

知能情報メディア実験
「複数話者の同時発話音声からの個別音声の抽出」最終レポート

情報科学類, 3 年, 2 クラス, 学籍番号 : 201311403 山岡 洸瑛, Yamaoka Kouei

2015/8/5

目次

1	はじめに	2
2	課題 1: 瞬時混合 ICA の実装	2
2.1	ステップ 1	2
2.2	ステップ 2	2
2.3	ステップ 3	4
2.4	ステップ 4	4
2.5	ステップ 5	5
3	課題 2-1: 畳み込み混合の分析	5
3.1	0.	5
3.2	1.	7
3.3	2.	8
4	課題 2-2: IVA による音源分離	10
4.1	1.	10
4.2	2.	13
5	課題 2-3: 分離性能の客観・主観評価	13
5.1	1.	13
5.2	2.	13
6	終わりに	15

1 はじめに

知能情報メディア実験「複数話者の同時発話音声からの個別音声の抽出」の最終レポートとして、与えられた各課題の結果と考察をまとめる。このレポート内において原則として、スカラーは x 、ベクトルは太字で \mathbf{x} 、行列は大文字の太字で \mathbf{X} と表し、行列の i 行 j 列成分を X_{ij} と表すとする。その他 f と m に関する行列を $\mathbf{X}[f, m]$ とすることもある。このとき、 $f = f_i, m = m_j$ の場合 $\mathbf{X}[f_i, m_j]$ と表す。また単位行列を \mathbf{I} とする。ソースコードについては、単純な計算のみのものや、与えられた式をそのまま用いたものなどは示さない。独自に作成したものや、説目に必要なもののみ示す。また作成した音源については TA の杉本さんに確認していただいた。

2 課題 1: 瞬時混合 ICA の実装

2.1 ステップ 1

2 つの音声サンプルを x_1, x_2 とし、式 (1) の混合行列を用いて (2) の式で混合した。行列の形から 1 行 2 列成分と 2 行 1 列成分の絶対値が小さいため、出力 y_1, y_2 について、 y_1 は x_1 の、 y_2 は x_2 の音声が他方に比べて振幅が大きくなることが予想される。

2.1.1 結果

作成した混合を再生すると、確かに混合されていた。また、予想通り y_1 においては x_1 の音が x_2 よりも大きく聞こえ、 y_2 についても x_2 の方が大きく聞こえた。混合の結果を図 1 に示す。左側がサンプル音源、右側が混合した音声である。この図からも y_1 は x_1 の、 y_2 は x_2 の影響が大きいことが読み取れる。

2.1.2 考察

結果より、音声をベクトル化した数値で扱うことができるという事がわかった。また、音声ベクトルに対して演算を行うことが可能という事もわかった。使用した混合行列の値を式 (3) のように変更して混合したところ、予想通り y_1 では x_2 の、 y_2 では x_1 の音が大きくなった。これより、瞬時混合は式 (4) のように単純な線形結合で表され、かつ各々の係数によって音声の大きさが決まるという事が確認できた。これは音源の数が増えても \mathbf{H} のサイズが変わるだけであり、本質的には変わらないと言える。

$$\mathbf{H} = \begin{pmatrix} 0.9 & 0.4 \\ -0.3 & 1 \end{pmatrix} \quad (1)$$

$$\begin{pmatrix} y_1 \\ y_2 \end{pmatrix} = \mathbf{H} * \begin{pmatrix} x_1 \\ x_2 \end{pmatrix} \quad (2)$$

$$\mathbf{H} = \begin{pmatrix} 0.1 & 0.9 \\ 0.9 & 0.1 \end{pmatrix} \quad (3)$$

$$\begin{cases} y_1 = H_{11} * x_1 + H_{12} * x_2 \\ y_2 = H_{21} * x_1 + H_{22} * x_2 \end{cases} \quad (4)$$

2.2 ステップ 2

混合行列 \mathbf{H} の逆行列を用いて分離信号を作成し、それが元の音源と一致することを確認する。分離は混合音声 y_1, y_2 を用いて式 (5) のように行った。分離した音声を s_1, s_2 とした。

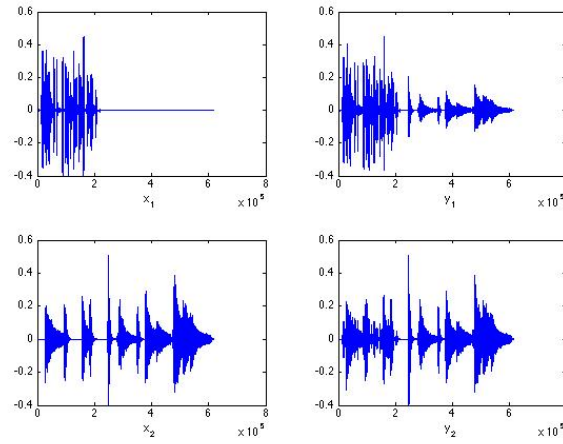


図 1: 瞬時混合による混合

2.2.1 結果

実際に分離を行い，元の音源と一致することが確認できた．分離した音声プロットしたものを図 2 に示す．この図が図 1 の右側とほぼ一致することから視覚的にも分離できていることが確認できる．しかし同時に僅かに変わっている部分があることから知覚できない程度の誤差が生じてしまっていることも読み取れる．

2.2.2 考察

この実験より，行列によって混合された音源は，その逆行列によって分離できるということが確認できた．これより，混合行列が不明であってもその逆行列の近似的な値がわかれば分離できると予想される．生じる誤差は知覚できない程度のものなので無視してよいと言える．この誤差は計算機で計算しているために生じる丸め誤差などの影響だと思われる．

$$\begin{pmatrix} s_1 \\ s_2 \end{pmatrix} = H^{-1} * \begin{pmatrix} y_1 \\ y_2 \end{pmatrix} \quad (5)$$

