



EMORY
UNIVERSITY
SCHOOL OF
MEDICINE

Department of
Biomedical Informatics

BMI 500: Introduction to Biomedical Informatics

Lecture 11. An Introduction to Blind Source Separation and Independent Component Analysis

Reza Sameni

Department of Biomedical Informatics, Emory University, Atlanta, GA, USA

30th October, 2020

Overview

- 1 Background & Motivation
- 2 Classical Independent Component Analysis Algorithms
- 3 Iterative BSS Methods
- 4 Algebraic BSS Methods
- 5 Further Reading & Advanced Topics
- 6 Appendix: BSS Lab Session Notes

Outline

- 1 Background & Motivation
- 2 Classical Independent Component Analysis Algorithms
- 3 Iterative BSS Methods
- 4 Algebraic BSS Methods
- 5 Further Reading & Advanced Topics
- 6 Appendix: BSS Lab Session Notes

Introduction

- **Blind source separation (BSS):** the problem of recovering random variables or random processes from an unknown linear or nonlinear mixture of such sources, using minimal assumptions on the sources and mixtures:

$$\mathbf{x}(t) = \mathbf{f}(\mathbf{s}(t), t) \quad (1)$$

- Although the problem seems unsolvable, various evidences in biology and nature prove it to be possible
- The problem was first studied in signal processing in the 1980s
- In this lecture, a short and broad introduction is presented on BSS and one of its most popular solutions, known as **independent component analysis (ICA)**

Evidence of BSS in nature: the cocktail party effect

- **Cocktail party effect or selective attention:** the ability to focus attention on a single speaker in a crowd
- From the signal processing perspective, people can understand a speaker even in very low quality environments (even negative SNR)
- Various properties of speech are used in selective attention



Evidence of BSS in nature: the cocktail party effect

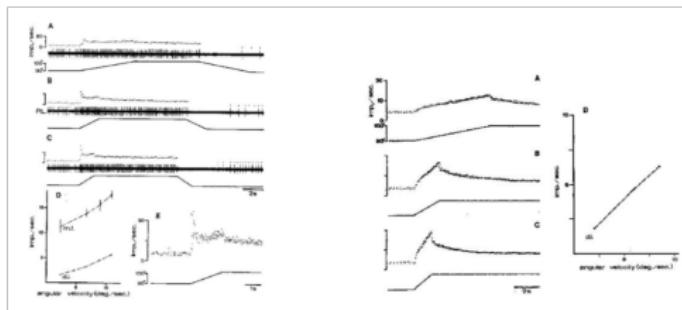
(continued)

Various factors of speech influence selective attention and intelligibility:

- Signal to noise ratio (SNR)
- Non-stationarity and sparsity of speech (time domain)
- Speaker voice tone (frequency domain)
- Facial gestures and lip reading (multimodality)
- Speech language, context, sentence structure, etc. (priors)
- Physiological and anatomic factors (two ears and their shapes improve selective listening)

Evidence of BSS in nature: motion decoding in vertebrates

Observation by Jean-Pierre Roll
(Roll, 1981): The human brain
is able to learn the position
(joint stretch) and velocity of
body extremities from mixtures
of these variables received from
the nervous system.



This observation inspired BSS research in signal processing:

$$\begin{aligned} f_1(t) &= a_{11}p(t) + a_{12}v(t) & p(t) &: \text{joint stretch} \\ f_2(t) &= a_{21}p(t) + a_{22}v(t) & v(t) &: \text{joint velocity} \end{aligned} \tag{2}$$

History of BSS in signal processing

A legendary article: J. Hérault, C. Jutten, and B. Ans. Détection de grandeurs primitives dans un message composite par une architecture de calcul neuromimétique en apprentissage non supervisé. In Actes du Xème colloque GRETSI, pages 1017–1022, Nice, France, 20–24 May 1985.

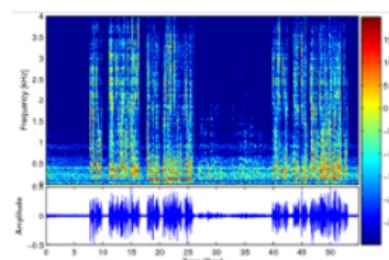
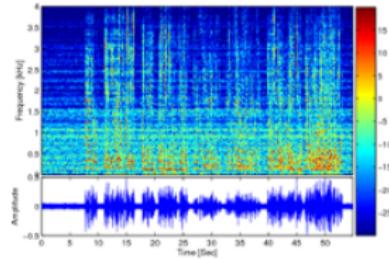
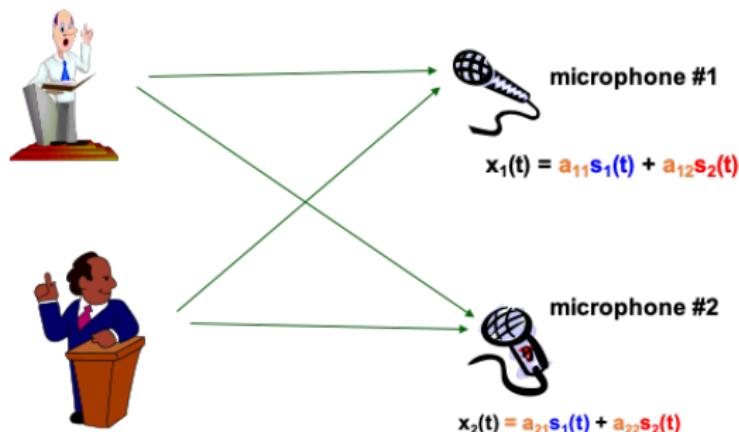
History of BSS in signal processing

A legendary article: J. Hérault, C. Jutten, and B. Ans. Détection de grandeurs primitives dans un message composite par une architecture de calcul neuromimétique en apprentissage non supervisé. In Actes du Xème colloque GRETSI, pages 1017–1022, Nice, France, 20–24 May 1985.



An analog source separation machine (Courtesy of Christian Jutten). See the following for notations and algorithm: C. Jutten and J. Hérault. "Blind separation of sources, part I: An adaptive algorithm based on neuromimetic architecture." *Signal processing* 24, no. 1 (1991): 1–10. APA

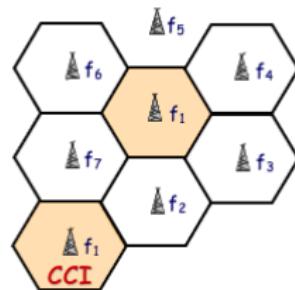
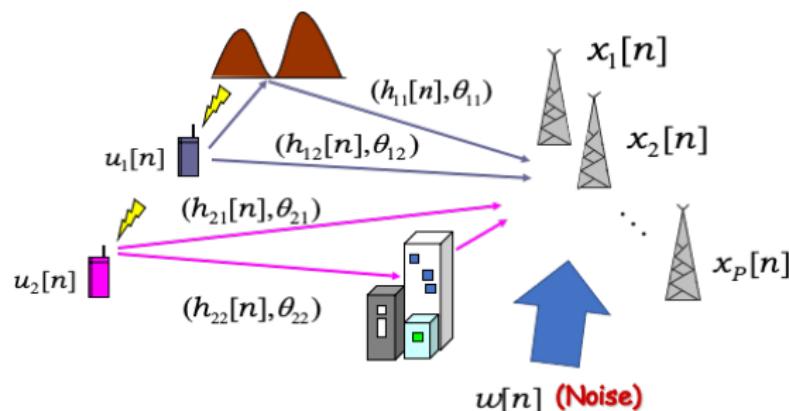
Applications: speech separation and enhancement



Sound propagates at 330m/s in air; so it takes 3ms for sound to propagate 1m. Therefore a **linear instantaneous mixture** is a good approximation of the audio scenario in some cases. For farther distances, **convulsive mixtures** are required.

Applications: multi-path telecommunication channels

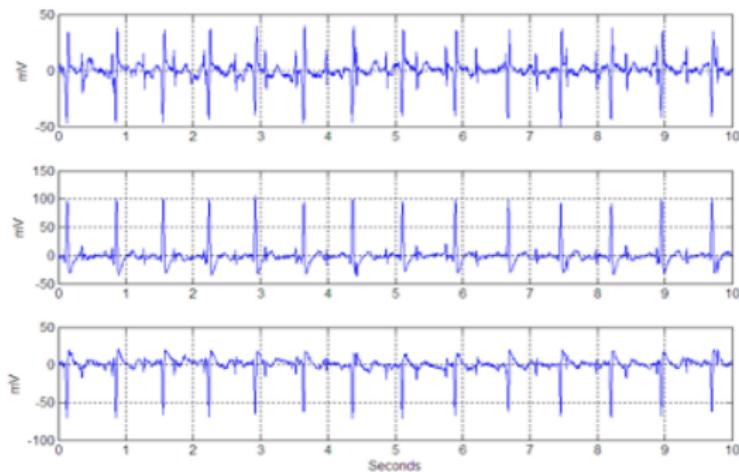
In wireless communications, the desired signal is degraded by interference and noise received from unwanted paths.



