

Independent Component Analysis

Xing Yi Liu, Weiyu Huang, Landon Choi

March 10, 2021

Section 1

Introduction

What is ICA?

- A method of dimension reduction used to separate data into independent basis vectors
- Deciphers the latent non-Gaussian variables which make up a multivariate mixed signal
- **Applications of ICA:** Cocktail Party Problem, EEG (brain activity), feature extraction, image processing

Cocktail Party Problem



Figure 1: Cocktail Party

Blind Source Separation

- Given m recordings of mixed sound sources, blind source separation seeks to uncover the underlying individual sources without any information other than the mixed data.
- Typically, we assume that if there are m recordings of a mixed sound, there are also m independent sound sources.

Graphical Illustration

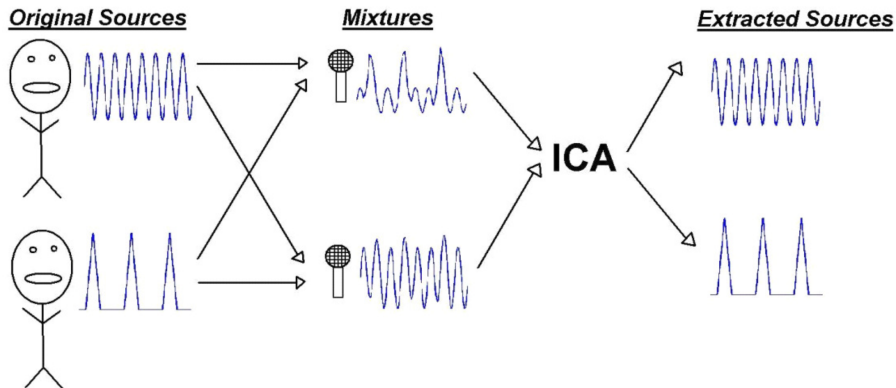


Figure 2: Illustration of ICA, $m = 2$

Section 2

Mathematical Formulation

Notation

Let the observed data be $\{\mathbf{x}_i\}_{i=1}^N$, where

$$\mathbf{x}_i = \begin{pmatrix} x_{i1} \\ x_{i2} \\ \vdots \\ x_{im} \end{pmatrix} = \begin{pmatrix} a_{11}s_{i1} + a_{12}s_{i2} + \dots + a_{1m}s_{im} \\ a_{21}s_{i1} + a_{22}s_{i2} + \dots + a_{2m}s_{im} \\ \vdots \\ a_{m1}s_{i1} + a_{m2}s_{i2} + \dots + a_{mm}s_{im} \end{pmatrix} = \mathbf{A}\mathbf{s}_i.$$

Then,

$$\mathbf{W}\mathbf{x}_i = \mathbf{s}_i,$$

where $\mathbf{W} = \mathbf{A}^{-1}$. Our goal is to estimate \mathbf{W} by assuming the distribution of \mathbf{s}_i , and approximate \mathbf{s}_i .

Preprocessing

Preprocessing involves centering and whitening the observation data.

Let $\mathbf{x} \in \{\mathbf{x}_i\}_{i=1}^N$. Let

$$\tilde{\mathbf{x}} = \mathbf{x} - \bar{\mathbf{x}}$$

represent the centered data. Our goal is to obtain whitened data \mathbf{x}_w satisfying

$$\mathbb{E}(\mathbf{x}_w \mathbf{x}_w^T) = \mathbf{I}.$$

To achieve this, perform eigen-decomposition on the covariance matrix of $\tilde{\mathbf{x}}$:

$$\mathbb{E}(\tilde{\mathbf{x}} \tilde{\mathbf{x}}^T) = \mathbf{V} \mathbf{D} \mathbf{V}^{-1},$$

where \mathbf{V} is the matrix whose columns are eigenvectors of $\mathbb{E}(\tilde{\mathbf{x}} \tilde{\mathbf{x}}^T)$, and \mathbf{D} is the diagonal matrix of eigenvalues $\lambda_1, \dots, \lambda_m$.

Preprocessing

The whitened signal can be obtained by

$$\mathbf{x}_w = \mathbf{V}\mathbf{D}^{-\frac{1}{2}}\mathbf{V}^T\tilde{\mathbf{x}} = \mathbf{A}_w\mathbf{s},$$

where $\mathbf{A}_w = \mathbf{V}\mathbf{D}^{-1/2}\mathbf{V}^T\mathbf{A}$ is orthogonal.

\mathbf{A}_w has only $\frac{m(m-1)}{2}$ degrees of freedom, less than the m^2 parameters of \mathbf{A} . Whitening thus allows us to estimate less parameters.

