

FEATURE EXTRACTION AND SELECTION METHODS

1

Feature extraction and selection methods

- The task of the feature extraction and selection methods is to obtain the most ***relevant information*** from the original data and represent that information in a ***lower dimensionality*** space.

2

Selection methods

- When the cost of the acquisition and manipulation of all the measurements is high we must make a selection of features.
- The goal is to select, **among all the available features**, those that will perform better.
- Example: which features should be used for classifying a student as a good or bad one:
 - Available features: marks, height, sex, weight, IQ.
 - Feature selection would choose **marks** and **IQ** and would discard height, weight and sex.
- We have to choose P variables in a set of M variables so that the separability is maximal.

3

Extraction methods

- The goal is to **build**, using the available features, those that will perform better.
- Example: which features should be used for classifying a student as a good or bad one:
 - Available features: marks, height, sex, weight, IQ.
 - Feature extraction may choose **marks + IQ²** as the best feature (in fact, it is a combination of two features).
- The goal is to transform the origin space X in a new space Y to obtain new features that work better. This way, we can compress the information.

4

Principal Component Analysis

- PCA = Karhunen-Loeve transform = Hotelling transform
- PCA is the most popular feature extraction method
- PCA is a linear transformation
- PCA is used in face recognition systems based on appearance

5

Principal Component Analysis

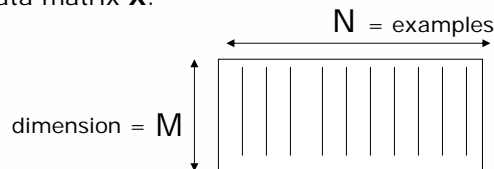
- PCA has been successfully applied to human face recognition.
- PCA consists on a transformation from a space of high dimension to another with more reduced dimension.
- If the data are highly correlated, there is redundant information.
 - PCA decreases the amount of redundant information by decorrelating the input vectors.
 - The input vectors, with high dimension and correlated, can be represented in a lower dimension space and decorrelated.
 - PCA is a powerful tool to compress data.

6

PCA by Maximizing Variance (I)

We will derive PCA by maximizing the variance in the direction of principal vectors.

Let us suppose that we have N M -dimensional vectors \mathbf{x}_j aligned in the data matrix \mathbf{X} .



Let \mathbf{u} be a direction (a vector of length 1). The projection of the j -th vector \mathbf{x}_j onto the vector \mathbf{u} can be calculated in the following way:

$$p_j = \vec{u}^T \cdot \vec{x}_j = \sum_{i=1}^M u_i x_{ij}$$

7

PCA by Maximizing Variance (II)

We want to find a direction \mathbf{u} that maximizes the variance of the projections of all input vectors \mathbf{x}_j , $j=1, \dots, N$.

The function to maximize is:

$$J^{PCA}(\vec{u}) = \sigma^2(p_j) = \frac{1}{N} \sum_{j=1}^N (p_j - \bar{p})^2 = \dots = \vec{u}^T \mathbf{C} \vec{u}$$

where \mathbf{C} is the covariance matrix of the data matrix \mathbf{X} .

$$\mathbf{C} = \frac{1}{N} \hat{\mathbf{X}} \cdot \hat{\mathbf{X}}^T \quad \hat{\mathbf{X}} = \mathbf{X} - \boldsymbol{\mu} \cdot \mathbf{1}_{1 \times N}$$

$$\boldsymbol{\mu} = [\mu_1, \dots, \mu_m]^T$$

Using the technique of Lagrange multipliers, the solution to this maximization problem is to compute the **eigenvectors** and the **eigenvalues** of the covariance matrix \mathbf{C} .

MORE INFO in [PCA.pdf](#)

8

PCA by Maximizing Variance (III)

The largest eigenvalue equals the maximal variance, while the corresponding eigenvector determines the direction with the maximal variance.

By performing singular value decomposition (SVD) of the covariance matrix \mathbf{C} we can diagonalize \mathbf{C} :

$$\mathbf{C} = \mathbf{U}\mathbf{\Lambda}\mathbf{U}^T$$

in such a way that the orthonormal matrix \mathbf{U} contains the eigenvectors $\mathbf{u}_1, \mathbf{u}_2, \dots, \mathbf{u}_N$ in its columns and the diagonal matrix $\mathbf{\Lambda}$ contains the eigenvalues $\lambda_1, \lambda_2, \dots, \lambda_N$ on its diagonal.

