## Pollution Sources Detection via Principal Component Analysis and Rotation

#### Vanessa Kuentz<sup>1</sup>

in collaboration with : Marie Chavent <sup>1</sup> Hervé Guégan<sup>2</sup> Brigitte Patouille <sup>1</sup> Jérôme Saracco<sup>1,3</sup>

> <sup>1</sup>IMB, Université Bordeaux 1, Talence, France (e-mail: vanessa.kuentz@math.u-bordeaux1.fr)

> > <sup>2</sup>ARCANE-CENBG, Gradignan, France

<sup>3</sup>GREThA, Université Bordeaux 4, Pessac, France



# Outline

### Factor Analysis

- Notations
- Recall about PCA
- Factor Analysis model
- Link between FA and PCA

### 2 Rotation of factors

- Motivation
- Criterion
- Remark on rotation in PCA
- Application to air pollution sources detection
  - A three steps process
  - Statistical treatment
  - Some confirmation of the results



#### Issue

- Air pollution = small particles + liquid droplets
- High concentrations of particles = danger to human health, especially fine particles
- Development of air quality control strategies
   → identification of pollution sources
- Receptor modeling, commonly used = Principal Component Analysis
- Part of scientific project initiated by French Ministry of Ecology



Notations Recall about PCA Factor Analysis model Link between FA and PCA

X = (x<sub>i</sub><sup>j</sup>)<sub>n,p</sub> numerical data matrix where n objects are described on p < n variables x<sup>1</sup>,..., x<sup>p</sup>

- $\tilde{X} = (\tilde{x}_i^j)_{n,p}$  standardized data matrix:  $\tilde{x}_i^j = \frac{x_i^j \bar{x}^j}{s^j}$  $(\bar{x}^j, s^j$  sample mean and standard deviation of  $x^j$ )
- *R* sample correlation matrix of  $x^1, \ldots, x^p$ :

 $R = \tilde{X}' M \tilde{X}$  where  $M = \frac{1}{m} I_n$  (with m = n or n - 1)

We have : R = Z'Z with  $Z = M^{1/2}\tilde{X}$ 



Notations Recall about PCA Factor Analysis model Link between FA and PCA

- Replace p variables by q ≤ p uncorrelated components
- The **SVD** of Z (of rank  $r \le p$ ) is  $Z = U\Lambda^{1/2} V'$  with :
  - $\Lambda$  : diagonal matrix (r, r) of nonnull eigenvalues  $\lambda_1, \ldots, \lambda_r$  of Z'Z (in decreasing order)
  - U: orthonormal matrix (n, r) of eigenvectors of ZZ' associated with first r eigenvalues
  - *V* : orthonormal matrix (p, r) of eigenvectors of Z'Z associated with first *r* eigenvalues

• Decomposition of 
$$\tilde{X}$$
:  $\tilde{X} = M^{-1/2} U \Lambda^{1/2} V'$ 



< ロ > < 団 > < 豆 > < 豆 >

Notations Recall about PCA Factor Analysis model Link between FA and PCA

• q = r: In PCA, equation  $\tilde{X} = M^{-1/2} U \Lambda^{1/2} V'$  is written :

$$\tilde{X} = \Psi V'$$
 with  $\Psi = M^{-1/2} U \Lambda^{1/2}$ 

We have  $\Psi = \tilde{X}V (VV' = I_p)$ Columns of  $\Psi$ : principal components  $\psi^k, k = 1, ..., q$ 

 q < r : In practice, user retains only first q < r eigenvalues of Λ</li>

Approximation of  $\left| ilde{X} : \hat{ ilde{X}}_q = \Psi_q V_q' \right|$ 

 $(\Psi_q, V_q : \text{matrices } \Psi \text{ and } V \text{ reduced to first } q \text{ columns})$ 



- FA : underlying common factors + dimension reduction
- Model is written :

$$egin{array}{rcl} ilde{\mathbf{x}} &=& A\mathbf{f} &+& \mathbf{e} \ (p imes 1) && (p imes q)(q imes 1) && (p imes 1) \end{array}$$

with :

- $ilde{\mathbf{x}} = ( ilde{x}^1, \dots, ilde{x}^{
  ho})'$  random vector of  $\mathbb{R}^{
  ho}$
- A : loading (or pattern) matrix
- $\mathbf{f} = (f^1, \dots, f^q)'$  random vector of q common factors
- $\mathbf{e} = (e^1, \dots, e^p)'$  random vector of p unique factors
- $\tilde{x}^{j}$  : linear combination of common factors + unique factor



Notations Recall about PCA Factor Analysis model Link between FA and PCA

- FA model : various estimation methods
  - Principal factor method
  - Maximum likelihood method
  - Principal component method

- ...

