#### K. Gibert (1,2)

(1)Department of Statistics and Operation Research

(2) Knowledge Engineering and Machine Learning group Universitat Politècnica de Catalunya, Barcelona

> Master Oficial en Enginyeria Informàtica Universitat Politècnica de Catalunya

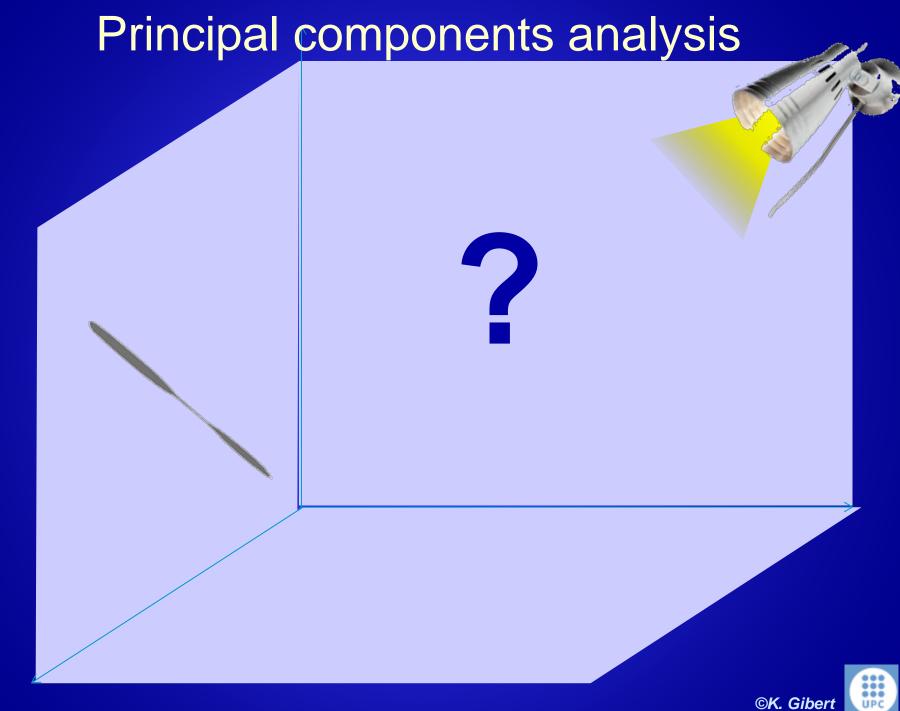
Find the isomorph transformation from original space

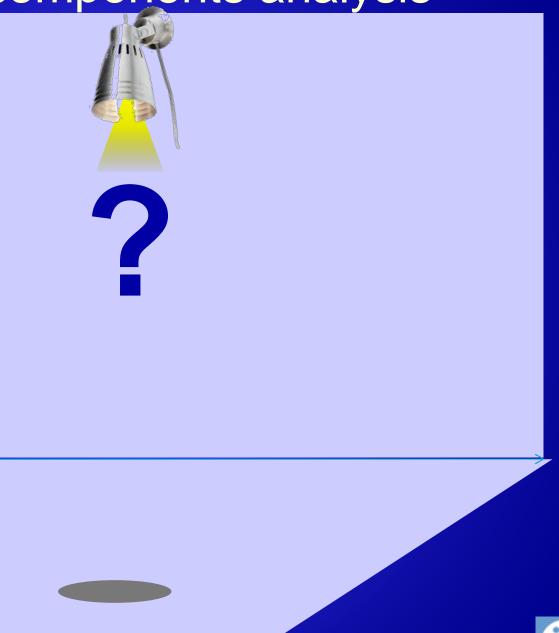
keeps the adjacency relationships among variables

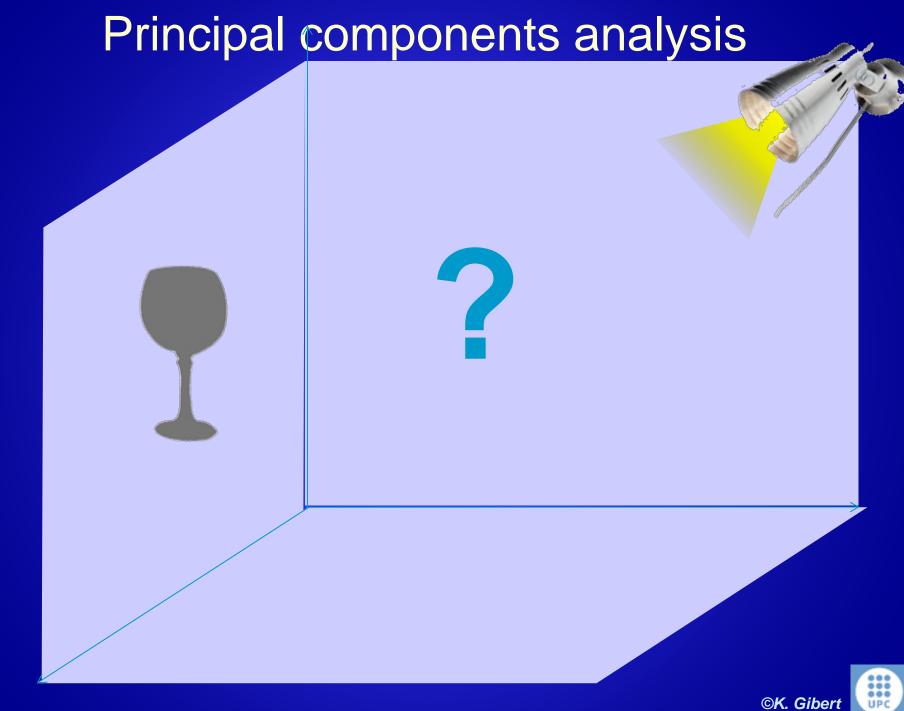
- Results expressed in a ficticious space
- Might produce interpretation problems
- Methods
  - PCA (Principal components analysis)
  - Simple correspondence analysis
  - Multiple correspondence analysis

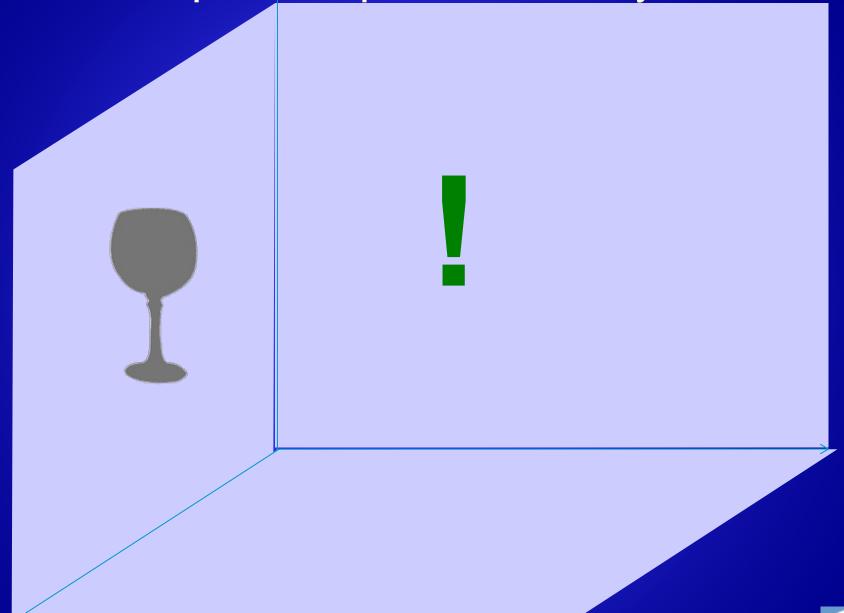
- Principal Components Analysis
  - Only numerical variables
  - Find the most informative projection planes (factorial planes)

Example "Copas"



