

音環境の変動に頑健な音源分離システムの開発

石橋 孝昭* 中島 栄俊**

Blind Source Separation under a Dynamic Acoustic Environment

Takaaki Ishibashi*, Hidetoshi Nakashima**

Independent component analysis can estimate unknown source signals from their mixtures under the assumption that the source signals are statistically independent. However, in a real environment, the separation performance is often deteriorated because the number of the source signals is different from that of the sensors. In this paper, we propose an estimation method for the number of the sources based on the joint distribution of the observed signals under two-sensor configuration. From several simulation results, it is found that the number of the sources is coincident to that of peaks in the histogram of the distribution. The proposed method can estimate the number of the sources even if it is larger than that of the observed signals. And we propose a new blind source separation method based on the estimated number on the sources under a dynamic acoustic environment. The proposed methods have been verified by several experiments.

キーワード：ブラインド信号分離，独立成分分析，原信号数推定，目的信号抽出，音環境の変動

Keywords: blind source separation, independent component analysis, source number estimation, target source extraction, dynamic acoustic environment

1. はじめに

人間と機械のコミュニケーションにおいて，ユーザーに負担が無く，高速に情報伝達できるインターフェースは音声であり，多くの音声認識システムが実用化されている．しかし，現存の音声認識は，雑音のある環境では認識率が低下するため，雑音を含めて学習させる認識方法や，認識の前処理である雑音除去の研究が進められているが，音環境の変動で性能が低下する問題がある．

その原因の一つは，既存の音源分離では音源の情報は未知と設定しているにもかかわらず，音源数や音源の特徴が分かっている条件下で研究が進められているからである．このことが，実用化するときの大きな問題として立ちはだかり，音源分離はなかなか実用化できない．

そのため，音源分離の前処理として，これまでに固定された環境下で音源の数を推定する研究はいくつか報告され

ているが，処理時間がかかる問題がある．また，時々刻々と変動する音源数をリアルタイムで推定する研究の報告は見当たらない．さらに，音環境の変動を考慮した分離アルゴリズムに対する研究報告も無いようである．

そこで本研究では，音環境の変動と音源数を推定し，その結果に基づいて音源分離や雑音除去を行って，目的信号を抽出することのできるシステムの構築を目的とする．

2. ブラインド信号分離

ブラインド信号分離 (BSS: Blind Source Separation) は，目的信号と，目的信号に雑音混合される過程を未知として，混合前の原信号を推定する技術である． N 個の音源 $\mathbf{s}(t) = [s_1(t), \dots, s_n(t), \dots, s_N(t)]^T$ を M 個のマイクロホンで観測するとき，混合信号 $\mathbf{x}(t) = [x_1(t), \dots, x_m(t), \dots, x_M(t)]^T$ は，混合行列 A によって

$$\mathbf{x}(t) = A\mathbf{s}(t) \dots \dots \dots (1)$$

のように表現できると考える．このとき， $\mathbf{x}(t)$ と M のみが既知で， $\mathbf{s}(t)$ ， N ， A は未知である．

$\mathbf{x}(t)$ のみを用いて， $\mathbf{s}(t)$ と A を推定することがブラインド信号分離の目的である． A の逆関数 W を推定すれば，分離信号 $\mathbf{u}(t) = [u_1(t), \dots, u_n(t), \dots, u_N(t)]^T$ は

$$\mathbf{u}(t) = W\mathbf{x}(t) \dots \dots \dots (2)$$

* 情報通信エレクトロニクス工学科
〒 861-1102 熊本県合志市須屋 2659-2
Department of Information, Communication and Electronic Engineering
2659-2, Suya, Koshi, Kumamoto 861-1102
** 制御情報システム工学科
〒 861-1102 熊本県合志市須屋 2659-2
Department of Control and Information Systems Engineering
2659-2, Suya, Koshi, Kumamoto 861-1102

