Unsupervised Learning Dimensionality Reduction - Latent Variable Analysis







- Autoencoders
- Nonlinear PCA, kernel PCA
- Sparse PCA
- Independent Component Analysis

590

Dimensionality Reduction Neural Networks

◆□▶ ◆□▶ ◆臣▶ ◆臣▶ 臣 のへで

Autoencoders

An autoencoder is a network with one input layer, one or more hidden layers and one output layer. This type of network aims at providing an internal representation (the layer in the middle) by learning how to predict the input from itself: $x \approx g(x)$.



- An auto-encoder is learned by backpropagation of the gradient. When the autoencoder is sparse, the loss considered in general is quadratic and penalized by an additive regularisation term, so that the activity of each unit of the hidden layer is limited in average.
- The autoencoder learns internal representations of complex data

Alternative Techniques for Dimensionality Reduction

◆□▶ ◆□▶ ◆三▶ ◆三▶ 三三 のへで

Example 1 - Finance

Analysis of interest rates



▲ロト ▲園ト ▲ヨト ▲ヨト 三国 - のくで

• Variables = 18 maturities = 1M, 3M, 6M, 9M, 1Y, 2Y, ..., 30Y

• Observations = History (monthly) over 8 years = 96 bonds

◆□▶ ◆□▶ ◆三▶ ◆三▶ 三三 のへで

Last-FM - collaborative webradio



▲□▶ ▲□▶ ▲□▶ ▲□▶ ▲□▶ ▲□▶

• 28302 artists and their "tags"

 Variables = 735 tags = trance, techno, ambient, alternative, rap metal, rock, ...

900

• Observations = 2840 users

Example 3 - Face Identification



◆□▶ ◆□▶ ◆三▶ ◆三▶ 三三 のへで

◆□▶ ◆□▶ ◆三▶ ◆三▶ ◆□▶

• Variables = 256 x 256 pixels

• Observations = 64 images

- Multivariate data
- Need for interpretation
- Variability explained by combinations of the original variables

▲□▶ ▲圖▶ ▲필▶ ▲필▶ - 필.

590

- Dimension = number of variables = p
- Sample size = number of observations = n

◆□▶ ◆□▶ ◆臣▶ ◆臣▶ 臣 のへで

• Table $n \times p$ of quantitative variables

Graphical representation



 \Rightarrow Cloud of *n* points in \mathbb{R}^p

・ロト ・日・ ・ヨト ・

1

≣ ▶

590

- Reduction of the dimension
- Visualization of the cloud in 2D or 3D

◆□▶ ◆□▶ ◆臣▶ ◆臣▶ 臣 のへで

• Explaining the variability

Principal Component Analysis (PCA)

◆□▶ ◆□▶ ◆三▶ ◆三▶ ◆□▶

Rationale of PCA

 \rightarrow Project the cloud onto the "right" axis



(日)、<部)、<注)、<注)、</p>

1

590

Rationale of PCA

 \rightarrow Project the cloud onto the "right" axis



Statistical framework : data arrays

- Observations : $X_i \in \mathbb{R}^p$, $1 \leq i \leq n$
- Variable $j : X_{1j}, \ldots, X_{nj}$
- Matrix $n \times p$ of data $X = (X_1, \ldots, X_n)^T$

$$X = (X_{ij})_{i,j} = \begin{pmatrix} X_{11} & \dots & X_{1p} \\ \vdots & \ddots & \vdots \\ X_{n1} & \dots & X_{np} \end{pmatrix}$$

◆ロト ◆御ト ◆注ト ◆注ト 注 のへで

Empirical covariance matrix

• Barycenter

$$ar{X} = rac{1}{n}\sum_{i=1}^n X_i \in \mathbb{R}^p$$

• Empirical covariance matrix $(p \times p)$

$$S = (s_{kj})_{k,j} = \frac{1}{n} \sum_{i=1}^{n} X_i X_i^{T} - \bar{X} \bar{X}^{T}$$

