

## 独立成分分析による残響環境下での混合音声の分離

半田 晶寛<sup>†</sup> 堤 憲亮<sup>†</sup> Leandro Di Persia<sup>‡</sup> 柳田 益造<sup>†</sup>

<sup>†</sup>同志社大学工学部 〒610-0394 京都府京田辺市多々羅都谷 1-3

<sup>‡</sup> Universidad Nacional de Entre Rios Argentina

E-mail: <sup>†</sup> {dte0733, bta0014}@mail4.doshisha.ac.jp, myanagid@mail.doshisha.ac.jp <sup>‡</sup> ldipersia@hotmail.com

### 1. はじめに

現在の音声認識システムは、マイクに接近した位置からの音声に対してはそれなりに高い認識精度を有する一方で、マイクから離れた位置からの音声に対しては、周囲の雑音や部屋の残響の影響を受けて認識精度は著しく低下してしまう[1]。実環境で音声認識を行うためには音声認識システムに入力される音声に対して何らかの処理を行うことが考えられる。その一つがブラインド音源分離(BSS: Blind Source Separation)である。

本稿では防音室やリビングルームにおいて、様々な条件を変えて音を収録し、その収録音に対してICAを用いた周波数領域BSSを行い、音声認識システムJuliusでの認識率の改善でその動作を評価する。

### 2. 研究の目的

ICAを用いたBSSは、線形混合に対しては有効な手法であるといえるが、残響のある実環境において効果を挙げるには到っていない。

そこで本研究では2音源に対する周波数領域BSSの実環境での有効性検証を目的として、残響時間の異なる収録場所、音源(スピーカ)と観測点(マイク)との位置関係、目的音源と雑音のパワー比、雑音の種類、を変えて収録した音に対して、JADE[2, 3]の結果を初期値に用いたFastICA[4]による周波数領域BSSを行い、音声認識率の向上効果を検証する。

### 3. 分離実験

収録場所は防音室(490×400×290(H)cm<sup>3</sup>)及びリビングルーム(380×230×218(H)cm<sup>3</sup>)の2ヶ所である。音を収録した環境(防音室)を図1に示す。正面壁からスピーカまでの距離が100cm、横壁からマイクaまでの距離が200cmである。リビングルームにおけるマイク・スピーカの配置は防音室と同じとした。また防音室に関しては標準状態に加えて、壁の1面(図1中の横壁)に反射板を設置した環境でも収録を行った。残響時間は防音室の標準状態、反射板1面、リビングに関してそれぞれ白色ノイズで130ms, 150ms, 440msである。信号の観測には、

<sup>†</sup>Separation of Mixed Speech in Reverberant Environment Using Independent Component Analysis

Akihiro HANDA<sup>†</sup> Kensuke TSUTSUMI<sup>†</sup>  
Leandro Di Persia<sup>‡</sup> and Masuzo YANAGIDA<sup>†</sup>  
<sup>†</sup> Faculty of Engineering, Doshisha University  
<sup>‡</sup> Universidad Nacional de Entre Rios, Argentina

素子間隔5cmの2素子アレイ(小野測器M1-1233)を使用した。音源信号としては、防音室で録音し、Juliusでの認識率が100%の男女各20フレーズ、計40フレーズの音声を使用した。雑音はコンピュータの冷却ファンノイズ、テレビCMノイズ、目的音声の話者から見て異性によるスピーチノイズの3種類である。スピーカの出力比(設定S/N)は0dBと6dBの2種類とした。

<防音室>

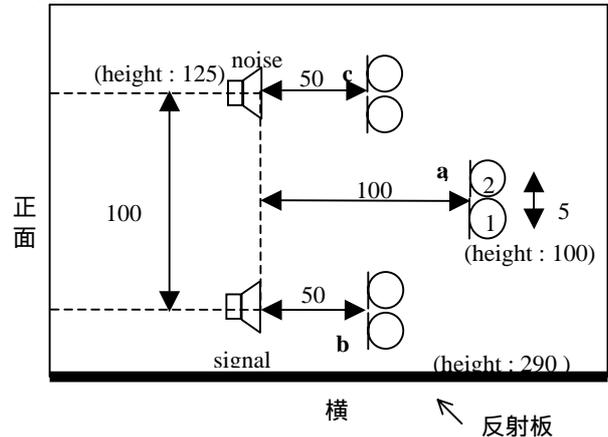


図1: 収録環境 (標準状態ではなし)

