# Factorial Methods

## 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*
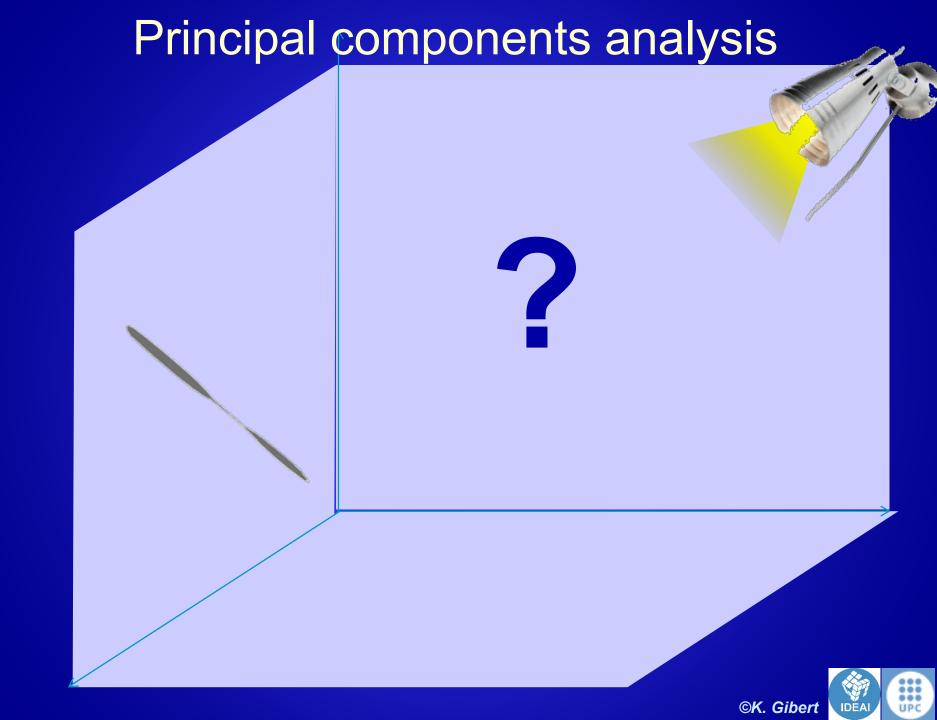
# Factorial Methods

- Find the isomorph transformation from original space
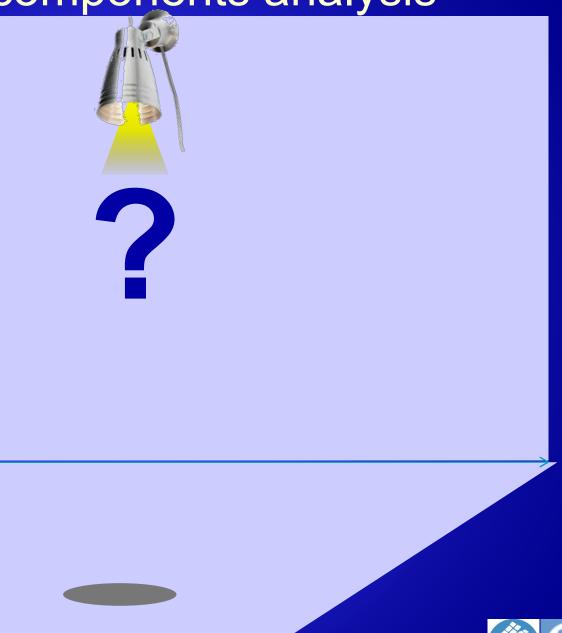
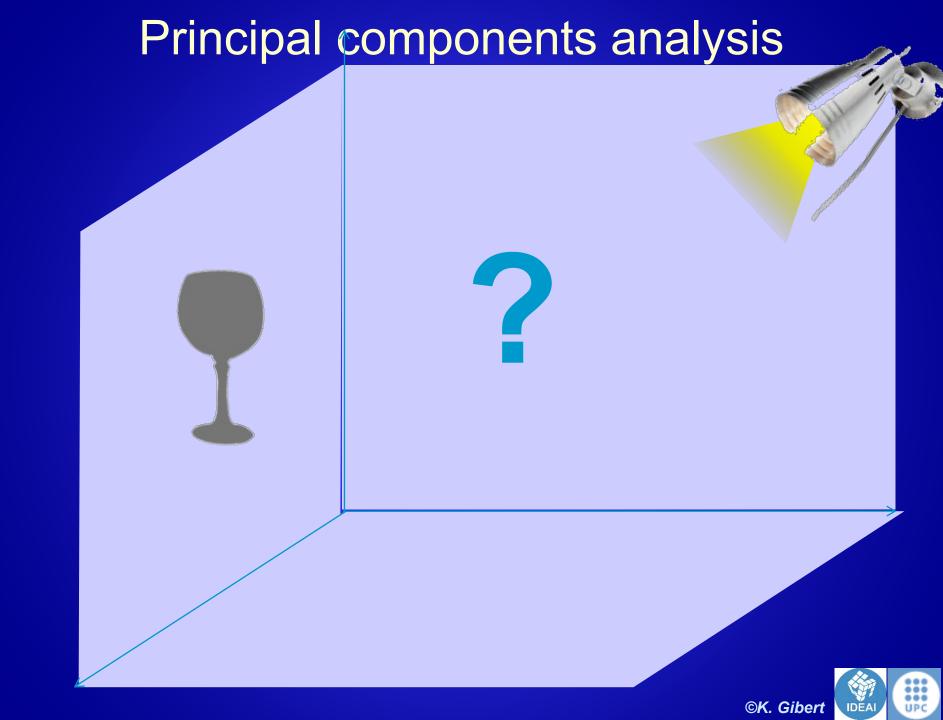  *keeps the adjacency relationships among variables*

- Results expressed in a ficticious space

- Might produce interpretation problems

- Methods
  - PCA (Principal components analysis)
  - Simple correspondence analysis
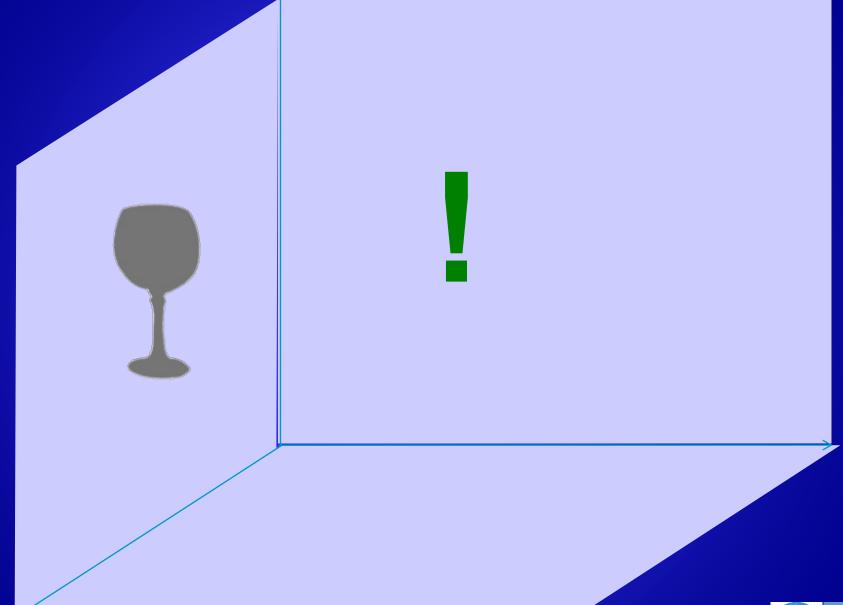  - Multiple correspondence analysis

