

# Independent Component Analysis

## P9120 Group 3

Ruoyu Ji   Haokun Yuan   Apoorva Srinivasan   Tianchen Xu

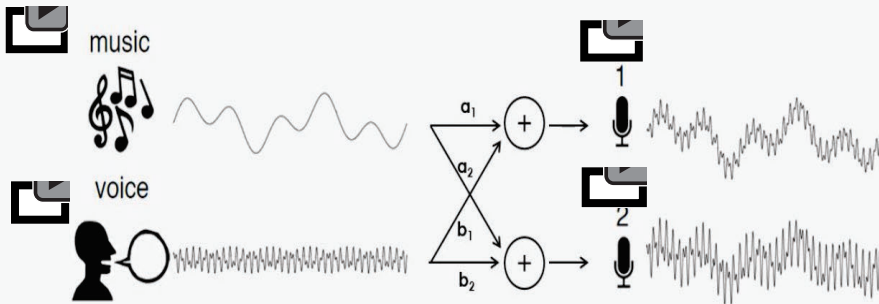
Biostatistics Department  
Columbia University

October 24, 2019

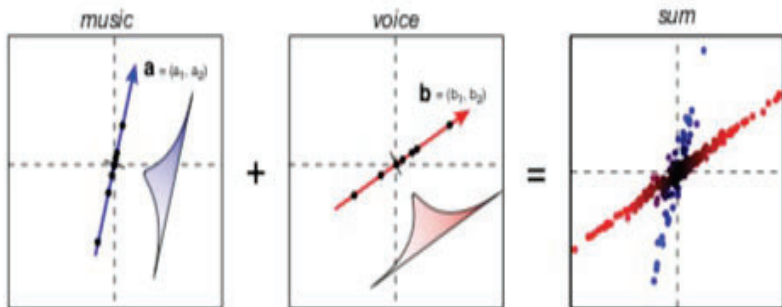
# Table of Contents

- ① Cocktail Party Problem
- ② ICA Theory
- ③ Direct Approach to ICA
- ④ Applications

# Motivating Example: Cocktail Party Problem



# Example data



# Mathematical Formulation of Independent Component Analysis

$$\mathbf{x} = \mathbf{A}\mathbf{s}$$
$$\mathbf{A} = \begin{bmatrix} a_1 & a_2 \\ b_1 & b_2 \end{bmatrix}$$

## Assumptions:

- 1) The mixing matrix is invertible
- 2) The sources are statistically independent
- 3) The independent components have non-Gaussian distributions
- 4) Data has been centered.  $\mathbf{x}$  and  $\mathbf{s}$  are zero-mean vectors

## Goal:

In this setting, the goal of ICA is to find an unmixing matrix  $\mathbf{W}$  that is appropriate  $\mathbf{A}^{-1}$  so that,  $\hat{\mathbf{s}} \approx \mathbf{s}$

$$\hat{\mathbf{s}} = \mathbf{W}\mathbf{x}$$

Seems impossible: Find two unknowns  $\mathbf{A}$  and  $\mathbf{s}$  by only observing their matrix product  $\mathbf{x}$  ????

# Singular Value Decomposition(SVD)

- Divide and conquer!
- Lets focus on **A** first.

$$\mathbf{A} = \mathbf{U}\mathbf{\Sigma}\mathbf{V}^T$$

Where,

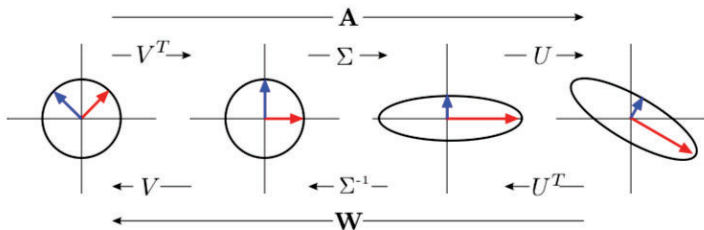
$\mathbf{U}$  = orthogonal matrix(eigenvector of  $\mathbf{A}^T\mathbf{A}$ )

$\mathbf{\Sigma}$  = Diagonal matrix non-negative diagonal entries(eigenvalue)

$\mathbf{V}$  = orthogonal matrix(eigenvector of  $\mathbf{A}\mathbf{A}^T$ )

$$\mathbf{W} = \mathbf{A}^{-1} = \mathbf{V}\mathbf{\Sigma}^{-1}\mathbf{U}^T$$

# Graphical Depiction of SVD



# Finding un-mixing matrix $W$

1. Examine the covariance of the data  $\mathbf{x}$  in order to calculate  $\mathbf{U}$  and  $\mathbf{\Sigma}$
2. Return to the assumption of independence of  $\mathbf{s}$  to solve for  $\mathbf{V}$ .