CCI: Co-channel Interference

Applications: noninvasive fetal electrocardiogram extraction

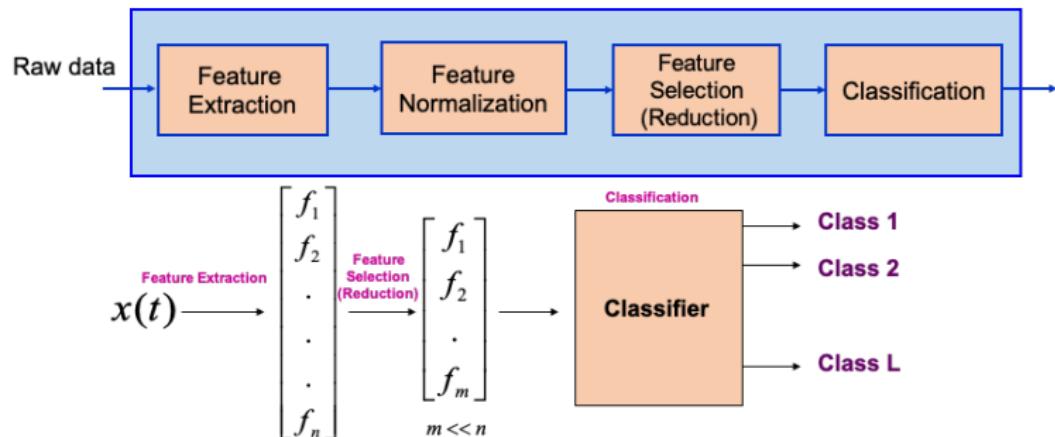
The noninvasive extraction of fetal ECG (fECG) from an array of electrodes placed on the maternal abdominal is a challenging biomedical problem



(Courtesy of Danilo Pani)

Applications: feature extraction and classification

In classification, **statistically independent** features are preferred over correlated or uncorrelated ones, as there is no redundancy between the features. How to make a set of features **independent**?



Outline

- 1 Background & Motivation
- 2 Classical Independent Component Analysis Algorithms
- 3 Iterative BSS Methods
- 4 Algebraic BSS Methods
- 5 Further Reading & Advanced Topics
- 6 Appendix: BSS Lab Session Notes

Mathematical problem definition

Data Model: In the most general case, blind source separation (BSS) can be mathematically formulated as follows:

$$\mathbf{x}(t) = \mathbf{f}(\mathbf{s}(t), t) \quad (3)$$

where

- t denotes time
- $\mathbf{x}(t)$ are N channels of observations (**known**)
- $\mathbf{s}(t)$ are M channels of sources (**unknown**)
- $\mathbf{f}(\cdot, t)$ is a (generally) nonlinear and time-variant transform (**unknown**)

Objective: to estimate $\mathbf{s}(t)$ from $\mathbf{x}(t)$, using minimal assumptions on the sources and mixtures (this is what the prefix “blind” stands for)

Blind source separation solutions

- The problem is apparently ill-posed and without any other prior assumptions there are no unique solutions
- Various cases of BSS with proven solutions include:
 - **Linear mixtures:**
 - Statistically independent and non-Gaussian sources
 - Signals with temporal correlation
 - Sparse sources (in various domains)
 - Temporally nonstationary sources
 - (Pseudo-)periodic sources
 - **Nonlinear mixtures:**
 - Post-nonlinear sources
 - Signals with temporal correlation
 - Sparse sources

Blind source separation solutions

(continued)

Herein, we study two approaches used for linear BSS without noise:

- **Iterative methods** (such as fastICA)
- **Algebraic methods** (such as AMUSE, JADE and Π CA)

which represent two major classes of BSS methods

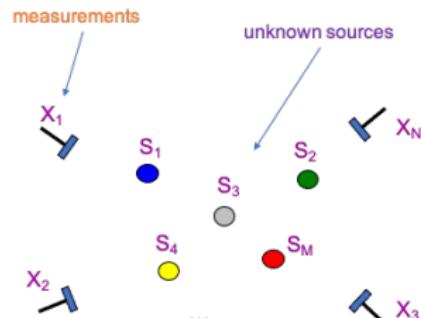
Outline

- 1 Background & Motivation
- 2 Classical Independent Component Analysis Algorithms
- 3 Iterative BSS Methods
- 4 Algebraic BSS Methods
- 5 Further Reading & Advanced Topics
- 6 Appendix: BSS Lab Session Notes

Linear mixtures of independent sources

Data Model: The most classical case of BSS with a proven solution is for linear mixtures of **statistically independent** and **non-Gaussian** sources (random variables of random processes): $x_i = a_{i1}s_1 + a_{i2}s_2 + \dots + a_{iM}s_M$

$$\begin{bmatrix} x_1 \\ \vdots \\ x_N \end{bmatrix} = \begin{bmatrix} a_{11} & a_{12} & \cdots & a_{1M} \\ \vdots & \ddots & & \vdots \\ a_{N1} & a_{N2} & \cdots & a_{NM} \end{bmatrix} \cdot \begin{bmatrix} s_1 \\ \vdots \\ s_M \end{bmatrix}$$



- In matrix form: $\mathbf{x} = \mathbf{As}$
- The noisy case: $\mathbf{x} = \mathbf{As} + \mathbf{n}$

Linear mixtures of independent sources

(continued)

- Different cases:

- $N = M$: Determined

- $N > M$:

Over-determined

- $N < M$:

Under-determined
(degenerate)

$$\begin{bmatrix} x_1 \\ \vdots \\ x_N \end{bmatrix} = \begin{bmatrix} a_{11} & a_{12} & \cdots & a_{1M} \\ \vdots & \ddots & & \vdots \\ a_{N1} & a_{N2} & \cdots & a_{NM} \end{bmatrix} \cdot \begin{bmatrix} s_1 \\ \vdots \\ s_M \end{bmatrix}$$

- The measurements and the sources can generally be **random processes** rather than **random variables**, i.e., $x(t)$ and $s(t)$.
- Independent component analysis (ICA)** is one of the possible solutions of BSS, which assumes mixtures of statistically independent sources
- How does it work?** we need to see how linear transforms alter the **topology** of random variables

Linear transforms of Gaussian random variables

Example: Suppose that s_1 and s_2 be independent Gaussian random variables with $s_1 \sim N(0, 1)$ and $s_2 \sim N(0, 2)$

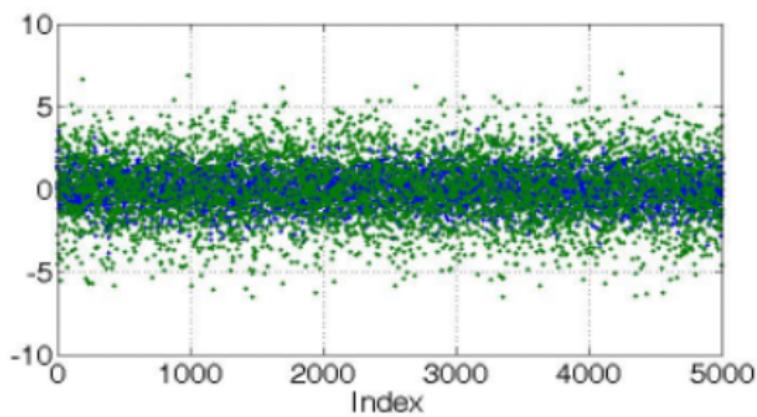


Figure: Sample plot

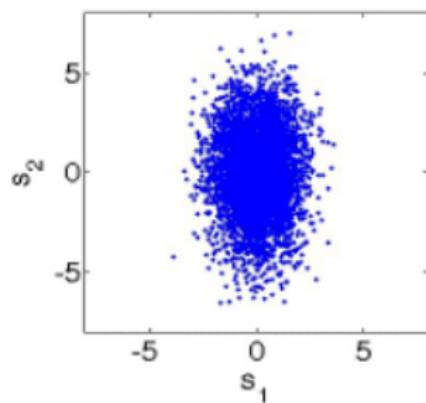


Figure: Scatter plot

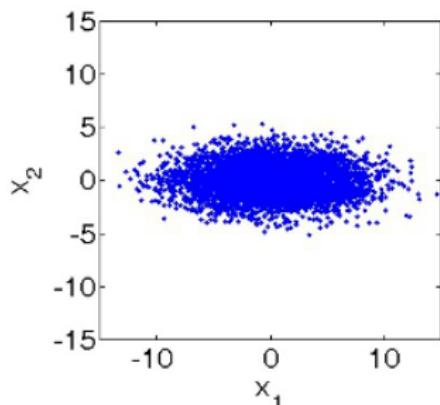
Linear transforms of Gaussian random variables

(continued)