©*K. Gibert*

# Factorial Methods

- **Principal Components Analysis**

  - Only numerical variables

  - Find the most informative projection planes

    *(factorial planes)*

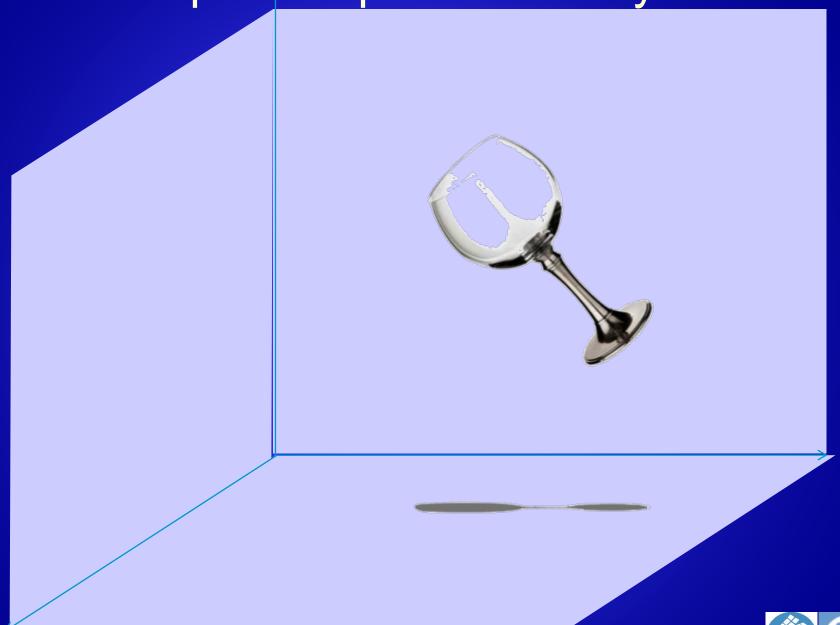  Example "Copas"

# Principal components analysis



?

# Principal components analysis



?

# Principal components analysis

?

# Principal components analysis

# Principal components analysis

# Principal components analysis

# Principal components analysis

# Principal components analysis

# Principal components analysis

# Principal components analysis

# Principal components analysis

# Principal components analysis

Factorial Plane: 2 factorial axes

Factorial axis: Linear combination of original variables

# Principal components analysis



Factorial Axis:
$$PC_\alpha = u_{\alpha 1}X_1 + u_{\alpha 2}X_2 + \ldots u_{\alpha p}X_p$$

©*K. Gibert*

# Principal components analysis

- 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

©K. Gibert

# Principal components analysis

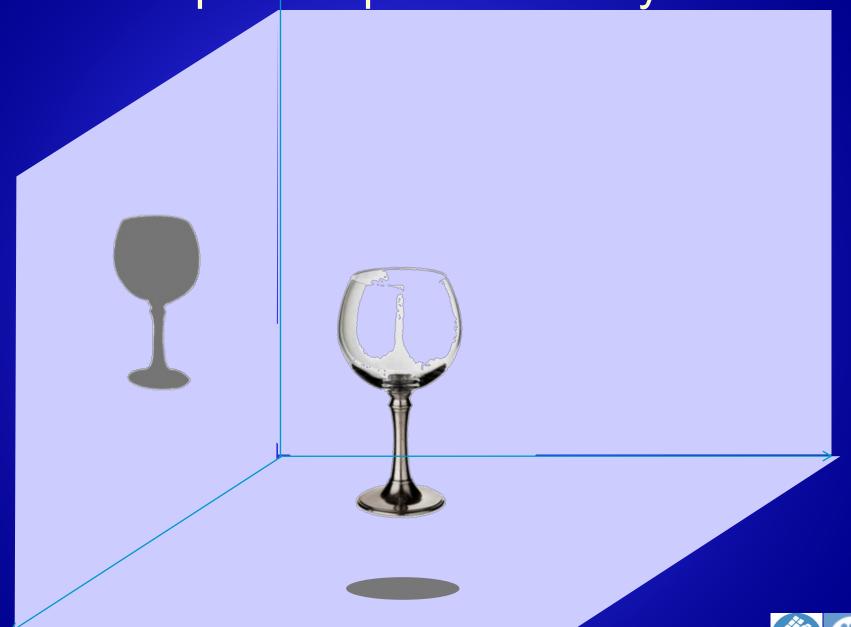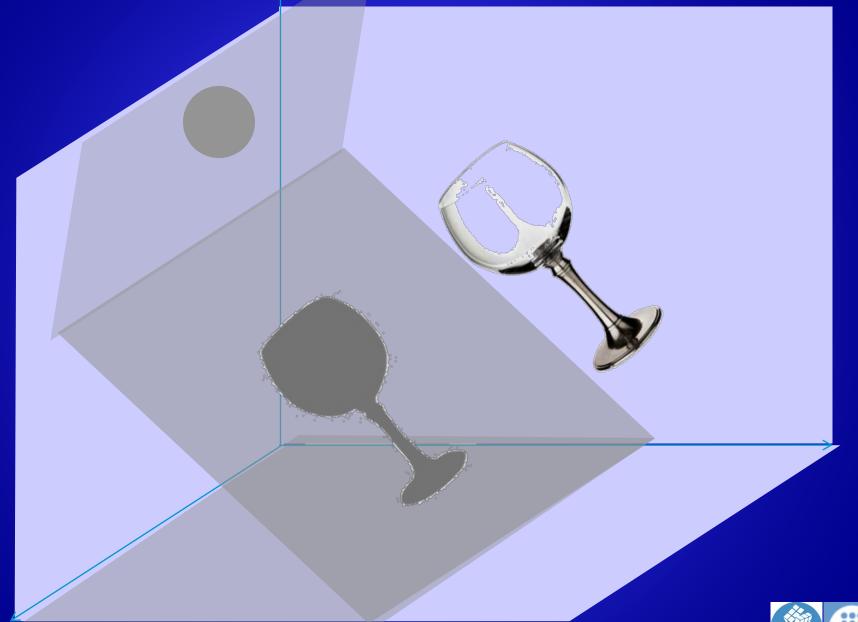- Find the most informative projection planes of data cloud

  *(factorial planes)*

# Factorial Methods

- 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

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

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

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

https://www.facebook.com/grupoajau/posts/1875392015926077
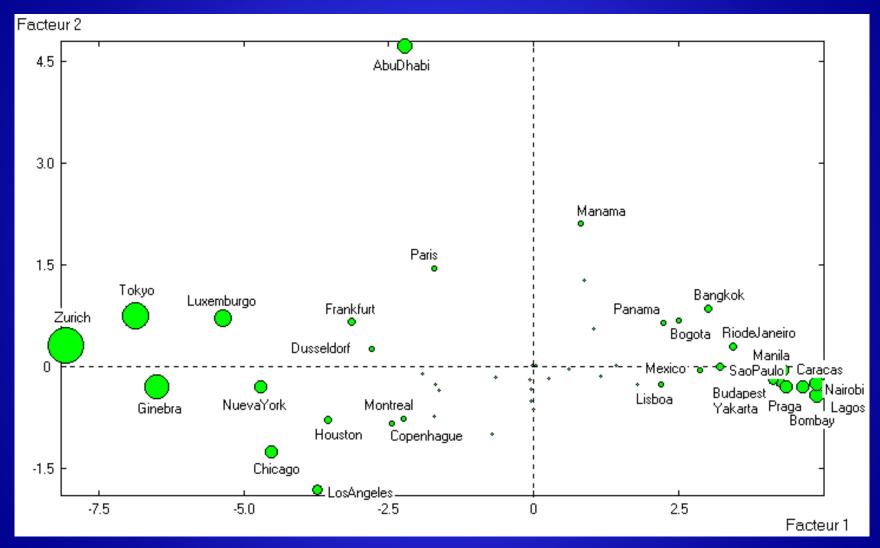
# Factorial Methods

– Output: K factors rotating original X variables

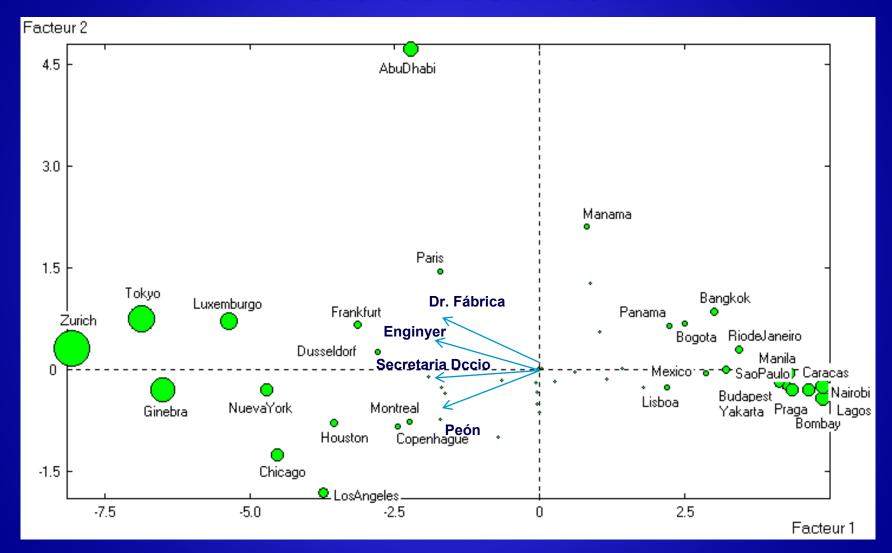–  Factors: Linear combinations of original variables

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.