# Examining the covariance of the data

- As a reminder, the covariance is the expected value of the outer product of individual data points  $\langle \mathbf{x} \mathbf{x}^T \rangle$ .

$$\begin{aligned}\langle \mathbf{x} \mathbf{x}^T \rangle &= \langle (\mathbf{A} \mathbf{s}) (\mathbf{A} \mathbf{s})^T \rangle \\ &= \langle (\mathbf{U} \Sigma \mathbf{V}^T \mathbf{s}) (\mathbf{U} \Sigma \mathbf{V}^T \mathbf{s})^T \rangle \\ &= \mathbf{U} \Sigma \mathbf{V}^T \langle \mathbf{s} \mathbf{s}^T \rangle \mathbf{V} \Sigma \mathbf{U}^T\end{aligned}$$

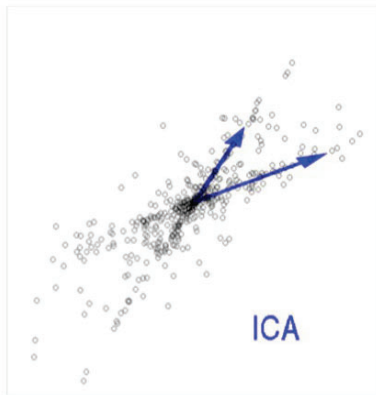
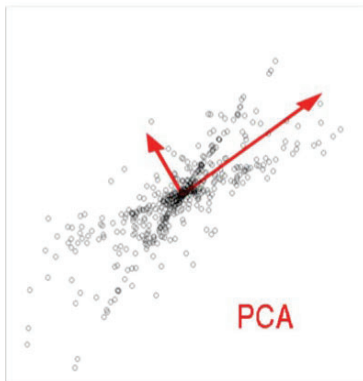
Assumption: covariance of the source  $\mathbf{s}$  is *whitened*  $\langle \mathbf{s} \mathbf{s}^T \rangle = \mathbf{I}$ ,

$$\langle \mathbf{x} \mathbf{x}^T \rangle = \mathbf{U} \Sigma^2 \mathbf{U}^T.$$

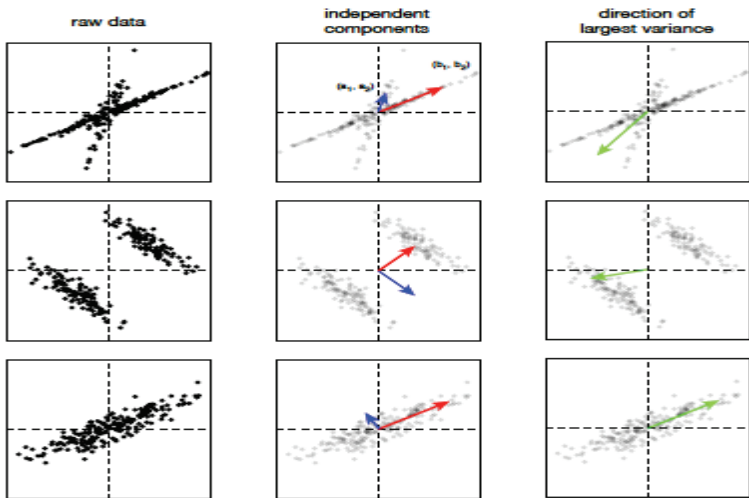
# PCA vs ICA

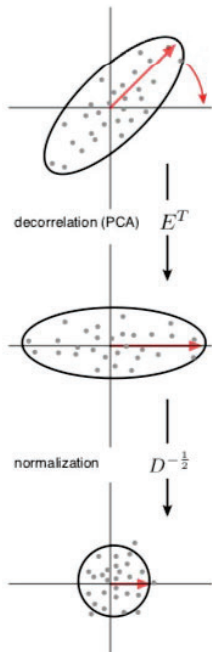
- Similarity:
  1. Feature extraction
  2. Dimension reduction
- Differences:
  - PCA just removes correlations, **not** higher order dependence ,  
ICA removes correlations, **and** higher order dependence
  - Two different sound signals need not be orthogonal as in PCA,  
even if they are independent

