

プログラミング基礎演習レポート 2017-2

長谷川禎彦

注意点

レポートは以下の2つを提出すること。

- プログラムのソースコード (Python ソース (.py) または Jupyter notebook (.ipynb))
- レポート本体 (doc, docx, pdf, odt). 英語でも良い.

以下の点に留意せよ. なお, その他の注意点は第一回レポートに準ずる.

- 提出〆切り: 2018/2/13 の 23:59
- 提出方法: 作成したソースファイル (ファイル名は自由. 一つのプログラムを複数ファイルに分割しても良い) とレポート本体 (doc, docx, pdf, odt 等) を「学籍番号.zip」としてまとめる. 作成した zip ファイルをホームページの提出フォームから提出する (「ファイル」と書かれた所から, ソースとレポート本体の zip ファイルを添付する). その際, 課題の選択を「レポート2」とすること (フォームからは 12/20 以降に提出可能となる).
- なお, 全ての問題において大枠として題意を満たしていれば, 方法・内容とも自由に変更して良い. 問題の拡張や手法の改良を行って良い. なお, Python の Pandas, Numpy, Scipy, Matplotlib 等 (これらに限らず) ライブラリは基本的に自由に用いよ. ただし, ICA それ自体を行うライブラリを用いてはならない (例えば, scikit-learn は利用不可).

カクテルパーティ効果とは, たくさんの人が異なる会話をしている時でも, 自分が興味のある人の会話は聞き取ることができるというよく知られた効果である. このように, 人間には複数の音源から必要な情報だけを抽出する機構が備わっている. これを模したアルゴリズムが独立成分分析 (ICA: Independent Component Analysis) である. 本レポートでは, ICA を Python を用いて実装してみる.

観測データを $\mathbf{x} \in \mathbb{R}^n$ で表し, 信号源データ (未知) を $\mathbf{y} \in \mathbb{R}^n$ とする (どちらも列ベクトル). 実際のデータは $\mathbf{X} = \{\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_N\}$ のように N 個のデータベクトルによって与えられ, それに対応して信号源は $\mathbf{Y} = \{\mathbf{y}_1, \mathbf{y}_2, \dots, \mathbf{y}_N\}$ とする. ここでは, 観測データ \mathbf{x} は平均が $\mathbf{0} \in \mathbb{R}^n$ になるように調整している (つまり $\mathbf{x} \leftarrow \mathbf{x} - \mathbb{E}[\mathbf{x}]$). なお, $\mathbb{E}[\mathbf{x}]$ は平均 (期待値) を表し, $\mathbb{E}[\mathbf{x}] = \frac{1}{N} \sum_{i=1}^N \mathbf{x}_i$ とする.

観測データは, 信号源の重ね合わせのため, 以下の線形関係があると仮定する

$$\mathbf{x} = \mathbf{A}\mathbf{y}. \quad (1)$$

ここで $\mathbf{A} \in \mathbb{R}^{n \times n}$ の行列である. これより, もし $\mathbf{W} = \mathbf{A}^{-1}$ が分かれば, 信号源 \mathbf{y} は

$$\mathbf{y} = \mathbf{W}\mathbf{x}, \quad (2)$$

によって復元可能である. つまり, ICA は行列 \mathbf{W} を求める問題に他ならない.

中心極限定理より, 任意の独立な確率変数の和は正規分布に収束する. 中心極限定理は無限の和の極限において成立する. しかし, たった二つの確率変数 q_1 と q_2 を足した場合でも, 足した後の $q_1 + q_2$ の方が, q_1 や q_2 よりも正規分布に近い分布となる場合が多い. 逆に言えば, 正規分布から遠い分

