

APPLIED MACHINE LEARNING

Independent Component Analysis (ICA)

Algorithm



ICA: Algorithm

Consider a N -dimensional observation vector $x \in \mathbb{R}^N$

Assume that x was generated by a linear variable model

$$x = As$$

A : is an **unknown** $N \times N$ mixing matrix

$s \in \mathbb{R}^N$ are **unknown** latent random variables,
referred to as the **sources**.

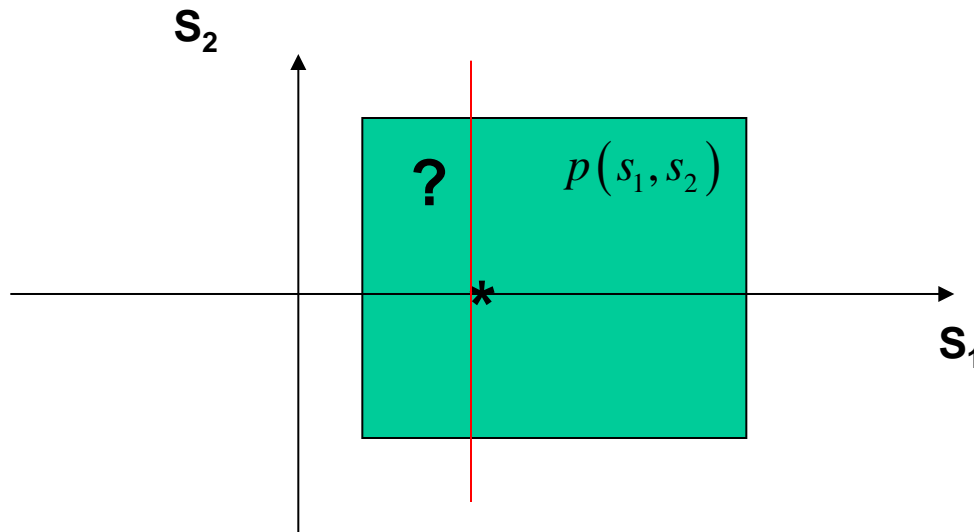
ICA estimates both the mixing matrix A and the sources s knowing only x .



1st ambiguity of ICA: x is product of A and s

ICA: Hypotheses

- Two hypotheses:
- 1: $s := \{s_1, \dots, s_N\}$ are **statistically independent**.
 - 2: The distribution of s is likely to be **non-Gaussian**.