のように得られる。

ブラインド信号分離問題を解く方法の一つに独立成分分析がある。独立成分分析は、混合信号が互いに統計的に独立な音源の重ね合わせであるという仮定の下で、混合信号から独立成分を分解する統計的手法であり、様々なアルゴリズムが提案されている。独立成分分析は音源のパワーに依存せず、非定常信号も分離できる特徴を持つ。

2.1 独立成分分析

独立成分分析 (ICA: Independent Component Analysis) は、原信号が統計的に独立という仮定の下で原信号を推定する統計的手法で、エンジン音のような定常信号だけでなく話者音声や音楽などの非定常信号も分離できるという特徴を持っている。また、原信号の推定だけでなく、観測信号の隠れた構造を見つけるための特徴抽出としても利用されている^{(1)~(6)}。

独立成分分析では、各信号を確率変数として扱う。また、原信号 $s_n(t)$ は互いに統計的に独立と仮定される。すなわち、 $s_n(t)$ の同時確率密度関数 $p(s_1(t), \dots, s_N(t))$ は

$$p(s_1(t), \dots, s_N(t)) = \prod_{n=1}^N p(s_n(t)) \dots \dots \dots (3)$$

のように、 $s_n(t)$ の周辺確率密度関数 $p(s_n(t))$ の積で表現できると仮定する。この仮定は、多くの場合は非現実的な仮定でなく、実際には厳密に独立である必要はないとされている。一方、式 (1) によって変換された混合信号 $x_m(t)$ の同時確率密度関数 $p(x_1(t), \dots, x_M(t))$ は

$$p(x_1(t), \dots, x_M(t)) \neq \prod_{m=1}^M p(x_m(t)) \dots \dots \dots (4)$$

のように独立性が失われる。したがって、式 (2) によって推定される分離信号 $u_n(t)$ の同時確率密度関数 $p(u_1(t), \dots, u_N(t))$ が

$$p(u_1(t), \dots, u_N(t)) = \prod_{n=1}^N p(u_n(t)) \dots \dots \dots (5)$$

のように、周辺確率密度関数 $p(u_n(t))$ の積で表現できて統計的に独立になれば、 $u_n(t)$ は原信号 $s_n(t)$ とみなすことができる。このような、原信号の独立性のみを手掛かりとして、原信号を推定する統計的手法を独立成分分析という。

ブラインド信号分離では原信号が未知であるため、その確率密度関数も未知である。したがって、独立性の評価方法を工夫することで、多くの独立成分分析のアルゴリズムが提案されている。具体的には、原信号の非ガウス性を利用し全ての原信号を同時に復元する方法として、分離信号の相互情報量を最小化する方法、最尤法を用いた方法、高次キュムラントを用いる方法などがある⁽¹⁾。また、原信号を一つずつ復元する方法として、尖度を用いる方法、Negentropy に基づく方法などがある⁽²⁾。また、原信号の非定常性に基づく方法や信号の時間相関に基づく方法がある⁽⁷⁾⁽⁸⁾。その他にも、原信号の確率密度関数を推定して独立成分分析を

行う方法や幾何学的な観点から解く方法なども研究されている^{(9)~(11)}。その中から、Kullback-Leibler 情報量の最小化に基づく Natural Gradient と、Negentropy の最大化に基づく FastICA について以降で述べる。

2.2 Natural Gradient

Natural Gradient⁽¹⁾ は、独立性の評価に Kullback-Leibler 情報量

$$I_{KL}(\mathbf{u}(t)) = \int p(\mathbf{u}(t)) \log \frac{p(\mathbf{u}(t))}{\prod_{n=1}^N p(u_n(t))} d\mathbf{u}(t) (6)$$

$$= \sum_{n=1}^N \mathcal{H}(u_n(t)) - \mathcal{H}(\mathbf{u}(t)) \dots \dots \dots (7)$$