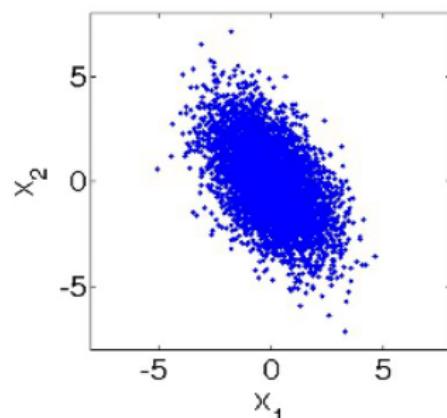
Example: Gaussian samples under transformations

$$\begin{bmatrix} x_1 \\ x_2 \end{bmatrix} = \begin{bmatrix} 4.0 & 0 \\ 0 & 0.75 \end{bmatrix} \begin{bmatrix} s_1 \\ s_2 \end{bmatrix}$$

$$\begin{bmatrix} x_1 \\ x_2 \end{bmatrix} = \begin{bmatrix} \cos(\frac{\pi}{6}) & -\sin(\frac{\pi}{6}) \\ \sin(\frac{\pi}{6}) & \cos(\frac{\pi}{6}) \end{bmatrix} \begin{bmatrix} s_1 \\ s_2 \end{bmatrix}$$



diagonal transforms stretch the distribution along the axes



orthonormal transforms rotate the distribution around the axes

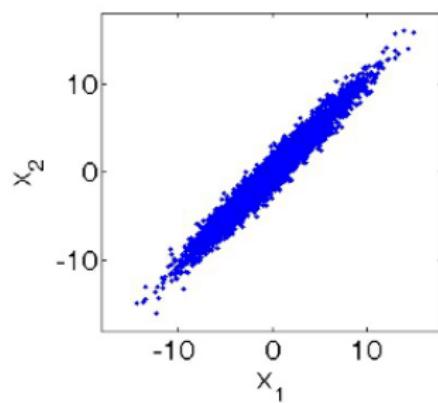
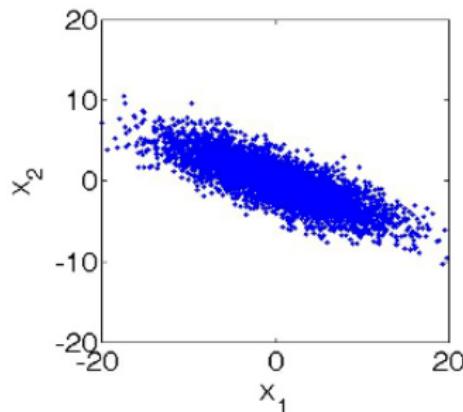
Linear transforms of Gaussian random variables

(continued)

Example: Gaussian samples under transformations

$$\begin{bmatrix} x_1 \\ x_2 \end{bmatrix} = \begin{bmatrix} 1.0 & -3.0 \\ -2.0 & 1.0 \end{bmatrix} \begin{bmatrix} s_1 \\ s_2 \end{bmatrix}$$

$$\begin{bmatrix} x_1 \\ x_2 \end{bmatrix} = \begin{bmatrix} 1.0 & 2.0 \\ 2.0 & 2.0 \end{bmatrix} \begin{bmatrix} s_1 \\ s_2 \end{bmatrix}$$



Arbitrary transforms (rotation + scaling), rotate and stretch the distribution

Singular value decomposition (SVD)

Reminders from linear algebra

Reminder (Strang, 1988)

Singular value decomposition of an arbitrary matrix:

$$\mathbf{A} = \mathbf{U}\Sigma\mathbf{V}^T \quad (4)$$

where \mathbf{U} and \mathbf{V} are **orthonormal**, and Σ is a **diagonal** matrix.

- Therefore, any **linear transform** of the form $\mathbf{x} = \mathbf{As}$ can be decomposed into a series of **rotations** (unitary transforms) and **scalings** (diagonal transforms).
- In an **inverse problem**, recovering (estimating) \mathbf{s} from \mathbf{x} can be done in a sequence of rotations and scalings, with appropriate assumptions and algorithms.

Eigenvalue decomposition (EVD)

Reminders from linear algebra

A special case of SVD is when the matrix \mathbf{A} is symmetric.

Reminder

- For a zero-mean random vector \mathbf{x} , the covariance matrix $\mathbf{C}_x = \mathbb{E}\{\mathbf{x}\mathbf{x}^T\}$ can be decomposed as follows:

$$\mathbf{C}_x = \mathbf{U}\Lambda\mathbf{U}^T \tag{5}$$

where $\mathbf{U} = [\mathbf{u}_1, \mathbf{u}_2, \dots, \mathbf{u}_N]$ is an orthonormal matrix with the eigenvectors of \mathbf{C}_x as its columns, and $\Lambda = \text{diag}(\lambda_1, \dots, \lambda_N)$ is a diagonal matrix of the eigenvalues of \mathbf{C}_x .

- Note for later use: $\text{tr}(\mathbf{C}_x) = \sum_{k=1}^N \mathbb{E}(x_k^2) = \sum_{k=1}^N \lambda_k$

Eigenvalue decomposition (EVD)

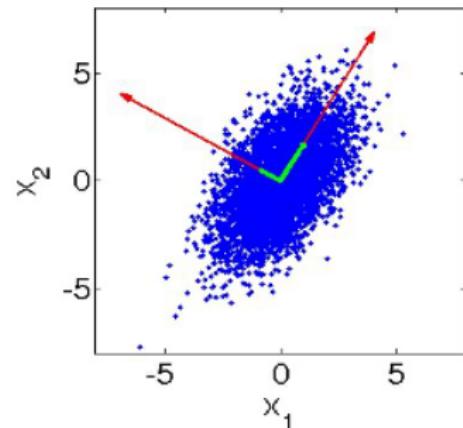
(continued)

Example: 2-dimensional case for zero-mean random variables

$$\mathbf{C}_x = \begin{bmatrix} \mathbb{E}\{x_1^2\} & \mathbb{E}\{x_1 x_2\} \\ \mathbb{E}\{x_2 x_1\} & \mathbb{E}\{x_2^2\} \end{bmatrix}$$

$$\mathbf{U} = [\mathbf{u}_1, \mathbf{u}_2]$$

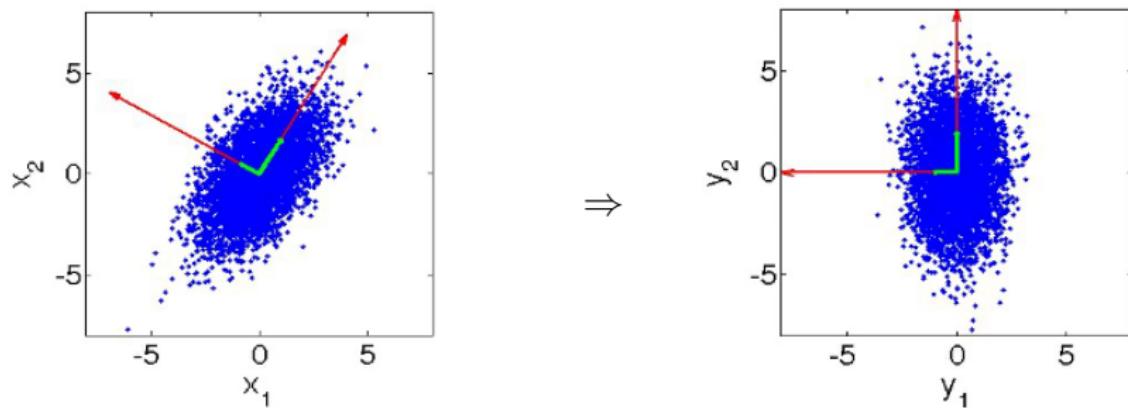
$$\Lambda_x = \begin{bmatrix} \lambda_1 & 0 \\ 0 & \lambda_2 \end{bmatrix}$$



Principal component analysis (PCA)

- Principal component analysis (PCA), also known as the (discrete) Karhunen-Loëve transform or the Hotelling transform, uses EVD to diagonalize the covariance matrix (decorrelates the data).
- If \mathbf{U} is the eigen-matrix of $\mathbf{C}_x = \mathbb{E}\{\mathbf{x}\mathbf{x}^T\}$, the transformation $\mathbf{y} = \mathbf{U}^T \mathbf{x}$ results in:

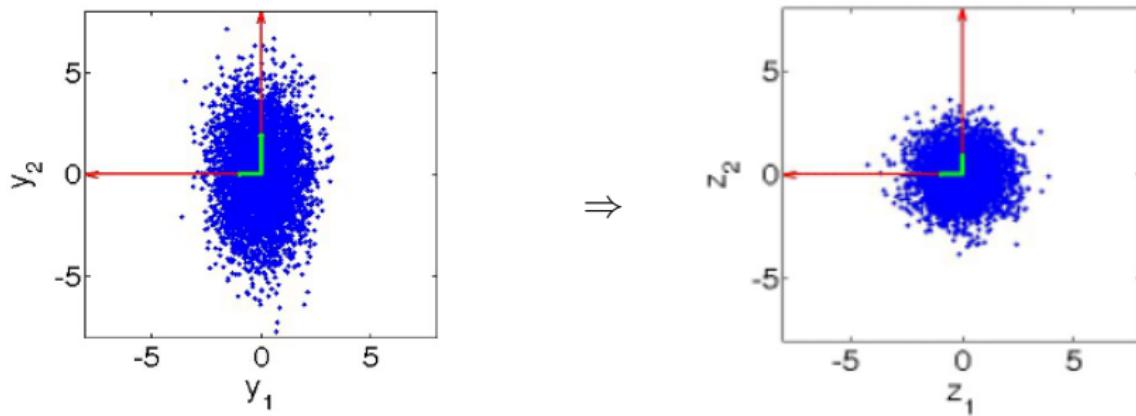
$$\mathbf{C}_y = \mathbb{E}\{\mathbf{y}\mathbf{y}^T\} = \mathbf{U}^T \mathbf{C}_x \mathbf{U} = \Lambda$$



