

# Lecture 16

# Independent Component Analysis

Jochen Braun

Otto-von-Guericke-Universität Magdeburg,  
Cognitive Biology Group

Theoretical Neuroscience II  
SS 2020

Credits: <https://research.ics.aalto.fi/ica/icademo/>  
Hyvarinen & Oja, 2000

## 16. Independent component analysis (ICA)

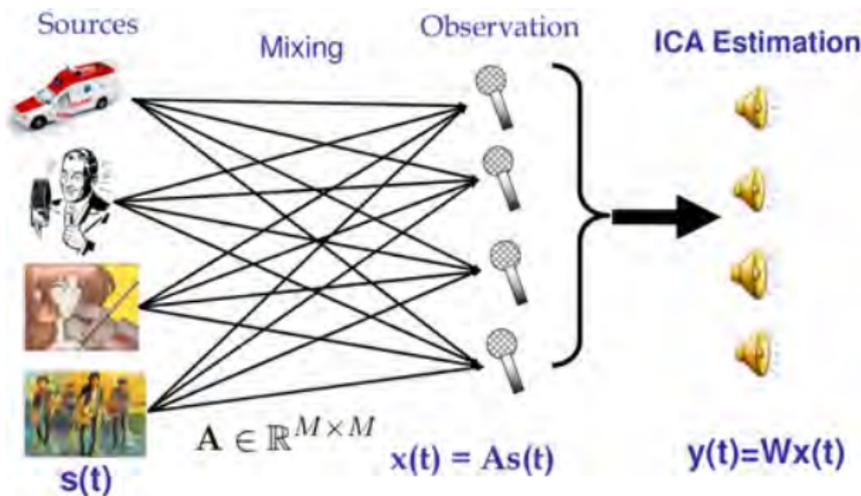
*Typical natural causes have rich statistics (more complex than Gaussian). **Blind source separation** decomposes observed signal mixtures into independent and maximally non-Gaussian sources (causes). Statistical dependencies may remain even after signals have been whitened (scaled to zero mean and unit variance) and decorrelated (by Principal Component Analysis). **ICA** decomposes signals such as to eliminate these remaining dependencies (in higher moments), maximizing both independence and non-Gaussianity. This heuristic approach can be justified by information theory. **Fast ICA** is a convenient computational implementation. ICA has many successful applications, for example, in functional brain imaging.*

# Overview

1. **Blind source separation**
2. **ICA demo**
3. **Theoretical justification**
4. **Processing steps of ICA**
5. **Imaging applications**

# 1. Blind source separation

Consider a 'cocktail party' situation, with a cacophony of sounds, musical performances, and voices ...



Assume that multiple microphones record various mixtures (linear combinations)  $x(t)$  of these independent sources  $s(t)$ .

## Other examples

The blind source separation problem occurs in many contexts:

- ▶ Electrical recordings of brain activity (EEG). Each scalp electrode measures a mixture of potential sources (dipoles) inside the brain.
- ▶ To better understand neural processes, it would be useful to identify independent potential sources.
- ▶ Natural signals and images may be a mixture of prototypical components (sound or image features).
- ▶ To efficiently represent and analyze natural signals and images, it would be useful to have a dictionary of such components.

## Mixing matrix

Observed time-courses  $\mathbf{x}(t) = \mathbf{x}_i$  are a linear mixture of latent sources  $\mathbf{s}(t) = \mathbf{s}_i$ , where  $i = 1 \dots n$  indexes time-samples:

$$\mathbf{x}_i = \mathbf{A} \cdot \mathbf{s} = \begin{pmatrix} \mathbf{a}_1 \cdot \mathbf{s}_i \\ \mathbf{a}_2 \cdot \mathbf{s}_i \\ \vdots \\ \mathbf{a}_n \cdot \mathbf{s}_i \end{pmatrix} = \begin{pmatrix} a_{11}s_{1i} + \dots + a_{1m}s_{mi} \\ a_{21}s_{1i} + \dots + a_{2m}s_{mi} \\ \vdots \\ a_{n1}s_{1i} + \dots + a_{nm}s_{mi} \end{pmatrix}$$

where  $n$  is the number of observed time-courses,  $m$  is the number of latent sources, and  $k$  is the number of time-samples.

$$\mathbf{A} \in \mathcal{R}^{n \times m}, \quad \mathbf{x} \in \mathcal{R}^{n \times k}, \quad \mathbf{s} \in \mathcal{R}^{m \times k}$$

Typically, we assume that the number of observations  $n$  is at least as large as the number of sources  $m$ .

# Several ambiguities

As neither mixing matrix  $\mathbf{A}$  nor sources  $\mathbf{s}$  are known, the situation is intrinsically ambiguous:

- ▶ We cannot know the absolute variance (power) of any source  $\mathbf{s}_j$ , because we observe only  $a_{1j}\mathbf{s}_{jk}$ ,  $a_{2j}\mathbf{s}_{jk}$ ,  $\dots$ ,  $a_{nj}\mathbf{s}_{jk}$ .
- ▶ Accordingly, we assume that all sources have unit variance (power).
- ▶ We cannot order or rank the sources in any way, because any permutation yields the same observations.
- ▶ Accordingly, we accept any ordering of sources as equivalent.

# Unmixing matrix

We seek an unmixing matrix  $\mathbf{W}$  which renders the recovered sources  $\mathbf{y}(t) = \mathbf{y}_i$  as independent as possible:

$$\mathbf{y}_i = \mathbf{W} \cdot \mathbf{x} = \begin{pmatrix} \mathbf{w}_1 \cdot \mathbf{x}_i \\ \mathbf{w}_2 \cdot \mathbf{x}_i \\ \vdots \\ \mathbf{w}_m \cdot \mathbf{x}_i \end{pmatrix} = \begin{pmatrix} w_{11}x_{1i} + \dots + w_{1n}x_{ni} \\ w_{21}x_{1i} + \dots + w_{2n}x_{ni} \\ \vdots \\ w_{m1}x_{1i} + \dots + w_{mn}x_{ni} \end{pmatrix}$$

$$\mathbf{W} \in \mathcal{R}^{m \times n}, \quad \mathbf{x} \in \mathcal{R}^{n \times k}, \quad \mathbf{y} \in \mathcal{R}^{m \times k}$$

The number of recovered sources  $m$  cannot be larger than the number of observations  $n$ !