を用いて、分離信号の Kullback-Leibler 情報量が最小になるように分離行列 W を推定する逐次更新アルゴリズムである。ここに、 $\mathcal{H}(u_n(t))$ は Entropy で

$$\mathcal{H}(u_n(t)) = - \int p(u_n(t)) \log p(u_n(t)) du \dots \dots \dots (8)$$

のように定義され、情報のあいまいさを表す。

$I_{KL}(\mathbf{u}(t))$ の最小値を最急降下法によって解いた後、ユークリッド空間からリーマン空間へ拡張すれば

$$W + \Delta W = W - \eta E[\varphi(\mathbf{u}(t))\mathbf{u}(t)^T - I]W \dots (9)$$

の逐次更新式が得られる。ここに、 η は学習係数、 I は単位行列である。このとき、 $\varphi(\mathbf{u}(t)) = [\varphi_1(u_1(t)), \dots, \varphi_N(u_N(t))]^T$ は、 $u_n(t) = s_n(t)$ の確率密度関数が未知であるため、求められない。しかし、 $u_n(t)$ の分布がラプラス分布のようなスーパーガウシアン、すなわち、裾が長くピークが尖った分布の場合

$$\varphi_n(u_n(t)) = \tanh(u_n(t)) \dots \dots \dots (10)$$

を選択し、一方、一様分布のようなサブガウシアン、すなわち、裾が短くピークが滑らかである分布の場合

$$\varphi_n(u_n(t)) = u_n(t)^3 \dots \dots \dots (11)$$

を選択することで、安定して分離できる。

2.3 FastICA

FastICA⁽²⁾ は、独立性の評価に Negentropy を用いて、分離信号の Negentropy が最大となるような分離荷重ベクトル \mathbf{w}_n を一つずつ推定する方法である。Negentropy $J(u_n(t))$ は、常に非負の値をとり、確率変数がガウス分布であるときのみ 0 となる統計量で

$$J(u_n(t)) = \mathcal{H}(u_g) - \mathcal{H}(u_n(t)) \dots \dots \dots (12)$$

と定義される。ここに、 u_g は $u_n(t)$ と等しい平均と分散を持つガウス分布の確率変数である。しかし、 $J(u_n(t))$ を非ガウス性の尺度として利用するには、確率変数 $u_n(t)$ の確率密度関数を推定する必要があり、その計算量は膨大となる。そこで、 $J(u_n(t))$ は

$$J(u_n(t)) \propto \sum_{i=1}^p \alpha_i \{E[F_i(u_n(t))] - E[F_i(\nu)]\}^2 \quad (13)$$

のように近似される。ここに、 α_i は正の定数、 ν は平均が 0 で分散が 1 のガウス変数、 $F_i(\cdot)$ は 2 次的でない関数である。

FastICA アルゴリズムは、この $J(u_n(t))$ を用いて、 $u_n(t)$ の分散が 1 という制約条件の下で、その最大値を求める条件付き最適化問題となる。この最適化問題を解くと、FastICA アルゴリズムの荷重更新式は

$$\mathbf{w}_n^+ = E[\tilde{\mathbf{x}}(t)f(\mathbf{w}_n^T \tilde{\mathbf{x}}(t))] - E[f'(\mathbf{w}_n^T \tilde{\mathbf{x}}(t))]\mathbf{w}_n \quad (14)$$

$$\mathbf{w}_n = \frac{\mathbf{w}_n^+}{\|\mathbf{w}_n^+\|} \dots \dots \dots \quad (15)$$

のように導出できる。ここに、 $\tilde{\mathbf{x}}(t)$ は白色化後の $\mathbf{x}(t)$ 、 $f(\cdot)$ は $F(\cdot)$ の導関数、 $f'(\cdot)$ は $f(\cdot)$ の導関数である。収束条件は、更新前後の分離荷重 \mathbf{w}_n の向きが一致して

$$|\mathbf{w}_{n,\text{old}}^T \mathbf{w}_{n,\text{new}}| \simeq 1 \dots \dots \dots \quad (16)$$

