# 周波数・時間・到来方向のスパースネスを統合した

# 近接分離最適化による IVA \*

☆牛島隆裕 (大分大), 太刀岡勇気 (デンソーIT ラボ), 上ノ原進吾, 古家賢一 (大分大)

# 1 はじめに

音源分離とは、複数の音源が混合されて観測され た信号から、単一の音源を含む信号を復元する技術で ある。音源分離が適用される場面として、雑音に頑健 な音声認識を実現するために使用者の音声のみを取 り出して雑音を抑圧する音声強調が例に挙げられる。

音源分離に用いられる手法は様々存在しており、例 として独立成分分析 [1]、独立ベクトル分析 (IVA)[2]、 同期同時対角化 [3]、独立低ランク行列分析 [4]、マル チチャネル非負値行列因子分解 [5] が存在している。 これらの手法に用いる最適化手法として補助関数法 による最適化が提案されており、高速でよい解に収束 することが挙げられている。

近年、補助関数法以外の最適化手法を用いた音源分 離手法として、近接分離最適化によるブラインド音源 分離 (BSS) が提案されている [6]。近接分離最適化に よる BSS は音源の生成モデルを表す音源項と分離行 列のスケールを調整する項を独立に設定することが できる利点がある。音源項の近接写像を求めることが 可能ならば、容易に音源項を変更した結果を検証可 能である。[6] では例として音源項の生成モデルにラ プラス球対称確率分布を用いたラプラス IVA が提案 されている。また異なる応用として、音源分離に必須 の前処理である白色化を行った場合、スパース構造が 崩れて発生する雑音の影響を軽減するため、スパー ス構造を復元する重みを導入したスパース IVA[7] が 提案されている。従来法の検討 [7] では音声データに 対する有効性が示されているが、音源によっては分離 性能が低下することを本稿の実験により示す。

本稿では、種々の音源での分離性能の向上のため、 従来のスパース IVA で考えられている周波数方向の スパース性の他に、時間方向と到来方向のスパース性 も考慮したスパース IVA を提案する。そしてそれら の複数のスパースネスを組み合わせることによって、 従来のスパース IVA に比べて様々な音源に対応した スパースネス統合 IVA の提案を行う。実験では従来 法で検討されてきた音声データのほかに音楽データ に対し検討を行い提案法の有効性を示す。

# 2 近接分離最適化による BSS

# 2.1 音源分離の概要

音源信号  $s_{ij}$ 、観測信号  $x_{ij}$  及び分離信号  $y_{ij}$  を短時間フーリエ変換したスペクトログラムにより表す。

$$\mathbf{s}_{ij} = (s_{ij,1}, s_{ij,2}, \dots, s_{ij,N})^T \tag{1}$$

$$\boldsymbol{x}_{ij} = (x_{ij,1}, x_{ij,2}, \dots, x_{ij,M})^T$$
(2)

$$\boldsymbol{y}_{ij} = (y_{ij,1}, y_{ij,2}, \dots, y_{ij,N})^T \tag{3}$$

ただし、各インデックスを周波数ビンi ( $1 \le i \le I$ )、 時間フレームj ( $1 \le j \le J$ )、分離対象となる音源の ID をn ( $1 \le n \le N$ )、及び観測チャネル番号をm( $1 \le m \le M$ )として表す。また、 $^{T}$  は転置を表す。 混合系の音源分離では、音源から観測マイクへの混 合の過程を表す混合行列 $A_i \in C^{M \times N}$ を用いて観測 信号、音源との関係を次式で表せる。

$$\boldsymbol{x}_{ij} = A_i \boldsymbol{s}_{ij} \tag{4}$$

M = Nかつ  $A_i$ が正則である場合、 $W_i = A_i^{-1}$ となる分離行列  $W_i \in \mathbb{C}^{N \times M}$ が存在し、観測信号から音源信号への分離を次式で表せる。

$$\boldsymbol{y}_{ij} = W_i \boldsymbol{x}_{ij} \tag{5}$$

IVA では各音源が互いに独立でありかつ何らかの確 率分布に従って生成される場合、観測信号 $x_{ij}$ から分 離行列 $W_i$ を推定できる。分離行列 $W_i$ の推定は以下 の最小化問題を解くことで達成される。

$$\underset{\{W_i\}_{i=1}^{I}}{\operatorname{Minimize}} \quad \sum_{i=1}^{I} \mathcal{P}\left(W_i \boldsymbol{x}_{ij}\right) - \sum_{i=1}^{I} \log \left|\det\left(W_i\right)\right| \quad (6)$$

