

モノラル音響信号に対する音源分離のための 無限相関テンソル分解

吉井 和佳^{1,2,a)}

概要：本稿では、モノラルの音響信号に対して音源分離を行うための究極の因子分解法について述べる。従来は、非負値行列分解 (NMF) を用いて、混合音の各時刻におけるパワースペクトルを、少数の基底スペクトルの線形和で近似することがよく行われてきた。近年では、複素スペクトルに含まれる位相情報を適切に取り扱うため、半正定値テンソル分解 (PSDTF) が提案されている。PSDTF は、NMF の本質的な拡張となっており、混合音の各時刻における複素スペクトルから計算される局所的な共分散行列を、少数の基底共分散行列の線形和で近似する。この方法により、複素スペクトログラムにおける周波数方向の相関を考慮することができるようになったが、依然として時間方向の独立性が仮定されていた。複素スペクトログラムの周波数軸と時間軸を入れ替えて PSDTF を利用することも可能であるが、時間方向の相関は考慮できても、周波数方向の独立性を仮定せざるを得なかった。この問題を解決するため、本研究では、周波数方向の相関と時間方向との相関を同時に取り扱うことができる相関テンソル分解 (CTF) を提案する。CTF は、PSDTF のさらなる拡張となっており、複素スペクトログラムを構成するすべての時間周波数ビン間の共分散行列を、周波数方向の共分散行列と時間方向の共分散行列とのクロネッカー積の和で近似する。この方法により、すべての時間周波数ビン間の相関を考慮したウィナーフィルタを用いて、音源信号の複素スペクトログラムを一挙に推定することが可能となる。CTF は、NMF や PSDTF と同様に EM アルゴリズムあるいは補助関数法を用いた最尤推定や、ガンマ過程に基づくノンパラメトリックベイズ拡張 (基底数の無限化) が可能である。実験の結果、高品質な音源分離ができることを確かめた。

1. はじめに

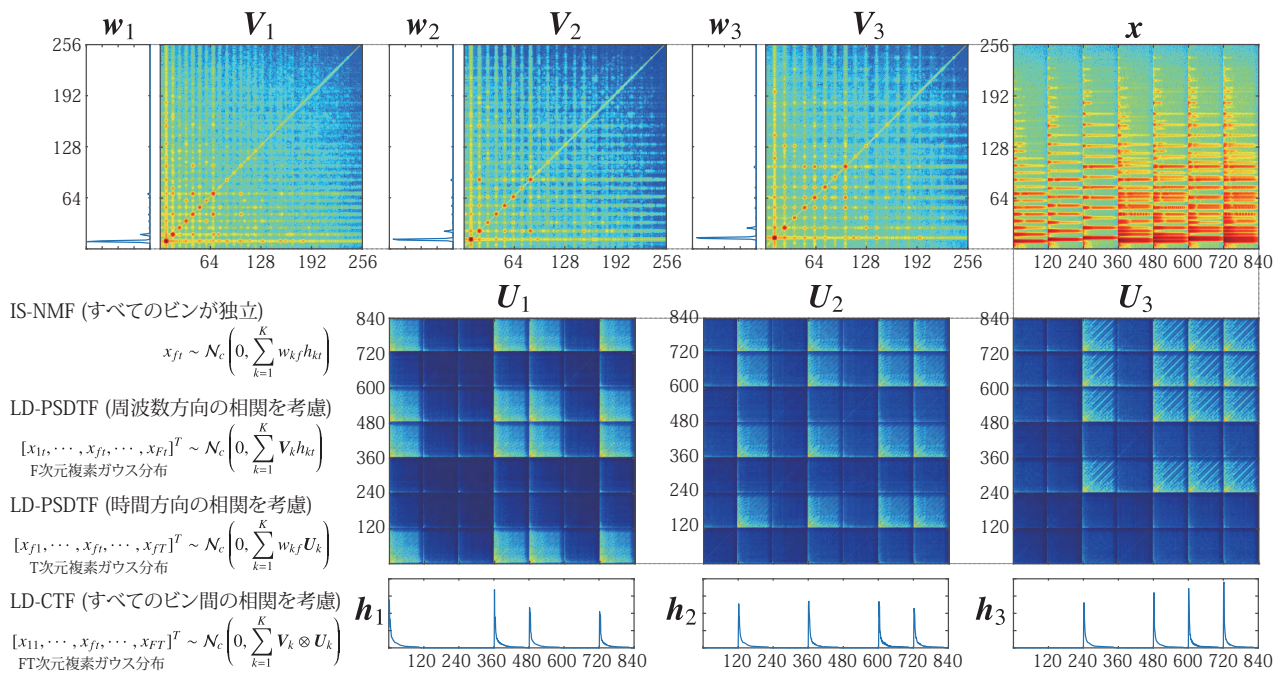
音楽音響信号に対する音源分離は、音楽情報処理分野における基盤技術のひとつである。音源分離が実現できれば、歌唱者の声質や性別、楽器構成などの楽曲の内容に基づいて、ユーザの好みに合う楽曲を検索できるようになる [1]。また、能動的音楽鑑賞 [2] として、混合音中に含まれる既存の楽器パートを自分好みに編集しながら音楽鑑賞を楽しむシステム [1-4] を実現することができる。このとき、分離音の品質の向上は重要な課題である。

モノラル音響信号に対する音源分離を行うには、非負値行列分解 (Nonnegative Matrix Factorization, NMF) [5] を用いるのが主流である。NMF は、入力となる非負値行列 (非負値スペクトログラム) を二つの非負値行列 (基底スペクトルの集合と対応する音量ベクトルとの集合) の積で近似することができる。NMF には多くの変種が存在するが、混合音の複素スペクトログラム中の時間周波数ビン

はすべて独立であり (現実には成立しない)、それぞれが異なる複素ガウス分布に従うという仮定のもとでは、混合音のパワースペクトログラムに対して Itakura-Saito (IS) ダイバージェンスに基づく NMF (IS-NMF) を適用することが理論的に妥当である。IS-NMF の結果に基づくウィナーフィルタを用いると、時間周波数ビンごとに独立に、混合音の複素成分を音源信号の複素成分の和に分解することができる。このとき、混合音と音源信号の複素スペクトログラムの位相は同一にならざるを得ず、復元される時間領域の音源信号の品質には限界があった。

時間領域信号に対応する「無矛盾な」複素スペクトログラムを推定する試みはいくつか存在する。Griffin ら [6] は、与えられた振幅スペクトログラムの位相を復元するため、その振幅スペクトログラムにできるだけ近い振幅スペクトログラムをもつ時間領域信号を推定できる反復短時間フーリエ変換 (STFT) 法を提案している。Le Roux ら [7] は、与えられた複素スペクトログラムの無矛盾性を評価する関数を提案し、それを最大化するアルゴリズムを導出している [8]。一方、亀岡ら [9] は、混合音の複素スペクトログラムを取り扱うことができる複素 NMF を提案し、無矛盾性に関する評価関数 [10] を組み込む拡張を行っている。ここ

¹ 京都大学 大学院情報学研究所 知能情報学専攻
Yoshida-hoftachi, Sakyo, Kyoto, Kyoto 606-8501, Japan
² 理研 革新知能統合研究センター (AIP) 音響情景理解チーム
15F, 1-4-1 Nihonbashi, Chuo, Tokyo 103-0027, Japan
a) yoshii@kuis.kyoto-u.ac.jp



*混合音の複素スペクトログラム $x \in \mathbb{C}^{FT}$ に対し、IS-NMF は、周波数方向の非負値ベクトル群 $W = \{w_1, w_2, w_3 \in \mathbb{R}_+^F\}$ および時間方向の非負値ベクトル群 $H = \{h_1, h_2, h_3 \in \mathbb{R}_+^T\}$ を推定する。LD-PSDTF は、周波数方向の半正定値行列群 $V = \{V_1, V_2, V_3 \in \mathcal{S}_+^F\}$ および時間方向の非負値ベクトル群 H 、あるいは周波数方向の非負値ベクトル群 W および時間方向の半正定値行列群 $U = \{U_1, U_2, U_3 \in \mathcal{S}_+^T\}$ を推定する。LD-CTF は、周波数方向の半正定値行列群 V および時間方向の半正定値行列群 U を推定する（図中では $F = 256, T = 840$ ）。

図 1 時間周波数領域での音源分離における IS-NMF, LD-PSDTF, LD-CTF の比較。

で、複素スペクトログラムが無矛盾であることが、対応する時間領域信号が高品質であることを必ずしも意味しないことに注意されたい。このことは、周波数領域における位相復元は容易ではないことを示唆している。

近年、複素スペクトログラムに含まれる位相情報を取り扱うため、半正定値テンソル分解 (Positive Semidefinite Tensor Factorization, PSDTF) [11, 12] と呼ばれる因子分解法が提案されている。NMF は、非負値ベクトルを少数の非負値基底ベクトルの線形和で近似するのに対し、PSDTF は、半正定値行列を少数の半正定値基底行列の線形和で近似する。PSDTF にも様々な変種が考えられるが、LogDet (LD) ダイバージェンスに基づく PSDTF (LD-PSDTF) が IS-NMF の自然な拡張となっている。音源分離において、IS-NMF は、各時刻の混合音のパワースペクトルを、各音源の典型的なパワースペクトルパターンの線形和で近似するのに対し、LD-PSDTF は、各時刻の混合音の複素スペクトルから計算される共分散行列を、各音源の複素スペクトルの共分散行列の線形和で近似する。ここで、共分散行列の対角成分がパワースペクトルに対応することに注目されたい。LD-PSDTF では、周波数間の相関を考慮するウィナーフィルタを用いて、時間ごとに独立に、混合音の複素スペクトルを音源信号の複素スペクトルの和に分解することができる。実は、このような周波数領域での分解は、時間領域での分解に対応しており、高品質な時間領域の音源信号を推定できる理由となっている。