を満たすことである。ここに、添え字の old と new はそれぞれ更新前後を表す。ここで得られた分離荷重 \mathbf{w}_n を用いて、 $u_n(t) = \mathbf{w}_n^T \tilde{\mathbf{x}}(t)$ と計算すれば、分離信号 $u_n(t)$ は原信号の一つが推定されることになる。

二つ以上の原信号を推定する場合、上述の FastICA アルゴリズムを適用して、原信号を一つ一つ推定すればよい。このとき、二つ以上の分離荷重が同じ値をとる可能性があるため、それぞれの分離荷重を直交化して無相関化する。具体的には、二つ以上の分離荷重が得られたとき、グラムシュミットの直交化法により

$$\mathbf{w}_n = \mathbf{w}_n - \sum_{i=1}^{n-1} \mathbf{w}_i^T \mathbf{w}_n \mathbf{w}_i \dots \dots \dots \quad (17)$$

のように、それぞれの分離行列を直交化させて、再度、式 (15) に代入し規格化する。このようにして、全ての分離荷重 \mathbf{w}_n ($n = 1, \dots, N$) を推定すれば、分離行列 W は $W = [\mathbf{w}_1, \dots, \mathbf{w}_N]^T$ のように求められる。

3. 音環境の変動に頑健な音源分離システム

独立成分分析を用いたブラインド信号分離は、一般に、音源数 N はマイクロホン数 M と等しい仮定の下で導出されたアルゴリズムが多く、 $N \neq M$ のときに分離性能が低下するという問題がある。したがって、ブラインド信号分離を行うには、音源数を知るための前処理が必要となる。以降では、観測信号から音源数 N を推定し、目的音声抽出する方法について報告する。

3.1 音環境の変動検出

音源は図 1 のような話者音声で、観測信号は複数話者の同時発話となる場合を考える。また、二つのマイクロホン ($M = 2$) による観測信号を用いる。

まず、音源が存在しない場合、観測信号 $x_1(t)$ と $x_2(t)$ は

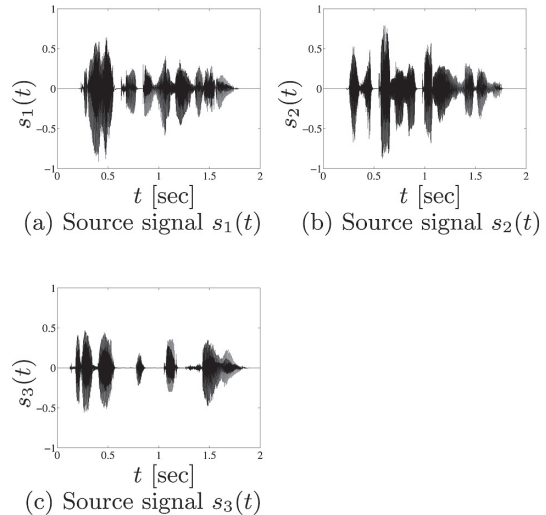


図 1 Source signals.

共に振幅が 0 の信号となることは明らかである。したがって、全ての観測信号の振幅が 0 と判断できるときには、音源数は $N = 0$ と容易に推定できる。

次に、音源が一つ (図 1 の $s_1(t)$ のみが存在) の場合、観測信号 $x_1(t)$ と $x_2(t)$ は図 2 のように得られる。観測信号は音源からマイクロホンまでの減衰率が異なるだけなので、振幅のスケールが異なる波形となる。

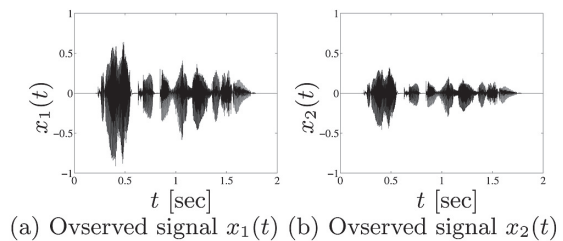


図 2 Observed signals ($N = 1$).

