

# Propuestas Robustas para el Análisis de Correlaciones Canónicas

Stella Maris Donato<sup>1</sup>    Jorge G. Adrover<sup>2</sup>

<sup>1</sup>Universidad de Buenos Aires

<sup>2</sup>Universidad Nacional de Córdoba, CIEM y CONICET.

22 de setiembre de 2016

- Revisión de Correlaciones Canónicas.

- Revisión de Correlaciones Canónicas.
- Estimadores propuestos.

- Revisión de Correlaciones Canónicas.
- Estimadores propuestos.
- Estudio de simulación.

- Revisión de Correlaciones Canónicas.
- Estimadores propuestos.
- Estudio de simulación.
- Conclusiones.

- Revisión de Correlaciones Canónicas.
- Estimadores propuestos.
- Estudio de simulación.
- Conclusiones.
- Referencias.

# Revisión de Correlaciones Canónicas

Dados dos grupos de variables, el primero conformado por  $p$  variables y representado a través del vector aleatorio  $\mathbf{x} = (x_1, x_2, \dots, x_p)^t$ , y el segundo formado por  $q$  variables y representado a través del vector aleatorio  $\mathbf{y} = (y_1, y_2, \dots, y_q)^t$ . Es conveniente considerarlos conjuntamente, para lo cual se define el vector  $\mathbf{z}$

$$\mathbf{z} = \begin{pmatrix} \mathbf{x} \\ \mathbf{y} \end{pmatrix} = (x_1, x_2, \dots, x_p, y_1, y_2, \dots, y_q)^t,$$

que tiene vector de medias

$$\boldsymbol{\mu} = E(\mathbf{z}) = \begin{pmatrix} E(\mathbf{x}) \\ E(\mathbf{y}) \end{pmatrix} = \begin{pmatrix} \boldsymbol{\mu}_x \\ \boldsymbol{\mu}_y \end{pmatrix}$$

y matriz de covarianzas

$$\boldsymbol{\Sigma} = E[(\mathbf{z} - \boldsymbol{\mu})(\mathbf{z} - \boldsymbol{\mu})^t] = \begin{pmatrix} \boldsymbol{\Sigma}_{xx} & \boldsymbol{\Sigma}_{xy} \\ \boldsymbol{\Sigma}_{yx} & \boldsymbol{\Sigma}_{yy} \end{pmatrix}.$$

Clásicamente, el primer par de vectores canónicos  $\boldsymbol{\alpha}_1 = (\alpha_{11}, \alpha_{12}, \dots, \alpha_{1p})^t$  y  $\boldsymbol{\beta}_1 = (\beta_{11}, \beta_{12}, \dots, \beta_{1q})^t$  son los que resuelven

$$(\boldsymbol{\alpha}_1, \boldsymbol{\beta}_1) = \arg \max_{(\mathbf{a}, \mathbf{b}) \in \mathcal{A}_1} \text{Corr}(\mathbf{a}^t \mathbf{x}, \mathbf{b}^t \mathbf{y}),$$

donde

$$\mathcal{A}_1 = \left\{ (\mathbf{a}, \mathbf{b}) \in \mathbb{R}^p \times \mathbb{R}^q : \begin{array}{l} \text{Cov}(\mathbf{a}^t \mathbf{x}, \mathbf{a}^t \mathbf{x}) = \mathbf{a}^t \boldsymbol{\Sigma}_{xx} \mathbf{a} = 1, \\ \text{Cov}(\mathbf{b}^t \mathbf{y}, \mathbf{b}^t \mathbf{y}) = \mathbf{b}^t \boldsymbol{\Sigma}_{yy} \mathbf{b} = 1 \end{array} \right\}.$$

Si el rango de  $\boldsymbol{\Sigma}_{xy} > 1$ , se buscan vectores canónicos de orden superior  $\boldsymbol{\alpha}_k = (\alpha_{k1}, \alpha_{k2}, \dots, \alpha_{kp})^t$  y  $\boldsymbol{\beta}_k = (\beta_{k1}, \beta_{k2}, \dots, \beta_{kq})^t$ .