実際の音響信号の複素スペクトログラムには、周波数方向だけではなく、時間方向にも相関が存在する。理論上は、定常信号に対してフーリエ変換を用いると、周波数間は独立となることが知られている。そこで、実際の非定常な音響信号に対しては、局所的な定常性を仮定して STFT を行うのが一般的である。しかし、厳密には、窓関数を用いて切り出された短時間信号は定常ではないため、対応する複素スペクトルにおいては、倍音関係や隣接関係にある周波数ビン間には強い相関が発生することが避けられない。また、音響信号は時系列データであるという性質上、時間方向にも相関が存在する。しかし、LD-PSDTF では、周波数方向の相関を考慮することができるが、時間方向の独立性が仮定されていた。混合音の複素スペクトログラムの周波数軸と時間軸を入れ替えて LD-PSDTF を利用することも可能であるが、時間方向の相関は考慮できても、周波数方向の独立性を仮定せざるを得なかった。

本稿では、NMF, PSDTF に続く究極の因子分解技法として、時間周波数上の相関を完全に取り扱える相関テンソル分解 (Correlated Tensor Factorization, CTF) を提案する。CTF も NMF や PSDTF と同様に様々な変種を考慮ことができ、本稿では、LD ダイバージェンスに基づく CTF (LD-CTF) を取り上げる。図 1 に IS-NMF, LD-PSDTF, LD-CTF の比較を示す。LD-CTF は、半正定値行列を、少数の半正定値行列のペアのクロネッカー積の和で近似する、すなわち、混合音の複素スペクトログラム中の全ての時間

周波数ビン間の超大規模な共分散行列を、各音源に対応する周波数方向の共分散行列と時間方向の共分散行列のクロネッカー積の和で近似する。この結果を用いると、混合音の複素スペクトログラムに対し、全ての時間周波数ビン間の相関を考慮したウィナーフィルタを適用することで、音源の複素スペクトログラムを一挙に得ることができる。

LD ダイバージェンス最小化に基づく LD-CTF は、全ての時間周波数ビン上の超高次元の多変量複素ガウス分布の共分散行列パラメータ（周波数方向の共分散行列と時間方向の共分散行列の少数のペア）の最尤推定を行うことと等価である。本稿では、IS-NMF や LD-PSDTF と同様に、収束性の保証された Expectation-Maximization (EM) アルゴリズムおよび Minorization-Maximization (MM) アルゴリズムを導出する。実際には、適切な音源数を事前に設定することが困難である場合が多い。そこで、観測データに合わせて音源数を自動的に調整するため、ガンマ過程に基づくノンパラメトリックベイズモデル（無限モデル）の定式化を行う。さらに、パラメータの事後分布を推定するため、補助関数を用いた変分ベイズ法を導出する。

理論的には、LD-CTF は究極の分解技法ではあるが、数百万次元を超えるような巨大行列の計算は現実的ではないため、計算量削減のための工夫が必要不可欠である。まず、周波数方向および時間方向の共分散行列をともにブロック対角行列に限定した LD-CTF の近似モデルを提案する。これは、局所的な時間周波数領域に限定して完全な相関を考慮することに相当する。また、これまで時間周波数領域における LD-CTF について議論してきたが、複素スペクトログラムの周波数軸あるいは時間軸に対して、任意の線形変換を施した空間（例：観測音響信号の各フレーム信号を並べただけの時間時間領域や複素スペクトログラムの時間軸に対してさらにフーリエ変換を施した周波数周波数領域など）においても等価な分解が得られる。このことから、周波数方向の共分散行列群を同時にほぼ対角化できるような線形変換と、また、時間方向の共分散行列群を同時にほぼ対角化できるような線形変換をそれぞれ見つけることができれば、変換後の空間における IS-NMF が時間周波数領域での LD-CTF の良い近似になると考えられる。

2. 因子分解と音源分離

本章では、従来の因子分解技法である NMF や PSDTF について説明する。特に、音源分離を行ううえで、IS-NMF [13] や LD-PSDTF [11] の理論的な裏付けについて説明する。

2.1 非負値行列分解 (NMF)

NMF では、非負値行列 $X \in \mathbb{R}_+^{F \times T}$ を二つの非負値行列 $W \in \mathbb{R}_+^{K \times F}$ および $H \in \mathbb{R}_+^{K \times T}$ の積 Y で近似する。

$$X \approx Y \stackrel{\text{def}}{=} W^T H \quad (1)$$

ただし、 $K \ll \min(F, T)$ は基底数、 $Y \in \mathbb{R}_+^{F \times T}$ は低ランクな再構成行列を表す。行列積表現である式 (1) は、各要素ごとに書き直すことができる。

$$x_{ft} \approx y_{ft} \stackrel{\text{def}}{=} \sum_{k=1}^K w_{kf} h_{kt} \quad (2)$$

ここで、 $y_{kft} = w_{kf} h_{kt}$ と定義すると、 $y_{ft} = \sum_k y_{kft}$ が成立する。観測値 x_{ft} と再構成値 y_{ft} との間の誤差 $C(x_{ft}|y_{ft})$ を評価する尺度のひとつに、Bregman ダイバージェンス [14] が知られている。

$$C_\phi(x_{ft}|y_{ft}) = \phi(x_{ft}) - \phi(y_{ft}) - \phi'(y_{ft})(x_{ft} - y_{ft}) \quad (3)$$

ここで、 ϕ は厳密に凸な関数である。常に $C_\phi(x_{ft}|y_{ft}) \geq 0$ が成り立ち、 $x_{ft} = y_{ft}$ であるときに限り $C_\phi(x_{ft}|y_{ft}) = 0$ となる。その特別な場合として、 $\phi(x) = x \log x - x$ の場合の Kullback-Leibler (KL) ダイバージェンスや、 $\phi(x) = -\log x$ の場合の IS ダイバージェンスがよく知られている。

$$C_{\text{KL}}(x_{ft}|y_{ft}) = x_{ft} \log \frac{x_{ft}}{y_{ft}} - x_{ft} + y_{ft} \quad (4)$$

$$C_{\text{IS}}(x_{ft}|y_{ft}) = \frac{x_{ft}}{y_{ft}} - \log \frac{x_{ft}}{y_{ft}} - 1 \quad (5)$$

NMF のコスト関数 $C_\phi(X|Y)$ は次式で与えられる。

$$C_\phi(X|Y) = \sum_{f=1}^F \sum_{t=1}^T C_\phi(x_{ft}|y_{ft}) \quad (6)$$

このコスト関数を最小化する W および H を求めるため、乗法更新 (Multiplicative Update, MU) 型の MM アルゴリズムが提案されている [15]。これは、ある確率モデルの最尤推定を行うことと等価である。一方、ガンマ事前分布を導入してベイズ推定を行うこともできる [16, 17]。

2.2 IS-NMF を用いた音源分離

与えられた混合音の複素スペクトログラム $S \in \mathbb{C}^{F \times T}$ (F は周波数ビン数、 T はフレーム数) を K 個の音源信号の複素スペクトログラムの和に分解することを考える。音源 k の複素スペクトログラムを $Z_k \in \mathbb{C}^{F \times T}$ とし、周波数領域での瞬時混合過程を仮定すると、以下が成り立つ。

$$S = \sum_{k=1}^K Z_k \quad (7)$$

IS-NMF では、潜在変数 z_{kft} が、 y_{kft} を分散パラメータとする複素ガウス分布に従うことを仮定する。

$$z_{kft}|y_{kft} \sim \mathcal{N}_c(0, y_{kft}) \quad (8)$$

ここで、 $s_{ft} = \sum_k z_{kft}$ かつ $y_{ft} = \sum_k y_{kft}$ であり、複素ガウス分布の再生性から、 s_{ft} も複素ガウス分布に従う。

$$s_{ft}|y_{ft} \sim \mathcal{N}_c(0, y_{ft}) \quad (9)$$

これは, s_{ft} の位相に関わらず, 時刻 t ・周波数 f におけるパワー $x_{ft} \stackrel{\text{def}}{=} |s_{ft}|^2$ が指数分布に従うことと等価である.

$$x_{ft}|y_{ft} \sim \text{Exponential}(y_{ft}) \quad (10)$$

最終的に, 観測される複素スペクトログラム S に対する対数尤度関数を導出できる.

$$\begin{aligned} \log p(S|Y) &= \sum_{f=1}^F \sum_{t=1}^T \log p(s_{ft}|y_{ft}) \\ &= \sum_{f=1}^F \sum_{t=1}^T \left(-\frac{x_{ft}}{y_{ft}} - \log y_{ft} \right) \\ &\stackrel{c}{=} -C_{\text{IS}}(X|Y) \end{aligned} \quad (11)$$

したがって, 対数尤度関数の最大化 $p(S|Y)$ は, IS ダイバージェンス $C_{\text{IS}}(X|Y)$ の最小化と等価である.