©K. Gibert

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

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

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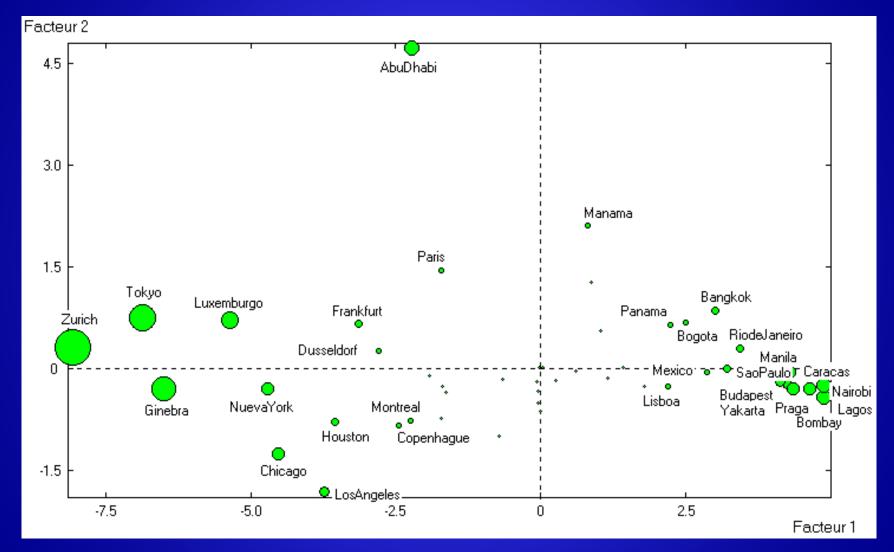## *BCN* Quarters x *Profession of inhabitants*

# Factorial Methods

– 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)
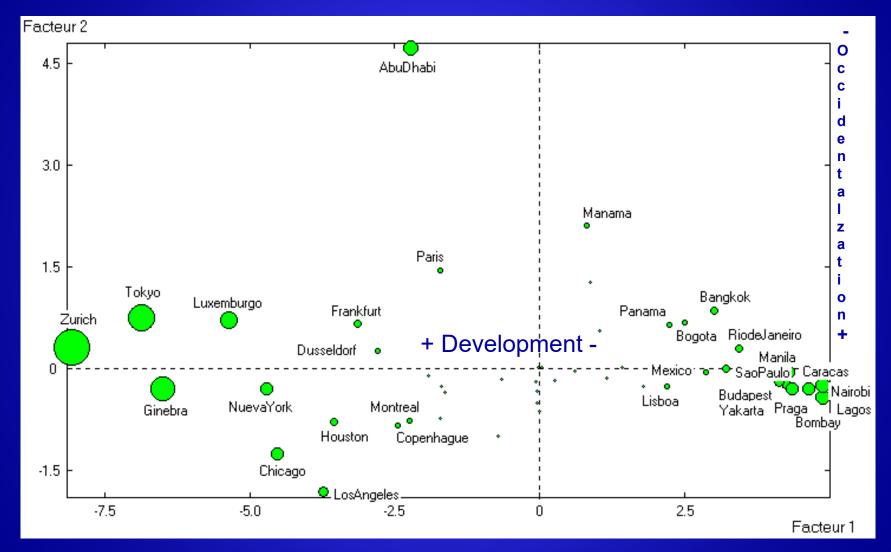
– As a preprocessing method for multidimensionality reduction

©K. Gibert    IDEAI    UPC

# Factorial Methods

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

# Factorial Methods

- **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= $\mathbb{I}_p$)

  *Si les dades estan centrades l'angle entre dues variables projectades coincideix amb la correlació entre elles*

  Matrix $M^{1/2} X'DXM^{1/2}$

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



|         | Workload | Distance to work | Salary |
|---------|----------|------------------|--------|
| Smith   | 1.0      | 0.2              | 1.2    |
| Johnson | 2.0      | 0.0              | 0.3    |
| Williams| -1.0     | 0.1              | -1.0   |
| Jones   | -2.0     | 0.2              | -0.1   |
| Davis   | 0.0      | -0.4             | -0.4   |

# Factorial Methods

*Given triplet <X,M,D>, diagonalize $M^{1/2} X'DXM^{1/2}$*

| Data | Factorial Method | X | M | D |
|---|---|---|---|---|
| Continuous variables | PCA | Centered data matrix | $\mathbb{I}_p$ | $\mathbb{I}_n$ |
| Contingency table ($n_{ij}$) | CA | $F=(n_{ij}/n_i)$ | $diag(1/f_j)$ | $diag(f_i)$ |
| | | $G=(n_{ij}/n_j)$ | $diag(1/f_i)$ | $diag(f_j)$ |
| Categorical variables | MCA | $F=(f_{ij}/(f_i/\sqrt{f_j}))$ | $\mathbb{I}_p$ | $diag(f_i)$ |
| | | Burt table | $\mathbb{I}_{n+p}$ | $diag(n_{ij})$ |

©*K. Gibert*

# Factorial Methods

- **Principal Components Analysis**

  *$M^{\frac{1}{2}}X'DXM^{\frac{1}{2}}$ catches well the data structures*

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

  *$M = \mathbb{I}_p : M^{\frac{1}{2}}X'DXM^{\frac{1}{2}} = X'DX$*

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

  *X standardized and D diagonal : X'DX = Corr(X)*

  *(preferred, big variabilities do not dominate analysis)*

  *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 $u_{\alpha} = (u_{\alpha 1} \ldots u_{\alpha p})$*

# Factorial Methods

- ## 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* $u_\alpha = (u_{\alpha 1} .... u_{\alpha p})$

$u^{-1}X'DXu = \lambda$ *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)

# Factorial Methods

- Given  <X,M,D>

*Diagonalize correlations matrix (with normalized data X'DX)*

*Get r eigen values $\lambda_\alpha$ and sort decreasingly*

$$\{\lambda_\alpha\}_{\alpha=1:r} \qquad \lambda_1 \geq \lambda_2 \geq \lambda_3 \geq \ldots \geq \lambda_r$$

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

$$|u_\alpha| = 1$$

$$u_\alpha u_{\alpha'} = 0$$

$\{u_\alpha\}_{\alpha=1:r}$  *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*

# Factorial Methods

- Given  <X,M,D>

  *Diagonalize Correlations matrix X'DX*

  *Get r eigen values $\lambda_\alpha$ and sort decreasingly*

  $$\{\lambda_\alpha\}_{\alpha=1:r} \qquad \lambda_1 \geq \lambda_2 \geq \lambda_3 \geq \ldots \geq \lambda_r$$

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

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

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

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

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

# Factorial Methods

- Given  <X,M,D>

In general *Diagonalize $M^{½}X'DXM^{½}$*

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

$$\{\lambda_\alpha\}_{\alpha=1:r} \qquad \lambda_1 \geq \lambda_2 \geq \lambda_3 \geq \ldots \geq \lambda_r$$

*Corresponding eigenvectors $u*_\alpha = (u*_{\alpha1} \ldots u*_{\alpha p})$*

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

$$u*_\alpha = M^{-½} u_\alpha$$

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

$$|u*_\alpha|_M = 1: \quad u*'_\alpha M u*_\alpha = u'_\alpha M^{-½} M M^{-½} u_\alpha = 1$$

$$u*_\alpha M u*_{\alpha'} = 0: \quad u*_\alpha M u*_{\alpha'} = u'_\alpha M^{-½} M M^{-½} u_{\alpha'} = 0$$