La definición es recursiva para  $k = 2, 3, \dots, \min(p, q)$ . Entonces, dado los primeros  $k - 1$  vectores canónicos  $(\boldsymbol{\alpha}_1, \boldsymbol{\beta}_1), \dots, (\boldsymbol{\alpha}_{k-1}, \boldsymbol{\beta}_{k-1})$ , se define

$$(\boldsymbol{\alpha}_k, \boldsymbol{\beta}_k) = \arg \max_{(\mathbf{a}, \mathbf{b}) \in \mathcal{A}_k} \text{Corr}(\mathbf{a}^t \mathbf{x}, \mathbf{b}^t \mathbf{y}), \quad (1)$$

donde

$$\mathcal{A}_k = \left\{ (\mathbf{a}, \mathbf{b}) \in \mathbb{R}^p \times \mathbb{R}^q : \begin{array}{ll} \text{Cov}(\mathbf{a}^t \mathbf{x}, \mathbf{a}^t \mathbf{x}) = 1, & \text{Cov}(\mathbf{a}^t \mathbf{x}, \boldsymbol{\alpha}_j^t \mathbf{x}) = 0, \\ \text{Cov}(\mathbf{b}^t \mathbf{y}, \mathbf{b}^t \mathbf{y}) = 1, & \text{Cov}(\mathbf{b}^t \mathbf{y}, \boldsymbol{\beta}_j^t \mathbf{y}) = 0, \\ & j = 1, 2, \dots, k - 1 \end{array} \right\}$$

Este planteo tiene una solución conocida

- Los vectores canónicos  $(\alpha_k, \beta_k)$ ,  $k = 1, \dots, \min(p, q)$  son los autovectores correspondientes a los autovalores  $\lambda_1 \geq \dots \geq \lambda_{\min(p,q)}$  de las matrices

$$\Sigma_{xx}^{-1} \Sigma_{xy} \Sigma_{yy}^{-1} \Sigma_{yx} \quad y \quad \Sigma_{yy}^{-1} \Sigma_{yx} \Sigma_{xx}^{-1} \Sigma_{xy}, \quad (2)$$

respectivamente

Este planteo tiene una solución conocida

- Los vectores canónicos  $(\boldsymbol{\alpha}_k, \boldsymbol{\beta}_k)$ ,  $k = 1, \dots, \min(p, q)$  son los autovectores correspondientes a los autovalores  $\lambda_1 \geq \dots \geq \lambda_{\min(p,q)}$  de las matrices

$$\Sigma_{\mathbf{xx}}^{-1} \Sigma_{\mathbf{xy}} \Sigma_{\mathbf{yy}}^{-1} \Sigma_{\mathbf{yx}} \quad y \quad \Sigma_{\mathbf{yy}}^{-1} \Sigma_{\mathbf{yx}} \Sigma_{\mathbf{xx}}^{-1} \Sigma_{\mathbf{xy}}, \quad (2)$$

respectivamente

- Y las correlaciones canónicas están dadas por

$$[\text{Corr}(\boldsymbol{\alpha}_k^t \mathbf{x}, \boldsymbol{\beta}_k^t \mathbf{y})]^2 = \lambda_k \quad (3)$$

# Revisión de Correlaciones Canónicas

El análisis canónico se puede plantear en forma alternativa buscando las matrices  $A^* \in \mathbb{R}^{r \times p}$  y  $B^* \in \mathbb{R}^{r \times q}$  ( $r \leq \min(p, q)$ ), y el vector  $\mathbf{a}^* \in \mathbb{R}^r$  tales que:

$$(A^*, B^*, \mathbf{a}^*) = \underset{(A, B, \mathbf{a}) \in \mathcal{C}}{\operatorname{arg\,min}} E \left( \|A\mathbf{x} - B\mathbf{y} - \mathbf{a}\|^2 \right) \quad (4)$$

donde

$$\mathcal{C} = \left\{ (A, B, \mathbf{a}) : A \in \mathbb{R}^{r \times p}, B \in \mathbb{R}^{r \times q}, \mathbf{a} \in \mathbb{R}^r, A\Sigma_{xx}A^t = I = B\Sigma_{yy}B^t \right\}. \quad (5)$$

Cuya solución es la siguiente (Ver Seber (1984) y Brillinger (1975)):

- $A^*$  y  $B^*$  son matrices que tienen en sus filas los autovectores correspondientes a los autovalores  $\lambda_1 \geq \dots \geq \lambda_r$  de las matrices en (2),

