

Lecture 3:

Linear Methods of Machine Learning:

2) Independent component analysis (ICA)

- 1. Independent Component Analysis (ICA) - introduction**
- 2. ICA - Cocktail party problem**
- 3. ICA - task definition**
- 4. ICA versus PCA**
- 5. ICA as maximization of nongaussianity**
- 6. ICA as Likelihood maximization**
- 7. ICA as minimization of Mutual Information**

Slide from Lecture 1

All original (p -dimensional) features (components of high-dimensional vectors) are actually be governed by a few q , $q < p$, simple variables (“hidden causes”, “latent variables”, etc.)

Panoramic image $X \in \mathbb{R}^p$ ($p = 2^{12} \div 2^{18}$)
captured by mobile robot at position $\theta \in \mathbb{R}^3$ is
determined by latent variable θ : $X = \varphi(\theta)$



- all captured panoramic images (features) $\mathbf{M} = \{X\}$ captured by mobile robot at all possible positions $\Theta = \{\theta\}$ is 3-dimensional surface $\mathbf{M} = \{X = \varphi(\theta): \theta \in \Theta \subset \mathbb{R}^3\} \subset \mathbb{R}^p$ in high dimensional space \mathbb{R}^p
- any 3-dimensional one-to-one parameterization $Z = F(X) \in \mathbb{R}^3$ of the surface \mathbf{M} gives the new reduced features

Feature space is q -dimensional surface in \mathbb{R}^p ‘parameterized’ by latent variables

Original multiple features $\{X\}$ are indirect measurements of small number of an underlying sources which typically cannot be directly measured

Independent component analysis (ICA)

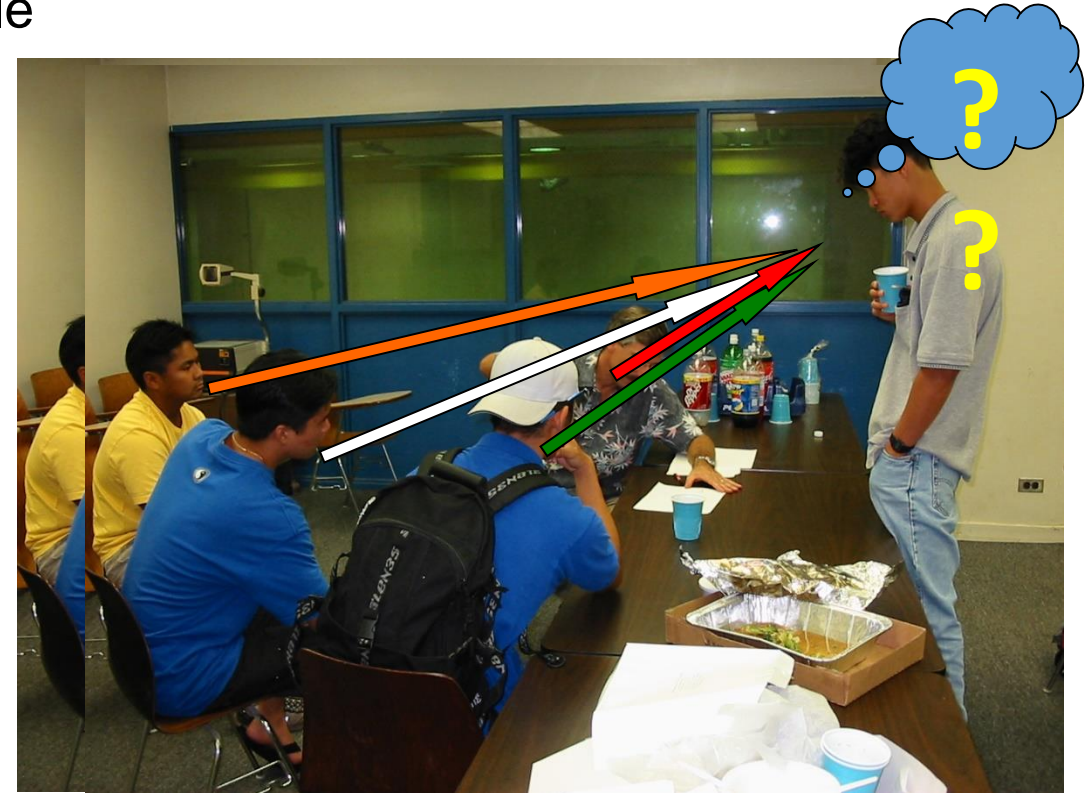
is a method for finding underlying **statistically independent** latent variables (factors, hidden causes, components) (Jeanny Hérault and Bernard Ans, 1984)

Most popular example: the **Cocktail party problem**

called also **Blind Source Separation** (BSS)

Cocktail party problem

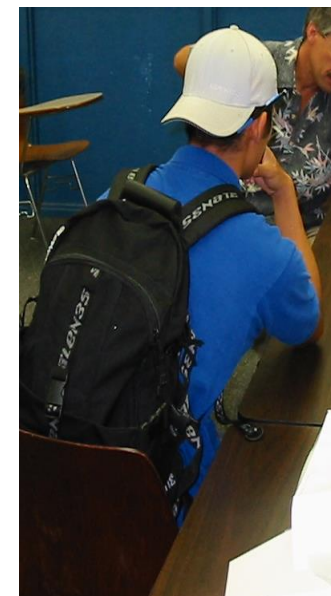
Suppose you are in a crowded room with many people



How do you understand what any one person is saying?



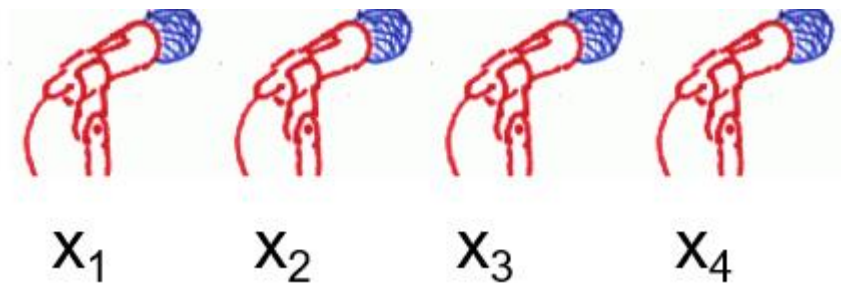
Humans can separate multiple signals **with only two ears/sensors**



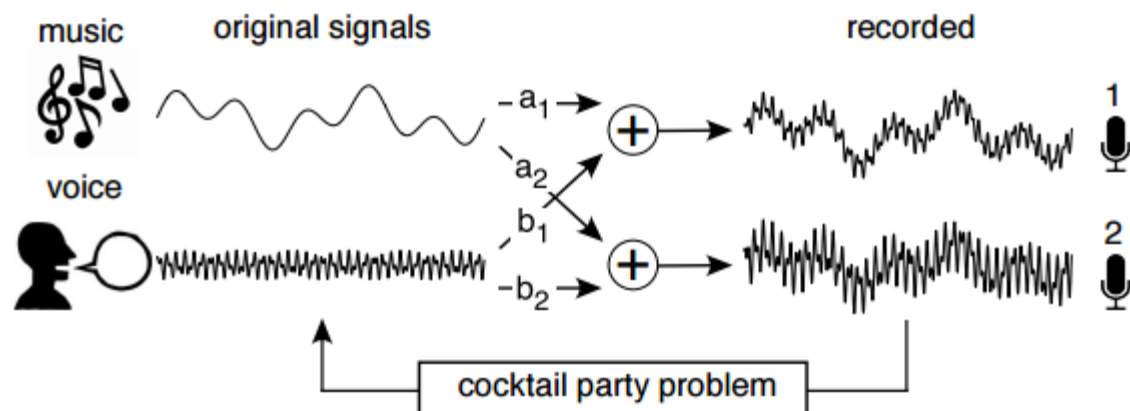
Computer can do the same but **with using many ears/sensors**

at least, as many ears/sensors as message signals

‘Human-computer’ has four ears:
a few microphones should be used



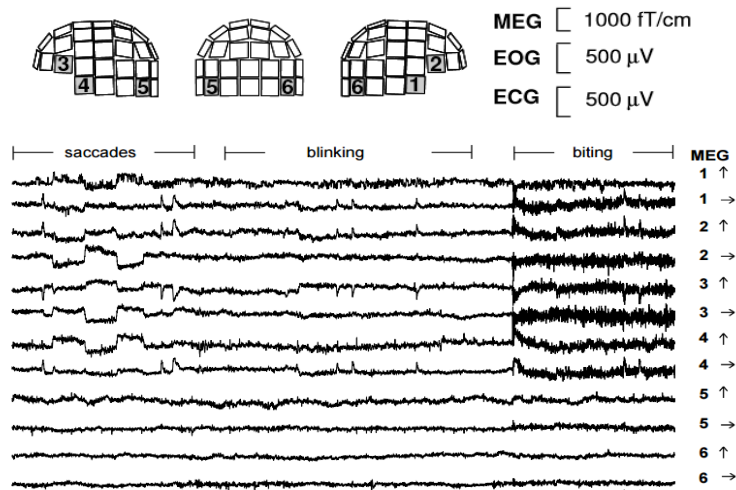
- **Data** (observed signals): mixed signals from **independent sources** obtained by **multiple sensors**
- **Analysis goal**: to find (to distinguish) the independent signals



The goal is to recover the original sources (i.e. music and voice) solely using the microphone recordings

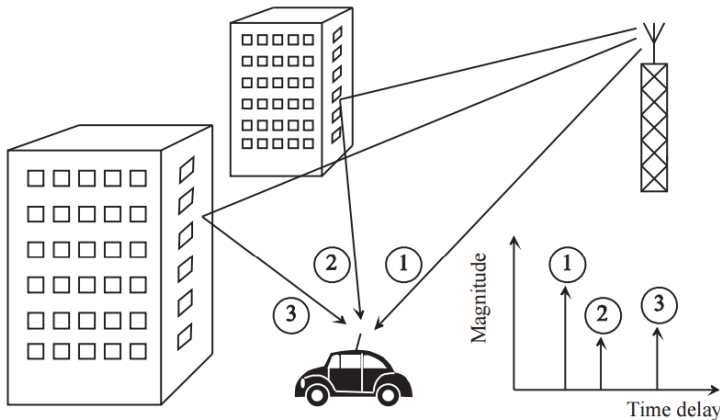


Removing blur from an image due to camera motion: identifying the original image and the motion path from a single blurry image



MagnetoEncephaloGrams (MEG) signals:

- from the frontal, temporal and occipital areas (multiple sensors)
- data contains several types of artifacts (signal sources) - ocular and muscle activity, the cardiac cycle, environmental magnetic disturbances, etc.