### 4. S/N の計算

受音段階及び分離後のS/Nは、暫定的に以下の式により計算している。

$$SNR = 10 \log_{10} \frac{\sum_n s_n^2}{\sum_n (s_n - \alpha y_n)^2} [\text{dB}]$$

ここで、 $s_n$ は音源信号、 $y_n$ は評価対象信号の各標本値である。また $\alpha$ は振巾に関する未定係数であり、分母が最小となるような値を最急降下法で求めた。ただし、上式を用いて実環境で収録した音のS/Nの計算をするためには、音源信号と、評価対象信号のサンプル点を正確に同期させるか、もしくは正確な補間値を求める必要があり、これらは共に非常に困難である。また振幅についても用いて調整しているが、信号と雑音の間に相関があると正確な値が計算できないなどの問題がある。このようなことから、本稿でのS/Nの値はあくまでも参考値である。

## 5. 実験結果と考察

図2は受音段階と分離音の認識率, 図3はS/Nの比較である. この図での値は, 3種類の雑音について平均した値である. 図2, 3の横軸のc0は図1におけるcの位置にあるマイクロフォンで, スピーカ出力比が0dBの音を収録した状態であること表示. 他も同様である.

収録環境別の認識率の改善に関するt検定, 収録環境別のS/Nに関するt検定の結果を表1の上下段に示す. \*, \*\*はそれぞれ $p < 0.05$ ,  $p < 0.01$ で有意水準であることを示す. なお, 表内の数値は分離前後の認識率[%], S/N[dB]の改善であり, 上述した有意差が確認できれば, 数値の前部のセルに「\*」記号を記入した.

### 5.1 認識率の改善

標準状態の防音室においては表中の有意差から, 認識率が改善されたと言える. 反射板を用いた場合には, b6以外で悪くとも有意水準5%の有意差が得られたことから, 概ね改善されたと言える. しかし, リビングルームでは有意水準5%でも全ての状態で有意差がなく, 認識率に関して改善の効果があるとは言えない.

またマイクロフォンの位置とスピーカ出力比の違いによる比較では, 防音室に関しては, 表1よりほぼすべての状態で認識率に対して改善の効果が見られると言える. ただし, b6のようにもともとS/Nが良く, 認識率の高いものに対しては, 改善が見られるとは言えない.

### 5.2 S/Nの改善

受音段階のS/Nが悪い時ほど改善の割合が大きくなっている. 表2のt検定の結果で示されているように, 受音段階のS/Nが高い場合はどの収録環境に関しても有意水準5%でも有意差がなく, S/Nが改善されるとは言えない.

表 1: 収録環境毎の認識率の向上(上段)とS/Nの改善(下段)に関するt検定

| [%]<br>[dB] | 防音室 |         | 防音室<br>(反射板1面) |         | リビング<br>ルーム |         |
|-------------|-----|---------|----------------|---------|-------------|---------|
| c0          | **  | 0 8     | *              | 0 4     | **          | 0 0     |
|             | **  | 0.4 2.8 | **             | 0.4 2.0 | **          | 0.2 0.4 |
| a0          | **  | 0 16    | *              | 1 8)    | **          | 0 1     |
|             | **  | 2.0 3.9 | **             | 1.9 3.4 | **          | 0.4 0.5 |
| a6          | **  | 6 22    | *              | 6 17    | **          | 1 1     |
|             | **  | 3.1 4.1 | **             | 3.0 3.7 | **          | 0.5 0.6 |
| b0          | **  | 8 18    | **             | 7 17    | **          | 2 1     |
|             | **  | 4.3 4.9 | **             | 4.2 4.5 | **          | 0.7 0.8 |
| b6          | *   | 16 29   |                | 15 23   |             | 6 5     |
|             |     | 5.0 4.8 |                | 4.9 4.4 |             | 0.8 0.8 |

## 6. まとめ

収録環境の違いによる比較では, 残響時間が150ms程度までなら音声認識率, S/N共に改善の効果があることが確認された. 一方で, リビングルー

ムにおける分離処理に関しては, S/Nが改善されても, 認識率は改善されなかった.