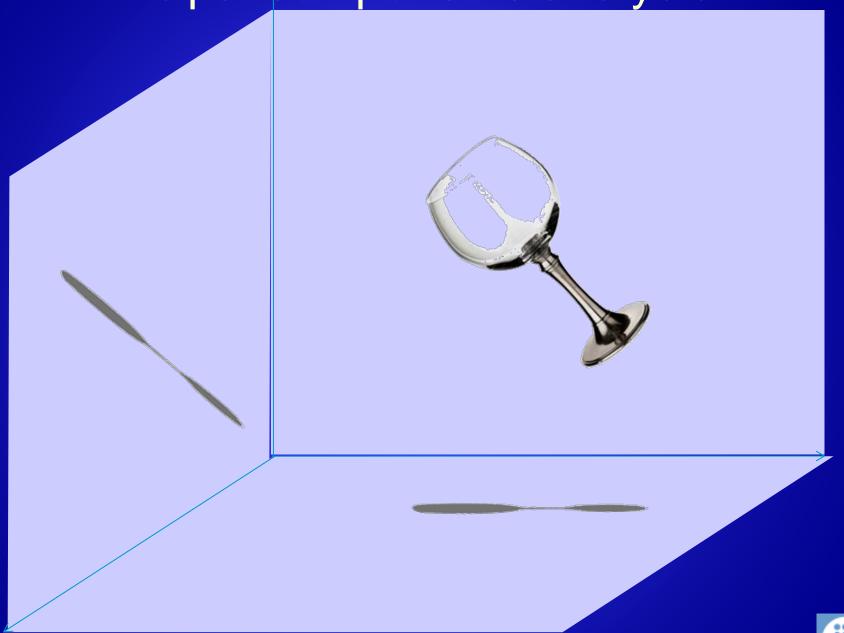




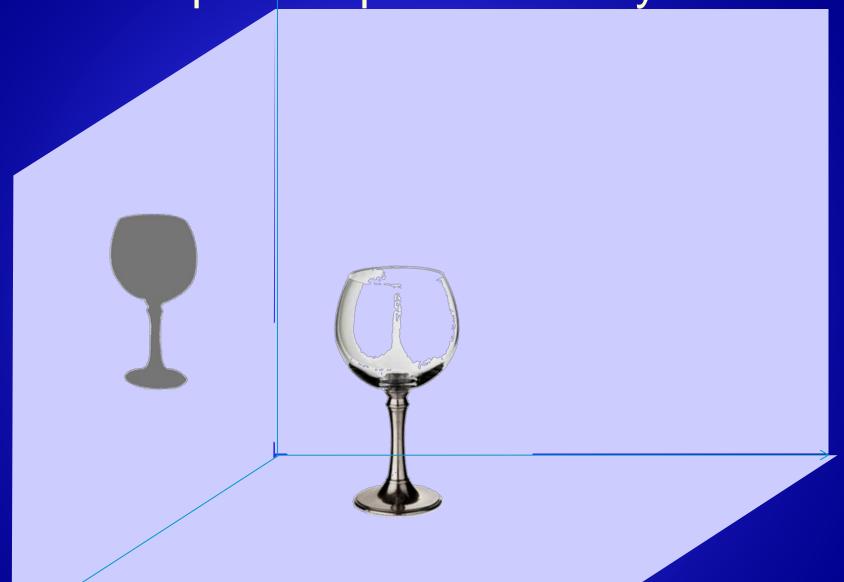


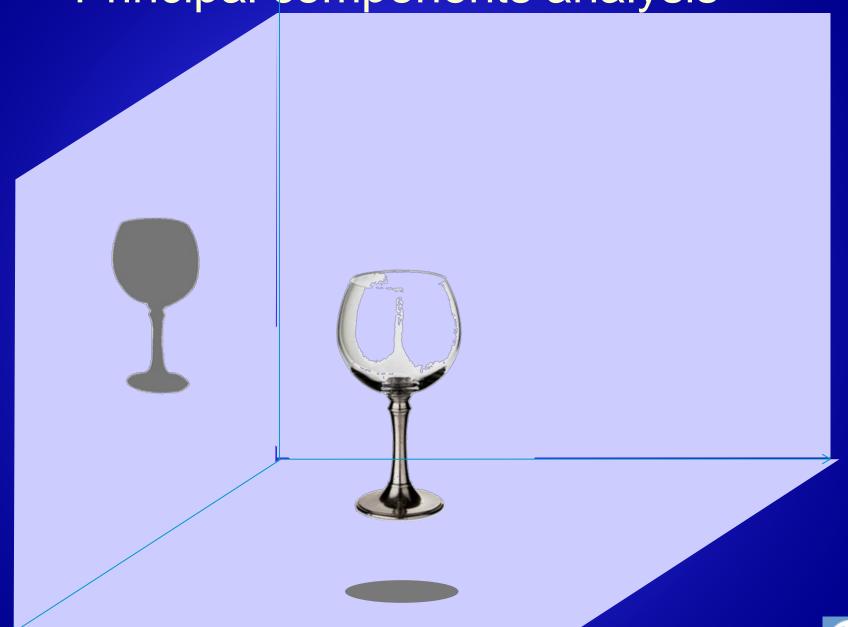


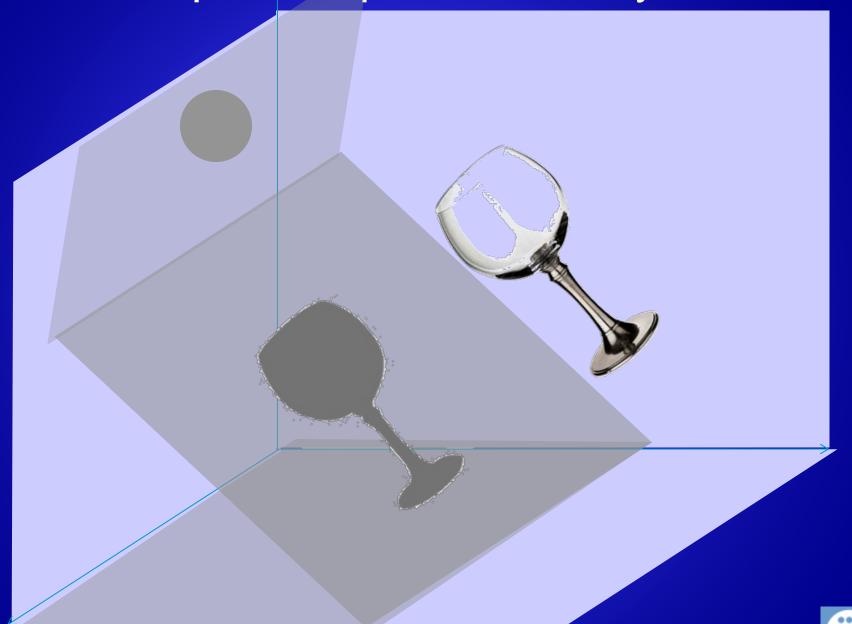


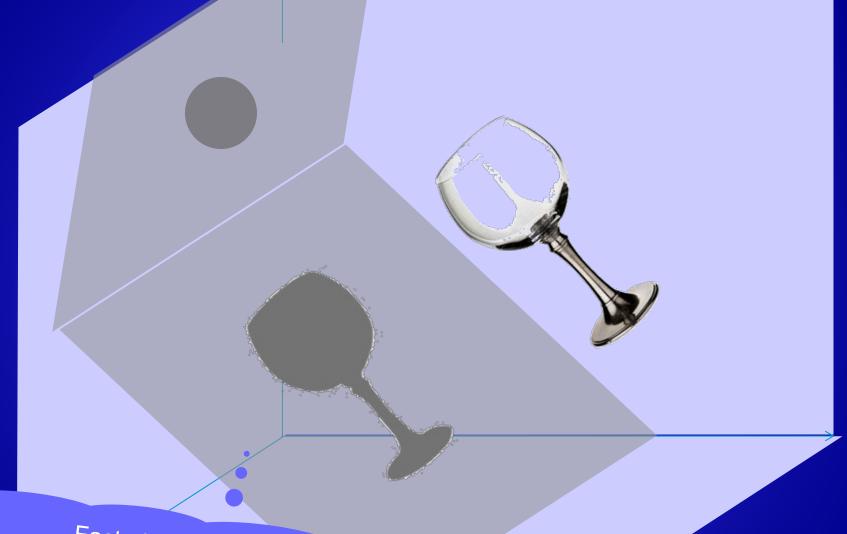








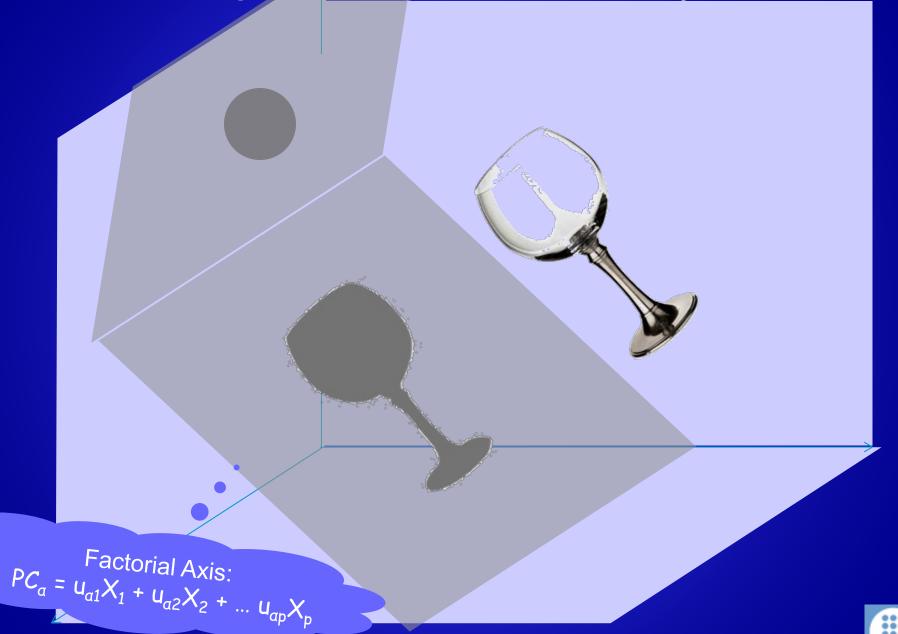




Factorial Plane: 2 factorial axes

Factorial axis: Linear combination of original variables

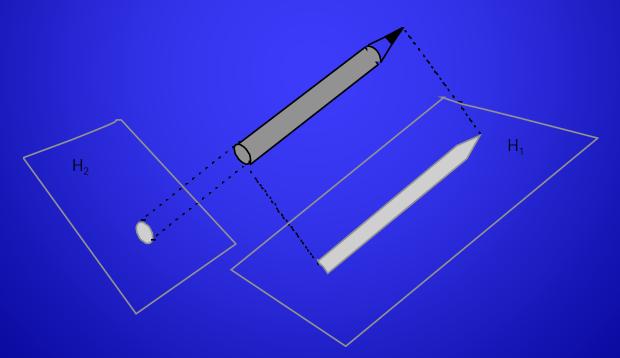




## Purpose:

 To project the cloud of points upon a subspace (plane) retaining as much original cloud information.