Telecommunications:

multipath propagation in urban environment

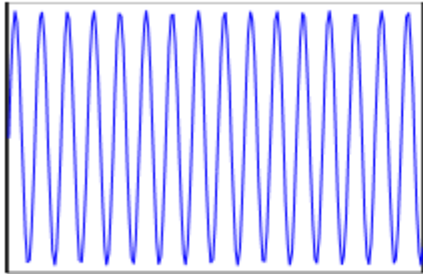
- Series of brain images with multiple voxels caused by different neuronal processes occurring within the brain
- Seismographic data from multiple seismometers

Formal task definition

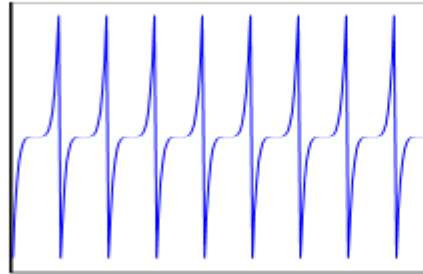
Hidden sources: q 'hidden' sources (latent variables, independent components):

s_j – signal from j^{th} source, $j = 1, 2, \dots, q$

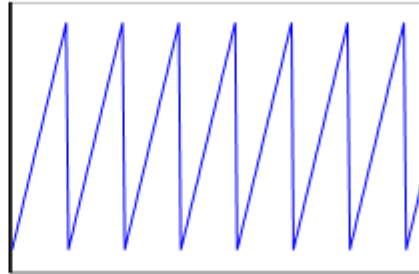
Example: $\{s_j = s_j(t)\}$ – time-dependent processes



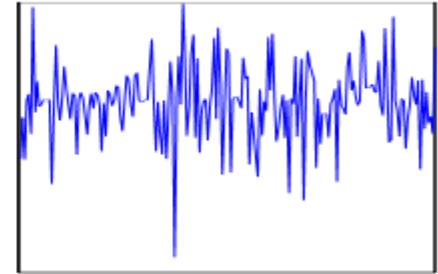
$$s_1 = s_1(t)$$



$$s_2 = s_2(t)$$

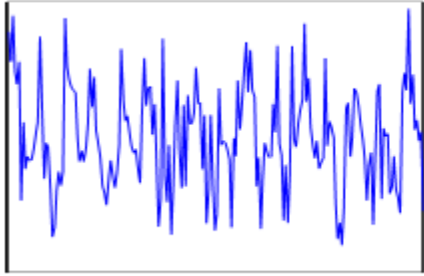


$$s_3 = s_3(t)$$

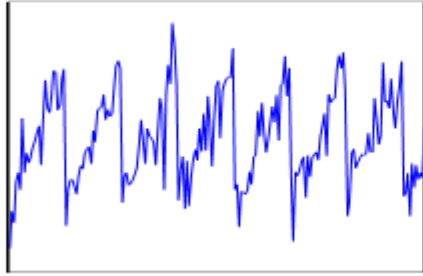


$$s_4 = s_4(t)$$

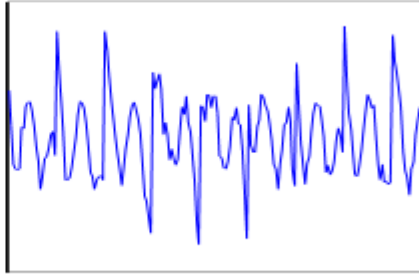
Data: m sensors: x_i – data (signal) from i^{th} sensor, $i = 1, 2, \dots, m$



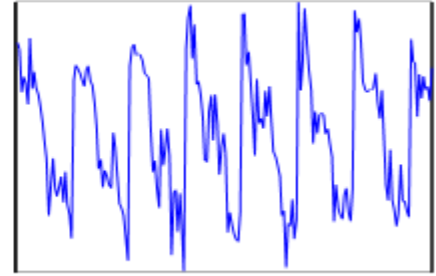
$$x_1 = x_1(t)$$



$$x_2 = x_2(t)$$



$$x_3 = x_3(t)$$



$$x_4 = x_4(t)$$

Data model - only **linear mixtures** of the source signals are observed:

$$X_i = \sum_{j=1}^q a_{ij} s_j, i = 1, 2, \dots, m$$

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

$$\mathbf{S} = \begin{pmatrix} s_1 \\ \dots \\ s_q \end{pmatrix} \in \mathbb{R}^q$$

- source signal vector

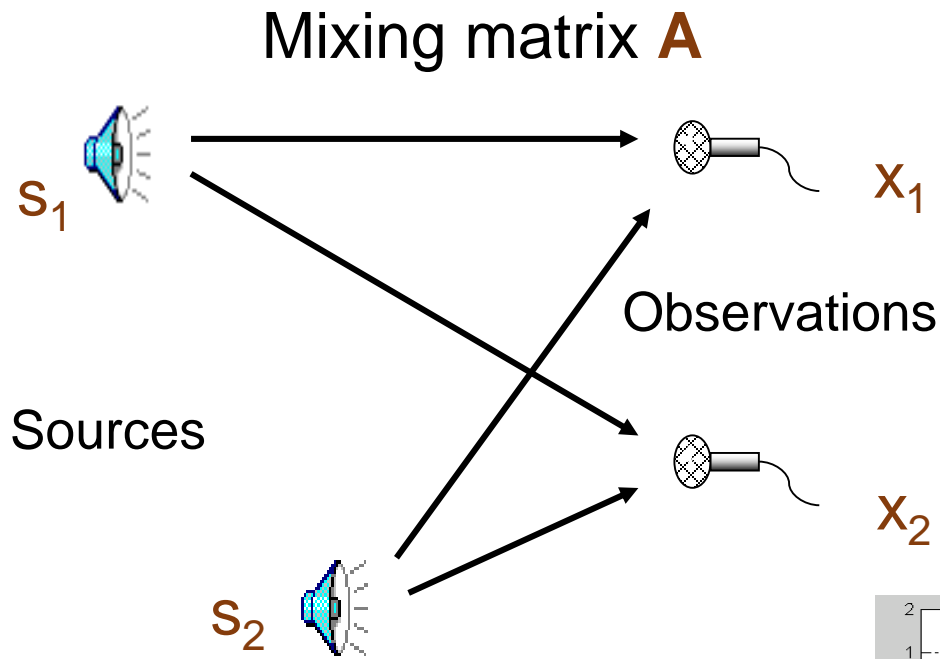
$$\mathbf{X} = \begin{pmatrix} x_1 \\ \dots \\ x_m \end{pmatrix} \in \mathbb{R}^m$$

- data (observed signal, multiplexed signal) vector

$$\mathbf{A} = \|a_{ij}\|$$

- **m×q** mixing matrix

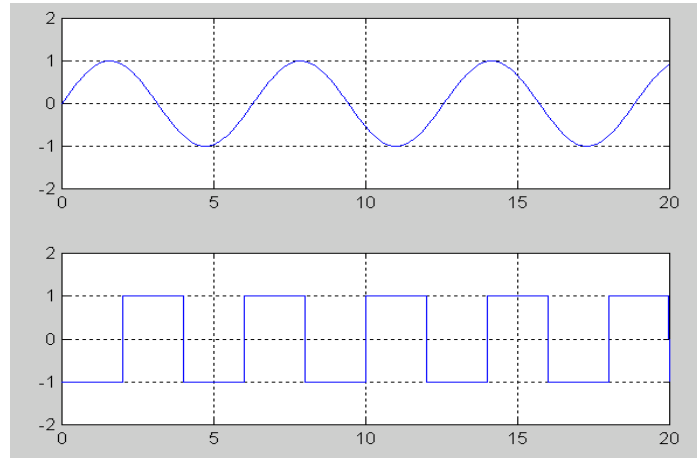
Cocktail party problem



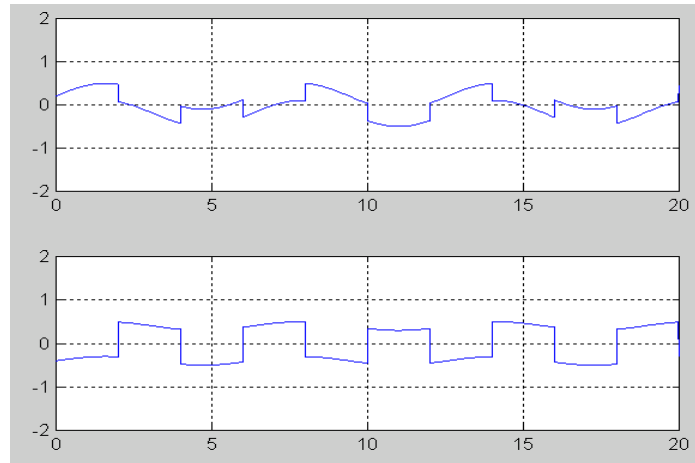
$$\mathbf{x} = \mathbf{A}\mathbf{s}$$

$$x_1(t) = a_{11}s_1 + a_{12}s_2$$

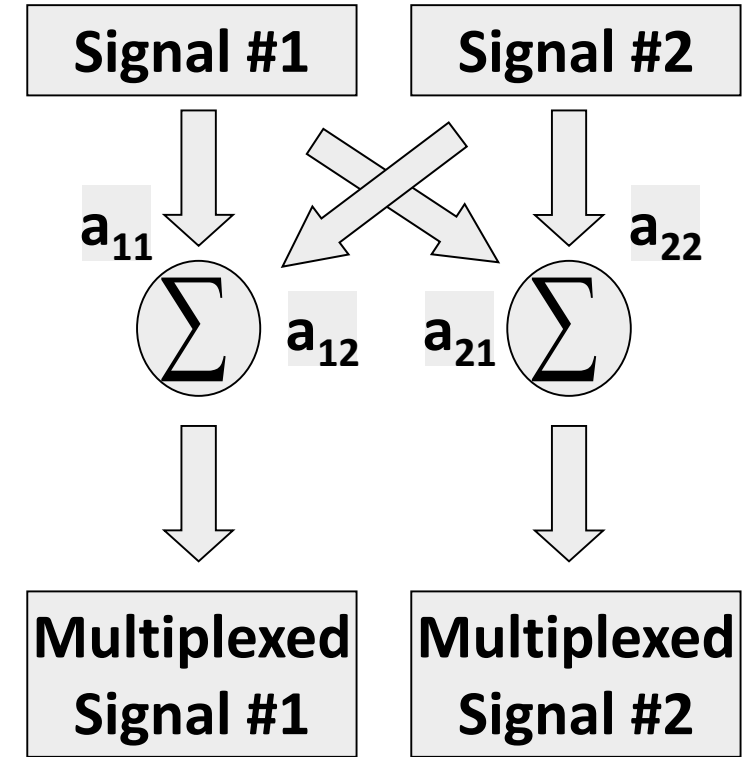
$$x_2(t) = a_{21}s_1 + a_{22}s_2$$



Two Independent sources



Mixture at two Mics



$\{a_{ij}\}$ depend on the distances
of the microphones from the speakers

ICA problem: to recover

- unknown **independent** 'source signals' $\{s_j\}$
- unknown mixing matrix A

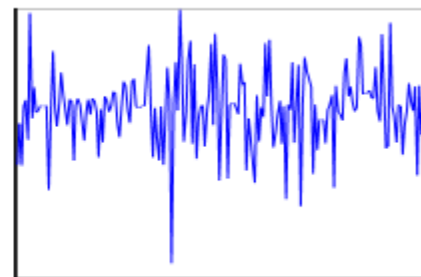
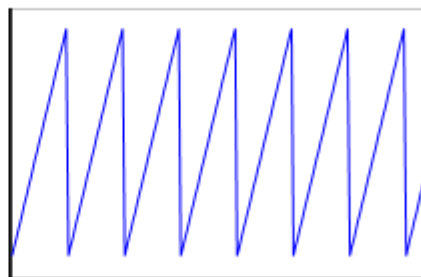
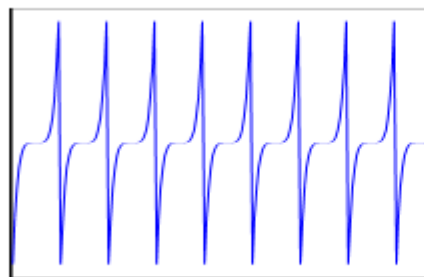
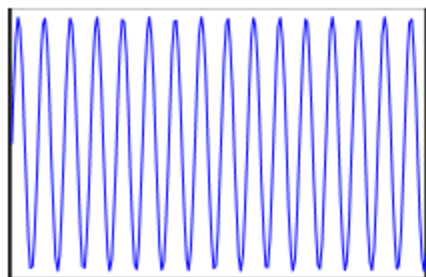
only from mixed data $\{x_i\}$ (unsupervised approach)

Desired **linear** solution – separated (estimated) signals: $\hat{S} = W \times X$

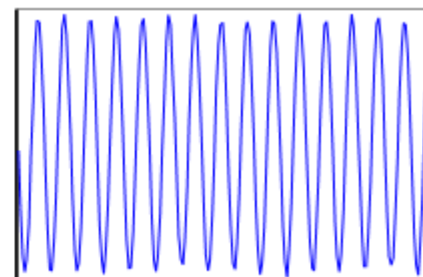
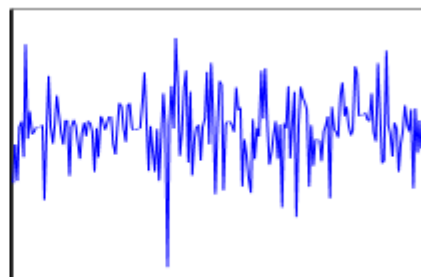
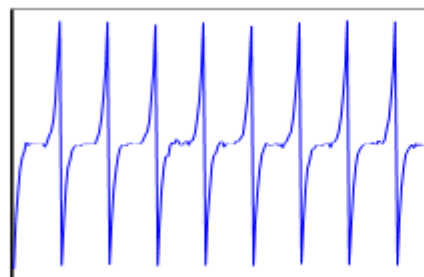
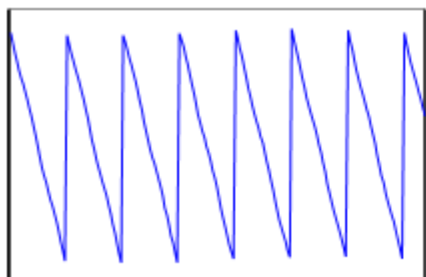
W - $q \times m$ **unmixing** matrix