第1項 $\mathcal{P}(W_i x_{ij})$ は音源の生成モデルを表す確率 分布から導かれる実数値関数であり以下音源項と呼 ぶ。第2項は分離行列のスケールを調整するための 正規化項である。

# 2.2 近接分離最適化ブラインド音源分離

近接分離最適化ブラインド音源分離は,式(6)の最 小化問題を主-双対近接分離法によって解いていく手 法である。Fig. 1に最適化の流れを示している。ま ず、式(6)を主-双対近接分離法が適用できるような 形に変形すると以下の式のようになる。

$$\underset{\mathbf{w}}{\text{Minimize}} \quad \mathcal{P}(\mathbf{X}\mathbf{w}) + \mathcal{L}(\mathbf{w}) \tag{7}$$

$$\mathcal{L}(\mathbf{w}) = -\sum_{i=1}^{I} \sum_{m=1}^{M} \log \sigma_m(\mathfrak{M}(\mathbf{w})_i)$$
(8)

行列 X は観測信号のスペクトログラム *x* を *IMJ*× *IMN* の形状に拡張したブロック対角行列である。ベク トル w は *I*×*M*×*N* の形状の分離行列 W を1×*IMN* 

<sup>\*</sup>IVA based on proximal splitting optimization integrating sparseness of frequency, time and direction of arrival by Takahiro Ushijima (Oita University), Yuuki Tachioka (Denso IT Laboratory), Shingo Uenohara, Ken'ichi Furuya (Oita University)



Fig. 1 近接分離最適化音源分離の流れ

の形にベクトル化したものである。また、 $\mathfrak{M}(\mathbf{w})$ は **w**を分離行列 W の形に戻すための処理であり、 $W_i$ と一致する。 $\sigma_m(\mathfrak{M}(\mathbf{w})_i)$ は $W_i$ の m 番目の特異値で ある。 $\mathcal{L}(\mathbf{w})$ は分離行列のスケールの正規化項を近接 写像が適用できる形に変形したものであり、これは特 異値の近接写像の計算を行うことで求めることが可 能である。 $\mathcal{P}(\mathbf{X}\mathbf{w})$ は音源の生成モデルを表す音源項 である。

2.3 スパース IVA[7]

スパース IVA は近接分離最適化の枠組みにおい て、音源項を工夫することで、白色化により失われ た周波数方向のスパース性を復元させる手法である。 これは混合  $L_{2,1}$  ノルムと  $L_1$  ノルムの和の近接写像  $\operatorname{prox}_{\lambda_1 \|\cdot\|_{2,1}+\lambda_2 \|\cdot\|_1}[\mathbf{z}]$  にスパース構造を復元する重み を適用している。以下の式で定式化され Fig. 1 中の  $\mathcal{M}[\mathbf{z}]$  に使用される。

$$(\mathcal{M}[\boldsymbol{z}])_{ijn}^{Freq} = \zeta_{ijn}^{z,\kappa} \\ \times \Xi_{\kappa} \left[ \left( 1 - \frac{\lambda_1}{\left( \Sigma_{i=1}^{I} \left( \Theta_{\eta}^{Freq}[\boldsymbol{x}] \right)_i \left| \zeta_{ijn}^{z,\kappa} z_{ijn} \right|^2 \right)^{1/2}} \right)_+ \right] (9)$$

 $\zeta_{ijn}^{z,\kappa}$  が  $L_1$  ノルムの近接写像、( $\Xi_{\kappa}[\cdot]$ ) の引数部分 が  $L_{2,1}$  ノルムの近接写像に対応している。ただし、 ( $\cdot$ )<sub>+</sub> = max(0, $\cdot$ )、 $\lambda_1, \lambda_2$  は閾値パラメータ、 $\kappa \ge 1$ は バイアス低減処理 ( $\Xi_{\kappa}[z]$ )<sub>ijn</sub> に用いるパラメータで あり、以下の式のように表される。

$$\left(\Xi_{\kappa}\left[\boldsymbol{z}\right]\right)_{ijn} = \left(\kappa z_{ijn}/\max_{ijn}\left\{z_{ijn}\right\}\right) \qquad (10)$$

 $\zeta_{iin}^{z,\kappa}$ は要素毎の firm thresholding を表している。

$$\zeta_{ijn}^{z,\kappa} = \Xi_{\kappa} \left[ \left( 1 - \lambda_2 / z_{ijn} \right)_{-} \right]$$
(11)