◆□▶ ◆□▶ ◆臣▶ ◆臣▶ 臣 のへで

• Direction of the projection $a \in \mathbb{R}^p$

• Sample (1D) =
$$(a^T X_1, \ldots, a^T X_n)$$

• Maximize the empirical variance in a :

$$s_a^2 = a^T S a$$

◆□▶ ◆□▶ ◆三▶ ◆三▶ 三三 のへで

Solution (linear algebra):
eigenvector g₍₁₎ of the largest eigenvalue l₁

- Eigenvalues : $l_1 \ge \ldots \ge l_p$
- Orthonormal basis of eigenvectors $g_{(1)}, \ldots, g_{(p)}$
- Reduction of the matrix $S = GLG^T$ where
 - $L = diag(l_1, \ldots, l_p)$ diagonal matrix $p \times p$
 - G orthogonal matrix $p \times p$

$$G = (g_{(1)}, \ldots, g_{(p)}) = (g_{kj})_{k,j}$$

◆□▶ ◆□▶ ◆臣▶ ◆臣▶ 臣 のへで

• Principal components : for all vector $z \in \mathbb{R}^p$

$$y_j(z) = g_{(j)}^T(z-\bar{X}) \;, \quad 1 \leq j \leq p$$

• The matrix $n \times p$

$$Y = (y_j(X_i))_{1 \le i \le n, \ 1 \le j \le p}$$

▲ロト ▲圖ト ▲国ト ▲国ト 三国 - のへで

replaces the matrix X with the original data.

• Empirical linear correlation between variable k and the PC y_j :

$$ilde{r}_{kj} = g_{kj} \sqrt{rac{l_j}{s_{kk}}}$$
 (définition)

• Property:

$$\sum_{j=1}^p ilde{r}_{kj}^2 = 1$$

◆□▶ ◆□▶ ◆三▶ ◆三▶ 三三 のへで

• Part of the empirical variance of the *k*-th variable explained by the first 2 PC's (*y*₁, *y*₂) :

$$\tilde{r}_{k1}^2 + \tilde{r}_{k2}^2$$

• We have :

$$l_1 + l_2 = \sum_{k=1}^{p} s_{kk} (\tilde{r}_{k1}^2 + \tilde{r}_{k2}^2)$$

▲ロト ▲圖ト ▲国ト ▲国ト 三国 - のへで

• 2D Visualization : correlation disk

Correlation disk

• Point $(\tilde{r}_{k1}, \tilde{r}_{k2})$ corresponds to variable k



▲ロト ▲圖ト ▲国ト ▲国ト 三国 - のへで

• Part of the empirical variance of the cloud of points explained by the PC y_j:

$$v_j = \frac{l_j}{Tr(S)}$$

◆□▶ ◆□▶ ◆臣▶ ◆臣▶ 臣 のへで

où
$$Tr(S) = \sum_{j=1}^{p} I_j$$

• Visualization : scree-graph

Scree-graph

• Axes = index j of the PC and part of the variance v_i



Results of the PCA - Last-FM (1)

Projection of the cloud of points onto (PC1, PC2)



▲ロト ▲圖ト ▲国ト ▲国ト 三国 - のへで

Results of the PCA - Last-FM (2)

Projection of the cloud of points onto (CP3, CP4)



・ロト ・日ト ・ヨト ・ヨー うへで

Results of the PCA - Faces (1)



Data

Results of the PCA - Faces (2)

"Clean faces"



Results of the PCA - Faces (3)

Partial reconstruction (sub-column of the matrix Y)



▲□▶ ▲□▶ ▲□▶ ▲□▶ ▲□ ● ● ●

Results of the PCA - Faces (4)

Projection of other images



< □ > < □ > < □ > < □ > < □ > < □ >

590

1

A few remarks

- PCA = linear tool
- Orthogonality of principal components
- In practice :

Reduction of the matrix $R = (r_{kj})_{k,j}$ of correlations

$$r_{kj} = \frac{s_{kj}}{\sqrt{s_{kk}s_{jj}}}$$

◆□▶ ◆□▶ ◆三▶ ◆三▶ 三三 のへで

• Computational difficulties :

Reduction of S in very high dimension
• When the clouds are of the form of ellipsoïds

◆□▶ ◆□▶ ◆臣▶ ◆臣▶ 臣 のへで

- Implicit model = gaussian model
- Information carried by order statistics
- Absence of outliers

Failure of PCA



▲ロト ▲圖ト ▲国ト ▲国ト 三国 - のへで

Positive kernels

◆□▶ ◆□▶ ◆臣▶ ◆臣▶ 臣 の�?