パラメータ Y (W および H) が推定できれば, 式 (8) および式 (9) から, 観測変数 S が与えられたもとで, 潜在変数 Z_k の事後分布を求めることができる.

$$p(z_{kft}|s_{ft}) = \mathcal{N}_c \left(z_{kft} \left| y_{kft} y_{ft}^{-1} s_{ft}, y_{ft} - y_{kft}^2 y_{ft}^{-1} \right. \right) \quad (12)$$

この処理はウィナーフィルタと呼ばれ, Z_k と S の位相は同一となる. これは NMF に基づく音源分離の性質であり, 音源信号の合成品質に限界がある根本的な原因のひとつとなっている. 最後に, 重畳加算合成法 [18] を用いれば, $\mathbb{E}[Z_k|S]$ から音源信号を復元することができる.

2.3 半正定値テンソル分解 (PSDTF)

PSDTF では, 観測データとして, 特別な形式を持つ 3 階のテンソル $X = [X_1, \dots, X_t] \in \mathbb{C}^{F \times F \times T}$ が与えられるものとする. ここで, 各要素 $X_t \in S_+^F$ は半正定値行列であるとする. PSDTF の目標は, それぞれの半正定値行列 X_t を K 個の半正定値行列 $\{V_k \in S_+^F\}_{k=1}^K$ (基底行列とよぶ) の線形和で近似することである.

$$X_t \approx Y_t \stackrel{\text{def}}{=} \sum_{k=1}^K h_{kt} V_k \quad (13)$$

ここで, $Y_{kt} = h_{kt} V_k$ と定義すると, $Y_t = \sum_k Y_{kt}$ が成立する. $h_{kt} \geq 0$ は n 番目の要素 X_t における基底行列 V_k の重みである. 観測行列 X_t と再構成行列 Y_t との間の誤差 $C_\phi(X_t|Y_t)$ を評価する尺度として, Bregman 行列ダイバージェンス [14] が利用できる.

$$\begin{aligned} C_\phi(X_t|Y_t) &= \phi(X_t) - \phi(Y_t) \\ &\quad - \text{tr}(\nabla \phi(Y_t)^T (X_t - Y_t)) \end{aligned} \quad (14)$$

ここで, ϕ は微分可能で厳密に凸な関数である. 式 (14) は式 (3) の自然な多次元拡張となっており, 常に $C_\phi(X_t|Y_t) \geq 0$ が成り立ち, $X_t = Y_t$ であるときに限り $C_\phi(X_t|Y_t) = 0$ となる. その特別な場合として, $\phi(Z) = \text{tr}(X \log X - X)$

の場合の von Neumann (vN) ダイバージェンスや, $\phi(Z) = -\log|Z|$ の場合の LogDet (LD) ダイバージェンスなどが知られている [19].

$$C_{\text{vN}}(X_t|Y_t) = \text{tr}(X \log X - X \log Y - X + Y) \quad (15)$$

$$C_{\text{LD}}(X_t|Y_t) = -\log|X_t Y_t^{-1}| + \text{tr}(X_t Y_t^{-1}) - M \quad (16)$$

PSDTF のコスト関数 $C_\phi(X|Y)$ は次式で与えられる.

$$C_\phi(X|Y) = \sum_{t=1}^T C_\phi(X_t|Y_t) \quad (17)$$

LD-PSDTF に関しては, このコスト関数を最小化する V および H を求めるため, 乗法更新 (Multiplicative Update, MU) 型の MM アルゴリズムが提案されている [11,12]. これは, ある確率モデルの最尤推定を行うことと等価であり, 共役事前分布を導入してベイズ推定を行うこともできる.

2.4 LD-PSDTF を用いた音源分離

LD-PSDTF では, 音源信号 k の時刻 t における複素ベクトル (潜在変数) を $z_{kt} = [z_{k1t}, \dots, z_{kFt}]^T \in \mathbb{C}^F$ とすると, z_{kt} は $Y_{kt} \in S_+^F$ を共分散行列パラメータとする多変量複素ガウス分布に従うことを仮定する.

$$z_{kt}|Y_{kt} \sim \mathcal{N}_c(\mathbf{0}, Y_{kt}) \quad (18)$$

ここで, 混合音の時刻 t における複素スペクトル (観測変数) を $s_t = [s_{1t}, \dots, s_{Ft}]^T \in \mathbb{C}^F$ とすると, $s_t = \sum_k z_{kt}$ かつ $Y_t = \sum_k Y_{kt}$ であることと, 複素ガウス分布の再生性から, s_t も多変量複素ガウス分布に従う.

$$s_t|Y_t \sim \mathcal{N}_c(\mathbf{0}, Y_t) \quad (19)$$

最終的に, 混合音の時刻 t における局所的な共分散行列を $X_t \stackrel{\text{def}}{=} s_t s_t^H$ とすると, 観測される複素スペクトログラム S に対する対数尤度関数を導出できる.

$$\begin{aligned} \log p(S|Y) &= \sum_{t=1}^T \log p(s_t|Y_t) \\ &= \sum_{t=1}^T -\log|Y_t| - \text{tr}(X_t Y_t^{-1}) \\ &\stackrel{c}{=} -C_{\text{LD}}(X|Y) \end{aligned} \quad (20)$$

したがって, 対数尤度関数 $p(S|Y)$ の最大化は, LD ダイバージェンス $C_{\text{LD}}(X|Y)$ の最小化と等価である.

パラメータ Y (V および H) が推定できれば, 式 (18) および式 (19) から, 観測変数 S が与えられたもとで, 潜在変数 Z_k の事後分布を求めることができる.

$$p(z_{kt}|s_t) = \mathcal{N}_c \left(z_{kt} \left| Y_{kt} Y_t^{-1} s_t, Y_t - Y_{kt} Y_t^{-1} Y_{kt} \right. \right) \quad (21)$$

このウィナーフィルタを用いると, S から Z_k の位相をある程度正しく復元することができる.

3. 相関テンソル分解

本章では、相関テンソル分解 (Correlated Tensor Factorization, CTF) とよぶ新しい因子分解法について説明する。CTF は NMF や PSDTF の本質的な拡張となっており、最尤推定やベイズ推定が可能である。

3.1 定式化

CTF では、観測データとして、半正定値行列 $\mathbf{X} \in S_+^{FT}$ が与えられるものとする。CTF の目標は、半正定値行列 \mathbf{X} を半正定値行列の集合 $\{\mathbf{V}_k \in S_+^F\}_{k=1}^K$ と対応する半正定値行列の集合 $\{\mathbf{U}_k \in S_+^{FT}\}_{k=1}^K$ とのクロネッカー積の和で近似することである。

$$\mathbf{X} \approx \mathbf{Y} \stackrel{\text{def}}{=} \sum_{k=1}^K \mathbf{V}_k \otimes \mathbf{U}_k \quad (22)$$

ここで、 $\mathbf{Y}_k = \mathbf{V}_k \otimes \mathbf{U}_k$ と定義すると、 $\mathbf{Y} = \sum_k \mathbf{Y}_k$ が成立する。これらすべての半正定値行列が対角行列のとき、式 (22) で定義される CTF は、式 (2) で定義される NMF に帰着する。また、 $\{\mathbf{V}_k\}_{k=1}^K$ あるいは $\{\mathbf{U}_k\}_{k=1}^K$ のいずれかが対角行列の集合であるとき、式 (22) で定義される CTF は、式 (13) で定義される PSDTF に帰着する。CTF では、PSDTF と同様に、観測行列 \mathbf{X} と再構成行列 \mathbf{Y} との間の誤差 $\mathcal{C}_\phi(\mathbf{X}|\mathbf{Y})$ を評価する尺度として、Bregman 行列ダイバージェンス [14] を利用する。

$$\mathcal{C}_\phi(\mathbf{X}|\mathbf{Y}) = \phi(\mathbf{X}) - \phi(\mathbf{Y}) - \text{tr}(\nabla\phi(\mathbf{Y})^T(\mathbf{X} - \mathbf{Y})) \quad (23)$$

CTF では、 $\mathcal{C}_\phi(\mathbf{X}|\mathbf{Y})$ を最小化する \mathbf{V} および \mathbf{U} を求めることが目標となる。本稿では特に、モノラル音響信号の音源分離を行うのに適した LD ダイバージェンス $\mathcal{C}_{\text{LD}}(\mathbf{X}|\mathbf{Y})$ に基づく CTF (LD-CTF) について議論する。

3.2 LD-CTF を用いた音源分離

LD-CTF では、音源信号 k の複素スペクトログラム \mathbf{Z}_k の全ての時間周波数ピンを並べたベクトル (潜在変数) を $\mathbf{z}_k = [z_{k11}, \dots, z_{k1T}, \dots, z_{kF1}, \dots, z_{kFT}]^T \in \mathbb{C}^{FT}$ とし、 \mathbf{z}_k が $\mathbf{Y}_k \in S_+^{FT}$ を共分散行列パラメータとする多変量複素ガウス分布に従うことを仮定する。

$$\mathbf{z}_k | \mathbf{Y}_k \sim \mathcal{N}_c(\mathbf{0}, \mathbf{Y}_k) \quad (24)$$