- Commonly used method : PCA
  - Easy implementation
  - Reasonable results
  - Default method in statistical softwares

8/22



(I)

Notations Recall about PCA Factor Analysis model Link between FA and PCA

• q = r: In FA (with PCA estimation), equation  $\tilde{X} = M^{-1/2}U\Lambda^{1/2}V'$ is written :

$$ilde{X} = FA'$$

with :

 $-F = M^{-1/2}U$ : factor scores matrix  $-A = V\Lambda^{1/2}$ : loading (or pattern) matrix

We have  $F = \tilde{X} V \Lambda^{-1/2} (V' V = I_q)$ Columns of F: common factors  $f^k, k = 1, ..., q$ 

● *q* < *r* :

In practice, user retains only first q < r eigenvalues of  $\Lambda$ 

Approximation of  $\tilde{X}$ :  $\begin{vmatrix} \hat{X}_q = F_q A'_q \end{vmatrix}$ ( $F_q$ ,  $A_q$ : matrices F and A reduced to first q columns)



Notations Recall about PCA Factor Analysis model Link between FA and PCA

We have :

 $\rightarrow$  factors  $f^k$  correspond to standardized principal components :

$$f^k = rac{\psi^k}{\sqrt{\lambda_k}}, k = 1, \dots, q$$



Motivation Criterion Remark on rotation in PCA

- Loading matrix  $A_q$ :  $a_j^k = \operatorname{corr}(x^j, f^k)$ 
  - $\rightarrow$  identification of groups of correlated elements
- Problem :

Intermediate values  $\rightarrow$  difficult association variables/factors

- Objective :
  - Column of  $A_q$ : values close to 0 or 1
  - Row of  $A_q$ : only one value close to 1
- Solution :
  - Non-uniqueness of FA solution : rotation of factors
  - Orthogonal transformation matrix  $T (TT' = T'T = I_q)$
  - $\check{A}_q = A_q \times T$ : correlation of variables to rotated factors  $\check{f}^k$ - How to determine T?
- MAR A STATE

Motivation Criterion Remark on rotation in PCA

- Several criteria : most used is varimax
- Maximizes empirical variance of squared elements of each column of Ă<sub>q</sub> = (ă<sub>j</sub><sup>α</sup>)<sub>ρ,q</sub>:

$$\sum_{\alpha=1}^{q} \left\{ \frac{\sum_{j=1}^{p} (\check{a}_{j}^{\alpha})^{4}}{p} - \left( \frac{\sum_{j=1}^{p} (\check{a}_{j}^{\alpha})^{2}}{p} \right)^{2} \right\}$$

with  $\check{a}_j^{\alpha} = a_j t^{\alpha}$  and u.cs.  $t^{\alpha} (t^{\alpha})' = 1, t^l (t^k)' = 0, k \neq l$ .



・ロト ・ 日 ・ ・ ヨ ・ ・ ヨ ・ ・

Motivation Criterion Remark on rotation in PCA

- Possible rotation in PCA as in FA
- Must apply *T* to the "good" matrices → keep property of non correlation of components
- One must not write :

$$\tilde{X} = \overbrace{M^{-1/2}U\Lambda^{1/2}}^{\Psi_q} \mathbf{TT'} V'$$

but :

$$\tilde{X} = \overbrace{M^{-1/2}U}^{F_q} TT' \Lambda^{1/2} V'$$

 $\rightarrow$  rotation of standardized principal components = factors of FA



- Part of scientific project of French Ministry of Ecology
- Improve air quality : identification + quantification of pollution sources
- Three steps process :
  - **Collecting** of fine particles in French urban site (Anglet) : n = 61 samples (Dec. 2005)
  - **Measurements** of p = 16 chemical species : PIXE method (Particle Induced X-ray Emission)

- **Statistical treatment** of  $X = (x_i^j)_{n,p}$ : concentration of *j*th chemical compound in *i*th sample