Spherering

- In the previous example, let's go a step further and normalize the covariance matrix via $\mathbf{z} = \Lambda^{-1/2}\mathbf{y}$, resulting in:

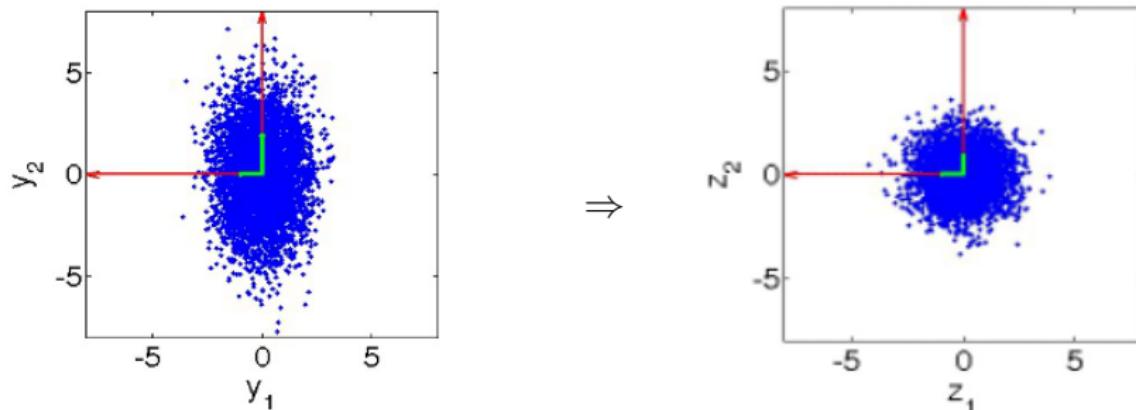
$$\mathbf{C}_z = \mathbb{E}\{\mathbf{zz}^T\} = \Lambda^{-1/2}\mathbf{C}_y\Lambda^{-1/2} = \mathbf{I}$$



Spherering

- In the previous example, let's go a step further and normalize the covariance matrix via $\mathbf{z} = \Lambda^{-1/2}\mathbf{y}$, resulting in:

$$\mathbf{C}_z = \mathbb{E}\{\mathbf{zz}^T\} = \Lambda^{-1/2}\mathbf{C}_y\Lambda^{-1/2} = \mathbf{I}$$



Note: There is no further information in the second-order statistics of \mathbf{z} , i.e., multiplication of \mathbf{z} by any orthonormal (rotation) matrix preserves its whiteness.

Uncorrelatedness vs independence of random variables

Reminders from statistics

Reminder (Hyvarinen et al., 2001)

- **Uncorrelatedness** of random variables x and y means:

$$\mathbb{E}\{xy\} = \mathbb{E}\{x\}\mathbb{E}\{y\}$$

Uncorrelatedness vs independence of random variables

Reminders from statistics

Reminder (Hyvarinen et al., 2001)

- **Uncorrelatedness** of random variables x and y means:
 $\mathbb{E}\{xy\} = \mathbb{E}\{x\}\mathbb{E}\{y\}$
- **Independence** of x and y means: $f_{XY}(x, y) = f_X(x)f_Y(y)$ resulting in
 $\mathbb{E}\{g(x)h(y)\} = \mathbb{E}\{g(x)\}\mathbb{E}\{h(y)\}$ for (almost) arbitrary $g(\cdot)$ and $h(\cdot)$

Uncorrelatedness vs independence of random variables

Reminders from statistics

Reminder (Hyvärinen et al., 2001)

- **Uncorrelatedness** of random variables x and y means:
 $\mathbb{E}\{xy\} = \mathbb{E}\{x\}\mathbb{E}\{y\}$
- **Independence** of x and y means: $f_{XY}(x, y) = f_X(x)f_Y(y)$ resulting in
 $\mathbb{E}\{g(x)h(y)\} = \mathbb{E}\{g(x)\}\mathbb{E}\{h(y)\}$ for (almost) arbitrary $g(\cdot)$ and $h(\cdot)$
- Hence, independence is far stronger than uncorrelatedness and relies on **higher order statistics**.

Uncorrelatedness vs independence of random variables

Reminders from statistics

Reminder (Hyvarinen et al., 2001)

- **Uncorrelatedness** of random variables x and y means:
 $\mathbb{E}\{xy\} = \mathbb{E}\{x\}\mathbb{E}\{y\}$
- **Independence** of x and y means: $f_{XY}(x, y) = f_X(x)f_Y(y)$ resulting in
 $\mathbb{E}\{g(x)h(y)\} = \mathbb{E}\{g(x)\}\mathbb{E}\{h(y)\}$ for (almost) arbitrary $g(\cdot)$ and $h(\cdot)$
- Hence, independence is far stronger than uncorrelatedness and relies on **higher order statistics**.
- For Gaussian random variables, uncorrelatedness is equivalent with independence. Therefore, their higher order statistics do not contain any additional information.



Non-gaussianity of the sources is a key point for conventional ICA algorithms.

Linear transforms of non-Gaussian random variables

Example: Suppose that $s_1 \sim U(-0.5, 0.5)$ and $s_2 \sim U(-1.0, 1.0)$ be zero-mean independent uniformly distributed random variables.

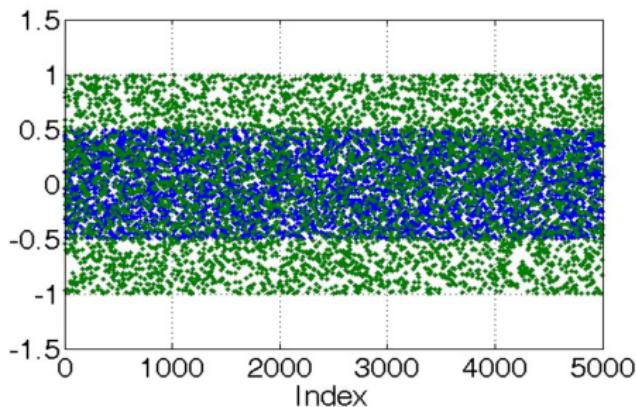


Figure: Sample plot

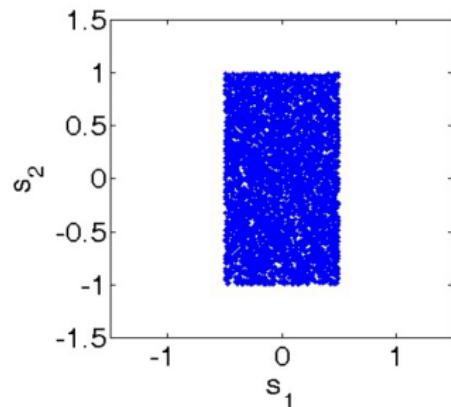


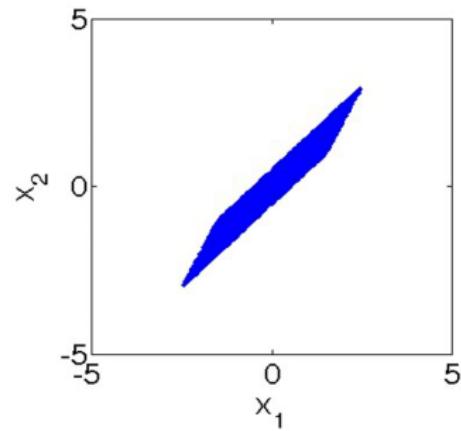
Figure: Scatter plot

Linear transforms of non-Gaussian random variables

(continued)

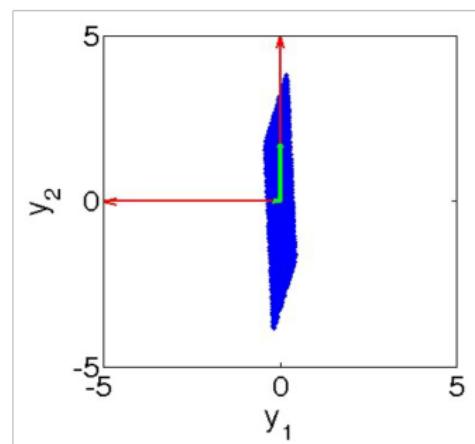
Example: An arbitrary linear transform of (s_1, s_2) will stretch and rotate the mixture

$$\begin{bmatrix} x_1 \\ x_2 \end{bmatrix} = \begin{bmatrix} 1.0 & 2.0 \\ 2.0 & 2.0 \end{bmatrix} \begin{bmatrix} s_1 \\ s_2 \end{bmatrix}$$

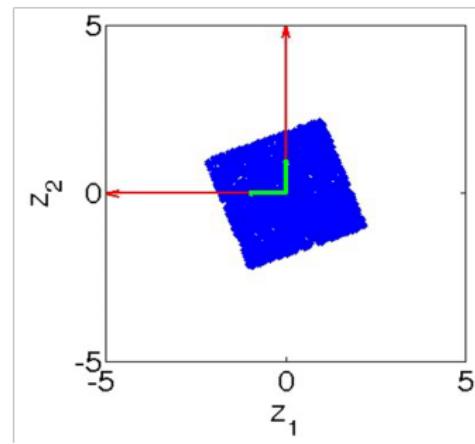


Linear transforms of non-Gaussian random variables

(continued)



After PCA



After spherering

- There is still some **mutual information** in the joint distribution of z_1 and z_2 , which can be used for their separation.
- Therefore, **spherering** is only half the way to **independence** for non-Gaussian variables.

