

# PCA and Classification in Symbolic Context

SDA 2011, Beijing, China

Richard Emilion

MAPMO Lab., University of Orléans, France

October 27, 2011

- **I - Symbolic PCA**
- **II - Unsupervised Classification of multihistograms**
- **III - Conclusion**

# I - Principal Component Analysis in symbolic context

## PCA in $\mathbb{R}^p$ PCA in Hilbert space

- $x_i = (x_{i,1}, \dots, x_{i,p}) \in \mathbb{R}^p, \alpha_i, i = 1, \dots, n$ : cloud of weighted vectors.

Affine line in  $\mathbb{R}^p$  minimizing the distance between the cloud and its orthogonal projection on this line ?

Affine plane containing the above line minimizing the the distance between the cloud and its orthogonal projection on this plane ? etc

Unique solution: line =  $(g, v_1)$ , plane =  $(g, v_1, v_2)$ , etc ...

$g$  : barycenter of the cloud,  $v_1, v_2, \dots, v_p$  orthogonal eigenvectors of the positive definite symmetric correlation matrix  $C$  of the  $p$  column vectors  $y_j = (x_{i,j}) - g_j \in \mathbb{R}^n$  associated to the eigenvalues  $\lambda_1 \geq \lambda_2, \dots \geq \lambda_p \geq 0$  of  $C$ , resp.

- Linear method (linear algebra, linear space)
- Dimension reduction method (from  $p$  to 1 or 2)
- Graphical representation: trends, clusters, correlations, outliers

$x_i$  : a class of data

- Examples of classes of data: a district of individuals, an image of pixels ...
- Description of the classes: mean, variance, quantiles, interval, histogram, ...
- How to perform PCA on units which are classes of data: symbolic PCA

```

5.5735987e+00 1.4157128e+01 7.4254032e+00 7.8532778e+00 7.5762937e-01 3.8035638e-01 6.7518185e-01 6.5644447e-01 1.7278023e-02 1.0632455e-02 1.5002581e-
02 1.4248458e-02 5.1229317e-01 3.9179368e-01 4.8897378e-01 4.8391289e-01
6.1905987e+00 1.2295111e+01 4.8628032e+00 6.3413111e+00 7.5432557e-01 5.1148637e-01 8.0429611e-01 1.7734111e-01 1.1990523e-02 1.8805481e-
02 1.6899133e-02 5.3950466e-01 4.3054111e-01 5.4164947e-01 5.1545141e-01
6.0917328e+00 1.4381984e+01 5.7864848e+00 7.1289111e+00 4.4727087e-01 4.8929128e-01 7.6326437e-01 7.0573524e-01 1.6085864e-02 1.1190351e-02 1.6042472e-
02 1.5662340e-02 5.1677307e-01 4.1876964e-01 5.1904722e-01 5.0723535e-01
7.2379032e+00 1.4232025e+01 6.0395917e+00 7.9484848e+00 7.2139380e-01 4.3552378e-01 7.7804689e-01 4.8429359e-01 1.0425859e-02 1.6748419e-02
1.4011378e+02 4.9206217e-01 4.6253358e-01 5.2327469e-01 4.7811064e-01
7.0443548e+00 1.3495117e+01 5.8689515e+00 7.9491111e+00 7.3093930e-01 4.8226071e-01 7.7673769e-01 6.9440662e-01 1.5556407e-02 1.1222809e-02 1.6313471e-
02 1.1743643e-02 3.2419338e-01 3.9894211e-01 4.7972939e-01 4.4480638e-01
7.0907258e+00 1.4154444e+01 6.1391129e+00 7.5962358e+00 7.0753766e-01 4.2166936e-01 7.4714313e-01 6.8063437e-01 1.3495053e-02 7.8633382e-03 1.3976225e-
02 1.4326804e+00 1.1943808e+01 7.0705444e+00 7.4815517e+00 8.1285211e-01 5.2523111e-01 7.1818246e-01 6.9485475e-01 1.4940593e-02 9.5631111e-02 1.2234977e-
02 1.4979388e+00 1.6341892e+01 5.0108088e+00 4.4807048e+00 7.5185757e-01 5.3149461e-01 7.7418510e-01 6.8068353e-01 1.6875427e-01 1.2371132e-02 1.7982898e-02
1.5551352e+00 5.2423434e+01 4.4466999e+00 5.2464678e-01 4.9935504e-01
5.0413388e+00 1.1187111e+01 5.4861288e+00 4.8284688e+00 8.0452850e-01 6.1186514e-01 8.1493920e-01 7.8456874e-01 1.7894627e-02 1.2152945e-02 1.6924617e-
02 1.5982311e-02 5.5145777e-01 4.3681148e-01 5.2977591e-01 5.2036284e-01
4.8981958e+00 1.3232781e+01 6.0282358e+00 9.4146848e+00 7.1640818e-01 4.5175021e-01 7.2974382e-01 4.2635287e-01 1.7564979e-02 1.2573185e-02 1.7538072e-
02 1.6322082e-02 5.1428811e-01 4.2187194e-01 5.0351302e-01 4.7984284e-01
4.5794371e+00 1.4377723e+01 7.4546371e+00 8.5067638e+00 7.3295128e-01 4.1775577e-01 6.9722379e-01 6.5523366e-01 1.5416172e-02 1.0489204e-02 1.4013412e-
02 1.3204897e-02 3.0546755e-01 3.9810334e-01 4.7358150e-01 4.4745589e-01
5.7641129e+00 1.1940911e+01 5.7004048e+00 6.4274714e+00 6.9320533e-01 3.6640350e-01 6.9665704e-01 4.6977070e-01 1.7108463e-02 1.2612058e-02 1.8432323e-
02 1.5843170e-02 5.0458888e-01 4.8002747e-01 5.2028290e-01 4.7681385e-01
5.8528228e+00 1.4585723e+01 7.4578411e+00 4.8189751e+00 7.4487389e-01 3.7084999e-01 6.7487035e-01 4.9964804e-01 2.0362719e-02 1.3140921e-02 1.8564846e-
02 1.8787237e-02 5.3864767e-01 4.1252214e-01 4.9404602e-01 5.1225403e-01
5.7318548e+00 1.2407100e+01 6.3891129e+00 7.1675338e+00 7.2328772e-01 4.0786841e-01 6.9984087e-01 6.5703073e-01 1.9139597e-01 1.3180114e-02 1.6375329e-02
1.7589453e+00 5.2333989e+01 4.0534811e-01 5.0779203e-01 4.9299118e-01
6.5382329e+00 1.1238282e+01 5.2419358e+00 4.8021851e+00 6.3525047e-01 3.7187496e-01 7.0420173e-01 6.1335867e-01 1.5923950e-02 1.2512999e-02 1.7216804e-02
1.5844178e+02 4.9585453e-01 4.1844848e-01 5.0738087e-01 4.80308489e-01
6.1280242e+00 1.2484877e+01 4.5473795e+00 4.8438084e+00 6.9389402e-01 2.4951703e-01 6.6792739e-01 4.6940412e-01 1.7567455e-02 1.2757686e-02 1.7587721e-
02 1.6264384e+00 5.0997006e-01 4.1384282e-01 4.9533899e-01 4.7610232e-01
7.0322641e+00 1.3182641e+01 4.3508011e+00 7.2121261e+00 6.6811623e-01 3.7559673e-01 6.9593636e-01 4.5385737e-01 1.4274478e-02 1.071136e-02 1.4514300e-
02 1.1251241e-02 4.9209192e-01 4.0044796e-01 4.8283711e-01 4.6452648e-01
6.8974184e+00 1.2847034e-01 4.4838713e+00 7.1941172e+00 6.6975572e-01 3.8224160e-01 6.8729191e-01 4.3470974e-01 1.3925415e-02 1.0646798e-02 1.4478491e-02
1.1154457e-02 4.9150766e-01 4.1037016e-01 4.3287314e-01 4.4452667e-01