ICA allows to recover source signals:

Source signals



ICA-recovered source signals



ICA has two related interpretations:

- **dimensional reduction:** if each source can be identified, a practitioner might choose to selectively delete or retain a single source (e.g. a person's voice, above)
- **filtering operation:** some aspect of the data is selectively removed or retained = is equivalent to projecting out some aspect (or dimension) of the data = a prescription for dimensional reduction

ICA was initially developed to deal with problems closely related to the **Cocktail party problem**

Now **ICA** has many other applications too:

- Blind source separation
- Image denoising
- Medical signal processing – fMRI, ECG, EEG
- Modelling of the hippocampus and visual cortex
- Feature extraction, face recognition
- Compression, redundancy reduction
- Watermarking
- Clustering
- Time series analysis
- Topic extraction
- Scientific Data Mining

Desired solution – separated (estimated) signals: $\hat{\mathbf{S}} = \mathbf{W} \times \mathbf{X}$

\mathbf{W} - $q \times m$ **unmixing** matrix

Mixing matrix \mathbf{A} is known: \mathbf{W} is solution of linear equation $\mathbf{X} = \mathbf{A} \times \mathbf{S}$:

$q = m$ and $q \times q$ mixing matrix \mathbf{A} is invertible (nonsingular: $\text{Rank}(\mathbf{A}) = q$): $\mathbf{W} = \mathbf{A}^{-1}$

$q < m$ and $\text{Rank}(\mathbf{A}) = q$:

$$\mathbf{W} = \mathbf{A}^+ \text{ - the solution}$$

$$\mathbf{A} = \mathbf{U} \times \mathbf{\Lambda} \times \mathbf{V}^T \text{ - SVD}$$

$$\mathbf{X} = (\mathbf{U} \times \mathbf{\Lambda} \times \mathbf{V}^T) \times \mathbf{S}$$

$$\mathbf{S} = (\mathbf{V} \times \mathbf{\Lambda}^{-1} \times \mathbf{U}^T) \times \mathbf{X}$$

$\mathbf{\Lambda}^{-1}$: non-zero diagonal elements are replaced by inverse

$$\mathbf{A}^+ = \mathbf{V} \times \mathbf{\Lambda}^{-1} \times \mathbf{U}^T$$

- (left) pseudoinverse Moore-Penrose matrix to $m \times q$ matrix \mathbf{A}

Mixing matrix **A** is unknown:

- what is possible to estimate in the ICA model
- when it is possible (in which cases)
- how to estimate if it is possible

What is possible to estimate: ICA limitations

1) ICA can't determine the order of the Independent components

Data model: $\mathbf{X} = \mathbf{A} \times \mathbf{S}$

\mathbf{B} - arbitrary $q \times q$ permutation matrix

\mathbf{B}^{-1} – inverse permutation matrix: $\mathbf{B}^{-1} \times \mathbf{B} = \mathbf{I}_q$

$q = 3$:

$$\mathbf{B} = \begin{pmatrix} 0 & 1 & 0 \\ 0 & 0 & 1 \\ 1 & 0 & 0 \end{pmatrix}$$

$$\mathbf{B}^{-1} = \begin{pmatrix} 0 & 0 & 1 \\ 1 & 0 & 0 \\ 0 & 1 & 0 \end{pmatrix}$$

$$\mathbf{X} = \mathbf{A} \times \mathbf{S} = \mathbf{A} \times (\underbrace{\mathbf{B}^{-1} \times \mathbf{B}}_{\mathbf{I}_q}) \times \mathbf{S} = (\mathbf{A} \times \mathbf{B}^{-1}) \times (\mathbf{B} \times \mathbf{S}) = \mathbf{A}^* \times (\mathbf{B} \times \mathbf{S})$$

\nearrow
 $= \mathbf{A} \times \mathbf{B}^{-1}$

$$\mathbf{B} \times \mathbf{S} = \begin{pmatrix} 0 & 1 & 0 \\ 0 & 0 & 1 \\ 1 & 0 & 0 \end{pmatrix} \times \begin{pmatrix} s_1 \\ s_2 \\ s_3 \end{pmatrix} = \begin{pmatrix} s_2 \\ s_3 \\ s_1 \end{pmatrix}$$

$$\mathbf{X} = \mathbf{A} \times \begin{pmatrix} s_1 \\ s_2 \\ s_3 \end{pmatrix} = \mathbf{A}^* \times \begin{pmatrix} s_2 \\ s_3 \\ s_1 \end{pmatrix}$$

So, we can call any Independent component as the first one

2) ICA can't determine the variances (energies) of the Independent components

Data model: $\mathbf{X} = \mathbf{A} \times \mathbf{S}$

\mathbf{B} - arbitrary $q \times q$ diagonal matrix

\mathbf{B}^{-1} – inverse diagonal matrix

$q = 3$:

$$\mathbf{B} = \begin{pmatrix} b_1 & 0 & 0 \\ 0 & b_2 & 0 \\ 0 & 0 & b_3 \end{pmatrix} \quad \mathbf{B}^{-1} = \begin{pmatrix} b_1^{-1} & 0 & 0 \\ 0 & b_2^{-1} & 0 \\ 0 & 0 & b_3^{-1} \end{pmatrix}$$

$$\mathbf{X} = \mathbf{A} \times \mathbf{S} = \mathbf{A} \times (\underbrace{\mathbf{B}^{-1} \times \mathbf{B}}_{\mathbf{I}_q}) \times \mathbf{S} = (\mathbf{A} \times \mathbf{B}^{-1}) \times (\mathbf{B} \times \mathbf{S}) = \mathbf{A}^* \times (\mathbf{B} \times \mathbf{S})$$

\uparrow
 $= \mathbf{A} \times \mathbf{B}^{-1}$

$$\mathbf{B} \times \mathbf{S} = \begin{pmatrix} b_1 & 0 & 0 \\ 0 & b_2 & 0 \\ 0 & 0 & b_3 \end{pmatrix} \times \begin{pmatrix} s_1 \\ s_2 \\ s_3 \end{pmatrix} = \begin{pmatrix} b_1 \times s_1 \\ b_2 \times s_2 \\ b_3 \times s_3 \end{pmatrix}$$
$$\mathbf{X} = \mathbf{A} \times \begin{pmatrix} s_1 \\ s_2 \\ s_3 \end{pmatrix} = \mathbf{A}^* \times \begin{pmatrix} b_1 \times s_1 \\ b_2 \times s_2 \\ b_3 \times s_3 \end{pmatrix}$$

So, we can find Independent component up to scaling

The source signals can principally be only recovered up to permutation and scaling

Mixing matrix **A** is unknown:

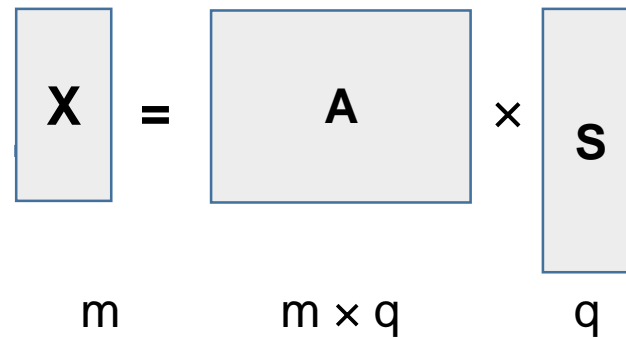
- what is possible to estimate in the ICA model
- **when it is possible (in which cases)**
- how to estimate if it is possible

When mixing matrix and source signals can be identified (estimated):

Requirement 1: the number q of separated (source) signals cannot be larger than the number m of observed variables (inputs)

Understandable!

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



Under $m < q$, even with known matrix \mathbf{A} , we can not find the vector \mathbf{S} from the vector \mathbf{X} **by the only way**

Further, for simplicity, we will assume sometimes that:

- $m = q$
- $q \times q$ mixing matrix \mathbf{A} is invertible

Requirement 2: source signals $\{s_i\}$ are **mutually statistically independent**

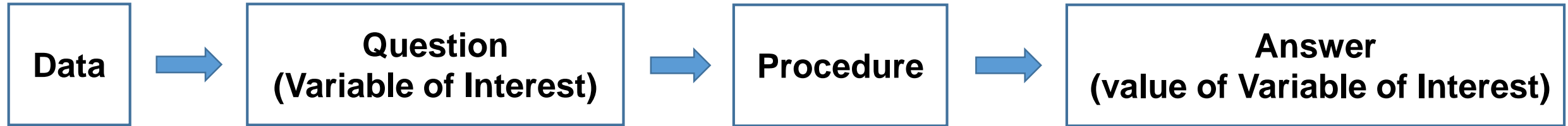
Without any requirements - trivial solution: $\hat{\mathbf{S}} = \mathbf{X}$, \mathbf{A} and \mathbf{W} - unit $q \times q$ matrices

Independent sources:

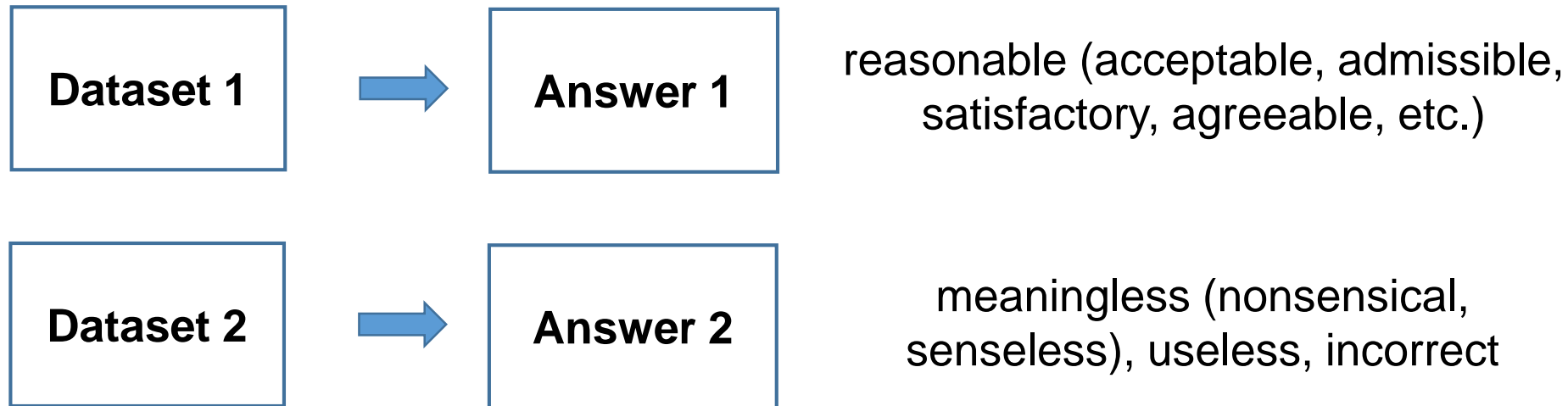
- a rather natural ‘default’ assumption when we do not want to postulate any specific dependencies between the components
- a **physical interpretation** of independence is also sometimes possible: if the components are created by physically separate and non-interacting entities

Strict mathematical definition requires the concept of a Mathematical data generative model

Data analysis process



Applying the same question and the same procedure to different datasets:



**We need some assumptions about data (data model) to predict a quality of answer:
which quality we can expect if the data satisfy the assumptions**

We need some assumptions about data (data model) to predict a quality of answer

Inverse task: given data model and selected quality criteria, to construct 'the best procedure'

Physical data model: follows from 'physical aspects' which describes data generation process and determines data properties

Math data model: which can generate (simulate) data using certain math procedures (**generative model**) similar to the data obtained from real sources

Generative model:

- **Geometrical model:** Data space - data support - 'a place where data lie'
- **Extracting model:** how data are extracted (are selected) from the Data space