From principal component analysis (PCA) to independent component analysis (ICA)

$$\mathbf{x} \xrightarrow[\text{(rotation)}]{\text{PCA}} \mathbf{y} \xrightarrow[\text{(scaling)}]{\text{Sphering}} \mathbf{z} \xrightarrow[\text{(rotation)}]{\text{ICA}} \mathbf{s}$$

- Multiplication of \mathbf{z} by any orthonormal matrix preserves its whiteness (second order statistics) and will only change its higher order statistics.
- ICA can use the additional information in the higher-order statistics to separate the sources
- How to find the ICA rotation transform?

The central limit theorem

A reminder from statistics

Reminder (Papoulis and Pillai, 2002)

- The summation of N independent identically distributed (IID) random variables, tends to Gaussian as N grows, regardless of their original distributions.
- So the summation of two random variables is “more Gaussian” than the original distribution.



In terms of ICA **non-Gaussian** is **independent**. In practice, we need measures of “independence” and “gaussianity” to quantify these terms for random variables.

How does classical ICA work?

A simplified perspective

Data model & assumptions:

- $\mathbf{x} = \mathbf{As} = \sum_{i=1}^N \mathbf{a}_i s_i$, where $\mathbf{s} = [s_1, \dots, s_N]$ is the source vector and $\mathbf{A} = [\mathbf{a}_1, \dots, \mathbf{a}_N]$ is the **mixing matrix**
- The sources are statistically independent: $f_{\mathbf{S}}(\mathbf{s}) = \prod_{i=1}^N f_{\mathbf{S}_i}(s_i)$
- The sources are non-Gaussian distributed (except at most one)

How does classical ICA work?

A simplified perspective

Data model & assumptions:

- $\mathbf{x} = \mathbf{As} = \sum_{i=1}^N \mathbf{a}_i s_i$, where $\mathbf{s} = [s_1, \dots, s_N]$ is the source vector and $\mathbf{A} = [\mathbf{a}_1, \dots, \mathbf{a}_N]$ is the **mixing matrix**
- The sources are statistically independent: $f_{\mathbf{S}}(\mathbf{s}) = \prod_{i=1}^N f_{\mathbf{s}_i}(s_i)$
- The sources are non-Gaussian distributed (except at most one)

According to the central limit theorem, a linear mixture of the observations

$$y = \sum_{i=1}^N b_i x_i = \mathbf{b}^T \mathbf{x} = \mathbf{b}^T \mathbf{As} = \mathbf{b}^T \sum_{i=1}^N \mathbf{a}_i s_i = \sum_{i=1}^N w_i s_i = \mathbf{w}^T \mathbf{s}$$

is “more Gaussian” than any of the components of \mathbf{s} , except when only one of the elements of \mathbf{w} is nonzero, or $y = \hat{s}_p$ (an estimate of one of the sources).

How does classical ICA work?

A simplified perspective

Data model & assumptions:

- $\mathbf{x} = \mathbf{As} = \sum_{i=1}^N \mathbf{a}_i s_i$, where $\mathbf{s} = [s_1, \dots, s_N]$ is the source vector and $\mathbf{A} = [\mathbf{a}_1, \dots, \mathbf{a}_N]$ is the **mixing matrix**
- The sources are statistically independent: $f_{\mathbf{S}}(\mathbf{s}) = \prod_{i=1}^N f_{\mathbf{s}_i}(s_i)$
- The sources are non-Gaussian distributed (except at most one)

According to the central limit theorem, a linear mixture of the observations

$$y = \sum_{i=1}^N b_i x_i = \mathbf{b}^T \mathbf{x} = \mathbf{b}^T \mathbf{As} = \mathbf{b}^T \sum_{i=1}^N \mathbf{a}_i s_i = \sum_{i=1}^N w_i s_i = \mathbf{w}^T \mathbf{s}$$

is “more Gaussian” than any of the components of \mathbf{s} , except when only one of the elements of \mathbf{w} is nonzero, or $y = \hat{s}_p$ (an estimate of one of the sources).

ICA Algorithm: update \mathbf{b} in the direction of minimum gaussianity of y .

How does classical ICA work?

(continued)

ICA algorithm overview:

- ① Apply PCA and spherling to remove second order correlations (optional but very helpful stage)
- ② Define a measure of independence: Kurthosis (4th Order Cumulant), Negentropy, Mutual Information, etc.
- ③ Find the mixture $y = \mathbf{b}^T \mathbf{z}$, which maximize this measure to find independent sources

How does classical ICA work?

(continued)

ICA algorithm overview:

- ① Apply PCA and spherling to remove second order correlations (optional but very helpful stage)
- ② Define a measure of independence: Kurthosis (4th Order Cumulant), Negentropy, Mutual Information, etc.
- ③ Find the mixture $y = \mathbf{b}^T \mathbf{z}$, which maximize this measure to find independent sources

Question: Is there any guarantee that independent components obtained from this algorithm are the originally mixed components?

How does classical ICA work?

(continued)

ICA algorithm overview:

- ① Apply PCA and spherling to remove second order correlations (optional but very helpful stage)
- ② Define a measure of independence: Kurthosis (4th Order Cumulant), Negentropy, Mutual Information, etc.
- ③ Find the mixture $y = \mathbf{b}^T \mathbf{z}$, which maximize this measure to find independent sources

Question: Is there any guarantee that independent components obtained from this algorithm are the originally mixed components?

Answer: Yes, to some extent!

Limitations of classical ICA

- Requires the same number of sources as sensors

Limitations of classical ICA

- Requires the same number of sources as sensors
- Does not work for gaussian data (at most one of the sources can be gaussian)

Limitations of classical ICA

- Requires the same number of sources as sensors
- Does not work for gaussian data (at most one of the sources can be gaussian)
- Can estimate the sources, up to a “linear copy”:
 - ① Can't estimate the energy nor sign of the signal:

$$\mathbf{x} = \mathbf{As} = (\mathbf{A}\alpha)(\alpha^{-1}\mathbf{s}) = \tilde{\mathbf{A}}\tilde{\mathbf{s}}$$

where α is an arbitrary diagonal matrix or scalar

Limitations of classical ICA

- Requires the same number of sources as sensors
- Does not work for gaussian data (at most one of the sources can be gaussian)
- Can estimate the sources, up to a “linear copy”:
 - ① Can't estimate the energy nor sign of the signal:

$$\mathbf{x} = \mathbf{As} = (\mathbf{A}\alpha)(\alpha^{-1}\mathbf{s}) = \tilde{\mathbf{A}}\tilde{\mathbf{s}}$$

where α is an arbitrary diagonal matrix or scalar

- ② Can't estimate the order of the signals:

$$\mathbf{x} = \mathbf{As} = (\mathbf{AP})(\mathbf{P}^{-1}\mathbf{s})$$

where \mathbf{P} is an arbitrary permutation matrix.

A fast converging ICA algorithm

The FastICA algorithm

- ① Given \mathbf{x} , pre-whiten the data by PCA and spherling to obtain \mathbf{z}
- ② Select a non-quadratic nonlinear function $f(\cdot)$, with its first derivative $g(\cdot)$, second derivative $\dot{g}(\cdot)$, and an initial weight vector \mathbf{w} .
- ③ Repeat until convergence:
 - ① Let $\mathbf{w} \leftarrow \mathbb{E}\{\mathbf{z}g(\mathbf{w}^T\mathbf{z})\} - \mathbb{E}\{\dot{g}(\mathbf{w}^T\mathbf{z})\}\mathbf{w}$, where $\mathbb{E}\{\cdot\}$ denotes expectation
 - ② Let $\mathbf{w} \leftarrow \mathbf{w}/\|\mathbf{w}\|$

A fast converging ICA algorithm

The FastICA algorithm

- ① Given \mathbf{x} , pre-whiten the data by PCA and spherizing to obtain \mathbf{z}
- ② Select a non-quadratic nonlinear function $f(\cdot)$, with its first derivative $g(\cdot)$, second derivative $\dot{g}(\cdot)$, and an initial weight vector \mathbf{w} .
- ③ Repeat until convergence:
 - ① Let $\mathbf{w} \leftarrow \mathbb{E}\{\mathbf{z}g(\mathbf{w}^T\mathbf{z})\} - \mathbb{E}\{\dot{g}(\mathbf{w}^T\mathbf{z})\}\mathbf{w}$, where $\mathbb{E}\{\cdot\}$ denotes expectation
 - ② Let $\mathbf{w} \leftarrow \mathbf{w}/\|\mathbf{w}\|$