# PCA vs ICA



# PCA vs ICA





# Whitening

- Whitening is an operation that removes all linear dependencies in a data set (i.e. second-order correlations) and normalizes the variance along all dimensions.
- In our problem whitening simplifies the ICA problem down to finding a single rotation matrix  $V$ .

$$\hat{s} = Vx_w$$

- Where,

$$x_w = (D^{-1/2} E^T)x$$

# Estimating $\mathbf{V}$

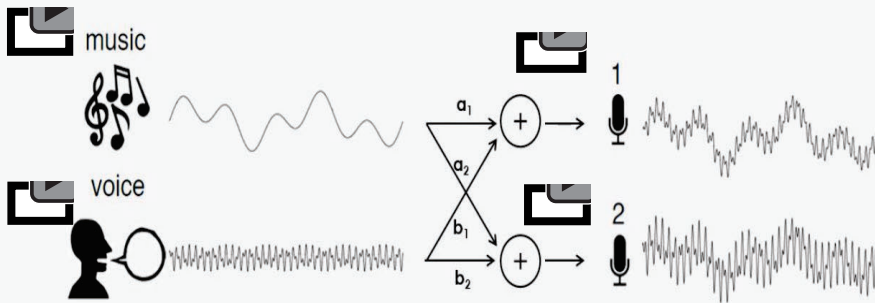
- The goal of ICA can now be stated succinctly. Find a rotation matrix  $\mathbf{V}$  such that  $\hat{\mathbf{s}}$  is statistically independent.
- Using an optimization technique from information theory, we estimate

$$\mathbf{V} = \underset{\mathbf{V}}{\operatorname{argmin}} \sum_i H[(\mathbf{V}\mathbf{x}_w)_i]$$

# Putting it all together

- We have just found  $\mathbf{U}$ ,  $\mathbf{V}$  and  $\mathbf{\Sigma}$  !!
- Original formulation:  $\mathbf{x} = \mathbf{A}\mathbf{s}$
- We needed to estimate sampled source signals using samples of observed data  $\hat{\mathbf{s}} = \mathbf{W}\mathbf{x}$
- where,  $\mathbf{W} = \mathbf{A}^{-1} = \mathbf{V}\mathbf{\Sigma}^{-1}\mathbf{U}^T$

# Results





# Notation - Factor Analysis Model

Let's introduce the notation first. Suppose observed  $\mathbf{X} \in \mathbb{R}^{N \times p} (N > p)$  can be decomposed to two matrices  $\mathbf{S} \in \mathbb{R}^{N \times p}$ ,  $\mathbf{A} \in \mathbb{R}^{p \times p}$ :

$$\mathbf{X} = \mathbf{S}\mathbf{A}^T$$

$$\Leftrightarrow \begin{pmatrix} \boxed{x_{11}} & \boxed{x_{12}} & \cdots & \boxed{x_{1p}} \\ \boxed{x_{21}} & \boxed{x_{22}} & \cdots & \boxed{x_{2p}} \\ \vdots & \vdots & \ddots & \vdots \\ \boxed{x_{Np}} & \boxed{x_{N2}} & \cdots & \boxed{x_{Np}} \end{pmatrix} = \begin{pmatrix} \boxed{s_{11}} & \boxed{s_{12}} & \cdots & \boxed{s_{1p}} \\ \boxed{s_{21}} & \boxed{s_{22}} & \cdots & \boxed{s_{2p}} \\ \vdots & \vdots & \ddots & \vdots \\ \boxed{s_{Np}} & \boxed{s_{N2}} & \cdots & \boxed{s_{Np}} \end{pmatrix} \begin{pmatrix} \boxed{a_{11}} & \cdots & \boxed{a_{p1}} \\ \boxed{a_{12}} & \cdots & \boxed{a_{p2}} \\ \vdots & \ddots & \vdots \\ \boxed{a_{1p}} & \cdots & \boxed{a_{pp}} \end{pmatrix}$$