Let \mathcal{X} be the space of the observations.

Positive kernel A function $k : \mathcal{X} \times \mathcal{X} \to \mathbb{R}$ is a positive kernel iff • k is symmetric: k(x, x') = k(x', x), $\forall x, x' \in \mathcal{X}$ • k is positive:

$$\sum_{i=1}^{n}\sum_{j=1}^{n}c_{i}c_{j}k(x_{i},x_{j})\geq 0, \quad \forall c_{i}\in\mathbb{R}, \quad \forall x_{i}\in\mathcal{X}, \quad \forall n\geq 1$$

・ロト ・四ト ・ヨト ・ヨト 三日

590

Theorem of Mercer

For all positive kernel k on \mathcal{X} there exists a Hilbert space \mathcal{H} and a mapping Φ ('feature map') such that:

$$k(x,x') = <\Phi(x), \Phi(x')>, \quad \forall x,x'\in \mathcal{X}$$

<ロト <部ト <注入 <注下 = 正

590

where <, > represents the inner product on \mathcal{H} .



 \bullet The theorem of Mercer is not constructive: it does not provide $\mathcal H,$ nor Φ

▲□▶ ▲□▶ ▲臣▶ ▲臣▶ 臣 のへで



 \bullet The theorem of Mercer is not constructive: it does not provide $\mathcal H,$ nor Φ

▲□▶ ▲圖▶ ▲目▶ ▲目▶ 目 のへで

• In practice:

 $\bullet\,$ The theorem of Mercer is not constructive: it does not provide ${\cal H},$ nor $\Phi\,$

◆□▶ ◆□▶ ◆臣▶ ◆臣▶ 臣 のへで

- In practice:
 - \mathcal{H} is a space of high (possibly infinite) dimension

 $\bullet\,$ The theorem of Mercer is not constructive: it does not provide ${\cal H},$ nor $\Phi\,$

◆□▶ ◆□▶ ◆三▶ ◆三▶ 三三 のへで

- In practice:
 - \mathcal{H} is a space of high (possibly infinite) dimension
 - Φ is a non-linear mapping

- \bullet The theorem of Mercer is not constructive: it does not provide ${\cal H},$ nor Φ
- In practice:
 - \mathcal{H} is a space of high (possibly infinite) dimension
 - Φ is a non-linear mapping
- ${\cal H}$ is a space used for representing the data, referred to as "feature space"

◆□▶ ◆舂▶ ◆注▶ ◆注▶ 三注.

590

- \bullet The theorem of Mercer is not constructive: it does not provide ${\cal H},$ nor Φ
- In practice:
 - \mathcal{H} is a space of high (possibly infinite) dimension
 - Φ is a non-linear mapping
- \mathcal{H} is a space used for representing the data, referred to as "feature space"
- The kernel trick consists in avoiding to specify ${\cal H}$ and Φ when we know they do exist!

・ロト ・ 御 ト ・ 臣 ト ・ 臣 ト … 臣

590

• Euclidean norm on \mathbb{R}^m : $\forall u \in \mathbb{R}^m$, $||u|| = \sqrt{\langle u, u \rangle}$ where \langle , \rangle inner product on \mathbb{R}^m

◆□▶ ◆□▶ ◆三▶ ◆三▶ ◆□▶

• Euclidean norm on \mathbb{R}^m : $\forall u \in \mathbb{R}^m$, $||u|| = \sqrt{\langle u, u \rangle}$ where \langle , \rangle inner product on \mathbb{R}^m

◆□▶ ◆□▶ ◆三▶ ◆三▶ 三三 のへで

• Euclidean distance: $d(u, v) = ||u - v|| = \sqrt{\langle u, u \rangle + \langle v, v \rangle - 2 \langle u, v \rangle}$

• Euclidean norm on \mathbb{R}^m : $\forall u \in \mathbb{R}^m$, $||u|| = \sqrt{\langle u, u \rangle}$ where \langle , \rangle inner product on \mathbb{R}^m

