

# EST-46114 Métodos Multivariados y Datos Categóricos

Sesion 06 - Analisis de Componentes Principales - Parte 2/3

Juan Carlos Martinez-Ovando

Maestria en Ciencia de Datos





# Objetivos

- Estudiaremos el procedimiento inferencial asociado con PCA (extension del SPCA).
- Realizaremos una ilustracion del procedimiento inferencial.

# Preambulo

PCA es usado en muchas ocasiones para producir estadísticas e indicadores oficiales, e.g.:

- Estadísticas demográficas (Mexico)
- Incidencia criminal (Nigeria)
- Aprovechamiento escolar (Colombia, Mexico)
- Habilidades sociales (Argentina)
- Analisis de rendimientos financieros (Internacional)
- Entre muchos otros. . .

# Motivacion

- Reconociendo que muchos indicadores oficiales (incluyendo los contruidos con PCA), **Manski (2015)** abogo por la inclusion de medidas de error (a.k.a. dispersion) en la publicacion de indicadores estadisticos.
- Esto reconociendo el aspecto inferencial asociado con su construccion.

Revisen la [rectura complementaria de esta sesion](#).

## PCA | Recapitulacion

Recodremos, a partir de un conjunto de datos, matriz  $\mathbf{X}$  de dimension  $(n \times p)$ , PCA se contruye a partir de la asociacion de la estructura espectral de  $\mathbf{\Sigma}$ , la matriz de varianzas y covarianzas subyacente a  $\mathbf{X}$ .

Es decir, los componentes principales de  $\mathbf{X}$  corresponden a la transformacion

$$\mathbf{c}_{j(n \times 1)} = \mathbf{X} \mathbf{v}_j,$$

para  $j = 1, \dots, p$ , donde

$$(\mathbf{e}_j, \mathbf{v}_j)_{j=1}^p,$$

con los eigenvalores asociados con  $\mathbf{X}$  y  $\mathbf{\Sigma}$ .

- Cuando la matriz  $\mathbf{X}$  de dimension  $(n \times p)$  es una muestra de datos, la informacion poblacional acerca de las interacciones de las  $p$  variables es incompleta.
- Informacion incompleta poblacional involucra que PCA construido a partir de  $\mathbf{X}$  o  $\mathbf{\Sigma}$  sea en realidad construido a partir de  $\hat{\mathbf{\Sigma}}$  (un estimador puntual de  $\mathbf{\Sigma}$ ).
- Aun cuando  $\mathbf{X}$  refiera a un censo poblacional, la estructura de codependencia entre las  $p$ -dimensiones de estudio es incierta (i.e. la codependencia puede no ser lineal).
- En los escenarios anteriores, se descarta el aspecto temporal/dinamico del PCA implementado.

## PCA | Ejemplo SPCA Swiss 1/2

En el caso de los datos Swiss, realizamos SPCA como:

```
## Importance of components:
##
##          PC1      PC2      PC3      PC4      PC5      PC6
## Standard deviation  1.7888 1.0901 0.9207 0.66252 0.45225 0.34765
## Proportion of Variance 0.5333 0.1981 0.1413 0.07315 0.03409 0.02014
## Cumulative Proportion 0.5333 0.7313 0.8726 0.94577 0.97986 1.00000

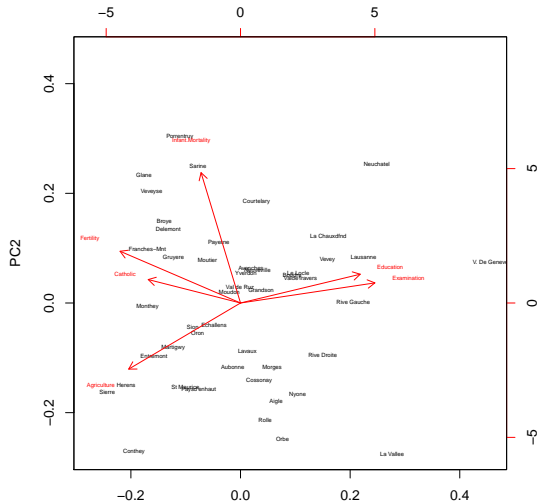
## [1] "center"      "rotation" "scale"      "sdev"      "x"

##          PC1      PC2      PC3
## Fertility    -0.46    0.32   -0.17
## Agriculture  -0.42   -0.41    0.04
## Examination   0.51    0.13   -0.09
## Education     0.45    0.18    0.53
## Catholic     -0.35    0.15    0.81
## Infant.Mortality -0.15   0.81   -0.16
```



## PCA | Ejemplo SPCA Swiss 2/2

```
biplot(swiss.pca, cex = 0.4)
```



**Planteamiento** Pensemos que cada renglon de la matriz  $\mathbf{X}$  es una observacion de

$$\mathbf{X}_{j\cdot} \sim N_p(\mathbf{X}|\boldsymbol{\mu}, \boldsymbol{\Lambda}),$$

para  $j = 1, \dots, n$ .

- Adoptamos el supuesto de homogeneidad e independencia estocastica en la enunciaci3n anterior.