# Independence

Two variables  $y_{1,2}$  are independent if the joint probability density is separable

$$p(y_1, y_2) = p(y_1) p(y_2)$$

or, equivalently, if the mutual information is zero

$$I_m = \int \int p(y_1, y_2) \ln \frac{p(y_1, y_2)}{p(y_1) p(y_2)} dy_1 dy_2 = 0$$

or, equivalently, if expectations of all transformations are separable

$$E \{h_1(y_1)h_2(y_2)\} = E \{h_1(y_1)\} E \{h_2(y_2)\}, \quad \forall h_1(y), h_2(y)$$

# Uncorrelatedness

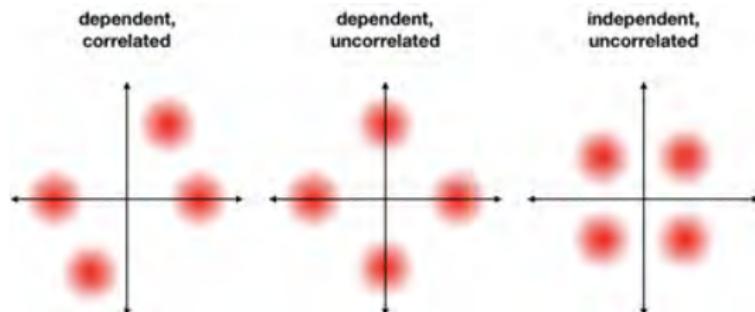
Uncorrelatedness is a far weaker constraint than independence.

Two variables  $y_{1,2}$  are uncorrelated if the covariance is zero:

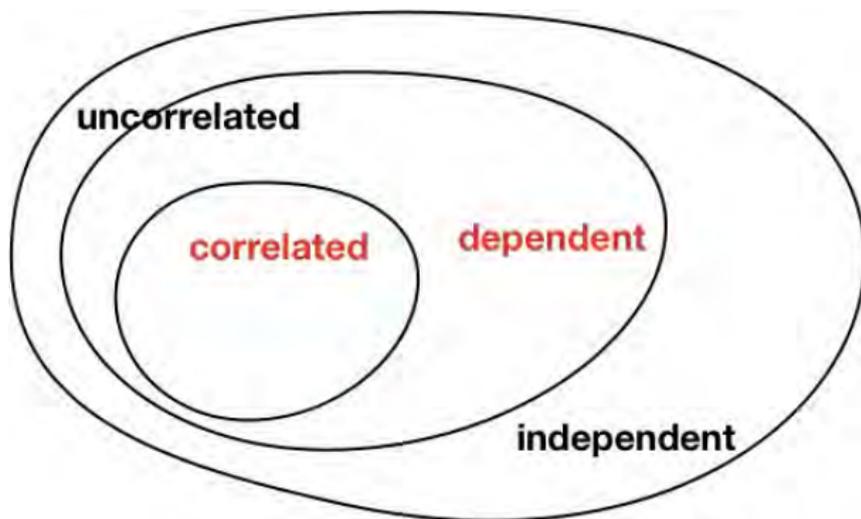
$$E \{y_1 y_2\} - E \{y_1\} E \{y_2\} = 0 \quad \Leftrightarrow \quad E \{y_1 y_2\} = E \{y_1\} E \{y_2\}$$

This does not imply separability of higher moments

$$E \{y_1^k y_2^k\} \neq E \{y_1^k\} E \{y_2^k\} \quad k > 1$$



# Correlation and dependence



# Central limit theorem

Random variables may be 'Gaussian' or 'non-Gaussian' or anything in between.

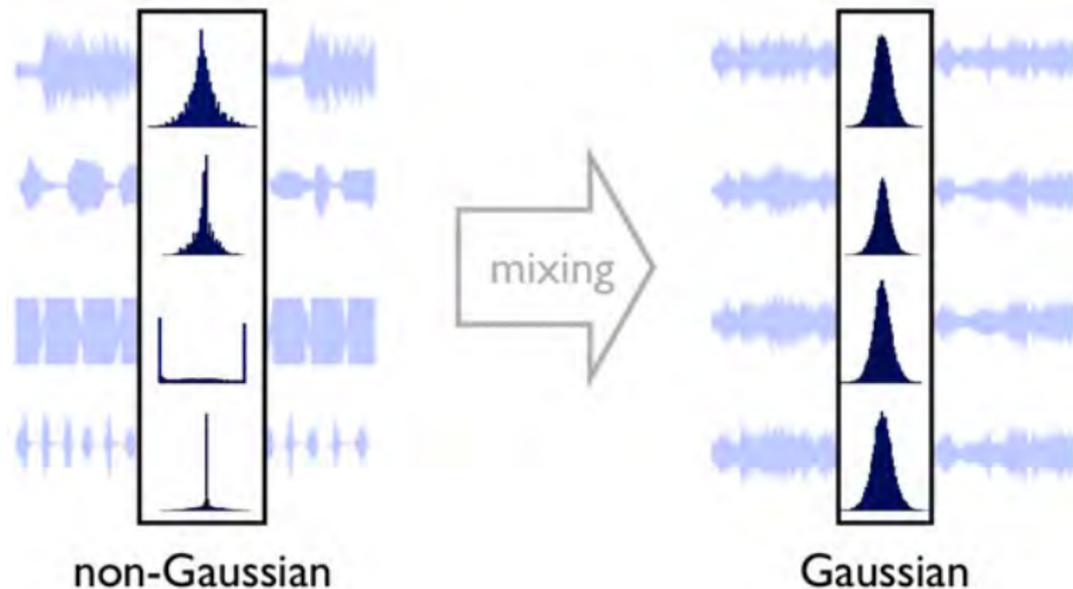


sources



mixtures

Linear combinations of *independent* 'non-Gaussian' are *more* 'Gaussian' than the original variables.



## Independent = maximally non-Gaussian

Given unknown mixing matrix  $\mathbf{x} = \mathbf{A} \cdot \mathbf{s}$  and linear combination (attempted partial unmixing):  $y = \mathbf{w}^T \cdot \mathbf{x}$ .