またS/Nと認識率の関係についても, 全体としてはS/Nが改善された時には認識率も改善されており, 一定の相関が見られる. しかし今後, より確かなS/Nの評価法を開発し, 再び検討を行う必要がある.

## 7. 終わりに

様々に条件を変えた環境で収録した音に, 独立成分分析による周波数領域BSSの処理を行い, その音声認識への有効性を検証した.

その結果, 残響時間が150ms程度までなら効果が得られることが確認された. ただし, 残響時間としては調査した対象が150msから440msに飛んでいるので, その間の値の残響時間での音声に対する認識率に関しても調べる必要がある.

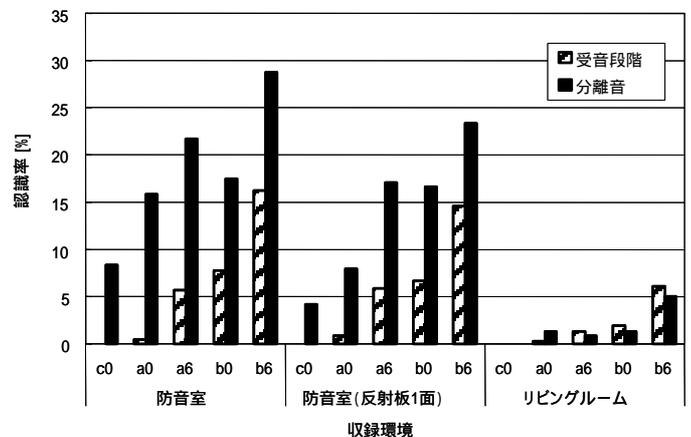


図 2: 収録環境毎の認識率の比較

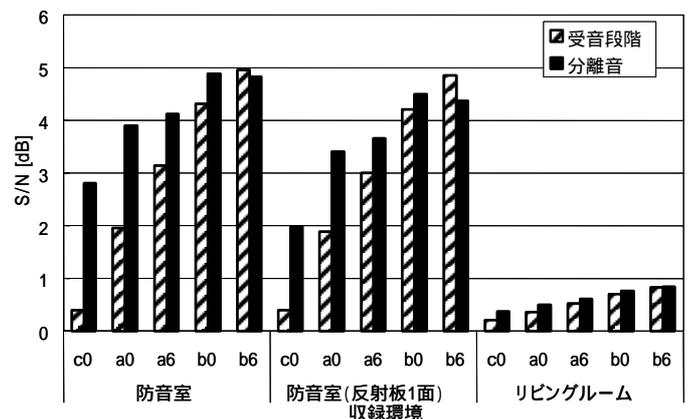


図 3: 収録環境毎のS/Nの比較

### 文献

- [1]中村哲, 実音響環境に頑健な音声認識を目指して, 電子情報通信学会 技術報告, SP2002-12, pp.31-36, 2002.
- [2]甘利俊一, 村田昇. SGC ライブラリ 18:独立成分分析 - 多変量データ解析の新しい方法, サインズ社, 東京, 2002.
- [3]J.F.Cardoso, and A.Souloumiac. Blind beamforming for non Gaussian signals, IEEE Proceeding-F, vol. 140, no 6, pp.362-370, December 1993.
- [4]横田康成. 信号処理 ~第5部 独立成分解析~, <http://www.ykt.info.gifu-u.ac.jp/dsp/sp5.pdf>