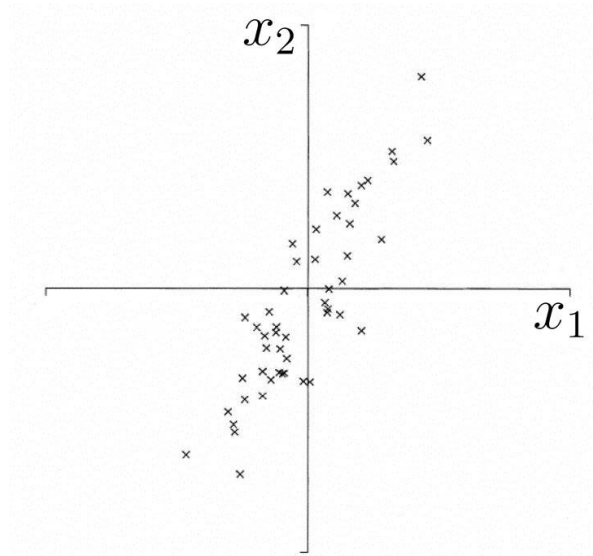
Probabilistic extracting model: data (training data and new, out-of-sample, data) are selected **randomly** from Data space independently of each other according to the certain **unknown probability distribution**

Probability distribution: determines statistical data properties

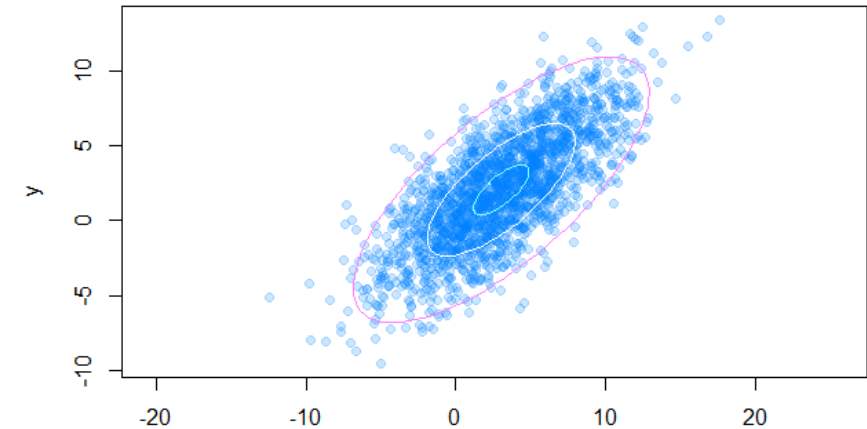
Statistical approach: given training dataset,

- to select appropriate **Probability distribution** which 'well describes' the training data
- selected Probability distribution is considered as **Statistical data model** and used for constructing 'best' data analysis procedures

Example: training data looks like a single cloud



A sample from two-dimensional Gaussian (normal) distribution looks similar



Thus, we can assume that we deal with random two-dimensional Gaussian (normal) random variables

Independent Component Analysis uses **probabilistic data model** (assuming data are random) and, thus, **probabilistic language** for strict formulating the requirements to data and describing the principles underlying the algorithms **is used**

Source signals $\{s_i\}$ are **mutually statistically independent** means that their joint **density function**

$$p(\mathbf{s}) = p(s_1, s_2, \dots, s_q)$$

is factorizable:

$$p(s_1, s_2, \dots, s_q) = p(s_1) \times p(s_2) \times \dots \times p(s_q)$$

Independence implies: conditional distribution of any source signal (say, s_1) under given all remaining signals $\{s_2, s_3, \dots, s_q\}$ doesn't depend on them:

$$p(s_1 | s_2, \dots, s_q) = p(s_1)$$

Thus, ICA searches for the rotation W of the observed data such that estimated source signal $\{\hat{s}_j\}$ are mutually statistically independent too: $p(\hat{s}_1, \hat{s}_2, \dots, \hat{s}_q) = p(\hat{s}_1) \times p(\hat{s}_2) \times \dots \times p(\hat{s}_q)$

This implies, for example, that $p(\hat{s}_1 | \hat{s}_2, \dots, \hat{s}_q) = p(\hat{s}_1)$

Probability theory: independency of random variables implies noncorrelatedness: $\text{cov}(s_i, s_j) = 0$

Thus, the estimated source signal $\{\hat{s}_j\}$ should be at least uncorrelated

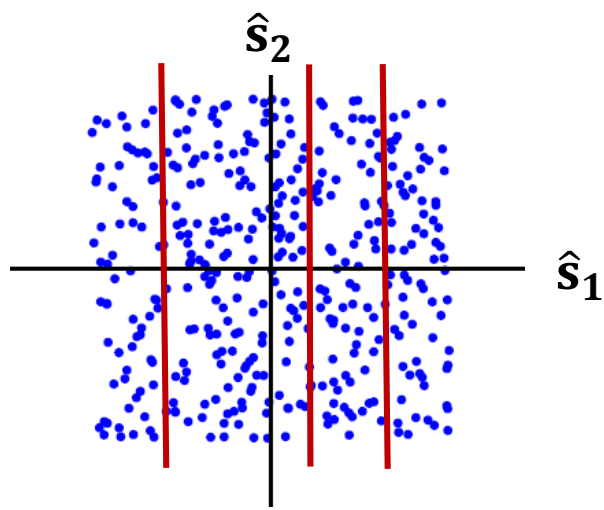
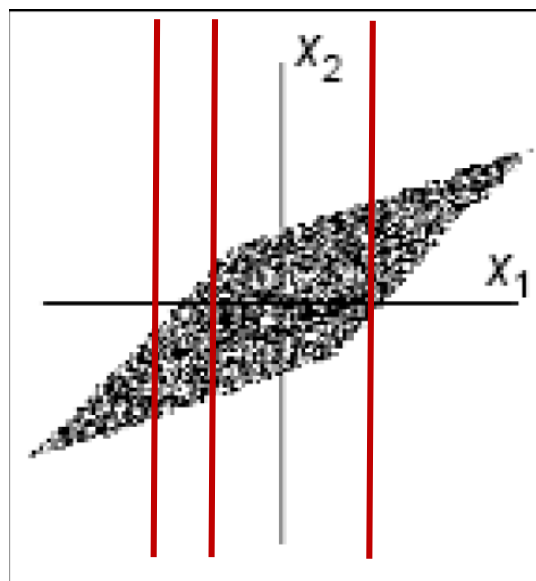
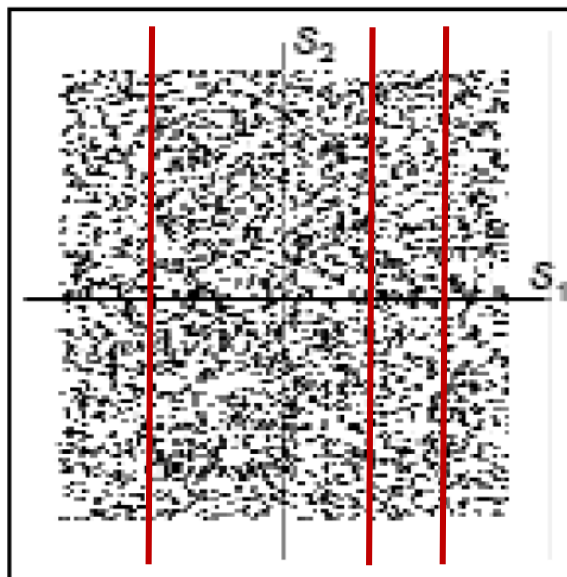
PCA transforms correlated multiplexed signals $\{x_i\}$ to uncorrelated quantities.

Can PCA give the solution to the ICA task?

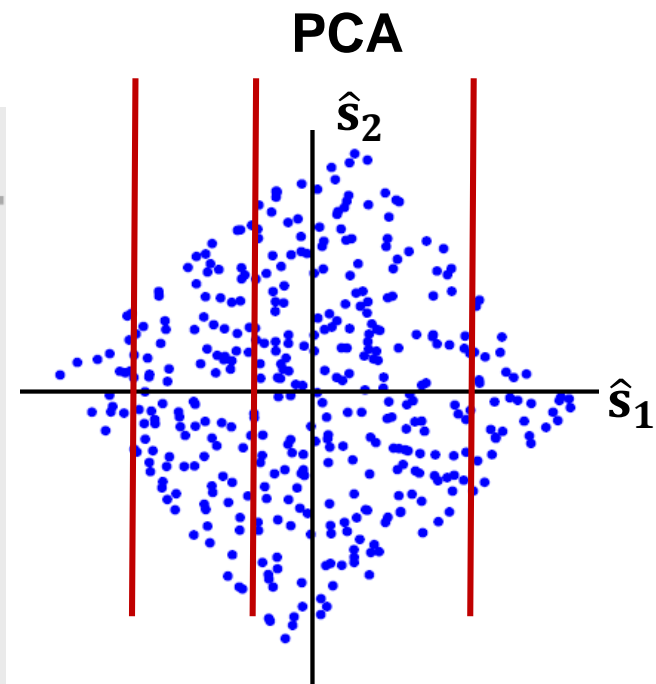
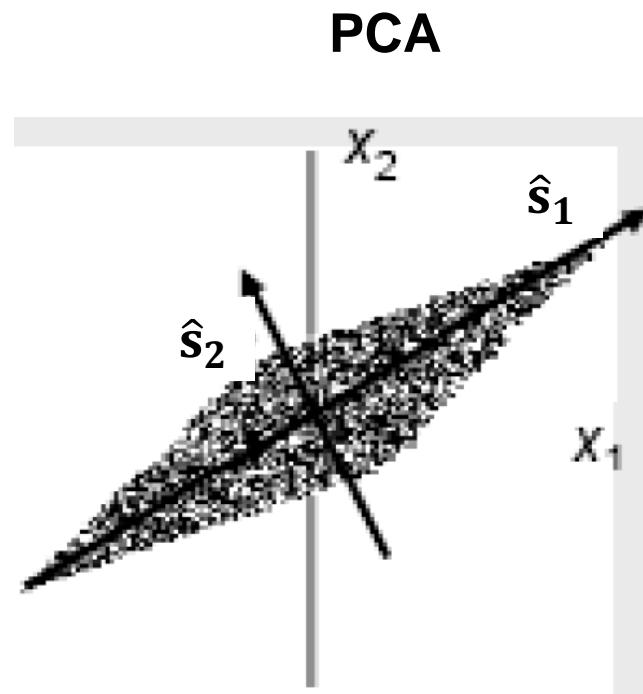
Answer - NO

Illustrative example: two components s_1 and s_2 with uniform distributions

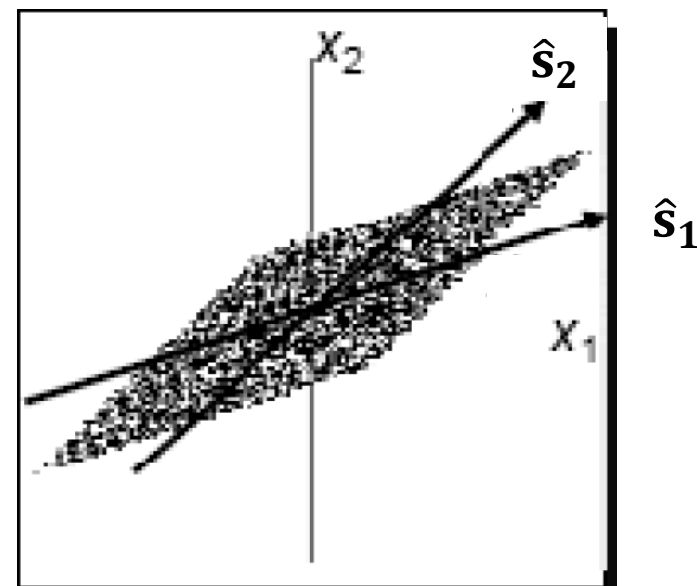
Mixing matrix
 $\mathbf{A} = \begin{pmatrix} 1 & 2 \\ -1 & 2 \end{pmatrix}$



ICA



ICA



Probability theory: independency of random variables implies their noncorrelatedness

But the converse is not true in general

Example: $\xi \sim N(0, 1)$

$$\eta_1 = \xi$$

$$\mathbf{M}\eta_1 = 0$$

η_1 and η_2 are uncorrelated:

$$\eta_2 = \xi^2 - 1$$

$$\mathbf{M}\eta_2 = 0$$

$$\text{cov}(\eta_1, \eta_2) = \mathbf{M}(\xi \times (\xi^2 - 1)) = \mathbf{M}(\xi^3 - \xi) = 0$$

These random variables are connected by non-random functional dependency

Probability theory: in Gaussian case, a noncorrelatedness implies an independency

Maybe, the PCA will solve **the ICA problem for Gaussian signals?**

Answer - NO

A surprising and unexpected fact: **no method** can solve **the ICA problem for Gaussian signals**

The **fundamental restriction** in ICA is that the independent components **must be nongaussian** (perhaps, **with the exception of only one** component) for **ICA to be possible**

Why Gaussian variables **are forbidden** (at most one of the sources can be Gaussian)?