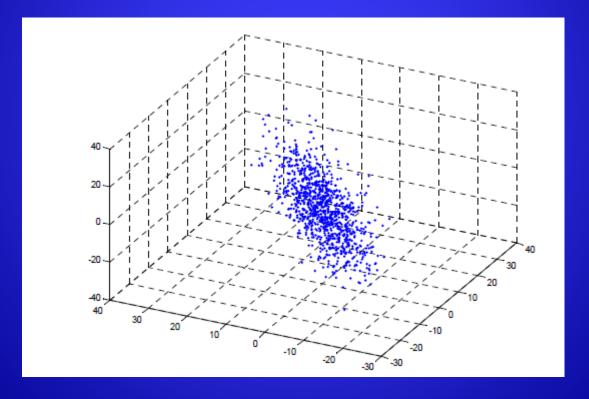
(see video)



Course DM: Multivariate Visualisation. T. Aluja



•Find the most informative projection planes of data cloud (factorial planes)

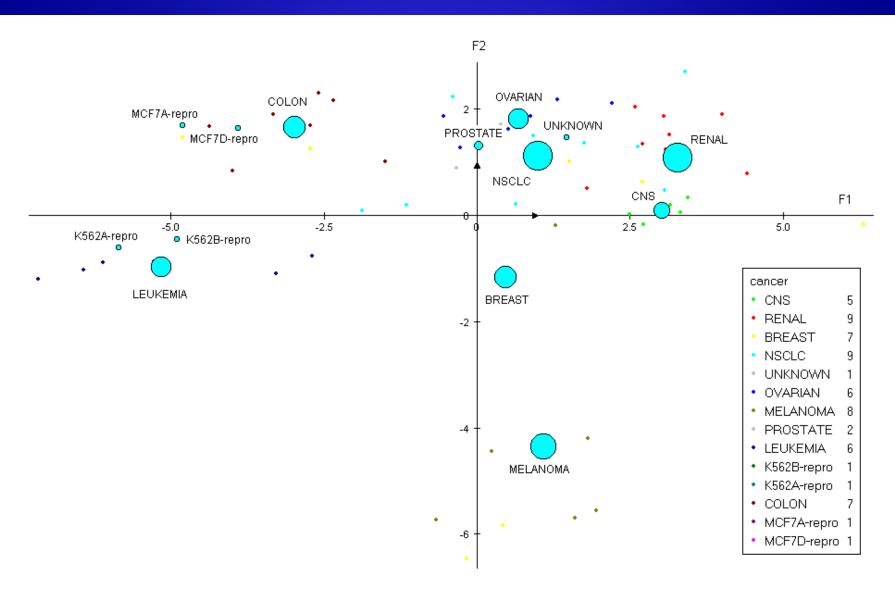


- Output: K factors rotating original X variables
- Factors: Linear combinations of original variables

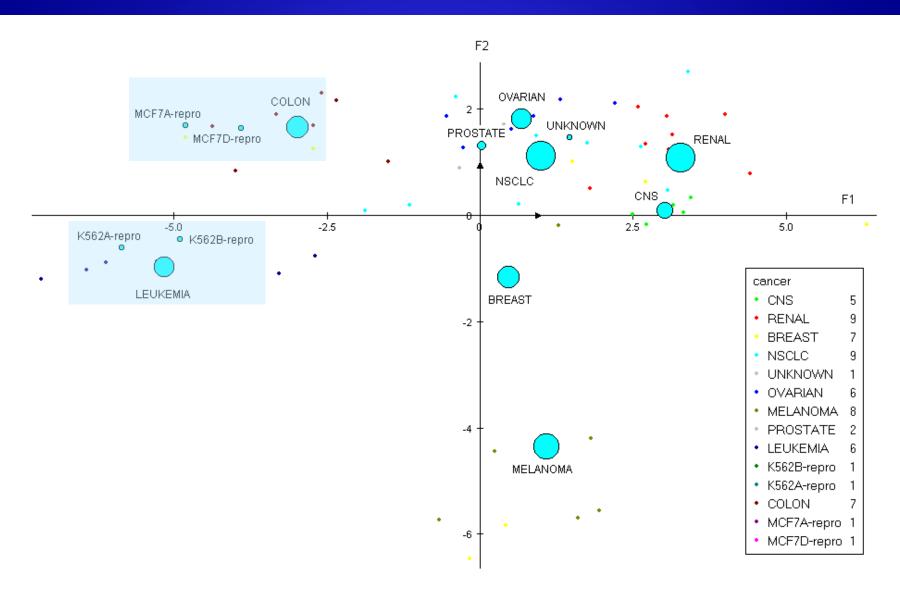
#### Several uses:

 As an associative data mining method: analyze relationships among variables
 Project variables and modalities and find associations

#### Microarray data: 64 cancers 6830 gen cromotografy



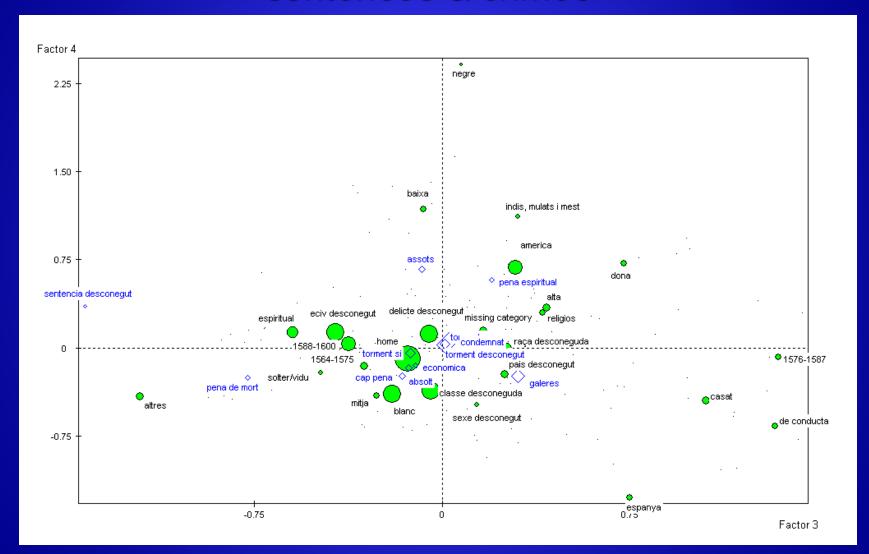
#### Microarray data: 64 cancers 6830 gen cromotografy