Examples of suitable nonlinear function

- $f(u) = \log(\cosh(u))$, $g(u) = \tanh(u)$ and $\dot{g} = 1 - \tanh^2(u)$
- $f(u) = -e^{-u^2/2}$, $g(u) = ue^{-u^2/2}$ and $\dot{g}(u) = (1 - u^2)e^{-u^2/2}$

Outline

- 1 Background & Motivation
- 2 Classical Independent Component Analysis Algorithms
- 3 Iterative BSS Methods
- 4 Algebraic BSS Methods
- 5 Further Reading & Advanced Topics
- 6 Appendix: BSS Lab Session Notes

Degrees of freedom in linear algebraic transforms

Problem restatement

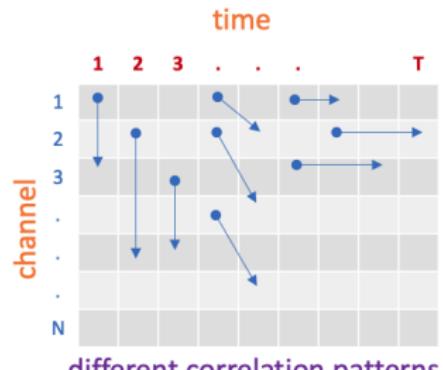
Let us restate BSS from a totally different viewpoint:

Degrees of freedom in linear transforms

- Consider the observations $\mathbf{x}(t) \in \mathbb{R}^{N \times T}$. We seek linear transforms of the form $\mathbf{y}(t) = \mathbf{W}\mathbf{x}(t)$ where $\mathbf{W} \in \mathbb{R}^{N \times N}$ is an unknown **separating matrix**.
- \mathbf{W} acts as a **spatial filter** (N instantaneous linear combinations of the multichannel observations) that changes the properties of $\mathbf{y}(t)$
- The problem of linear source separation can be viewed as a **matrix design** problem with N^2 **degrees of freedom (DoF)**, to obtain $\mathbf{y}(t)$ with desired properties

Second and higher order statistics matrices

- Examples of matrices containing second and higher order statistics of a multichannel random process $\mathbf{x}(t)$ includes:
 - $\mathbf{C} = \mathbb{E}\{\mathbf{x}(t)\mathbf{x}^T(t)\}$: The covariance matrix
 - $\mathbf{C}(\tau) = \mathbb{E}\{\mathbf{x}(t)\mathbf{x}^T(t + \tau)\}$: The lagged covariance matrix
 - $\mathbf{C}_k = \mathbb{E}\{\mathbf{x}(t)\mathbf{x}^T(t)x_k(t)\}$: Slices of higher-order statistics
- This can more generally be extended to **tensors** of higher order statistics.
- These matrices carry average inter-channel similarities (correlation) over time
- Diagonalization of a number or all of these matrices results in the **decorrelation** and **independence** of the corresponding channels



different correlation patterns

Source separation by temporal decorrelation

- Consider $\mathbf{x}(t)$ with covariance matrix $\mathbf{C}_0 = \mathbb{E}\{\mathbf{x}(t)\mathbf{x}(t)^T\}$ and lagged-covariance matrices $\mathbf{C}_1 = \mathbb{E}\{\mathbf{x}(t)\mathbf{x}(t + \tau_1)^T\}$, $\mathbf{C}_2 = \mathbb{E}\{\mathbf{x}(t)\mathbf{x}(t + \tau_2)^T\}$, etc.
- For $\mathbf{y}(t) = \mathbf{Wx}(t)$:

$$\mathbf{A}_0 = \mathbb{E}\{\mathbf{y}(t)\mathbf{y}(t)^T\} = \mathbb{E}\{\mathbf{Wx}(t)\mathbf{x}(t)^T \mathbf{W}^T\} = \mathbf{WC}_0 \mathbf{W}^T$$

...

$$\mathbf{A}_k = \mathbb{E}\{\mathbf{y}(t)\mathbf{y}(t + \tau_k)^T\} = \mathbb{E}\{\mathbf{Wx}(t)\mathbf{x}(t + \tau_k)^T \mathbf{W}^T\} = \mathbf{WC}_k \mathbf{W}^T$$

- These matrices show the inter-channel correlations of the output $\mathbf{y}(t)$, at different time lags.
- For temporally correlated data, diagonalizing all such matrices results in channel-wise independence of $\mathbf{y}(t)$.

Source separation by temporal decorrelation

- Consider $\mathbf{x}(t)$ with covariance matrix $\mathbf{C}_0 = \mathbb{E}\{\mathbf{x}(t)\mathbf{x}(t)^T\}$ and lagged-covariance matrices $\mathbf{C}_1 = \mathbb{E}\{\mathbf{x}(t)\mathbf{x}(t + \tau_1)^T\}$, $\mathbf{C}_2 = \mathbb{E}\{\mathbf{x}(t)\mathbf{x}(t + \tau_2)^T\}$, etc.
- For $\mathbf{y}(t) = \mathbf{Wx}(t)$:

$$\mathbf{A}_0 = \mathbb{E}\{\mathbf{y}(t)\mathbf{y}(t)^T\} = \mathbb{E}\{\mathbf{Wx}(t)\mathbf{x}(t)^T \mathbf{W}^T\} = \mathbf{WC}_0 \mathbf{W}^T$$

...

$$\mathbf{A}_k = \mathbb{E}\{\mathbf{y}(t)\mathbf{y}(t + \tau_k)^T\} = \mathbb{E}\{\mathbf{Wx}(t)\mathbf{x}(t + \tau_k)^T \mathbf{W}^T\} = \mathbf{WC}_k \mathbf{W}^T$$

- These matrices show the inter-channel correlations of the output $\mathbf{y}(t)$, at different time lags.
- For temporally correlated data, diagonalizing all such matrices results in channel-wise independence of $\mathbf{y}(t)$.

Question: Can we find a matrix \mathbf{W} , which simultaneously diagonalizes all or a number of the above matrices?

Degrees-of-freedom (DoF) in algebraic transforms

Reminder: eigenvalue decomposition (EVD)

Using EVD, a symmetric matrix $\mathbf{A} \in \mathbb{R}^{N \times N}$ can be decomposed as:

$$\mathbf{C}_0 = \mathbf{U}\Lambda\mathbf{U}^T \quad (6)$$

where

- \mathbf{C}_0 is a symmetric matrix with $N(N + 1)/2$ DoF
- Λ is a diagonal matrix with N DoF
- \mathbf{U} is the orthonormal eigen-matrix with $N(N - 1)/2$ DoF

Therefore

$$\text{DoF}(\mathbf{C}_0) = \text{DoF}(\Lambda) + \text{DoF}(\mathbf{U})$$

Degrees-of-freedom (DoF) in algebraic transforms

(continued)

Reminder: generalized eigenvalue decomposition (GEVD)

Two symmetric matrices \mathbf{C}_0 and \mathbf{C}_1 can be diagonalized using GEVD as follows:

$$\begin{aligned}\mathbf{W}\mathbf{C}_0\mathbf{W}^T &= \mathbf{I} \\ \mathbf{W}\mathbf{C}_1\mathbf{W}^T &= \Lambda\end{aligned}\tag{7}$$

where

- \mathbf{C}_0 and \mathbf{C}_1 are symmetric matrices, each with $N(N + 1)/2$ DoF
- Λ is a diagonal matrix with N DoF
- \mathbf{W} is a square matrix with N^2 DoF
- \mathbf{I} is the identity matrix with no DoF

Therefore

$$\text{DoF}(\mathbf{C}_0) + \text{DoF}(\mathbf{C}_1) = \text{DoF}(\Lambda) + \text{DoF}(\mathbf{W})$$

Degrees-of-freedom (DoF) in algebraic transforms

(continued)

Reminder: joint diagonalization of multiple matrices

No more than two symmetric matrices can be simultaneously diagonalized, except when they belong to the same eigen-space:

$$\begin{aligned} \mathbf{W} \mathbf{C}_0 \mathbf{W}^T &= \mathbf{I} \\ \mathbf{W} \mathbf{C}_1 \mathbf{W}^T &= \Lambda_1 \\ &\dots \\ \mathbf{W} \mathbf{C}_K \mathbf{W}^T &= \Lambda_K \end{aligned} \tag{8}$$

 For more than two matrices only **approximate joint diagonalization (AJD)** may be achieved.

Source separation by matrix diagonalization

Based on the above properties, two classes of algebraic source separation methods are common:

- ① Exact diagonalization of two “well-chosen” matrices:
 - Advantage: The corresponding statistics are fully separated
 - Disadvantage: Requires prior information for selecting the matrices
 - Major methods of this class: AMUSE, Π CA, etc.
- ② Approximate joint diagonalization of more than two matrices:
 - Advantage: It is easier to select the matrices
 - Disadvantage: The corresponding statistics are only separated approximately; measures of approximate diagonalization are not unique
 - Major methods of this class: JADE, SOBI, etc.