▲ロト ▲圖ト ▲国ト ▲国ト 三国 - のへで

- Euclidean distance: $d(u, v) = ||u - v|| = \sqrt{\langle u, u \rangle + \langle v, v \rangle - 2 \langle u, v \rangle}$
- Non-linear transform : Φ : $\mathbb{R}^d \to \mathbb{R}^m$ avec m > d

• Euclidean norm on \mathbb{R}^m : $\forall u \in \mathbb{R}^m$, $||u|| = \sqrt{\langle u, u \rangle}$ where \langle , \rangle inner product on \mathbb{R}^m

▲ロト ▲圖ト ▲国ト ▲国ト 三国 - のへで

- Euclidean distance: $d(u, v) = ||u - v|| = \sqrt{\langle u, u \rangle + \langle v, v \rangle - 2 \langle u, v \rangle}$
- Non-linear transform : Φ : $\mathbb{R}^d \to \mathbb{R}^m$ avec m > d

• Kernel:
$$k(x, x') = < \Phi(x), \Phi(x') >$$

- Euclidean norm on \mathbb{R}^m : $\forall u \in \mathbb{R}^m$, $||u|| = \sqrt{\langle u, u \rangle}$ where \langle , \rangle inner product on \mathbb{R}^m
- Euclidean distance: $d(u, v) = ||u - v|| = \sqrt{\langle u, u \rangle + \langle v, v \rangle - 2 \langle u, v \rangle}$
- Non-linear transform : Φ : $\mathbb{R}^d \to \mathbb{R}^m$ avec m > d

• Kernel:
$$k(x,x') = <\Phi(x), \Phi(x') >$$

Image distance:

$$d_{\Phi}(x,x') = \|\Phi(x) - \Phi(x')\| = \sqrt{k(x,x) + k(x',x') - 2k(x,x')}$$

 \Rightarrow the distance induced by Φ involves the kernel only



• No algorithmic complication when replacing the original inner product by another similarity measure

▲□▶ ▲□▶ ▲臣▶ ▲臣▶ 臣 のへで



- No algorithmic complication when replacing the original inner product by another similarity measure
- Turn a problem initially non-linear into a linear problem by sending the data to a space of higher dimension

◆□▶ ◆□▶ ◆三▶ ◆三▶ 三三 のへで



- No algorithmic complication when replacing the original inner product by another similarity measure
- Turn a problem initially non-linear into a linear problem by sending the data to a space of higher dimension

◆□▶ ◆□▶ ◆三▶ ◆三▶ 三三 のへで



- No algorithmic complication when replacing the original inner product by another similarity measure
- Turn a problem initially non-linear into a linear problem by sending the data to a space of higher dimension

Example

Let $f(x, y) = ax^2 + bx + c - y = 0$ be a quadratic surface decision (parabolic in \mathbb{R}^2).



- No algorithmic complication when replacing the original inner product by another similarity measure
- Turn a problem initially non-linear into a linear problem by sending the data to a space of higher dimension

Example

Let $f(x, y) = ax^2 + bx + c - y = 0$ be a quadratic surface decision (parabolic in \mathbb{R}^2).

Key role of the transformation:

$$\begin{array}{cccc} \flat & : & \mathbb{R}^2 & \rightarrow & \mathbb{R}^4 \\ & x & \mapsto & \left(x^2, x, 1, y\right)^T \end{array}$$

◆□▶ ◆□▶ ◆三▶ ◆三▶ 三三 のへで

Example (continued)

One may write:

$$g(x^2, x, 1, y) = ax^2 + bx + c - y = 0$$

où g(u, v, w, y) = au + bv + cw - y.

The equation g(u, v, w, y) = 0 defines a decision surface that is linear in \mathbb{R}^4 .

◆□▶ ◆□▶ ◆三▶ ◆三▶ 三三 のへで

Example (continued)

One may write:

$$g(x^2, x, 1, y) = ax^2 + bx + c - y = 0$$

où g(u, v, w, y) = au + bv + cw - y.

The equation g(u, v, w, y) = 0 defines a decision surface that is linear in \mathbb{R}^4 .

A non-linear problem in a certain space can be formulated sometimes as a linear problem in a 'larger' space.