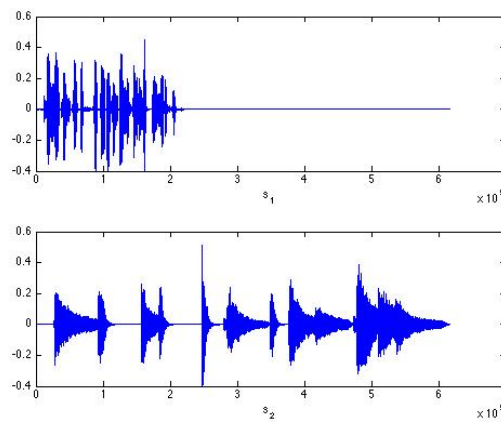


図 2: 逆行列計算による分離

2.3 ステップ 3

瞬時混合 ICA を用いて分離行列 W を学習する．インフォマックスという手法を用いて，式 (6) と式 (7) のような勾配法の更新を繰り返すことによって学習する．ここで活性化関数としては式 (8) を用いた．音源が 2 つの場合， I は 2×2 , $i = 1, 2$ であり， N は音源のフレーム数である． W の初期値は式 (9) とし，ステップサイズ $\mu = 0.1$ ，活性化関数で用いる適当な正の値 $\gamma = 0.5$ とした．

2.3.1 結果

学習した結果推定された分離行列 W を得た．確認のため $W * H$ を計算したところ，1 行 2 列成分と 2 行 1 列成分が共に 0 に近くなったため，近似的な値は求まっていると言える．ただし 1 行 1 列成分，2 行 2 列成分が大きく，単位行列からは大きく離れていた，また，学習は勾配法で行ったため，収束条件を W のすべての要素について更新前の値と更新後の値の差の絶対値が 10^{-4} 以下となること，とし反復計算を行った結果，今回のパラメータの場合 179 回の反復で終了した．

2.3.2 考察

インフォマックスという手法を用いると分離行列が求まることが確認できた．この分離行列で実際に分離できるのかどうかはわからないが，結果で確認したように正しい逆行列の近似にはなっているようなので，分離できることが期待される．ただし，単位行列に近づかなかったことや， W の要素が非常に大きいことから，元の音源よりも音量が大きくなっていると予想され，これがスケーリングの問題であると言える．

$$\begin{pmatrix} y_1 \\ y_2 \end{pmatrix} = W * \begin{pmatrix} x_1 \\ x_2 \end{pmatrix} \quad (6)$$

$$W \leftarrow W + \mu \left(I - \frac{1}{N} \sum_{n=0}^{N-1} \psi(y_i[n]) y_i[n]^t \right) W \quad (7)$$

$$\psi(x) = (\tanh(\gamma x_1), \dots, \tanh(\gamma x_i))^t \quad (8)$$

$$W = \begin{pmatrix} 0.3 & 0.7 \\ 0.5 & 1 \end{pmatrix} \quad (9)$$

2.4 ステップ 4

ステップ 1 で作成した混合音声と，ステップ 3 で作成した分離行列を用いて実際に分離を行った．この時スケーリング解決を行った．また分離性能を上げるために数回分離行列学習時のパラメータを調整した．スケーリング解決及び分離は式 (10) で行った． $y_i[n]$, W はステップ 3 で得た値で， $Z[n]$ はスケーリング解決を行った分離信号である． $\text{diag}()$ は引数を対角要素とする対角行列を表す．

2.4.1 結果

分離し，スケーリング解決を行った音声を聞くと元の音源と全く同じように聞こえた．サンプル音源 1 の裏でサンプル音源 2 が少し聞こえる，などという事はなく，完全に分離されていた．分離した音声をプロットし，図 3 に示す．この図からも分離していることが分かる．スケーリングの解決もできているようで，音源 1 の方はほぼ完璧である．しかし，音源 2 の方は分離できているものの振幅が小さくなってしまっている．

2.4.2 考察

ステップ3で考察したように分離行列を推定することで分離することができた．パラメータの調整は必要だが，瞬時混合の場合，正確な逆行列でなくても，非常に綺麗に分離できるといえる．スケーリングに問題は残っているが，実際に音を聞くと気になるほどの違いはわからないため，大きな問題にはならないと言える．ここまでの計算は音源の数を2つとして行なってきたが，音源の数が変わったとしても同じように計算できることは使用した計算式からわかる．これより瞬時混合ならば，混合音声を分離することができると言える．

$$\mathbf{Z}[n] = \mathbf{W}^{-1} * \text{diag}(\mathbf{y}_i[n]) \quad (10)$$

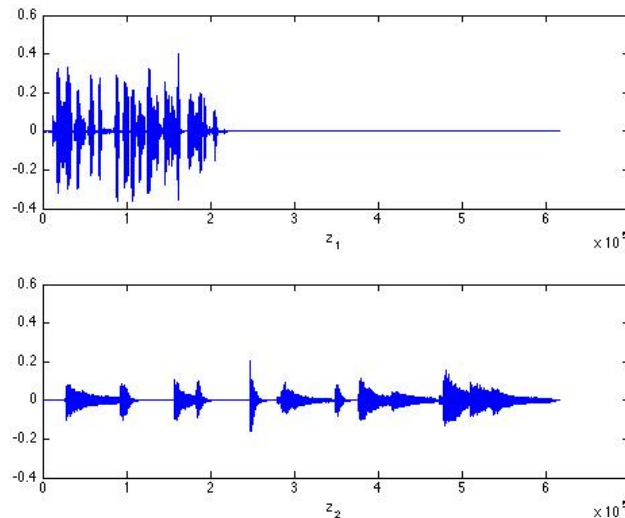


図 3: インフォマックスによる分離

2.5 ステップ5

混合行列が与えられていない未知の瞬時混合波形に対してICAを行い音源分離を行え，という課題であったが未知の瞬時混合音声取得出来なかったため，混合行列 H の値を変えて混合音声を作成し， H が未知であるとして音源分離を行った．混合行列を作成した翌週に分離を行うことで混合行列を未知であるとした．実際に混合行列の値は覚えていなかったため，未知であったとしても良いと考える．

2.5.1 結果

ステップ4同様に分離できた．しかしステップ4と違い混合行列の情報がないとしたためパラメータの調節に時間がかかってしまった．これは混合行列が未知であるため，分離行列との積を求めることで確認する，という方法がとれなかったため，目安がなかったことが原因である．時間がかかってしまったものの分離できたためステップ4の考察で述べたように瞬時混合音声であれば分離できると言える．

3 課題2-1: 畳み込み混合の分析

3.1 0.