original  $p$  variables

transformed  $p$  variables

transformation

**FA:**

*Score matrix*

*Loading matrix*

**ICA:**

*Source matrix*

*Mixing matrix*

## Notation - ICA Model

$$\mathbf{X} = \mathbf{S}\mathbf{A}^T$$

$$\Leftrightarrow \begin{pmatrix} x_{11} & x_{12} & \cdots & x_{1p} \\ x_{21} & x_{22} & \cdots & x_{2p} \\ \vdots & & & \vdots \\ x_{Np} & x_{N2} & \cdots & x_{Np} \end{pmatrix} = \begin{pmatrix} s_{11} & s_{12} & \cdots & s_{1p} \\ s_{21} & s_{22} & \cdots & s_{2p} \\ \vdots & & & \vdots \\ s_{Np} & s_{N2} & \cdots & s_{Np} \end{pmatrix} \begin{pmatrix} a_{11} & \cdots & a_{p1} \\ a_{12} & \cdots & a_{p2} \\ \vdots & & \vdots \\ a_{1p} & \cdots & a_{pp} \end{pmatrix}$$

---

$$\mathbf{X} = \mathbf{A}\mathbf{S}$$

$$\Leftrightarrow \begin{pmatrix} X_1 \\ X_2 \\ \vdots \\ X_p \end{pmatrix} = \begin{pmatrix} a_{11} & \cdots & a_{1p} \\ a_{21} & \cdots & a_{2p} \\ \vdots & & \vdots \\ a_{p1} & \cdots & a_{pp} \end{pmatrix} \begin{pmatrix} S_1 \\ S_2 \\ \vdots \\ S_p \end{pmatrix}$$

where  $\mathbf{X}$  and  $\mathbf{S}$  are  $p$  dimensional r.v.s.

# Rotation Problem

The model is unidentifiable even we constrain  $S$  to be orthonormal:

$$\begin{aligned} X &= \mathbf{A}S \\ &= \mathbf{A}\mathbf{R}^T \mathbf{R}S \\ &= \mathbf{A}^* \mathbf{S}^* \end{aligned}$$

where  $\mathbf{R}$  is any orthonormal matrices.

- *Quartimax*: To maximize the sum of all loadings raised to power 4 in  $S$ . It thus minimizes the number of factors needed to explain a variable.
- *Varimax*: To maximize variance of the squared loadings in each factor in  $S$ . As the result, each factor has only few variables with large loadings by the factor.
- .....

# Rotation - ICA Assumption

- The starting point for ICA is the very simple assumption that the components  $S_i$  are statistically independent.
- It will be shown that we must also assume that the independent component must have non-Gaussian distributions. However, in the basic model we do not assume these distributions known (if they are known, the problem is considerably simplified.)
- Then, after estimating the matrix  $\mathbf{A}$ , we can compute its inverse, say  $\mathbf{W}$ , and obtain the independent component simply by:

$$\mathbf{S} = \mathbf{A}^{-1}\mathbf{X} = \mathbf{W}\mathbf{X}.$$

# Ambiguities of ICA

- ❶ The variances of the independent components  $S_i$  cannot be determined.
  - Since both  $S$  and  $\mathbf{A}$  are unknown, any scalar multiplier of source  $S_i$  can be cancelled by dividing the corresponding column of  $\mathbf{A}$  with the same scalar value.
  - The most natural way to assume that each 'source' has unit variance  $\text{Var}(S_i^2) = 1$ .
- ❷ The order of the independent components cannot be determined.
  - Again, since  $S$  and  $\mathbf{A}$  are unknown, order of the terms in the model can be changed freely, and we can call any of the independent components the first one.

# Model Estimation

- Distribution of a sum of independent random variables tends toward a Gaussian distribution.
- Thus, a sum of two independent random variables usually has a distribution that is closer to gaussian than any of the two original random variables.