Intuitively:

- **a correlatedness** is determined by the **second moments only**
- **an independency is more stronger property** than correlatedness which **is essential for** estimation of **the ICA model** and determined by ‘higher-order information’ about random variables (their higher-order moments)
- the **Gaussian distributions** are “**too simple**”: their higher order moments are completely determined by first two moments

A more strict explanation

two Gaussian source signals \mathbf{s}_1 and \mathbf{s}_2 \rightarrow observed multiplexed signals \mathbf{x}_1 and \mathbf{x}_2 are Gaussian too

Centering

$$\mathbf{a} = \mathbf{M}\mathbf{X} \text{ and } \mathbf{X}' = \mathbf{X} - \mathbf{a} \rightarrow \mathbf{M}\mathbf{X}' = 0$$

If $\hat{\mathbf{S}}' = \mathbf{W} \times \mathbf{X}'$ - the solution to the centered ICA problem \rightarrow

$\hat{\mathbf{S}} = \hat{\mathbf{S}}' + \mathbf{W} \times \mathbf{a}$ - the solution to the original ICA problem

Whitening

$\mathbf{cov}(\mathbf{X}) = \mathbf{E} \times \mathbf{D} \times \mathbf{E}^T$ – eigenvalue decomposition of positive semidefinite covariance matrix:

\mathbf{E} - orthogonal, \mathbf{D} - diagonal

$\mathbf{V} = \mathbf{E} \times \mathbf{D}^{-1/2} \times \mathbf{E}^T$ - whitening matrix

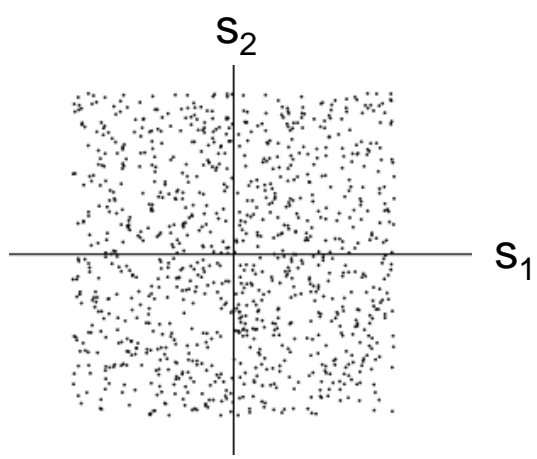
$\mathbf{X}_w = \mathbf{V} \times \mathbf{X} \rightarrow$

$\mathbf{cov}(\mathbf{X}_w) = \mathbf{V} \times \mathbf{cov}(\mathbf{X}) \times \mathbf{V}^T = (\mathbf{E} \times \mathbf{D}^{-1/2} \times \mathbf{E}^T) \times (\mathbf{E} \times \mathbf{D} \times \mathbf{E}^T) \times (\mathbf{E} \times \mathbf{D}^{-1/2} \times \mathbf{E}^T)^T = \dots = \mathbf{I}_2$ - unit matrix

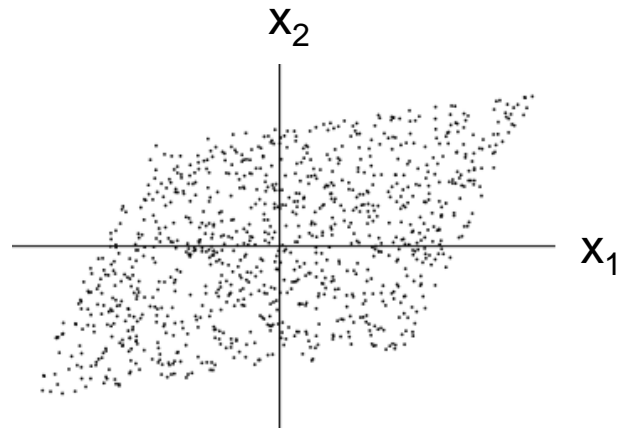
If $\hat{\mathbf{S}}' = \mathbf{W}' \times \mathbf{X}_w$ - the solution to the whitened ICA problem \rightarrow

$\hat{\mathbf{S}} = \mathbf{W} \times \mathbf{X}$, where $\mathbf{W} = \mathbf{W}' \times \mathbf{V}$, - the solution to the original ICA problem

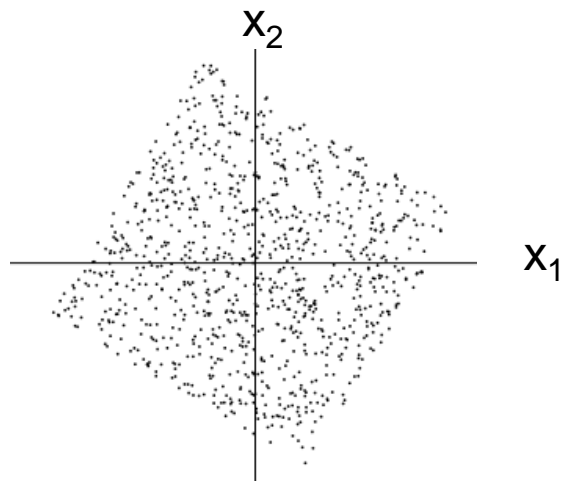
Without loss of generality, we can assume that $\mathbf{X} \sim \mathcal{N}(0, \mathbf{I}_2)$



Dataset consisting of the independent uniformly distributed signals s_1 and s_2



Dataset consisting of the observed mixtures x_1 and x_2



Dataset consisting of the **whitened** observed mixtures x_1 and x_2

Source signals $\hat{\mathbf{S}}$ can be found up to scaling

→ we can find Gaussian vector \mathbf{S} with independent components as vector with unit covariance matrix \mathbf{I}_2

Model $\mathbf{X} = \mathbf{A} \times \mathbf{S}$, $\mathbf{X} \sim \mathcal{N}(0, \mathbf{I}_2)$ and $\mathbf{S} \sim \mathcal{N}(0, \mathbf{I}_2)$, implies that mixing matrix \mathbf{A} is orthogonal

In Gaussian model, mixing matrix \mathbf{A} and, thus, source vector \mathbf{S} , **cannot be identified** (estimated) – any orthogonal matrix \mathbf{A} transform Gaussian source vector $\mathbf{S} \sim \mathcal{N}(0, \mathbf{I}_2)$ with independent components to the random vector \mathbf{X} **with the same distribution** $\mathcal{N}(0, \mathbf{I}_2)$

ICA is essentially impossible for the Gaussian sources

If **just one** of the independent components is Gaussian, the **ICA model can still be identified**

Fortunately, signals measured by physical sensors are usually quite nongaussian

- R1: number of source signals is equal to number of observed variables (at least not less)
- R2: source signals $\{s_i\}$ are **mutually statistically independent**
- R3: all source signals $\{s_i\}$, with the exception of at most one, must be non-Gaussian

Proven result: under requirements R1 – R3, mixing matrix and components can be identified

ICA model in which some of the components are Gaussian, some non-Gaussian:

- we can estimate all the non-Gaussian components
- the Gaussian components cannot be separated from each other: these components will be arbitrary linear combinations of some constructed Gaussian components
- in the case of just one Gaussian component, we can estimate the model (the single Gaussian component does not have any other Gaussian components that it could be mixed with)

Mixing matrix **A** is unknown:

- what is possible to estimate in the ICA model
- when it is possible (in which cases)
- **how to estimate if it is possible**

ICA preprocessing steps: centering and whitening

Centering

$$\mathbf{a} = \mathbf{M}\mathbf{X} \text{ and } \mathbf{X}' = \mathbf{X} - \mathbf{a} \rightarrow \mathbf{M}\mathbf{X}' = 0$$

$$\hat{\mathbf{S}}' = \mathbf{W} \times \mathbf{X}' \text{ - solution to the centered ICA } \rightarrow \hat{\mathbf{S}} = \hat{\mathbf{S}}' + \mathbf{W} \times \mathbf{a} \text{ - solution to the original ICA}$$

Without loss of generality, we can assume that $\mathbf{M}\mathbf{X} = 0$

When centered observed variables are used, $\mathbf{M}\mathbf{S} = \mathbf{M}(\mathbf{W} \times \mathbf{X}) = 0$

ICA model: $\mathbf{X} = \mathbf{A} \times \mathbf{S}$

$\mathbf{S} \in \mathbb{R}^q$ - source signal, $\mathbf{X} \in \mathbb{R}^m$ - data, $\mathbf{A} = \parallel a_{ij} \parallel$ - $m \times q$ mixing matrix ($m \geq q$, $\text{Rank}(\mathbf{A}) = q$)

$\text{cov}(\mathbf{X}) = \mathbf{E} \times \mathbf{D} \times \mathbf{E}^T$ – Eigenvalue decomposition of positive semidefinite covariance $m \times m$ matrix:

\mathbf{E} – orthogonal $m \times m$ matrix, \mathbf{D} – diagonal $m \times m$ matrix at most q nonzero diagonal elements

Whitening

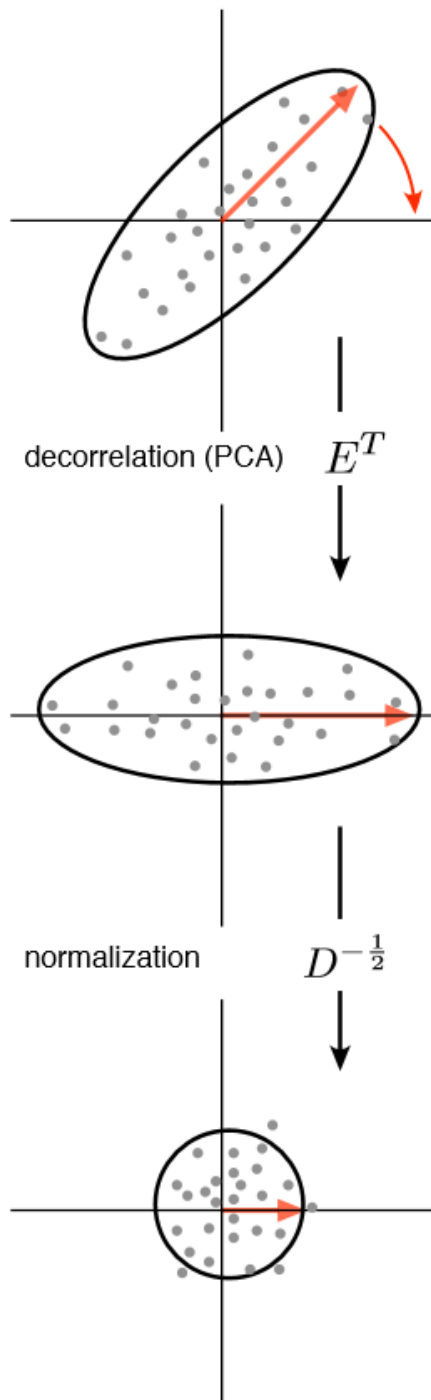
$\mathbf{V} = \mathbf{E}_q \times \mathbf{D}_q^{-1/2} \times \mathbf{E}_q^T$ - whitening $m \times m$ matrix

\mathbf{E}_q - orthogonal $m \times q$ matrix consisting of first q columns of $m \times m$ matrix \mathbf{E}

\mathbf{D}_q - diagonal $q \times q$ matrix with positive eigenvalues of matrix $\text{cov}(\mathbf{X})$: $\text{Rank}(\text{cov}(\mathbf{X})) = q$

In fact, ‘PCA technique’ has been applied to the observed data!

