# Data analysis Principal component analysis

Olivier Roustant

February 14, 2022

**Olivier Roustant** 

Principal Component Analysis (PCA): Outline

# Figures only!



**3** Variations (metric, weights)

4 Results interpretation

## 5 Conclusion and further readings

#### The aim: To reduce dimension



This is a 2*D* cloud of points, centered at 0. Can you find a 1D axis 'containing' the maximum of information?

#### Figures only

## Inertia



Total inertia: 5.402

Total inertia: mean square of distances to the center.

```
Olivier Roustant
```

#### Figures only

## Inertia



Projected inertia on u: 4.351

x1

Projected inertia: inertia of projections. How much do we lose?

Projected inertia on u: 4.351



Projected inertia: For what axis is it maximal?

**Olivier Roustant** 

Projected inertia on u: 4.798



Projected inertia: For what axis is it maximal?

**Olivier Roustant** 

Projected inertia on u: 4.993



Projected inertia: For what axis is it maximal?

**Olivier Roustant** 

Projected inertia on u: 4.911



#### Projected inertia: For what axis is it maximal?

**Olivier Roustant** 



Projected inertia on u: 4.562

### Projected inertia: For what axis is it maximal?

Olivier Roustant





#### Projected inertia: For what axis is it maximal?

**Olivier Roustant** 





Projected inertia: For what axis is it maximal?

**Olivier Roustant** 





Projected inertia: For what axis is it maximal?

**Olivier Roustant** 



Projected inertia on u: 1.698

#### Projected inertia: For what axis is it maximal?

Olivier Roustant

Projected inertia on u: 1.051



Projected inertia: For what axis is it maximal?

**Olivier Roustant** 

Projected inertia on u1: 4.999



Projected inertia: Maximal for the largest eigenvalue of the covariance matrix

Olivier Roustant

## Maximizing the projected inertia, recursion





The second largest eigenvalue maximizes the projected inertia in the orthogonal of the first

ivie	or F	Rou	sta	ant
	51 1	iuu	510	21.11

## Maximizing the projected inertia, summary





Projected points on the first two 'principal components'

OI	ivier	Ro	ust	ant
· • ·				~

## Maximizing the projected inertia, summary

Individuals factor map (PCA)



Representation with package FactoMineR. Percentages are inertia ratio w.r.t. total inertia

-		_		
$-\alpha$	IN/IOT		1101	201
<u> </u>		I NU		
0		110	uu	

## Theory

#### Notations and assumption

• **X**: a matrix of size  $n \times p$ , representing the data:



• **g**: center of gravity (empirical mean),  $\mathbf{g} = \bar{\mathbf{x}} = \frac{1}{n} \sum_{i=1}^{n} \mathbf{x}_i (\in \mathbb{R}^p)$ .

$$\mathbf{g} \| \overline{\mathbf{x}^1} \dots \overline{\mathbf{x}^j} \dots \overline{\mathbf{x}^p}$$

We assume that g = 0, i.e. the data have been centered.

#### Notations and assumption

- The rows of X lie in ℝ<sup>ρ</sup>, and form the indivuals space.
   It is an Euclidean space, equipped with the usual ℓ<sup>2</sup> norm ||.||.
- The columns of X lie in ℝ<sup>n</sup>, and form the variables space.
   It is an Euclidean space. Instead of choosing the usual ℓ<sup>2</sup> norm, we rescale it by 1/n. Indeed, as the data are centered, it corresponds to the empirical covariance:

$$\langle \mathbf{x}^j, \mathbf{x}^k \rangle_{\mathbb{R}^n} := \frac{1}{n} \sum_{i=1}^n x_i^j x_i^k = \widehat{\operatorname{cov}}(\mathbf{x}^j, \mathbf{x}^k).$$

Notice that **orthogonal variables = uncorrelated variables**.  $\Gamma$  denotes the  $p \times p$  empirical covariance matrix:

$$\Gamma = \left(\widehat{\operatorname{cov}}(\mathbf{x}^j, \mathbf{x}^k)\right)_{1 \le j, k \le p} = \frac{1}{n} \mathbf{X}^\top \mathbf{X} = \frac{1}{n} \sum_{i=1}^n \mathbf{x}_i \mathbf{x}_i^\top.$$