## Complementariamente

El desconocimiento acerca de  $(\boldsymbol{\mu}, \boldsymbol{\Lambda})$  es expresado como:

$$(\boldsymbol{\mu}, \boldsymbol{\Lambda}) \sim \text{N-Wi}(\boldsymbol{\mu}, \boldsymbol{\Lambda}|\mathbf{m}_0, s_0, a_0, \mathbf{B}_0),$$

para ciertos valores de  $\mathbf{m}_0, s_0, a_0, \mathbf{B}_0$ .

### Aprendizaje

Así, al consolidar las fuentes de información tenemos:

$$(\boldsymbol{\mu}, \boldsymbol{\Lambda} | \mathbf{X}) \sim \text{N-Wi}(\boldsymbol{\mu}, \boldsymbol{\Lambda} | \mathbf{m}_n, s_n, a_n, \mathbf{B}_n),$$

con  $\mathbf{m}_n, s_n, a_n, \mathbf{B}_n$  expresados analíticamente en la presentación de la Sesión 04.

### PCA inferencial

Con base en el aprendizaje realizado, se define para un componente principal  $c_{j\cdot (n \times 1)}$  como

$$\begin{aligned} c_{j\cdot} &= \mathbf{X} \mathbf{v}_j \\ &= \mathbf{X} \mathbf{v}_j(\mu, \Lambda), \end{aligned}$$

donde  $\mathbf{v}_j(\mu, \Lambda)$  **hace explicito** que el eigenvector asociado esta en funcion del *parametro desconocido*  $\Lambda^{-1}$ .

- Por consiguiente, al ser definidos  $(\mathbf{e}_j, \mathbf{v}_j)_{j=1}^p$  en funcion de  $\mathbf{\Lambda}$  (*aleatoria*), los define como **desconocidos y aleatorios** tambien.
- La distribucion de  $(\mathbf{e}, \mathbf{V}) = (\mathbf{e}_j, \mathbf{v}_j)_{j=1}^p$  esta definida por  $\text{N-Wi}(\mu, \mathbf{\Lambda} | \mathbf{m}_n, s_n, a_n, \mathbf{B}_n)$ .

Siguiendo el ultimo punto, el calculo de

$$\mathbb{P}(\mathbf{e}, \mathbf{V} | \mathbf{m}_n, s_n, a_n, \mathbf{B}_n),$$

es *difícil de obtener analiticamente*.

- Press (2005) sugiere emplear distribuciones asintoticas basadas en el CLT (vea Teoremas 9.3.2 y 9.3.3). Eso solo para la distribucion de los eigenvalores  $\mathbf{e} = (e_1, \dots, e_p)$ , mas no para los eigenvectores.
- En la practica, esta distribucion puede aproximarse via muestras de Monte Carlo de  $N\text{-Wi}(\boldsymbol{\mu}, \boldsymbol{\Lambda} | \mathbf{m}_n, s_n, a_n, \mathbf{B}_n)$  de acuerdo al siguiente algoritmo.

## Algoritmo

0. Actualizar la distribución de los parámetros  $(\mu, \Lambda)$  con base en  $\mathbf{X}$ . Fijar el tamaño de la muestra simulada  $M$ .

Para  $m = 1, \dots, M$ ,

1. Generar los pseudodatos aleatorios  $(\mu^{(m)}, \Lambda^{(m)})$  como

$$\begin{aligned}\Lambda^{(m)} &\sim \text{Wi}(\Lambda | a_n, \mathbf{B}_n), \\ \mu^{(m)} | \Lambda^{(m)} &\sim \text{N}(\mu | \mathbf{m}_n, s_n \Lambda^{(m)}).\end{aligned}$$

2. Generar la descomposición espectral

$$(\mathbf{e}_j^{(m)}, \mathbf{v}_j^{(m)})_{j=1}^p,$$

$$\text{de } \Sigma^{(m)} = (\Lambda^{(m)})^{-1}.$$

3. Generar los componentes principales  $\mathbf{C}^{(m)} = (\mathbf{c}_{1\cdot}^{(m)}, \dots, \mathbf{c}_{p\cdot}^{(m)})$  como

$$\mathbf{c}_{j\cdot}^{(m)} = \mathbf{X} \mathbf{v}_j^{(m)}.$$

## Comentarios

- Las muestras

$$\left( \mathbf{e}_j^{(m)}, \mathbf{v}_j^{(m)} \right)_{j=1, m=1}^{p, M}$$

son iid de

$$\mathbb{P}(\mathbf{e}, \mathbf{V} | \mathbf{m}_n, s_n, a_n, \mathbf{B}_n).$$

- Las muestras

$$\left( \mathbf{C}^{(m)} \right)_{m=1}^M = (\mathbf{c}_{1\cdot}^{(m)}, \dots, \mathbf{c}_{p\cdot}^{(m)})_{m=1}^M$$