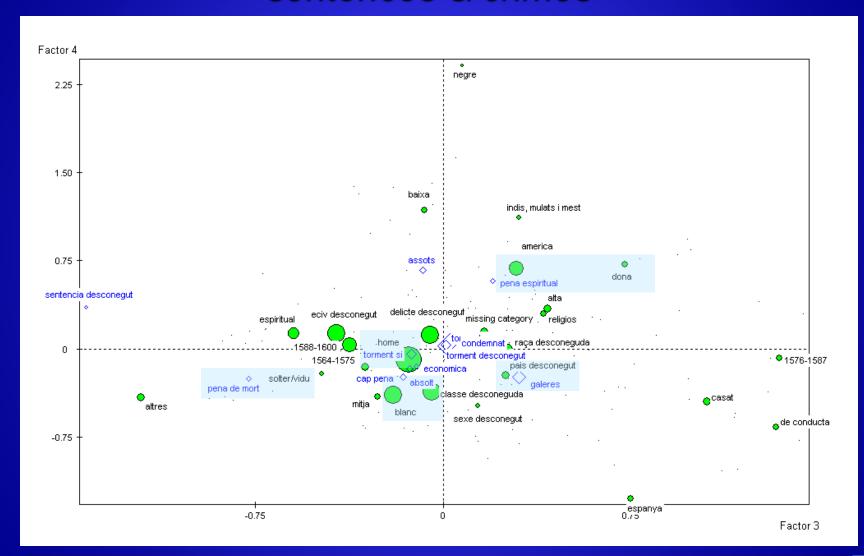
## Spanish inquisition 1567-1600

### sentences & crimes



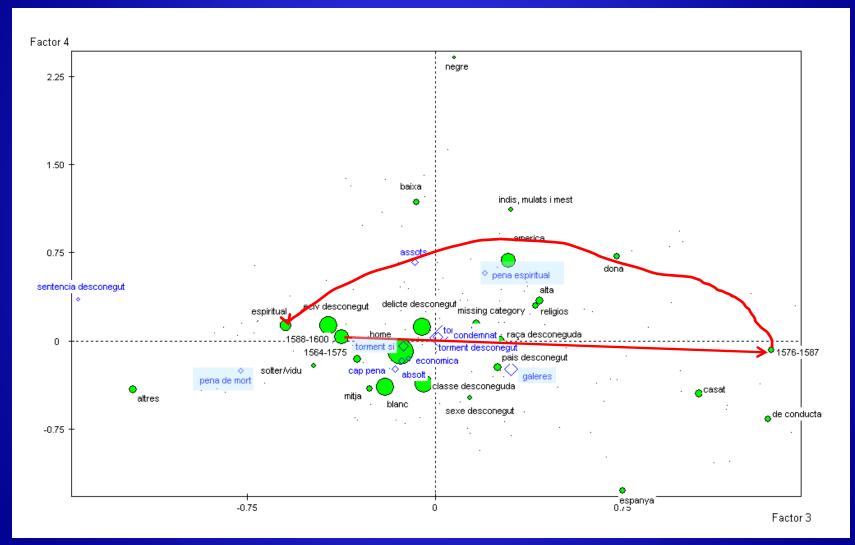
## Spanish inquisition 1567-1600

## sentences & crimes



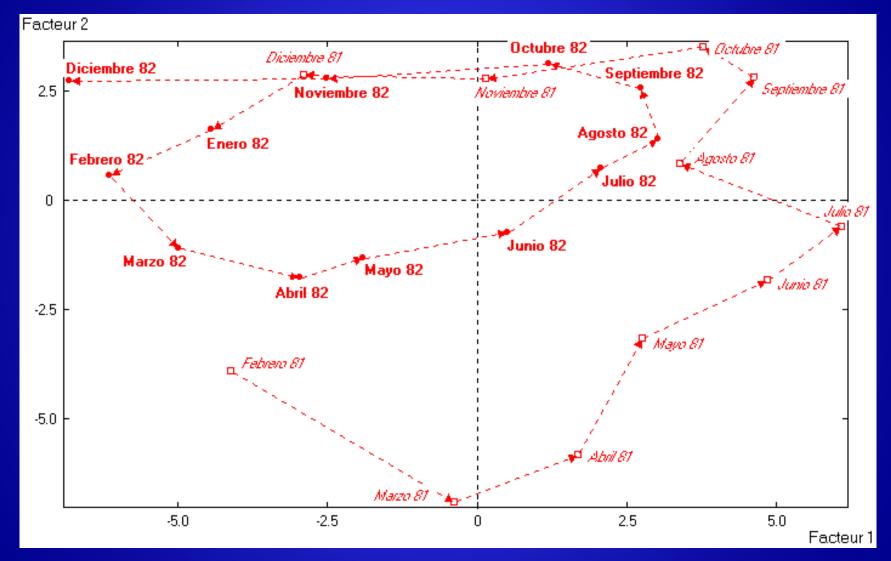
## Spanish inquisition 1564-1600

### sentences & crimes





#### Monitoring of the inner temperatures of Lascaux cave (France)



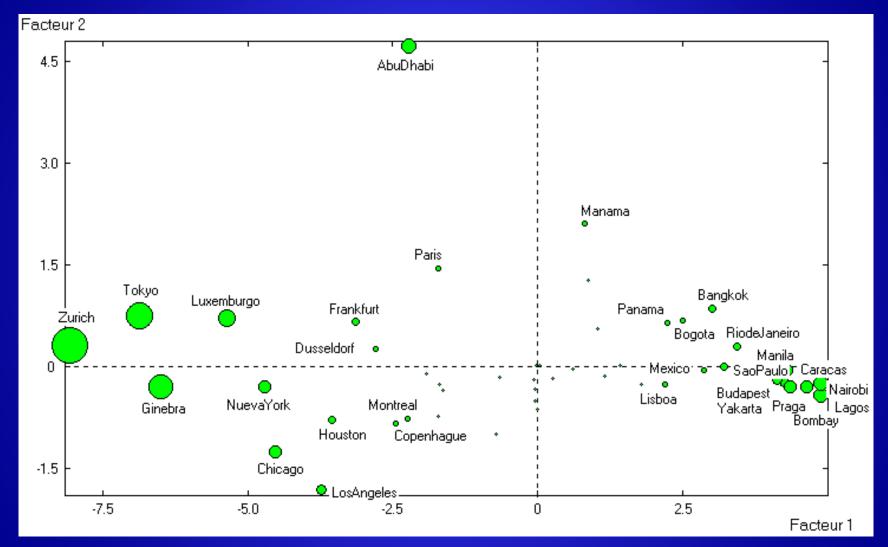
- Output: K factors rotating original X variables
- Factors: Linear combinations of original variables

#### Several uses:

- As an associative data mining method to analyze relationships among variables
   Project variables and modalities and find associations
- As a preprocessing method for elicitation of latent variables
   Project active and illustrative variables/individuals on first/second factorial plane and interpret factors (find latent variables)



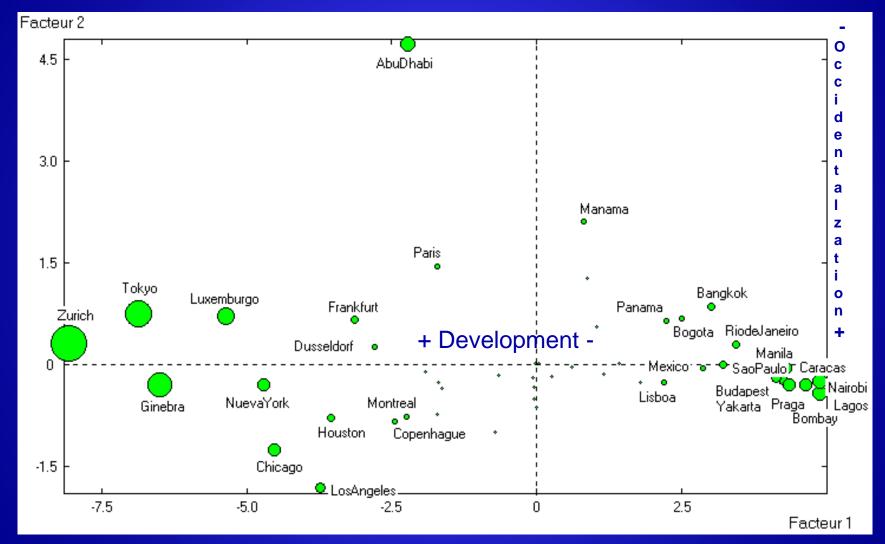
# Visualisation of international cities according their salaries. USB 1994.







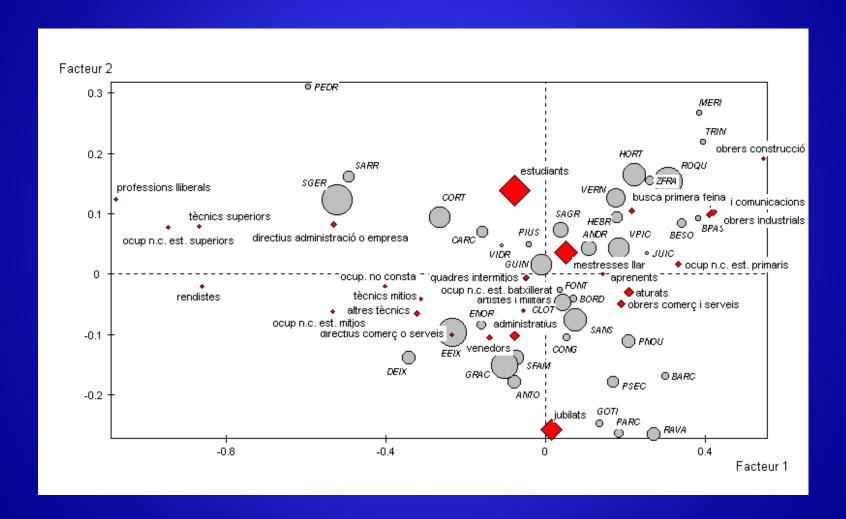
# Visualisation of international cities according their salaries. USB 1994.



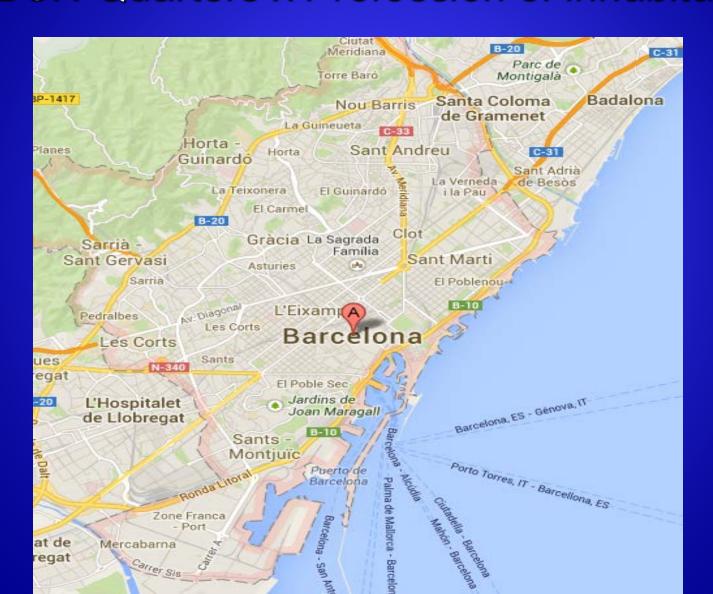




## Visualization of the table BCN Quarters x Profession of inhabitants

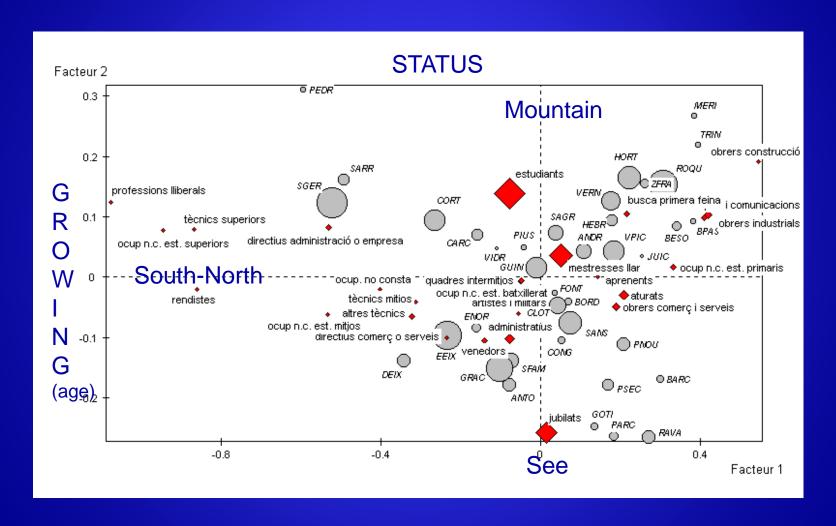


## Visualization of the table BCN Quarters x Profession of inhabitants





## Visualization of the table BCN Quarters x Profession of inhabitants



- Output: K factors rotating original X variables
- Factors: Linear combinations of original variables

#### Several uses:

- As an associative data mining method to analyze relationships among variables
   Project variables and modalities and find associations
- As a preprocessing method for elicitation of latent variables
   Project active and illustrative variables/individuals on first/second factorial plan and interpret factors (find latent variables)
- -As a preprocessing method for multidimensionality reduction