```

Figure: Each image: 400 pixels, Each pixel: 16 measurements

One image  $x_i$  = One matrix  $k = 400 \times p = 16$

$x_{i,1}$  image  $i$  measurement : column vector  $k \times 1 \in \mathbb{R}^k$

...

$x_{i,p}$  image  $i$  measurement  $p$  : column vector  $k \times 1 \in \mathbb{R}^k$

- E. Diday et. al.:
  - Columns of raw data are summarized by intervals, histograms, ...
  - extended PCA
  - a lot of raw informations are lost: multivariate dist., correlations
  - basic vector space is changed and not precised
  - consistency problems (results depend on the bins)
- R. Emilion:
  - Hilbert space PCA considering both **raw data** (if available) **and classes**
  - Summarize just at the end the results of PCA for visualization purposes of classes

More generally, our approach of the symbolic case: higher level math. representation using as far as possible both huge raw datasets nowadays available, classes, power of computers. Sometimes raw data are not available (private, cost, confidentiality, security ...)

- $\mathbb{H}$ : Hilbert space, i.e. vector space with an inner product  $\langle \cdot, \cdot \rangle$ , complete.

Examples : matrices, square integrable functions, product of Hilbert spaces, ...

- $x_i = (x_{i,1}, \dots, x_{i,p}) \in \mathbb{H}^p, \alpha_i, i = 1, \dots, n$ : cloud of weighted vectors.

$$g = \sum_{i=1}^n \alpha_i x_i \in \mathbb{H}^p$$

: barycenter of the cloud

$v_1, v_2, \dots, v_p$  orthogonal eigenvectors of the positive definite symmetric correlation,  $p \times p$  real matrix  $C$  of the  $p$  column vectors  $(x_{i,j}) - g_j \in \mathbb{H}^n$  associated to the eigenvalues

$\lambda_1 \geq \lambda_2, \dots, \geq \lambda_p \geq 0$  of  $C$ , resp.