このときの観測信号の同時分布を図 3(a) に示す。グラフの横軸は $x_1(t)$ 、縦軸は $x_2(t)$ である。図から、スケールのみが異なるため同時分布は直線になることが分かる。また、この分布の方位に対するヒストグラムを図 3(b) に示す。すなわち、各観測時刻におけるプロット点の角度 θ を

$$\theta(t) = \tan^{-1} \frac{x_2(t)}{x_1(t)} \dots \dots \dots \quad (18)$$

を計算して、その角度に対するヒストグラムを作成する。横軸は $-\pi/2 \sim \pi/2$ rad の同時分布の方位、縦軸はその方位に対する頻度を示す。同時分布が一方向のみに存在するため、ヒストグラムは一つの尖ったピークを持つ。

音源が二つ (図 1 の $s_1(t)$ と $s_2(t)$ が存在) の場合、観測信号は図 4 となる。複数の音源が混合したため、二つの観測信号の波形は全く異なる形状になることが分かる。

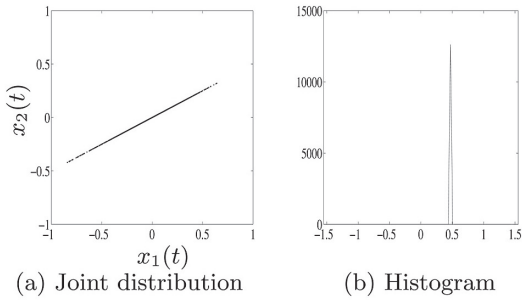


図3 Joint distribution and histogram ($N = 1$).

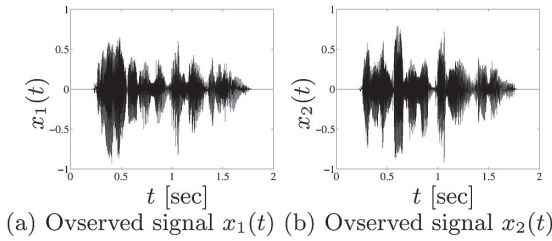


図4 Observed signals ($N = 2$).

観測信号の同時分布を図5(a), そのヒストグラムを(b)に示す。(a)の同時分布から二本の直線上に多く分布していることが分かる。このことは、(b)のヒストグラムが二つの尖ったピークを持つことから明確である。

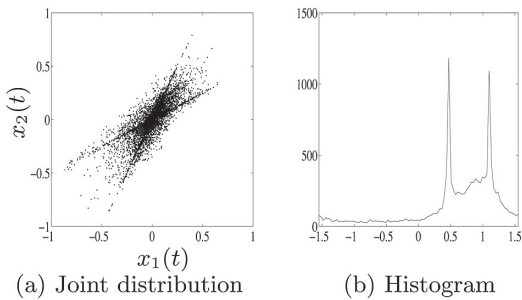


図5 Joint distribution and histogram ($N = 2$).

最後に、 $N = 3$ (図1の $s_1(t)$, $s_2(t)$, $s_3(t)$ が存在) の観測信号は図6のように得られる。この同時分布を図7(a)に、そのヒストグラムを(b)に示す。(a)の同時分布では特徴を見出すことが困難であるが、(b)のヒストグラムでは三つの尖ったピークの存在を確認できる。

3.2 音源数推定

上述の事実から、観測信号の同時分布から作成されるヒストグラムは、音源数と等しい数の尖ったピークを持つことが分かる。したがって、このピーク数を求めることで音源数を推定できる。ピーク数を求めるために、本稿では、ヒ

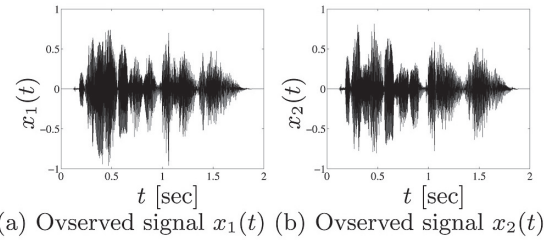


図6 Observed signals ($N = 3$).