ここで、混合音の複素スペクトログラム \mathbf{S} (観測変数) を同様に展開し、 $\mathbf{s} = [s_{11}, \dots, s_{1T}, \dots, s_{F1}, \dots, s_{FT}]^T \in \mathbb{C}^{FT}$ とすると、 $\mathbf{s} = \sum_k \mathbf{z}_k$ かつ $\mathbf{Y} = \sum_k \mathbf{Y}_k$ であることと、複素ガウス分布の再生性から、次式が成立する。

$$\mathbf{s} | \mathbf{Y} \sim \mathcal{N}_c(\mathbf{0}, \mathbf{Y}) \quad (25)$$

したがって、 $\mathbf{X} \stackrel{\text{def}}{=} \mathbf{s}\mathbf{s}^H$ とすると、観測される複素スペクトログラム \mathbf{S} に対する対数尤度関数は、次式で与えられる。

$$\begin{aligned} \log p(\mathbf{S}|\mathbf{Y}) &= \log p(\mathbf{s}|\mathbf{Y}) \\ &= -\log |\mathbf{Y}| - \text{tr}(\mathbf{X}\mathbf{Y}^{-1}) \\ &\stackrel{c}{=} -\mathcal{C}_{\text{LD}}(\mathbf{X}|\mathbf{Y}) \end{aligned} \quad (26)$$

したがって、対数尤度関数 $p(\mathbf{S}|\mathbf{Y})$ の最大化は、LD ダイバージェンス $\mathcal{C}_{\text{LD}}(\mathbf{X}|\mathbf{Y})$ の最小化と等価である。

パラメータ \mathbf{Y} (\mathbf{V} および \mathbf{U}) が推定できれば、式 (24) および式 (25) から、観測変数 \mathbf{S} が与えられたもとで、潜在変数 \mathbf{Z}_k の事後分布を求めることができる。

$$p(\mathbf{z}_k | \mathbf{s}) = \mathcal{N}_c\left(\mathbf{z}_k \mid \mathbf{Y}_k \mathbf{Y}^{-1} \mathbf{s}, \mathbf{Y} - \mathbf{Y}_k \mathbf{Y}^{-1} \mathbf{Y}_k\right) \quad (27)$$

このウィナーフィルタでは、時間周波数ビン間の完全な相関が考慮されている。また、 $\mathbf{Z}_k \stackrel{\text{def}}{=} \mathbf{z}_k \mathbf{z}_k^H$ と定義しておく。

3.3 補助関数法

LD ダイバージェンス $\mathcal{C}_{\text{LD}}(\mathbf{X}|\mathbf{Y})$ を \mathbf{V} および \mathbf{U} に関して直接最小化することは困難であるため、補助関数法 [15] を用いて、 $\mathcal{C}_{\text{LD}}(\mathbf{X}|\mathbf{Y})$ を間接的に最小化することを考える。いま、 θ に関して最小化したい目的関数 $\mathcal{F}(\theta)$ に対して、

$$\mathcal{F}(\theta) \leq \mathcal{F}^+(\theta, \phi) \quad (28)$$

を満たす上限関数 $\mathcal{F}^+(\theta, \phi)$ を $\mathcal{F}(\theta)$ の補助関数と呼ぶ。ここで、 ϕ は補助パラメータであり、任意の値に対して不等式が成立する必要がある。このとき、以下の反復更新則

$$\phi^{\text{new}} \leftarrow \text{argmin}_\phi \mathcal{F}^+(\theta^{\text{old}}, \phi) \quad (29)$$

$$\theta^{\text{new}} \leftarrow \text{argmin}_\theta \mathcal{F}^+(\theta, \phi^{\text{new}}) \quad (30)$$

を用いると、 $\mathcal{F}(\theta)$ は単調非増加となる (式 (29) により $\mathcal{F}^+(\theta, \phi)$ が最大化されて $\mathcal{F}(\theta)$ に一致する)。このアルゴリズムの収束性は保証されており、IS-NMF や LD-PSDTF のベイズ学習でも同様の手法が利用されている [11, 17]。対数尤度関数 $p(\mathbf{S}|\mathbf{Y})$ を最大化する場合には、下限関数を補助関数として、下限関数を逐次最大化すればよい。

3.4 Expectation-Maximization アルゴリズム

LD-CTF の確率モデル (3.2 節) を考えることで、補助関数法の一つである EM アルゴリズムを用いた最尤推定を行うことができる。対数尤度関数 $p(\mathbf{S}|\mathbf{Y})$ を直接最大化することは困難であるので、下限関数を設計する。

$$\begin{aligned} \log p(\mathbf{S}|\mathbf{Y}) &= \log \int q(\mathbf{Z}) \frac{p(\mathbf{S}, \mathbf{Z}|\mathbf{V}, \mathbf{U})}{q(\mathbf{Z})} d\mathbf{Z} \\ &\geq \int q(\mathbf{Z}) \log \frac{p(\mathbf{S}, \mathbf{Z}|\mathbf{V}, \mathbf{U})}{q(\mathbf{Z})} d\mathbf{Z} \\ &= \int q(\mathbf{Z}) \log p(\mathbf{S}, \mathbf{Z}|\mathbf{V}, \mathbf{U}) d\mathbf{Z} - \int q(\mathbf{Z}) \log q(\mathbf{Z}) d\mathbf{Z} \\ &\stackrel{\text{def}}{=} \mathcal{L}(q(\mathbf{Z}), \mathbf{V}, \mathbf{U}) \end{aligned} \quad (31)$$

ここで、 $q(\mathbf{Z})$ は潜在変数 \mathbf{Z} に関する変分事後分布である。

まず, E ステップにおいて, V および U が既知のもとで, $\mathcal{L}(q(Z), V, U)$ を最大化する $q(Z)$ を求める. 具体的には, $q(Z)$ が真の事後分布 $p(Z|S, V, U)$ と一致するとき (式 (27) 参照), $\mathcal{L}(q(Z), V, U)$ は最大化される.

$$q(Z) = p(Z|S, V, U) = \prod_{k=1}^K p(z_k|S, V_k, U_k) \\ = \prod_{k=1}^K \mathcal{N}_c(z_k|Y_k Y^{-1}x, Y_k - Y_k Y^{-1}Y_k) \quad (32)$$

この結果から潜在変数 Z に関する期待値を計算しておく.

$$\mathbb{E}[z_k|X, V, U] = Y_k Y^{-1}x \quad (33)$$

$$\mathbb{E}[Z_k|X, V, U] = Y_k - Y_k Y^{-1}Y_k + Y_k Y^{-1}X Y^{-1}Y_k \\ = Y_k (I + (Y^{-1}X - I) Y^{-1}Y_k) \quad (34)$$

次に, M ステップにおいて, $q(Z)$ が既知のもとで, $\mathcal{L}(q(Z), V, U)$ を最大化する V および U を求める. まず, $\mathcal{L}(q(Z), V, U)$ のうち, V および U を含む項を書き下す.

$$\mathcal{L}(q(Z), V, U) \\ \stackrel{c}{=} \int q(Z) \log p(X|Z) p(Z|V, U) dZ \\ \stackrel{c}{=} \sum_{k=1}^K \mathbb{E}_{q(z_k)} [\log \mathcal{N}_c(z_k|0, V_k \otimes U_k)] \\ \stackrel{c}{=} -T \sum_{k=1}^K \log |V_k| - F \sum_{k=1}^K \log |U_k| \\ - \sum_{k=1}^K \text{tr}((V_k^{-1} \otimes U_k^{-1}) \mathbb{E}_{q(z_k)}[Z_k]) \quad (35)$$

$\mathcal{L}(q(Z), V, U)$ を V_k^{-1} で偏微分する.

$$\frac{\partial \mathcal{L}}{\partial V_k^{-1}} = T V_k^T - (I_{F,F} \otimes \mathbf{1}_T^T) \\ \left((I_{F,F} \otimes U_k^{-1}) \odot \mathbb{E}_{q(z_k)}[Z_k^T] \right) (I_{F,F} \otimes \mathbf{1}_T) \quad (36)$$

ここで, $\mathbf{1}_*$ は全要素が 1 のベクトルあるいは行列を表す (添え字はデータのサイズ). これを 0 とおいて V_k について解き, 式 (34) を代入すると V_k の更新則を得る.

$$V_k \leftarrow \frac{1}{T} (I_{F,F} \otimes \mathbf{1}_T^T) \left((I_{F,F} \otimes U_k^{-T}) \odot (Y_k (I + (Y^{-1}X - I) Y^{-1}Y_k)) \right) (I_{F,F} \otimes \mathbf{1}_T) \quad (37)$$

同様に, U_k の更新則も求められる.

$$U_k \leftarrow \frac{1}{F} (\mathbf{1}_F^T \otimes I_{T,T}) \left((V_k^{-T} \otimes \mathbf{1}_{T,T}) \odot (Y_k (I + (Y^{-1}X - I) Y^{-1}Y_k)) \right) (\mathbf{1}_F \otimes I_{T,T}) \quad (38)$$

EM アルゴリズムでは, 対数尤度関数 $p(S|Y)$ が収束するまで式 (37) および式 (38) を反復する. このとき, V_k と U_k にはスケール不定性があるので, 反復のたびに, V_k の対角成分 (IS-NMF における基底スペクトル w_k に相当) の和が 1 となるように, 両者のスケールの調整を行う. この処理で, $\mathcal{L}(q(Z), V, U)$ の値は影響を受けない.