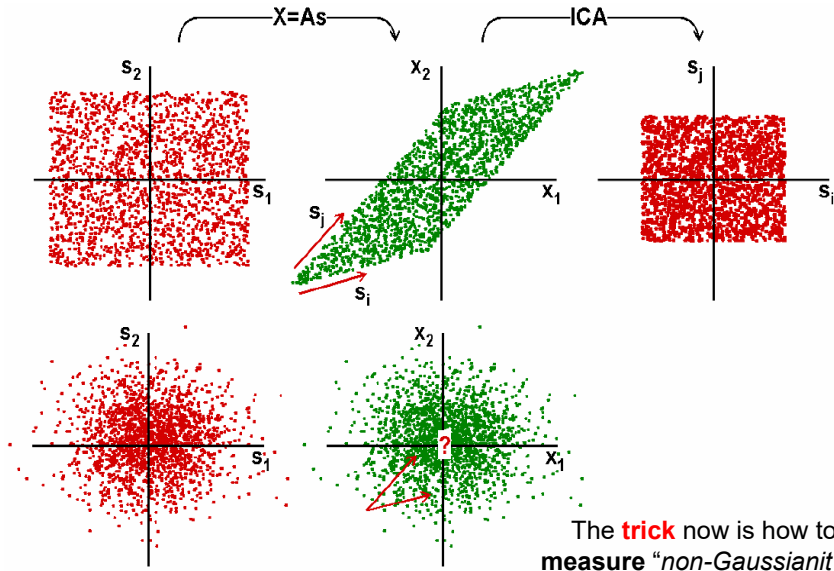
## Problem Formulation

To find a linear mapping  $\mathbf{W}$  such that the unmixed sequences  $u$ ,

$$u = \mathbf{W}^T \mathbf{X}$$

where  $\mathbf{W}$  is a row of  $\mathbf{W}$ , are maximally statistically independent.

# Necessary Non-Gaussianity



# Measurement of Non-Gaussianity

- *Kurtosis*

- Defined by:  $\text{kurt}(y) = E(y^4) - 3(E(y^2))^2$ .
- The kurtosis for a Gaussian is zero.
- Kurtosis is very sensitive to outliers when its value has to be estimated from a measured sample.

- *Differential entropy*

- Defined by:  $H(y) = - \int g(y) \ln g(y) dy$  where  $g(y)$  is the density of  $y$ .
- The Gaussian random variable has the largest entropy among all random variables of equal variance, which means that entropy can be used to measure non-gaussianity.

- *Negentropy*

- Defined by:  $J(y) = H(z) - H(y)$  where  $z$  is a normally distributed variable.



# Approximation of Negentropy

- *Classical approximation:*

$$J(y) = \frac{1}{12} E(y^3)^2 + \frac{1}{48} \text{kurt}(y)^2$$

- *New approximation based on maximum-entropy principle:*

$$J(y) = \sum_{i=1}^p k_i [E\{G_i(y)\} - E\{G_i(z)\}]^2$$

where  $k_i$  are some positive constants, and  $z$  is a Gaussian variable of zero mean and unit variance. The functions  $G_i$  are some nonquadratic functions such as:

$$G(u) = \frac{1}{a} \ln \cosh(au), \quad G(u) = -e^{-\frac{u^2}{2}}, \dots \quad (1 \leq a \leq 2)$$

# Mutual Information

- Defined by:

$$I(y_1, y_2, \dots, y_p) = \sum_{i=1}^p H(y_i) - H(y)$$

where  $y_i$  are the components of the random variable  $y$ .

- It is the natural measure of the dependence between random variables. Its value is always nonnegative, and zero if and only if the variables are statistically dependent.
- Let  $u = \mathbf{W}X$ , then

$$I(u_1, u_2, \dots, u_p) = \sum_{i=1}^p H(u_i) - H(X) - \ln |\det(\mathbf{W})|.$$

This is equivalent to minimizing the sum of the entropies of the separate components of  $u$ .

# Use of Negentropy

- Recall *Negentropy*:

Defined by:  $J(y) = H(z) - H(y)$  where  $z$  is a normally distributed. It has many approximation forms.

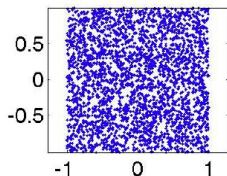
- Entropy and negentropy differ only by a constant and sign. Therefore, finding an invertible transformation  $\mathbf{W}$  that minimizes the mutual information is roughly equivalent to finding directions in which negentropy is maximized.

# Preprocessing for ICA - Centering

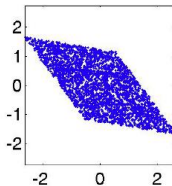
- The most basic and necessary preprocessing is to center the data matrix  $\mathbf{X}$  that is, subtract the mean vector, to make the data a zero mean variable. With this,  $S$  can be considered to be zero mean, as well.
- After estimating the mixing matrix  $A$  the mean vector of  $S$  can be added back to the centered estimates of  $S$  to complete the estimation.
- The mean vector of  $S$  is given by  $\mathbf{A}^{-1}\mu$  where  $\mu$  is the mean vector of the data matrix  $\mathbf{X}$ .