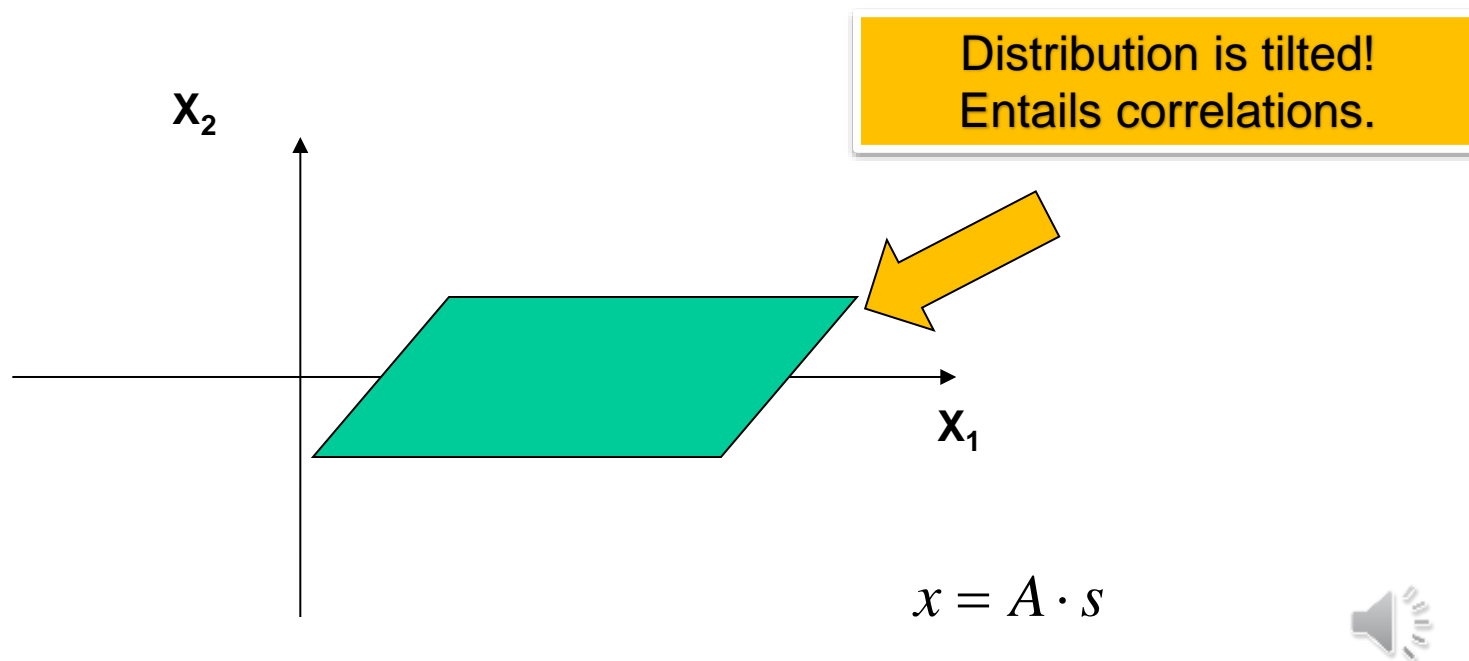


Here is an example of two independent sources s_1 and s_2 , whose distribution is non-Gaussian (here uniform distribution).



ICA: Statistical dependency

Observables X_1 and X_2 result from a linear combination of the independent sources. The distribution of the mixed observables is not statistically independent.



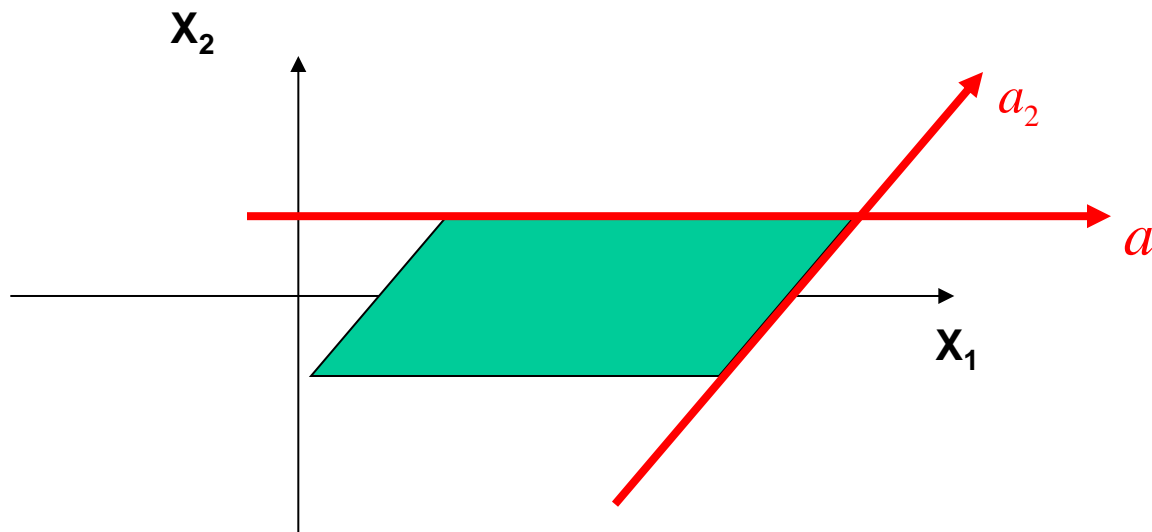
Approach: Use the correlation to find A , then compute $s = A^{-1}x$

ICA: Finding mixing matrix - intuition

Idea:

Finds the columns of the matrix A .