$\mathbf{X}_w = \mathbf{V} \times \mathbf{X}$ - whitened data vector $\rightarrow \text{cov}(\mathbf{X}_w) = \mathbf{I}_q$ - unit matrix



If $\hat{\mathbf{S}}' = \mathbf{W}' \times \mathbf{X}_w$ - the solution to the whitened ICA problem \rightarrow
 $\hat{\mathbf{S}} = \mathbf{W} \times \mathbf{X}$, $\mathbf{W} = \mathbf{W}' \times \mathbf{V}$, - the solution to the original ICA problem

Without loss of generality, we can assume that $\text{cov}(\mathbf{X}_w) = \mathbf{I}_q$

Source signals $\hat{\mathbf{S}}$ can be found up to scaling

$$\text{Cov}(\mathbf{S}) = \mathbf{M}(\mathbf{S} \times \mathbf{S}^T) = \mathbf{I}_q$$

→ we can find \mathbf{S} with independent components and **unit variances**

$$\text{Model } \mathbf{X} = \mathbf{A} \times \mathbf{S} \quad \rightarrow \quad \mathbf{X}_w = \mathbf{V} \times \mathbf{A} \times \mathbf{S} \quad \rightarrow \quad \mathbf{X}_w = \mathbf{A}_w \times \mathbf{S}, \text{ where } \mathbf{A}_w = \mathbf{V} \times \mathbf{A}$$

$$\text{cov}(\mathbf{X}_w) = \text{cov}(\mathbf{A}_w \times \mathbf{S}) = \mathbf{A}_w \times \text{cov}(\mathbf{S}) \times (\mathbf{A}_w)^T = \mathbf{A}_w \times (\mathbf{A}_w)^T$$

$$\text{cov}(\mathbf{X}_w) = \mathbf{I}_q \quad \rightarrow \quad \mathbf{A}_w \times (\mathbf{A}_w)^T = \mathbf{I}_q \quad \rightarrow \quad \mathbf{A}_w - \text{orthogonal } m \times q \text{ matrix}$$

We can restrict our search for the mixing matrix to the space of **orthogonal** $m \times q$ matrices:

Case $m = q$:

- original mixing matrix \mathbf{A} depends on q^2 parameters which must be estimated
- orthogonal mixing matrix \mathbf{A}_w depends on $q(q - 1)/2$ parameters only

whitening solves half of the problem of ICA by reducing the complexity of the problem

Fundamental approaches to the solution of the ICA problem:

- maximization the nongaussianity of the marginals
- maximum likelihood approach
- minimization of mutual information between components

$$\hat{\mathbf{S}} = \mathbf{W} \times \mathbf{X}$$

$$\begin{array}{c} \boxed{\hat{\mathbf{S}}} \\ q \end{array} = \begin{array}{c} \boxed{\mathbf{W}} \\ q \times q \end{array} \times \begin{array}{c} \boxed{\mathbf{X}} \\ q \end{array}$$

$$\begin{array}{c} \boxed{\mathbf{W}} \\ q \times q \end{array} = \begin{array}{c} \boxed{W_1} \\ \boxed{W_2} \\ \boxed{\dots} \\ \boxed{W_q} \\ q \end{array}$$

W_1, W_2, \dots, W_q - rows of \mathbf{W}

Recovered signals: $\hat{s}_i = W_i \times \mathbf{X}$

Assumption: source signals $\{s_i\}$ are **random, mutually statistically independent**, and all signals, with the exception of at most one, are **non-Gaussian**

Data vector \mathbf{X} is random \rightarrow we assume (after centering and whitening) that

$$\mathbf{MX} = \mathbf{0} \quad \rightarrow \quad \mathbf{cov}(\mathbf{X}) = \mathbf{I}_q \text{ - unit matrix}$$

$$\mathbf{MX} = \mathbf{0} \rightarrow \mathbf{MS} = \mathbf{0}$$

Source signals can be found up to scaling \rightarrow we assume that $\mathbf{cov}(\mathbf{S}) = \mathbf{I}_q$

$$\mathbf{cov}(\mathbf{X}) = \mathbf{I}_q, \mathbf{cov}(\mathbf{S}) = \mathbf{I}_q \quad \rightarrow \quad \text{mixing matrix } \mathbf{A} \text{ is orthogonal}$$

Source signals can be found up to permutation and scaling \rightarrow we will search:

- independent recovered signals $\hat{\mathbf{S}} = \mathbf{W} \times \mathbf{X}$ satisfying relations $\mathbf{M}\hat{\mathbf{S}} = \mathbf{0}$, $\mathbf{cov}(\hat{\mathbf{S}}) = \mathbf{I}_q$
- orthogonal unmixing matrix $\mathbf{W} \rightarrow$ the rows $\mathbf{W}_1, \mathbf{W}_2, \dots, \mathbf{W}_q$ of \mathbf{W} - **orthonormal** vectors

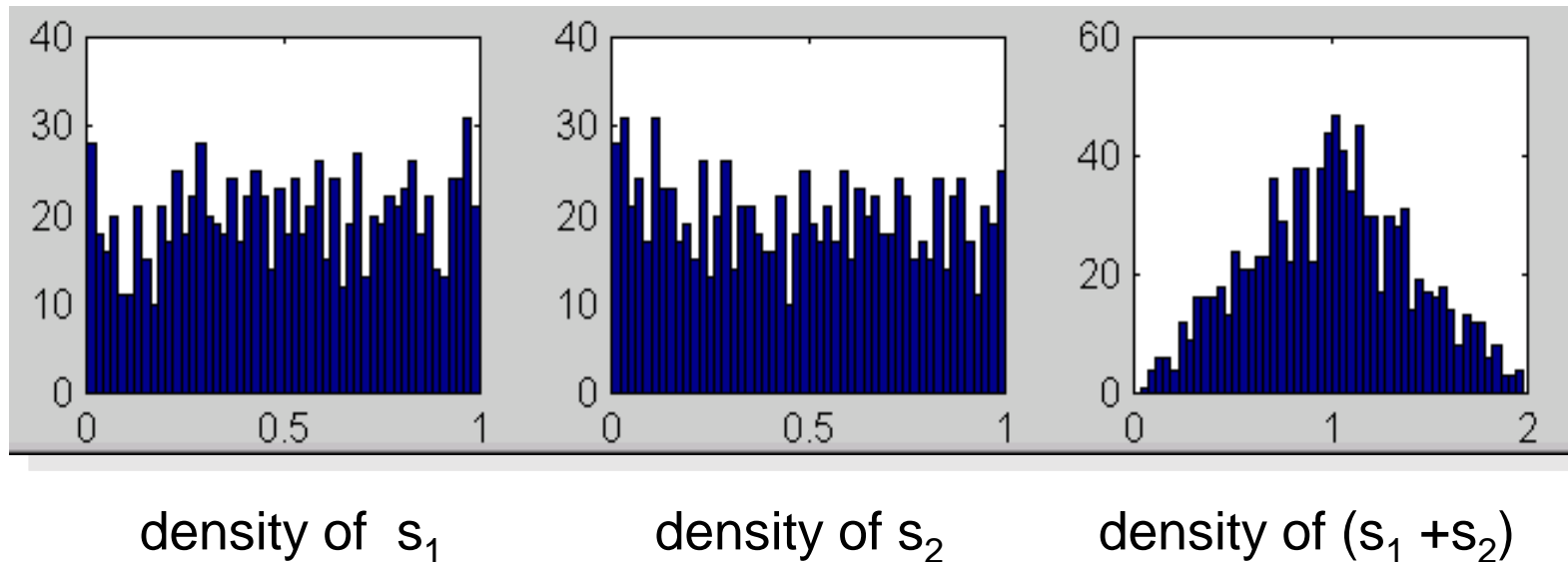
Fundamental approaches to the solution of the ICA problem:

1) Maximum nongaussianity principle: find maximally 'nongaussian' (with respect to **chosen measure of nongaussianity**) **independent** signals $\{\hat{s}_i = W_i \times \mathbf{X}\}$

Independency as maximum nongaussianity: motivation

Probability theory (Central Limit Theorem) - *informally*: The distribution of a sum of independent random variables tends toward a Gaussian distribution

A sum of a few independent random variables **is more gaussian** than individual random variables



$$\hat{\mathbf{S}} = \mathbf{W} \times \mathbf{X}$$

\mathbf{W}_k - k -th row of desired unmixing matrix \mathbf{W} \rightarrow

$\hat{\mathbf{S}}_k = \mathbf{W}_k \times \mathbf{X}$ - k -th component of the estimated signal $\hat{\mathbf{S}}$

$\hat{\mathbf{S}}_k = \mathbf{W}_k \times \mathbf{X} = \mathbf{W}_k \times \mathbf{A} \times \mathbf{S}$ - a weighted linear combination of independent signals $\{\mathbf{S}_i\}$

$\hat{\mathbf{S}}_k$ is 'more Gaussian' than source signals $\{\mathbf{S}_i\}$ **AND becomes least gaussian when** equals to one of $\{\mathbf{S}_i\}$

Starting from the first row, we could take \mathbf{W}_1 as a vector which **maximizes 'non-gaussianity'** of $\mathbf{W}_1 \times \mathbf{X}_w$

After choosing $\hat{\mathbf{S}}_1 = \mathbf{W}_1 \times \mathbf{X}$, we will search for next estimated signal $\hat{\mathbf{S}}_2 = \mathbf{W}_2 \times \mathbf{X}_w$ (\mathbf{W}_2 - 2nd row of desired matrix \mathbf{W}) by similar way taking into account that $\hat{\mathbf{S}}_2$ and $\hat{\mathbf{S}}_1$ are independent variables

\rightarrow they must be uncorrelated

and so on

$f(\xi)$ - a chosen measure of non-Gaussianity of random variable ξ , $\mathbf{M}\xi = 0$ and $\mathbf{M}\xi^2 = 1$

$\mathbf{b} \in \mathbb{R}^q$ - arbitrary row of desired unmixing matrix \mathbf{W} which determines estimated signal $\hat{s} = \mathbf{b} \times \mathbf{X}$

$F(\mathbf{b}) = \mathbf{M}f(\mathbf{b} \times \mathbf{X})$ - objective function defined by chosen measure of non-Gaussianity which should be maximized under constraint $\mathbf{b}^T \times \mathbf{b} = 1$ with taking into account that $\mathbf{cov}(\mathbf{X}_w) = \mathbf{I}_q$

X_1, X_2, \dots, X_T - training dataset consisting of centering and whitened observed signals

$F_n(\mathbf{b})$ - some estimator of the objective function $F(\mathbf{b})$ based on training dataset

For example: $F_n(\mathbf{b}) = \frac{1}{n} \sum_{t=1}^T F(\mathbf{b}^T \times X_t)$

ICA: $F_n(\mathbf{b}) \rightarrow \max$ under constraint $\mathbf{b}^T \times \mathbf{b} = 1$ with taking into account that $\frac{1}{n} \sum_{t=1}^T (X_t \times X_t^T) = \mathbf{I}_q$

We look for the rows of matrix \mathbf{W} one after another, starting from the first row \mathbf{W}_1 by maximizing an objective function $F(\mathbf{W}_1)$ under constraint $\mathbf{W}_1^T \times \mathbf{W}_1 = 1$

After choosing $\hat{\mathbf{s}}_1 = \mathbf{W}_1 \times \mathbf{X}$, we will search for next estimated signal $\hat{\mathbf{s}}_2 = \mathbf{W}_2 \times \mathbf{X}$ by maximizing $F(\mathbf{W}_2)$ under constraints $\mathbf{W}_2^T \times \mathbf{W}_2 = 1$

and under additional constraint $\mathbf{W}_2 \times \mathbf{W}_1^T = 0$ which follows from independence of $\hat{\mathbf{s}}_2$ and $\hat{\mathbf{s}}_1$:

$$\mathbf{Cov}(\hat{\mathbf{s}}_2, \hat{\mathbf{s}}_1) = \mathbf{M}(\hat{\mathbf{s}}_2 \times \hat{\mathbf{s}}_1) = \mathbf{M}\{(\mathbf{W}_2 \times \mathbf{X}) \times (\mathbf{W}_1 \times \mathbf{X})\} = \mathbf{W}_2 \times \mathbf{M}(\mathbf{X} \times \mathbf{X}^T) \times \mathbf{W}_1^T = \mathbf{W}_2 \times \mathbf{W}_1^T = 0$$

and so on

ICA: used quantitative measures of non-Gaussianity

ξ - random variable, $\mathbf{M}\xi = 0$, $\mathbf{M}\xi^2 = 1$

1) **Kurtosis:** $\mathbf{Kurt}(\xi) = \mathbf{M}\xi^4 - 3(\mathbf{M}\xi^2)^2$, $\mathbf{Kurt}(\xi_{\text{Gauss}}) = 0$ for Gauss random variable ξ_{Gauss}