Compare our attempt at unmixing  $\mathbf{w}^T$  with the ground truth  $\mathbf{A}^T$  (true unmixing):

$$\mathbf{z} = \mathbf{A}^T \cdot \mathbf{w} \quad \Rightarrow \quad y = \mathbf{w}^T \cdot \mathbf{x} = \mathbf{w}^T \cdot \mathbf{A} \cdot \mathbf{s} = \left( \mathbf{A}^T \cdot \mathbf{w} \right)^T \cdot \mathbf{s} = \mathbf{z}^T \cdot \mathbf{s}$$

As a linear combination of  $s_i$ ,  $y$  is more Gaussian than any  $s_i$ . To be maximally non-Gaussian,  $y$  must be equal to one  $s_i$ .

Thus, by choosing  $\mathbf{w}$  such as to maximize non-Gaussianity of

$$y = \mathbf{w}^T \cdot \mathbf{x} = \mathbf{z}^T \cdot \mathbf{s}$$

we can identify one independent component!

## Procedure of blind source separation

- ▶ Assume that sources are *independent* and have *non-Gaussian* statistics (dependent and/or Gaussian sources would not be distinguishable).
- ▶ Assume that mixing matrix is square and orthogonal (number of observations equals number of independent sources).
- ▶ Find first unmixing vector  $\mathbf{w}$  such as to maximize non-Gaussianity of  $\mathbf{w}^T \cdot \mathbf{x}$ .
- ▶ This identifies first independent component  $s_1$  ( $-s_1$  is also a solution).
- ▶ Find second unmixing vector  $\mathbf{v}$ , orthogonal to  $\mathbf{w}$ , such as to maximize non-Gaussianity of  $\mathbf{v}^T \cdot \mathbf{x}$ .
- ▶ This identifies second independent component  $s_2$  ( $-s_2$  is also a solution).
- ▶ Repeat until the number of 'independent components' equals the number of observations.

## Points to note

- ▶ Given a linear mixture of unknown signals, 'blind source separation' seeks an 'unmixing' matrix.
- ▶ Some ambiguities cannot be resolved: signal power, dependent signals, Gaussian signals.
- ▶ Important distinction between dependence (all moments) and correlatedness (only second moment).
- ▶ Generally, *mixing* signals makes statistics *more* Gaussian.
- ▶ Thus, a matrix that makes statistics *less* Gaussian, is presumably an *unmixing* matrix.
- ▶ 'Less Gaussian' implies non-zero higher moments.
- ▶ Abstract procedure of 'blind source separation'

## 2. ICA Demo

[cis.legacy.ics.tkk.fi/aapo/papers/IJCNN99\\_tutorialweb/](http://cis.legacy.ics.tkk.fi/aapo/papers/IJCNN99_tutorialweb/)

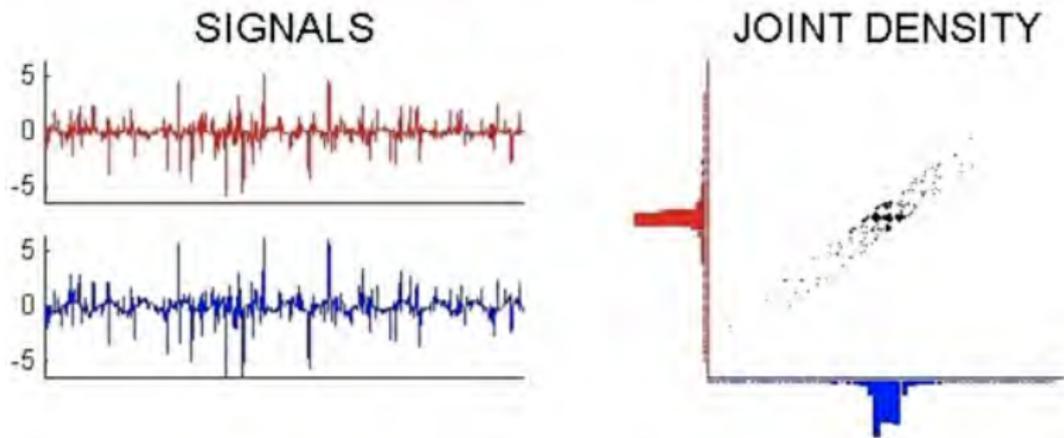
Here, we demonstrate ICA for solving the blind source separation problem. We are given two linear mixtures of two source signals which we know to be independent of each other, i.e. observing the value of one signal does not give any information about the value of the other. We seek to determine the source signals given only the mixtures.

Putting this into mathematical notation, we model the problem by

$$\mathbf{x} = \mathbf{A} \cdot \mathbf{s}$$

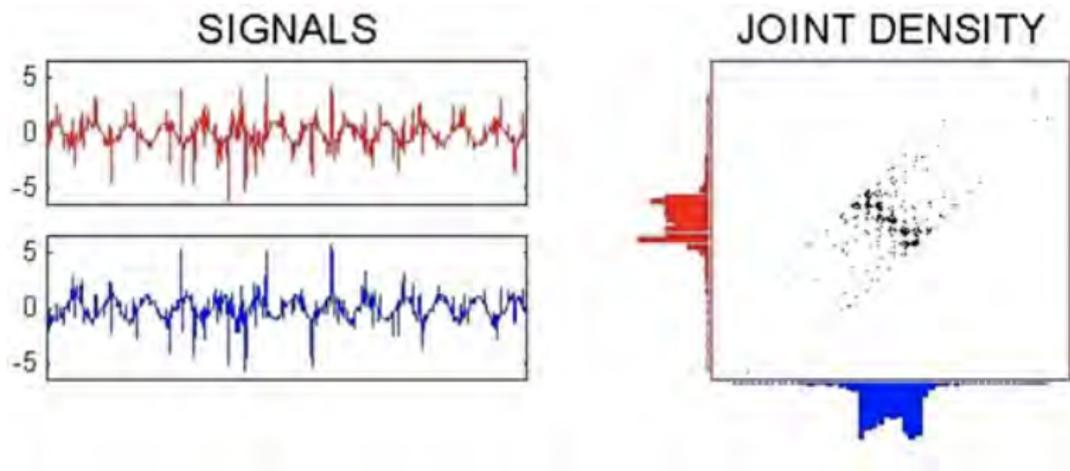
where  $\mathbf{s}$  is a two-dimensional random vector containing the independent source signals,  $\mathbf{A}$  is the two-by-two mixing matrix, and  $\mathbf{x}$  contains the observed (mixed) signals.