The vectors embeds correlation across variables



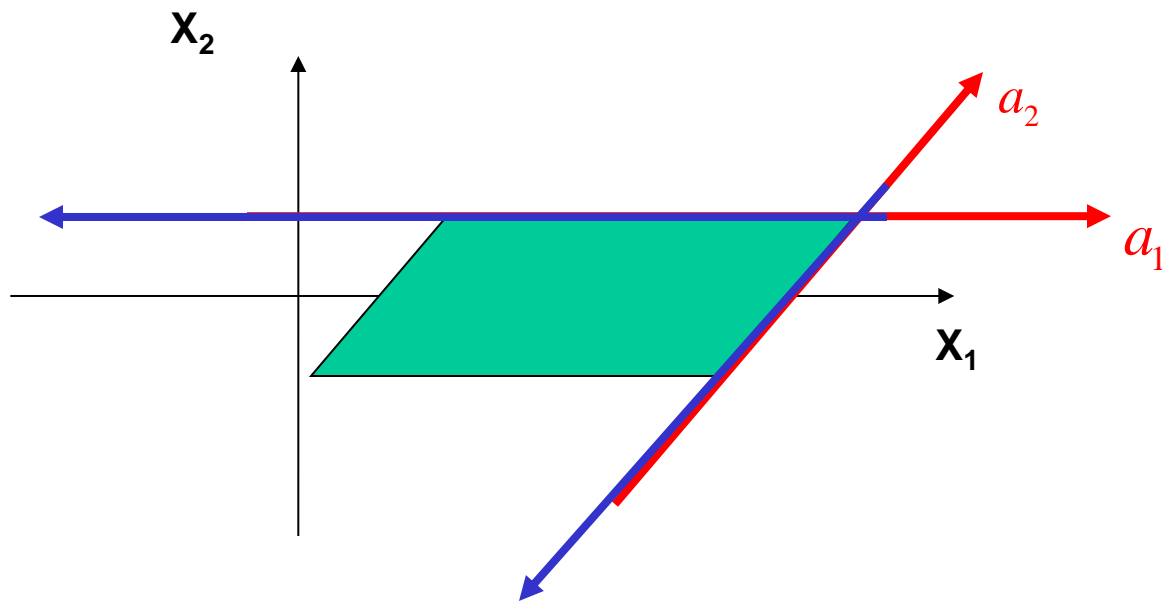
$$x = A \cdot s$$

$$A = [a_1 \ a_2]$$

Once projected onto these vectors, the distribution is statistically independent



ICA: Ambiguity

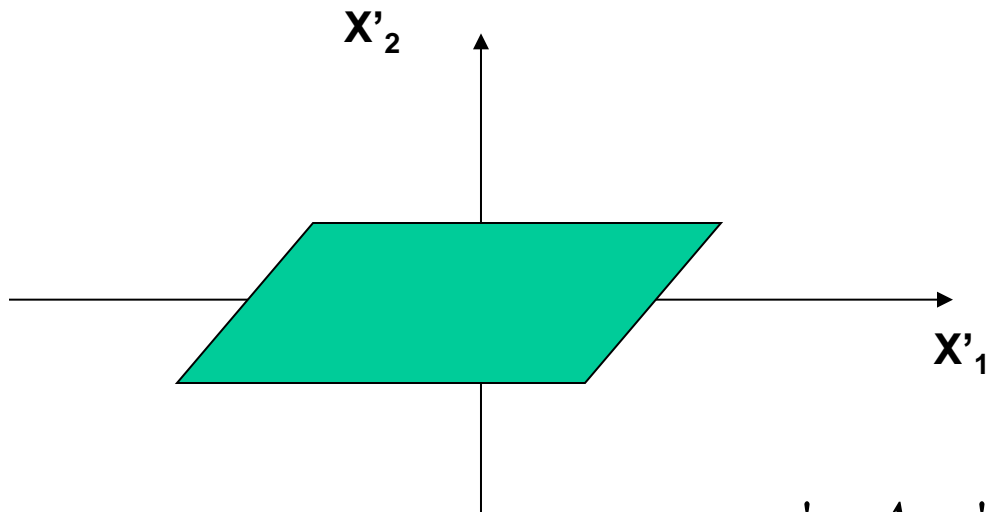


2nd ambiguity of ICA: cannot determine the sign of the vectors

ICA: Preprocessing - Centering

ICA first pre-processing step starts

by centering the distribution such that: $x' = x - E\{x\}$
 $\Rightarrow E\{x'\} = 0$