ただし、 $(\cdot)_{-} = \min(1, \cdot)$ である。 $\Theta_{\eta}^{Freq}$ は周波数帯 域毎のスパース性を復元するためのピークを計測す るための重みである。この重みは、白色化によって励 起した雑音部分の周波数帯域のパワーを軽減する効 果がある。

$$\left(\Theta_{\eta}^{Freq}\left[\boldsymbol{x}\right]\right)_{i} = \Upsilon_{\eta} \left[\frac{\left(\Sigma_{m=1}^{M}\Sigma_{j=1}^{J}\left|\boldsymbol{x}_{ijm}\right|^{2}\right)^{\frac{1}{2}}}{\Sigma_{m=1}^{M}\Sigma_{j=1}^{J}\left|\boldsymbol{x}_{ijm}\right|}\right] \quad (12)$$

Υ<sub>η</sub>[·] は L<sub>1</sub> 正規化を表しており以下の形で定式化さ れる。

$$\Upsilon_{\eta}\left[\boldsymbol{\xi}\right] = \boldsymbol{\xi}_{\eta} / \left( \left\| \boldsymbol{\xi}_{\eta} \right\|_{1} / I \right) \qquad \xi_{\eta} = \left( \boldsymbol{\xi} - \eta \right)_{+}$$

 $\xi$ は  $\Upsilon_{\eta}$  への引数、 $\xi_{\eta}$  は  $\xi$  を閾値処理した後の値と なる、元々スパース性が低い雑音の支配的な周波数帯 域をカットする効果がある。

# 3 提案法

#### 3.1 概要

スパース IVA では周波数ビン毎のスパース性の計 算を行なっている。提案法では周波数方向の他に二つ のスパース性の導入を行う。一つ目は信号のパワーの 時間変動のスパース性、二つ目は音源の到来方向に 対するスパース性である。

#### 3.2 時間方向のスパース性

時間方向のスパース性は以下のように、各時間フ レームにおける周波数ビンでの観測信号のパワーの 総和を計算している。

$$\left(\Theta_{\eta}^{Time}\left[\boldsymbol{x}\right]\right)_{j} = \Upsilon_{\eta} \left[\frac{\left(\Sigma_{m=1}^{M}\Sigma_{i=1}^{I}\left|\boldsymbol{x}_{ijm}\right|^{2}\right)^{\frac{1}{2}}}{\Sigma_{m=1}^{M}\Sigma_{i=1}^{I}\left|\boldsymbol{x}_{ijm}\right|}\right]$$
(13)

提案法では、 $\Theta_{\eta}^{Time}$ を式 (9) の  $\Theta_{\eta}^{Freq}$ の代わりに用いる。

#### 3.3 到来方向のスパース性

到来方向はマイク間の位相差によって求める。到来 方向に基づくバイナリマスクは以下のように求めて いる。

$$\mathcal{D}_{i,j,n} = \begin{cases} 1 & (|\theta_{ij} - \theta_n| > \theta_c) \\ \epsilon & else \end{cases}$$
(14)

ただし  $\mathcal{D}_{i,j,n}$  は音源 n を強調するバイナリマスク、  $\theta_{i,j}$  は周波数ビン i と時間フレーム j における推定到 来角度、 $\theta_n$  は音源 n の到来角度、 $\theta_c$  は到来方向を許 容するための定数である。後の実験では $\theta_{i,j}$  のヒスト グラムの最頻値から到来方向  $\theta_n$  を求めている。

到来方向のスパース性を扱う際、本稿では二通り のマスクを提案する。一つ目は到来方向分離マスク で式 (9) に代わって以下のような重みを用いる。

$$\Theta_{n,n}^{DOAsep} = \mathcal{D}_{i,j,n} \tag{15}$$

これは各音源の到来方向に対するバイナリマスクを 用いた重みであり、式 (9) に使用することで音源 n ご とに強調するようなマスク ( $\mathcal{M}[z]$ )<sup>DOAsep</sup>を生成す る。二つ目は到来方向共通マスクで以下のような重 みを用いる。

$$\Theta_{\eta}^{DOAshare} = \max_{n} \mathcal{D}_{i,j,n}$$
(16)

これは直接音を通し、それ以外の音を排除するマス クである。例として2方向から音源が到来する場合、 到来方向分離マスクでは1方向の音源のみを強調し 他の到来方向の音源を排除するマスクを音源数だけ 持つが、到来方向共通マスクでは到来する2方向を どちらとも通す1つのマスクを持っている。これは、 直接音の分離は音源分離アルゴリズムに任せ、残響成 分や雑音成分をマスクにより除去することを狙って いる。

#### 3.4 スパースネスの統合

従来の周波数方向、提案した時間方向と到来方向 のスパース性をそれぞれ音源項として式(7)に加え、 複数の音源項がある場合の近接分離最適化による音 源分離[6]として統合することで多様な音源への対応 を試みる。スパースネスを統合した際の最適化式を 以下に表す。