#### Notations and assumption

• Inertia: mean squared distance of the data to their center (here 0),

$$U = \frac{1}{n} \sum_{i=1}^{n} \|\mathbf{x}_i\|^2$$

Projected inertia on a subspace F ⊆ ℝ<sup>ρ</sup>. Same definition for the projected points onto F (we denote by Π<sub>F</sub> the projection operator):

$$\mathcal{I}_F = \frac{1}{n} \sum_{i=1}^n \|\Pi_F(\mathbf{x}_i)\|^2$$

#### Theory

#### **Properties of inertia**

## Link with variance, and inertia decomposition.

Consider a 1*D* axis spanned by a unit vector **a**, and denote  $\mathcal{I}_{a} = \mathcal{I}_{\mathbb{R}a}$ . Then:

$$\mathcal{I}_{\mathbf{a}} = \mathbf{a}^{\top} \Gamma \mathbf{a}, \quad \text{and} \quad \mathcal{I} = \mathcal{I}_{\mathbf{a}} + \mathcal{I}_{\mathbf{a}^{\perp}}$$

Moreover,  $\mathcal{I}_a$  and  $\mathcal{I}$  are interpreted in terms of variances:

- *I*<sub>a</sub> is the empirical variance of the projected points onto Ra,
- *I* is the sum of the empirical variances of the *p* variables:

$$\mathcal{I}_{\mathbf{a}} = \frac{1}{n} \sum_{i=1}^{n} \langle \mathbf{x}_i, \mathbf{a} \rangle^2, \qquad \mathcal{I} = \sum_{j=1}^{p} \hat{\sigma}_j^2, \quad \text{with} \quad \hat{\sigma}_j^2 = \frac{1}{n} \sum_{i=1}^{n} (x_i^j)^2$$

Remark: The empirical variances are computed here by dividing by n the sum of squares, contrarily to unbiased statistical estimates (division by n - 1).

**Properties of inertia (proofs)** 

Left to exercise.

#### Main result

### Theorem (principal component analysis)

As the covariance matrix  $\Gamma$  is real symmetric, it admits a spectral decomposition in orthogonal eigenspaces. Denote  $\lambda_1 \ge \cdots \ge \lambda_p \ge 0$  the eigenvalues, and  $\mathbf{v}_1, \ldots, \mathbf{v}_p$  orthogonal eigenvectors. Then:

- $\mathbf{v}_1$  maximizes  $\mathcal{I}_{\mathbf{a}}$  over  $\mathbf{a}$ , which is then equal to  $\lambda_1$ .
- $\mathbf{v}_2$  maximizes  $\mathcal{I}_{\mathbf{a}}$  over  $\mathbf{a}$  in  $(\mathbf{v}_1)^{\perp}$ , which is then equal to  $\lambda_2$ .
- v<sub>3</sub> maximizes *I*<sub>a</sub> over a in (v<sub>1</sub>, v<sub>2</sub>)<sup>⊥</sup>, which is then equal to λ<sub>3</sub>.
  ...

Furthermore the inertia (called *total inertia*) is decomposed:

$$\mathcal{I} = \mathcal{I}_{\mathbf{v}_1} + \dots + \mathcal{I}_{\mathbf{v}_p} = \lambda_1 + \dots + \lambda_p$$

Main result (proof)

Left to exercise. Hint: Use the decomposition of **a** in the basis of eigenvectors.

#### **Principal components**

- The eigenvectors  $\mathbf{v}_1, \ldots, \mathbf{v}_p$  define a new orthonormal basis in  $\mathbb{R}^p$ .
- The change of variables is defined by:

$$\mathbf{C} = \mathbf{XP}, \quad \text{with} \quad \mathbf{P} = [\mathbf{v}_1, \dots, \mathbf{v}_p].$$

The  $n \times p$  matrix **C** is called matrix of principal components. The columns of **C** are called principal variables. They contain the coordinates of the individuals in the new space.

#### **Principal components**

 Principal variables are linear combinations of the original variables, with coefficients given by the eigenvectors:

$$\mathbf{C}^j = \mathbf{X}\mathbf{P}_j = \sum_{k=1}^p (\mathbf{v}_j)_k \mathbf{x}^k$$

• Principal variables are uncorrelated and  $\widehat{var}(\mathbf{C}^k) = \lambda_k$ :

$$\left(\widehat{\operatorname{cov}}(\mathbf{C}^{j},\mathbf{C}^{k})\right)_{1\leq j,k\leq p} = \frac{1}{n}\mathbf{C}^{\top}\mathbf{C} = \mathbf{P}^{\top}\mathbf{\Gamma}\mathbf{P} = \operatorname{diag}(\lambda_{1},\ldots,\lambda_{p}).$$

## Remark: singular value / spectral decomposition

PCA can be done with **Singular Value Decomposition (SVD)**, which decomposes a rectangular matrix  $n \times m$  or rank r as

$$\mathbf{X} = \mathbf{U} \Lambda^{1/2} \mathbf{V}^{\top},$$

where  $\Lambda$  is the diagonal matrix containing the *r* non-zero eigenvalues of  $\mathbf{X}^{\top}\mathbf{X}$  (or  $\mathbf{X}\mathbf{X}^{\top}$ ), ranked by decreasing order, and **U** (resp. **V**) is an orthogonal matrix for  $\|.\|_{\mathbb{R}^n}$  (resp. for  $\|.\|_{\mathbb{R}^m}$ ) containing the eigenvectors of  $\mathbf{X}\mathbf{X}^{\top}$  (resp.  $\mathbf{X}^{\top}\mathbf{X}$ ).