$$x' = A \cdot s'$$



$$s = s' + \mu_s, \quad \mu_s = A^{-1} \cdot E\{x\}$$

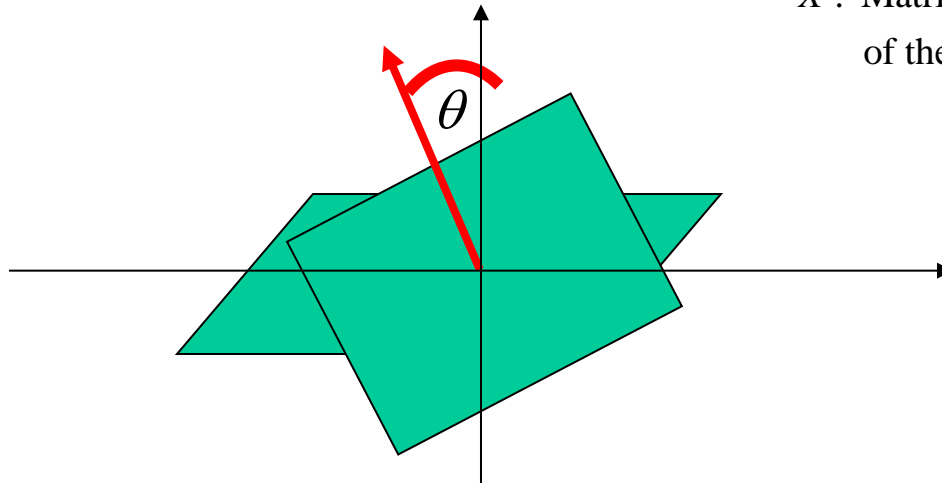
ICA: Preprocessing - Whitening

The joint distribution $p(x)$ of two variates x_1, x_2 is said to be **white** if the variates are **uncorrelated** and their **variance is equal to unity**.

$$\text{Uncorrelated: } E\{x_1, x_2\} = E\{x_1\} E\{x_2\}$$

$$\text{White: } E\{XX^T\} = I$$

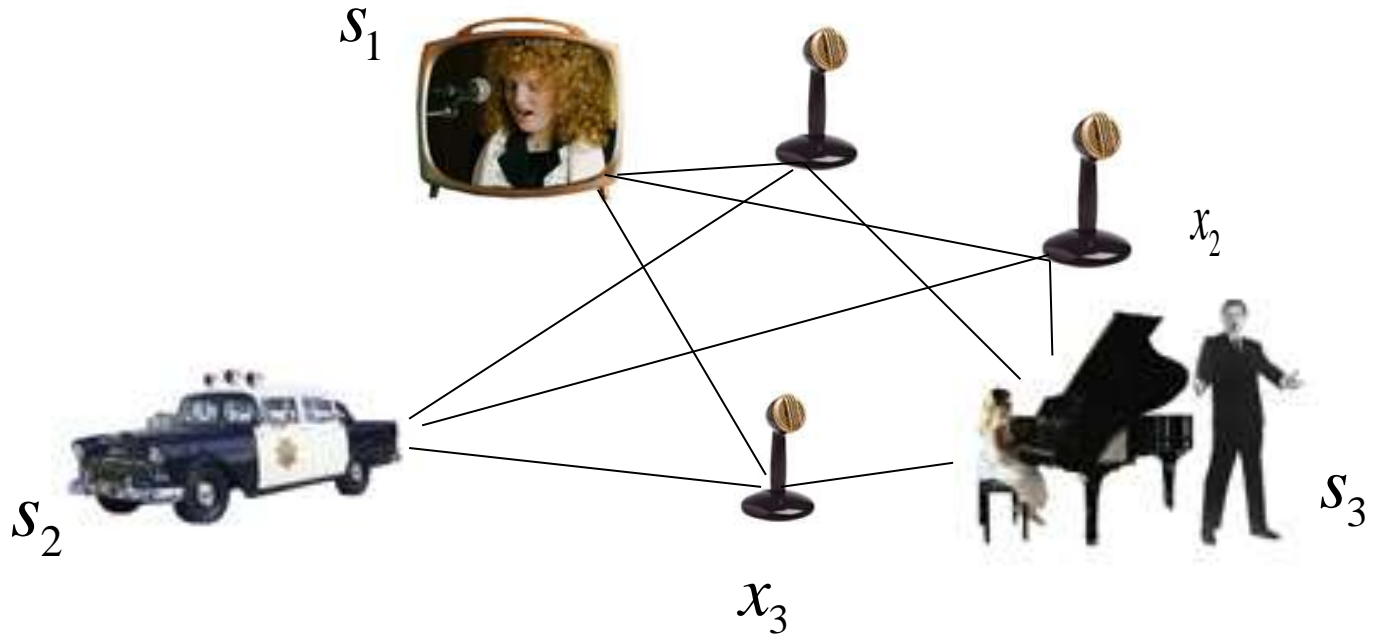
X : Matrix composed of all instances of the observations x