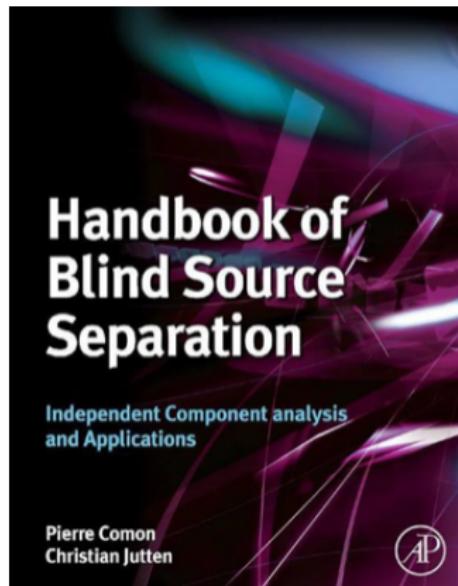
Outline

- 1 Background & Motivation
- 2 Classical Independent Component Analysis Algorithms
- 3 Iterative BSS Methods
- 4 Algebraic BSS Methods
- 5 Further Reading & Advanced Topics
- 6 Appendix: BSS Lab Session Notes

Textbooks on BSS and ICA

Further reading

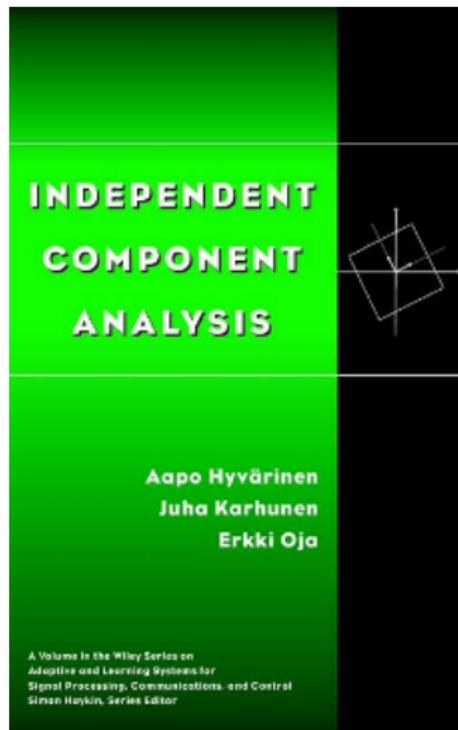
Comon, P., & Jutten, C. (Eds.). (2010). Handbook of Blind Source Separation: Independent component analysis and applications. Academic press. (Comon and Jutten, 2010)



Textbooks on BSS and ICA

Further reading

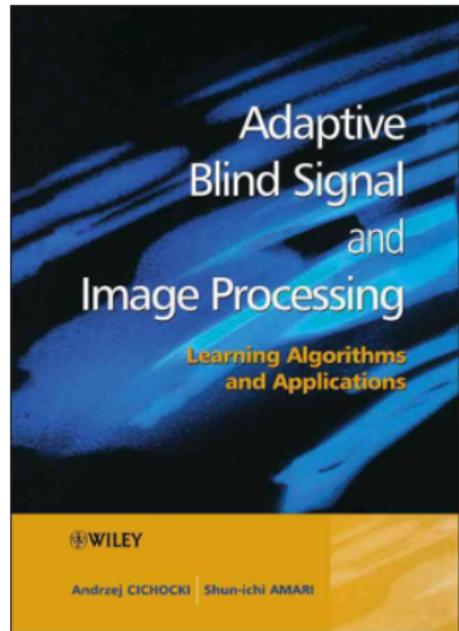
Hyvärinen, A., Karhunen, J., & Oja, E. (2004). Independent component analysis, John Wiley & Sons.(Hyvarinen et al., 2001)



Textbooks on BSS and ICA

Further reading

Cichocki, A., & Amari, S. I. (2002). Adaptive blind signal and image processing: learning algorithms and applications (Vol. 1). John Wiley & Sons. (Cichocki and Amari, 2003)



Source codes and algorithms to consider

- FastICA. Online available:
<http://research.ics.aalto.fi/ica/fastica/>
- Joint Approximation Diagonalization of Eigen-matrices (JADE):
http://www2.iap.fr/users/cardoso/compsep_classic.html
- Second order blind identification (SOBI)
- The Open-Source Electrophysiological Toolbox (OSET): www.oset.ir

- Semi-blind source separation: temporal structure, sparsity, periodicity, etc.
- Noisy mixtures; integration of filtering and source separation
- Approximate joint diagonalization methods
- Nonlinear mixtures
- Sparse component analysis
- Distributed (non-punctual) sources
- Convulsive mixtures
- Online ICA and time-variant mixtures
- Curse of dimensionality and implementation issues
- BSS using classical estimation techniques (MAP, ML, CRLB)
- Tensor decomposition for BSS

Thank you! Questions?

Outline

- 1 Background & Motivation
- 2 Classical Independent Component Analysis Algorithms
- 3 Iterative BSS Methods
- 4 Algebraic BSS Methods
- 5 Further Reading & Advanced Topics
- 6 Appendix: BSS Lab Session Notes

BSS lab session

preface

Algorithms:

- PCA
- GEVD: AMUSE, NICA, NSCA
- JADE
- FastICA

Case Studies:

- Audio source separation
- Noninvasive fetal ECG extraction
- EOG artifact cancellation from EEG
- Device noise and artifact cancellation from biomedical signals

BSS lab session

requirements

Clone the OSET from either repositories:

- <https://gitlab.com/rsameni/OSET.git> (or its image repository on GitHub: <https://github.com/alphanumericslab/OSET.git>)

Clone the tutorial for efficient Matlab coding (optional):

- <https://github.com/rsameni/TechReport-EfficientMatlabCodeTutorial.git>

Clone the lab sample codes (**publicly available after the lecture**):

- <https://github.com/rsameni/BSSLecture.git>

Principal component analysis algorithms

Reminder

By definition, the eigenvalues and eigenvectors of a symmetric matrix $\mathbf{A} \in \mathbb{R}^{N \times N}$ satisfy:

$$\mathbf{A}\mathbf{v} = \lambda\mathbf{v} \quad (9)$$

Principal component analysis algorithms

Reminder

By definition, the eigenvalues and eigenvectors of a symmetric matrix $\mathbf{A} \in \mathbb{R}^{N \times N}$ satisfy:

$$\mathbf{A}\mathbf{v} = \lambda\mathbf{v} \quad (9)$$

Exercise

- Calculate the eigenvalues and eigenvectors of $\mathbf{A} = \begin{bmatrix} 3 & 1 \\ 1 & 4 \end{bmatrix}$

How?

Principal component analysis algorithms

Reminder

By definition, the eigenvalues and eigenvectors of a symmetric matrix $\mathbf{A} \in \mathbb{R}^{N \times N}$ satisfy:

$$\mathbf{A}\mathbf{v} = \lambda\mathbf{v} \quad (9)$$

Exercise

- Calculate the eigenvalues and eigenvectors of $\mathbf{A} = \begin{bmatrix} 3 & 1 \\ 1 & 4 \end{bmatrix}$

How? Find the roots of the characteristic polynomial $\det(\mathbf{A} - \lambda\mathbf{I}) = 0$; then use them to obtain the eigenvectors from (9).

Note: Impose $\|\mathbf{v}\| = 1$ to obtain unique results.

Principal component analysis algorithms

Reminder

By definition, the eigenvalues and eigenvectors of a symmetric matrix $\mathbf{A} \in \mathbb{R}^{N \times N}$ satisfy:

$$\mathbf{A}\mathbf{v} = \lambda\mathbf{v} \quad (9)$$

Exercise

- ① Calculate the eigenvalues and eigenvectors of $\mathbf{A} = \begin{bmatrix} 3 & 1 \\ 1 & 4 \end{bmatrix}$

How? Find the roots of the characteristic polynomial $\det(\mathbf{A} - \lambda\mathbf{I}) = 0$; then use them to obtain the eigenvectors from (9).

Note: Impose $\|\mathbf{v}\| = 1$ to obtain unique results.

- ② Use the function `eig` in Matlab to check the result

Principal component analysis algorithms

Reminder

By definition, the eigenvalues and eigenvectors of a symmetric matrix $\mathbf{A} \in \mathbb{R}^{N \times N}$ satisfy:

$$\mathbf{A}\mathbf{v} = \lambda\mathbf{v} \quad (9)$$

Exercise

- ① Calculate the eigenvalues and eigenvectors of $\mathbf{A} = \begin{bmatrix} 3 & 1 \\ 1 & 4 \end{bmatrix}$
How? Find the roots of the characteristic polynomial $\det(\mathbf{A} - \lambda\mathbf{I}) = 0$; then use them to obtain the eigenvectors from (9).
Note: Impose $\|\mathbf{v}\| = 1$ to obtain unique results.
- ② Use the function `eig` in Matlab to check the result
- ③ Generate 10000 ensembles of a Gaussian random vector with covariance matrix equal to \mathbf{A}
- ④ How can we calculate the eigenvalues and eigenvectors numerically?