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

# Principal components analysis
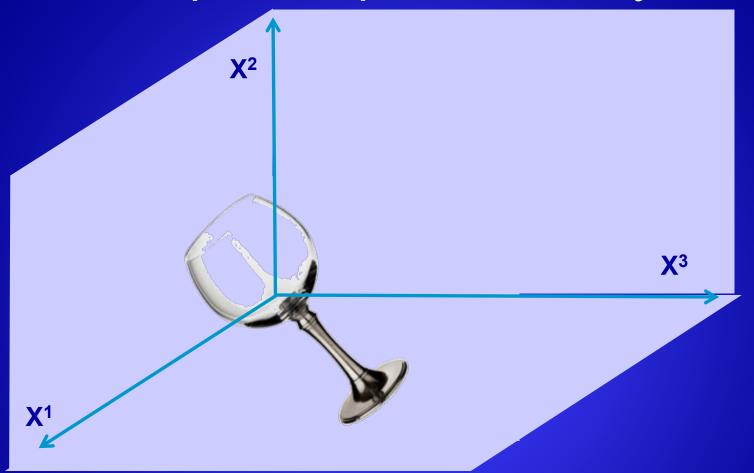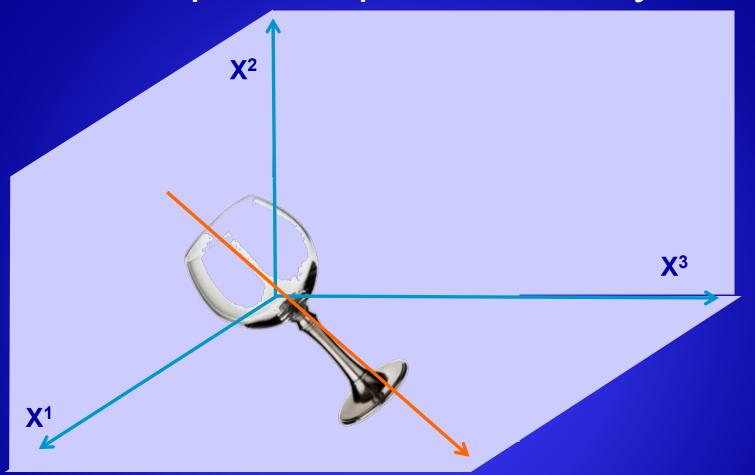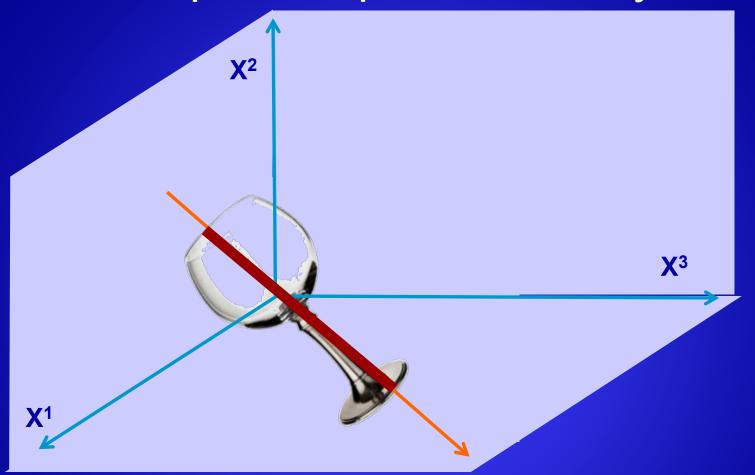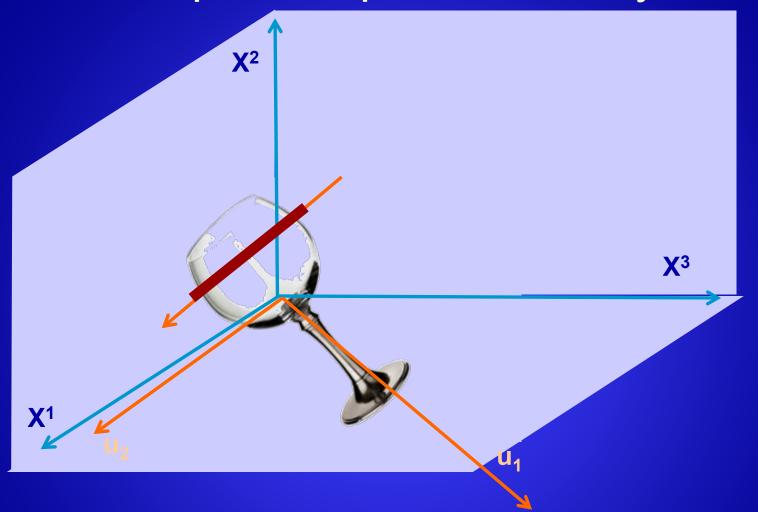
# Principal components analysis



**Centering X**

**(0,0,0)**

# Principal components analysis

# Principal components analysis

# Principal components analysis

# Principal components analysis

# Principal components analysis



©K. Gibert

# Principal components analysis



$X^2$　$u_2$　90°　90°　90°　$X^3$　$X^1$　$u_3$　$u_1$

©K. Gibert

# Principal components analysis

# Factorial Methods

- Given  <X,M,D>

*Diagonalize Correlations matrix X'DX*

*Get r eigen values $\lambda_\alpha$ and sort decreasingly*

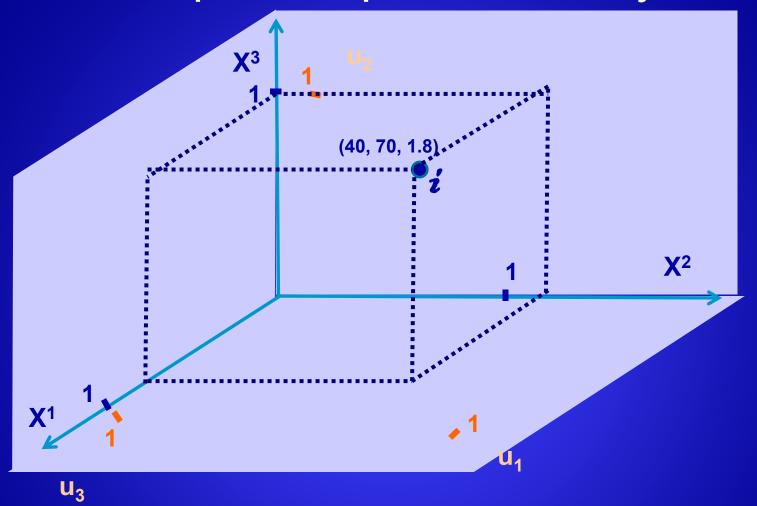$$\{\lambda_\alpha\}_{\alpha=1:r} \qquad \lambda_1 \geq \lambda_2 \geq \lambda_3 \geq_{.....} \geq \lambda_r$$

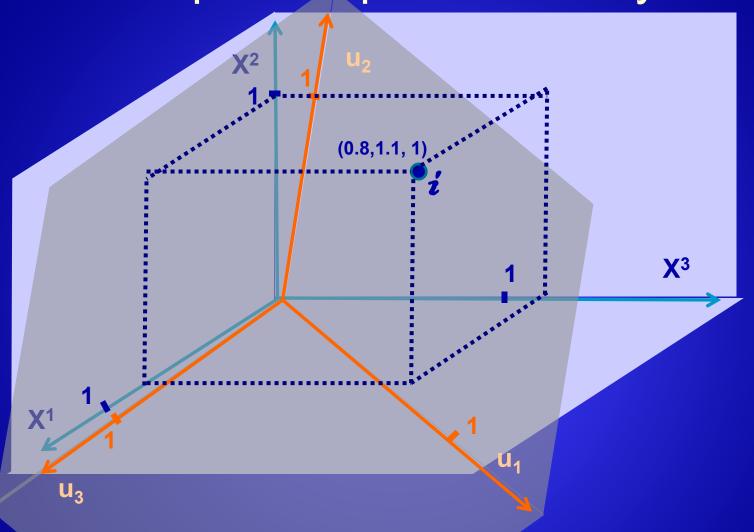*Corresponding eigenvectors $u_\alpha = (u_{\alpha 1} .... u_{\alpha p})$*