# Preprocessing for ICA - Whitening

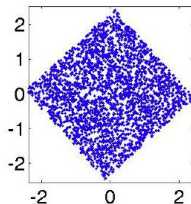
- By eigen-value decomposition (EVD):  $E(\mathbf{X}\mathbf{X}^T) = \mathbf{V}\mathbf{D}\mathbf{V}^T$ , where  $\mathbf{V}$  is the orthogonal matrix of eigenvectors and  $\mathbf{D}$  is the diagonal matrix of eigenvalues.
- transformation of  $\mathbf{X}$ :  $\tilde{\mathbf{X}} = \mathbf{V}\mathbf{D}^{-1/2}\mathbf{V}^T\mathbf{X}$ , then  $\text{Cov}(\tilde{\mathbf{X}}) = \mathbf{I}$ .
- Since  $\mathbf{S}$  has covariance  $\mathbf{I}$ , then for  $\tilde{\mathbf{X}} = \tilde{\mathbf{A}}\mathbf{S}$ , we have  $\tilde{\mathbf{A}}$  is orthogonal which only contains  $p(p-1)/2$  degrees of freedom instead of  $p^2$ .



original



mixed



whitened

# Preprocessing for ICA - Sphering

- Centering + Whitening = Sphering
- Sphering removes the first and second order statistics of the data. Both the mean and covariance are set to zero and the variance are equalized.
- Because it is a very simple and standard procedure, much simpler than any ICA algorithms, it is a good idea to reduce the complexity of the problem this way.
- We look at the eigenvalues of  $E(\mathbf{X}\mathbf{X}^T)$  and discard those that are too small. This has often the effect of reducing noise. Moreover, dimension reduction prevents over learning, which can sometimes be observed in ICA.

# A Direct Approach to ICA

Recall that many approaches to ICA , are based on minimizing the approximation on entropy.

# Product Density ICA

Independent Component by definition have a joint product density:

$$f_S(s) = \prod_{i=1}^p f_i(s_i)$$

And  $f_j$  can be represented as a tilted Gaussian Density:

$$f_j(s_j) = \phi(s_j) e^{g_j(s_j)}$$

The log-likelihood for the observed data  $X = \mathbf{A}S$ :

$$l(\mathbf{A}, \{g_j\}_1^p; \mathbf{X}) = \sum_{i=1}^N \sum_{j=1}^p [\log \phi_j(a_j^T x_i) + g_j(a_j^T x_i)]$$



# Penalized Density

$$\sum_{j=1}^p \left[ \frac{1}{N} \sum_{i=1}^N [\log \phi_j(a_j^T x_i) + g_j(a_j^T x_i)] - \int \phi(t) e^{g_j(t)} dt - \lambda_j \int \{g_j'''(t)\}^2(t) dt \right]$$

- The first enforces the density constraint  $\int \phi(t) e^{g_j(t)} dt = 1$  on any solution  $\hat{g}_j$ .
- The second is a roughness penalty, which guarantees that the solution  $\hat{g}_j$  is a quartic-spline with knots at the observed values of  $s_{ij} = a_j^T x_i$ .

# Product Density ICA Algorithm

---

**Algorithm 14.3** *Product Density ICA Algorithm: ProDenICA*

---

1. Initialize  $\mathbf{A}$  (random Gaussian matrix followed by orthogonalization).
  2. Alternate until convergence of  $\mathbf{A}$ :
    - (a) Given  $\mathbf{A}$ , optimize (14.91) w.r.t.  $g_j$  (separately for each  $j$ ).
    - (b) Given  $g_j$ ,  $j = 1, \dots, p$ , perform one step of a fixed point algorithm towards finding the optimal  $\mathbf{A}$ .
-

## Product Density ICA Algorithm 2a)

Given  $\mathbf{A}$  optimize log-likelihood function w.r.t.  $g_j$

$$\sum_{j=1}^p \left[ \frac{1}{N} \sum_{i=1}^N [\log \phi_j(a_j^T x_i) + g_j(a_j^T x_i)] - \int \phi(t) e^{g_j(t)} dt - \lambda_j \int \{g_j'''(t)\}^2(t) dt \right]$$

For simplicity, we focus on a single coordinate:

$$\frac{1}{N} \sum_{i=1}^N [\log \phi(s_i) + g(s_i)] - \int \phi(t) e^{g(t)} dt - \lambda_j \int \{g'''(t)\}^2(t) dt$$

Since the function involves integral, an approximation is needed.