The eigenvalues and the eigenvectors are arranged with respect to the descending order of the eigenvalues, thus $\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_N$. Therefore, the most of the variability of the input random vectors is contained in the first eigenvectors. Hence, the eigenvectors are called principal vectors.

9

Computing PCA

Steps to compute the PCA transformation of a data matrix \mathbf{X} :

- Center the data $\bar{\mathbf{X}}$
- Compute the covariance matrix \mathbf{C}
- Obtain the eigenvectors and eigenvalues of the covariance matrix $\mathbf{U}, \mathbf{\Lambda}$
- Project the original data in the eigenspace $\mathbf{P} = \mathbf{U}^T \cdot \bar{\mathbf{X}}$

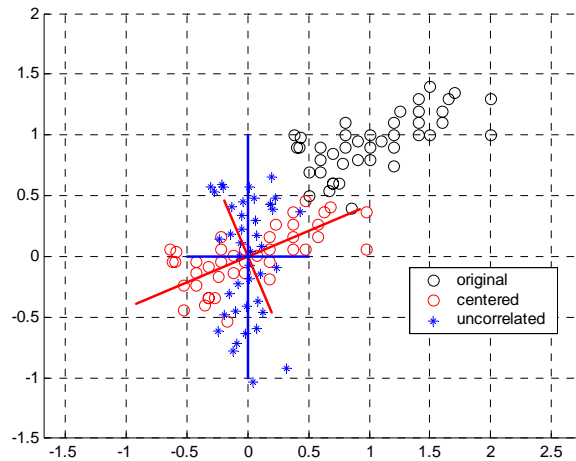
\mathbf{U} can be used as a linear transformation to project the original data of high dimension into a space of lower dimension.

Matlab code:

```
%number of examples
N=size(X,2);
%dimension of each example
M=size(X,1);
%mean
meanX=mean(X,2);
%centering the data
Xm=X-meanX*ones(1,N);
%covariance matrix
C=(Xm*Xm')/N;
%computing the eigenspace:
[U D]=eig(C);
%projecting the centered data
over the eigenspace
P=U'*Xm;
```

10

PCA of a bidimensional dataset



11

Computing PCA of a set of images

This approach to the calculation of principal vectors is very clear and widely used. However, if the size of the data vector M is very large, which is often the case in the field of computer vision, the covariance matrix \mathbf{C} becomes very large and eigenvalue decomposition of \mathbf{C} becomes unfeasible.

But, if the number of input vectors is smaller than the size of these vectors ($N < M$), PCA can be sped up using a method proposed by Murakami and Kumar which is also known as Turk and Pentland's trick.

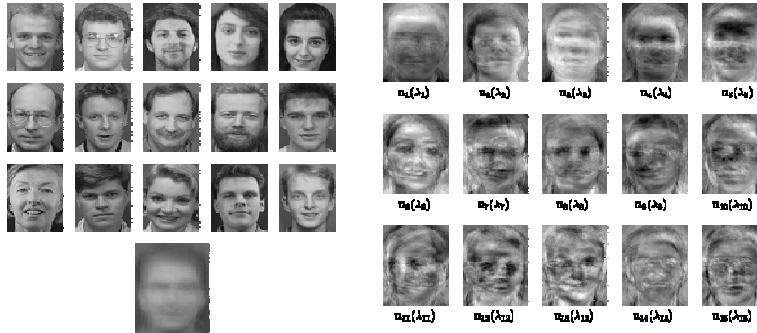
When we work on Matlab, we use `pc_evectors.m` to apply the Turk and Pentland's trick and compute the PCA transformation of a set of images

MORE INFO in [PCA.pdf](#)

12

Face recognition using PCA (I)

Eigenfaces for Recognition, Turk, M. & Pentland, A. ,
Journal of Cognitive Neuroscience, 3, 71-86, 1991.



13

Linear Discriminant Analysis (I)

- LDA = Fisher analysis
- LDA is a linear transformation
- LDA is also used in face recognition
- LDA seeks directions that are efficient for discrimination between classes

- In PCA, the subspace defined by the vectors is the one that better describes the conjunct of data.
- LDA tries to discriminate between the different classes of data.

