

Statistiques en grande dimension

Christophe Giraud^{1,2} et Tristan Mary-Huart^{3,4}

- (1) Université Paris-Sud
- (2) Ecole Polytechnique
- (3) AgroParistech
- (4) INRA - Le Moulon

M2 MathSV & Maths Aléa

High-dimensional data

Données en grande dimension

- **Données biotech:** mesure des milliers de quantités par "individu".
- **Images :** images médicales, astrophysique, video surveillance, etc. Chaque image est constituées de milliers ou millions de pixels ou voxels.
- **Marketing:** les sites web et les programmes de fidélité collectent de grandes quantités d'information sur les préférences et comportements des clients. Ex: systèmes de recommandation...
- **Business:** exploitation des données internes et externes de l'entreprise devient primordial
- **Crowdsourcing data :** données récoltées online par des volontaires. Ex: eBirds collecte des millions d'observations d'oiseaux en Amérique du Nord

Blessing?

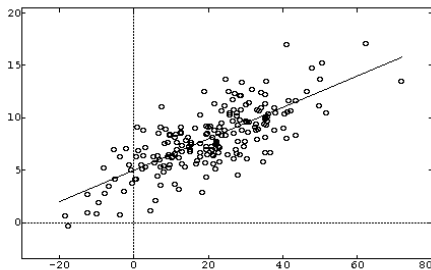
😊 we can sense thousands of variables on each "individual" : potentially we will be able to scan every variables that may influence the phenomenon under study.

😞 the curse of dimensionality : separating the signal from the noise is in general almost impossible in high-dimensional data and computations can rapidly exceed the available resources.

Renversement de point de vue

Cadre statistique classique:

- petit nombre p de paramètres
- grand nombre n d'expériences
- on étudie le comportement asymptotique des estimateurs lorsque $n \rightarrow \infty$ (résultats type théorème central limite)



Renversement de point de vue

Cadre statistique classique:

- petit nombre p de paramètres
- grand nombre n d'expériences
- on étudie le comportement asymptotique des estimateurs lorsque $n \rightarrow \infty$ (résultats type théorème central limite)

Données actuelles:

- inflation du nombre p de paramètres
- taille d'échantillon réduite: $n \approx p$ ou $n \ll p$

\implies penser différemment les statistiques!
(penser $n \rightarrow \infty$ ne convient plus)

Fléau de la dimension

Course 1 : fluctuations cumulate

Exemple : linear regression $Y = \mathbf{X}\beta^* + \varepsilon$ with $\mathbf{cov}(\varepsilon) = \sigma^2 I_n$. The Least-Square estimator $\hat{\beta} \in \operatorname{argmin}_{\beta \in \mathbb{R}^p} \|Y - \mathbf{X}\beta\|^2$ has a risk

$$\mathbb{E} \left[\|\hat{\beta} - \beta^*\|^2 \right] = \operatorname{Tr} \left((\mathbf{X}^T \mathbf{X})^{-1} \right) \sigma^2.$$

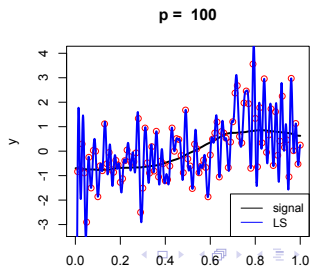
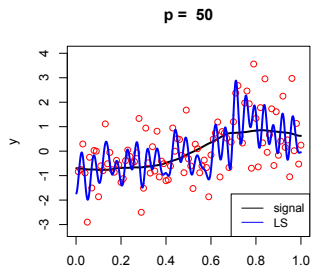
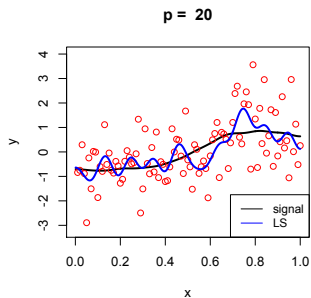
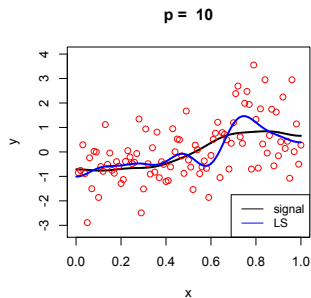
Illustration :