Data	Factorial Method
Continuous variables	PCA - Principal Component Analysis
Count variables	CA - (Simple) Correspondence Analysis
Categorical variables	MCA - Multiple Correspondence Analysis

## Principal Components Analysis

- Only numerical variables
- Find the most informative projection planes (factorial planes, maximize projected inertia)

#### Given <X,M,D>

- A data matrix X (nxp) centered
- A matrix of individuals weights D (nxn)
- Assume euclidean metrics to compare individuals (M= Ip)

#### Matrix M<sup>1/2</sup> X'DXM <sup>1/2</sup>

- Product of data with the two metrics
- Simetric,
- Semidefinite
- Catches relationships and opositions of data



Principal Components Analysis

M<sup>1/2</sup>X'DXM<sup>1/2</sup> catches well the data structures

$$Rang(M^{\frac{1}{2}}X'DXM^{\frac{1}{2}}) = r, r = rang(X)$$
 r pos

r positive vaps and p-r null vaps

Trace(
$$M^{\frac{1}{2}}X'DXM^{\frac{1}{2}}$$
)= $\sum_{\alpha=1}^{r}\lambda_{\alpha}$ 

 $(\lambda_{\alpha}$ , the r non null vaps)

$$M = I_D: M^{1/2}X'DXM^{1/2} = X'DX$$

X centered and D diagonal: X'DX= Cov(X)

Build variances and covariances matrix: X'DX

Diagonalize X'DX (i.e. solving the equation ) X'DXu=  $\lambda u$ provides eigen values  $\lambda_{\alpha}$  and

eigenvectors  $u_{\alpha} = (u_{\alpha 1} .... u_{\alpha p})$ 

Principal Components Analysis

```
Diagonalize X'DX (i.e. solving the equation ) X'DXu= \lambda_u (1) \det(X'DX-\lambda)=0 (find roots of characteristic polynomial) provides eigen values \lambda_{\alpha} (\alpha=1:r, r=rang(X)) substituting in (1) provides eigenvectors u_{\alpha}=(u_{\alpha 1}....u_{\alpha n})
```

 $u^{-1}X'DXu = \lambda$  is a diagonal matrix

(X'DX becomes diagonal when pre/post multiplied by u)

 $u^{-1}=u'$  in orthonormal basis:  $u'X'DXu=\lambda$ X'DX decompose in a product by a diagonal matrix X'DX =  $u\lambda u'$ 

 $X'DX = u\lambda u' = u\lambda^{1/2}\lambda^{1/2} u' = u\lambda^{1/2} \mathbb{I}\lambda^{1/2} u' = u\lambda^{1/2} u'u \lambda^{1/2} u' = A^{1/2}A^{1/2}$ X'DX decompose in a product of something by itself (A square root)

 $Trace(X'DX)=Trace(\lambda)$  (property of diagonalization)



Given <X,M,D>

Diagonalize covariances matrix (with centered data X'DX)

Get r eigen values  $\lambda_{\alpha}$  and sort decreasingly

$$\{\lambda_{\alpha}\}_{\alpha=1:r}$$
  $\lambda_{1} \geq \lambda_{2} \geq \lambda_{3} \geq \ldots \geq \lambda_{r}$ 

Corresponding eigenvectors  $u_{\alpha} = (u_{\alpha 1} \dots u_{\alpha p})$ 

$$|u_{\alpha}|=1$$
  
 $u_{\alpha}u_{\alpha'}=0$   
 $\{u_{\alpha}\}_{\alpha=1:r}$  orthonormal base for individuals

The subspace generated by  $\{u_{\alpha}\}_{\alpha=1:r}$  is the same as the subspace generated by the rows of X

Given <X,M,D>

In general *Diagonalize M<sup>1/2</sup>X'DXM<sup>1/2</sup>*)

Get r eigen values  $\lambda_{\alpha}$  and sort decreasingly (vaps are conserved!!!!)

$$\{\lambda_{\alpha}\}_{\alpha=1:r}$$
  $\lambda_{1} \geq \lambda_{2} \geq \lambda_{3} \leq \ldots \geq \lambda_{r}$ 

Corresponding eigenvectors  $u^*_{\alpha} = (u^*_{\alpha 1}....u^*_{\alpha p})^{n}$ 

by algebraic properties,  $u^*_{\alpha}$  can be found from  $u^*$ 

$$u^*_{\alpha} = M^{-1/2}u_{\alpha}$$

 $\{u^*_{\alpha}\}_{\alpha=1:r}$  orthonormal base for individuals

$$|u^*_{\alpha}|_{M}=1: u^{*'}_{\alpha}Mu^*_{\alpha}=u'_{\alpha}M^{-\frac{1}{2}}MM^{-\frac{1}{2}}u_{\alpha}=1$$
  
 $u^*_{\alpha}Mu^*_{\alpha'}=0: u^*_{\alpha}Mu^*_{\alpha'}=u'_{\alpha}M^{-\frac{1}{2}}MM^{-\frac{1}{2}}u_{\alpha'}=0$ 

Subspace generated by  $\{u^*_{\alpha}\}_{\alpha=1:r}$  = Subspace generated by X rows



Given <X,M,D>

Diagonalize Covariance matrix X'DX

Get r eigen values  $\lambda_{\alpha}$  and sort decreasingly

$$\{\lambda_{\alpha}\}_{\alpha=1:r}$$
  $\lambda_{1} \geq \lambda_{2} \geq \lambda_{3} \geq \ldots \geq \lambda_{r}$ 

Corresponding eigenvectors  $u_{\alpha} = (u_{\alpha 1} .... u_{\alpha p})$ 

for 
$$M = \mathbb{I}_p : u^*_{\alpha} = u_{\alpha}$$
; for  $M \neq \mathbb{I}_p : u^*_{\alpha} = M^{-1/2} u_{\alpha}$ 

 $\{u^*_{\alpha}\}_{\alpha=1:r}$  orthonormal base for individuals

 $u^*_{\alpha}$  are the principal factors of X : good rotation directions

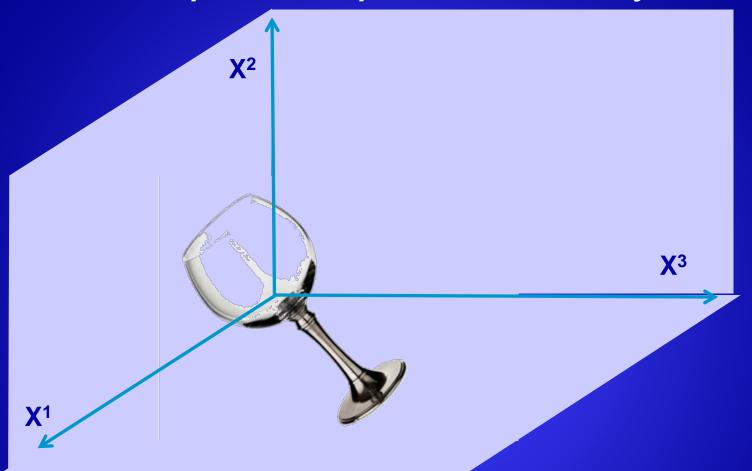
 $U^*=([u^*_1][u^*_2]....[u^*_r])$  is the basis for the projection space

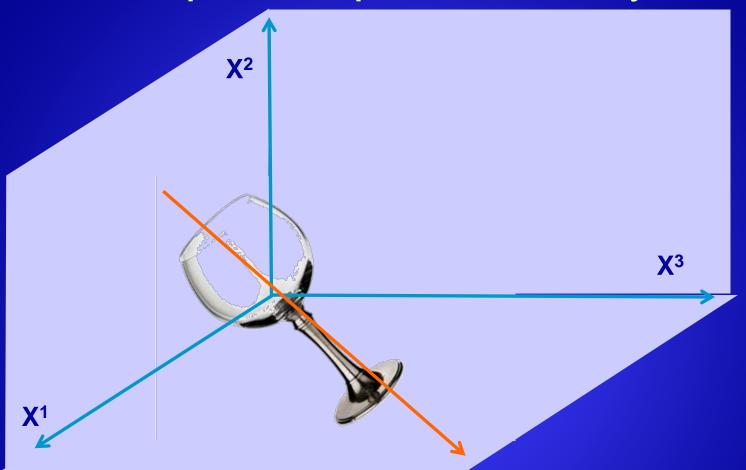


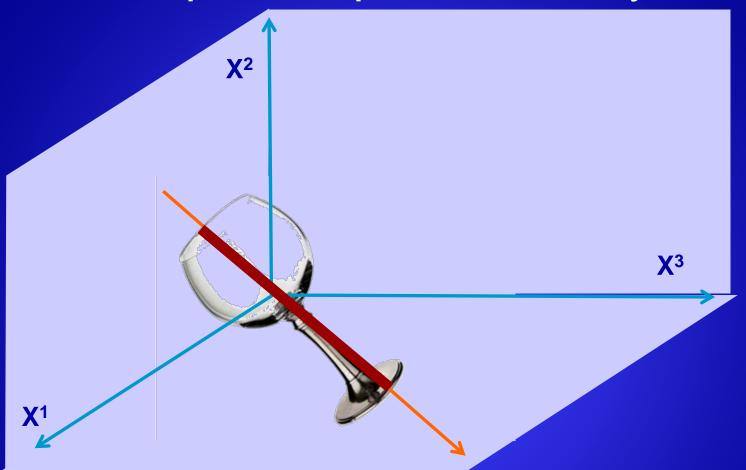
**Centering X** 

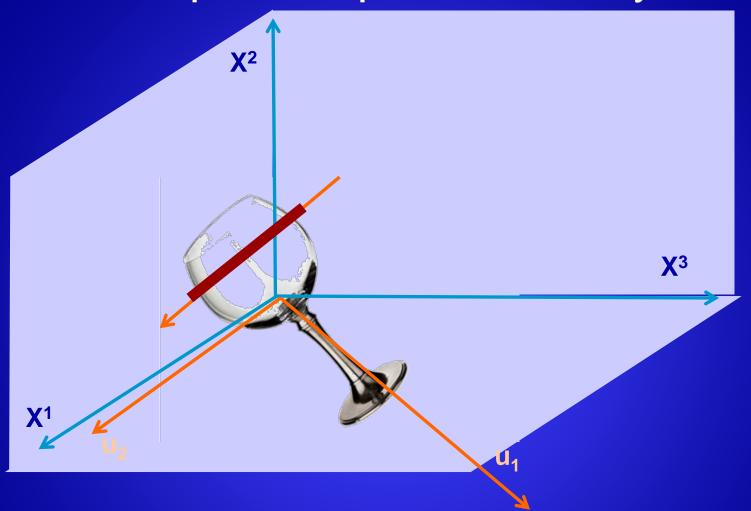
(0,0,0)