- FA model with PCA estimation
- Choice of number of factors : q = 5
- Orthogonal rotation of factors with varimax criterion
- Detection of groups of correlated chemical compounds
- Identification of air pollution sources



A three steps process Statistical treatment Some confirmation of the results

	f <sup>1</sup>	f <sup>2</sup>	f <sup>3</sup>	f <sup>4</sup>	f <sup>5</sup>			
Al2O3	0.67	-0.66	0.22	-0.19	0.02			
SiO2	0.65	-0.67	0.25	-0.20	-0.06			
Р	0.68	-0.64	0.22	-0.24	0.01			
SO4	0.59	0.45	-0.41	-0.22	0.18			
CI	-0.47	-0.28	0.19	0.68	0.34			
K	0.89	-0.15	-0.21	0.00	0.28			
Ca	0.64	-0.41	0.10	0.40	-0.19			
Mn	0.38	0.78	0.18	0.09	-0.25			
Fe2O3	0.79	0.32	0.04	0.27	-0.36			
Cu	0.80	0.25	-0.10	0.23	-0.36			
Zn	0.35	0.66	0.59	-0.07	0.25			
Br	0.75	-0.18	-0.13	0.35	0.36			
Pb	0.43	0.66	0.52	-0.09	0.30			
C-Org	0.60	0.30	-0.61	-0.02	0.22			

Table: Loading matrix before rotation  $A_5$ 

Lots of intermediate values  $\rightarrow$  difficult association variables/factors



A three steps process Statistical treatment Some confirmation of the results

	Ĭ <sup>1</sup>	Ĭ2	<i>f</i> 3	Ĭ <sup>4</sup>	ĭ⁵		
Al2O3	0.98	0.09	-0.04	0.07	-0.04		
SiO2	0.98	0.01	-0.06	0.10	-0.07		
Р	0.97	0.09	-0.02	0.07	-0.09		
SO4	-0.03	0.77	0.25	0.18	-0.35	Factor	Possible source
CI	-0.15	-0.27	-0.14	-0.18	0.88	Factor1	Soil dust
К	0.60	0.72	0.11	0.23	0.03	Factor2	Combustion
Ca	0.61	0.09	-0.11	0.56	0.27	Factor3	Industry
Mn	-0.28	0.12	0.60	0.58	-0.24	Factor4	Vehicle
Fe2O3	0.20	0.28	0.29	0.85	-0.11	Factor5	Sea
Cu	0.21	0.36	0.16	0.82	-0.15	Table: P	ollution sources
Zn	-0.03	0.05	0.98	0.13	-0.04	rabier i	
Br	0.49	0.62	0.10	0.28	0.39		
Pb	0.00	0.16	0.97	0.13	-0.05		
C-Org	-0.02	0.89	0.02	0.22	-0.16		

Table: Rotated loading matrix  $\check{A}_5$ 



17/22

Applied Stochastic Models and Data Analysis, Crete, 2007

- Rotated factors : columns of  $\breve{F}_5 = F_5 \times T = (\breve{f}_i^k)_{n,5}$
- Score *f*<sup>k</sup><sub>i</sub> = "relative" contribution of *k*th source to *i*th sample
- Confrontation of rotated factors with external parameters :
  - Meteorological data (wind, temperature, ...)
  - Periodicity day/night
  - Other chemical coumpounds



A three steps process Statistical treatment Some confirmation of the results



Figure: Evolution of vehicle pollution source

Stronger contribution during day than night  $\rightarrow$  confirmation of cars pollution source



A three steps process Statistical treatment Some confirmation of the results



Figure: Evolution of heating pollution and temperatures

Middle of period :

Increase in contribution Decrease in temperature

ightarrow confirmation of heating pollution

Applied Stochastic Models and Data Analysis, Crete, 2007

A three steps process Statistical treatment Some confirmation of the results



Figure: Map of sampling site and correlations wind/sea pollution

Strong correlation to N-W wind Location toward Atlantic Ocean

 $\rightarrow$  validation of sea pollution



A three steps process Statistical treatment Some confirmation of the results

## Conclusion

- FA (with PCA estimation) followed by rotation : identification of five sources of particulate emission
- No prior knowledge : number, composition ?
   → more complex sampling site
- Sources quantification : percentage of total fine dust mass of each source (JDS, Angers, 2007)