$$Y_i = \sum_{j=1}^p \beta_j^* \cos(\pi j i / n) + \varepsilon_i = f_{\beta^*}(i/n) + \varepsilon_i, \quad \text{for } i = 1, \dots, n,$$

with

- $\varepsilon_1, \dots, \varepsilon_n$ i.i.d with $\mathcal{N}(0, 1)$ distribution
- β_j^* independent with $\mathcal{N}(0, j^{-4})$ distribution

Curse 1 : fluctuations cumulate



Curse 2 : locality is lost

Observations $(Y_i, X^{(i)}) \in \mathbb{R} \times [0, 1]^p$ for $i = 1, \dots, n$.

Model: $Y_i = f(X^{(i)}) + \varepsilon_i$ with f smooth.

Local averaging: $\hat{f}(x) = \text{average of } \{Y_i : X^{(i)} \text{ close to } x\}$

Curse 2 : locality is lost

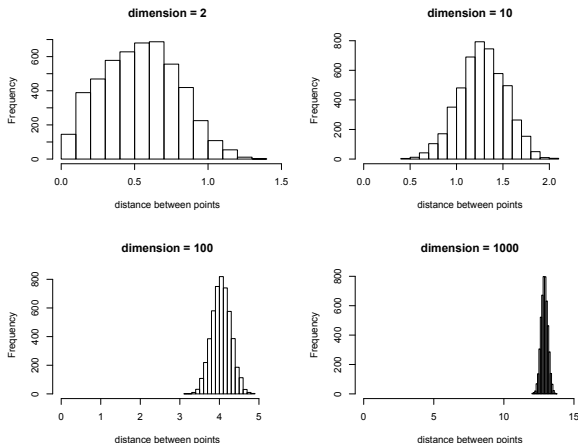


Figure: Histograms of the pairwise-distances between $n = 100$ points sampled uniformly in the hypercube $[0, 1]^p$, for $p = 2, 10, 100$ and 1000 .

Curse 2 : locality is lost

Number n of points x_1, \dots, x_n required for covering $[0, 1]^p$ by the balls $B(x_i, 1)$:

$$n \geq \frac{\Gamma(p/2 + 1)}{\pi^{p/2}} \underset{p \rightarrow \infty}{\sim} \left(\frac{p}{2\pi e}\right)^{p/2} \sqrt{p\pi}$$

p	20	30	50	100	200
n	39	45630	$5.7 \cdot 10^{12}$	$42 \cdot 10^{39}$	larger than the estimated number of particles in the observable universe

Some other curses

- Curse 3 : an accumulation of rare events may not be rare (false discoveries, etc)
- Curse 4 : algorithmic complexity must remain low

Low-dimensional structures in high-dimensional data

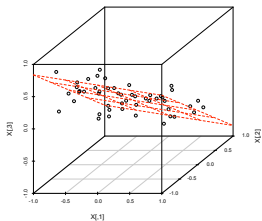
Hopeless?

Low dimensional structures : high-dimensional data are usually concentrated around low-dimensional structures reflecting the (relatively) small complexity of the systems producing the data

- geometrical structures in an image,
- regulation network of a "biological system",
- social structures in marketing data,
- human technologies have limited complexity, etc.