課題2-1で使用する畳み込み信号を自分で作成するという追加課題を行った．2音源2マイク環境下での畳み込みを行った．インパルス応答は与えられたものを使い，音源には課題1で使用したサンプル音源をそのまま使用

した．

音源の長さを N 、インパルス応答の長さを R として $K = N + R - 1$ とする．まず畳み込みを周波数領域で行うために音源とインパルス応答を離散フーリエ変換 (以下 DFT とする) する．DFT の方法としては高速フーリエ変換 (以下 FFT とする) を用いる．音源 1,2 を K 点 FFT したものを s_1, s_2 、音源 i からマイク j へのインパルス応答を K 点 FFT したものを imp_{ij} とする．また周波数領域でのインデックスを f とすると、畳み込みは式 (11) で表される．畳み込み混合信号を x_i として、 x_i を逆フーリエ変換 (以下 iFFT とする) することで時間領域の畳み込み混合信号を得る．

3.1.1 結果

畳み込み混合信号を得た．実際に聞いてみると音声に少しディレイがかかっているように聞こえ、畳み込みに成功していることが確認できた．畳み込みに使用したインパルス応答を図 4 に示す．ほとんど同じに見えるがよく見ると音源 1 からの第 1 インパルスの方が音源 2 からの第 1 インパルスよりも大きいなどの違いがある．また、畳み込み混合信号を図 5 に示す．左が元の音源で、右が畳み込み混合信号である．畳み込んでいるのかどうかはよく読み取れないが、振幅が $\frac{1}{10}$ 程度になっていることが分かる．これはインパルス応答の振幅が小さいことが原因だと思われる．これよりインパルス応答が反映されていると言えるため、畳み込みに成功していると言える．

3.1.2 考察

周波数領域における畳み込みを実装したところ非常に高速に動くことが確認できた．時間領域の畳み込みを MATLAB 関数の `conv()` で行ったところメモリが足りずに計算出来なかったことから、周波数領域での畳み込みの威力が分かる．また K 点 FFT をする理由は、畳み込みによって $R - 1$ だけ音声伸びるからである．これは実際に短い信号で畳み込みを行なってみると確認できる．以上より以降の課題で使用する畳み込み混合信号が作成できた．

$$\begin{pmatrix} x_1[f] \\ x_2[f] \end{pmatrix} = \begin{pmatrix} imp_{11}[f] & imp_{21}[f] \\ imp_{12}[f] & imp_{22}[f] \end{pmatrix} * \begin{pmatrix} s_1[f] \\ s_2[f] \end{pmatrix} \quad (11)$$

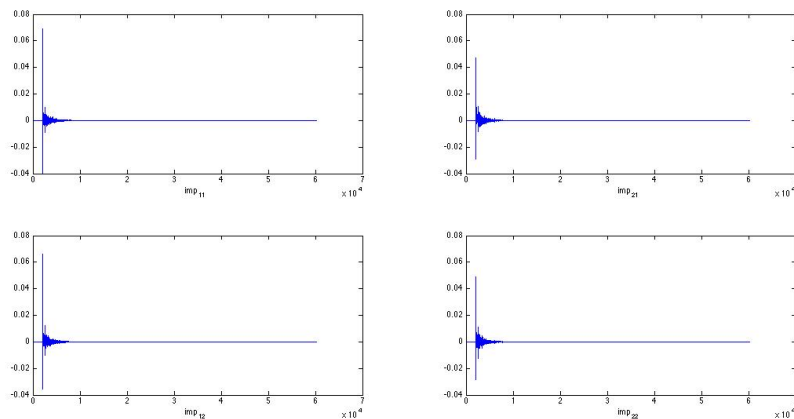


図 4: インパルス応答

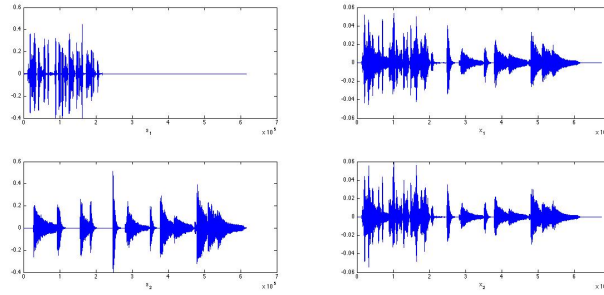


図 5: 畳み込み混合信号

3.2 1.

作成した畳み込み混合信号を短時間フーリエ変換 (以下 STFT とする) を用いて分析する．また逆 STFT(以下 iSTFT とする) によって元の波形に戻ることを確認する．

3.2.1 実装

フレーム長 p , フレームシフトを $ol = p * \frac{1}{o}$, 総フレーム数を al , 窓関数を w とする．これらを用いて p フレーム分 w をかけて取り出し FFT, ol だけずらして同様の計算を行う．これを繰り返すことで STFT が実装できる．iSTFT では切り出したフレームを iFFT して加算, ol だけずらして同様の計算を行うことで実装できる．しかし, STFT の時にフレームシフトや窓関数によってオーバーラップが生まれてしまい, その場合このままでは iSTFT は正しく行えていない．例えば $ol = \frac{1}{2}$ で窓関数がハミング窓の場合, オーバーラップはなくこのままで良いが, 窓関数が箱形窓の場合, オーバーラップは 2 となり, 前後 $\frac{p}{2}$ 以外は振幅が 2 倍になってしまう．これを解決するためにカウント用配列を用意する．これは切り出した回数をカウントする配列で, 信号の第 k インデックスを n 回切り出した, という情報を格納する．iSTFT 時にこの配列を用いて割ることでオーバーラップ問題を解決できる．これに関する図を図 6 に示した．以上より STFT, iSTFT が実装できた．ソースコードをリスト 1, 2 に示す．ただし, X は STFT した信号, Y は iSTFT した信号で, $count$ がカウント用配列である．なお, MATLAB の都合上 $ifft()$ を用いてもごく小さな値が虚数として残ってしまうことがあるため, $real()$ 関数を用いて虚数を消してある．

