathbf{a}_{fd} には無響室で 5° 間隔で録音したインパルス応答を用いた。ギブスサンプリングの試行回数は 200 回とし, はじめの 180 回は burn-in として棄却した。評価尺度として, Signal-to-distortion ratio (SDR), Signal-to-inference ratio (SIR), Signal-to-artificial ratio (SAR) を用いた。SDR は総合的な分離性能, SIR は目的音以外の音の除去性能, SAR は分離音の歪みの少なさを表す尺度である。

4.2 実験結果

表 2, 3, 4 に実験結果を示す。それぞれの条件において最も数値が大きくなったものを太字で示した。SIR は LDA-LDA が最も高い性能を示したが, SDR と SAR は NMF-NMF が最も高い性能を示した。また, 提案法では事前分布のパラメータとして使用するマイクロホンアレイのインパルス応答を必要とするが, 無響室で録音したインパルス応答を用いても残響時間 400ms の環境下での混合音の分離ができたことから, 未知の環境下においても使用するマイクロホンアレイが同じであれば提案法により音源分離が可能であると考えられる。

5. おわりに

本稿では、音源モデル・重畳過程のそれぞれに対して、因子モデルと混合モデルを用いたモデル化を行うことで複数の音源分離法を提案した。それぞれのモデル化による音源分離法を比較すると、SDR, SAR の観点では音源モデルと重畳過程が共に因子モデルによりモデル化された手法が最も性能が高く、SIR の観点では音源モデルと重畳過程が共に混合モデルでモデル化された手法が最も性能が高いことを確認した。ただし、音源モデルが混合モデル、重畳過程が因子モデルでモデル化された手法はまだ実現できておらず、今後はそのモデル化を行う方法について考える必要がある。

付 録

複素ウィシャート分布と一般化逆ガウス分布の確率密度関数は以下のとおりである。

$$\mathcal{W}_c(\mathbf{G}|\nu, \mathbf{G}^0) = \frac{|\mathbf{G}|^{\nu-M} \exp(-\text{tr}(\mathbf{G}(\mathbf{G}^0)^{-1}))}{|\mathbf{G}^0|^{\nu} \pi^{M(M-1)/2} \prod_{m=0}^{M-1} \Gamma(\nu-m)} \quad (\text{A.1})$$

$$\text{GIG}(y|\gamma, \rho, \tau) = \frac{\exp\{(\gamma-1)\log y - \rho y - \tau/y\} \rho^{\gamma/2}}{2\tau^{\gamma/2} \mathcal{K}_\gamma(2\sqrt{\rho\tau})} \quad (\text{A.2})$$

ただし \mathcal{K}_γ は第 2 種変形ベッセル関数である。また、 MGIG_c の確率密度関数は以下の比例式で表される。

$$\text{MGIG}_c(\mathbf{X}|\gamma, \mathbf{R}, \mathbf{U}) \propto |\mathbf{X}|^{\gamma-M} \exp\{-\text{tr}(\mathbf{R}\mathbf{X} + \mathbf{U}\mathbf{X}^{-1})\} \quad (\text{A.3})$$

謝辞 本研究の一部は、JSPS 科研費 24220006, 15K12063 の支援を受けた。

文 献

- [1] A. Hyvärinen, J. Karhunen, and E. Oja, Independent component analysis, John Wiley & Sons, 2004.
- [2] I. Lee, T. Kim, and T. Lee, “Fast fixed-point independent vector analysis algorithms for convolutive blind source separation,” *Signal Processing*, vol.87, no.8, pp.1859–1871, 2007.
- [3] A. Hyvärinen, “Fast and robust fixed-point algorithms for independent component analysis,” *IEEE TNN*, vol.10, no.3, pp.626–634, 1999.
- [4] O. Yilmaz and S. Rickard, “Blind separation of speech mixtures via time-frequency masking,” *IEEE TSP*, vol.52, no.7, pp.1830–1847, 2004.
- [5] N. Ito, S. Araki, and T. Nakatani, “Permutation-free convolutive blind source separation via full-band clustering based on frequency-independent source presence priors,” *IEEE ICASSP*, pp.3238–3242, 2013.
- [6] H. Sawada, S. Araki, and S. Makino, “Underdetermined convolutive blind source separation via frequency bin-wise clustering and permutation alignment,” *IEEE TASLP*, vol.19, no.3, pp.516–527, 2011.
- [7] L. Drude, C. Boeddeker, and R. Haeb-Umbach, “Blind speech separation based on complex spherical k-mode clustering,” *IEEE ICASSP*, pp.141–145, 2016.
- [8] M.I. Mandel, R.J. Weiss, and D.P.W. Ellis, “Model-based expectation-maximization source separation and localization,” *IEEE TASLP*, vol.18, no.2, pp.382–394, 2010.
- [9] T. Otsuka, K. Ishiguro, H. Sawada, and H.G. Okuno, “Bayesian nonparametrics for microphone array processing,” *IEEE TASLP*, vol.22, pp.493–504, 2014.
- [10] K. Itakura, Y. Bando, E. Nakamura, K. Itoyama, and K. Yoshii, “A unified Bayesian model of time-frequency clustering and low-rank approximation for multi-channel source

separation,” *EUSIPCO*, pp.2280–2284, 2016.

- [11] P. Smaragdis and J.C. Brown, “Non-negative matrix factorization for polyphonic music transcription,” *IEEE WASPAA*, pp.177–180, 2003.
- [12] H. Sawada, H. Kameoka, S. Araki, and N. Ueda, “Multichannel extensions of non-negative matrix factorization with complex-valued data,” *IEEE TASLP*, vol.21, no.5, pp.971–982, 2013.
- [13] D.M. Blei, P.R. Cook, and M.D. Hoffman, “Bayesian non-parametric matrix factorization for recorded music,” *ICML*, pp.439–446, 2010.
- [14] K. Yoshii, K. Itoyama, and M. Goto, “Student’s t nonnegative matrix factorization and positive semidefinite tensor factorization for single-channel audio source separation,” *IEEE ICASSP*, pp.51–55, 2016.
- [15] N. Ono, “Stable and fast update rules for independent vector analysis based on auxiliary function technique,” *IEEE WASPAA*, pp.189–192, 2011.
- [16] S. Araki, F. Nesta, E. Vincent, Z. Koldovský, G. Nolte, A. Ziehe, and A. Benichoux, “The 2011 signal separation evaluation campaign (SiSEC2011):-audio source separation,” *Latent Variable Analysis and Signal Separation*, pp.414–422, Springer, 2012.
- [17] K. Itou, M. Yamamoto, K. Takeda, T. Takezawa, T. Matsuoaka, T. Kobayashi, K. Shikano, and S. Itahashi, “The design of the newspaper-based Japanese large vocabulary continuous speech recognition corpus,” *ICSLP*, pp.3261–3264, 1998.