14

Linear Discriminant Analysis (I)

- We have a conjunct of N vectors of dimension M in the data matrix $M \times N$.
- We have C classes and k vectors per class.
- We want to find the transformation matrix W that better describes the subspace that discriminates between classes, after projecting the data in the new space.
- The objective is to make maximum the distance between classes S_b and minimizing S_w .

$$P = W \cdot X \quad W \text{ are the eigenvectors of } S_w^{-1} S_b$$

"between class" scatter matrix

$$S_b = \sum_j^C (\mu_j - \mu)(\mu_j - \mu)^T$$

"within-class" scatter matrix

$$S_w = \sum_{j=1}^C \sum_{i=1}^K (x_i^j - \mu_j)(x_i^j - \mu_j)^T = \sum_{j=1}^C \hat{X}^j \hat{X}^{jT}$$

$$X = [X^1 | X^2 | \dots | X^j | \dots | X^C]$$

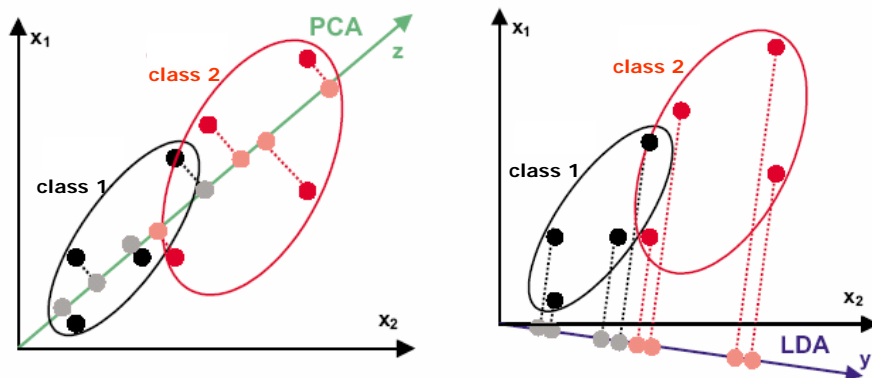
$$X^1 = [x_1^1, x_2^1, \dots, x_i^1, \dots, x_K^1]$$

$$X^j = [x_1^j, x_2^j, \dots, x_i^j, \dots, x_K^j]$$

$$X^C = [x_1^C, x_2^C, \dots, x_i^C, \dots, x_K^C]$$

Linear Discriminant Analysis (II)

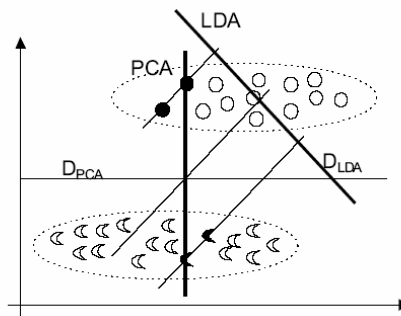
- The figure shows the effect of LDA transform in a conjunct of data composed of 2 classes.



16

Linear Discriminant Analysis (III)

Limitations of LDA



- LDA works better than PCA when the training data are well representative of the data in the system.
- If the data are not representative enough, PCA performs better.

17

Independent Component Analysis (I)

•Independent Component Analysis

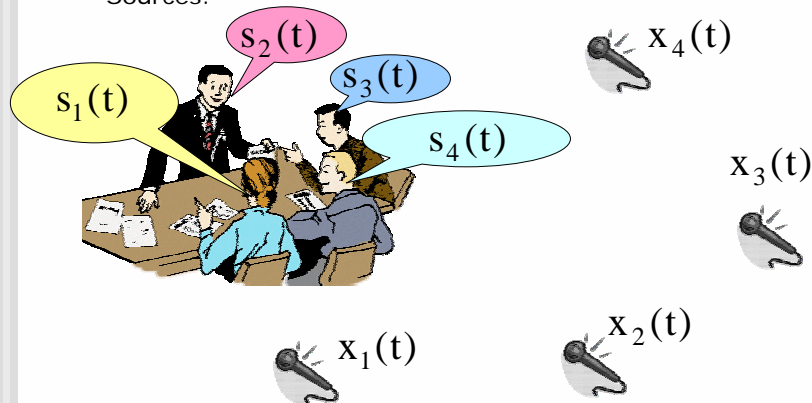
- ICA is a statistical technique that represents a multidimensional random vector as a linear combination of nongaussian random variables ('independent components') that are as independent as possible.
- ICA is somewhat similar to PCA.
- ICA has many applications in data analysis, source separation, and feature extraction.