In the frequent case when p = r (e.g. n > p), we have:

$$\mathbf{V} = \mathbf{P}, \qquad \Lambda = \operatorname{diag}(\lambda_1, \ldots, \lambda_n).$$

(In the general case, **V** contains the *r* columns of **P** corresponding to non-zero eigenvalues.) Further, due to our definition of the scalar product in  $\mathbb{R}^n$ , we have  $\frac{1}{n}\mathbf{U}^{\top}\mathbf{U} = I_p$ . Then, you can recover all the formulas of the textbook, e.g.:

$$\mathbf{C} = \mathbf{X}\mathbf{P} = \mathbf{U}\Lambda^{1/2}\mathbf{P}^{\top}\mathbf{P} = \mathbf{U}\Lambda^{1/2}.$$

## Variations (metric, weights)

## Changing the metric in the individuals space

Consider a new norm on  $\mathbb{R}^{p}$ , called **metric**, defined by a positive definite matrix **M**, of size *p*:

$$\|\mathbf{x}\|_M^2 = \mathbf{x}^\top \mathbf{M} \mathbf{x}.$$

Let **R** be an invertible matrix s.t.  $\mathbf{R}^{\top}\mathbf{R} = \mathbf{M}$  (e.g. square root, Choleski decomposition). Then, the map

$$\mathbf{\mathsf{R}}: \frac{(\mathbb{R}^{\rho}, \|.\|_{M})}{\mathbf{\mathsf{x}}} \xrightarrow{\rightarrow} \frac{(\mathbb{R}^{\rho}, \|.\|)}{\mathbf{\mathsf{Rx}}}$$

is an isometry, and thus preserves distances and orthogonality.

Indeed:  $\|\mathbf{R}\mathbf{x}\|^2 = (\mathbf{R}\mathbf{x})^\top (\mathbf{R}\mathbf{x}) = \mathbf{x}^\top \mathbf{M}\mathbf{x} = \|\mathbf{x}\|_M^2$ .

**Olivier Roustant** 

## Changing the metric in the individuals space

Due to the isometry property, we deduce immediately:

## PCA with / without metric

**v** max. projected inertia for original data  $\mathbf{x}_1, \dots, \mathbf{x}_n$  with metric  $\|.\|_M$   $\Leftrightarrow$  **Rv** max. proj. inertia for transformed data  $\mathbf{Rx}_1, \dots, \mathbf{Rx}_n$  with  $\|.\|$   $\Rightarrow$  **Rv** is an eigenvector of  $\frac{1}{n} \sum_{i=1}^n (\mathbf{Rx}_i) (\mathbf{Rx}_i)^\top = \mathbf{R} (\frac{1}{n} \mathbf{X}^\top \mathbf{X}) \mathbf{R}^\top$   $\Rightarrow$ **v** is an eigenvector of  $(\frac{1}{n} \mathbf{X}^\top \mathbf{X}) \mathbf{M} = \Gamma \mathbf{M}$ 

## Changing the metric in the individuals space

Recall that the data are assumed to be centered.

Example. Standardize (centered) data.

$$\mathbf{M} = \operatorname{diag}\left(\frac{1}{\hat{\sigma}_1^2}, \dots, \frac{1}{\hat{\sigma}_p^2}\right)$$

Then we can choose  $\mathbf{R} = \text{diag}\left(\frac{1}{\hat{\sigma}_1}, \dots, \frac{1}{\hat{\sigma}_p}\right)$ . Thus doing PCA with the metric **M** is equivalent to doing usual PCA on the standardized data.

## Changing the weights in the variable space

In the standard formulation, each individual  $\mathbf{x}_1, \ldots, \mathbf{x}_n$  has weight  $\frac{1}{n}$ .

Obviously, one can use positive weights  $\omega_1, \ldots, \omega_n$  that sum to one. It can be useful if some individuals have more importance.

This can be viewed as an isometric transformation in the space  $\mathbb{R}^n$  by the diagonal matrix containing the square roots of  $\omega_i$ .

The theory is immediately adapted, by modifying the definitions, e.g.:

$$\mathcal{I} = \sum_{i=1}^{n} \omega_i \|\mathbf{x}_i\|^2, \qquad \Gamma = \sum_{i=1}^{n} \omega_i \mathbf{x}_i \mathbf{x}_i^\top.$$

## Link between the notations slides / textbook

In this presentation, we started from the simplest case (centered data, standard Euclidean metric, same weights), and explained how to obtain the general formula.

