

Machine learning 2

Dimension reduction

Pierre Chainais



1 Dimension reduction

- Motivations
- Basic ideas
- Principal Component Analysis (PCA)
- Fisher discriminant analysis (FLD or FDA)

Motivations

Data in large dimension

p.3



visages

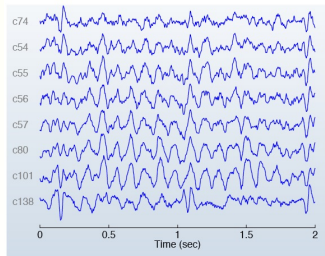


genetics

Zambian President Levy Mwanawasa has won a second term in office in an election his challenger Michael Sata accused him of rigging, official results showed on Monday.

According to media reports, a pair of hackers said on Saturday that the Firefox Web browser, commonly perceived as the safer and more customizable alternative to market leader Internet Explorer, is critically flawed. A presentation on the flaw was shown during the ToorCon hacker conference in San Diego.

documents



electroencephalograms (ECG)

Why do we need dimension reduction ?

(of the features vector)

- ▶ Compression / compact representation
- ▶ Statistics / efficiency : simpler \Rightarrow + robust
- ▶ Visualization : 2D, or 3D maximum...
- ▶ Detection of anomalies : normal average / extremes

Methods for the reduction of dimension :

- ▶ Projection on characteristic directions
- ▶ Clustering : cf. vectorial quantization, unsupervised classif. (later)
- ▶ Selection of variables (stat. tests, regularization...)
- ▶ Non linear method (kernels...)

Why do we need dimension reduction ?

(of the features vector)

- ▶ Compression / compact representation
- ▶ Statistics / efficiency : simpler \Rightarrow + robust
- ▶ Visualization : 2D, or 3D maximum...
- ▶ Detection of anomalies : normal average / extremes

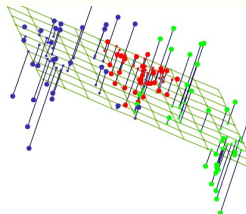
Methods for the reduction of dimension :

- ▶ **Projection on characteristic directions**
- ▶ Clustering : cf. vectorial quantization, unsupervised classif. (later)
- ▶ Selection of variables (stat. tests, regularization...)
- ▶ Non linear method (kernels...)

- ▶ **Best possible prediction** (classification or regression)
application : minimize error rate (pragmatic)
- ▶ **Discover structure**
application : interpretable characteristics, visualization
- ▶ **Estimation of density** $p(\mathbf{x})$, model the data,
applications : detection of anomalies, models of language...



face = image 19×19 , that is $\mathbf{x} \in \mathbb{R}^{361}$, $D = 361$



$$\mathbf{x} \in \mathbb{R}^{361} \Rightarrow \mathbf{z} \in \mathbb{R}^{10} ?$$

Idea : search for a projection $\mathbf{z} = \mathbf{U}^T \mathbf{x}$, $\mathbf{U} : D \times M$, $M \ll D$

Let N data points in dimension D : $\mathbf{x}_1, \dots, \mathbf{x}_N \in \mathbb{R}^D$

$$X = \begin{pmatrix} \mathbf{x}_1^T \\ \vdots \\ \mathbf{x}_N^T \end{pmatrix} \in \mathbb{R}^{N \times D}$$

One wants to reduce the dimension from D to M by choosing M orthogonal directions $\mathbf{u}_1, \dots, \mathbf{u}_M \in \mathbb{R}^D$:

$$U = \begin{pmatrix} | & & | \\ \mathbf{u}_1 & \dots & \mathbf{u}_M \\ | & & | \end{pmatrix} \in \mathbb{R}^{D \times M}$$

Projection of \mathbf{x} sur $\mathbf{z} = (z_1, \dots, z_M)^T = U^T \mathbf{x}$ with $U^T U = I$

\Rightarrow How to determine U ?