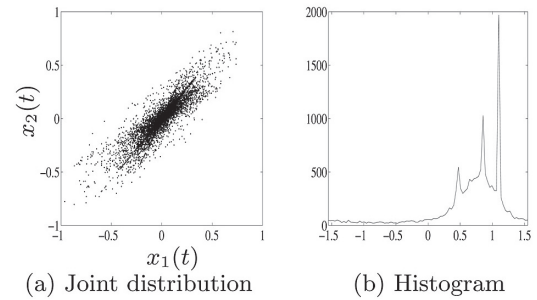


図7 Joint distribution and histogram ($N = 3$).

ストグラムの振幅 $h(k)$ の変動に着目し

$$h(k + 1) - h(k) \geq \epsilon \dots \dots \dots (19)$$

を満たす点の数を求めることで、音源数を推定する。ここに、 k はヒストグラムのフレーム番号、 ϵ は閾値である。

3.3 変動する音響環境下での音源分離

音源数を推定した後、その数によって図8のように音源分離や雑音除去の処理を行う。すなわち、音源数が0と推定された場合、何も出力しない。音源数が1と推定された場合、その音が目的音源と判断したときのみ出力する。音源数が2以上の場合、ブラインド信号分離によって音源分離し、分離信号から目的音源を選択し出力する。以上の処理によって、変動する音響環境下での音源分離が実現可能となる。

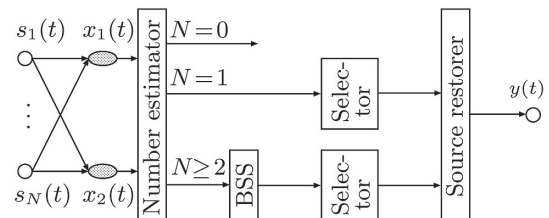


図8 Blind source separation based on the estimation for the number of the source signals.

また、独立成分分析の固有の問題である分離信号のスケー

ルの不定性については

$$\mathbf{v}_n(t) = W^{-1}[0, \dots, 0, u_n(t), 0, \dots, 0]^T \dots \dots (20)$$

のように計算することで解消する方法が提案されている⁽⁷⁾. $\mathbf{v}_n(t) = [v_{n1}(t), \dots, v_{nM}(t)]^T$ は, n 番目の音源から m 番目のマイクロホンまでの伝達関数 a_{mn} と音源 $s_n(t)$ の積で

$$v_{nm}(t) = a_{mn}s_n(t) \dots \dots \dots (21)$$

のように一意に決まることが証明されている. したがって, 分割スペクトルがマイクロホンで観測された原信号の推定値であることは明らかである.

さらに, 出力チャンネルの選択問題に関しては kurtosis に基づいて解決する. 一般的に, 話者音声は騒音に比べて尖度が高い特徴を持つ. この特徴を利用して, kurtosis を用いた話者音声の抽出が可能となる. kurtosis は 2 次のモーメントと 4 次のモーメントで

$$\kappa_4(u_n(t)) = E[u_n^4(t)] - 3\{E[u_n^2(t)]\}^2 \dots \dots \dots (22)$$

と定義される. $u_n(t)$ が図 9 の鎖線で示されるようなガウス分布の場合, 4 次のモーメント $E[u_n^4(t)]$ は $3\{E[u_n^2(t)]\}^2$ と等しくなり, kurtosis の値は 0 になる. 一方, $u_n(t)$ がガウス分布でない場合, kurtosis は正か負のどちらかの値をとる. 正の kurtosis を持つ確率変数は supergaussian と呼ばれ, 負の kurtosis を持つ確率変数は subgaussian と呼ばれる. supergaussian は, 図 9 の実線で示される分布のように, 裾が長く, ピークが尖った分布を持つ. 一方, subgaussian は, 図 9 の破線で示されるような平らな分布を持つ. したがって, 目的音声の話者音声で雑音が定常信号に近い騒音である場合, 分離信号から kurtosis に基づいて目的音声 $y(t)$ を抽出するには

$$y(t) = \begin{cases} u_1(t) & \text{if } \kappa_4(u_1(t)) > \kappa_4(u_2(t)) \\ u_2(t) & \text{if } \kappa_4(u_1(t)) < \kappa_4(u_2(t)) \end{cases} \dots \dots (23)$$