In the textbook, the general case is considered.

To obtain the formula of the slideshow from the textbook, you should simply use:

$$ar{\mathbf{X}} = \mathbf{X}, \qquad \mathbf{M} = \mathbf{I}_{p}, \qquad \mathbf{D} = rac{1}{n}\mathbf{I}_{n}, \qquad \mathbf{V} = \mathbf{P}, \qquad \mathbf{S} = \mathbf{\Gamma}$$

## **Results interpretation**

#### Example on a temperature dataset



Dataset: Temperature at n = 36 cities (individuals) for p = 12 months (variables).

### **Dimension reduction**

Here the variables are highly correlated, and a strong dimension reduction is expected. The decrease of inertia shows that 2 dimensions explain approx 99% of the variance.



## Interpretation of principal components

Remember that  $\mathbf{C}^{j} = \sum_{k=1}^{p} (\mathbf{v}_{j})_{k} \mathbf{x}^{k}$ . To interpret  $\mathbf{C}^{j}$ , look at  $\mathbf{v}_{j}$ . Here we can plot them as a function of time. •  $\mathbf{C}^{1} \propto (\mathbf{x}^{1} + \dots + \mathbf{x}^{12})$ , proportional to the annual temperature •  $\mathbf{C}^{2} \propto (\mathbf{x}^{5} + \dots + \mathbf{x}^{8}) - (\mathbf{x}^{1} + \mathbf{x}^{2} + \mathbf{x}^{11} + \mathbf{x}^{12})$ , contrast summer/winter



Coordinates of the first 2 eigenvectors in  $\mathbb{R}^{12}$ .

- (1)	N/IOF	Douc	tant
		11005	lani

## **Graphics for individuals**



#### Individuals factor map (PCA)

PCA: Projection on the first 2 principal axis.

 nt

#### **Graphics for variables**

- The principal variables C<sup>k</sup> are orthogonal with variance λ<sub>k</sub>. Thus, they define an orthonormal basis U<sup>k</sup> = C<sup>k</sup>/√λ<sub>k</sub>.
- Consider the coordinates  $a_{i,k}$  of the original variables in this basis

$$a_{j,k} = \operatorname{cov}(\mathbf{X}^j, \mathbf{U}^k).$$

We thus have,  $\|\mathbf{x}^j\|_{\mathbb{R}^n}^2 = \hat{\sigma}_j^2 = \sum_k a_{j,k}^2$ .

• The idea is to plot these coordinates for two principal components.

## Graphics for variables, case of unit variance

When the variables have been normalized (unit variance),

$$a_{j,k} = \operatorname{cor}(\mathbf{X}^{j}, \mathbf{U}^{k}) = \cos(\widehat{\mathbf{X}^{j}, \mathbf{U}^{k}})$$

and  $\sum_{k=1}^{p} a_{j,k}^{2} = 1$ .

- Thus the coordinates  $(a_{i,k})_k$  belong to a *p*-dimensional sphere.
- Further  $(a_{j,1}, a_{j,2})$  belongs to the unit disk:  $a_{j,1}^2 + a_{j,2}^2 \le 1$ . It is closed to the unit circle if  $a_{j,3}, \ldots, a_{j,p}$  are nearly zero. In that case,  $\mathbf{X}^j$  is well-represented by  $\mathbf{C}^1, \mathbf{C}^2$ . This is the circle of correlations for components (1, 2).

## Interpretation of principal components



Variables factor map (PCA)

Coordinates of the variables in the orthonormal basis of principal variables. We see again that Axis 1 weigths all months nearly equally, whereas Axis 2 exhibits a contrast summer / winter.

## Interpretation of principal components



Circle of correlation (normalized variables). Here all variables are well-represented by the first 2 principal components.

0				les mak
U	Ivier	RC	us	lanı

### **Graphics for variables**

#### Exercise

- Check that U<sup>k</sup> is the k-th column of the matrix U of the SVD decomposition of X (see slide 21).
- Check that the coordinate of  $\mathbf{X}^{j}$  onto  $\mathbf{U}^{k}$  is  $a_{j,k} = (v_{k})_{j}\sqrt{\lambda_{k}}$ . Explain the link between the circle of correlation and the plot of eigenvectors coordinates (slides 31 and 35).

## **Conclusion and further readings**

- PCA is a dimension reduction technique which finds uncorrelated variables, called principal variables, that are linear combination of the original ones, which approximate the best the data in the mean-square sense.
- PCA = spectral decomposition of the covariance matrix
  - Up to isometric transformations (metric, weights)
- Several graphs can be used to interpret principal components: projection of individuals, circle of correlation (normalized case).
  - Mind that what you visualize is only a projection. Several tools quantify the quality of the representation.
    - $\rightarrow$  See textbook page 29, 30.