## Product Density ICA Algorithm 2a)

Construct a fine grid  $L$  with values  $s_l^*$  in increments of  $\Delta$ .

$$y_l^* = \frac{\#s_i \in (s_l^* - \Delta/2, s_l^* + \Delta/2)}{N}$$

Likelihood function can be written as

$$\sum_{l=1}^L \left\{ y_l^* [\log \phi(s_l^*) + g(s_l^*)] - \Delta \phi(s_l^*) e^{g(s_l^*)} \right\} - \lambda_j \int \{g'''(t)\}^2(t) dt$$

## Product Density ICA Algorithm 2b)

Given  $\hat{g}_j$ , perform one step of fixed point algorithm towards finding the optimal  $\mathbf{A}$ .

Optimize log-likelihood function w.r.t.  $\mathbf{A}$  is equivalent to maximize:

$$C(\mathbf{A}) = \frac{1}{N} \sum_{i=1}^N \sum_{j=1}^p \hat{g}_j(\mathbf{a}_j^T \mathbf{x}_i)$$

$C(\mathbf{A})$  is the log-likelihood ratio between the fitted density and Gaussian and can be seen as a estimate of negentropy.

## 2b) Fixed Point Algorithm

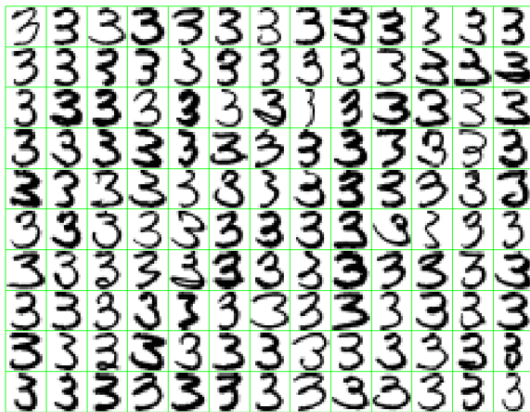
- For each  $j$  update

$$a_j \leftarrow E \left\{ X \hat{g}'_j(a_j^T x_i) - E[\hat{g}''_j(a_j^T x_i)] \right\}$$

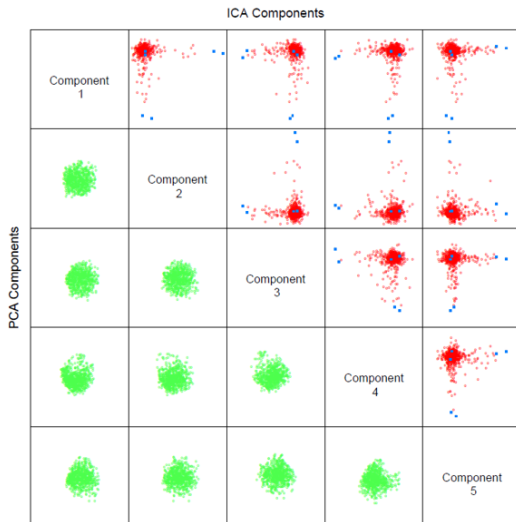
- Orthogonalize  $\mathbf{A}$  using  $(\mathbf{A}\mathbf{A}^T)^{-\frac{1}{2}}\mathbf{A}$   
Compute the SVD of  $\mathbf{A}$ ,  $\mathbf{A} = \mathbf{U}\mathbf{D}\mathbf{V}^T$ , and then replace  $\mathbf{A}$  with  $\mathbf{A} \leftarrow \mathbf{U}\mathbf{V}^T$

## Example: Handwritten Digits

- Digitized  $16 \times 16$  grayscale images
- Points in 256-dimensional space
- High-dimensional data



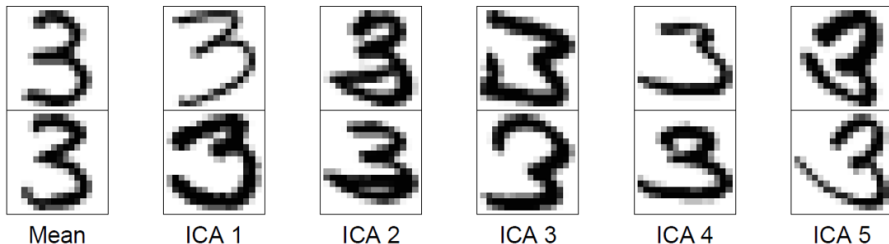
# Comparison: standardized first five ICA components vs standardized first five PCA components