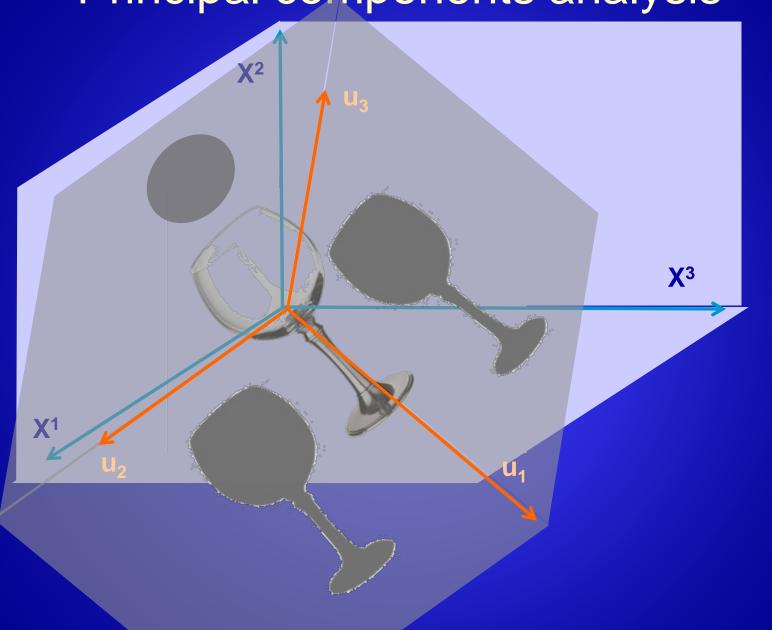


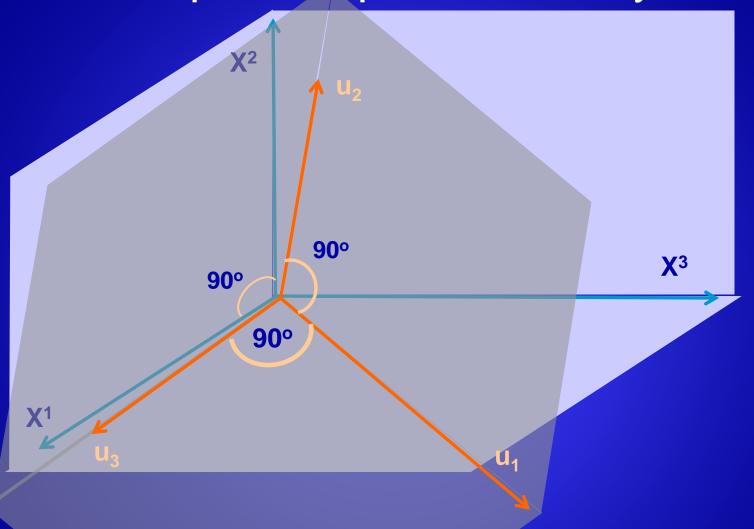


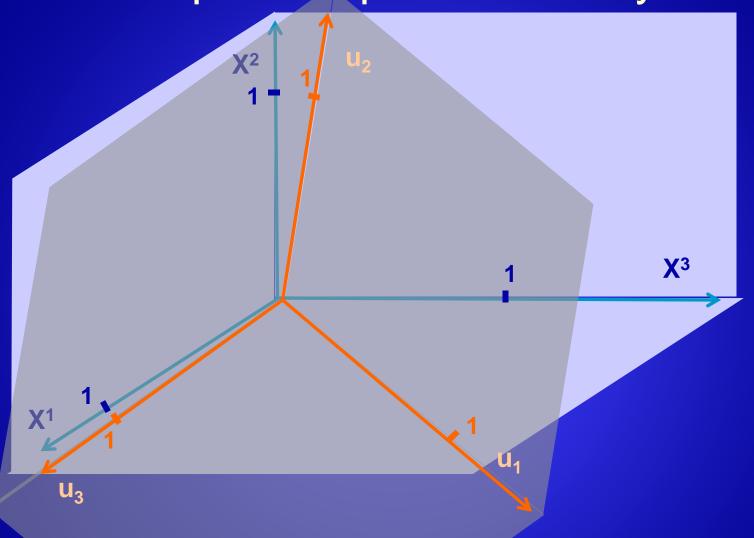












Given <X,M,D>

Diagonalize Covariance matrix X'DX

Get r eigen values  $\lambda_{\alpha}$  and sort decreasingly

$$\{\lambda_{\alpha}\}_{\alpha=1:r}$$
  $\lambda_{1} \geq \lambda_{2} \geq \lambda_{3} \geq \ldots \geq \lambda_{r}$ 

Corresponding eigenvectors  $u_{\alpha} = (u_{\alpha 1} .... u_{\alpha p})^{\prime}$ 

for 
$$M = \mathbb{I}_p : u^*_{\alpha} = u_{\alpha}$$
; for  $M \neq \mathbb{I}_p : u^*_{\alpha} = M^{-1/2} u_{\alpha}$ 

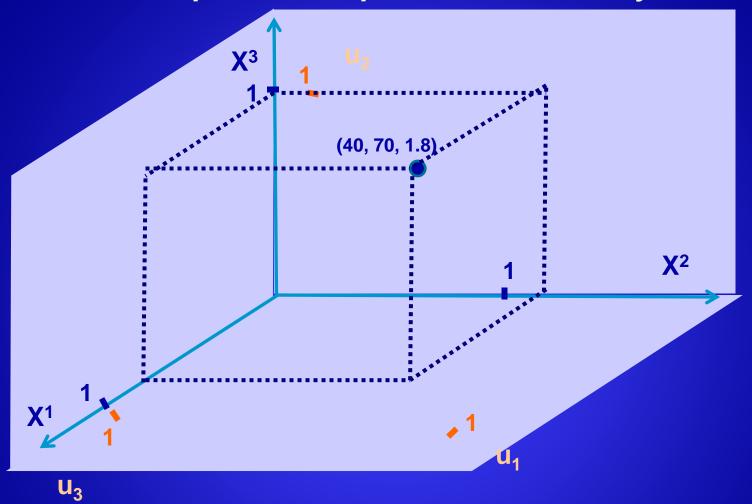
 $\{u^*_{\alpha}\}_{\alpha=1:r}$  orthonormal base for individuals

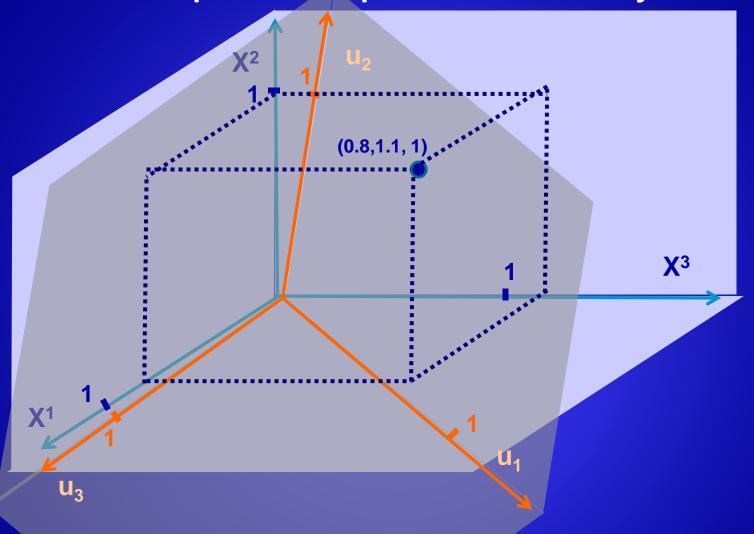
 $u^*_{\alpha}$  are the principal factors of X : good rotation directions

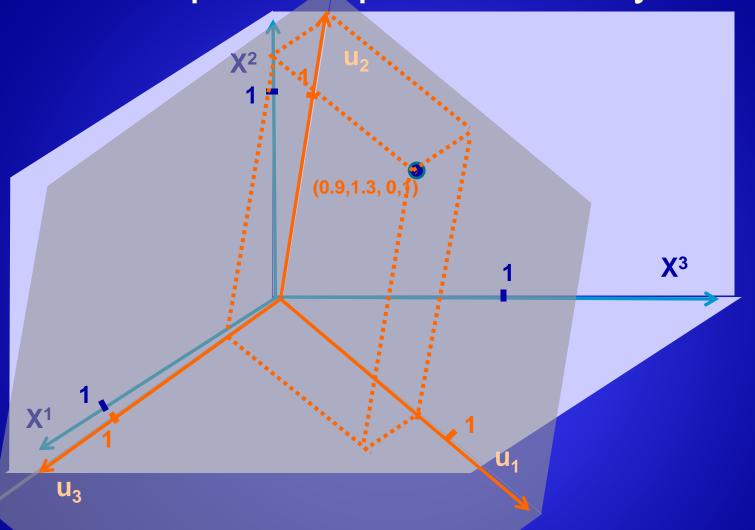
 $U^*=([u^*_1][u^*_2]....[u^*_r])$  is the basis for the projection space