This first plot (below) shows the signal mixtures on the left and the corresponding joint density plot on the right. That is, at a given time instant, the value of the top signal is the first component of  $x$ , and the value of the bottom signal is the corresponding second component. The plot on the right is then simply constructed by plotting each such point  $x$ . The marginal densities are also shown at the edge of the plot.



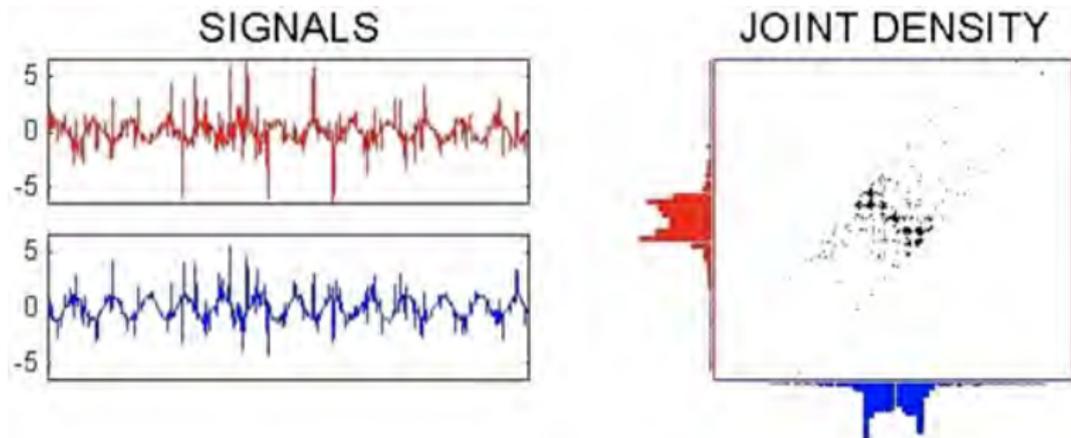
**Input signals and density**

The first step to remove any correlations in the data ('whitening' or PCA). We seek a linear transformation  $\mathbf{V}$  such that when  $\mathbf{y} = \mathbf{V} \cdot \mathbf{x}$  we obtain  $E(\mathbf{y} \cdot \mathbf{y}^T) = \mathbf{I}$ . To this end, we choose  $\mathbf{V} = \mathbf{C}_x^{-1/2}$  as the 'scaled eigenvectors' of the covariance matrix  $\mathbf{C}_x = (\mathbf{x} \cdot \mathbf{x}^T)$ , ensuring  $E(\mathbf{y} \cdot \mathbf{y}^T) = E(\mathbf{V} \cdot \mathbf{x} \mathbf{x}^T \cdot \mathbf{V}^T) = \mathbf{C}_x^{-1/2} \cdot \mathbf{C}_x^{-1/2} = \mathbf{I}$ . The figure shows the whitened signal  $\mathbf{y}$  and joint density  $p(\mathbf{y})$



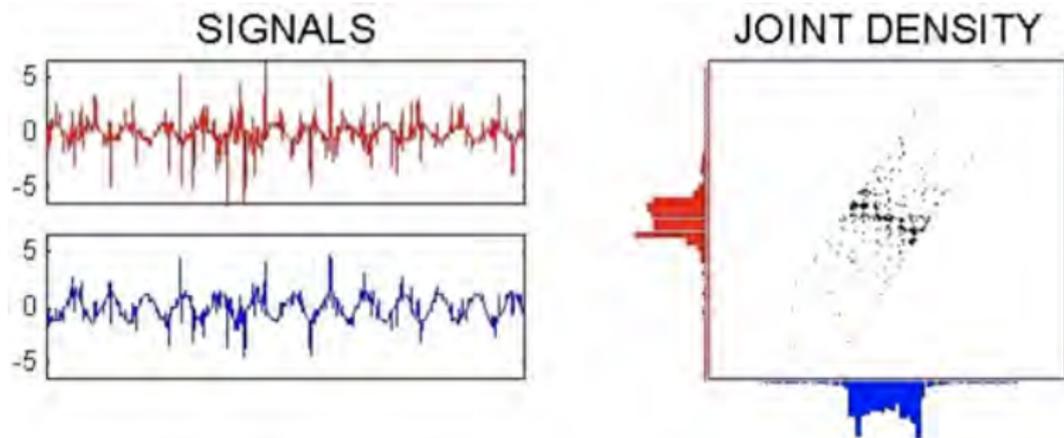
**Whitened signals and density**

After whitening, we seek to increase the degree of independence by performing a rotation (orthogonal transformation). The appropriate rotation maximizes the non-Gaussianity of the marginal densities.



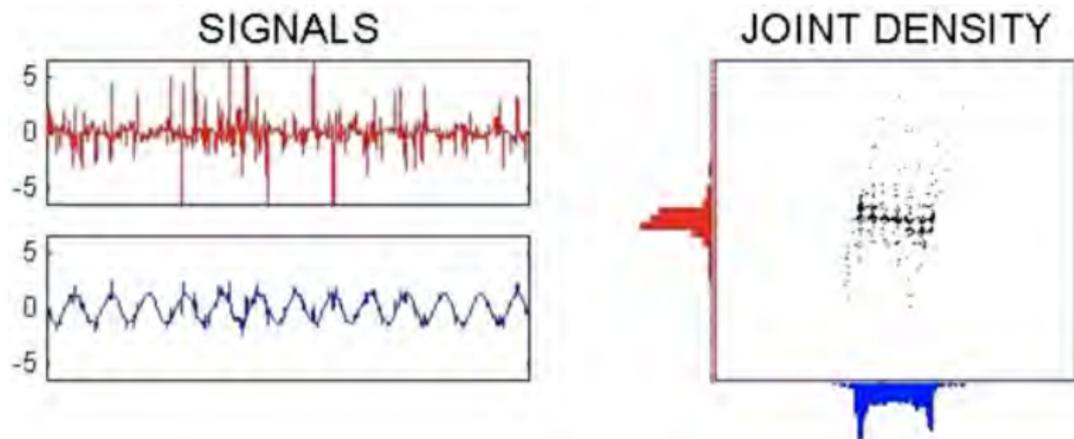
**Separated signals after 1 step of FastICA**

There are many algorithms for performing ICA, but a particularly efficient one is FastICA, which was developed by Hyvarinen and Oja (1999).