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

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

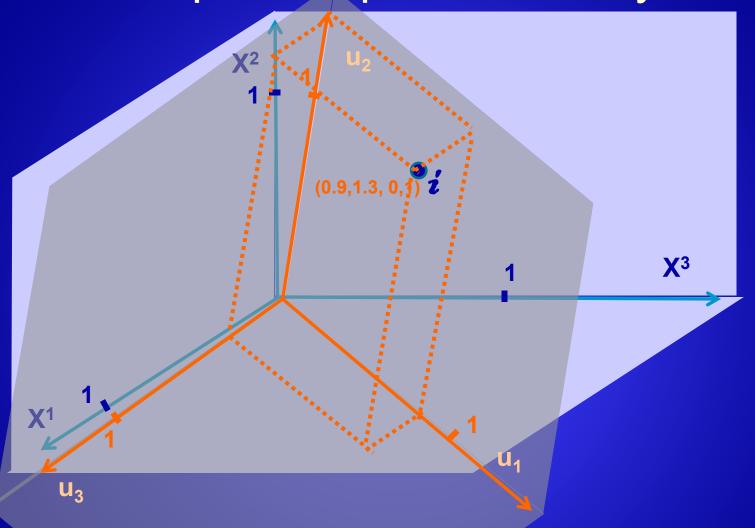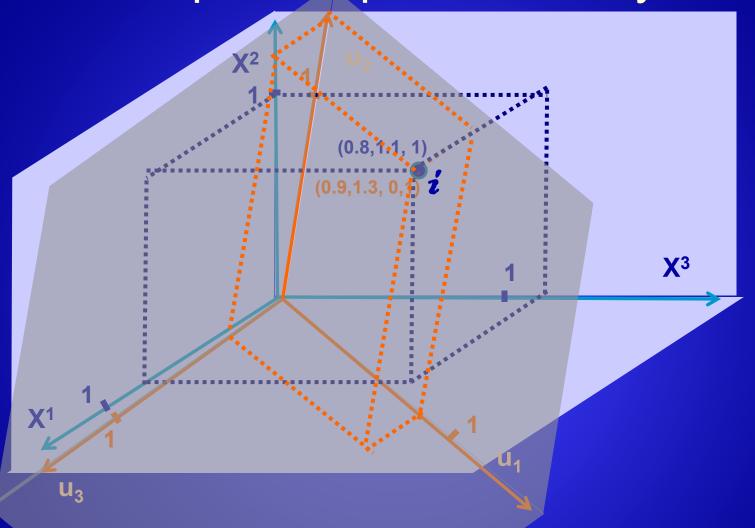*$u^*_\alpha$  are the principal factors of X : good rotation directions*

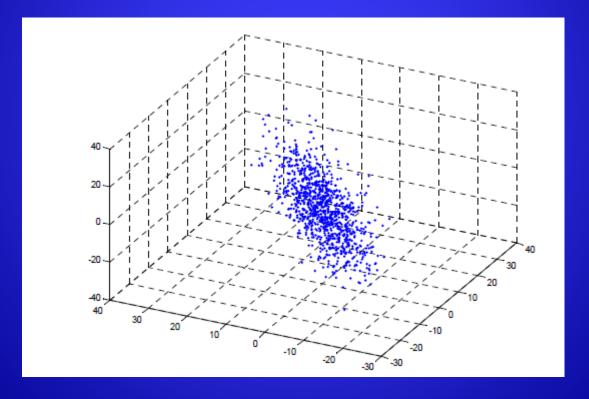*$U^* = ([u^*_1] [u^*_2] ..... [u^*_r])$ is the basis for the projection space*

*How is i expressed in rotated space?*

# Principal components analysis



©K. Gibert

# Principal components analysis



X² u₂

1

1

(0.8, 1.1, 1)

i′

X³

1

1

X¹

1

1

u₁

u₃

# Principal components analysis



$X^2$

$u_2$

$1$

$1$

$(0.9, 1.3, 0.1)$ $i$

$u_3$

$1$

$1$

$u_1$

# Principal components analysis

# Principal components analysis

# Principal components analysis

- Find the most informative projection planes of data cloud

  *(factorial planes)*

# Factorial Methods

- Given <X,M,D>

*Can we find coordinates in rotated space from original ones?*

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

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

*Projection of all individuals: $Pr(X) = U^*_k U^{*\prime}_k M X'$*

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

Projections expressed in original vectorial space

*The best possible projection over k dimensions*

©K. Gibert

# Factorial Methods

- 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

# Factorial Methods

- Given  <X,M,D>

*Diagonalize correlations matrix (with normalized data)*

*eigenvectors $u_\alpha = (u_{\alpha 1} .... u_{\alpha p})$  (direction of factor α, α =1:p)*

$\qquad u_{\alpha p}$ : contribution of variable p to the factor $\alpha$

$\qquad (u_1 ......... u_k)$  *ortonormal*

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

(Projected inertia)

$\{\lambda_\alpha\}_{\alpha =1:r} \qquad \lambda_1 \geq \lambda_2 \geq \lambda_3 \geq_{.....} \geq \lambda_r$

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

*Close objects project close
proximity linked with  association*

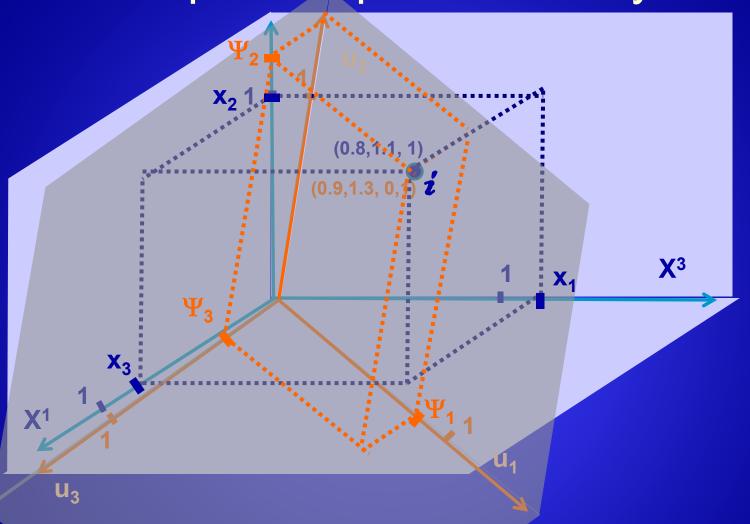# Factorial Methods

- Given <X,M,D>

  *eigenvectors u $_\alpha$ = (u$_{\alpha 1}$.... u$_{\alpha p}$)  (direction of factor k)*

  u$_{\alpha p}$: contribution of variable p to the factor $\alpha$

  *eigen values $\lambda_\alpha$ (quantity of information conserved by factor $\alpha$)*
  (Projected inertia)