- Change of coordinates from the canonical basis of  $\mathbb{R}^p$  to the  $(v_i)$
- Apply the same linear changes to the  $x_i$ 's to get new coordinates  $y_i$ 's: principal components.
- Dimension reduction (from  $\mathbb{H}^p$  to  $\mathbb{H}^1$  or  $\mathbb{H}^2$ )

- First and second component:  $n$  vectors in  $\mathbb{R}^k$
- If  $k$  is large: statistics on these vectors.

Example: if use the means of the components (and not the components of the means): each image is represented by a point in the plane.

## II- Classification of multihistograms

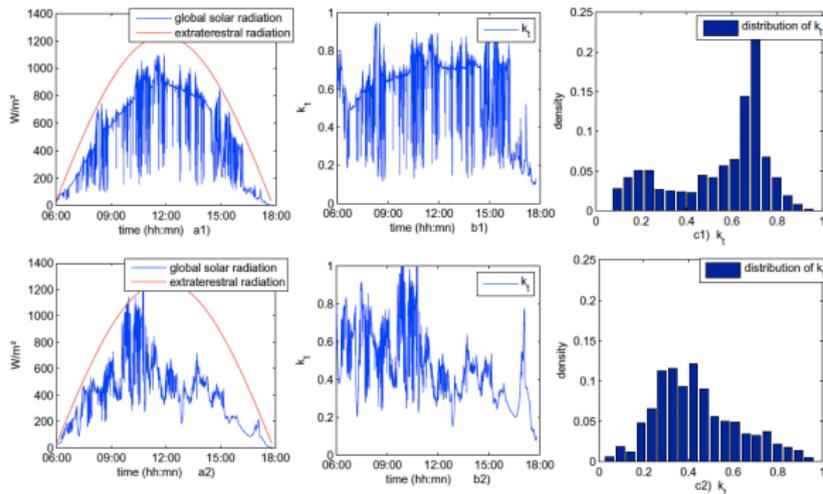
Example

Dirichlet mixtures

Consistency theorem

1058

*T. Soudhan et al / Solar Energy 83 (2009) 1056–1063*



**Figure:** Each solar day stochastic process: 12h x 3600 measurements, histogram

- Each histogram is a vector of probability
- histograms are not those of any parametric distributions
- Considered as outputs of a random variable: random distribution
- Estimating this random distribution as a mixture of Dirichlet distributions (Solar Energy 2009)
- The Classification by estimating this mixture is consistent if the true distribution is an outcome of a mixture of Dirichlet processes: proof uses the Martingale convergence theorem (submitted to SAM)

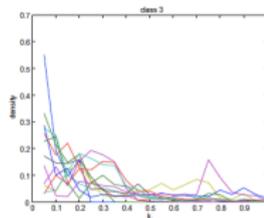
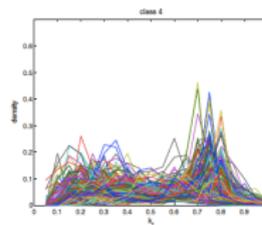
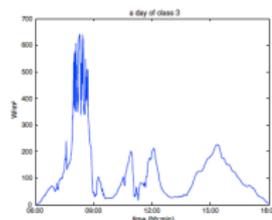
Fig. 6. Daily distributions of  $k_t$  in class 3.Fig. 8. Daily distributions of  $k_t$  in class 4.

Fig. 7. A day of class 3.

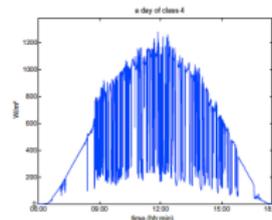


Fig. 9. A day of solar radiation in class 4.

Figure: 4 Classes of histograms were identified

- One day: Solar radiation, humidity, Temperature, ... stochastic processes
- Marginal histogram for each process: do not capture correlations
- Consider one multihistogram for each day
- Outputs of a Dirichlet mixture with Dirichlet's defined on a product space
- Same story

- Link between the two parts: Symbolic PCA and then classification of 2-dim histograms of the 2 best principal components
- SDA requires more complex mathematics
- What is a symbolic object:  
Symbolic Galois concepts (Diday-Emilion 1995)  
Concepts: representation of classes, First order logic (Diday).  
Both can be derived from Galois connections.
- Historically point. J. F. C. Kingmann (1975). J. Royal Stat. B, 37. Random discrete distributions.  
Spoke of many previous works where people were interested with 'Random objects which are themselves probability distributions'.