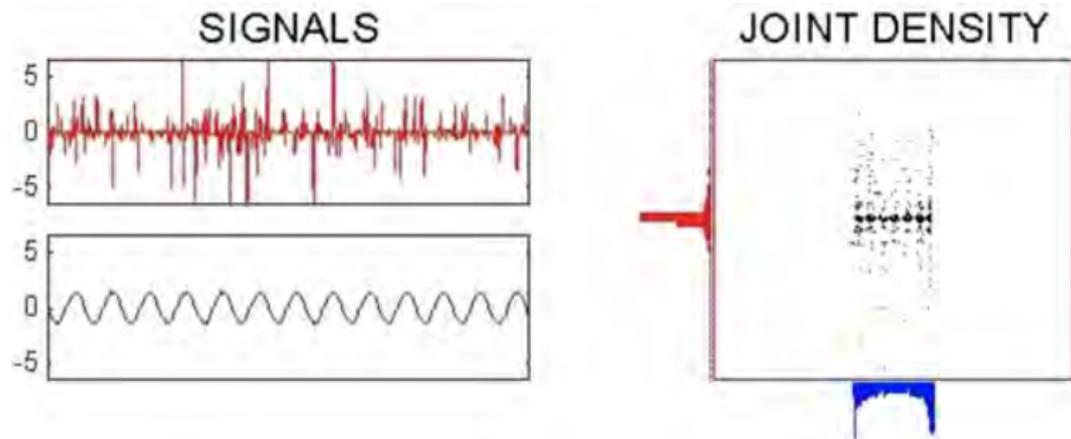
**Separated signals after 2 steps of FastICA**

.. the rotation continues ...



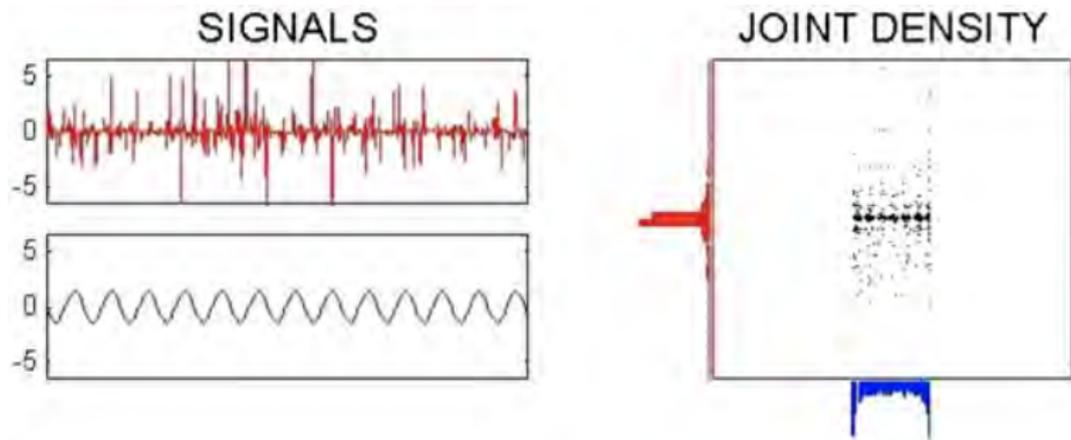
**Separated signals after 3 steps of FastICA**

... until it begins to converge ...



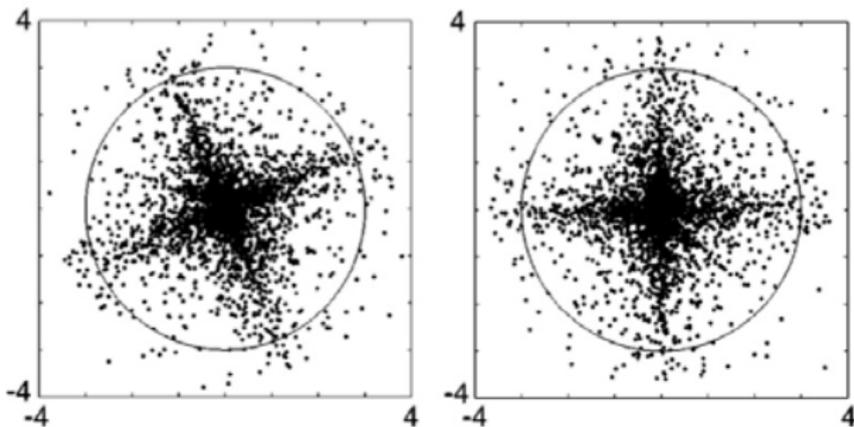
**Separated signals after 4 steps of FastICA**

Convergence! In this example, the two source signals were a sinusoid and impulse noise. The joint density can be seen to be the product of marginal densities,  $p(\mathbf{y}) = p(y_1)p(y_2)$ . Thus we have recovered two independent source signals.



**Separated signals after 5 steps of FastICA**

ICA rotates coordinates such as to maximize higher-order moments and independence!



## Points to note

- ▶ Fast ICA is a computational scheme for solving the blind source separation problem.
- ▶ Hyvarinen and Oja (2000) illustrated its operation graphically.
- ▶ First step is to whiten and decorrelate signals (perform principal component analysis, PCA).
- ▶ In general this, leaves residual dependencies in higher-order moments.
- ▶ Several steps of Fast ICA rotate coordinate system.
- ▶ After convergence, both higher-order moments and independence are maximal.

### 3. Theoretical justification

Several quantitative measures of nongaussianity are possible:

- ▶ Excess kurtosis
- ▶ Negentropy
- ▶ Approximate negentropy
- ▶ Mutual information
- ▶ Maximum likelihood
- ▶ Infomax principle

We assume a random variable  $y$  with zero-mean and unit-variance:

$$E\{y\} = 0, \quad E\{y^2\} = 1$$

## Excess kurtosis

$$\text{kurt}(y) = E \{y^4\} - 3 (E \{y^2\})^2 = E \{y^4\} - 3$$

$$\text{kurt}(y_1 + y_2) = \text{kurt}(y_1) + \text{kurt}(y_2), \quad \text{kurt}(\alpha y_1) = \alpha^4 \text{kurt}(y_1)$$

In terms of the transformation  $\mathbf{z} = \mathbf{A}^T \cdot \mathbf{w}$ , we seek the values  $z_{1,2}$  maximizing the kurtosis of  $y = z_1 s_1 + z_2 s_2$ :

$$|\text{kurt}(y)| = |\text{kurt}(z_1 s_1)| + |\text{kurt}(z_2 s_2)| = z_1^4 |\text{kurt}(s_1)| + z_2^4 |\text{kurt}(s_2)|$$

with

$$E \{y^2\} = z_1^2 + z_2^2 = 1$$

It can be shown that the maxima are where one element of  $\mathbf{z}$  is 0 and the other is  $\pm 1$ .