  $$\lambda_1 \geq \lambda_2 \geq \lambda_3 \geq_{.....} \geq \lambda_r$$

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

# Principal components analysis

# Factorial Methods

- Given  <X,M,D>

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

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

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

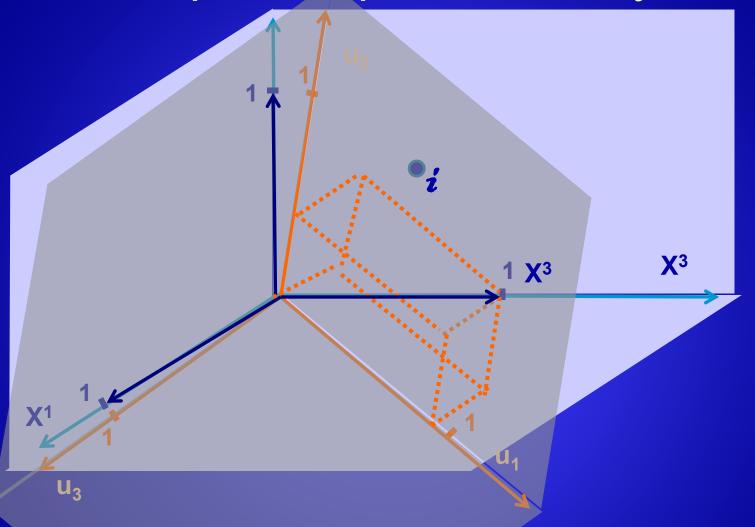Then $\quad \Psi_\alpha^{'} \, D\Psi_\alpha = \lambda_\alpha$
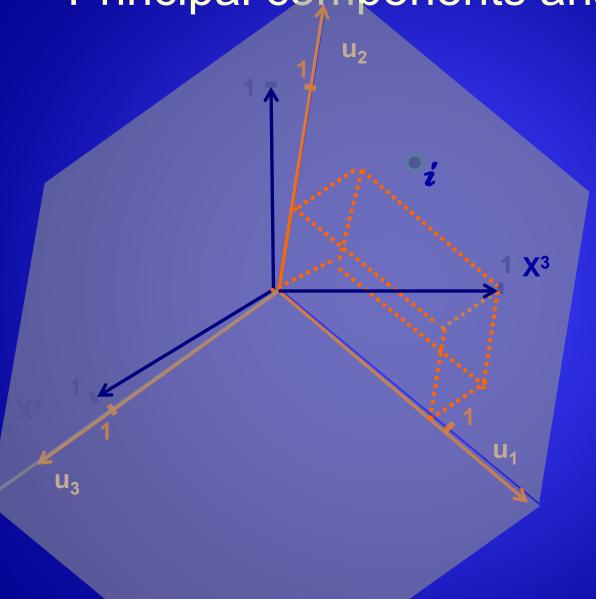
Illustrative points z also projectable

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

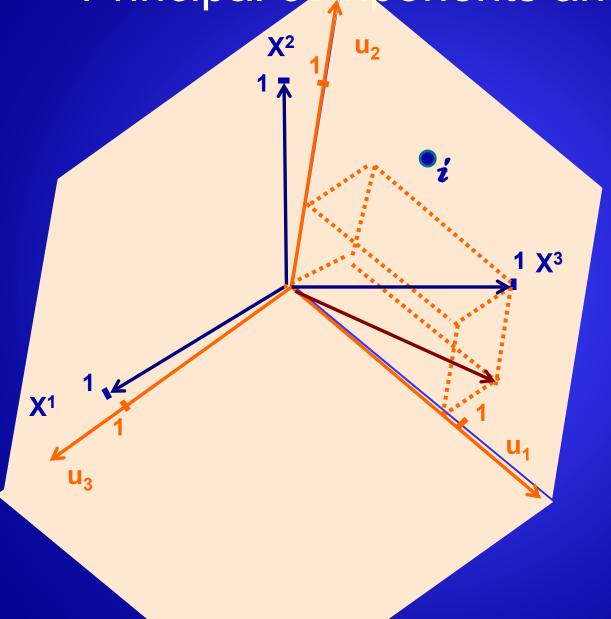*Factors are linear combinations of original variables*

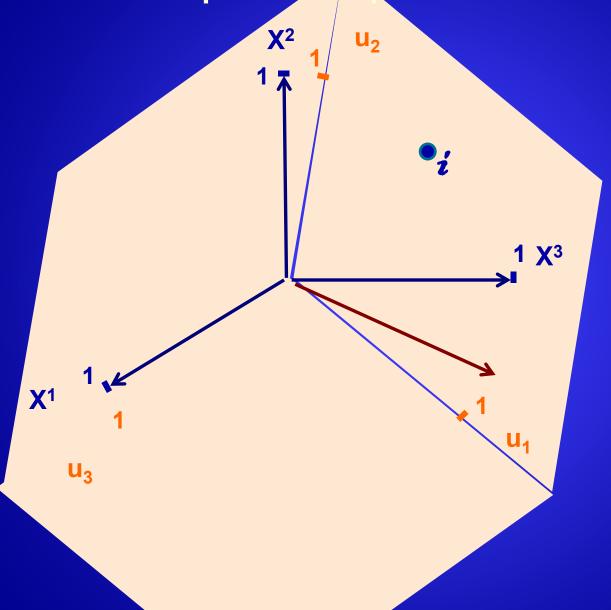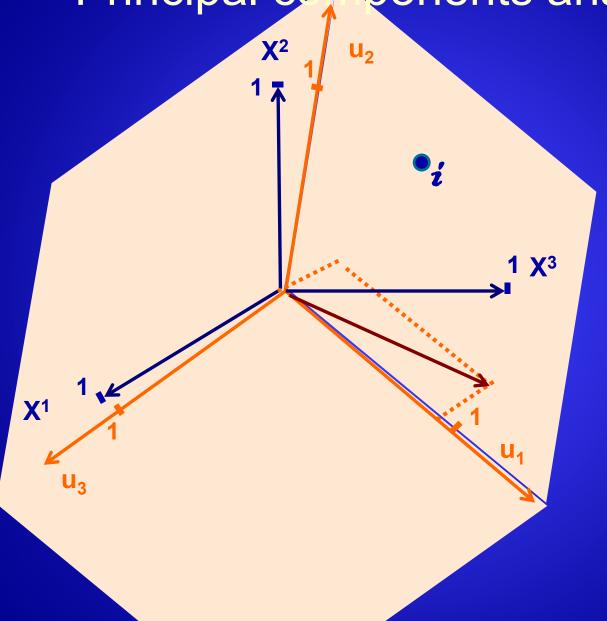Original variables project as VECTORS over factorial space
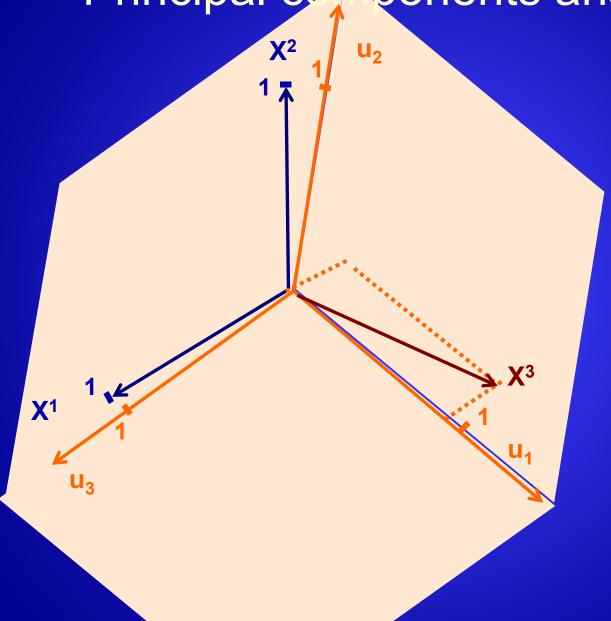angle and lenght  important

# Principal components analysis

# Principal components analysis



©K. Gibert

# Principal components analysis

# Principal components analysis

# Principal components analysis

# Principal components analysis



©K. Gibert

# Principal components analysis

# Principal components analysis

# Principal components analysis



©K. Gibert

# Principal components analysis

# Principal components analysis



©K. Gibert

# Principal components analysis

# Principal components analysis



Map of projected variables

Angles linked with Association

Small angles : correlation

©K. Gibert

# Principal components analysis



*Bicing trips Washington*

©K. Gibert

# Principal components analysis

| Variables | Meaning |
|---|---|
| Start.date | Date of the beginning of the trip |
| End.date | Date of the arrival |
| Durada.Trajecte | Transit's total duration |
| Capacity.S | Bike capacity of the origin station |
| Capacity.E | Bike capacity of the destination station |
| Elevation | Difference in altitude between the stations of arrival and origin |
| Start.long | Starting station's longitude according to the CSR WGS84 |
| End.long | Ending station's longitude according to the CSR WGS84 |
| Temperature | Air temperature |
| Rel.humidity | Air relative humidity |
| Wind.speed | Wind speed |
| Atm.pressure | Atmospheric pressure |