900

From non-linear to linear



Input Space

Feature Space

▲□▶ ▲圖▶ ▲필▶ ▲필▶ 三월

990

Kernel PCA

◆□▶ ◆□▶ ◆臣▶ ◆臣▶ 臣 の�?



Goals of PCA

- Method for visualizing the data
- Effective reduction of the dimension of the data

◆□▶ ◆□▶ ◆三▶ ◆三▶ 三三 のへで

Goals of PCA

- Method for visualizing the data
- Effective reduction of the dimension of the data

 PCA consists in identifying the principal components of the data formed by

(日) (문) (문) (문) (문)

• the best direction for projecting the cloud of points i.e. that with maximal variance

Goals of PCA

- Method for visualizing the data
- Effective reduction of the dimension of the data

PCA consists in identifying the principal components of the data formed by

- the best direction for projecting the cloud of points i.e. that with maximal variance
- 2 next, the best direction for projecting that is orthogonal to the first

Goals of PCA

- Method for visualizing the data
- Effective reduction of the dimension of the data

PCA consists in identifying the principal components of the data formed by

- the best direction for projecting the cloud of points i.e. that with maximal variance
- 2 next, the best direction for projecting that is orthogonal to the first
- 3 and, so on so forth, until the *n*-th



▲□▶ ▲□▶ ▲臣▶ ▲臣▶ 臣 のへで

• Orthogonal projection of a vector x onto direction $w \in \mathbb{R}^d$:

$$p_w(x) = \frac{\langle x, w \rangle}{\|w\|}$$

▲□▶ ▲□▶ ▲臣▶ ▲臣▶ 臣 のへで

• Orthogonal projection of a vector x onto direction $w \in \mathbb{R}^d$:

$$p_w(x) = \frac{\langle x, w \rangle}{\|w\|}$$

• Empirical variance of the cloud of points along direction w:

$$\mathbb{V}(p_w) = \frac{1}{n} \sum_{i=1}^n \frac{\langle x_i, w \rangle^2}{\|w\|^2}$$

◆□▶ ◆□▶ ◆臣▶ ◆臣▶ 臣 のへで

• Orthogonal projection of a vector x onto direction $w \in \mathbb{R}^d$:

$$p_w(x) = \frac{\langle x, w \rangle}{\|w\|}$$

• Empirical variance of the cloud of points along direction w:

$$\mathbb{V}(p_w) = \frac{1}{n} \sum_{i=1}^n \frac{\langle x_i, w \rangle^2}{\|w\|^2}$$

• Empirical covariance matrix $\Sigma = \frac{1}{n} \sum_{i=1}^{n} x_i x_i^T$

• Orthogonal projection of a vector x onto direction $w \in \mathbb{R}^d$:

$$p_w(x) = \frac{\langle x, w \rangle}{\|w\|}$$

• Empirical variance of the cloud of points along direction w:

$$\mathbb{V}(p_w) = \frac{1}{n} \sum_{i=1}^n \frac{\langle x_i, w \rangle^2}{\|w\|^2}$$

- Empirical covariance matrix $\Sigma = \frac{1}{n} \sum_{i=1}^{n} x_i x_i^{T}$
- We then have :

$$\mathbb{V}(p_w) = \frac{w^T \Sigma w}{\|w\|^2}$$

▲ロト ▲圖ト ▲国ト ▲国ト 三国 - のへで

First principal component

$$\underset{w}{\operatorname{arg\,max}} \mathbb{V}(p_w) = \frac{w^T \Sigma w}{\|w\|^2}$$

< □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ = □

500
First principal component

$$\operatorname*{arg\,max}_{w} \mathbb{V}(p_w) = \frac{w^T \Sigma w}{\|w\|^2}$$

Solution

The principal components are the eigenvectors of Σ sorted by decreasing order of magnitude of the corresponding eigenvalues.

First principal component

$$\operatorname*{arg\,max}_{w} \mathbb{V}(p_w) = \frac{w^T \Sigma w}{\|w\|^2}$$

Solution

The principal components are the eigenvectors of Σ sorted by decreasing order of magnitude of the corresponding eigenvalues.

Remark : the matrix Σ is PSD, hence diagonalizable in an orthonormal basis.