3.5 Minorization-Maximization アルゴリズム

対数尤度関数 $p(S|Y)$ に対して EM アルゴリズムとは異なる補助関数 (下限関数) を設計することにより, MM アルゴリズムに基づく最尤推定を行うことができる. EM アルゴリズムは, 確率モデルの最尤推定にのみ適用できるが, MM アルゴリズムは, 補助関数さえ設計できれば, 一般的な目的関数の最適化問題に適用できる.

まず, 半正定値行列を入力にとる関数の凸性および凹性に基づく不等式を導出しておく. まず, $f(Z) = \log |Z|$ ($Z \in S_+^M$) が凹関数であることに着目すると, $f(Z)$ に対して 1 次のテイラー展開を行うことで, 次式を得る.

$$\log |Z| \leq \log |\Omega| + \text{tr}(\Omega^{-1}Z) - M \quad (39)$$

ここで, Ω は任意の半正定値行列である. 等号成立条件は, $\Omega = Z$ で与えられる. 次に, 任意の半正定値行列 $A \in S_+^M$ に対して $g(Z) = \text{tr}(Z^{-1}A)$ は凸関数であることに着目すると, 澤田らの提案する不等式 [20] を適用可能である.

$$\text{tr} \left(\left(\sum_{k=1}^K Z_k \right)^{-1} A \right) \leq \sum_{k=1}^K \text{tr} \left(Z_k^{-1} \Phi_k A \Phi_k^H \right) \quad (40)$$

ここで, $\{Z_k \in S_+^M\}_{k=1}^K$ は任意の半正定値行列の集合であり, $\{\Phi_k \in \mathbb{C}^{M \times M}\}_{k=1}^K$ は $\sum_k \Phi_k = I_{M,M}$ を満たす補助変数である. ここで, $I_{M,M}$ は $M \times M$ の単位行列である. 等号成立条件は, $\Phi_k = Z_k (\sum_{k'} Z_{k'})^{-1}$ で与えられる.

式 (39) および式 (40) を用いると, 対数尤度関数 $p(S|Y)$ に対する下限関数を導出できる.

$$\log p(S|Y) \\ \stackrel{c}{=} -\log |Y| - \text{tr}(XY^{-1}) \\ \stackrel{c}{\geq} -\log |\Omega| - \sum_{k=1}^K \text{tr}(Y_k \Omega^{-1}) - \sum_{k=1}^K \text{tr} \left(Y_k^{-1} \Phi_k X \Phi_k^H \right) \\ \stackrel{\text{def}}{=} \mathcal{L}(\Omega, \Phi, V, U) \quad (41)$$

ここで, Ω は半正定値行列であり, $\{\Phi_k\}_{k=1}^K$ は $\sum_k \Phi_k = I_{FT,FT}$ を満たす補助変数である. 等号が成立する, すなわち, $\mathcal{L}(\Omega, \Phi, V, U)$ を最大化するときの条件 (補助パラメータの更新則) は次式で与えられる.

$$\Omega = Y \quad (42)$$

$$\Phi_k = Y_k Y^{-1} \quad (43)$$

次に, Ω および Φ が既知のもとで, $\mathcal{L}(\Omega, \Phi, V, U)$ を最大化する V および U を求める. 行列変数微分のチェーン則に注意しながら, $\mathcal{L}(\Omega, \Phi, V, U)$ を V_k で偏微分する.

$$\frac{\partial \mathcal{L}}{\partial \text{vec}^T(V_k)} \\ = -\text{vec}^T \left((I_{F,F} \otimes \mathbf{1}_T^T) \left((I_{F,F} \otimes U_k) \odot \Omega^{-T} \right) (I_{F,F} \otimes \mathbf{1}_T) \right) \\ + \text{vec}^T \left((I_{F,F} \otimes \mathbf{1}_T^T) \left((I_{F,F} \otimes U_k^{-1}) \odot \Phi_k^C X^T \Phi_k^T \right) (I_{F,F} \otimes \mathbf{1}_T) \right) \left(V_k^{-T} \otimes V_k^{-1} \right) \quad (44)$$

$$\begin{aligned}
&= -\text{vec}^T \left((\mathbf{I}_{F,F} \otimes \mathbf{1}_T^T) \left((\mathbf{1}_{F,F} \otimes \mathbf{U}_k) \odot \boldsymbol{\Omega}^{-T} \right) (\mathbf{I}_{F,F} \otimes \mathbf{1}_T) \right) \\
&\quad + \text{vec}^T \left(\mathbf{V}_k^{-T} (\mathbf{I}_{F,F} \otimes \mathbf{1}_T^T) \left((\mathbf{1}_{F,F} \otimes \mathbf{U}_k^{-1}) \odot \boldsymbol{\Phi}_k^C \mathbf{X}^T \boldsymbol{\Phi}_k^T \right) \right. \\
&\quad \left. (\mathbf{I}_{F,F} \otimes \mathbf{1}_T) \mathbf{V}_k^{-T} \right) \quad (45)
\end{aligned}$$

これを0とおいて \mathbf{V}_k について解き, 式 (42) および式 (43) を代入すると \mathbf{V}_k の更新則を得る.

$$\begin{aligned}
\mathbf{P}_k &\stackrel{\text{def}}{=} (\mathbf{I}_{F,F} \otimes \mathbf{1}_T^T) \left((\mathbf{1}_{F,F} \otimes \mathbf{U}_k^T) \odot \mathbf{Y}^{-1} \right) \\
&\quad (\mathbf{I}_{F,F} \otimes \mathbf{1}_T) \quad (46)
\end{aligned}$$

$$\begin{aligned}
\mathbf{Q}_k &\stackrel{\text{def}}{=} (\mathbf{I}_{F,F} \otimes \mathbf{1}_T^T) \left((\mathbf{1}_{F,F} \otimes \mathbf{U}_k^T) \odot \mathbf{Y}^{-1} \mathbf{X} \mathbf{Y}^{-1} \right) \\
&\quad (\mathbf{I}_{F,F} \otimes \mathbf{1}_T) \quad (47)
\end{aligned}$$

$$\begin{aligned}
\mathbf{V}_k &\leftarrow \mathbf{P}_k^{-\frac{1}{2}} \left(\mathbf{P}_k^{\frac{1}{2}} \mathbf{V}_k \mathbf{Q}_k \mathbf{V}_k \mathbf{P}_k^{\frac{1}{2}} \right)^{\frac{1}{2}} \mathbf{P}_k^{-\frac{1}{2}} \\
&= \mathbf{P}_k^{-1} \circ (\mathbf{V}_k \mathbf{Q}_k \mathbf{V}_k) \quad (48)
\end{aligned}$$

ここで, \circ は半正定値行列同士の幾何平均を表す [21]. 同様に, \mathbf{U}_k の更新則も求められる.

$$\begin{aligned}
\mathbf{R}_k &\stackrel{\text{def}}{=} (\mathbf{1}_F^T \otimes \mathbf{I}_{T,T}) \left((\mathbf{V}_k^T \otimes \mathbf{1}_{T,T}) \odot \mathbf{Y}^{-1} \right) \\
&\quad (\mathbf{1}_F \otimes \mathbf{I}_{T,T}) \quad (49)
\end{aligned}$$

$$\begin{aligned}
\mathbf{S}_k &\stackrel{\text{def}}{=} (\mathbf{1}_F^T \otimes \mathbf{I}_{T,T}) \left((\mathbf{V}_k^T \otimes \mathbf{1}_{T,T}) \odot \mathbf{Y}^{-1} \mathbf{X} \mathbf{Y}^{-1} \right) \\
&\quad (\mathbf{1}_F \otimes \mathbf{I}_{T,T}) \quad (50)
\end{aligned}$$

$$\begin{aligned}
\mathbf{U}_k &\leftarrow \mathbf{R}_k^{-\frac{1}{2}} \left(\mathbf{R}_k^{\frac{1}{2}} \mathbf{U}_k \mathbf{S}_k \mathbf{U}_k \mathbf{R}_k^{\frac{1}{2}} \right)^{\frac{1}{2}} \mathbf{R}_k^{-\frac{1}{2}} \\
&= \mathbf{R}_k^{-1} \circ (\mathbf{U}_k \mathbf{S}_k \mathbf{U}_k) \quad (51)
\end{aligned}$$

3.6 ノンパラメトリックベイズモデル

本節では, 基底数を $K \rightarrow \infty$ としたノンパラメトリックベイズモデルについて説明する. 無限モデルにおいては, 無限個の基底のうちで, 観測データを表現するのに実質的には少数の基底しか利用されないようなスパースな学習を行う必要がある. いま, 基底の重みを表す無限次元の非負値ベクトル $\boldsymbol{\theta} = [\theta_1, \dots, \theta_\infty]^T$ を式 (25) に導入する.

$$\mathbf{s} | \mathbf{Y} \sim \mathcal{N}_c(\mathbf{0}, \mathbf{Y}) = \mathcal{N}_c \left(\mathbf{0}, \sum_{k=1}^K \theta_k (\mathbf{V}_k \otimes \mathbf{U}_k) \right) \quad (52)$$

ここで, $\mathbf{Y}_k = \theta_k (\mathbf{V}_k \otimes \mathbf{U}_k)$ かつ $\mathbf{Y}_k = \sum_k \mathbf{Y}_k$ とした. このとき, 無限次元の $\boldsymbol{\theta}$ のうちで, 一部の要素のみが有意に大きな値をもち, それ以外はほとんどゼロとなるような学習を行うには, $\boldsymbol{\theta}$ に対する事前分布としてガンマ過程を用いるのが自然である [17]. 最も簡単なガンマ過程の近似方法は, K を十分に大きくとり, 各要素 θ_k がガンマ事前分布に従うことを仮定することである.

$$\theta_k \sim \mathcal{G}(\alpha c / K, \alpha) \quad (53)$$

ここで, $\alpha \geq 0$ および $c \geq 0$ は正の超パラメータであり, $\mathbb{E}_{\text{prior}}[\theta_k] = c/K$ および $\mathbb{E}_{\text{prior}}[\sum_k \theta_k] = c$ となっている.