$f(\xi) = |\mathbf{Kurt}(\xi)|$ - a measure of non-Gaussianity

$$F(\mathbf{b}) = f(\mathbf{b} \times \mathbf{X}) = \mathbf{M}(\mathbf{b} \times \mathbf{X})^4 - 3(\mathbf{M}(\mathbf{b} \times \mathbf{X})^2)^2 = \frac{1}{T} \sum_{t=1}^T (\mathbf{b} \times \mathbf{X}_t)^4 - 3 \times \left(\frac{1}{T} \sum_{t=1}^T (\mathbf{b} \times \mathbf{X}_t)^2 \right)^2$$

$\mathbf{X}_1, \mathbf{X}_2, \dots, \mathbf{X}_T$ - training dataset consisting of centering and whitened observed signals

$$\frac{1}{T} \sum_{t=1}^T (\mathbf{b} \times \mathbf{X}_t)^2 = \frac{1}{T} \sum_{t=1}^T \{ (\mathbf{b} \times \mathbf{X}_t) \times (\mathbf{X}_t^T \times \mathbf{b}^T) \} = \mathbf{b} \times \left(\frac{1}{T} \sum_{t=1}^T \{ \mathbf{X}_t \times \mathbf{X}_t^T \} \right) \times \mathbf{b}^T = \mathbf{b} \times \mathbf{b}^T = 1$$

$$F(\mathbf{b}) = \frac{1}{T} \sum_{t=1}^T (\mathbf{b} \times \mathbf{X}_t)^4 - 3 \rightarrow \max \quad \text{under constraint } \mathbf{b} \times \mathbf{b}^T = 1$$

$\mathbf{cov}(\mathbf{X}) = \mathbf{I}_q$



$$F(\mathbf{b}) = |\text{Kurt}(\mathbf{b} \times \mathbf{X})| = \frac{1}{T} \sum_{t=1}^T (\mathbf{b} \times X_t)^4 - 3 \rightarrow \max \quad \text{under constraint } \mathbf{b} \times \mathbf{b}^T = 1$$

Gradient descent method is usually used

Let $\mathbf{w} = \mathbf{b}^T$ - desired vector-column of unmixing matrix \mathbf{W} and $F(\mathbf{b}) = F(\mathbf{w})$

$$\nabla F(\mathbf{w}) = \frac{\partial}{\partial \mathbf{w}} |\text{Kurt}(\mathbf{w}^T \times \mathbf{X})| = 4 \times \text{sign}(\text{Kurt}(\mathbf{w}^T \times \mathbf{X})) \times \{ \mathbf{M}(\mathbf{X} \times (\mathbf{w}^T \times \mathbf{X})^3) - 3\mathbf{w} \times \underbrace{\|\mathbf{w}\|^2}_{=1} \}$$

Iterations in the gradient descent method: $\mathbf{w} \leftarrow \mathbf{w} + \Delta \mathbf{w}$, $\Delta \mathbf{w} \propto \nabla F(\mathbf{w})$

- since $\|\mathbf{w}\|^2 = 1$, iterations $\mathbf{w} \leftarrow \mathbf{w} + \Delta \mathbf{w}$ must be complemented by projecting on the unit sphere after every step: $\mathbf{w} \leftarrow \mathbf{w} / \|\mathbf{w}\|$
- the term $-3\mathbf{w}$ in the gradient $\nabla F(\mathbf{w})$ can be omitted: it changes the norm of \mathbf{w} in the gradient algorithm which will be normalized

Iterations in the gradient descent method:

- $\mathbf{w} \leftarrow \mathbf{w} + \Delta \mathbf{w}, \quad \Delta \mathbf{w} \propto \text{sign}(\text{Kurt}(\mathbf{w}^T \times \mathbf{X})) \times \mathbf{M}(\mathbf{X} \times (\mathbf{w}^T \times \mathbf{X})^3)$
- $\mathbf{w} \leftarrow \mathbf{w} / \|\mathbf{w}\| \quad \mathbf{M}(\mathbf{X} \times (\mathbf{w}^T \times \mathbf{X})^3) = \frac{1}{T} \sum_{t=1}^T X_t \times (\mathbf{w}^T \times X_t)^3$

A fast fixed-point algorithm using kurtosis- FastICA algorithm (kurtosis version)

Math: at a stable point of the gradient algorithm, the gradient $\nabla F(\mathbf{w})$ must point in the direction of \mathbf{w}

- the gradient $\nabla F(\mathbf{w})$ must be equal to \mathbf{w} multiplied by some scalar constant - only in such a case, adding the gradient to \mathbf{w} does not change its direction

Iterations in the FastICA method:

- $\mathbf{w} \leftarrow \mathbf{M}(\mathbf{X} \times (\mathbf{w}^T \times \mathbf{X})^3) - 3\mathbf{w}$
- $\mathbf{w} \leftarrow \mathbf{w} / \|\mathbf{w}\|$

Kurtosis-based measure of non-Gaussianity has some drawbacks in practice, when kurtosis value has to be estimated from a measured sample: **kurtosis can be very sensitive to outliers**

Kurtosis is not a robust measure of nongaussianity

2) Negentropy: $p_\xi(y)$ - density of random variable ξ

Entropy (measure of uncertainty): $H(\xi) = -M(\log_2 p_\xi(\xi)) = -\int p_\xi(y) \times \log_2 p_\xi(y) dy$

Probability theory: which random variable ξ has maximal entropy?

The answer depends on available information about ξ (its density p_ξ)

Let available information has a form: for given functions $h_i(\xi)$, $i = 1, 2, \dots, m$:

$$Mh_i(\xi) = \int p_\xi(y) \times h_i(y) dy = k_i$$

Example: $h_1(\xi) = \xi$, $k_1 = 0$; $h_2(\xi) = \xi^2$, $k_2 = 1$

Under these constraints, a density with maximum entropy is

$$p_0(y) = A \times \exp\{\sum_{i=1}^m \alpha_i \times h_i(y)\}$$

Example: $p_0(y) = A \times \exp\{\alpha_1 \times y + \alpha_2 \times y^2\}$ - Gaussian density

Negentropy: $J(\xi) = H(\xi_{\text{Gauss}}) - H(\xi) \geq 0$

$H(\xi_{\text{Gauss}})$ - entropy of Gaussian random variable $\xi_{\text{Gauss}} \sim$ such that $\text{Var}(\xi_{\text{Gauss}}) = \text{Var}(\xi)$

$J(\xi_{\text{Gauss}}) = 0 \quad \rightarrow \quad f(\xi) = J(\xi)$ - a measure of distance from Gaussianity

$\xi = \mathbf{b} \times \mathbf{X} \quad \rightarrow \quad \mathbf{cov}(\xi) = \mathbf{M}(\mathbf{b} \times \mathbf{X})^2 = 1 \quad \rightarrow$

$\xi_{\text{Gauss}} \sim N(0, 1) \quad \rightarrow \quad H(\xi_{\text{Gauss}}) = \frac{1}{2}(1 + \log_2(2\pi))$

$J(\mathbf{b} \times \mathbf{X}) = \frac{1}{2}\{1 + \log_2(2\pi)\} - H(\mathbf{b} \times \mathbf{X}) = \frac{1}{2}\{1 + \log_2(2\pi)\} + \mathbf{M} \log_2(p_s(\mathbf{b} \times \mathbf{X}))$

$= \frac{1}{2}\{1 + \log_2(2\pi)\} + \frac{1}{T} \sum_{t=1}^T \{\log_2(p_s(\mathbf{b} \times X_t))\}$ p_s - density function of random signal \mathbf{s}

$F(\mathbf{b}) = \frac{1}{T} \sum_{t=1}^T \{\log_2(p_s(\mathbf{b} \times X_t))\} \rightarrow \max$ under constraint $\mathbf{b} \times \mathbf{b}^T = 1$

The trouble: computationally very difficult, density function p_s of random signal \mathbf{s} is unknown and should be estimated \rightarrow simpler approximations of the negentropy/entropy are used

Approximation 1 based on Gram-Charlier expansion for density estimation

Probability theory: $p_\xi(y)$ - unknown density of random variable ξ $H_3(y), H_4(y)$ - Chebyshev-Hermite polynomials of order 3 and 4
 $\varphi(y) = (2\pi)^{-1/2} \times \exp(-y^2/2)$ - density of $N(0, 1)$

$$p_\xi(y) \approx p^*(y) = \varphi(y) \times \left\{ 1 + k_3(\xi) \times \frac{H_3(y)}{3!} + k_4(\xi) \times \frac{H_4(y)}{4!} \right\} \quad K_3(\xi) = M\xi^3, K_4(\xi) = M\xi^4 - 3$$

Negentropy approximation:

$$H^*(\xi) = - \int p^*(y) \times \log_2 p^*(y) dy \approx H(\xi) \quad J^*(\xi) = \frac{1}{12} (M\xi^3)^2 + \frac{1}{48} (\mathbf{Kurt}(\xi))^2 \approx J(\xi)$$

$M\xi^3 = 0$ for random variables with symmetric distributions $\rightarrow J^*(\xi) = \frac{1}{48} (\mathbf{Kurt}(\xi))^2$

\rightarrow maximization of $f(\xi) = J^*(\xi)$ is equivalent to maximization of $|\mathbf{Kurt}(\xi)|$

$$p_{\xi}(y) \approx p^*(y) = \varphi(y) \times \left\{ 1 + k_3(\xi) \times \frac{H_3(y)}{3!} + k_4(\xi) \times \frac{H_4(y)}{4!} \right\}$$

Approximation 2 replaces polynomial functions $H_3(y)$ and $H_4(y)$ by other functions and is based on ‘maximum entropy’ density estimation

Assumption: $M\xi = 0$, $M\xi^2 = 1$ and we have observed (in practice, estimated) a number of expectations $Mh_i(\xi)$ of functions $\{h_i(\xi)\}$:

$$Mh_i(\xi) = \int p_{\xi}(y) \times h_i(y) dy = k_i, \quad i = 1, 2, \dots, m,$$

which:

- form orthonormal system: $\int \varphi(y) \times h_i(y) \times h_j(y) dy = \delta_{ij}$ (Kronecker symbol)
- orthogonal to all polynomials of second degree: $\int \varphi(y) \times h_i(y) \times y^k dy = 0, \quad k = 0, 1, 2$

We can choose arbitrary not-orthonormal functions $\{h_i(\xi)\}$ and then use standard Gram-Schmidt orthonormalization technique

Approximative maximum entropy density: $p_{\xi}(y) \approx p^{**}(y) = \varphi(y) \times \{1 + \sum_{i=1}^m k_i \times h_i(y)\}$

with nonquadratic functions $\{h_i(y)\}$

$$H^{**}(\xi) = - \int p^{**}(y) \times \log_2 p^{**}(y) dy \approx H(\xi)$$

Negentropy approximation: $J^{**}(\xi) = \sum_{i=1}^m k_i \times \{\mathbf{M}h_i(\xi) - \mathbf{M}h_i(\xi_{\text{Gauss}(0,1)})\}^2 \approx J(\xi)$

$\{k_i\}$ - positive constants, $\xi_{\text{Gauss}(0,1)} \sim N(0, 1)$

How many and how to choose the functions $\{h_i(\xi)\}$?