PCA fulfills both optimization criteria simultaneously :

- 1 **best least square approximation** of some data set $(\mathbf{x}_n)_{1 \leq n \leq N}$ par $M < D$ **orthogonal components** denoted by \mathbf{u}_j :

- Decomposition : $\mathbf{z} = U^T \mathbf{x}$, $z_j = \mathbf{u}_j^T \mathbf{x}$, $1 \leq j \leq M$
- Reconstruction : $\mathbf{x}^{app} = U\mathbf{z} = \sum_{j=1}^M z_j \mathbf{u}_j$

$$U = \operatorname{argmin}_{U^T U = I} \sum_{n=1}^N \|\mathbf{x}_n - \underbrace{UU^T \mathbf{x}_n}_{\mathbf{x}_n^{app}}\|^2$$

- 2 **projection on components z_j with maximal variances**

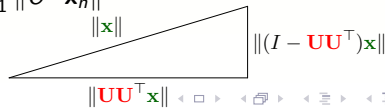
Intuition : large dispersion \Rightarrow significant

For centered inputs \mathbf{x} ,

$$U = \operatorname{argmax}_{U^T U = I} \|\underbrace{U^T \mathbf{x}}_{\mathbf{z}}\|^2 =$$

$$\operatorname{argmax}_{U^T U = I} \frac{1}{N} \sum_{n=1}^N \|U^T \mathbf{x}_n\|^2$$

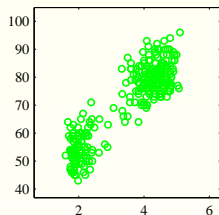
équivalence 1 \Leftrightarrow 2 :



Principal Component Analysis (PCA)

p.10

Identification of the 1st principal component \mathbf{u}_1



$$\mathbf{u}_1 = \operatorname{argmax}_{\|\mathbf{u}\|=1} \frac{1}{N} \sum_{n=1}^N \|\mathbf{u}_1^T \mathbf{x}_n\|^2 \quad (1)$$

$$= \operatorname{argmax}_{\|\mathbf{u}\|=1} \frac{1}{N} \|\mathbf{X} \mathbf{u}_1\|^2 \quad (2)$$

$$= \operatorname{argmax}_{\|\mathbf{u}\|=1} \mathbf{u}_1^T \left(\frac{1}{N} \mathbf{X}^T \mathbf{X} \right) \mathbf{u}_1 \quad (3)$$

$$= \operatorname{eigenvector}(\max \lambda_{\mathbf{X}^T \mathbf{X}}) \quad (4)$$

$$\mathbf{X} = \begin{pmatrix} \mathbf{x}_1^T \\ \vdots \\ \mathbf{x}_N^T \end{pmatrix} \in \mathbb{R}^{N \times D}$$

\mathbf{u}_1 = 1st eigenvector of $\mathbf{X}^T \mathbf{X}$.

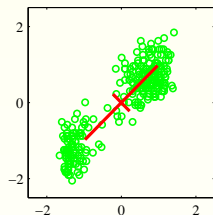
\mathbf{u}_m = m-th eigenvector $\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_M$.

Remark : $\frac{1}{N} \mathbf{X}^T \mathbf{X}$ is a covariance matrix.

Principal Component Analysis (PCA)

p.11

Identification of the 1st principal component \mathbf{u}_1



$$\mathbf{u}_1 = \operatorname{argmax}_{\|\mathbf{u}\|=1} \frac{1}{N} \sum_{n=1}^N \|\mathbf{u}_1^T \mathbf{x}_n\|^2 \quad (1)$$

$$= \operatorname{argmax}_{\|\mathbf{u}\|=1} \frac{1}{N} \|\mathbf{X} \mathbf{u}_1\|^2 \quad (2)$$

$$= \operatorname{argmax}_{\|\mathbf{u}\|=1} \mathbf{u}_1^T \left(\frac{1}{N} \mathbf{X}^T \mathbf{X} \right) \mathbf{u}_1 \quad (3)$$

$$= \operatorname{eigenvector}(\max \lambda_{\mathbf{X}^T \mathbf{X}}) \quad (4)$$

$$\mathbf{X} = \begin{pmatrix} \mathbf{x}_1^T \\ \vdots \\ \mathbf{x}_N^T \end{pmatrix} \in \mathbb{R}^{N \times D}$$

\mathbf{u}_1 = 1st eigenvector of $\mathbf{X}^T \mathbf{X}$.