のように, kurtosis の値が大きい推定信号を話者音声と判断することになる.

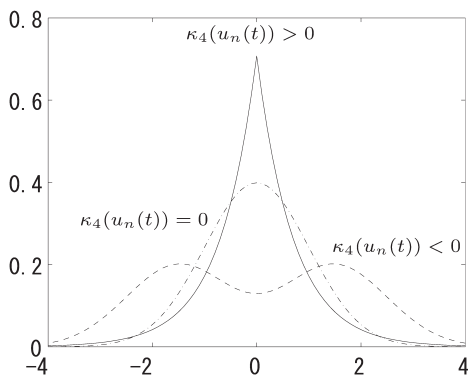


図 9 Probability density functions.

4. 実験および考察

本研究手法の有効性を検証するために実験を行った. 新聞記事読み上げ音声コーパス⁽¹²⁾ から 10 パターン (男性 5 パターン, 女性 5 パターン) の約 2 秒間の話者音声と, Ambient Noise Database⁽¹³⁾ から 2 パターンの駅構内の約 2 秒間の騒音を, それぞれ音源 1 と音源 2 として使用した. また, 目的音声は音源 1 とした. さらに, この音源については, 図 10(a) のように, 時刻とともに音源数が変動するように, それぞれ, 3 秒間の無音区間を追加した. すなわち, 混合信号の約 0~1 秒は $N = 0$, 約 1~2 秒は $N = 1$ で目的音声のみが存在し, 約 2~3 秒は $N = 2$, 約 3~4 秒は $N = 1$ で雑音のみが存在し, 約 4~5 秒は $N = 0$ とした. 音声データについては, サンプリング周波数を 8000Hz, 分解能を 16bit とした. 音源数の推定やブラインド信号分離の処理単位時間は 0.5 秒とした. ブラインド信号分離については FastICA アルゴリズムを用い, 初期荷重をノルムが 1 の乱数, 非線形関数を $f(|u_n(t)|^2) = 1 - 2/(e^{2|u_n(t)|^2} + 1)$, 最大繰返し回数を 100, 誤差を 0.000001 とした.

シミュレーションによる観測信号は (b) である. 観測信号から音源の波形や発話区間を判断することは困難であることが分かる. この観測信号のみを用いて目的信号を抽出した結果が (c) である. 提案法は目的信号 $s_1(t)$ を抽出可能であり, 処理単位ごとのスケールの不定性や出力チャンネル選択の問題についても解消されていることが分かる. 全ての実験パターンにおいて同様の結果が得られ, 提案法の有効性が確認された.

5. まとめ

本論文では音環境の変動に頑健な音源分離システムの開発を目的として, まず, 観測信号の同時分布を用いて音環境の変動を検出した. また, 同時分布から作成されるヒストグラムに基づいて原信号を推定する方法を提案した. さらに, 推定した原信号数に基づいた変動環境下での音源分離法を提案した. シミュレーションにより, 提案法は目的音声を抽出することが可能であり, その有効性が確認された.

音源分離や雑音除去に限らず, 音響信号処理の分野における大きな問題に, 高残響下で高精度に機能する処理方法の確立が挙げられる. 本研究においても, 高残響下で機能するアルゴリズムを開発することが今後の研究の方向である. また, リアルタイムブラインド信号分離の開発や, 実用化へ向けたアプリケーションの開発についても研究を進める予定である.