3.2.2 結果

STFT を行い iSTFT によって元の音源に戻ることが確認できた．音源信号をプロットした図を図 7 に示す．左側が元の音源 s_1, s_2 で右側が STFT, iSTFT した音源である．これより左右の信号がほぼ一致していることから, 視覚的にも成功していることが確認できた．

3.2.3 考察

図 7 の左右を見比べると僅かに異なっている部分が稀にある．これは計算機で計算しているために丸め誤差などの僅かな誤差が生じた結果だと考えられる．左右の音源を聞いてみると, その差は知覚出来なかったため, この程度の誤差であれば無視しても構わないと言える．また今回の実装で STFT, iSTFT をそれぞれ関数化したため, 今後 STFT 分析は関数 1 つで行えるようになった．

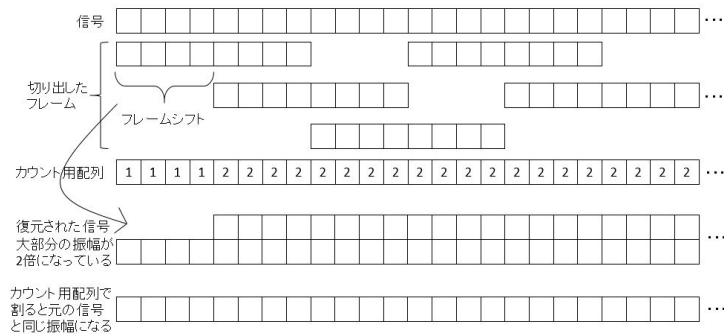


図 6: オーバーラップ問題の解決

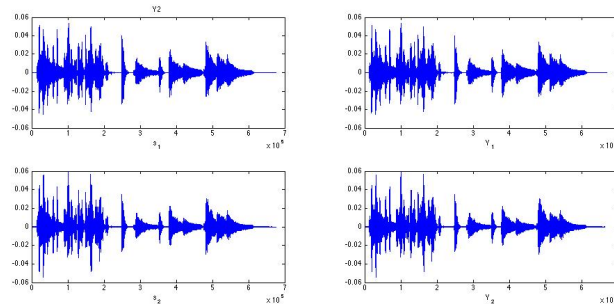


図 7: 元の音源と STFT, iSTFT 後の音源

リスト 1: stft.m

```

1 for k = 1:al
2     % 窓をかけ、fft
3     X(k,:) = s(1 + ol*(k-1) : p + ol*(k-1)) .* w';
4     X(k,:) = fft(X(k,:));
5     % 今のカウント用配列のインデックスに1加算
6     count(1 + ol*(k-1) : p + ol*(k-1)) = count(1 + ol*(k-1) : p + ol*(k-1)) + 1;
7 end

```

リスト 2: istft.m

```

1 for k = 1:al
2     % ifftし窓をとる
3     Y(1 + ol*(k-1) : ol*(k-1) + p) = Y(1 + ol*(k-1) : ol*(k-1) + p) + real(ifft(X(k,:))) ./ w';
4 end
5 Y = Y ./ count;

```

3.3 2.

与えられたインパルス応答を DFT し周波数毎の分離行列を作成する．また STFT した混合信号を分離行列を用いて分離し，それを iSTFT することで時間領域の分離信号を得る．各変数は 1. と同じとし，サンプル音源 1,2 を STFT したものを X_1, X_2 とする．ただし i 行に第 i フレームが格納されている．また，音源 i からマイク j へのインパルス応答を p 点 FFT したものを imp_{ij} とする．分離した信号を X_i と同じサイズの行列 Y_1, Y_2 に格納す

ると式 (12) で分離することができる．ただし周波数領域のインデックスは f で表している．計算した Y_i を iSTFT することで時間領域の分離信号を得ることができる．STFT, iSTFT は 1. で作成した関数を用いるが，iSTFT において 1 部変更を加えた．STFT した信号に分離行列をかけているため窓関数で割る必要がなくなったので，割らないように変更をした．

3.3.1 結果

分離に成功した．プロットした図を図 8 に示す．これより分離されていることが視覚的にも確認できる．しかし本来無音であるべき時間であっても僅かに振幅が存在することも確認できる．これは分離の過程で生じてしまったノイズである．これは分離した音声を聞いても聞き取れる程度には大きく．無視することはできない．

3.3.2 考察

時間周波数領域での分離に成功したため，ノイズを考慮しなければ今後 DFT の代わりに STFT を使うことができる．問題のノイズが生じた理由について考察する．フーリエ変換は切り出した範囲が周期関数であることを仮定して計算しているため，実装した STFT においては窓関数をかけることによってフレームの両端が綺麗につながるようにしている．しかし，切り出したフレームは実際には周期関数ではなく，その続きは今のフレームとは異なる信号であるため，周期関数を仮定して DFT すると実際とは異なる周波数領域の情報が得られてしまう．これを各フレームについて行なっているためすべてのフレームについてこのことが言える．STFT 分析した信号をそのまま iSTFT するのであれば各フレーム毎に FFT し iFFT したただけなので一致する．しかし，今回は STFT 分析した信号に分離行列をかけるという演算を行なっているため，誤差が生じたまま分離を行なっていることになる．そのため分離信号には誤差があり，それを iSTFT しても本来の時間領域の分離信号は得られず，誤差，つまりノイズのついた信号が得られてしまう．これが理由であると考えられる．