18

ICA – cocktail party problem

- Cocktail party problem

- ICA is a statistical technique for decomposing a complex dataset into independent sub-parts. Here we show how it can be applied to the problem of separation of Blind Sources.



19

ICA – cocktail party problem

- Cocktail party problem

Estimate the sources $s_i(t)$
from the mixed signals $x_i(t)$

20

ICA – cocktail party problem

- Linear model:

$$x_1(t) = a_{11}s_1(t) + a_{12}s_2(t) + a_{13}s_3(t) + a_{14}s_4(t)$$

$$x_2(t) = a_{21}s_1(t) + a_{22}s_2(t) + a_{23}s_3(t) + a_{24}s_4(t)$$

$$x_3(t) = a_{31}s_1(t) + a_{32}s_2(t) + a_{33}s_3(t) + a_{34}s_4(t)$$

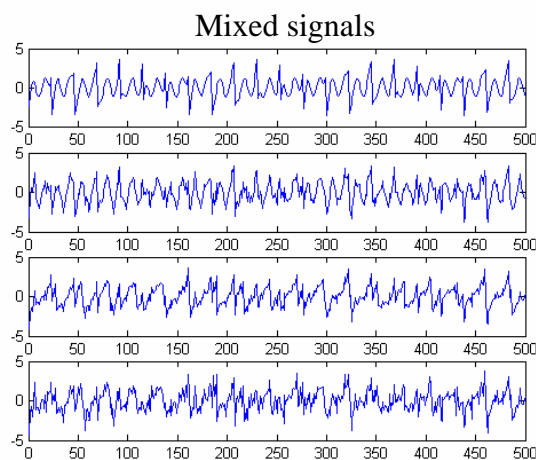
$$x_4(t) = a_{41}s_1(t) + a_{42}s_2(t) + a_{43}s_3(t) + a_{44}s_4(t)$$

- We can model the problem as $X=AS$

- S = 4D vector containing the independent source signals.
- A = mixing matrix.
- X = Observed signals.

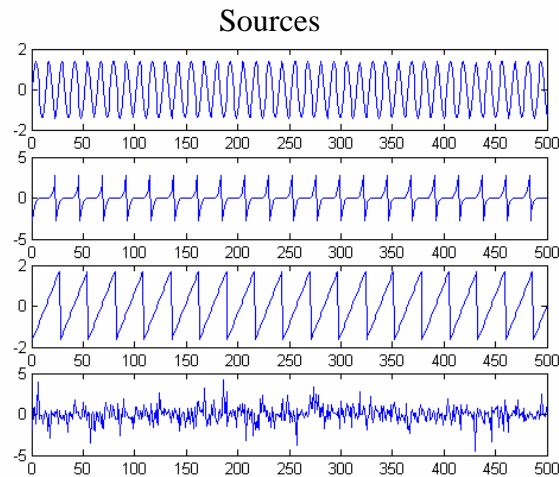
21

ICA – cocktail party problem



22

ICA – cocktail party problem



23

ICA – cocktail party problem

Estimate the sources $s_i(t)$
from the mixed signals $x_i(t)$

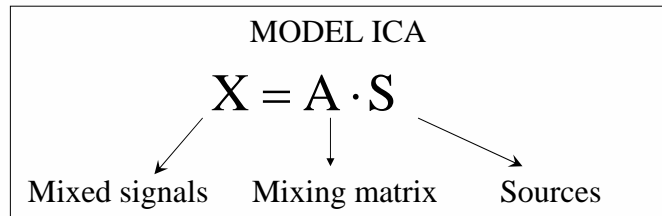


ICA: One possible solution is to assume that the sources are **independent**.