Maximum Likelihood Estimation

Since $\mathbf{x}_i = \mathbf{A}\mathbf{s}_i$,

$$\begin{aligned} p(\mathbf{x}_i) &= \frac{1}{|\det(\mathbf{A})|} p(\mathbf{s}_i) \\ &= |\det(\mathbf{W})| \prod_{j=1}^m p_i(\mathbf{w}_j^\top \mathbf{x}_i), \end{aligned}$$

where \mathbf{w}_j^\top is the j -th row (a row vector) of $\mathbf{W} = \mathbf{A}^{-1}$. Then, the likelihood is

$$L(\mathbf{W}) = |\det(\mathbf{W})|^N \prod_{i=1}^N \prod_{j=1}^m p_i(\mathbf{w}_j^\top \mathbf{x}_i)$$

$$\log(L(\mathbf{W})) = N \log |\det(\mathbf{W})| + \sum_{i=1}^N \sum_{j=1}^m \log(p_i(\mathbf{w}_j^\top \mathbf{x}_i))$$

$$\frac{1}{N} \log(L(\mathbf{W})) = \log |\det(\mathbf{W})| + \frac{1}{N} \sum_{i=1}^N \sum_{j=1}^m \log(p_i(\mathbf{w}_j^\top \mathbf{x}_i)).$$

Maximum Likelihood Estimation

$$\frac{\partial}{\partial \mathbf{W}} \frac{1}{N} \log(L(\mathbf{W})) = (\mathbf{W}^T)^{-1} + \frac{1}{N} \sum_{i=1}^N \mathbf{g}(\mathbf{W} \mathbf{x}_i) \mathbf{x}_i^T,$$

where

$$\mathbf{g}(\mathbf{a}) = \begin{pmatrix} g_1(a_1) \\ g_2(a_2) \\ \vdots \\ g_m(a_m) \end{pmatrix}.$$

The functions $g_i = (\log(p_i))' = \frac{p'_i}{p_i}$, $i = 1, \dots, m$ are called **score functions** and depend on the probability distribution p_i of the source s_i .

Bell-Sejnowski Algorithm

Require: $\{\mathbf{x}_n\}_{n=1}^N, \{p_i(s_i)\}_{i=1}^m$

- ➊ Calculate the sample mean $\bar{\mathbf{x}}$ of the data and center the data by $\tilde{\mathbf{x}}_i \leftarrow \mathbf{x}_i - \bar{\mathbf{x}}$ for $i = 1, \dots, N$.
- ➋ Calculate \mathbf{V}, \mathbf{D} such that $\mathbb{E}(\tilde{\mathbf{x}}_i \tilde{\mathbf{x}}_i^\top) = \mathbf{V} \mathbf{D} \mathbf{V}^{-1}$ is the eigen-decomposition of $\mathbb{E}(\tilde{\mathbf{x}}_i \tilde{\mathbf{x}}_i^\top)$.
- ➌ Preprocess the data by $\mathbf{x}_{wi} = \mathbf{V} \mathbf{D}^{-1/2} \mathbf{V}^\top \tilde{\mathbf{x}}_i$ for $i = 1, \dots, N$.
- ➍ Initialize $\mathbf{W} \in \mathbf{R}^{m \times m}$
- ➎ Update $\mathbf{W} \leftarrow \mathbf{W} - \eta \left((\mathbf{W}^\top)^{-1} + \frac{1}{N} \sum_{i=1}^N \mathbf{g}(\mathbf{W} \mathbf{x}_{wi}) \mathbf{x}_{wi}^\top \right)$, where $\eta > 0$.
- ➏ If stopping criterion is satisfied, proceed to step 7. Otherwise, return to step 5.
- ➐ Estimates for the separated latent variables will be given by $\mathbf{y}_i = \mathbf{W} \mathbf{x}_i$.

Non-Gaussianity

Kurtosis is a classical measure of non-Gaussianity:

$$\kappa(y) = \mathbb{E}(y^4) - 3(\mathbb{E}(y^2))^2.$$

A more robust measure of non-Gaussianity involves entropy:

$$H(\mathbf{y}) = - \int f(\mathbf{y}) \log f(\mathbf{y}) d\mathbf{y}.$$

- Since Gaussian distributions have the most entropy, can compare using entropy
- Brings definition of negentropy and its maximization

Negentropy

$$J(\mathbf{y}) = H(\mathbf{y}_{\text{Gauss}}) - H(\mathbf{y})$$

- $H(\mathbf{y}_{\text{Gauss}})$ is the entropy a Gaussian variable with the same covariance as \mathbf{y}
- But this definition of negentropy is hard to work with
- Approximization is used instead, which is the foundation of FastICA

$$J(y) \propto [\mathbb{E}\{G(y)\} - \mathbb{E}\{G(\nu)\}]^2$$

- $G(y)$ is a nonquadratic function, ν is a Gaussian distribution of mean 0 and unit variance

FastICA

FastICA is a fixed point algorithm that attempts to maximize negentropy by maximizing the term $\mathbb{E}\{G(y)\}$ or $\mathbb{E}\{G(\mathbf{w}_i^\top \mathbf{x})\}$ where \mathbf{w}_i is the i -th row of \mathbf{W} .