布を持つデータは、データが足し合わさっていないことを意味する。この性質を用いると、 \mathbf{y} ができるだけ正規分布から遠くなるように \mathbf{W} を決めれば、信号源 \mathbf{y} が計算できそうである。実は、尖度 (kurtosis) と呼ばれる統計量は、正規分布の時 0 になり、正規分布から乖離すると 0 から大きくずれる (尖度は正負の値をとれる)。平均が0の確率変数 $y \in \mathbb{R}$ の尖度は以下で定義される。

$$\text{kurtosis}[y] = \mathbb{E}[y^4] - 3\mathbb{E}[y^2]^2, \quad (3)$$

つまり、 $|\text{kurtosis}(y)|$ を最大化するように、 \mathbf{W} を求めればICAが実現できることがわかる。

現実にはこれだけではうまくいかない。以下に述べる白色化と呼ばれる操作を、ICAの前段階に挟む必要がある。今、観測データの共分散行列を

$$\mathbf{\Sigma} = \mathbb{E}[\mathbf{x}\mathbf{x}^T], \quad (4)$$

とする (\mathbf{x} は平均が0のように取っていることに注意せよ)。共分散行列 $\mathbf{\Sigma}$ は以下のように対角化させることができる

$$\mathbf{\Sigma} = \mathbf{E}\mathbf{D}\mathbf{E}^T, \quad (5)$$

ここで $\mathbf{D} = \text{diag}(\lambda_1, \lambda_2, \dots, \lambda_n)$, λ_i は $\mathbf{\Sigma}$ の固有値である。 \mathbf{E} は直交行列になっている。この時以下のように \mathbf{V} を導入する

$$\mathbf{V} = \mathbf{E}\mathbf{D}^{-1/2}\mathbf{E}^T, \quad (6)$$

ここで、 $\mathbf{D}^{-1/2} = \text{diag}\left(\frac{1}{\sqrt{\lambda_1}}, \frac{1}{\sqrt{\lambda_2}}, \dots, \frac{1}{\sqrt{\lambda_n}}\right)$ である。新しい確率変数 \mathbf{z} を以下で導入する

$$\mathbf{z} = \mathbf{V}\mathbf{x}, \quad (7)$$

この時 \mathbf{z} の共分散行列は単位行列となる、つまり $\mathbb{E}[\mathbf{z}\mathbf{z}^T] = \mathbf{I}$ 。共分散を単位行列にする操作を白色化という (\mathbf{z} は白色化されている)。ICAでは観測データ \mathbf{x} をまず白色化し、式(7)を用いて \mathbf{z} を計算する。 \mathbf{z} を用いると、式(2)は以下のようになる

$$\mathbf{y} = \mathbf{W}\mathbf{z}, \quad (8)$$

(この \mathbf{W} と式(2)の \mathbf{W} は当然一般には異なる)。式(8)において、信号源の一つに注目し、それを $y \in \mathbb{R}$ とする。この時、式(8)より

$$y = \mathbf{w}^T \mathbf{z}, \quad (9)$$

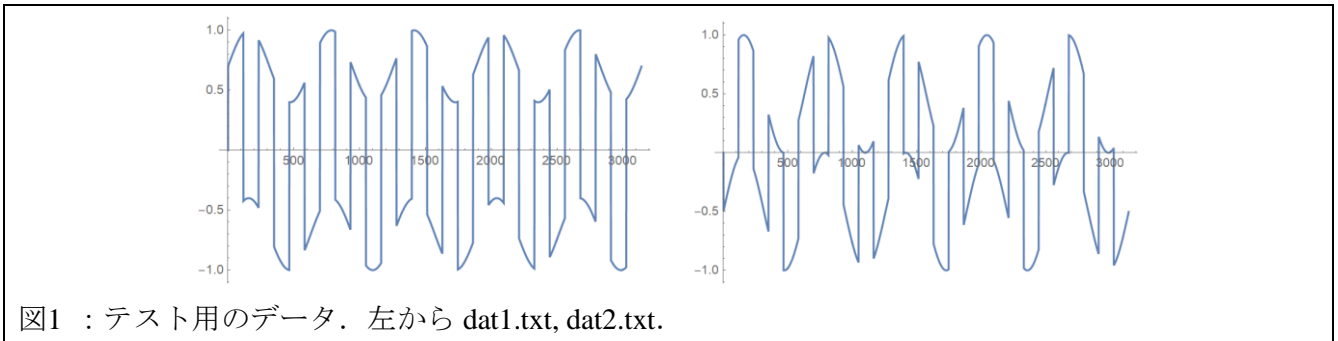
ここで $\mathbf{w} \in \mathbb{R}^n$ は列ベクトルである。 \mathbf{w}^T を縦に積んだものが \mathbf{W} になっている。計算の詳細は省くが、 y の尖度の絶対値を最大化させるには、以下の繰り返しアルゴリズムを適用することで可能である (詳細は参考文献参照)。