$m = 1$: $J^{**}(\xi) = \{\mathbf{M}h(\xi) - \mathbf{M}h(\xi_{\text{Gauss}(0,1)})\}^2$ h - arbitrary nonquadratic function

ξ - symmetric and $h(y) = y^4$ - kurtosis based approximation

$$m = 2: J^{**}(\xi) = k_1 \times \{Mh_1(\xi) - Mh_1(\xi_{\text{Gauss}(0,1)})\}^2 + k_2 \times \{Mh_2(\xi) - Mh_2(\xi_{\text{Gauss}(0,1)})\}^2$$

The following choices of $\{h_i(y)\}$ have proved very useful ($m = 2$)

$$h_1(y) = \frac{1}{\alpha_1} \log \cosh(\alpha_1 y), \quad 1 \leq \alpha_1 \leq 2; \quad h_2(y) = -\exp(-y^2/2)$$

Gradient descent algorithm - ($m = 1$)

$$F(w) = \{\mathbf{M}h(w^T \times \mathbf{X})) - \mathbf{M}h(\xi_{\text{Gauss}(0,1)})\}^2$$

$$h(y) = \frac{1}{\alpha_1} \log \cosh(\alpha_1 y)$$

$$\nabla F(w) = \mathbf{M}(\mathbf{X} \times g(w^T \times \mathbf{X}))$$

$$g(y) = h'(y) = \tanh(\alpha_1 y)$$

Iterations in the gradient descent method:

- $w \leftarrow w + \Delta w, \quad \Delta w \propto \mathbf{M}(\mathbf{X} \times g(w^T \times \mathbf{X}))$

- $w \leftarrow w / \|w\| \quad \mathbf{M}(\mathbf{X} \times g(w^T \times \mathbf{X})) = \frac{1}{T} \sum_{t=1}^T X_t \times g(w^T \times X_t)$

A fast fixed-point algorithm - FastICA ($m = 1$)

Iterations in the FastICA method:

- $w \leftarrow \mathbf{M}(\mathbf{X} \times h(w^T \times \mathbf{X})) - \mathbf{M}(g(w^T \times \mathbf{X})) \times w$
- $w \leftarrow w / \|w\|$

Fundamental approaches to the solution of the ICA problem:

2) Likelihood maximization

Probability theory: \mathbf{X} and \mathbf{S} - random vectors with densities $p_{\mathbf{X}}(\mathbf{x})$ and $p_{\mathbf{S}}(\mathbf{s})$, respectively

$$\mathbf{X} = \mathbf{A} \times \mathbf{S} \quad \rightarrow \quad p_{\mathbf{X}}(\mathbf{x}) = |\text{Det}(\mathbf{B})| \times p_{\mathbf{S}}(\mathbf{s}), \quad \mathbf{B} = \mathbf{A}^{-1}$$

$$\mathbf{S} = \mathbf{B} \times \mathbf{X} \quad \mathbf{B} = (b_1, b_2, \dots, b_q)^T$$

s_1, s_2, \dots, s_q - independent signals: $p_{\mathbf{S}}(\mathbf{s}) = p_{\mathbf{S}}(s_1, s_2, \dots, s_q) = p_1(s_1) \times p_2(s_2) \times \dots \times p_q(s_q)$

$$s_i = b_i^T \times \mathbf{X} \quad p_{\mathbf{X}}(\mathbf{x}) = |\text{Det}(\mathbf{B})| \times p_1(b_1^T \times \mathbf{X}) \times p_2(b_2^T \times \mathbf{X}) \times \dots \times p_q(b_q^T \times \mathbf{X})$$

X_1, X_2, \dots, X_T - training dataset

Likelihood: $L(\mathbf{B}, p_s) = \left\{ \prod_{t=1}^T \prod_{i=1}^q p_i(b_i^T x_t) \right\} \times |\text{Det}(\mathbf{B})|^T \rightarrow \max$

If density p_s of independent signals are known from some prior knowledge, likelihood depends only on unknown matrix \mathbf{B}

$\mathbf{W} = \arg \max L(\mathbf{B}, p_s)$ - desired unmixing matrix

Under unknown density p_s , it is approximated by some estimator p^* and $L(\mathbf{B}, p^*)$ depends only on \mathbf{B} :

$\mathbf{W} = \arg \max L(\mathbf{B}, p^*)$ - desired unmixing matrix

Maximum likelihood approach:

- a choice of 'appropriate' densities $p^*(s)$
- maximization $L(\mathbf{B}, p^*)$ over \mathbf{B}

Math: $\{p_i\}$ - chosen densities, $p^*(s_1, s_2, \dots, s_q) = p_1(s_1) \times p_2(s_2) \times \dots \times p_q(s_q)$.

Maximum likelihood estimator $\mathbf{W} = \arg \max L(\mathbf{B}, p^*)$ is locally consistent, if the assumed densities meet the conditions

$$\text{Index } \Psi(p_i) = \mathbf{M}\{s_i \times g_i(s_i) - (g_i(s_i))'\} > 0, \quad g_i(s) = (\log p_i(s))'$$

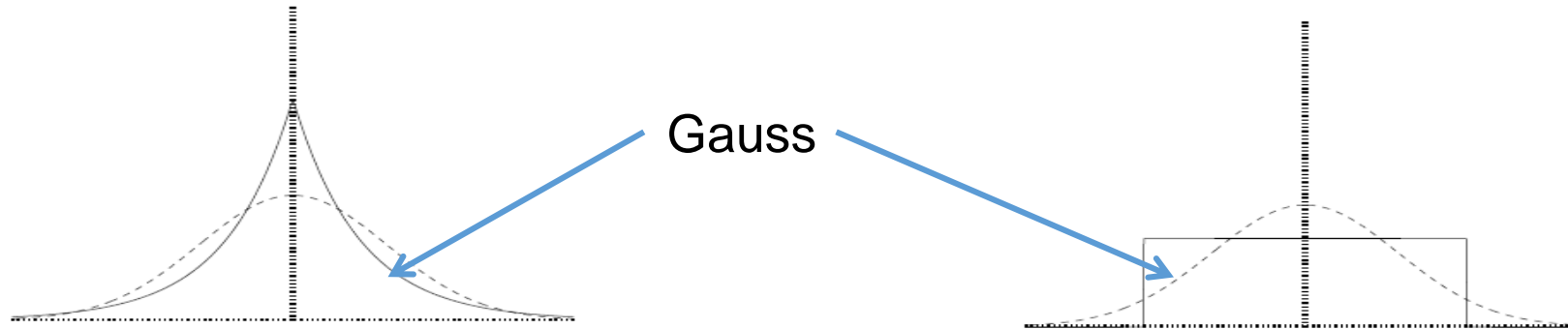
for all i .

Example:

$$\log(p_1(s)) = \alpha_1 - 2\log \cosh(s)$$

$$\log(p_2(s)) = \alpha_2 - [s^2/2 - \log \cosh(s)]$$

$p_1(s)$ - supergaussian density with positive kurtosis - with 'heavy' tails (the density is relatively large at zero and at large values of the variable, while being small for intermediate values)



$p_2(s)$ - subgaussian density with negative kurtosis - have typically a "flat" density which is rather constant near zero, and very small for larger values of the variable

$$\log(p_1(s)) = \alpha_1 - 2\log \cosh(s)$$

$$\log(p_2(s)) = \alpha_2 - [s^2/2 - \log \cosh(s)]$$

$$\Psi(p_1) = 2\mathbf{M}\{-\tanh(s) \times s + (1 - \tanh^2 s)\}$$

$$\Psi(p_2) = \mathbf{M}\{\tanh(s) \times s - (1 - \tanh^2 s)\}$$

$$\Psi(p_i) = -2\Psi(p_2)$$

Maximization $L(\mathbf{B}, p) = \{\prod_{t=1}^T \prod_{i=1}^q p(\mathbf{b}_i^T \times X_t)\} \times |\text{Det}(\mathbf{B})|^T$ over \mathbf{B} $p_i(s) = p(s)$

$$L(\mathbf{B}, p) \quad \rightarrow \quad \frac{1}{T} \log L(\mathbf{B}, p) = \sum_{t=1}^T \sum_{i=1}^q \log p(\mathbf{b}_i^T \times X_t) + \log |\text{Det} \mathbf{B}|$$

$$\nabla = \frac{1}{T} \frac{\partial}{\partial \mathbf{B}} \log L(\mathbf{B}, p) = \frac{1}{T} \sum_{t=1}^T \mathbf{g}(\mathbf{B} \times X_t) \times X_t^T + [\mathbf{B}^T]^{-1}$$

$$\mathbf{g}(\mathbf{B} \times X_t) = \begin{pmatrix} g(\mathbf{b}_1^T \times X_t) \\ \vdots \\ g(\mathbf{b}_q^T \times X_t) \end{pmatrix} \quad g(s) = (\log p(s))'$$

Iterations in the gradient descent maximization:

$$\mathbf{B} \leftarrow \mathbf{B} + \Delta \mathbf{B}, \quad \Delta \mathbf{B} \propto \frac{1}{T} \sum_{t=1}^T \mathbf{g}(\mathbf{B} \times X_t) \times X_t^T + [\mathbf{B}^T]^{-1}$$

There exists fast fixed-point algorithm for likelihood maximization - FastICA

Fundamental approaches to the solution of the ICA problem:

3) Minimization of Mutual Information

Probability theory: \mathbf{S} - random vector

Mutual information $I(\mathbf{s}_1, \mathbf{s}_2, \dots, \mathbf{s}_q)$ between $\mathbf{s}_1, \mathbf{s}_2, \dots, \mathbf{s}_q$ is:

$$I(\mathbf{s}_1, \mathbf{s}_2, \dots, \mathbf{s}_q) = \sum_{i=1}^q H(\mathbf{s}_i) - H(\mathbf{s}_1, \mathbf{s}_2, \dots, \mathbf{s}_q)$$

$$p_{\mathbf{S}}(\mathbf{s}) = p_{\mathbf{S}}(\mathbf{s}_1, \mathbf{s}_2, \dots, \mathbf{s}_q) \rightarrow$$

$$H(\mathbf{s}_1, \mathbf{s}_2, \dots, \mathbf{s}_q) = - \int p_{\mathbf{S}}(y_1, y_2, \dots, y_q) \times \log_2 p_{\mathbf{S}}(y_1, y_2, \dots, y_q) dy_1 dy_2 \dots dy_q$$

$$p_i(\mathbf{s}_i) \rightarrow H(\mathbf{s}_i) = - \int p_i(y) \times \log_2 p_i(y) dy$$

s_1, s_2, \dots, s_q - independent signals: $p_{\mathbf{s}}(\mathbf{s}) = p_{\mathbf{s}}(s_1, s_2, \dots, s_q) = p_1(s_1) \times p_2(s_2) \times \dots \times p_q(s_q)$

$$H(s_1, s_2, \dots, s_q) = \sum_{i=1}^q H(s_i) \quad \rightarrow \quad I(s_1, s_2, \dots, s_q) = 0$$

Probability theory: Kullback-Leibler divergence.

p_1, p_2 - densities $\rightarrow d(p_1, p_2) = \int p_1(y) \log \frac{p_1(y)}{p_2(y)} dy \geq 0$

- kind of a distance between the two probability densities, because $d(p_1, p_2)$ is always nonnegative, and zero if and only if $p_1(y) = p_2(y)$

$$p_1 = p_{\mathbf{s}}(s_1, s_2, \dots, s_q), p_2 = p_1(s_1) \times p_2(s_2) \times \dots \times p_q(s_q)$$

$$\rightarrow d(p_1, p_2) = I(s_1, s_2, \dots, s_q)$$

Mutual information $I(s_1, s_2, \dots, s_q)$ has minimal value when s_1, s_2, \dots, s_q - independent

$$\mathbf{S} = \mathbf{B} \times \mathbf{X} \quad \rightarrow \quad p_{\mathbf{X}}(\mathbf{x}) = |\text{Det}(\mathbf{B})| \times p_{\mathbf{S}}(\mathbf{s})$$

$$H(\mathbf{S}) = H(\mathbf{X}) + \log |\text{Det}(\mathbf{B})|$$

$$I(s_1, s_2, \dots, s_q) = \sum_{i=1}^q H(s_i) - H(\mathbf{X}) - \log |\text{Det}(\mathbf{B})|$$

$$\mathbf{B} \text{ - orthogonal matrix} \quad \rightarrow \quad \log |\text{Det}(\mathbf{B})| = 0$$

$$s_i = \mathbf{b}_i^T \times \mathbf{X}: \quad \sum_{i=1}^q H(\mathbf{b}_i^T \times \mathbf{X}) = - \mathbf{M} \sum_{i=1}^q \log_2 p(\mathbf{b}_i^T \times \mathbf{X}) \rightarrow \min$$

$$F(\mathbf{b}) = \frac{1}{T} \sum_{t=1}^T \sum_{i=1}^q \log_2 p(\mathbf{b}_i^T \times \mathbf{X}_t) \rightarrow \max \quad \text{under constraint } \mathbf{b} \times \mathbf{b}^T = \mathbf{1}$$

- Negentropy maximization

Mutual information minimization - Negentropy maximization