According to Kuhn-Tucker conditions and by also assuming that $\|\mathbf{w}_i\| = 1$ due to whitening earlier, the optima of \mathbf{w}_i is where

$$\mathbb{E}\{\mathbf{x}g(\mathbf{w}_i^\top \mathbf{x})\} - \beta \mathbf{w}_i = 0$$

where $g = G'$ and $\beta = \mathbb{E}\{\mathbf{w}_{i,\text{opt}}^\top \mathbf{x}g(\mathbf{w}_{i,\text{opt}}^\top \mathbf{x})\}$ is a constant.

FastICA

Therefore the overall algorithm will go as follows.

- ① Randomize a \mathbf{w}_i vector
- ② Calculate $\mathbf{w}_i^+ = \mathbb{E}\{\mathbf{x}g(\mathbf{w}_i^\top \mathbf{x})\} - \mathbb{E}\{g'(\mathbf{w}_i^\top \mathbf{x})\}$
- ③ Normalize $\mathbf{w}_i \leftarrow \frac{\mathbf{w}_i^+}{\|\mathbf{w}_i^+\|}$
- ④ If \mathbf{w}_i has not converged then repeat step 2

Note that computing the expectation of $\mathbf{x}g(\mathbf{w}_i^\top \mathbf{x})$ and $g'(\mathbf{w}_i^\top \mathbf{x})$ is difficult so it is often easier to compute the sample mean for all the values \mathbf{x}_n within the dataset.

Section 3

Experiments

Mixed Signals

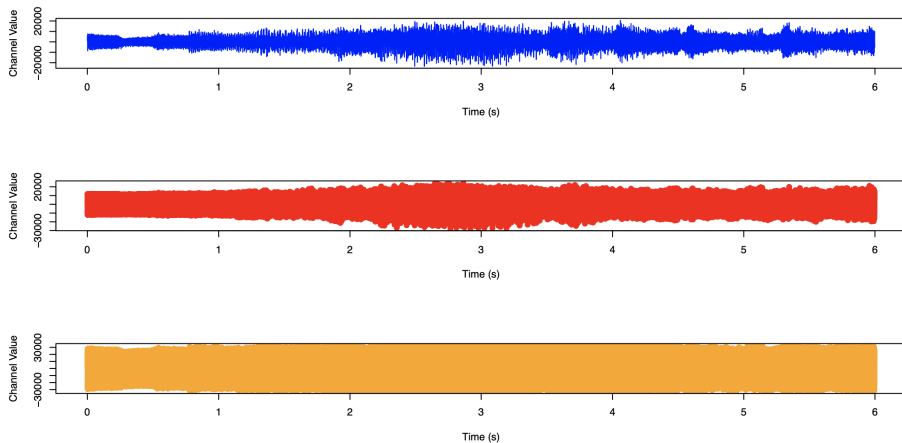


Figure 3: Recordings of Mixed Signals

Separated Sources

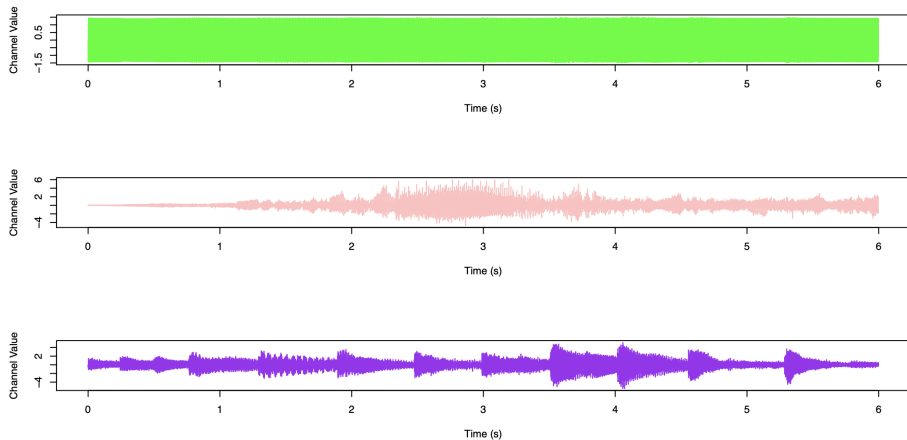


Figure 4: Separated Sources