1. \mathbf{w} に対して初期値を適当に選ぶ。正規化する $\mathbf{w} \leftarrow \frac{\mathbf{w}}{\|\mathbf{w}\|}$ 。
2. $\mathbf{w} \leftarrow \mathbb{E}[\mathbf{z}(\mathbf{w}^T \mathbf{z})^3] - 3\mathbf{w}$ を計算する
3. $\mathbf{w} \leftarrow \frac{\mathbf{w}}{\|\mathbf{w}\|}$ によって正規化する
4. 収束していなければ、2に戻る。収束していれば終了する。

このアルゴリズムを収束するまで実行すると、情報源 y の尖度の絶対値が最大化され、信号源の一つが求まる。異なる信号源を求めるには、1.において異なる初期値 \mathbf{w} でアルゴリズムを何度か実行すればよい (これが一番簡単)。ほかの方法としては、直交化を各ステップで使う方法もある (詳細は参考文献参照)。

1. 課題（必須課題）

ホームページの `dat1.txt` と `dat2.txt` は 2 つの信号源を異なる比率で混合したデータである (`report2_kadai1_data1.zip`, 図 1 参照). Python で ICA を実装し, `dat1.txt` と `dat2.txt` に ICA を適用し, 2 つの信号源を同定せよ.



2. 課題（自由課題）

ホームページに音声ファイルが入った圧縮ファイルが3つある. その内訳は以下のとおりである. これらのデータに作成した ICA を適用し, 音源を分離せよ. なお, 分離後の `wav` ファイルの提出は, 3 つのうちどれか一つでよい (`report2_kadai2_data3` を提出する場合は `mp3` 等で圧縮してサイズを小さくすること. `mp3` 圧縮は <https://online-audio-converter.com/ja/>などで可能).

- `report2_kadai2_data1.zip`
 - 二人の話者が同時に話したものを, (仮想的に) 異なる場所で録音したデータが入っている. `wav` 形式, 16bit, 8000Hz, モノラル
- `report2_kadai2_data2.zip`
 - 同じ声の 3 人の話者が同時に話したものを, (仮想的に) 異なる場所で録音したデータが入っている. `wav` 形式, 16bit, 8000Hz, モノラル
- `report2_kadai2_data3.zip`
 - 二つの曲を同時に演奏したものを, (仮想的に) 異なる場所で録音したデータが入っている. `wav` 形式, 16bit, 44100Hz, モノラル

Python での `wav` ファイルの読み込み・書き込みは `scipy.io.wavfile.read` や `scipy.io.wavfile.write` を用いると簡単にできる (他のライブラリでももちろん良い).

3. 課題（自由課題）

音声だけではなく, 画像に対しても ICA を適用することが可能である. ホームページに `image1.png` と `image2.png` が入ったファイルがある (`report2_kadai3_data1.zip`). これは, 二つの画像の異なる比率の重ね合わせである (図 2 参照). これに作成した ICA を適用し, 画像を分離せよ. なお画像はグレースケールの `png` フォーマットである.

4. 課題（自由課題）

課題で与えられたデータ以外 (の面白いデータ) に対して, 作成した ICA を適用してみよ.



図2 : テスト用のデータ. 左から image1.png, image2.png.

参考文献

- Hyvärinen, Aapo, Juha Karhunen, and Erkki Oja. Independent component analysis. Vol. 46. John Wiley & Sons, 2004.
- [上の翻訳] Hyvärinen, Aapo, Juha Karhunen, and Erkki Oja 著, 根本幾, 川勝真喜翻訳, 独立成分分析～信号解析の新しい世界～, 東京電機大学出版局