このとき, $K \rightarrow \infty$ とすれば, ガンマ過程が得られる.

$$G \sim \text{GaP}(\alpha, G_0) \quad (54)$$

ここで, G_0 はある空間 Θ 上に定義された基底測度であり, $G(\Theta) = c$ を満たす (確率測度に限らない). サンプルされる G は Θ 上の離散測度となることが知られており, $\mathbb{E}[G] = G_0$ となっている. 空間 Θ の微小区間への分割を $\{\Theta_1, \dots, \Theta_\infty\}$ とすると, $G(\Theta_k) = \theta_k$ を満たす. α は集中度と呼ばれ, α が小さくなるほど $\boldsymbol{\theta}$ はスパースになる.

また, パラメータ \mathbf{U} および \mathbf{V} に関しては, 複素ウィシャート事前分布を用いるのが都合がよい.

$$\mathbf{V}_k \sim \mathcal{W}_c(\mathbf{R}_0^V, \gamma_0^V) \quad (55)$$

$$\mathbf{U}_k \sim \mathcal{W}_c(\mathbf{R}_0^U, \gamma_0^U) \quad (56)$$

いま, 観測データ \mathbf{S} が与えられたもとの, パラメータの事後分布 $p(\boldsymbol{\theta}, \mathbf{V}, \mathbf{U} | \mathbf{S}) = p(\mathbf{S}, \boldsymbol{\theta}, \mathbf{V}, \mathbf{U}) / p(\mathbf{S})$ を計算したい. しかし, 周辺尤度 $p(\mathbf{S}) = \iiint p(\mathbf{S}, \boldsymbol{\theta}, \mathbf{V}, \mathbf{U}) d\boldsymbol{\theta} d\mathbf{V} d\mathbf{U}$ の計算は解析的に行えないため, 変分ベイズ法を用いて $p(\boldsymbol{\theta}, \mathbf{V}, \mathbf{U} | \mathbf{S})$ を近似的に求める. まず, 因子分解が可能な関数形をもつ変分事後分布 $q(\boldsymbol{\theta}, \mathbf{V}, \mathbf{U})$ を考える.

$$q(\boldsymbol{\theta}, \mathbf{V}, \mathbf{U}) = \prod_{k=1}^K q(\theta_k) q(\mathbf{V}_k) q(\mathbf{U}_k) \quad (57)$$

そのうえで, 真の事後分布 $p(\boldsymbol{\theta}, \mathbf{V}, \mathbf{U} | \mathbf{S})$ に対する KL ダイバージェンスを最小化するような $q(\boldsymbol{\theta}, \mathbf{V}, \mathbf{U})$ を求めたい. これは, 対数周辺尤度 $\log p(\mathbf{S})$ の変分下限 \mathcal{L} を最大化することと等価であることが知られている.

$$\begin{aligned}
&\log p(\mathbf{S}) \\
&\geq \mathbb{E}[\log p(\mathbf{S} | \boldsymbol{\theta}, \mathbf{V}, \mathbf{U})] \\
&\quad + \mathbb{E}[\log p(\boldsymbol{\theta})] + \mathbb{E}[\log p(\mathbf{V})] + \mathbb{E}[\log p(\mathbf{U})] \\
&\quad - \mathbb{E}[\log q(\boldsymbol{\theta})] - \mathbb{E}[\log q(\mathbf{V})] - \mathbb{E}[\log q(\mathbf{U})] \\
&\stackrel{c}{\geq} \mathbb{E}[\log \mathcal{L}(\boldsymbol{\Omega}, \boldsymbol{\Phi}, \boldsymbol{\theta}, \mathbf{V}, \mathbf{U})] \\
&\quad + \mathbb{E}[\log p(\boldsymbol{\theta})] + \mathbb{E}[\log p(\mathbf{V})] + \mathbb{E}[\log p(\mathbf{U})] \\
&\quad - \mathbb{E}[\log q(\boldsymbol{\theta})] - \mathbb{E}[\log q(\mathbf{V})] - \mathbb{E}[\log q(\mathbf{U})] \\
&\stackrel{\text{def}}{=} \mathcal{L}(\boldsymbol{\Omega}, \boldsymbol{\Phi}, q(\boldsymbol{\theta}), q(\mathbf{V}), q(\mathbf{U})) \quad (58)
\end{aligned}$$

ここで, 第一項に関して, 不等式 (41) では, $\mathbf{Y}_k = \mathbf{V}_k \otimes \mathbf{U}_k$ であったのを, $\mathbf{Y}_k = \theta_k (\mathbf{V}_k \otimes \mathbf{U}_k)$ とした不等式を用いた. このとき, 変分下限 $\mathcal{L}(\boldsymbol{\Omega}, \boldsymbol{\Phi}, q(\boldsymbol{\theta}), q(\mathbf{V}), q(\mathbf{U}))$ を最大化するには, $\mathcal{L}(\boldsymbol{\Omega}, \boldsymbol{\Phi}, q(\boldsymbol{\theta}), q(\mathbf{V}), q(\mathbf{U}))$ が収束するまで各変数の最適化を反復する. まず, 補助パラメータ $\boldsymbol{\Omega}$ および $\boldsymbol{\Phi}$ に関する最適化を行う. その後, 各パラメータに関する変分事後分布を順番に更新すればよい.

$$q(\boldsymbol{\theta}) \propto p(\boldsymbol{\theta}) \exp(\mathbb{E}_{q(\mathbf{V}, \mathbf{U})}[\log \mathcal{L}(\boldsymbol{\Omega}, \boldsymbol{\Phi}, \boldsymbol{\theta}, \mathbf{V}, \mathbf{U})]) \quad (59)$$

$$q(\mathbf{V}) \propto p(\mathbf{V}) \exp(\mathbb{E}_{q(\boldsymbol{\theta}, \mathbf{U})}[\log \mathcal{L}(\boldsymbol{\Omega}, \boldsymbol{\Phi}, \boldsymbol{\theta}, \mathbf{V}, \mathbf{U})]) \quad (60)$$

$$q(\mathbf{U}) \propto p(\mathbf{U}) \exp(\mathbb{E}_{q(\boldsymbol{\theta}, \mathbf{V})}[\log \mathcal{L}(\boldsymbol{\Omega}, \boldsymbol{\Phi}, \boldsymbol{\theta}, \mathbf{V}, \mathbf{U})]) \quad (61)$$

まず, $\mathcal{L}(\Omega, \Phi, q(\theta), q(V), q(U))$ を最大化する補助パラメータ Ω および Φ を求める.

$$\Omega = \sum_k \mathbb{E}[Y_k] \quad (62)$$

$$\Phi_k = (\mathbb{E}[Y_k^{-1}])^{-1} \left(\sum_{k'} (\mathbb{E}[Y_{k'}^{-1}])^{-1} \right)^{-1} \quad (63)$$

次に, $\mathcal{L}(\Omega, \Phi, q(\theta), q(V), q(U))$ を最大化する変分事後分布 $q(\theta), q(V), q(U)$ を求める. 各分布は一般化逆ガウス (Generalized Inverse Gaussian, GIG) 分布あるいは複素行列 GIG (Matrix GIG, MGIG) 分布で与えられる [11, 22].

$$\begin{aligned} q(\theta_k) &= \text{GIG}(\theta_k | \alpha c / K, \rho_k^\theta, \tau_k^\theta) \\ q(V_k) &= \text{MGIG}(V_k | \gamma_0^V, R_k^V, T_k^V) \\ q(U_k) &= \text{MGIG}(U_k | \gamma_0^U, R_k^U, T_k^U) \end{aligned} \quad (64)$$

ここで, 各分布のパラメータは以下で与えられる.

