

Statistical Learning with High-dimensional Data



Pr. Charles BOUVEYRON

Professor of Statistics
Chair of the Institut 3IA Côte d'Azur
Université Côte d'Azur & Inria

charles.bouveyron@univ-cotedazur.fr
@cbouveyron

Outline

1. Introduction
2. Reminder on the learning process
3. Learning in high-dimensions
4. Dimension reduction
5. Clustering and classification

Dimension reduction

A common phantasm about dimension reduction:

- believe that dimension reduction helps for classification,
- **this is not true** because, most of the time, dimension reduction implies an information loss which would be discriminative.

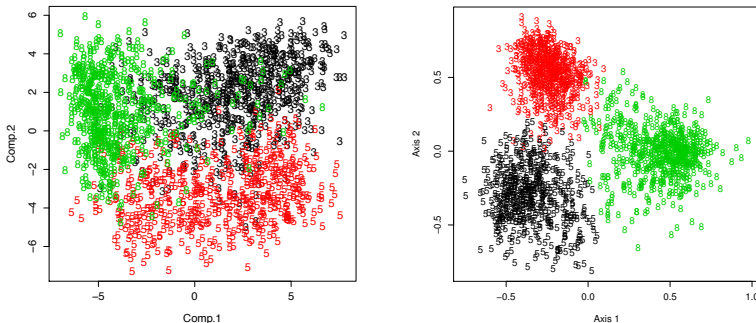
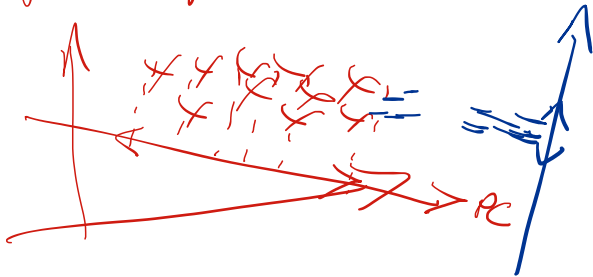


Figure: Projection of the 256-dimensional USPS data with PCA (left, unsupervised) and FDA (right, supervised).

PCA: principal component analysis

The goal of PCA is to create d new variables, which are linear combinations of the original var, such that the variance of the projected data is maximum.



PCA: the principle

$$\max_U \sqrt{X^t X} U$$

covariance matrix of the centered data

The solution: \hat{U} is a matrix made of the d eigenvectors associated to the largest eigenvalues

$$\bar{X}^t \bar{X} = S = Q^t \Delta Q \quad \text{where } \Delta = \text{diag}(\lambda_1, \dots, \lambda_d)$$

\uparrow SPD matrix

PCA: the principle

The recipe :

- (i) Compute $S = \overline{X^T X}$
- (ii) eigendecomposition of S
- (iii) keep the vectors q_j such that $\lambda_j, j=1 \dots d$ are the largest eigenvalues.

PCA: projection

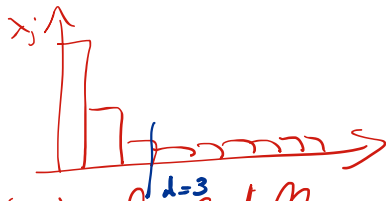
$$y = X \times U$$

$n \times d$ $m \times p$ $p \times d$

\uparrow scores.

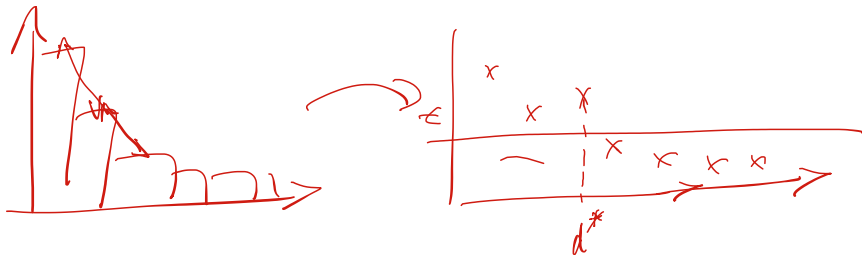
PCA: how many axes?

- (i) The 90% rule: we keep p variables such that
- $$\frac{\sum_{j=1}^d \lambda_j}{\sum_{j=1}^p \lambda_j} \geq 0.9$$
- (ii) The knee test: looking at the eigenvalue scree



- (iii) the scree-test of Cattell: looking at the differences between eigenvalues, we retain d such that all diff after d are smaller than a threshold

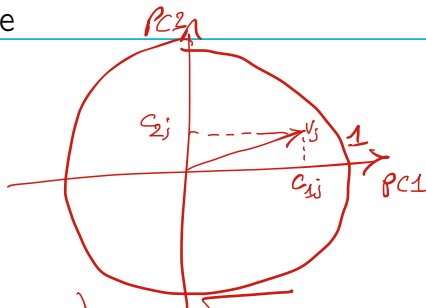
PCA: how many axes?