# Negentropy

Negentropy measures how much less entropy a variable has than a comparable Gaussian variable. It is based on differential entropy

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

A Gaussian variable has the largest differential entropy of all variables with constant variance. Negentropy is defined as the entropy difference to a Gaussian variable with identical covariance:

$$J(\mathbf{y}) = H(\mathbf{u}_{gauss}) - H(\mathbf{y})$$

Negentropy is well justified by statistical theory. However, negentropy is computationally difficult and requires the (high-dimensional) density  $p(\mathbf{y})$ .

## Approximate negentropy

Improved approximations of negentropy were developed by Hyvarinen (1998):

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

where  $k_i$  are positive constants,  $\nu$  is a Gaussian variable of zero-mean and unit variance, and functions  $G_i$  are suitable non-quadratic functions.

Note that  $G(y) = y^4$  retrieves a kurtosis-based approximation

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

Suitable choices of  $G$  yield robust estimates of negentropy, e.g.:

$$G_1(u) = \frac{1}{a_1} \log \cosh(a_1 u), \quad G_2(u) = -\exp(-u^2/2)$$

## Mutual information

The mutual information between  $m$  random variables  $y_i$  is

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

It is equivalent to the Kullback-Leibler divergence between the joint density  $p(\mathbf{y})$  and the product of marginal densities  $\prod_{i=1}^m p(y_i)$ . For an invertible transform  $\mathbf{y} = \mathbf{W} \cdot \mathbf{x}$ , it can be shown that

$$I(y_1, y_2, \dots, y_m) = \sum_{i=1}^m H(y_i) - H(\mathbf{x}) - \log |\det \mathbf{W}|$$

Mutual information and negentropy differ only by a constant and the sign:

$$I(y_1, y_2, \dots, y_m) = C - \sum_i J(y_i)$$

ICA of a random vector  $\mathbf{x}$  is an invertible transformation  $\mathbf{y} = \mathbf{W} \cdot \mathbf{x}$  which minimizes the mutual information of the transformed (unmixed) components  $y_i$ .

Finding this invertible transformation is roughly equivalent to finding the 1-D subspaces with maximum negentropy.

Rigorously, ICA estimation by minimization of mutual information is equivalent to maximizing the sum of nongaussianities of estimates, where the estimates are constrained to be uncorrelated.

# Points to note

- ▶ ICA is based on a heuristic approach: maximizing non-gaussianity (higher order moments).
- ▶ Several related theoretical arguments justify this approach.
- ▶ We outlined excess kurtosis, negentropy, and mutual information.

## 4. Processing steps of ICA

The processing steps are as follows:

- ▶ Centering
- ▶ Whitening
- ▶ Filtering
- ▶ Fast ICA (repeatedly)
- ▶ Decorrelation (repeatedly)

# Centering

Observations  $\mathbf{x}$  must have zero mean. If necessary, this can be ensured by subtracting the mean:

$$\mathbf{x} = \mathbf{x}_{ori} - \mathbf{m}, \quad \mathbf{m} = E\{\mathbf{x}_{ori}\}$$

After estimating the mixing matrix  $\mathbf{A}$  and zero-mean sources  $\mathbf{s}$ , we can add the mean back to obtain non-zero-mean sources

$$\mathbf{s}_{ori} = \mathbf{s} + \mathbf{A}^T \cdot \mathbf{m}$$

## Whitening

Observations  $\mathbf{x}$  must be 'whitened' by rotation to the principle component axes. Recall eigenvalue decomposition of covariance matrix

$$E \left\{ \mathbf{x} \cdot \mathbf{x}^T \right\} = \mathbf{E} \cdot \mathbf{D} \cdot \mathbf{E}^T, \quad \tilde{\mathbf{x}} = \mathbf{E} \cdot \mathbf{D}^{-1/2} \cdot \mathbf{E}^T \cdot \mathbf{x}$$

$$\Rightarrow \quad E \left\{ \tilde{\mathbf{x}} \cdot \tilde{\mathbf{x}}^T \right\} = \mathbf{I}$$

Whitening transforms the mixing matrix  $\mathbf{A}$  to a new one,  $\tilde{\mathbf{A}}$ , which is orthogonal:

$$\tilde{\mathbf{x}} = \mathbf{E} \cdot \mathbf{D}^{-1/2} \cdot \mathbf{E}^T \cdot \mathbf{x} = \mathbf{E} \cdot \mathbf{D}^{-1/2} \cdot \mathbf{E}^T \cdot \mathbf{A} \cdot \mathbf{s} = \tilde{\mathbf{A}} \cdot \mathbf{s}$$

$$\Rightarrow \quad E \left\{ \tilde{\mathbf{x}} \cdot \tilde{\mathbf{x}}^T \right\} = \tilde{\mathbf{A}} \cdot E \left\{ \mathbf{s} \cdot \mathbf{s}^T \right\} \cdot \tilde{\mathbf{A}}^T = \tilde{\mathbf{A}} \cdot \tilde{\mathbf{A}}^T = \mathbf{I}$$

# Fast ICA (repeatedly)

Seek unit vector  $\mathbf{w}$  that maximizes nongaussianity of projected observations  $\mathbf{w}^T \cdot \mathbf{x}$ , as measured by negentropy  $J(\mathbf{w}^T \cdot \mathbf{x})$ .

Denote by  $g_{1,2}$  the derivatives of non-quadratic functions  $G_{1,2}$ , for example

$$g_1(u) = \tanh(u), \quad g_2(u) = u \exp(-u^2/2)$$

1. Choose an initial (random) vector  $\mathbf{w}$ .
2. Rotate  $\mathbf{w}^+ = E \{ \mathbf{x} g(\mathbf{w}^T \cdot \mathbf{x}) \} - E \{ g'(\mathbf{w}^T \cdot \mathbf{x}) \}$
3. Renormalize  $\mathbf{w} = \mathbf{w}^+ / \|\mathbf{w}^+\|$ .
4. Repeat until converged.