\mathbf{u}_m = m-th eigenvector $\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_M$.

Remark : $\frac{1}{N} \mathbf{X}^T \mathbf{X}$ is a covariance matrix.

Principal Component Analysis (PCA)

p.12

How to compute a PCA?

Method 1 : decomposition in eigen values and vectors

$\mathbf{U} = (\mathbf{u}_j)_{1 \leq j \leq M} = M$ first eigen vectors of the covariance matrix

$$\mathbf{C} = \frac{1}{N} \mathbf{X}^T \mathbf{X}.$$

Cost is high : $O(ND^2)$. Remark : Karhunen-Loève in sig. proc.

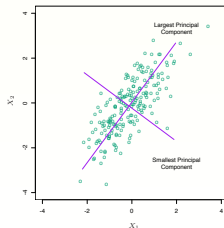
Method 2 : singular value decomposition **SVD**

$$\mathbf{X} = \underbrace{\mathbf{V}}_{N \times N} \underbrace{\mathbf{\Sigma}}_{N \times D} \underbrace{\mathbf{U}^T}_{D \times D}, \text{ where } \mathbf{\Sigma} \text{ contains a diagonal block } D \times D.$$

Cost of the M first eigenvectors : $O(NDM)$.

Connection :

$$\mathbf{C} = \mathbf{U} \mathbf{\Sigma}^2 \mathbf{U}^T$$



Principal Component Analysis (PCA)

Example : face recognition

- ▶ D = number of pixels per image, for instance $19 \times 19 = 361$,
- ▶ $\mathbf{x}_n \in \mathbb{R}^D$ is an image of a face,
- ▶ x_{ni} = intensity of the i -th pixel of image n ,

$$\begin{array}{ccc} X_{D \times N}^T & \simeq & U_{D \times M} \quad \times \quad Z_{M \times N} \\ \begin{array}{c} \text{img 1} \quad \dots \quad \text{img N} \end{array} & \simeq & \begin{array}{c} \text{img 1} \quad \text{img 2} \quad \text{img 3} \quad \text{img 4} \quad \text{img 5} \end{array} \quad \times \quad \left(\begin{array}{c|c|c} | & & | \\ \mathbf{z}_1 & \dots & \mathbf{z}_N \\ | & & | \end{array} \right) \end{array}$$

Interest :

- ▶ extraction of generical characteristics,
- ▶ use the \mathbf{z}_j for classification (K-NN,...),
- ▶ reduction of dimension \Rightarrow speed

[Turk & Pentland 1991]

Principal Component Analysis (PCA)

Example : semantic document analysis

- ▶ D = number of words in the vocabulary,
- ▶ $\mathbf{x}_n \in \mathbb{R}^D$ counts the frequency of words in document n ,
- ▶ x_{ni} = frequency of the i -th word in document n ,

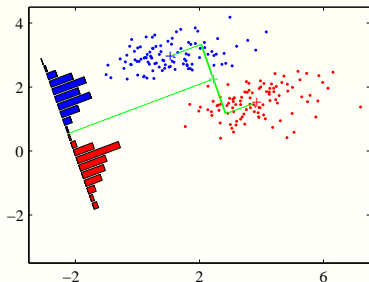
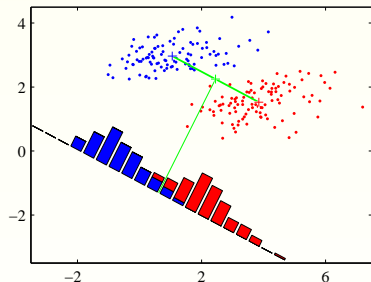
$$\begin{pmatrix} \text{stocks : } 2 \cdots 0 \\ \text{exchange : } 4 \cdots 1 \\ \text{the : } 10 \cdots 12 \\ \vdots \\ \text{wins : } 0 \cdots 3 \\ \text{game : } 1 \cdots 3 \\ \text{Interest : } \end{pmatrix} \begin{matrix} \overset{X_{D \times N}^T}{\simeq} \\ \underset{\simeq}{\simeq} \end{matrix} \begin{pmatrix} 0.4 & \cdots & -0.001 \\ 0.8 & \cdots & 0.03 \\ 0.01 & \cdots & 0.04 \\ \vdots & \cdots & \vdots \\ 0.002 & \cdots & 2.3 \\ 0.003 & \cdots & 1.9 \end{pmatrix} \begin{matrix} \times \\ \times \end{matrix} \begin{matrix} Z_{M \times N} \\ \left(\begin{array}{c|c|c} | & & | \\ \mathbf{z}_1 & \cdots & \mathbf{z}_N \\ | & & | \end{array} \right) \end{matrix}$$