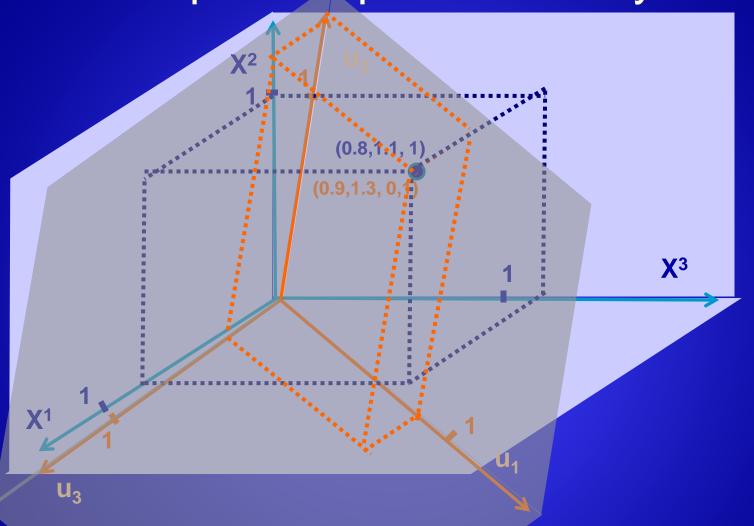
How is i expressed in rotated space?











Given <X,M,D>

Can we find coordinates in rotated space from original ones?

The projection matrix  $P = U^*_k U^{*'}_k M$ 

Projection of a single individual:  $Pr(i) = U_k^* U_k^{*\prime} M x_i$ 

Projection of all individuals:  $Pr(X) = U_k^* U_k^{*'} M X'$ 

Get a matrix with projections in ROWS:  $Pr(X)' = XMU^*_k U^{*'}_k$ 

Projections expressed in original vectorial space

The best possible projection over k dimensions



Given <X,M,D>

Matrix  $XMU^*_k U^{*'}_k$  provides the best possible k-projection of X

Silver-Smidth norm: 
$$|X|^2_{MD} = \sum_{\alpha=1}^r \lambda_{\alpha}$$

Measures variability, information contained in X

Property: 
$$||XMU^*_{k}U^{*'}_{k}||^2_{MD} = ||X||^2_{MD}$$

#### Any other k-projection of X

- Provides smallest values of Silver-Smidth norm
- Has less variability
- Keeps smallest information from X

Given <X,M,D>

Diagonalize covariances matrix (with centered data)

```
eigenvectors u_{\alpha} = (u_{\alpha 1} .... u_{\alpha p}) (direction of factor \alpha, \alpha = 1:p)
u_{p \alpha} : \text{contribution of variable p to the factor } \alpha
(u_{1} ..... u_{k}) \text{ ortonormal}
```

eigen values  $\lambda_k$  (quantity of information converved by factor k) (Projected inertia)

$$\{\lambda_{\alpha}\}_{\alpha=1:r}$$
  $\lambda_{1} \geq \lambda_{2} \geq \lambda_{3} \leq \ldots \geq \lambda_{r}$ 

 $\Sigma_{\forall \alpha} \lambda_{\alpha}$  = Total inertia of X (information in data)

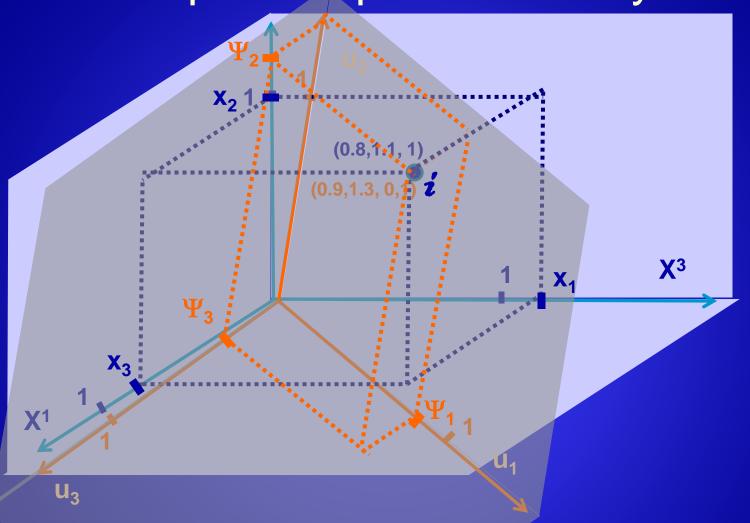
Given <X,M,D>

eigenvectors  $u_{\alpha} = (u_{1\alpha} .... u_{p\alpha})$  (direction of factor k)  $u_{p\alpha}$ : contribution of variable p to the factor  $\alpha$ 

eigen values  $\lambda_{\alpha}$  (quantity of information converved by factor  $\alpha$ ) (Projected inertia)

$$\lambda_1 \geq \lambda_2 \geq \lambda_3 \geq \ldots \geq \lambda_r$$

 $\Sigma_{\forall \alpha}$   $\lambda_{\alpha}$  = Total inertia of X (information of data)



Given <X,M,D>

$$i = (x_{1i}, ..., x_{pi})$$

Points in projected space:  $i = (\Psi_{1i}, ..., \Psi_{\alpha i}, ..., \Psi_{ri})$  (often r=p)

$$\Psi_{\alpha i} = x_{1i} u_{1\alpha} + x_{2i} u_{2\alpha} + \dots + x_{pi} u_{p\alpha} \qquad \qquad \psi_{\alpha} = X u_{\alpha}$$

Then 
$$\Psi'_{\alpha}$$
  $D\Psi_{\alpha} = \lambda_{\alpha}$ 

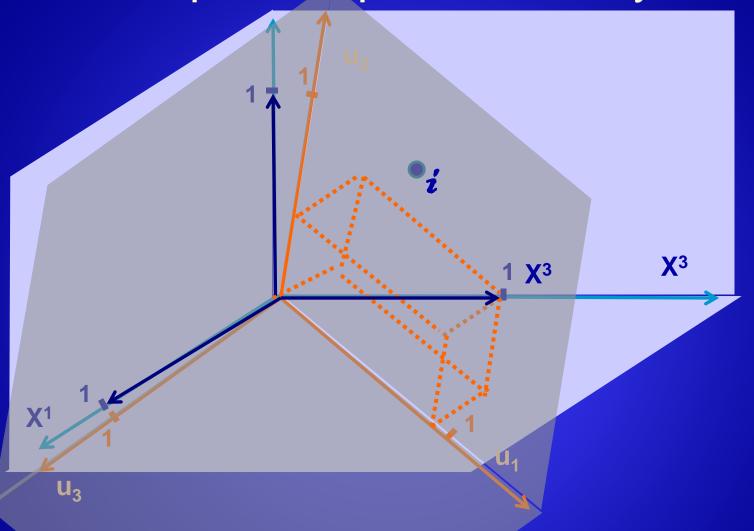
Illustrative points z also projectable

$$\Psi_{\alpha z} = x_{1z} u_{1\alpha} + x_{2z} u_{2\alpha} + \dots + x_{pz} u_{p\alpha}$$

Factors are linear combinations of original variables

Original variables project as VECTORS over factorial space angle and lenght important





- Principal Components Analysis
  - Output: K factors rotating original X variables
  - Factors: Linear combinations of original variables

#### Several uses:

- As an associative data mining method to analyze relationships among variables
   Project variables and modalities and find associations
- As a preprocessing method for elicitation of latent variables
   Project active and illustrative variables/individuals on first/second factorial plane and interpret factors (find latent variables)
- -As a preprocessing method for multidimensionality reduction

Select more informative factors  $\kappa << p$  (accumulate 80% inertia)

Reduce data matrix to selected factors

Alternative, keep variables mainly contributing to selected factors

(smaller angles with factorial axis)

©K. Gibert

Given <X,M,D>

Diagonalize Covariance matrix X'DX

Get r eigen values  $\lambda_{\alpha}$  and sort decreasingly

$$\{\lambda_{\alpha}\}_{\alpha=1:r}$$
  $\lambda_{1} \geq \lambda_{2} \geq \lambda_{3} \geq \ldots \geq \lambda_{r}$ 

Corresponding eigenvectors  $u_{\alpha} = (u_{\alpha 1} .... u_{\alpha p})$ 

for 
$$M = \mathbb{I}_p : u^*_{\alpha} = u_{\alpha}$$
; for  $M \neq \mathbb{I}_p : u^*_{\alpha} = M^{-1/2} u_{\alpha}$ 