$$\rho_k^\theta = 2\alpha + 2\text{tr}(\mathbb{E}[V_k \otimes U_k] \Omega^{-1}) \quad (65)$$

$$\tau_k^\theta = 2\text{tr}(\mathbb{E}[V_k^{-1} \otimes U_k^{-1}] \Phi_k X \Phi_k^H) \quad (66)$$

$$\begin{aligned} R_k^V &= (R_0^V)^{-1} + \mathbb{E}[\theta_k] (I_{F,F} \otimes \mathbf{1}_T^T) \\ &\quad \left((\mathbf{1}_{F,F} \otimes \mathbb{E}[U_k^T]) \odot \Omega^{-1} \right) (I_{F,F} \otimes \mathbf{1}_T) \end{aligned} \quad (67)$$

$$\begin{aligned} T_k^V &= \mathbb{E}[\theta_k^{-1}] (I_{F,F} \otimes \mathbf{1}_T^T) \\ &\quad \left((\mathbf{1}_{F,F} \otimes \mathbb{E}[U_k^{-T}]) \odot \Phi_k X \Phi_k^H \right) (I_{F,F} \otimes \mathbf{1}_T) \end{aligned} \quad (68)$$

$$\begin{aligned} R_k^U &= (R_0^U)^{-1} + \mathbb{E}[\theta_k] (\mathbf{1}_F^T \otimes I_{T,T}) \\ &\quad \left((\mathbb{E}[V_k^T] \otimes \mathbf{1}_{T,T}) \odot \Omega^{-1} \right) (\mathbf{1}_F \otimes I_{T,T}) \end{aligned} \quad (69)$$

$$\begin{aligned} T_k^U &= \mathbb{E}[\theta_k^{-1}] (\mathbf{1}_F^T \otimes I_{T,T}) \\ &\quad \left((\mathbb{E}[V_k^{-T}] \otimes \mathbf{1}_{T,T}) \odot \Phi_k X \Phi_k^H \right) (\mathbf{1}_F \otimes I_{T,T}) \end{aligned} \quad (70)$$

MGIG 分布に従う確率変数の各種期待値を計算するには, ウィシャート分布を用いた重点サンプリングを行う必要がある [23].

3.7 計算量削減

これまで説明した通り, 理論上は LD-CTF の最尤推定とベイズ推定が導出できるが, 計算量が極めて膨大であり, 現実的には実行不可能である. 音源分離問題において, F と T はそれぞれ周波数ビン数とフレーム数に対応しており, $FT \times FT$ という巨大な行列を取り扱う必要があることから, LD-CTF の計算量は $O(KF^3T^3)$ である.

計算量を削減する一つの方法は, 半正定値行列 V_k や U_k をブロック対角行列に制限することである (図 2). 音源分離においては, 時間周波数領域を矩形領域に分割し, 各ブロックでは完全な相関を考慮するが, ブロック間は独立であるとみなすことを意味する. ここで, V_k と U_k がどちらも対角行列の場合には IS-NMF, いずれかが対角行列の場合には LD-PSDTF となることに注意されたい. ブロック対角 LD-CTF は, 制約のない LD-CTF と LD-PSDTF の間の中間的な表現力のモデルになっている.

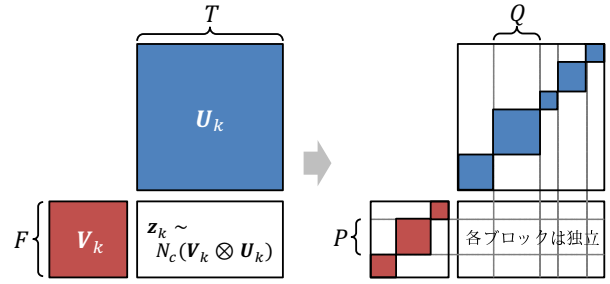


図 2 ブロック対角行列に制限した LD-CTF.

簡単のため, $V_k \in S_+^F$ は, 同一サイズ P の半正定値行列を対角上に I 個に並べた半正定値行列であり, $U_k \in S_+^T$ は, 同一サイズ Q の半正定値行列を対角上に J 個に並べた半正定値行列であるとする. ただし, $F = PI$ および $T = QJ$ とする. 実際には, すべての k について共通であれば, 各ブロックのサイズは可変でよい. また, 必ずしも隣接する行あるいは列を同じブロックにいれる必要はなく, 行および列をそれぞれ I 個と J 個のブロックにクラスタリングしておけばよい. このとき, 計算量は $O(KIJP^3Q^3)$ となり, 大幅な計算量の削減が可能となる.

優れた近似精度と劇的な高速化を同時に達成するための方法として, V および U に対してそれぞれ同時対角化を行うことが考えられる. 本稿では, 時間周波数領域における複素スペクトログラム S の確率モデルを考えていた.

$$S|V, U \sim \mathcal{N}_c \left(\mathbf{0}, \sum_{k=1}^K V_k \otimes U_k \right) \quad (71)$$

ここで, S は行列であるが, 全ての要素を並べたベクトルが多変量複素ガウス分布に従うと考える. S に任意の線形変換 A および B を両側から適用すると, 次式が成立する.

$$ASB^H|V, U \sim \mathcal{N}_c \left(\mathbf{0}, \sum_{k=1}^K AV_k A^H \otimes BU_k B^H \right)$$

このとき, すべての k について, $AV_k A^H$ および $BU_k B^H$ が対角行列となれば, 変換後の空間では, LD-CTF は IS-NMF に帰着し, 高速な実行が可能となる. したがって, 同時対角化を行う A および B の推定と, V および U の更新を反復的に繰り返す手法が有望である. 図 3 に, A と B が, ともに離散フーリエ変換 (DFT) 行列である場合の例を示す. もし, 時間領域信号が定常であれば, V_k は巡回行列となり, DFT 行列 D_F で対角されるため, $A = D_F$ とすることはよい選択となる. これが, 周波数方向には独立であるとして IS-NMF を適用する理由であったが, LD-PSDTF を用いれば分離精度が向上することから, D_F は同時対角化のための最適な変換ではないことが示唆されている. 一方, フレーム方向の時間軸の変換については, これまでほとんど議論されてこなかった. ガウス過程の観点から, 周波数方向・時間方向の依存性を一挙に解決する試みについてさらなる研究が期待される.

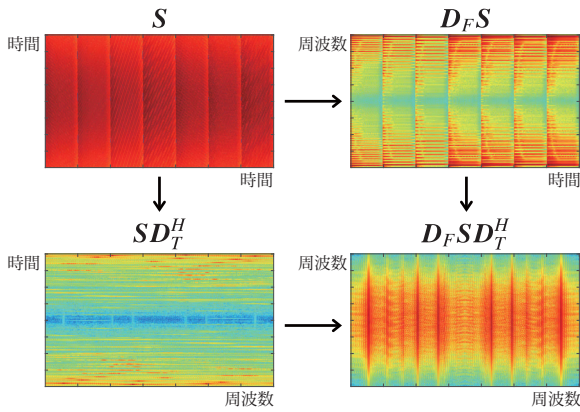


図 3 DFT による空間変換. $S \in \mathbb{C}^{F \times T}$ は混合音の複素スペクトログラム (ここでは F は窓幅) であり, $D_F \in \mathbb{C}^{F \times F}$ および $D_T \in \mathbb{C}^{T \times T}$ はそれぞれ離散フーリエ変換行列である.

4. 評価

LD-CTF を用いた音源分離実験について報告する. 本稿では, 3.5 節で述べた MM アルゴリズムに基づく最尤推定を用いて, LD-CTF の基本的な動作と性能を確認した.

4.1 実験条件

実験には, MIDI のピアノ音を用いた. 三つの音高 (C4, E4, G4) をもつ 1.2 秒間の音響信号を準備し, それらを 7 つの異なる組み合わせで重畳したもの (C4, E4, G4, C4+E4, C4+G4, E4+G4, C4+E4+G4) を連結して 8.4 秒の音響信号を合成した. サンプル周波数は 16[kHz] とした. 窓幅 512 点ガウス窓を用いて, 窓シフト長 160 点の STFT を行った ($F = 256, T = 840$).

合成した混合音を C4, E4, G4 に対応する音源信号に分離することを試みた ($K = 3$). ブロック対角 LD-CTF の設定は, $(P, Q) = (256, 1), (1, 840), (128, 10), (64, 20), (32, 40)$ とした. 比較のため, IS-NMF も評価した. このとき, 計算量を削減し, 局所解を回避するため, IS-NMF の結果を, LD-CTF の初期値とした. ブロック対角 LD-CTF は, $(P, Q) = (1, 1)$ のとき IS-NMF に, $(P, Q) = (256, 1)$ あるいは $(P, Q) = (1, 840)$ のときそれぞれ周波数方向あるいは時間方向の相関を考慮する LD-PSDTF と等価である. 各手法に対して, 反復回数は 100 回とした. BSS Eval Toolbox [24] を用いて, Source-to-Distortion Ratio (SDR), Source-to-Interferences Ratio (SIR) および Sources-to-Artifacts Ratio (SAR) で評価した.