Advantage: Need only to compute the angle of rotation θ



N-dimensional observation vector $x \in \mathbb{R}^N$ and sources $s \in \mathbb{R}^N$, $N = 3$.



But at each time step of the recording, we obtain a new observation of both mixtures and sources.

If we do T measurements, the total dimension of the dataset is

$$X: N \times T$$



ICA: Preprocessing - Whitening

Determine a matrix V , s.t.:

$$X'' = VX'$$

$$E\{X'' X''^T\} = I$$

X : Matrix composed of all instances
of the observations x

Use the eigenvalue decomposition

$$E\{XX^T\} = EDE^T$$

$$D = \text{diag}\{\lambda_1, \dots, \lambda_N\}, \quad E = \{e^1, \dots, e^N\}$$

$$V = D^{-\frac{1}{2}} E^T$$

$$E\{X'' X''^T\} = VE\{XX^T\}V^T = D^{-\frac{1}{2}} E^T E\{XX^T\} ED^{-\frac{1}{2}}$$

$$E\{X'' X''^T\} = D^{-\frac{1}{2}} E^T EDE^T ED^{-\frac{1}{2}} = I$$

Whitening is essentially decorrelation followed by scaling.

Whitening is bringing us closer to unmix the observables and find the sources.

However, it only guarantees to find projections that decorrelate the data.

It does not guarantee to find projections such that the dataset becomes statistically independent. 

ICA: Determining the sources

The distribution of s is **non-Gaussian**.

Find sources how joint distribution optimizes
a measure of non-gaussianity
Negentropy or Kurtosis



ICA: Optimizing – Fast ICA

Estimating the negentropy is difficult as it requires
to estimate the distribution of the data

→ Minimizes Kurtosis instead

$$\text{kurtosis}(x) = E\{x^4\}; \text{kurtosis}(s) = E\{s^4\} - 3(E\{s^2\})^2$$

Iterative process: candidate source $y = s$?

Look for an optimum of $J(y) \propto \left[E\{(w^T x)^4\} - 3E\{(w^T x)^2\} \right]^2$

Since data is white and uncorrelated $E\{xx^T\} = I \Rightarrow E\{y^2\} = 1$

Search optimum for: $J(y) = E\{(w^T x)^4\} - 3E\{\|w\|^4\}$



Fast – ICA optimization

For any non-linear transformation G , we would have:

$$E\{G\{s_1 \cdot s_2\}\} = E\{G\{s_1\}\} \cdot E\{G\{s_2\}\}$$

$$\text{Set to minimize: } J(y) = E\left\{G(w^T y)^4\right\} - 3E\left(\|w\|^4\right)$$

It is important to choose well the G functions.