Convergence means that old and new  $\mathbf{w}$  are almost identical (dot product almost unity).

For derivation see Hyvarinen & Oja (2000).

# Decorrelation (repeatedly)

To estimate several independent components, we need to iterate with several orthogonal projection vectors  $\mathbf{w}_1, \mathbf{w}_2, \dots, \mathbf{w}_n$ .

To prevent different vectors from converging to the same maxima, we must decorrelate the outputs  $\mathbf{w}_1 \cdot \mathbf{x}, \mathbf{w}_2 \cdot \mathbf{x}, \dots, \mathbf{w}_n \cdot \mathbf{x}$  after every iteration, such as to keep them orthogonal!

There are several methods for achieving this.

# Advantages of Fast ICA

According to its inventor, the advantages of Fast ICA include:

- ▶ Convergence is faster (cubic or quadratic) than for gradient descent (linear).
- ▶ No step-size parameter.
- ▶ Finds directly independent components of practically any non-Gaussian distribution, using any non-linearity  $g$ .
- ▶ Method can be optimized by choosing non-linearity  $g$ .
- ▶ Independent components may be estimated one by one, reducing computational load in an exploratory analysis.
- ▶ Matlab implementation is available.

## 5. Example applications

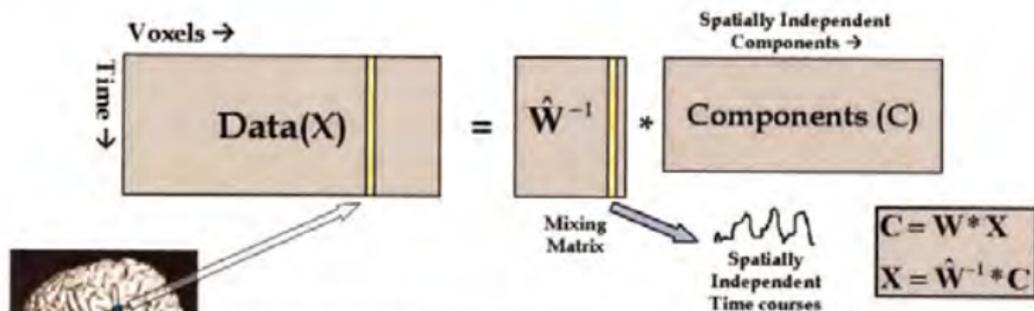
ICA is useful in many brain recording and brain imaging applications, for example