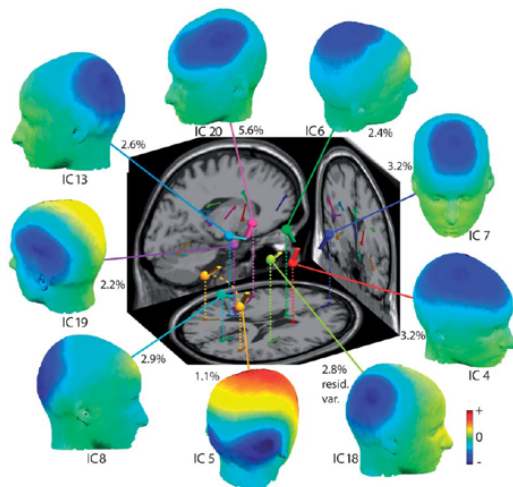
# Nature of each ICA component

- extreme digits vs central digits



# Example: EEG (electroencephalographic) Time Courses

- Domain: Brain Dynamics



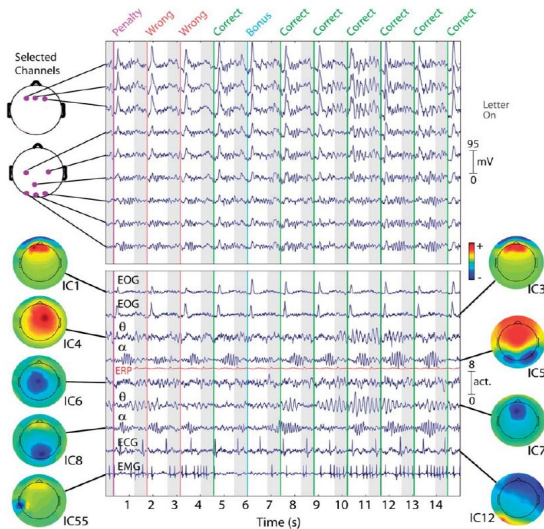
# Experiment

- Subjects wear a cap embedded with a lattice of 100 EEG electrodes, which record brain activity at different locations on the scalp.
- Data are from a single subject performing a standard “two-back” learning task over a 30-minute period, i.e., the subject is presented with a letter (B, H, J, C, F, or K) at roughly 1500-ms intervals, and responds by pressing one of two buttons to indicate whether the letter presented is the same or different from that presented two steps back.
- Depending on the answer, the subject earns or loses points, and occasionally earns bonus or loses penalty points for infrequent correct or incorrect trials, respectively.

- Goal: untangle the components of signals in multi-channel EEG data.
- Key assumption: signals recorded at each scalp electrode are a mixture of independent potentials arising from different cortical activities, as well as non-cortical artifact domains.
- ICA method

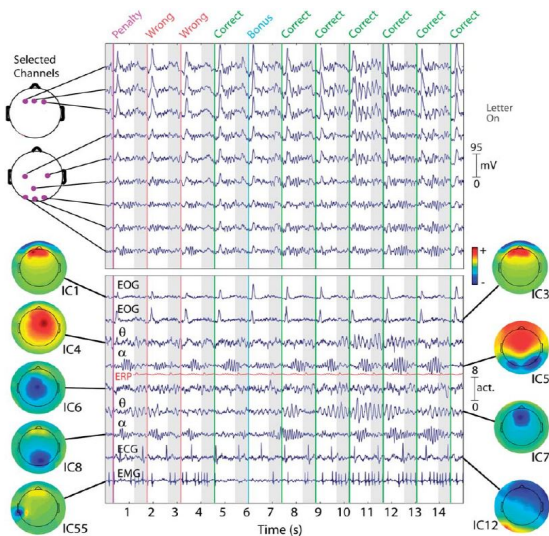
The top half:  
EEG data from 9 (of 100)  
electrodes/channels over the  
course of 15 (of 1917) seconds.

The bottom half:  
the “activities”/time courses  
of 9 (of 100) independent  
components during the same  
period.



The top half:  
time-course data show spatial  
correlation in EEG signals – the  
signals of nearby sensors look  
very similar.

The bottom half:  
ICA component activities are  
temporally distinct, i.e.,  
maximally independent over  
time.



Thank you! :)