- Non-quadratic functions, that do not grow too fast.
- Asymmetric derivative

Good candidates are :

$$G_1(y) = \frac{1}{a} \log(\cosh(a \cdot y)), \quad 1 \leq a \leq 2$$

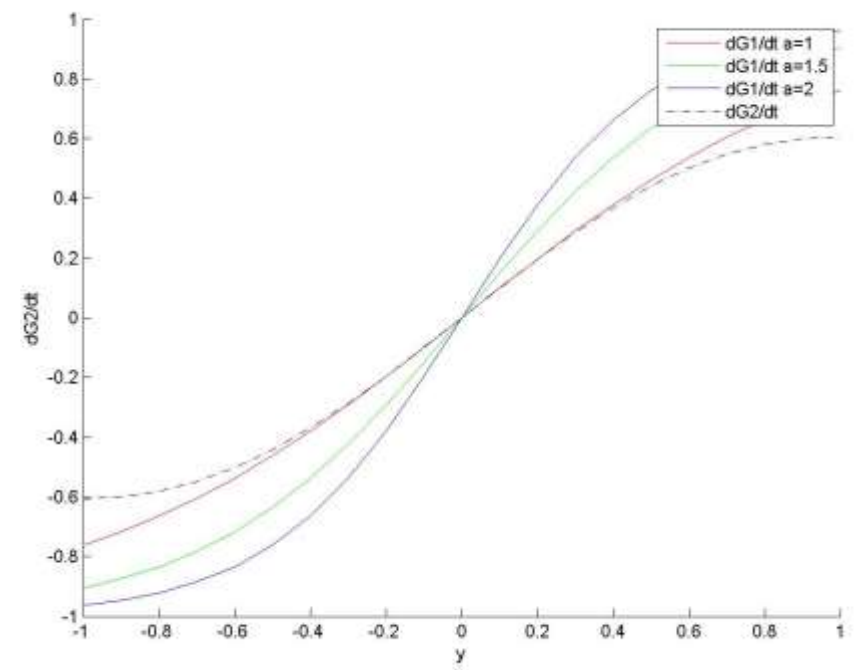
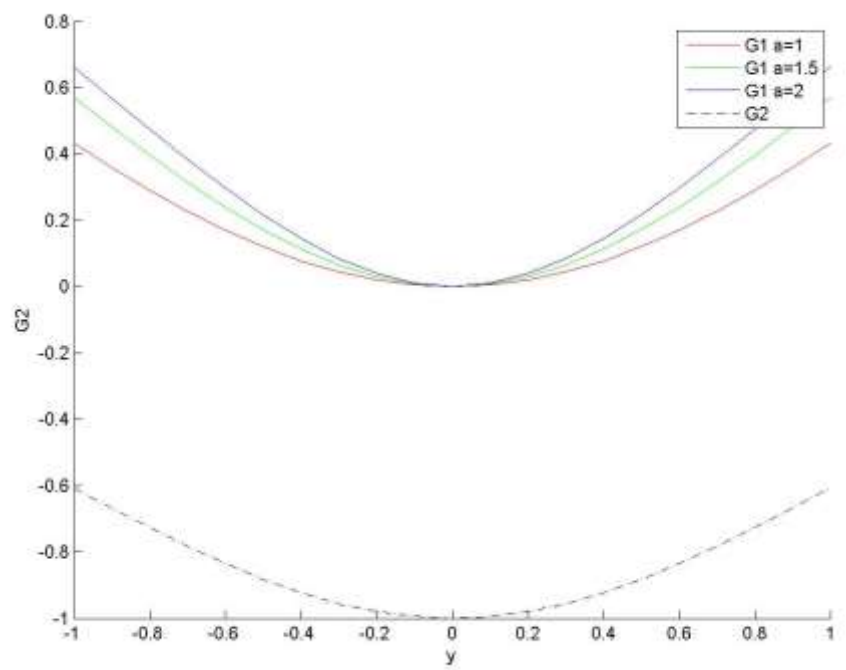
$$G_2(y) = -e^{-\frac{1}{2}y^2}$$