- ▶ measure of similarity between documents :
 $\mathbf{z}_1^T \mathbf{z}_2$ in place of $\mathbf{x}_1^T \mathbf{x}_2$,
- ▶ searching for information ($\mathbf{z}_j \simeq$ themes... ?)

Fisher discriminant analysis (FLD or FDA)

p.15

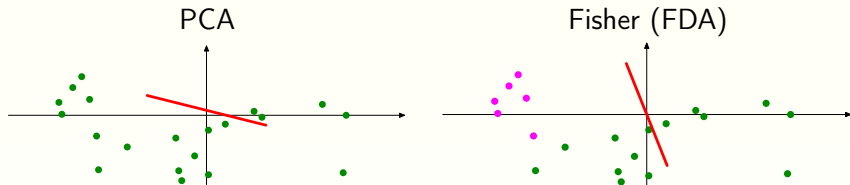
Objective : the most discriminant supervised projection



Fisher discriminant analysis (FLD or FDA)

p.16

Objective : the most discriminant supervised projection

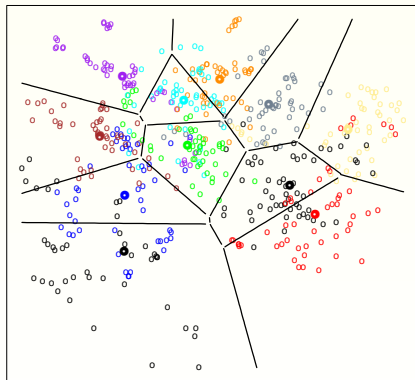


Fisher discriminant analysis (FLD or FDA)

p.17

Example : vowel recognition by LDA

11 vowels $\Rightarrow K = 11$ classes described by $D = 10$ features
(cf. time frequency analysis)



Classification based on $M=2$ Fisher discriminant components

Fisher discriminant analysis (FLD or FDA)

p.18

Supervised dimension reduction

Idea : maximize the projected ratio $\frac{\text{variance inter-class}}{\text{variance intra-class}}$

K= 2 classes :

projection $y = \mathbf{w}^T \mathbf{x}$ on \mathbf{w} such that $J(\mathbf{w}) = \frac{(\mu_2 - \mu_1)^2}{s_1^2 + s_2^2}$ maximum

$$\begin{cases} \mu_k &= \frac{1}{N_k} \sum_{n \in \mathcal{C}_k} y_n = \mathbf{w}^T \mathbf{m}_k \\ s_k^2 &= \sum_{n \in \mathcal{C}_k} (y_n - \mu_k)^2 \end{cases}$$

$$J(\mathbf{w}) = \frac{(\mu_2 - \mu_1)^2}{s_1^2 + s_2^2}$$

$$\mathbf{m}_k = \frac{1}{N_k} \sum_{n \in \mathcal{C}_k} \mathbf{x}_n$$

\Leftrightarrow

$$S_{inter} = (\mathbf{m}_2 - \mathbf{m}_1)(\mathbf{m}_2 - \mathbf{m}_1)^T$$

$$J(\mathbf{w}) = \frac{\mathbf{w}^T S_{inter} \mathbf{w}}{\mathbf{w}^T S_{intra} \mathbf{w}}$$

$$\begin{aligned} S_{intra} &= \sum_{n \in \mathcal{C}_1} (\mathbf{x}_n - \mathbf{m}_1)(\mathbf{x}_n - \mathbf{m}_1)^T \\ &\quad + \sum_{n \in \mathcal{C}_2} (\mathbf{x}_n - \mathbf{m}_2)(\mathbf{x}_n - \mathbf{m}_2)^T \end{aligned}$$