これを解決する方法は考えつかなかった．この問題は各フレーム毎の誤差が原因であるため，フレーム長を長くすればするほどノイズは減少すると予想される．実際にフレーム長を長くして実行した結果を図 9 に示す．外形自体が変わってしまっている部分があるが，ノイズ部分は減少していることが読み取れる．音声のノイズも減少していると感じた．これより解決はできていないが軽減はできたと言える．しかし STFT には不確定性原理があり，フレーム長を長くすると周波数分解能は良いが時間分解能は低くなってしまふ．フレーム数を最大にすればこの問題は最も軽減されるが，それは DFT であり，STFT する意味がなくなってしまう．これらよりこの問題の根本的な解決方法が見つからなかった．

$$\begin{pmatrix} Y_1[f] \\ Y_2[f] \end{pmatrix} = \begin{pmatrix} imp_{11} & imp_{21} \\ imp_{12} & imp_{22} \end{pmatrix}^{-1} * \begin{pmatrix} X_1[f] \\ X_2[f] \end{pmatrix} \quad (12)$$

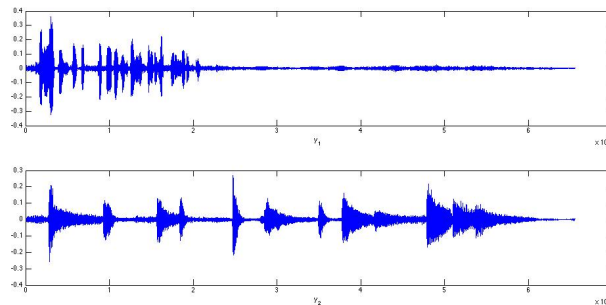


図 8: 時間領域の分離信号

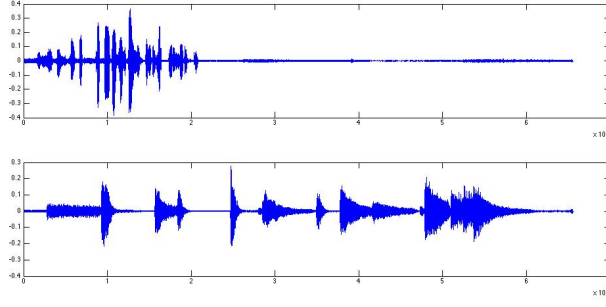


図 9: フレーム長を 4 倍にした場合

4 課題 2-2: IVA による音源分離

4.1 1.

IVA を実装し分離を行う．最大周波数を $\frac{L}{2}$ とする．分離行列 W の学習には課題 1 ステップ 3 と同様にインフォマックスを用いる．その時の分離信号 Y_i の更新は式 (13) であり，勾配法の更新は式 (14) である．ただしこのときの活性関数は式 (15) である．活性関数で使用している $normY_i[m]$ は式 (16) で計算される．またスケーリングの解決は式 (17) で行う．ここで X_i, Y_i は第 i 音源の時間周波数領域表現で，周波数インデックス f , 第 m フレームについて $X_i[f, m]$ などと表す．音源数が増えた場合も同様の式で計算できる．

$$\begin{pmatrix} Y_1[f, m] \\ Y_2[f, m] \end{pmatrix} = W * \begin{pmatrix} X_1[f, m] \\ X_2[f, m] \end{pmatrix} \quad (13)$$

$$W[f] \leftarrow W[f] + \mu \left(I - \frac{1}{M} \sum_{m=0}^M \psi(Y_i[f, m]) Y_i[f, m]^H \right) W[f] \quad (14)$$

$$\psi(X[f, m]) = \gamma \left(\frac{Y_1[f, m]}{\|normY_1[m]\|_2}, \dots, \frac{Y_J[f, m]}{\|normY_J[m]\|_2} \right)^t \quad (15)$$

$$\|normY_j[m]\|_2 = \sqrt{\sum_{f=0}^{\frac{L}{2}} |Y_j[f, m]|^2} \quad (16)$$

$$Z = W[f]^{-1} * diag([Y_1[f, m], Y_2[f, m]]) \quad (17)$$

4.1.1 実装

式 (14) ~ 式 (17) を実装する．しかしこのまま実装すると f, m についてループ計算が必要になるため for 文を多用することになり計算時間が非常に長くなってしまふ．そのため式変形と MATLAB 関数を利用して計算速度を速くする．

まず式 (13) については音源毎に f, m についての 2 重ループを計算することになるが，MATLAB 関数を使うことで複数の音源を同時にかつ f に関するループをなくして計算することができる．つまり音源数 $\times f \times m$ 回の計算を m 回の計算まで減らすことができる．まず X_i, Y_i について，行が m , 列が f となっているが，これを 2 行 m 列 f ページの 3 次元配列に変更し， i 行目に X_i, Y_i を格納する．この 3 次元配列を新たに X, Y とする．また W を 2 行 2 列 f ページの 3 次元配列とする．これを第 f ページについて式に表すと式 (18) となる．この行列計算はこのままだと実行できない． $W[f]$ を対角要素とする対角行列を用意し， X を転置することで計算はできるが，メモリを大量に消費するため，ここは for 文を使うことにする．第 m 列に着目すると式 (19) となり，for 文で実装できる．

$$\begin{pmatrix} Y_1[f, 1] & \cdots & Y_1[f, M] \\ Y_2[f, 1] & \cdots & Y_2[f, M] \end{pmatrix} = \begin{pmatrix} W_{11}[f] & W_{12}[f] \\ W_{21}[f] & W_{22}[f] \end{pmatrix} * \begin{pmatrix} X_1[f, 1] & \cdots & X_1[f, M] \\ X_2[f, 1] & \cdots & X_2[f, M] \end{pmatrix} \quad (18)$$

$$\begin{aligned} \begin{pmatrix} Y_1[f, m] \\ Y_2[f, m] \end{pmatrix} &= \begin{pmatrix} W_{11}[f] & W_{12}[f] \\ W_{21}[f] & W_{22}[f] \end{pmatrix} * \begin{pmatrix} X_1[f, m] \\ X_2[f, m] \end{pmatrix} \\ &= \begin{pmatrix} W_{11}[f] * X_1[f, m] + W_{12}[f] * X_2[f, m] \\ W_{21}[f] * X_1[f, m] + W_{22}[f] * X_2[f, m] \end{pmatrix} \end{aligned} \quad (19)$$