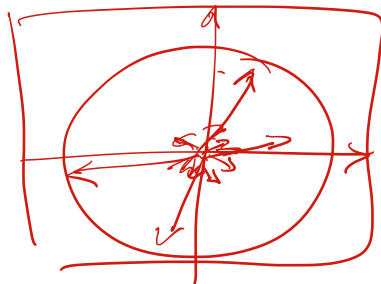
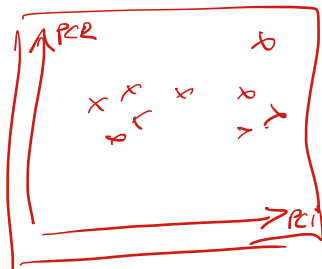
$$t = 0.1 \times \text{Max}(\text{diff})$$

(iv) The other solutions (statistical tests, Bayesian app, ...) are clearly more complex and perform only slightly better.

PCA: correlation circle



$$c_{ij} = \text{Corr}(PC_i, V_j) = \sqrt{\lambda_i} q_{ij}$$



PCA: analysis of the marathon data

See the R code ...

PPCA: a probabilistic version of PCA

PPCA was proposed by Tipping & Bishop (1996) to better establish the theory of PCA under a Gaussian distribution:

$$X \in \mathbb{R}^p$$

$$\bullet X = U^t Y + \epsilon$$

$$\bullet Y \sim N(\mu, I_d)$$

$$\epsilon \sim N(0, \sigma^2 I_p)$$

is called the
latent variable
 $Y \in \mathbb{R}^d$

$$\Rightarrow p(X|\theta) = N(U^t \mu, U^t U + \sigma^2 I_p) \quad \text{with } \theta = \{\mu, U, \sigma^2\}$$

PPCA: a probabilistic version of PCA

In this case, doing a (P)PCA is equivalent to estimate the parameters $\theta = \{\mu, \underline{U}, \sigma^2\}$ from the data.

Max Likelihood $\Rightarrow \hat{U}_{ML} =$ the eigenvectors of $\bar{X}^T \bar{X}$ associated to the largest eigen values.

$$\mu = \bar{X}$$

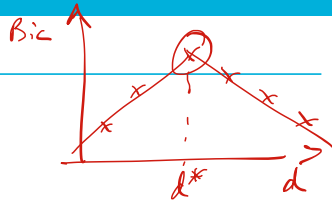
PPCA: why is it interesting?

- (i) justification of a very old result
- (ii) it allows to use model selection and all statistical tools within PCA
- (iii) it allows to propose new models based on this PPCA model

PPCA in practice...

PPCA = PCA in practice

except that we can use AIC or BIC
for choosing d , ...



AIC, BIC (...) are model selection criteria which
are based on penalized likelihood

$$\text{BIC}(\mathcal{M}) = \log(\hat{\sigma}^2) - \frac{\gamma(\mathcal{M})}{2} \log(n)$$

where $\gamma(\mathcal{M})$ is the nb of parameters in \mathcal{M} .

MDS: multi-dimensional scaling

The idea of MDS is to find a low-dim representation of the data that keeps the original topology of the data.

$$X \in \begin{matrix} \mathbb{R}^p \\ \text{Graph} \\ \dots \end{matrix} \longrightarrow Y \in \mathbb{R}^d$$

where it is possible to compute a distance



MDS: the principle

Given the distances between points in the original space, let say $d_{ij} \forall i, j = 1 \dots n$

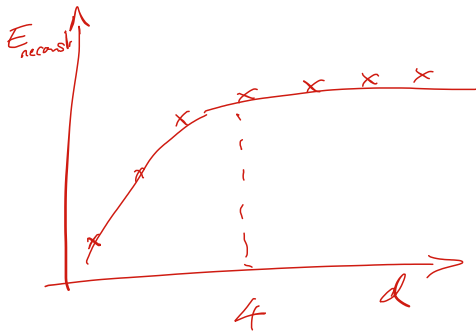
$$\min_g \sum_{i,j=1}^n \| \underline{d_{ij}} - \underline{\delta_{ij}} \|^2$$

where $\underline{\delta_{ij}} = \| \underline{y_i} - \underline{y_j} \|^2$

MDS in practice...

in R : `cmdscale(dist(X) , d)`

↑ it remains to
choose d



t-SNE: t-distributed stochastic neighbor embedding

t-SNE is a model-based version of MDS

(i) SNE : $\forall j=1 \dots n, j \neq i$ $x_j \sim N(x_i, \sigma_i^2)$

$p_{j|i}$ = probability that i chooses j as neighbor

\propto
Bayes $\frac{\exp(-\|x_i - x_j\|^2 / 2\sigma_i^2)}{\sum_{k \neq i} \exp(-\|x_i - x_k\|^2 / 2\sigma_i^2)}$

} Min KL div

$q_{j|i} = \frac{\exp(-\|y_i - y_j\|^2 / 2\sigma_i^2)}{\sum \quad \quad \quad}$

t-SNE: from SNE to T-SNE

$$\min_y \sum_{i \neq j} p_{j|i} \log \left(\frac{p_{j|i}}{q_{j|i}} \right)$$

⊖ very difficult to optimize

⊖ a very asymmetric modeling.

t-SNE: the principle

(ii) t-SNE:

- symmetry : $p_{ij} = p_{ji} = \frac{p_{j|i} + p_{i|j}}{2n}$

- distribution in the reduced space

$$\forall j \neq i, \quad y_j \sim T(y_i, 1)$$

$$q_{ij} = \frac{\left(1 + \|y_i - y_j\|^2\right)^{-1}}{\sum_{k \neq j} \quad \quad \quad}$$