Principal component analysis algorithms

(continued)

A simple (but inefficient) method for calculating the largest eigenvalue and its corresponding eigenvector is:

The power method for EVD

- ① Initialize by a random vector $\mathbf{v}_0 \in \mathbb{R}^{N \times N}$
- ② Repeat until convergence: $\mathbf{v}_{k+1} = \frac{\mathbf{A}\mathbf{v}_k}{\|\mathbf{A}\mathbf{v}_k\|}$
- ③ $\lambda = \|\mathbf{A}\mathbf{v}\| / \|\mathbf{v}\|$

Principal component analysis algorithms

(continued)

A simple (but inefficient) method for calculating the largest eigenvalue and its corresponding eigenvector is:

The power method for EVD

- ① Initialize by a random vector $\mathbf{v}_0 \in \mathbb{R}^{N \times N}$
- ② Repeat until convergence: $\mathbf{v}_{k+1} = \frac{\mathbf{A}\mathbf{v}_k}{\|\mathbf{A}\mathbf{v}_k\|}$
- ③ $\lambda = \|\mathbf{A}\mathbf{v}\| / \|\mathbf{v}\|$

Question: Does the algorithm necessarily converge?

Principal component analysis algorithms

(continued)

A simple (but inefficient) method for calculating the largest eigenvalue and its corresponding eigenvector is:

The power method for EVD

- ① Initialize by a random vector $\mathbf{v}_0 \in \mathbb{R}^{N \times N}$
- ② Repeat until convergence: $\mathbf{v}_{k+1} = \frac{\mathbf{A}\mathbf{v}_k}{\|\mathbf{A}\mathbf{v}_k\|}$
- ③ $\lambda = \|\mathbf{A}\mathbf{v}\| / \|\mathbf{v}\|$

Question: Does the algorithm necessarily converge?

Question: How to measure the convergence?

Principal component analysis algorithms

(continued)

A simple (but inefficient) method for calculating the largest eigenvalue and its corresponding eigenvector is:

The power method for EVD

- ① Initialize by a random vector $\mathbf{v}_0 \in \mathbb{R}^{N \times N}$
- ② Repeat until convergence: $\mathbf{v}_{k+1} = \frac{\mathbf{A}\mathbf{v}_k}{\|\mathbf{A}\mathbf{v}_k\|}$
- ③ $\lambda = \|\mathbf{A}\mathbf{v}\| / \|\mathbf{v}\|$

Question: Does the algorithm necessarily converge?

Question: How to measure the convergence?

Question: How can we find the other eigenvalues and eigenvectors?

Principal component analysis algorithms

(continued)

Reminder: rank-one expansion for EVD (Strang, 1988)

Consider the EVD of a symmetric matrix $\mathbf{A} = \mathbf{U}\Lambda\mathbf{U}^T$. It can be shown that \mathbf{A} can be decomposed as a summation of **rank-one matrices**:

$$\mathbf{A} = \lambda_1 \mathbf{u}_1 \mathbf{u}_1^T + \dots + \lambda_N \mathbf{u}_N \mathbf{u}_N^T$$

Principal component analysis algorithms

(continued)

Reminder: rank-one expansion for EVD (Strang, 1988)

Consider the EVD of a symmetric matrix $\mathbf{A} = \mathbf{U}\Lambda\mathbf{U}^T$. It can be shown that \mathbf{A} can be decomposed as a summation of **rank-one matrices**:

$$\mathbf{A} = \lambda_1 \mathbf{u}_1 \mathbf{u}_1^T + \dots + \lambda_N \mathbf{u}_N \mathbf{u}_N^T$$

The above property can be used to extract the other eigenvalues and eigenvectors in the power method:

The power method (continued)

- ① $\tilde{\mathbf{A}} = \mathbf{A} - \lambda_1 \mathbf{u}_1 \mathbf{u}_1^T$
- ② Apply the power method on $\tilde{\mathbf{A}}$
- ③ Repeat the same procedure for all eigenvalues and eigenvectors

Principal component analysis algorithms

(continued)

Reminder: rank-one expansion for EVD (Strang, 1988)

Consider the EVD of a symmetric matrix $\mathbf{A} = \mathbf{U}\Lambda\mathbf{U}^T$. It can be shown that \mathbf{A} can be decomposed as a summation of **rank-one matrices**:

$$\mathbf{A} = \lambda_1 \mathbf{u}_1 \mathbf{u}_1^T + \dots + \lambda_N \mathbf{u}_N \mathbf{u}_N^T$$

The above property can be used to extract the other eigenvalues and eigenvectors in the power method:

The power method (continued)

- ① $\tilde{\mathbf{A}} = \mathbf{A} - \lambda_1 \mathbf{u}_1 \mathbf{u}_1^T$
- ② Apply the power method on $\tilde{\mathbf{A}}$
- ③ Repeat the same procedure for all eigenvalues and eigenvectors

Drawbacks: slow convergence; error accumulation during rank-one updates

Generalized eigenvalue decomposition

Exercise

- Consider zero-mean multichannel electroencephalogram signals $\mathbf{x}(t) \in \mathbb{R}^{N \times N}$ with $t \in \mathcal{T} = \{1, \dots, T\}$ and covariance matrix $\mathbf{C}_0 = \mathbb{E}_t\{\mathbf{x}(t)\mathbf{x}(t)^T\}$
- Suppose that $\mathcal{U} = \{u\} \subset \mathcal{T}$ is a subset of time samples corresponding to event-related potentials (ERP), with $\mathbf{C}_1 = \mathbb{E}_u\{\mathbf{x}(u)\mathbf{x}(u)^T\}$
- We seek $y(t) = \mathbf{w}^T \mathbf{x}(t)$, such that it has the maximum energy ratio during the ERP time epochs, while its total energy remains fixed.
- How to formulate the problem?

Generalized eigenvalue decomposition

Exercise

- Consider zero-mean multichannel electroencephalogram signals $\mathbf{x}(t) \in \mathbb{R}^{N \times N}$ with $t \in \mathcal{T} = \{1, \dots, T\}$ and covariance matrix $\mathbf{C}_0 = \mathbb{E}_t\{\mathbf{x}(t)\mathbf{x}(t)^T\}$
- Suppose that $\mathcal{U} = \{u\} \subset \mathcal{T}$ is a subset of time samples corresponding to event-related potentials (ERP), with $\mathbf{C}_1 = \mathbb{E}_u\{\mathbf{x}(u)\mathbf{x}(u)^T\}$
- We seek $y(t) = \mathbf{w}^T \mathbf{x}(t)$, such that it has the maximum energy ratio during the ERP time epochs, while its total energy remains fixed.
- How to formulate the problem?

$$\mathbf{w}^* = \arg \max_{\mathbf{w}} \|\mathbf{w}^T \mathbf{x}(u)\|, \text{ s.t. } \|\mathbf{w}^T \mathbf{x}(t)\| = cte.$$

- Or equivalently

$$\mathbf{w}^* = \arg \max_{\mathbf{w}} \mathbf{w}^T \mathbf{C}_1 \mathbf{w}, \text{ s.t. } \mathbf{w}^T \mathbf{C}_0 \mathbf{w} = cte. \quad (10)$$

Generalized eigenvalue decomposition

(continued)

The Rayleigh quotient and generalized eigenvalue decomposition (GEVD) (Strang, 1988)

For real symmetric matrices \mathbf{C}_0 and \mathbf{C}_1 , the following problems are equivalent:

- $\mathbf{W}\mathbf{C}_0\mathbf{W}^T = \mathbf{I}$ and $\mathbf{W}\mathbf{C}_1\mathbf{W}^T = \Lambda$, where $\mathbf{W} = [\mathbf{w}_1, \dots, \mathbf{w}_N]$,
 $\Lambda = \text{diag}(\lambda_1, \dots, \lambda_N)$ and $\lambda_1 \geq \dots \geq \lambda_N$
- $\mathbf{w}_1 = \arg \max_{\mathbf{w}} r = \frac{\mathbf{w}^T \mathbf{C}_1 \mathbf{w}}{\mathbf{w}^T \mathbf{C}_0 \mathbf{w}}$ (known as the Rayleigh quotient)
- $\mathbf{w}_1 = \arg \max_{\mathbf{w}} \mathbf{w}^T \mathbf{C}_1 \mathbf{w}$, s.t. $\mathbf{w}^T \mathbf{C}_0 \mathbf{w} = \text{cte.}$

References

- Cichocki, A. and Amari, S., eds (2003). Adaptive Blind Signal and Image Processing. John Wiley & Sons Inc.
- Comon, P. and Jutten, C. (2010). Handbook of Blind Source Separation: Independent Component Analysis and Applications. Independent Component Analysis and Applications Series, Elsevier Science.
- Hyvärinen, A., Karhunen, J. and Oja, E. (2001). Independent Component Analysis. Wiley-Interscience.
- Papoulis, A. and Pillai, S. U. (2002). Probability, random variables, and stochastic processes. Fourth edition, McGraw-Hill.
- Roll, J.-P. (1981). Contribution de la proprioception musculaire à la perception et au contrôle du mouvement chez l'homme. PhD thesis, Uni. D'Aix-Marseille I. Thèse d'État.
- Strang, G. (1988). Linear Algebra and Its Applications. 3 edition, Brooks/Cole.