(平成 22 年 9 月 22 日受付)

文 献

- (1) A. Cichocki and S. Amari: Adaptive blind signal and image processing, learning algorithm and applications; *John Wiley & Sons, Ltd* (2002)
- (2) A. Hyvärinen, J. Karhunen and E. Oja: Independent component analysis; *John Wiley & Sons, Ltd* (2001)

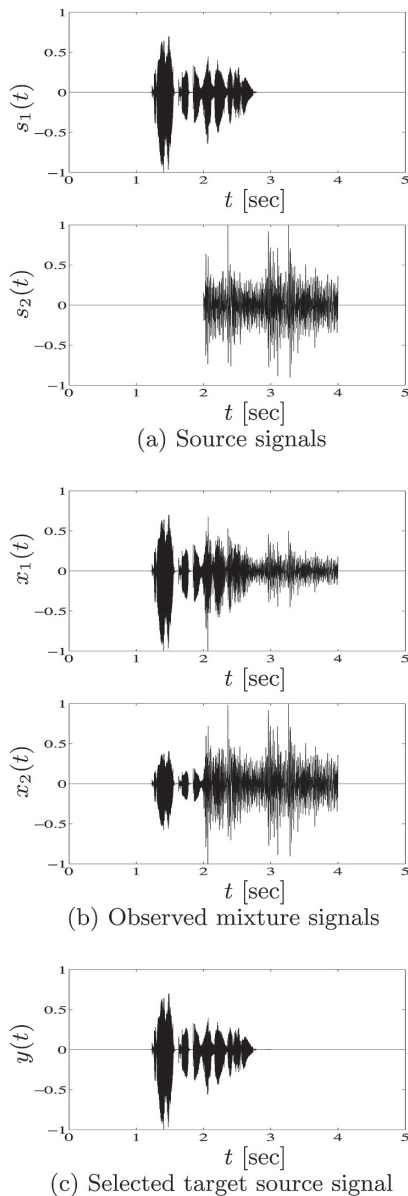


図 10 Experimental results on blind source separation under a dynamic acoustic environment.

- (11) H. Shindo and Y Hirai: Blind source separation by a geometrical method; *Proceedings of 2002 International Joint Conference on Neural Networks* pp. 1109-1114 (2002)
- (12) Acoustical Society of Japan: ASJ continuous speech corpus japanese newspaper article sentences; *JNAS Vols.1-16*, 1997.
- (13) NTT Advanced Technology Corporation: Ambient noise database for telephony 1996; 1996.

- (3) T. W. Lee: Independent component analysis; *Kluwer Academic Publishers* (1998)
- (4) 甘利俊一, 狩野裕, 佐藤俊哉, 松山裕, 竹内啓, 石黒真木夫: 多変量解析の展開 隠れた構造と因果を推理する; 株式会社岩波書店 (2002)
- (5) 甘利俊一, 村田昇: 独立成分分析 - 多変量データ解析の新しい方法; 株式会社サイエンス社 (2002)
- (6) 村田昇: 入門 独立成分分析; 東京電機大学出版局 (2004)
- (7) N. Murata, S. Ikeda and A. Ziehe: An approach to blind source separation based on temporal structure of speech signals; *Neurocomputing*, Vol. 41, Issue 1-4, pp. 1-24 (2001)
- (8) K. Matsuoka, M. Ohya and M. Kawamoto: A neural net for blind separation of nonstationary signals; *Neural Networks*, Vol. 8, No. 3, pp. 411-419 (1995)
- (9) R. Boscolo, H. Pan, V. P. Roychowdhury: Independent component analysis based on nonparametric density estimation; *IEEE Transactions on Neural Networks*, Vol. 15, No. 1, pp. 55-65 (2004)
- (10) M. R. Álvarez, F. Rojas, C. G. Puntonet, J. Ortega, F. Theis and E. W. Lang: A geometric ICA procedure based on a lattice of the observation space; *ICA2003*, pp. 1101-1106 (2003)