$$p(s_1, s_2, \dots, s_n) = p(s_1)p(s_2) \dots p(s_n)$$

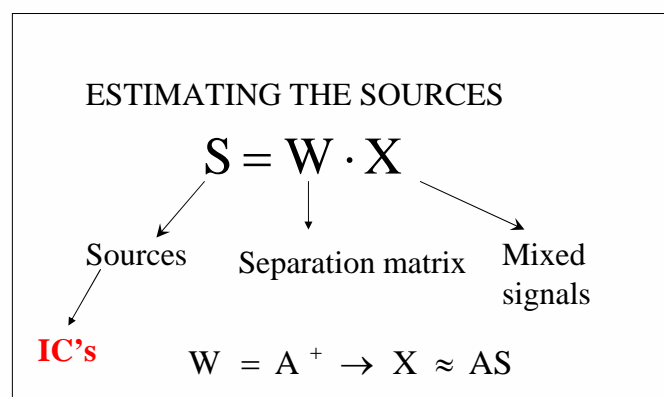
24

ICA – cocktail party problem



25

ICA – cocktail party problem



26

Computing IC's

- Typically, in ICA algorithms, \mathbf{W} is sought such that the rows of it have maximally non-gaussian distributions and are mutually uncorrelated.
- A simple way to do this is to first whiten the data and then seek orthogonal non-normal projections.
- We want to find arrows $\omega_i / s_i = \mathbf{w}_i^T \cdot \mathbf{x}$ have maximally non-gaussian distributions and mutually uncorrelated.

27

PCA, WHITENING, ICA (I)

- **PCA:**
uncorrelated data
(the covariance matrix of the PCA transformed data has the eigenvalues in its diagonal)
- **WHITENING:**
PCA + scaling
(the covariance matrix of the whiten data is the identity)
- **ICA:**
WHITENING + rotation

28

WHITENING

■ WHITENING:

PCA + scaling

$$\mathbf{X} \approx \sum_{j=1}^K \sigma_j \mathbf{u}_j \mathbf{v}_j^T = \tilde{\mathbf{U}} \tilde{\mathbf{\Lambda}} \tilde{\mathbf{V}}^T$$

$$\tilde{\mathbf{Z}} = \tilde{\mathbf{\Lambda}}^{-1} \tilde{\mathbf{U}}^T \mathbf{X} = \tilde{\mathbf{V}}^T$$

29

ICA (I)

■ ICA:

WHITENING + rotation

$$\mathbf{S} = \mathbf{R}^T \tilde{\mathbf{\Lambda}}^{-1} \tilde{\mathbf{U}}^T \mathbf{X} = \mathbf{R}^T \tilde{\mathbf{Z}}$$

\mathbf{R} is a rotation that maximizes the non-gaussianity of the projections

30

ICA (II)

- ICA model:

$$X \approx \widetilde{U}\widetilde{\Lambda}\widetilde{V}^T = \underbrace{\widetilde{U}\widetilde{A}}_A \underbrace{R^T\widetilde{V}^T}_S = AS$$

31

ICA (III)

FastICA

[FastICA](#) : is a free MATLAB program that implements the fast-fixed point algorithm.

32

PCA, WHITENING, ICA (II)

ICA transformation by FastICA

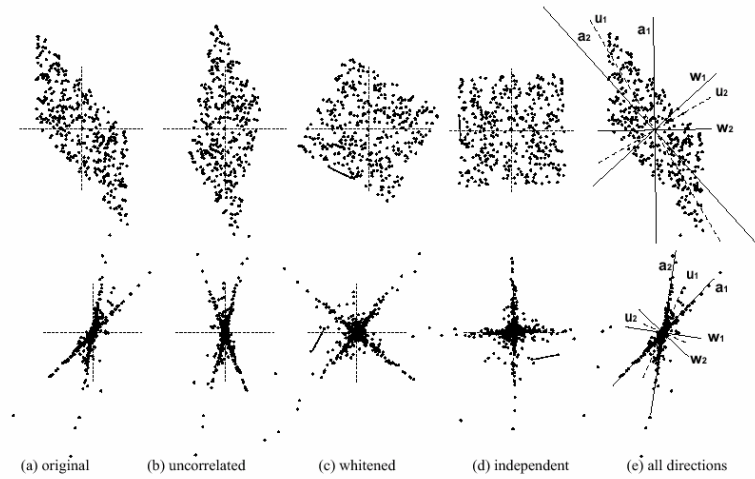
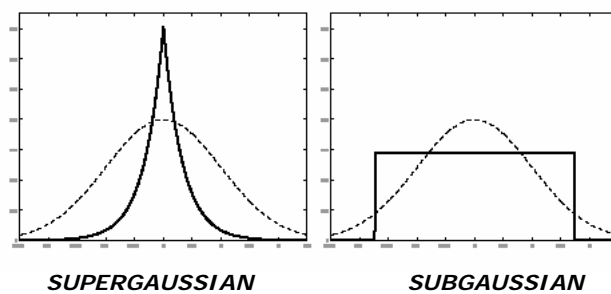