Trajectes de bicing Washington

# Principal components analysis
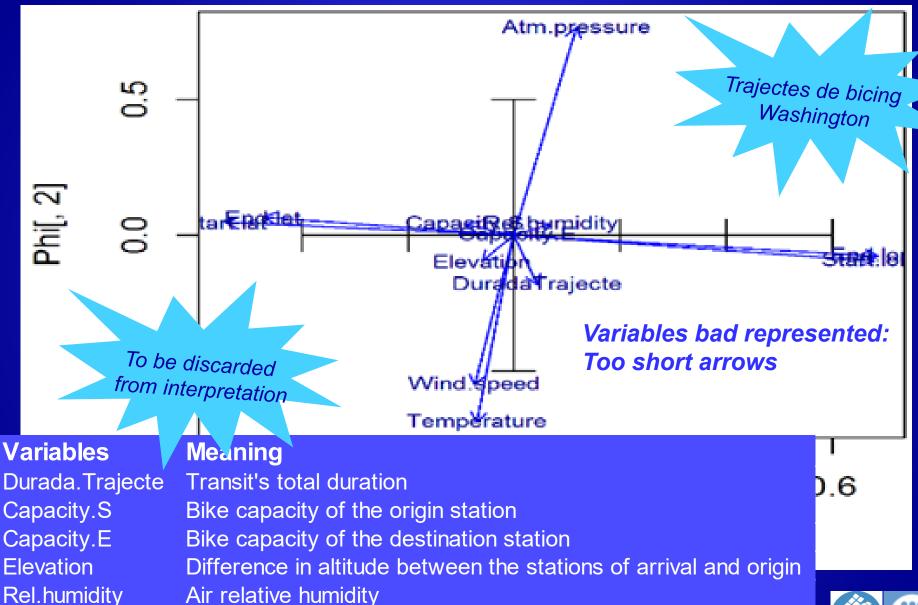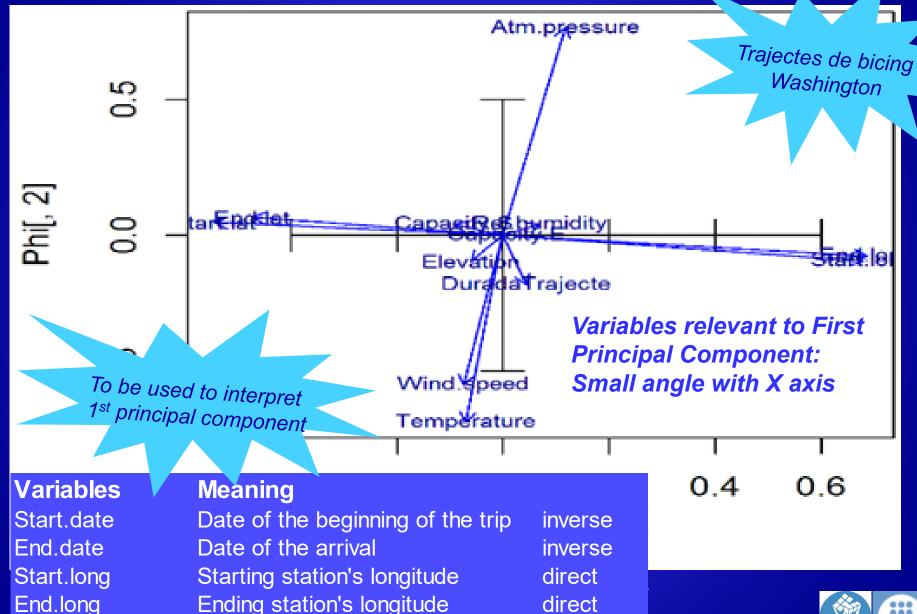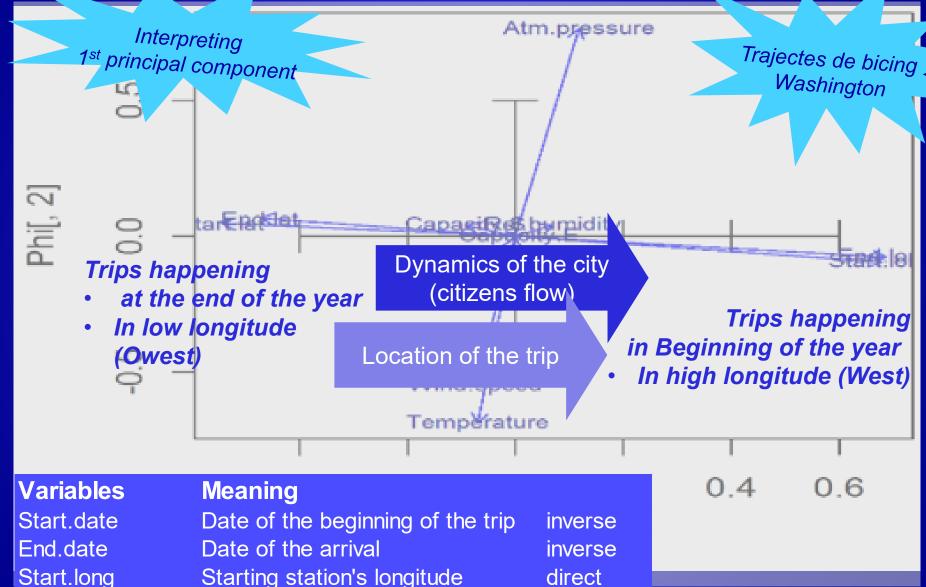
Process to interpret a factorial map

- Forget about variables bad represented in the factorial plan
- Which are the variables with relevant direct contribution to Factor in Axis X (eg. PCA1)?
- Which are the variables with relevant inverse contribution to Factor in Axis X (eg. PCA1)
- (later introduce info on qualitative variables as well)
- Analyze profiles opposed in two extremes of Axis X
- Induce a label for the Factor that represents the concept
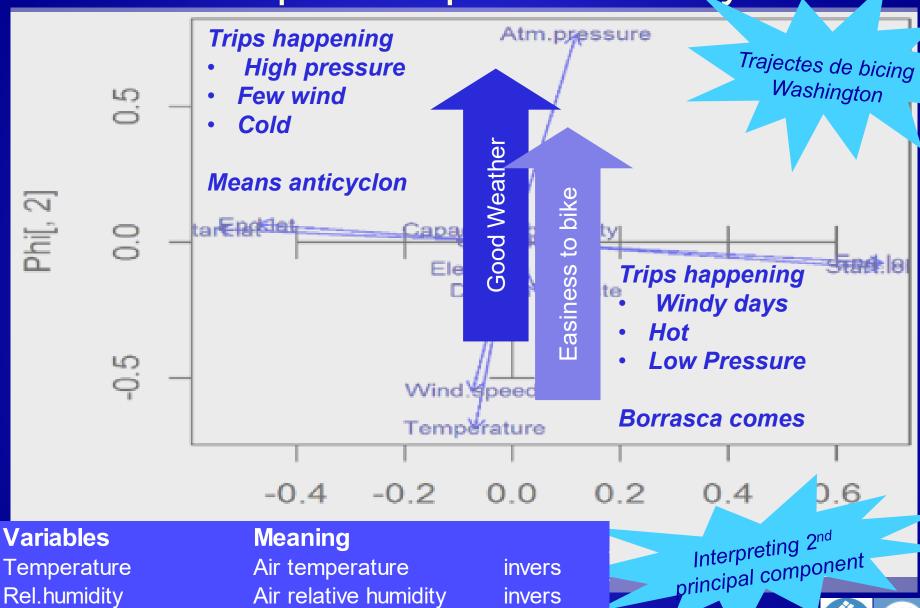
- Repeat with Factor in Axis Y

# Principal components analysis



Atm.pressure

*Trajectes de bicing Washington*

EndLat
CapaciRe.humidity
Capacity.E
Elevation
Durada Trajecte
StartLon

**Variables bad represented: Too short arrows**

*To be discarded from interpretation*

Wind.speed

Temperature

| Variables | Meaning |
|---|---|
| Durada.Trajecte | Transit's total duration |
| Capacity.S | Bike capacity of the origin station |
| Capacity.E | Bike capacity of the destination station |
| Elevation | Difference in altitude between the stations of arrival and origin |
| Rel.humidity | Air relative humidity |

# Principal components analysis



Trajectes de bicing Washington

Variables relevant to First Principal Component: Small angle with X axis

To be used to interpret 1st principal component

| Variables | Meaning | |
|-----------|---------|---|
| Start.date | Date of the beginning of the trip | inverse |
| End.date | Date of the arrival | inverse |
| Start.long | Starting station's longitude | direct |
| End.long | Ending station's longitude | direct |