- ▶ Spatial ICA
- ▶ Temporal ICA
- ▶ Artifact removal

# Spatial and temporal ICA

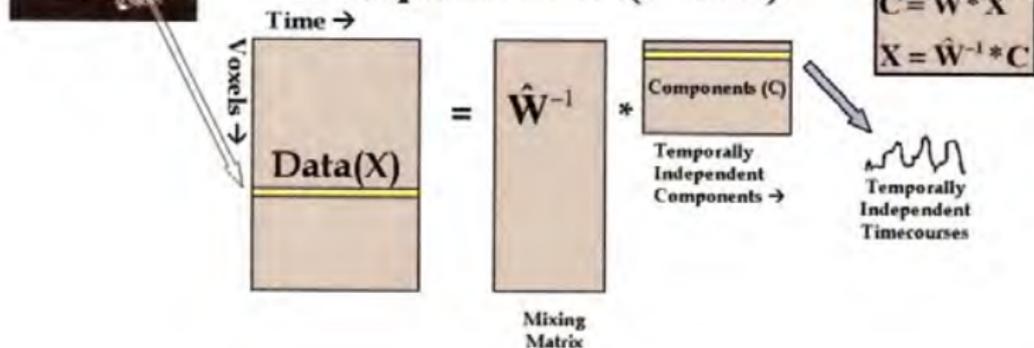
a

## Spatial ICA (SICA)

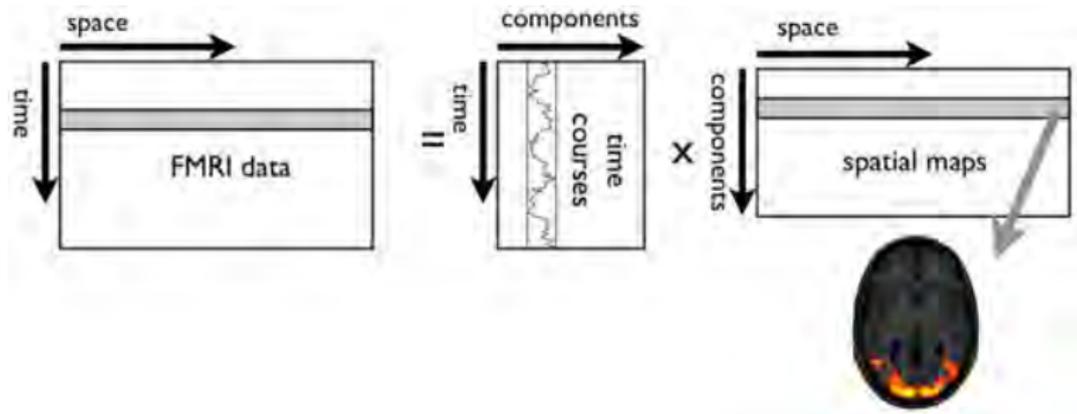


b

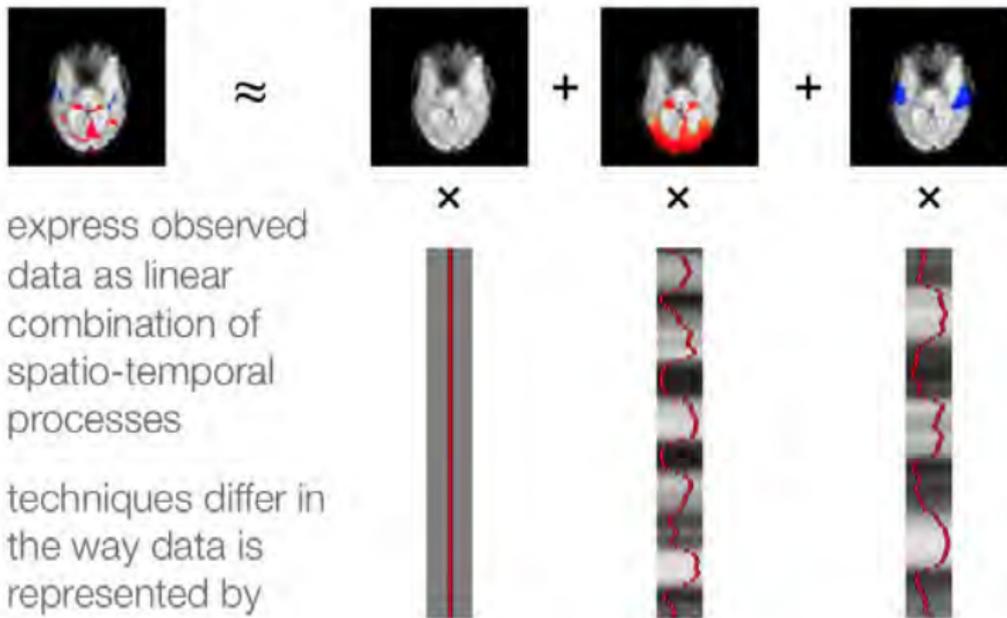
## Temporal ICA (TICA)



# Independent spatial components, modulated in time

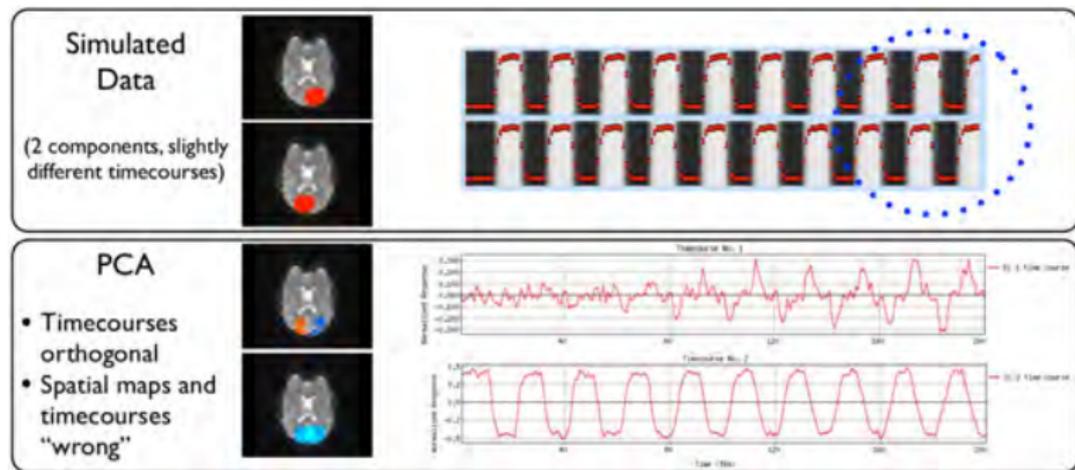


Caution: typical number of voxels is far greater than number of time-points!

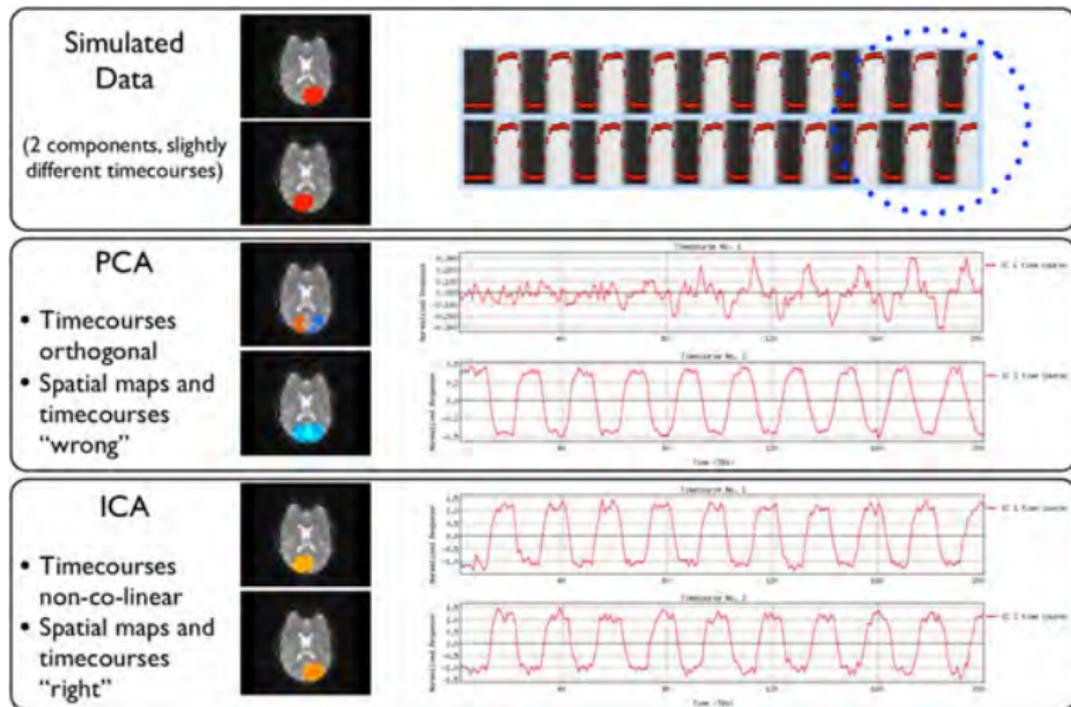


- express observed data as linear combination of spatio-temporal processes
- techniques differ in the way data is represented by components

# Comparison PCA and ICA



Orthogonality assumption decomposes phase-shifted time-courses into one time-course (first principal component) and one differential (secondary principal component). Associated spatial maps depict region of commonality and regions of differential.



ICA correctly identifies independent (but non-orthogonal) time-courses and spatial maps.

## Points to note

- ▶ ICA is used widely in signal processing.
- ▶ Functional brain imaging is a popular area of application.
- ▶ Spatial ICA, temporal ICA, probabilistic ICA (PICA).
- ▶ As all computational methods, ICA works well only when it is appropriate.

**The aim of this course was to enable you to judge for yourself which computations methods are, or are not, appropriate for the analysis of your observations!**

**Moral:**

**Trust yourself, not authority / literature.**