Dimension reduction :

- "unsupervised" (PCA)
- "estimation-oriented"

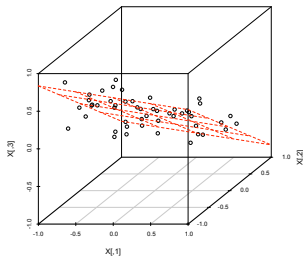


Principal Component Analysis

For any data points $X^{(1)}, \dots, X^{(n)} \in \mathbb{R}^p$ and any dimension $d \leq p$, the PCA computes the linear span in \mathbb{R}^p

$$V_d \in \operatorname{argmin}_{\dim(V) \leq d} \sum_{i=1}^n \|X^{(i)} - \operatorname{Proj}_V X^{(i)}\|^2,$$

where Proj_V is the orthogonal projection matrix onto V .

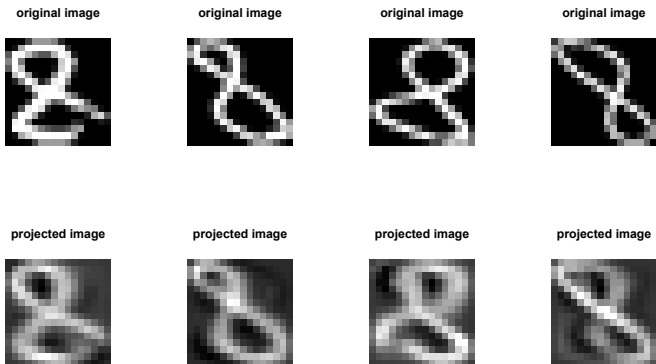


V_2 in dimension $p = 3$.

To do

Exercise 1.6.4

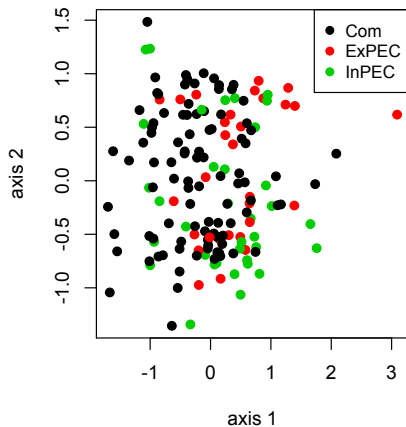
PCA in action



MNIST : 1100 scans of each digit. Each scan is a 16×16 image which is encoded by a vector in \mathbb{R}^{256} . The original images are displayed in the first row, their projection onto 10 first principal axes in the second row.

"Estimation-oriented" dimension reduction

PCA



LDA

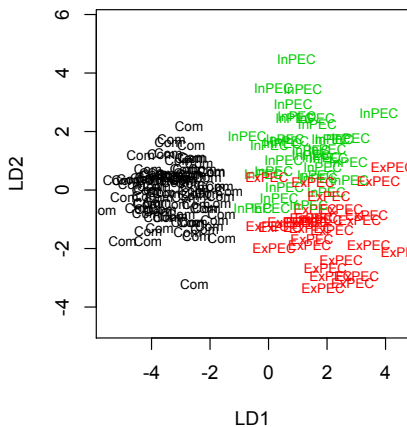


Figure: 55 chemical measurements of 162 strains of *E. coli*.

Left : the data is projected on the plane given by a PCA.

Right : the data is projected on the plane given by a LDA.

Résumé

Difficulté statistique

- données de très grande dimension
- peu de répétitions

Pour nous aider

Données issues d'un vaste système (plus ou moins dynamique et stochastique)

- existence de structures de faible dimension "effective"
- parfois: existence de modèles théoriques

La voie du succès

Trouver, à partir des données, ces structures "cachées" pour pouvoir les exploiter.

La voie du succès

Trouver, à partir des données, les structures cachées pour pouvoir les exploiter.

Exemples de structures

Regression Model

Regression model

$$Y_i = f(x^{(i)}) + \varepsilon_i, \quad i = 1, \dots, n \quad \text{with}$$

- $f : \mathbb{R}^p \rightarrow \mathbb{R}$
- $\mathbb{E}[\varepsilon_i] = 0$

Vectorial representation

The observations can be summarized in a vector form

$$Y = f^* + \varepsilon \in \mathbb{R}^n$$

with $f_i^* = f(x^{(i)})$, $i = 1, \dots, n$.