4.2 実験結果

表 1 に実験結果を示す. $(P, Q) = (256, 1)$ に対応する LD-CTF (周波数方向の相関を考慮した LD-PSDTF) は IS-NMF に対する明確な優位性を示した. LD-PSDTF を時間領域で実行した場合は, SDR は 26.6 [dB] であった [12]. 理論的には両者は等価であり, ほぼ同じ分離精度が得られ

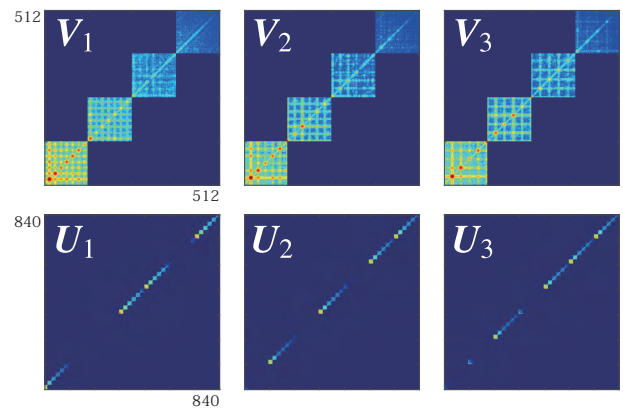


図 4 $(P, Q) = (64, 20)$ のときの LD-CTF の結果.

表 1 音源分離精度 [dB].

(P, Q)	IS-NMF	LD-PSDTF		LD-CTF		
	(1, 1)	(256, 1)	(840, 1)	(128, 10)	(64, 20)	(32, 40)
SDR	18.88	21.58	21.04	19.68	20.60	20.21
SIR	24.14	27.01	24.67	25.29	26.17	25.45
SAR	20.45	23.14	23.50	21.47	21.47	22.15

るはずであるが, 複素行列計算に関する計算誤差の影響が考えられる. 一方, $(P, Q) = (1, 840)$ に対応する LD-CTF (時間方向の相関を考慮した LD-PSDTF) でも優れた分離結果が得られることが本実験で初めて確認された. 図 1 を見ると, 周波数方向においては, 倍音の間には強い相関があり, 時間方向においては, 音量の大きいフレームの間には強い相関があることがわかる.

ブロック対角 LD-CTF は, IS-NMF より高精度な分離が可能であるはあるものの, LD-PSDTF には及ばなかった. 図 4 に実行結果の例を示す. 時間周波数領域を排他的なブロックに区切ることで, ブロックの境界付近で位相が不適切になる可能性が考えられる. また, 隣接する周波数を同じブロックにしているが, 倍音周波数を同じブロックに入れて, 倍音周波数間の相関は完全に扱うことで, 分離精度は向上することが期待できる.

5. おわりに

本稿では, NMF や PSDTF を特別な場合として含む, 相関テンソル分解 (CTF) と呼ぶ究極の因子分解法を提案した. 特に, 音源分離に適した LD ダイバージェンスに基づく CTF (LD-CTF) を取り上げ, 最尤推定のための EM アルゴリズムおよび MM アルゴリズム, さらにガンマ過程に基づくノンパラメトリックベイズモデルと変分ベイズ法を提案した. また, ブロック対角制約や同時対角化などの計算量を削減する方法について議論した.

CTF は一般の N 階のテンソルデータに対して適用でき, よく知られた Canonical Polyadic (CP) 分解のエレガントな数学的拡張でもある. テンソルの各モードにおける完全な相関を取り扱うことができる強力な枠組みであり, 音響信号処理分野以外への応用について検討していきたい.

謝辞: 本研究の一部は, JSPS 科研費 26700020, 16H01744 および JST ACCEL No. JPMJAC1602 の支援を受けた。

参考文献

- [1] Itoyama, K., Goto, M., Komatani, K., Ogata, T. and Okuno, H. G.: Query-by-Example Music Information Retrieval by Score-Informed Source Separation and Remixing Technologies, *EURASIP Journal on Advances in Signal Processing*, Vol. 2010 (2010). Article ID 172961.
- [2] Goto, M.: Active Music Listening Interfaces based on Signal Processing, *IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, Vol. 4, pp. 1441–1444 (2007).
- [3] Yoshii, K., Goto, M., Komatani, K., Ogata, T. and Okuno, H. G.: Drumix: An Audio Player with Real-time Drum-Part Rearrangement Functions for Active Music Listening, *IPSSJ Digital Courier*, Vol. 3, pp. 134–144 (2007).
- [4] Sturmel, N., Liutkus, A., Pinel, J., Girin, L., Marchand, S., Richard, G., Badeau, R. and Daudet, L.: Linear Mixing Models for Active Listening of Music Productions in Realistic Studio Conditions, *The 132nd Audio Engineering Society (AES) Convention* (2012).
- [5] Lee, D. and Seung, H.: Algorithms for Non-Negative Matrix Factorization, *Neural Information Processing Systems (NIPS)*, pp. 556–562 (2000).
- [6] Griffin, D. W. and Lim, J. S.: Signal Estimation from Modified Short-Time Fourier Transform, *IEEE Transactions on Acoustics, Speech, and Signal Processing*, Vol. 32, No. 2, pp. 236–243 (1984).
- [7] Roux, J. L., Kameoka, H., Ono, N. and Sagayama, S.: Explicit Consistency Constraints for STFT Spectrograms and Their Application to Phase Reconstruction, *Workshop on Statistical and Perceptual Audition (SAPA)*, pp. 23–28 (2008).
- [8] Roux, J. L., Kameoka, H., Ono, N. and Sagayama, S.: Fast Signal Reconstruction from Magnitude STFT Spectrogram Based on Spectrogram Consistency, *International Conference on Digital Audio Effects (DAFx)*, pp. 397–403 (2010).
- [9] Kameoka, H., Nishimoto, T. and Sagayama, S.: Complex NMF: A New Sparse Representation for Acoustic Signals, *IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pp. 45–48 (2009).
- [10] Roux, J. L., Vincent, E., Mizuno, Y., Kameoka, H., Ono, N. and Sagayama, S.: Consistent Wiener Filtering: Generalized Time-Frequency Masking Respecting Spectrogram Consistency, *International Conference on Latent Variable Analysis and Signal Separation (LVA/ICA)*, pp. 89–96 (2010).
- [11] Yoshii, K., Tomioka, R., Mochihashi, D. and Goto, M.: Infinite Positive Semidefinite Tensor Factorization for Source Separation of Mixture Signals, *International Conference on Machine Learning (ICML)*, pp. 576–584 (2013).
- [12] Yoshii, K., Tomioka, R., Mochihashi, D. and Goto, M.: Beyond NMF: Time-Domain Audio Source Separation without Phase Reconstruction, *International Society for Music Information Retrieval Conference (ISMIR)*, pp. 369–374 (2013).
- [13] Févotte, C., Bertin, N. and Durrieu, J.-L.: Nonnegative Matrix Factorization with the Itakura-Saito Divergence: With Application to Music Analysis, *Neural Computation*, Vol. 21, No. 3, pp. 793–830 (2009).
- [14] Bregman, L. M.: The Relaxation Method of Finding the Common Points of Convex Sets and Its Application to the Solution of Problems in Convex Programming, *USSR Computational Mathematics and Mathematical Physics*, Vol. 7, No. 3, pp. 200–217 (1967).
- [15] Nakano, M., Kameoka, H., Roux, J. L., Kitano, Y., Ono, N. and Sagayama, S.: Convergence-Guaranteed Multiplicative Algorithms for Non-Negative Matrix Factorization with Beta Divergence, *International Workshop on Machine Learning for Signal Processing (MLSP)*, pp. 283–288 (2010).
- [16] Cemgil, A. T.: Bayesian Inference for Nonnegative Matrix Factorisation Models, *Computational Intelligence and Neuroscience*, Vol. 2009, pp. 1–17 (2009).
- [17] Hoffman, M., Blei, D. and Cook, P.: Bayesian Nonparametric Matrix Factorization for Recorded Music, *International Conference on Machine Learning (ICML)*, pp. 439–446 (2010).
- [18] Allen, J. B. and Rabiner, L. R.: A Unified Approach to Short-Time Fourier Analysis and Synthesis, *IEEE*, Vol. 65, No. 11, pp. 1558–1564 (1977).
- [19] Kulis, B., Sustik, M. and Dhillon, I.: Low-rank Kernel Learning with Bregman Matrix Divergences, *Journal of Machine Learning Research (JMLR)*, Vol. 10, pp. 341–376 (2009).
- [20] Sawada, H., Kameoka, H., Araki, S. and Ueda, N.: Efficient Algorithms for Multichannel Extensions of Itakura-Saito Nonnegative Matrix Factorization, *International Conference on Acoustics, Speech, and Signal Processing (ICASSP)*, pp. 261–264 (2012).
- [21] Pusz, W. and Woronowicz, S. L.: Functional Calculus for Sesquilinear Forms and the Purification Map, *Reports on Mathematical Physics*, Vol. 8, No. 2, pp. 159–170 (1975).
- [22] Itakura, K., Bando, Y., Nakamura, E., Itoyama, K., Yoshii, K. and Kawahara, T.: Bayesian Multichannel Nonnegative Matrix Factorization for Audio Source Separation and Localization, *International Conference on Acoustics, Speech, and Signal Processing (ICASSP)*, pp. 551–555 (2017).
- [23] Yang, M., Li, Y. and Zhang, Z.: Multi-Task Learning with Gaussian Matrix Generalized Inverse Gaussian Model, *International Conference on Machine Learning (ICML)*, pp. 423–431 (2013).
- [24] Vincent, E., Gribonval, R. and Févotte, C.: Performance Measurement in Blind Audio Source Separation, *IEEE Transactions on Acoustics, Speech, and Signal Processing*, Vol. 14, No. 4, pp. 1462–1469 (2006).