これは行列ならば簡単に計算できる式だが、3次元配列の場合 MATLAB では計算することができない。そのためこの式を MATLAB で計算できるように式変形すると式 (20) となり、これは式 (19) と一致する。そしてこの式は 3次元配列のまま MATLAB で計算できる。これによって f に関する for 文の撤去に成功した。sum(), *, は MATLAB 表記で、 X^t は転置を表す。ソースコードはリスト 3 に示す。

$$\begin{aligned} &sum \left(\begin{pmatrix} W_{11}[f] & W_{12}[f] \\ W_{21}[f] & W_{22}[f] \end{pmatrix} * \begin{pmatrix} X_1[f, m] & X_1[f, m] \\ X_2[f, m] & X_2[f, m] \end{pmatrix}^t, 2 \right) \\ &= sum \left(\begin{pmatrix} W_{11}[f] * X_1[f, m] & W_{12}[f] * X_2[f, m] \\ W_{21}[f] * X_1[f, m] & W_{22}[f] * X_2[f, m] \end{pmatrix}, 2 \right) \\ &= \begin{pmatrix} W_{11}[f] * X_1[f, m] + W_{12}[f] * X_2[f, m] \\ W_{21}[f] * X_1[f, m] + W_{22}[f] * X_2[f, m] \end{pmatrix} \end{aligned} \quad (20)$$

次に式 (14) ~ (16) について考える。まず、式 (20) で計算した Y からそれぞれの信号を取り出し、m 行 f 列の Y_1, Y_2 とする。これを用いると、要素毎に 2 乗し、sum() を使い、sqrt() を使うことで式 (21) となり、式 (16) を得る。これより式 (16) には f, m どちらの for 文もいらず、音源毎の計算をすれば良い。

$$(\|normY_j[1]\|_2, \dots, \|normY_j[M]\|_2) = \left(\sqrt{\sum_{f=0}^{\frac{L}{2}} |Y_j[f, 1]|^2}, \dots, \sqrt{\sum_{f=0}^{\frac{L}{2}} |Y_j[f, M]|^2} \right) \quad (21)$$

式 (15), (16) を用いて式 (14) を音源を同時にかつ f に関する for 文のみで計算できるようにする。つまり計算回数が音源数 * f * m 回であったのを f 回に減らす。まず式 (14) 中の $\psi(Y_i[f, m])Y_i[f, m]^H$ について式変形をする。以下の式変形において、MATLAB 関数の sum(), 複素共役をとる conj(), 内積を計算する dot() を用いる。なお以下では音源数を J としている。

$$\begin{aligned} \psi(Y_i[f, m])Y_i[f, m]^H &= \gamma \begin{pmatrix} \frac{Y_1[f, m]}{normY_1[m]} \\ \vdots \\ \frac{Y_J[f, m]}{normY_J[m]} \end{pmatrix} \begin{pmatrix} conj(Y_1[f, m]), & \cdots, & conj(Y_J[f, m]) \end{pmatrix} \\ &= \gamma \begin{pmatrix} \frac{Y_1[f, m] * conj(Y_1[f, m])}{normY_1[m]}, & \cdots, & \frac{Y_1[f, m] * conj(Y_J[f, m])}{normY_1[m]} \\ \vdots & \ddots & \vdots \\ \frac{Y_J[f, m] * conj(Y_1[f, m])}{normY_J[m]}, & \cdots, & \frac{Y_J[f, m] * conj(Y_J[f, m])}{normY_J[m]} \end{pmatrix} \\ \sum_{m=1}^M \psi(Y_i[f, m])Y_i[f, m]^H &= \gamma \begin{pmatrix} \sum_{m=1}^M \frac{Y_1[f, m] * conj(Y_1[f, m])}{normY_1[m]}, & \cdots, & \sum_{m=1}^M \frac{Y_1[f, m] * conj(Y_J[f, m])}{normY_1[m]} \\ \vdots & \ddots & \vdots \\ \sum_{m=1}^M \frac{Y_J[f, m] * conj(Y_1[f, m])}{normY_J[m]}, & \cdots, & \sum_{m=1}^M \frac{Y_J[f, m] * conj(Y_J[f, m])}{normY_J[m]} \end{pmatrix} \\ &= \gamma \begin{pmatrix} \frac{dot(Y_1[f], conj(Y_1[f]))}{sum(normY_1)}, & \cdots, & \frac{dot(Y_1[f], conj(Y_J[f]))}{sum(normY_1)} \\ \vdots & \ddots & \vdots \\ \frac{dot(Y_J[f], conj(Y_1[f]))}{sum(normY_J)}, & \cdots, & \frac{dot(Y_J[f], conj(Y_J[f]))}{sum(normY_J)} \end{pmatrix} \end{aligned} \quad (22)$$

式 (22) において $dot(Y_j[f], conj(Y_j[f]))$ は、行列 $Y_j[f, m]$ の第 f 列を取り出したベクトル $Y_j[f]$ とその複素共役の内積であるからスカラーであり、 $sum(normY_1)$ はベクトルの要素の和なのでこちらもスカラーである。従っ

て $J * J$ の行列である．この行列を $\chi[f]$ とおくと，式 (14) は以下のように表される．これを更に式変形する．

$$\begin{aligned} W[f] &\leftarrow W[f] + \mu \left(I - \frac{1}{M} \chi[f] \right) W[f] \\ &= \left(I + \mu \left(I - \frac{1}{M} \chi[f] \right) \right) W[f] \end{aligned}$$