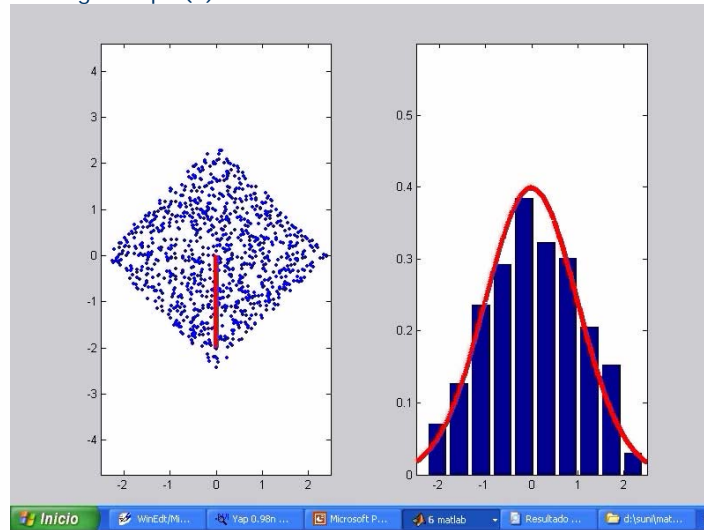
Fig. 2. Two artificial examples: a subgaussian dataset (1st. row) and a supergaussian dataset (2nd. row) (a), both transformed by PCA (b), whitening (c) and using the ICA model by means of FastICA (d). The last figure in each row (e) shows the original dataset with the ICA and PCA directions.

Non gaussianity (I)



No gaussianity (II)

generateNongExample(1)



35

ICA in CNS (Computational Neuroscience) (I)

- BSS applications with EEG and MEG signals.
 - The brains activity is measured through Electroencephalograms. Those signals are a mixture of different activities in the brain and other external noises.
 - ICA solves correctly the problem of extracting the original activity signals
- Modeling the performance of the neurons in area V1 of mammalian cortex.
 - Spikes
 - Receptive fields
 - Natural images
- Some studies propose that the behaviour of one kind of neurons can be computationally described through the ICA analysis of this natural inputs.

36

Spikes

- SPIKES: electrical signal in neurons

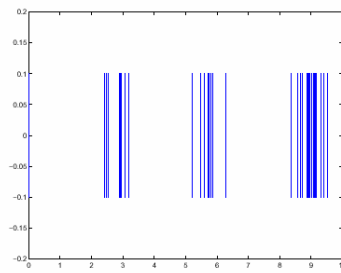


Figure 1: An inhomogeneous Poisson spike train

Receptive fields

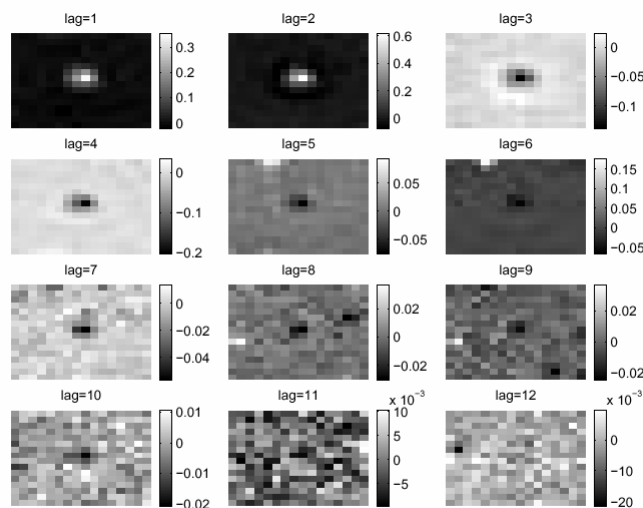
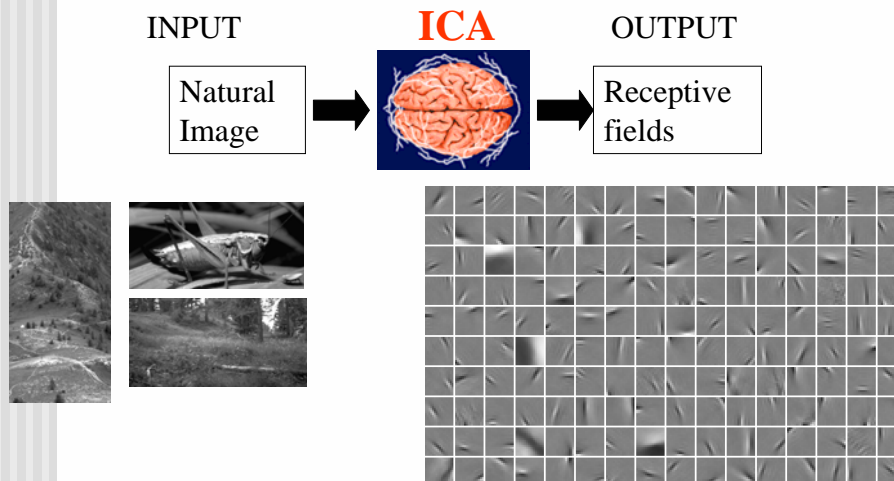