Fisher discriminant analysis (FLD or FDA)

p.19

Supervised dimension reduction

Idea : maximize the projected ratio $\frac{\text{variance inter-class}}{\text{variance intra-class}}$

K= 2 classes :

projection $y = \mathbf{w}^T \mathbf{x}$ on \mathbf{w} such that $J(\mathbf{w}) = \frac{(\mu_2 - \mu_1)^2}{s_1^2 + s_2^2}$ maximum

$$\begin{cases} \mu_k &= \frac{1}{N_k} \sum_{n \in \mathcal{C}_k} y_n = \mathbf{w}^T \mathbf{m}_k \\ s_k^2 &= \sum_{n \in \mathcal{C}_k} (y_n - \mu_k)^2 \end{cases}$$

$$J(\mathbf{w}) = \frac{(\mu_2 - \mu_1)^2}{s_1^2 + s_2^2}$$

$$\mathbf{m}_k = \frac{1}{N_k} \sum_{n \in \mathcal{C}_k} \mathbf{x}_n$$

\Leftrightarrow

$$S_{inter} = (\mathbf{m}_2 - \mathbf{m}_1)(\mathbf{m}_2 - \mathbf{m}_1)^T$$

$$J(\mathbf{w}) = \frac{\mathbf{w}^T S_{inter} \mathbf{w}}{\mathbf{w}^T S_{intra} \mathbf{w}}$$

$$\begin{aligned} S_{intra} &= \sum_{n \in \mathcal{C}_1} (\mathbf{x}_n - \mathbf{m}_1)(\mathbf{x}_n - \mathbf{m}_1)^T \\ &\quad + \sum_{n \in \mathcal{C}_2} (\mathbf{x}_n - \mathbf{m}_2)(\mathbf{x}_n - \mathbf{m}_2)^T \end{aligned}$$

Fisher discriminant analysis (FLD or FDA)

Supervised dimension reduction

$K \geq 2$ classes :

projection on M directions w_j : $y_j = \mathbf{w}_j^T \mathbf{x}$

$$\mathbf{y} = W^T \mathbf{x}$$

$$\mathbf{m}_k = \frac{1}{N_k} \sum_{n \in \mathcal{C}_k} \mathbf{x}_n$$

$$J(W) = \text{tr}[(WS_{intra}W^T)^{-1}(WS_{inter}W^T)]$$

$$\mathbf{m} = \frac{1}{N} \sum_{n=1}^N \mathbf{x}_n$$



$$S_k = \sum_{n \in \mathcal{C}_k} (\mathbf{x}_n - \mathbf{m}_k)(\mathbf{x}_n - \mathbf{m}_k)^T$$

$$S_{intra} = \sum_{k=1}^K S_k$$

$$S_{inter} = \sum_{k=1}^K N_k (\mathbf{m}_k - \mathbf{m})(\mathbf{m}_k - \mathbf{m})^T$$

Solution

$(\mathbf{w}_j)_{1 \leq j \leq M}$: M first eigenvectors
of $S_{intra}^{-1} S_{inter}$
[Fukunaga 1990]

- ▶ Fisher : dimension reduction + preparing classification
- ▶ LDA / QDA adapted : projection z_j = sum of independent r.v.

Central Limit Theorem : $z_j \simeq$ Gaussian !

Framework : $\mathbf{z} = U^T \mathbf{x}$, $\mathbf{x} \simeq U \mathbf{z}$

- ▶ **PCA** : maximize the variance of projected components,
- ▶ **FDA** : maximize the projected ratio $\frac{\text{variance inter-class}}{\text{variance intra-class}}$
- ▶ **CCA** : Canonical Correlation Analysis...
- ▶ possibility of random projections (compressed sensing !)...

Algorithm : eigenvalue decomposition