・ロト ・(部)・ ・(日)・ 「日

PCA (continued)

One searches for a vector v and a real number λ such that:

$$\Sigma v = \lambda v$$

Observe that we have:

$$\Sigma v = \frac{1}{n} \sum_{i=1}^{n} \langle x_i, v \rangle \ x_i$$

Thus:

$$v = \sum_{i=1}^{n} \left(\frac{\langle x_i, v \rangle}{n\lambda} \right) x_i = \sum_{i=1}^{n} \alpha_i x_i$$

PCA (continued)

One searches for a vector v and a real number λ such that:

$$\Sigma v = \lambda v$$

Observe that we have:

$$\Sigma v = \frac{1}{n} \sum_{i=1}^{n} \langle x_i, v \rangle \ x_i$$

Thus:

$$v = \sum_{i=1}^{n} \left(\frac{\langle x_i, v \rangle}{n\lambda}\right) x_i = \sum_{i=1}^{n} \alpha_i x_i$$

One uses

$$x_j^T \Sigma v = \lambda < x_j, v >, \quad \forall j$$

and one substitutes the expressions for Σ and v:

$$\frac{1}{n}\sum_{i=1}^{n}\alpha_{i}\left\langle x_{j},\sum_{k=1}^{n}< x_{k},x_{i}>x_{k}\right\rangle =\lambda\sum_{\substack{i=1\\ \ <\ \$$

• Denote by $K = (\langle x_i, x_j \rangle)_{i,j}$ the Gram matrix

- Denote by $K = (\langle x_i, x_j \rangle)_{i,j}$ the Gram matrix
- One may then write the system:

$$K^2 \alpha = n \lambda K \alpha$$

- Denote by $K = (\langle x_i, x_j \rangle)_{i,j}$ the Gram matrix
- One may then write the system:

$$K^2 \alpha = n\lambda K \alpha$$

 $\bullet\,$ To find $\alpha,$ one thus solves the problem

$$K\alpha = n\lambda\alpha$$

◆□▶ ◆□▶ ◆臣▶ ◆臣▶ 臣 のへで

• it is essentially tailored to the case of multivariate Gaussian data

- it is essentially tailored to the case of multivariate Gaussian data
 - in general, non-correlation does not imply independence of the principal directions

◆□▶ ◆□▶ ◆三▶ ◆三▶ 三三 のへで

- it is essentially tailored to the case of multivariate Gaussian data
 - in general, non-correlation does not imply independence of the principal directions

◆□▶ ◆□▶ ◆三▶ ◆三▶ 三三 のへで

alternative : Independent Component Analysis

- it is essentially tailored to the case of multivariate Gaussian data
 - in general, non-correlation does not imply independence of the principal directions

▲ロト ▲圖ト ▲国ト ▲国ト 三国 - のへで

- alternative : Independent Component Analysis
- it is tailored to detect linear structures

- it is essentially tailored to the case of multivariate Gaussian data
 - in general, non-correlation does not imply independence of the principal directions

- alternative : Independent Component Analysis
- it is tailored to detect linear structures
 - Not all clouds of points are ellipsoïds!!

- it is essentially tailored to the case of multivariate Gaussian data
 - in general, non-correlation does not imply independence of the principal directions

◆□▶ ◆舂▶ ◆逹▶ ◆逹▶ 三臣 ……

- alternative : Independent Component Analysis
- it is tailored to detect linear structures
 - Not all clouds of points are ellipsoïds!!
 - alternative : Kernel PCA

• Apply a transformation Φ that sends the cloud of points X to a space where the structure is linear

◆ロト ◆御ト ◆注ト ◆注ト 注 のへで

- Apply a transformation Φ that sends the cloud of points X to a space where the structure is linear
- The covariance matrix of $\Phi(X) = (\Phi(x_1), \dots, \Phi(x_n))^T$ is then

$$\Sigma = \frac{1}{n} \sum_{i=1}^{n} \Phi(x_i) \Phi(x_i)^{T}$$

◆□▶ ◆□▶ ◆臣▶ ◆臣▶ 臣 のへで

- Apply a transformation Φ that sends the cloud of points X to a space where the structure is linear
- The covariance matrix of $\Phi(X) = (\Phi(x_1), \dots, \Phi(x_n))^T$ is then

$$\Sigma = \frac{1}{n} \sum_{i=1}^{n} \Phi(x_i) \Phi(x_i)^{T}$$

◆□▶ ◆□▶ ◆三▶ ◆三▶ 三三 のへで

• Kernel trick : $K = (k(x_i, x_j))_{i,j} = (\Phi(x_i)^T \Phi(x_j))_{i,j}$

Kernel PCA (continued)