Low-dimensional x

Example 1: sparse piecewise constant regression

It corresponds to the case where $f : \mathbb{R} \rightarrow \mathbb{R}$ is piecewise constant with a small number of jumps.

This situation appears for example for CGH profiling.

Low-dimensional x

Example 2: sparse basis/frame expansion

We estimate $f : \mathbb{R} \rightarrow \mathbb{R}$ by expanding it on a basis or frame $\{\varphi_j\}_{j \in \mathcal{J}}$

$$f(x) = \sum_{j \in \mathcal{J}} c_j \varphi_j(x),$$

with a small number of nonzero c_j . Typical examples of basis are Fourier basis, splines, wavelets, etc.

The most simple example is the piecewise linear decomposition where $\varphi_j(x) = (x - z_j)_+$ where $z_1 < z_2 < \dots$ and $(x)_+ = \max(x, 0)$.

High-dimensional x

Example 3: sparse linear regression

It corresponds to the case where f is linear: $f(x) = \langle \beta, x \rangle$ where $\beta \in \mathbb{R}^p$ has only a few nonzero coordinates.

This model can be too rough to model the data.

Example: relationship between some phenotypes and some protein abundances.

- only a small number of proteins influence a given phenotype,
- but the relationship between these proteins and the phenotype is unlikely to be linear.

High-dimensional x

Example 4: sparse additive model

In the sparse additive model, we expect that $f(x) = \sum_k f_k(x_k)$ with most of the f_k equal to 0.

If we expand each function f_k on a frame or basis $\{\varphi_j\}_{j \in \mathcal{J}_k}$ we obtain the decomposition

$$f(x) = \sum_{k=1}^p \sum_{j \in \mathcal{J}_k} c_{j,k} \varphi_j(x_k),$$

where most of the vectors $\{c_{j,k}\}_{j \in \mathcal{J}_k}$ are zero.

Such a model can be hard to fit from a small sample since it requires to estimate a relatively large number of non-zero $c_{j,k}$.

High-dimensional x

In some cases the basis expansion of f_k can be sparse itself.

Example 5: sparse additive piecewise linear regression

The sparse additive piecewise linear model, is a sparse additive model $f(x) = \sum_k f_k(x_k)$ with sparse piecewise linear functions f_k . We then have two levels of sparsity :

- 1 most of the f_k are equal to 0,
- 2 the nonzero f_k have a sparse expansion in the following representation

$$f_k(x_k) = \sum_{j \in \mathcal{J}_k} c_{j,k} (x_k - z_{j,k})_+$$

In other words, the matrix $c = [c_{j,k}]$ of the sparse additive model has only a few nonzero columns, and this nonzero columns are sparse.

Reduction to a structured linear model

Reduction to a structured linear model

In all the above examples, we have a linear representation

$$f_i^* = \langle \alpha, \psi_i \rangle \quad \text{for } i = 1, \dots, n,$$

with a structured α .

Examples (continued)

Representation $f_i^* = \langle \alpha, \psi_i \rangle$

- Sparse piecewise constant regression: $\psi_i = e_i$ with $\{e_1, \dots, e_n\}$ the canonical basis of \mathbb{R}^n and $\alpha = f^*$ is piecewise constant.
- Sparse basis expansion: $\psi_i = [\varphi_j(x^{(i)})]_{j \in \mathcal{J}}$ and $\alpha = c$.
- Sparse linear regression: $\psi_i = x^{(i)}$ and $\alpha = \beta$.
- Sparse additive models: $\psi_i = [\varphi_j([x_k^{(i)}])]_{\substack{k=1, \dots, p \\ j \in \mathcal{J}_k}}$ and $\alpha = [c_{j,k}]_{\substack{k=1, \dots, p \\ j \in \mathcal{J}_k}} \cdot$