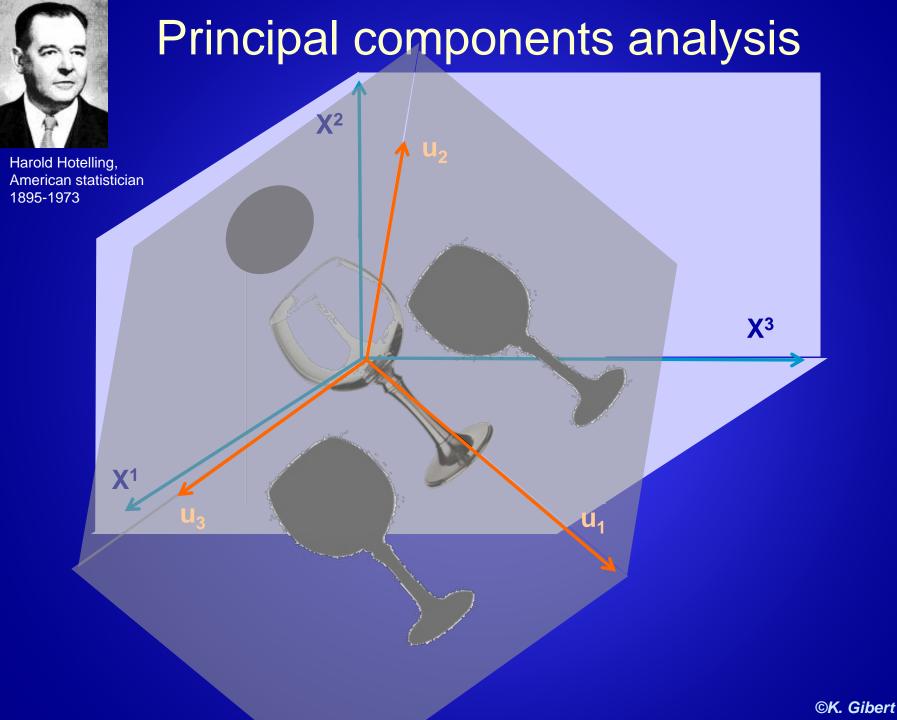
 $\{u^*_{\alpha}\}_{\alpha=1:r}$  orthonormal base for individuals

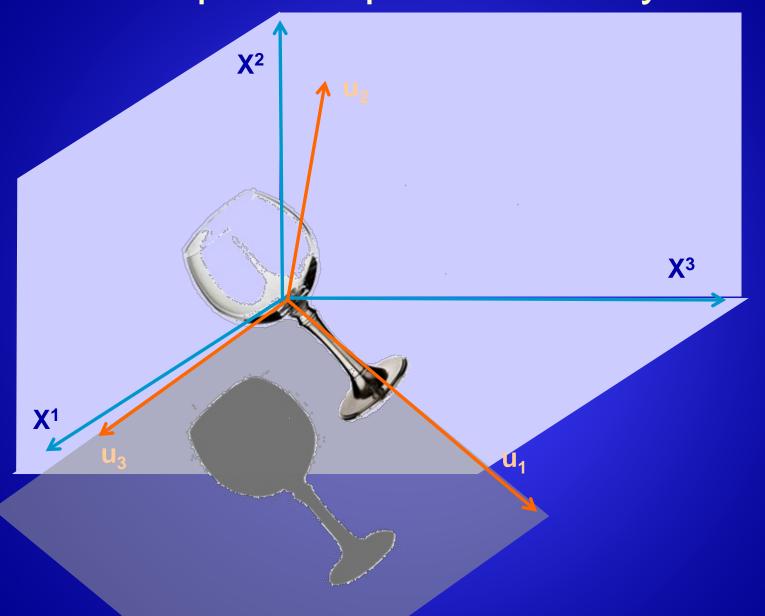
 $u^*_{\alpha}$  are the principal factors of X : good rotation directions

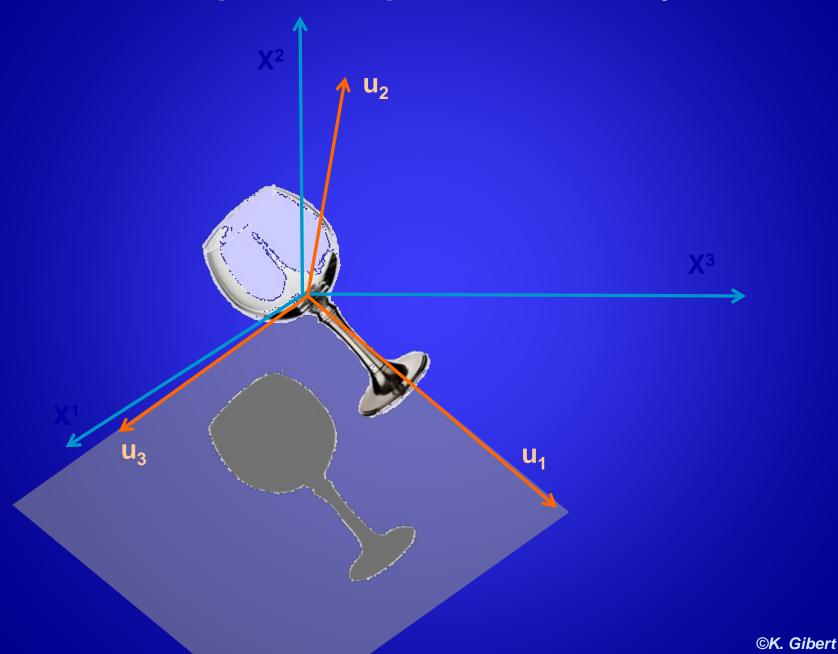
 $U^*=([u^*_1][u^*_2]....[u^*_r])$  is the basis for the projection space

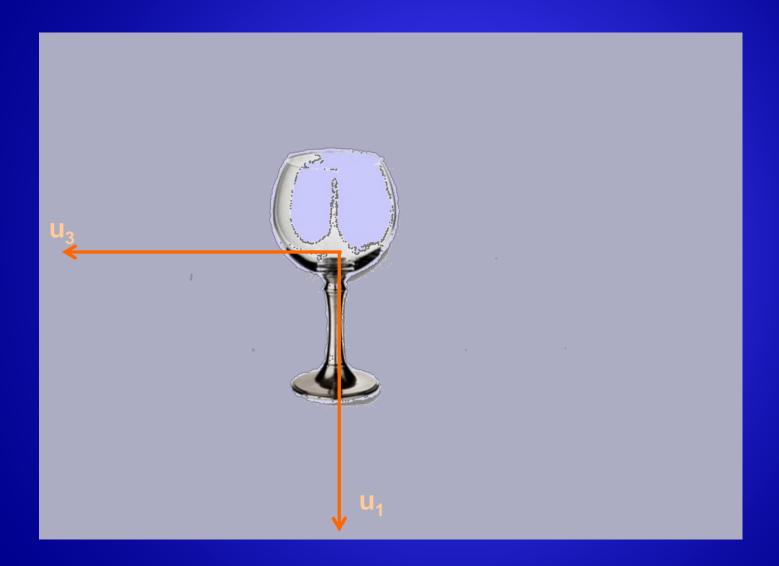
 $U_k^*=([u_1^*][u_2^*]....[u_k^*])$  is the basis for projecting in first k dimensions(k<r)



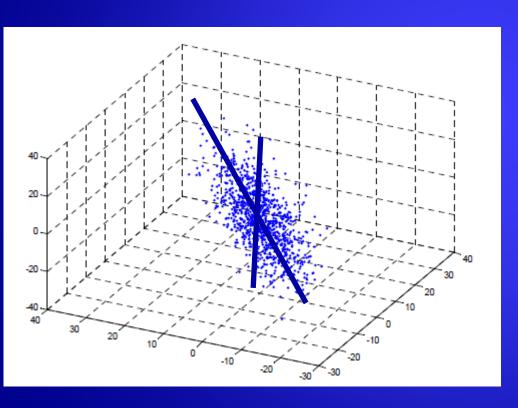


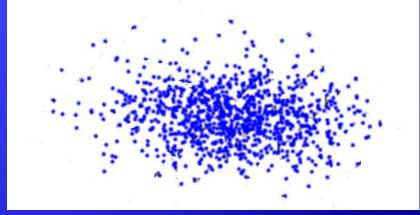




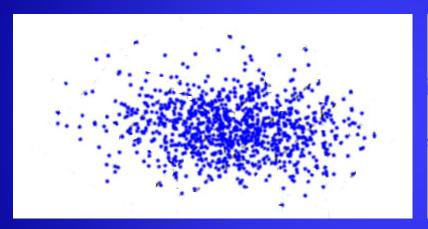


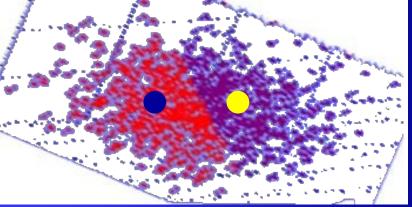
•Find the most informative projection planes of data cloud (factorial planes)



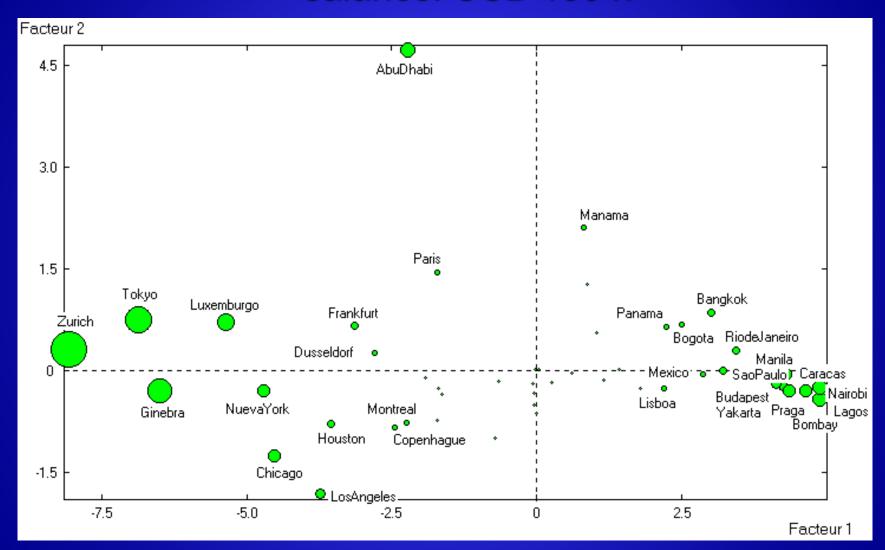


•Introduce qualitative information (projecting modalities)





# Visualisation of international cities according their salaries. USB 1994.





# Visualisation of international cities according their salaries. USB 1994.