• Principal "Directions" of the form:

$$p_i(x) = \sum_{j=1}^n \alpha_{i,j} k(x_j, x)$$

Kernel PCA (continued)

• Principal "Directions" of the form:

$$p_i(x) = \sum_{j=1}^n \alpha_{i,j} k(x_j, x)$$

the vector α_i = (α_{i,1},..., α_{i,n}) is solution to the optimisation problem:

$$\max_{\alpha} \frac{\alpha^{\mathsf{T}} \mathsf{K}^2 \alpha}{\alpha^{\mathsf{T}} \mathsf{K} \alpha}$$

◆□▶ ◆□▶ ◆臣▶ ◆臣▶ 臣 のへで

under the constraints: $\alpha_i^T K \alpha_j$ pour $j = 1, \ldots, i-1$

Kernel PCA (continued)

• Principal "Directions" of the form:

$$p_i(x) = \sum_{j=1}^n \alpha_{i,j} k(x_j, x)$$

the vector α_i = (α_{i,1},..., α_{i,n}) is solution to the optimisation problem:

$$\max_{\alpha} \frac{\alpha^{\mathsf{T}} \mathsf{K}^2 \alpha}{\alpha^{\mathsf{T}} \mathsf{K} \alpha}$$

under the constraints: $\alpha_i^T K \alpha_j$ pour $j = 1, \ldots, i-1$

• one solves the problem:

$$K\alpha = n\lambda\alpha$$

▲ロト ▲圖ト ▲国ト ▲国ト 三国 - のへで

Independent Component Analysis (ICA)