 $\begin{array}{ll} \underset{\mathbf{w}}{\text{Minimize}} & \lambda_{Freq} \mathcal{P}^{Freq}(\mathbf{X}\mathbf{w}) + \lambda_{Time} \mathcal{P}^{Time}(\mathbf{X}\mathbf{w}) \\ & + \lambda_{DOA} \mathcal{P}^{DOA}(\mathbf{X}\mathbf{w}) + \mathcal{L}(\mathbf{w}) \ (17) \end{array}$ 

ただし各スパース IVA の音源項を  $\mathcal{P}^{[\cdot]}$ 、音源項を適 用する大きさを表す重みを  $\lambda_{[\cdot]}$ 、 [·] の添え字は各ス パースネスを表し、周波数方向は *Freq*、時間方向は *Time*、到来方向は *DOA* で表す。

#### 4 評価実験

#### 4.1 スパースネスの違いによる分離性能の比較

#### 4.1.1 実験条件

本実験では2音源の混合音を2チャネルで音源分離 することで精度を確かめる。分離する音源は人間同士 の話し声から構成される音声データと、楽器や歌声で 構成される音楽データを対象とした。音声データの信 号はSiSEC[8]のUNDの女声データの内、-50度と45 度から到来する2つの信号を合成して MixtureA,-10 度と15度から到来する信号を合成して MixtureB と して使用する。音楽データの信号は RWCP 実環境音 声・音響データベースのインパルス応答 (E2A[10]) に、 SiSEC で提供されている楽曲 Bearlin-Roads と楽曲 Ultimate-Nz-Tour の各楽器の音源データを Table 2の ように畳み込んで生成した。初期値は分離行列 W に単 位行列、双対変数dにゼロベクトルを与えた。スパー ス IVA において閾値 η は ξ の最大値を超えるような値 に設定しまった場合、重み $\Theta_n$ の値がすべて零になって しまう恐れがあるため、ここでは ξの値を利用するこ とで、閾値を決定する。具体的には、 $\xi_i$ ( $i = 1, 2, \dots, I$ ) を昇順に並べ替えた $\hat{\xi}_i$ ( $i = 1, 2, \cdots, I$ )を導入し以下 のように閾値  $\eta$  を求めている。 $\eta = \hat{\xi}_k, k = I * \delta$ 。音 源分離性能の評価には、SDR[9] を用いた。値が高い ほど分離性能が高いことを表す。

#### 4.1.2 実験結果

分離実験の結果を Fig. 2 に示す。比較対象は従来 法となるラプラス IVA、周波数方向のスパース IVA、 そして提案法となる時間方向のスパース IVA、到来

#### Table 1 分離処理に用いるパラメータ

サンプリング周波数	16kHz
フレームサイズ	2048
シフトサイズ	1024
$\mu_1$ , $\mu_2$	1, 1
$\lambda_1,\lambda_2$	2, 0.01
$\kappa$	1.1
$収 東速 度 \alpha$	1
閾値パラメータδ	0.5

#### Table 2 使用した音楽データ

ID	Author/Song	Part 1	Part2
1	bearlin1	$Piano(90^\circ)$	$Vocal(10^{\circ})$
2	bearlin2	$Vocal(10^{\circ})$	$Ambient(150^\circ)$
3	bearlin3	$Piano(90^\circ)$	$Ambient(150^\circ)$
4	ultimate1	$Guitar(90^\circ)$	$Synth(10^{\circ})$
5	ultimate2	$\operatorname{Synth}(10^\circ)$	$Drum(150^\circ)$
6	ultimate3	$Guitar(90^\circ)$	$Drum(150^\circ)$



■ラプラスIVA(従来) ■周波数方向(従来) ⊠時間方向 ☑到来方向(分離) ■到来方向(共通)

#### Fig. 2 スパースネスの違いによる分離性能の比較

方向の分離マスクと共通マスクとしている。横軸は分 離対象の音源、縦軸は分離後の信号の平均 SDR を表 している。従来の周波数方向のスパース IVA は、音 声データに対しては有効であるものの、音楽データの 2曲に対しては一転して分離精度が他手法と比べて低 くなった。これは楽器音同士のパワーの偏りがあるに もかかわらず閾値処理を行ったため、片方のパワーの 小さい音源を削ってしまったためと考えられる。時間 方向マスクが最も改善量が大きい。到来方向マスクは 分離マスクに比べて共通マスクの方が良好な結果が 出た。到来方向分離マスクは各音源をそれぞれ強調 するような特徴を持っているが、マスクが強調しよう とする音源と音源分離アルゴリズムにより分離され た分離音の不整合により、分離が進まない信号があっ た。そのため到来方向共通マスクの方が優れている 結果となった。