©K. Gibert  IDEAI  UPC

# Principal components analysis

# Principal components analysis



**Trips happening**
- **High pressure**
- **Few wind**
- **Cold**

**Means anticyclon**

Good Weather

Easiness to bike

Atm.pressure

*Trajectes de bicing Washington*

**Trips happening**
- **Windy days**
- **Hot**
- **Low Pressure**

**Borrasca comes**

Wind.speed

Temperature

Phi[, 2]

*Interpreting 2nd principal component*

| Variables | Meaning | |
|---|---|---|
| Temperature | Air temperature | invers |
| Rel.humidity | Air relative humidity | invers |
| Atm.pressure | Atmospheric pressure | direct |

©K. Gibert

IDEAI    UPC

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

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

# Factorial Methods

- **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 $\kappa \ll p$ *(accumulate 80% inertia)*
    *Reduce data matrix to selected factors*
    *Alternative, keep variables mainly contributing to selected factors (smaller angles with factorial axis)*

©*K. Gibert*

# Factorial Methods

- Given  <X,M,D>

  *Diagonalize Correlations matrix X'DX*

  *Get r eigen values $\lambda_\alpha$ and sort decreasingly*

  $$\{\lambda_\alpha\}_{\alpha=1:r} \qquad \lambda_1 \geq \lambda_2 \geq \lambda_3 \geq_{.....} \geq \lambda_r$$

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

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

  $\{u^*_\alpha\}_{\alpha=1:r}$  *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)*
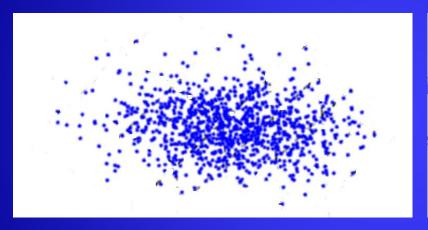
©*K. Gibert*

# Principal components analysis

Harold Hotelling,
American statistician
1895-1973

$X^2$

$u_2$

$X^3$

$X^1$

$u_3$

$u_1$

# Principal components analysis



$X^2$

$u_2$

$X^3$

$X^1$

$u_3$

$u_1$

# Principal components analysis



X²

u₂

X³

X¹

u₃

u₁

©K. Gibert

# Principal components analysis

# Principal components analysis

- Find the most informative projection planes of data cloud

*(factorial planes)*

# Principal components analysis

- Introduce qualitative information (projecting modalities)

Factor 2 - 13.41 %

0.050

0.025

0

-0.050

-0.150    -0.075    0    0.075    0.150

Factor 1 - 79.67 %

h40-49
Pvasco
Madrid
h50-59
Baleares
Catalunya
Asturias
Rioja
Navarra
Galicia
Cantabria
Oval
h60-69
Aragón
h30-39
Canarias
h20-29
Castil-Leo
h70-79
Andalusia
h>80
Murcia
Extremadura
Castil-Mancha

©*K. Gibert*

©*K. Gibert*

©K. Gibert

## *sentences & crimes*

©*K. Gibert*
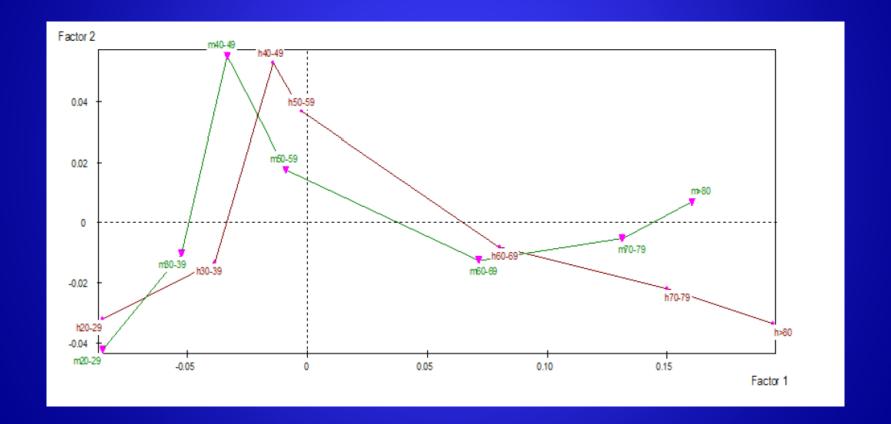
©*K. Gibert*

*sentences & crimes*

# Projecting qualitative variables

Respect the following principles:

- Choose a diferent color for each qualitative variable

- Use the color of the variable for all centroids corresponding to the modalities of the variable

- Include a legent with the list of variables and associated color

- Ensure that legend to not hide any centroide in the factorial map

- For ordinal variables link modalities with rows in the right order and use the color of the variable for the arrows

- Manipulate the size of the font to guarantee the màximum visibility of the map
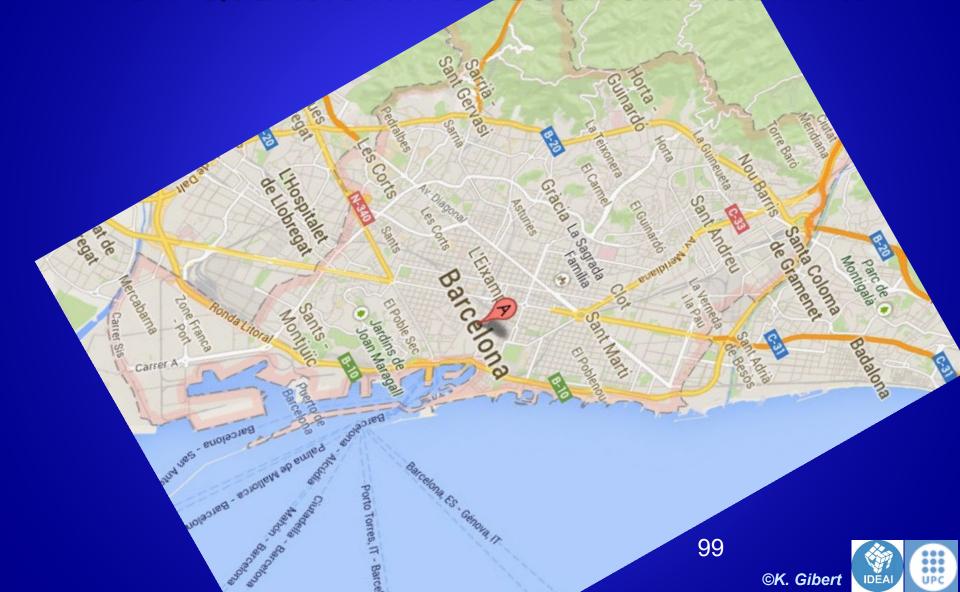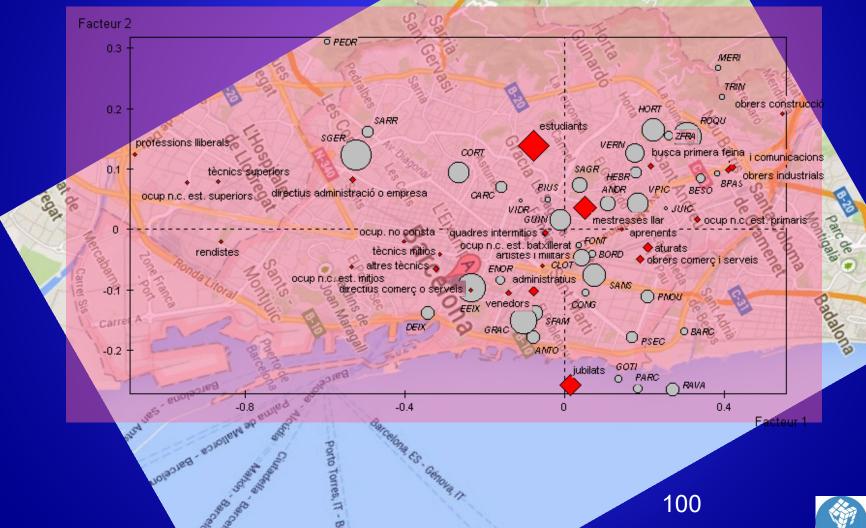
# Efecte guttmann

http://www.ugr.es/~gallardo/

©*K. Gibert*

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

©*K. Gibert*

©*K. Gibert*

# Visualization of the table
## BCN Quarters x *Professions*, *habitants*