son iid de

$$\mathbb{P}(\mathbf{C} | \mathbf{m}_n, s_n, a_n, \mathbf{B}_n).$$

## PCA | Swiss Inferencial 7/

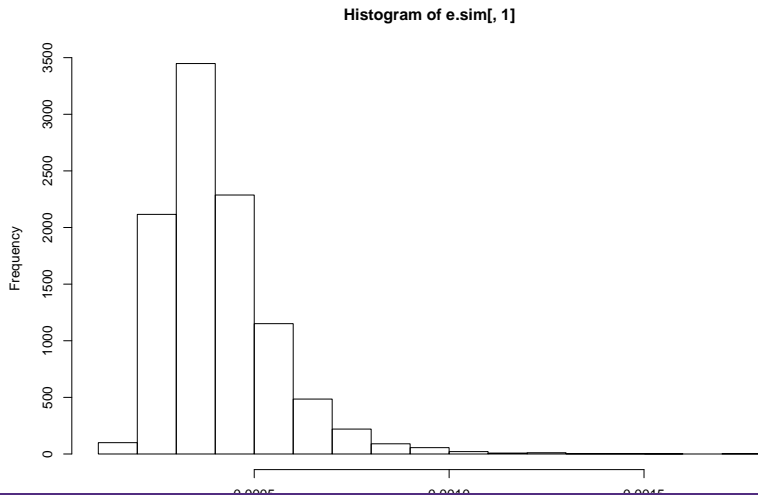
```
M <- 10000
mu.sim <- matrix(NA,nrow=M, ncol=ncol(data))
Lambda.sim <- array(NA,dim=c(M,ncol(data),ncol(data)))
e.sim <- matrix(NA,nrow=M, ncol=ncol(data))
V.sim <- array(NA,dim=c(M,ncol(data),ncol(data)))
C.sim <- array(NA,dim=c(M,nrow(data),ncol(data)))
m <- 1; X <- as.matrix(data)
for(m in 1:M){
  # Simulacion (mu,Lambda)
  Lambda.sim[m,,] <- rWishart(1, output[[3]], output[[4]])
  mu.sim[m,] <- mvrnorm(1, mu=output[[1]], Sigma=solve(output[[2]]*Lamb
  # Simulacion (e,V) + C
  eigen_aux <- eigen(solve(Lambda.sim[m,,]))
  e.sim[m,] <- eigen_aux$values
  V.sim[m,,] <- eigen_aux$vectors
  C.sim[m,,] <- X %*% V.sim[m,,]
}
```



# PCA | Swiss Inferencial 8/

Inferencia sobre el eigenvenctor  $e_1$

```
hist(e.sim[,1])
```



Inferencia sobre el eigenvenctor  $e_1$

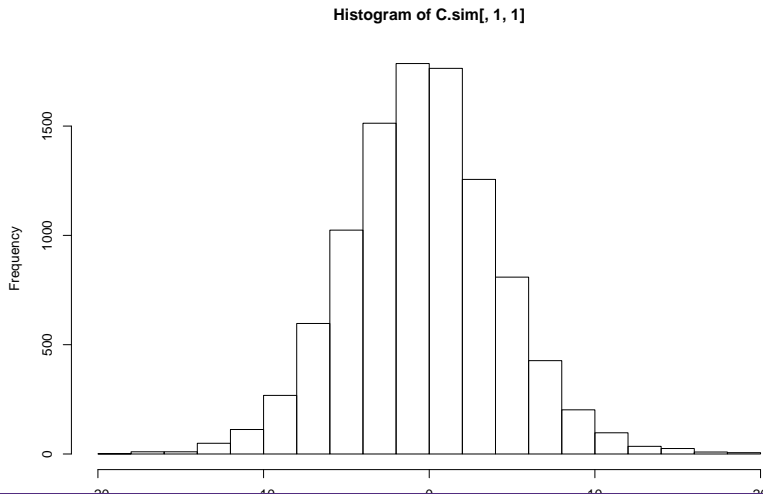
```
summary(e.sim[,2])
```

##	Min.	1st Qu.	Median	Mean	3rd Qu.	Max.
##	5.912e-05	1.277e-04	1.545e-04	1.630e-04	1.895e-04	4.855e-04

# PCA | Swiss Inferencial 10/

Inferencia sobre el componente principal uno  $c_{i,1}$  de la observacion  $i = 1$

```
hist(C.sim[,1,1])
```



Inferencia sobre el componente principal uno  $c_{i,1}$  de la observacion  $i = 1$

```
summary(C.sim[,1,1])
```

##	Min.	1st Qu.	Median	Mean	3rd Qu.	Max.
##	-19.1027	-3.3611	-0.3913	-0.3644	2.5175	19.4049

# Ejercicio

Consideren el archivo `est46114_s06_data.xls` con los tipos de cambio reales (respecto a USD) de varias economías (periodo 1970-2010).

1. Implementen el *procedimiento inferencial PCA* considerando distribuciones iniciales no informativas para  $\mu, \Lambda$ .
2. Reporten qué economía tiene el mayor peso esperado en la descomposición PCA.
3. Reporten qué economía tiene la mayor consistencia en estimación de los  $c_j$ s correspondientes.

# Lecturas complementarias

- Maski (2015) “Communicating Uncertainty in Official Economic Statistics”.  
est46114\_s06\_suplemento.pdf
- Press (2005) *Applied Multivariate analysis*, Caps. 3 y 5.

# Table of Contents