#### 4.2 スパースネス統合 IVA の実験

# 4.2.1 実験条件

上記のスパース性の違いによるスパース IVA の分離性能の比較を行った結果から、音声データに対しては従来法の周波数方向のスパース IVA、音楽デー



Table 3 実験に用いた音楽データ

Fig. 3 スパースネスを統合した際の分離結果

タに対しては時間方向のスパース IVA が有効である ことが分かった。よって 3.4 節の手法により 3 つの スパースネスを統合した手法により分離実験を行う。 ただし到来方向に対するスパース性については、安 定性から到来方向共通マスクのみを用いる。統合に は各音源項を適用する大きさを調整するための重み を設定する。音声データに対しては周波数方向のス パース IVA の重みを大きくした音声強調マスク、音 楽データに対しては時間方向のスパース IVA の重み を大きくした音楽強化マスクを用いて検討する。各 音源項を統合する重みは Table 3 に表している。

#### 4.2.2 実験結果

分離実験の結果を Fig. 3 に示す。比較対象は従来 法と音声強化マスク、音楽強化マスクである。縦軸は 分離後の信号の平均の SDR を表している。音声強化 マスクについては、従来の周波数方向のスパース IVA と比べて音声データの SDR がほぼ変わらない性能で、 楽器音が平均で約 1.4dB ほど向上した。音楽強化マ スクについては、音楽データに強いラプラス IVA に 対しても分離性能の若干の向上がみられ、さらに音声 データでは約 2.2dB 分離性能が向上した。音声強化 マスクや音楽強化マスクの双方で重みが最も大きい マスクの得意な音源に対する性能を維持しつつ、そ の他の音源に対して分離性能の向上ができた。これ は異なるスパースネスを統合することで、単一のス パースネスだけでは考慮されなかった別の特徴を持 つスパースネスの利用ができたためと考えられる。

# 5 まとめ

本稿では、まずラプラス IVA に関して周波数方向 のスパース性を導入したスパース IVA は音声データ には有効だが、音楽データでは性能が落ちることを 実験によって示した。そこで多様な音源に対応するた め、時間方向と到来方向の2つのスパース性を導入 する手法を提案した。分離実験の結果、平均的には時間方向のスパース性の導入が最も有効であることが分かった。また、周波数方向、時間方向、到来方向の3つのスパース性を統合し重みを工夫したスパースネス統合 IVA を用いることで音声・音楽の両データの性能が改善した。

# 参考文献

- T.-W. Lee, "Independent Component Analysis-Theory and Applications," Norwell, MA: Kluwer, 1998.
- [2] I. Lee *et al.*, "Fast fixedpoint independent vector analysis algorithms for convolutive blind source separation," Signal Processing 87(8), 2007
- [3] 澤田 宏, "同期同時対角化によるブラインド信号 分離,"第 32 回信号処理シンポジウム講演論文 集, pp.332-337, 2017
- [4] D. Kitamura *et al.*, "Dtermined blind source separation unifying independent vector analysis and nennegative matrix factorization," IEEE/ACM Transaction on Audio, Speech, and Language Processing, 24(9):1626-1641, 2016
- [5] H. Sawada *et al.*, "Multichannel Extensions of Non-Negative Matrix Factorization with Complex-Valued Data," IEEE Trans. ASLP, vol.21, no.5, pp.971-982, 2013.
- [6] Kohei Yatabe et al., "Determined blind source separation via proximal splitting algorithm," Proceedings of IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP 2018), pp. 776-780, 2018
- [7] Kohei Yatabe et al., "Time-frequency-maskingbased determined BSS with application to sparse IVA," Proceedings of IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP 2019), pp. 715-719, 2019.
- [8] S. Araki, et al., "The 2011 signal separation evaluation campaign (SiSEC2011): - audio source separation -," in Proc. Int. Conf. Latent Variable Anal. Signal Separation, pp. 414-422 2012.
- [9] E. Vincent *et al.*, "First Stereo Audio Source Separation Evaluation Campaigh: Data Algprithm and Results," Independent Component Analysis and Signal Separation(Springer, Bearlin, 2007), pp.552-559.
- [10] RWCP: "実環境音声・音響データベース (RWCP-SSD)" 音声資源コンソーシアム, http://research.nii.ac.jp/src/ RWCP-SSD.html,閲覧日:2019/12/19.