Figure 2: The spike-triggered average images for each of the 12 time steps

Simple experiment



NMF (I)

Non-negative matrix factorization (**NMF**) is a recently developed technique for finding parts, and it is based on linear representations of non-negative data.

Given a non-negative data matrix \mathbf{X} , NMF finds an approximate factorization

$$\mathbf{X} \approx \mathbf{W} \cdot \mathbf{H}$$

into non-negative factors \mathbf{W} and \mathbf{H} . The non-negativity constraints make the representation purely additive (allowing no subtractions), in contrast to many other linear representations such as PCA or ICA.

Motivation: In most real systems, the variables are non negative. PCA and ICA offer results complicated to interpret.

\mathbf{W} and \mathbf{H} are chosen as the matrix that minimize reconstruction error.

NMF (II)

NMF as a feature extraction method in faces



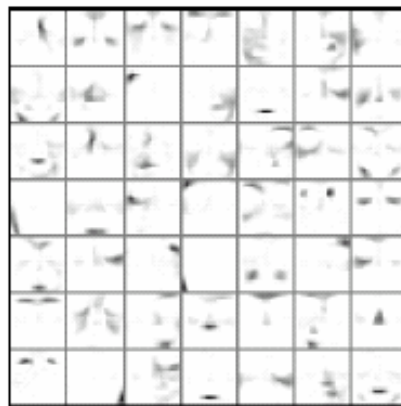
The importance of NMF is that it has capacity of obtaining significant features in collections of real biological data.

When applied to $X = \text{Faces}$, NFM generates base vectors that are intuitive features of the faces (eyes, mouth, nose...)

41

NMF (III)

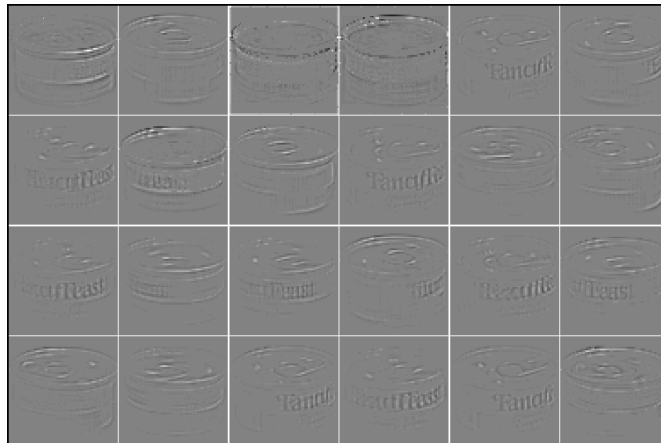
■ NMF → local features



42

NMF (IV)

■ NMF → local features



43

NMF (V)

- NMF presents features that make it adequate for applications in object recognition.
- It allows extracting local features as shown on the previous figure. Some images extract the text, others the top side, others the general shape of the object...
- It can be useful in presence of occlusions.
 - In this case, it is not possible to extract global features, but we can extract local ones.
- It can be useful also to identify objects in non-structured environments.
- At last, we can use it to extract categories of objects.

44