◆□▶ ◆□▶ ◆三▶ ◆三▶ 三三 のへで

The "cocktail-party" problem



<ロト <(四)> <(0)> <(0)> <(0)> <(0)> <(0)> <(0)> <(0)> <(0)> <(0)> <(0)> <(0)> <(0)> <(0)> <(0)> <(0)> <(0)> <(0)> <(0)> <(0)> <(0)> <(0)> <(0)> <(0)> <(0)> <(0)> <(0)> <(0)> <(0)> <(0)> <(0)> <(0)> <(0)> <(0)> <(0)> <(0)> <(0)> <(0)> <(0)> <(0)> <(0)> <(0)> <(0)> <(0)> <(0)> <(0)> <(0)> <(0)> <(0)> <(0)> <(0)> <(0)> <(0)> <(0)> <(0)> <(0)> <(0)> <(0)> <(0)> <(0)> <(0)> <(0)> <(0)> <(0)> <(0)> <(0)> <(0)> <(0)> <(0)> <(0)> <(0)> <(0)> <(0)> <(0)> <(0)> <(0)> <(0)> <(0)> <(0)> <(0)> <(0)> <(0)> <(0)> <(0)> <(0)> <(0)> <(0)> <(0)> <(0)> <(0)> <(0)> <(0)> <(0)> <(0)> <(0)> <(0)> <(0)> <(0)> <(0)> <(0)> <(0)> <(0)> <(0)> <(0)> <(0)> <(0)> <(0)> <(0)> <(0)> <(0)> <(0)> <(0)> <(0)> <(0)> <(0)> <(0)> <(0)> <(0)> <(0)> <(0)> <(0)> <(0)> <(0)> <(0)> <(0)> <(0)> <(0)> <(0)> <(0)> <(0)> <(0)> <(0)> <(0)> <(0)> <(0)> <(0)> <(0)> <(0)> <(0)> <(0)> <(0)> <(0)> <(0)> <(0)> <(0)> <(0)> <(0)> <(0)> <(0)> <(0)> <(0)> <(0)> <(0)> <(0)> <(0)> <(0)> <(0)> <(0)> <(0)> <(0)> <(0)> <(0)> <(0)> <(0)> <(0)> <(0)> <(0)> <(0)> <(0)> <(0)> <(0)> <(0)> <(0)> <(0)> <(0)> <(0)> <(0)> <(0)> <(0)> <(0)> <(0)> <(0)> <(0)> <(0)> <(0)> <(0)> <(0)> <(0)> <(0)> <(0)> <(0)> <(0)> <(0)> <(0)> <(0)> <(0)> <(0)> <(0)> <(0)> <(0)> <(0)> <(0)> <(0)> <(0)> <(0)> <(0)> <(0)> <(0)> <(0)> <(0)> <(0)> <(0)> <(0)> <(0)> <(0)> <(0)> <(0)> <(0)> <(0)> <(0)> <(0)> <(0)> <(0)> <(0)> <(0)> <(0)> <(0)> <(0)> <(0)> <(0)> <(0)> <(0)> <(0)> <(0)> <(0)> <(0)> <(0)> <(0)> <(0)> <(0)> <(0)> <(0)> <(0)> <(0)> <(0)> <(0)> <(0)> <(0)> <(0)> <(0)> <(0)> <(0)> <(0)> <(0)> <(0)> <(0)> <(0)> <(0)> <(0)> <(0)> <(0)> <(0)> <(0)> <(0)> <(0)> <(0)> <(0)> <(0)> <(0)> <(0)> <(0)> <(0)> <(0)> <(0)> <(0)> <(0)> <(0)> <(0)> <(0)> <(0)> <(0)> <(0)> <(0)> <(0)> <(0)> <(0)> <(0)> <(0)> <(0)> <(0)> <(0)> <(0)> <(0)> <(0)> <(0)> <(0)> <(0)> <(0)> <(0)> <(0)> <(0)> <(0)> <(0)> <(0)> <(0)> <(0)> <(0)> <(0)> <(0)> <(0)> <(0)> <(0)> <(0)> <(0)> <(0)> <(0)> <(0)> <(0)> <(0)> <(0)> <(0)> <(0)> <(0)> <(0)> <(0)> <(0)> <(0)> <(0)> <(0)> <(0)> <(0)>

- PCA = based on the notion of correlation
- Appropriate notion = notion of independence

- Notice that : X and Y independent $\Rightarrow cov(X, Y) = 0$
- Reciprocal false in general, except in the Gaussian case..

◆□▶ ◆□▶ ◆臣▶ ◆臣▶ 臣 のへで

• From PCA to ICA... (much more difficult !)

- $S = (S_1, \dots, S_d)^T$ unknown independent and non-Gaussian sources
- A mixing matrix $d \times d$ unknown
- $X = (X_1, \dots, X_d)^T$ observations (sensors), one assumes $\operatorname{Cov}(X) = \mathbf{I}$
- One has the system : $X = \mathbf{A}S$
- One searches for **A** orthogonal such that:

 $S = \mathbf{A}^T X$ has independent components

◆□▶ ◆□▶ ◆三▶ ◆三▶ 三三 のへで

Information theory

• Entropy of a r.v.
$$Z \sim p(z)$$
 :

$$H(Z) = -\mathbb{E}(\log(p(Z)))$$

• Consider the r.v. T with variance v, then

$$Z \sim \mathcal{N}(0,1) \quad o \quad \max_{T} H(T)$$

• Mutual information for $S = (S_1, \ldots, S_d)^T$:

$$I(S) = \sum_{i=1}^{d} H(S_i) - H(S)$$

▲□▶ ▲圖▶ ▲目▶ ▲目▶ 目 のへで

• Property of the entropy : if $S = \mathbf{A}^{\mathsf{T}} X$

$$H(S) = H(X) + \log(|\det(\mathbf{A})|)$$

• One then has the following optimization problem:

$$\rightarrow \min_{\mathbf{A}:\mathbf{A}^{\mathsf{T}}\mathbf{A}=\mathbf{I}} I(\mathbf{A}^{\mathsf{T}}X) = \sum_{i=1}^{d} H(S_i) - H(X)$$

▲ロト ▲圖ト ▲国ト ▲国ト 三国 - のへで

• Interpretation : deviation from the Gaussian behavior (minimization of the entropy of the components)