$$\frac{dG_1(y)}{dy} = g_1(y) = \tanh(a \cdot y) \quad 1 \leq a \leq 2$$

$$\frac{dG_2(y)}{dy} = g_2(y) = y \cdot e^{-\frac{1}{2}y^2}$$



Fast – ICA optimization



Good candidates are :

$$G_1(y) = \frac{1}{a} \log(\cosh(a \cdot y)), \quad 1 \leq a \leq 2$$

$$G_2(y) = -e^{-\frac{1}{2}y^2}$$

$$\frac{dG_1(y)}{dy} = g_1(y) = \tanh(a \cdot y) \quad 1 \leq a \leq 2$$

$$\frac{dG_2(y)}{dy} = g_2(y) = y \cdot e^{-\frac{1}{2}y^2}$$



Fast-ICA: Determining the sources iteratively

Idea: find each source one after the other one (iterative process)

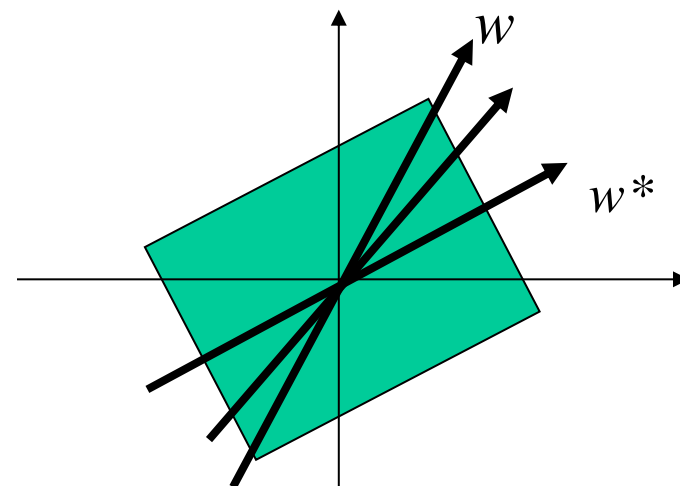
Let w be a candidate projection, such that:

$$y = w^T x = w^T A s = z s, \quad z = w^T A$$

y is more gaussian than s

(central limit theorem), unless $y = s$.

Search w , to seek optimum: $\frac{dJ(y = w^T x)}{dw} = 0$



=> iterate until one finds a w that maximizes the non-gaussianity of y .



Fast – ICA optimization

Let x'' be the centered and whitened projection of the data

1. Choose an initial (e.g. random) weight vector w .
2. Compute the quantity $w^+ = E \{ x G(w^T x) \} - E \{ g(w^T x) \} w$ (derivative of J)
3. Proceed to a normalization of the weight vector:

$$w = \frac{w^+}{\|w^+\|}$$

$$E \{ yy^T \} = 1 \Leftrightarrow w^T E \{ xx^T \} w = w^T w = 1$$

4. If the weight has not converged, that is if: $w^T(t-1) \cdot w(t) \neq 1$
go back to step 2.



ICA: iteration for several sources

To estimate several independent components, we must prevent the different vectors from converging to the same maximum.

→ *decorrelate* the vectors at each iteration.

1. Estimate the independent components one by one w_1, \dots, w_p
2. Run the one-unit fixed-point algorithm for w_{p+1}
3. after every iteration step subtract from w_{p+1} the other "projections"

$$w_{p+1}^T w_j w_j, \quad j = 1, \dots, p$$

1. Let $\mathbf{w}_{p+1} = \mathbf{w}_{p+1} - \sum_{j=1}^p \mathbf{w}_{p+1}^T \mathbf{w}_j \mathbf{w}_j$

2. Let $\mathbf{w}_{p+1} = \mathbf{w}_{p+1} / \sqrt{\mathbf{w}_{p+1}^T \mathbf{w}_{p+1}}$



ICA Estimation and Limitations

ICA Estimation

Centering (mean=0)

Whitening (variance = 1)

Fast-ICA for one component

Fast-ICA for several components +
decorrelation at each time step

Limitations of ICA

Reconstructed sources
may be inverted
(sign of vector unknown)

Source vectors cannot be ordered
according to importance
(unlike PCA)