ここで改めて $I + \mu \left(I - \frac{1}{M} \chi[f] \right)$ を $\chi[f]$ とおくと，最終的に式 (14) は式 (23) と表すことができる．

$$W[f] \leftarrow \chi[f] * W[f] \quad (23)$$

この計算は 3 次元配列であるため行えないが，第 i 列について 1 列ずつ行えば，式 (13) と同じ方法で実行できる．以上より全ての音源を同時に f のループのみで実行できるようになった．ソースコードはリスト 3 に示す．式 (22) を計算する関数 `calcm.m` はリスト 4 に示す．

これより分離行列の学習が実装できた．これを用いると，混合信号を STFT し分離行列を学習，学習した分離行列を用いて式 (13) のように分離，式 (17) でスケーリング解決を行い iSTFT することで時間領域の分離信号を得ることができる．STFT, iSTFT については実装済みなので省略する．

リスト 3: `bunri.m`

```
1 % 反復計算
2 for k=1:loop
3     % 式(13)の計算
4     for l = 1:fi
5         Y(:,l,:) = sum(W .* repmat(permute(X(:,l,:),[2,1,3]), 2, 1), 2);
6     end
7     % Yの計算終了 -----
8
9     % Wの更新 -----
10    % 準備
11    Y1 = squeeze(Y(1,:,:)).';
12    Y2 = squeeze(Y(2,:,:)).';
13    % 式(16)の計算
14    normY1 = sqrt(sum(abs(Y1).^2));
15    normY2 = sqrt(sum(abs(Y2).^2));
16
17    % X[f]の計算
18    for i = 1:L2
19        kai(:,i) = calcm(Y1(i,:), Y2(i,:), normY1, normY2, fi, mu, gamma);
20    end
21
22    % 式(23)の計算
23    W(:,1,:) = sum(kai .* repmat(permute(W(:,1,:),[2,1,3]), 2,1), 2);
24    W(:,2,:) = sum(kai .* repmat(permute(W(:,2,:),[2,1,3]), 2,1), 2);
25
26    % Wの更新終了 -----
27 end
28
29 % 最後の式(13)の計算
30 for l = 1:fi
31     Y(:,l,:) = sum(W .* repmat(permute(X(:,l,:),[2,1,3]), 2, 1), 2);
32 end
```

リスト 4: `calcm.m`

```
1 function kai = calcm(Y1, Y2, normY1, normY2, fi, mu, gamma)
2
3     ty1 = gamma * conj(Y1);
4     ty2 = gamma * conj(Y2);
5     kai = [sum((Y1.*ty1)./normY1), sum((Y1.*ty2)./normY1); ...
6            sum((Y2.*ty1)./normY2), sum((Y2.*ty2)./normY2)];
7     kai = eye(2) + mu * (eye(2) - (1/fi) * tmpY);
```

4.1.2 結果

分離に成功した． W の初期値はシード値 509 で生成した乱数とした． μ, γ は分離精度を上げるために変更しつつ調整した．分離音声 Y_1, Y_2 をプロットした図を図 10 に示す．この図からわかる通り，分離には成功したが完全に分離することは出来なかった．分離，というよりは片方の音源を強調していると言ったほうがいいかもしれない．どちらか一方の音が大きく聞こえ，他方の音が後ろで小さく聞こえている．課題 1 ステップ 4 のようには分離出来なかった．分離精度の向上が今後の課題と言える．

4.1.3 考察

まず計算回数を減らすためにした式変形によってどの程度計算速度が上がったのかを確認する．残った for 文は f に関するものと m に関するもの両方あるため，for 文の計算速度はトレードオフの関係になってしまっている．今回はフレーム長を 1024 として実行した．この時，式 (13) は 0.2 秒短縮され，式 (14) については 6 秒程度短縮された．フレーム長を長くすると式 (13) は更に短縮されるが，式 (14) は長くなってしまう．これよりフレーム長にかかわらず約 6 秒程度の短縮になったと言える．ただし混合信号の長さが短くなると短縮される時間も短くなる．これは f, m どちらも小さくなるため for 文自体にかかる時間が減るためである．このことから当然ではあるが，今回の短縮は混合信号が長いほど威力を発揮すると言える．

また特に式 (13) などはわかりやすくするため音源数が 2 個の場合で式変形しているが，音源数が J 個の場合も同じ方法で実装することができる．これより畳み込み混合信号であっても分離行列を学習し分離できるようになった．ここまでで音源数とマイクの数が適していれば実際の環境下で混合音声を分離できるようになったはずである．必要な情報は話者の数とマイクの数，話者の数以上の混合信号のみである．

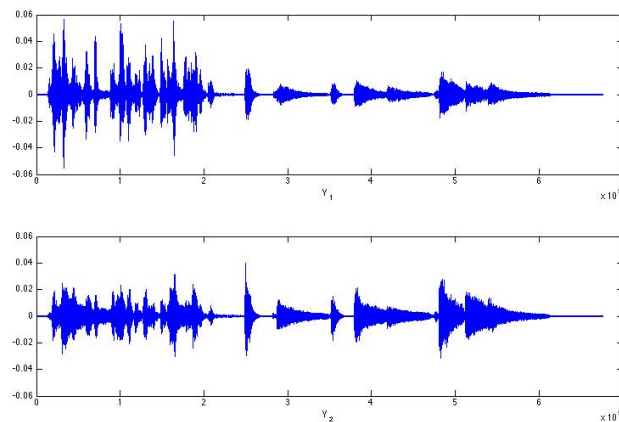


図 10: 分離信号

4.2 2.

この課題はやらない，という事なので行わなかった．配布するという関数も配布されなかったため実装出来なかった．

5 課題 2-3: 分離性能の客観・主観評価

5.1 1.

客観評価に使用する音源が配布されなかったため行わなかった．この課題もやらない，という事であった．

5.2 2.

実環境で録音した信号の分離を行う．録音した音声は 2 話者の場合と 3 話者の場合があり，その両方について分離を行った．分離には課題 2-2 で作成したプログラムをその儘使用した．3 話者の場合には同様の手法で音源数を 3 つに拡張したプログラムを使用した．分離した音源を用いて主観評価実験も行う．

5.2.1 結果

分離に成功した．ただし前の課題と同じくパラメータを調整したが，完全には分離できず1人が話している後ろで僅かに別の人間の話し声が残ってしまった．またノイズも大きく，それほど精度が高いとは言えないものとなった．しかし1人の話し声だけが強く聞こえるため成功と言えるだろう．2話者の場合でもノイズなどが大きかったにも関わらず，3話者のばあいはそのがさらに大きくなってしまい，精度は更に悪くなってしまった．混合信号によってはうまく分離できず，1人の声が少し大きくなったぐらいにしかないこともあり，3話者の分離は困難であった．例として2話者の場合と3話者の場合それぞれについて，混合信号と分離信号を図11，12に示す．どちらも左側が混合信号で右側が分離信号である．図の通り，ノイズが非常に大きくなってしまっているが，振幅が大きな位置が異なっていることから分離に成功していることが分かる．特に3話者の場合はほとんどノイズなのではないかと思うほどノイズが大きい，こちらも大きな振幅の位置で分離されていることが読み取れる．それでも分離精度が最も良いと感じたものを選んでプロットしたのである．このノイズ除去は重要な課題であると言える．

5.2.2 考察

今回の分離でノイズが大きくなってしまった一因に元の混合信号にもノイズが乗っていたという事が挙げられる．これに関してはこちらではどうしようもないため，改善は難しいだろう．改善方法としてバンドパスフィルタを通してみた．人間が通常会話で発する音の周波数は300～2000Hz程度の範囲に収まるため，それ以外の周波数をカットすることでノイズが減少するのではないか，という単純な考えから行なってみた．しかし実際にはノイズは減少せず，理由はわからなかったがノイズが大きくなってしまうこともあった．これより単純なバンドパスフィルタではノイズ除去は行えないことがわかった．ノイズの特性を調べてみるとはっきりしたことはわからなかったが，白色雑音も含まれているようで，これからも単純なフィルタでは除去できないことが分かる．

分離精度についてもパラメータを何度も調整したにも関わらず上がらないため，この分離方法ではそれほど精度を得られないのだと思われる．分離音声と混合音声を聞いていると，話者の声の大きさなど，個人差があるためそれによって精度が変わってくるように感じられた．例えば声が比較的大きい人と小さい人の2話者で分離すると，大きい人の方が分離精度が高く，小さい人の場合はそれほど高くなってしまっている．これより分離は混合信号に強く左右される，と言える．高精度の分離には限界があり，高精度の混合信号が必要なのではないだろうか．

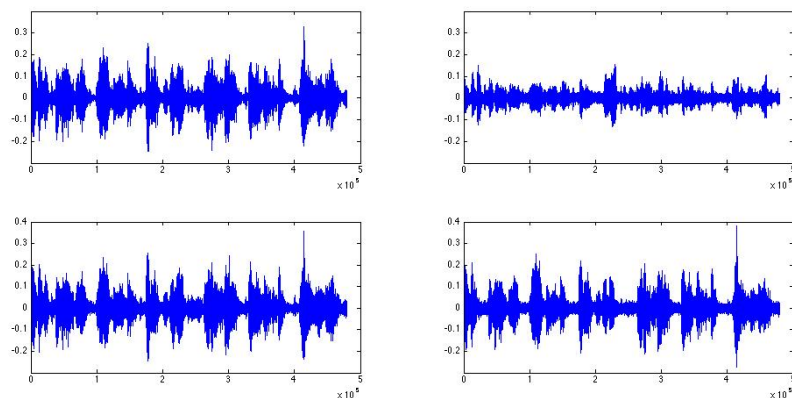


図 11: 2 話者の場合の分離

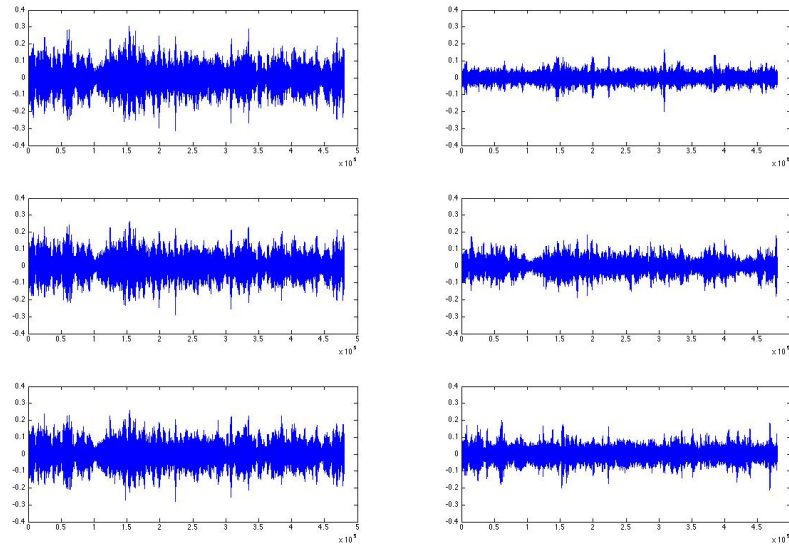


図 12: 3 話者の場合の分離

5.2.3 主観評価実験

分離した音源を用いて主観評価実験を行った。今回の実験は合計 5 人で行ったため、その 5 人の分離音声を集め、誰の音声かわからない状態でランダムに再生、それを音質、分離という 2 項目について 5 段階評価を行った。音質は分離にかかわらず、良いかどうか、分離は音質に関わらず、分離できているかできていないか、が評価基準であった。結果がまだわからないため、自分の分離精度がどの程度なのかはわからないが、中には自分のよりも明らかに精度よく分離できているものもあり、工夫次第でまだまだ精度をあげられることがわかった。

6 終わりに

この実験の結果、混合信号の分離を学ぶ過程で信号処理について理解が深まった。特に式で扱っただけのフーリエ変換を実際に行ったことで、利用方法を学ぶことができた。また混合信号分離についての知識と技術も手に入れることができた。分離精度の向上や、ノイズの除去など課題は多いが、この実験の目標は達成できたと言えるだろう。