# Revisión de Correlaciones Canónicas

El análisis canónico se puede plantear en forma alternativa buscando las matrices  $A^* \in \mathbb{R}^{r \times p}$  y  $B^* \in \mathbb{R}^{r \times q}$  ( $r \leq \min(p, q)$ ), y el vector  $\mathbf{a}^* \in \mathbb{R}^r$  tales que:

$$(A^*, B^*, \mathbf{a}^*) = \arg \min_{(A, B, \mathbf{a}) \in \mathcal{C}} E \left( \|A\mathbf{x} - B\mathbf{y} - \mathbf{a}\|^2 \right) \quad (4)$$

donde

$$\mathcal{C} = \left\{ (A, B, \mathbf{a}) : A \in \mathbb{R}^{r \times p}, B \in \mathbb{R}^{r \times q}, \mathbf{a} \in \mathbb{R}^r, A\Sigma_{xx}A^t = I = B\Sigma_{yy}B^t \right\}. \quad (5)$$

Cuya solución es la siguiente (Ver Seber (1984) y Brillinger (1975)):

- $A^*$  y  $B^*$  son matrices que tienen en sus filas los autovectores correspondientes a los autovalores  $\lambda_1 \geq \dots \geq \lambda_r$  de las matrices en (2),
- $\mathbf{a}^* = A^*\boldsymbol{\mu}_x - B^*\boldsymbol{\mu}_y$

# Revisión de Correlaciones Canónicas

El análisis canónico se puede plantear en forma alternativa buscando las matrices  $A^* \in \mathbb{R}^{r \times p}$  y  $B^* \in \mathbb{R}^{r \times q}$  ( $r \leq \min(p, q)$ ), y el vector  $\mathbf{a}^* \in \mathbb{R}^r$  tales que:

$$(A^*, B^*, \mathbf{a}^*) = \arg \min_{(A, B, \mathbf{a}) \in \mathcal{C}} E \left( \|A\mathbf{x} - B\mathbf{y} - \mathbf{a}\|^2 \right) \quad (4)$$

donde

$$\mathcal{C} = \left\{ (A, B, \mathbf{a}) : A \in \mathbb{R}^{r \times p}, B \in \mathbb{R}^{r \times q}, \mathbf{a} \in \mathbb{R}^r, A\Sigma_{xx}A^t = I = B\Sigma_{yy}B^t \right\}. \quad (5)$$

Cuya solución es la siguiente (Ver Seber (1984) y Brillinger (1975)):

- $A^*$  y  $B^*$  son matrices que tienen en sus filas los autovectores correspondientes a los autovalores  $\lambda_1 \geq \dots \geq \lambda_r$  de las matrices en (2),
- $\mathbf{a}^* = A^*\boldsymbol{\mu}_x - B^*\boldsymbol{\mu}_y$
- Y las correlaciones canónicas cumplen la ecuación (3).

- Los estimadores propuestos en este trabajo están basados en la idea de minimizar la distancia entre las variables canónicas y este problema de optimización tiene relación con Componentes Principales (PCA).

- Los estimadores propuestos en este trabajo están basados en la idea de minimizar la distancia entre las variables canónicas y este problema de optimización tiene relación con Componentes Principales (PCA).
- Sea  $A_o = A\Sigma_{xx}^{1/2}$ ,  $B_o = B\Sigma_{yy}^{1/2}$ ,  $D = (A_o - B_o) \in \mathbb{R}^{r \times m}$ ,  $m = p + q$ , y el vector aleatorio  $\mathbf{z} = (\mathbf{x}^t \Sigma_{xx}^{-1/2}, \mathbf{y}^t \Sigma_{yy}^{-1/2})^t$ .

- Los estimadores propuestos en este trabajo están basados en la idea de minimizar la distancia entre las variables canónicas y este problema de optimización tiene relación con Componentes Principales (PCA).
- Sea  $A_o = A\Sigma_{xx}^{1/2}$ ,  $B_o = B\Sigma_{yy}^{1/2}$ ,  $D = (A_o - B_o) \in \mathbb{R}^{r \times m}$ ,  $m = p + q$ , y el vector aleatorio  $\mathbf{z} = (\mathbf{x}^t \Sigma_{xx}^{-1/2}, \mathbf{y}^t \Sigma_{yy}^{-1/2})^t$ .
- Reformulando la expresión para los vectores estandarizados  $\Sigma_{xx}^{-1/2} \mathbf{x}$  y  $\Sigma_{yy}^{-1/2} \mathbf{y}$ , se obtiene

$$\min_{(A,B,\mathbf{a}) \in \mathcal{C}} E \left( \|A\mathbf{x} - B\mathbf{y} - \mathbf{a}\|^2 \right) = \min_{(D,\mathbf{a}) \in \mathcal{B}_{r,m}} E \left( \|D\mathbf{z} - \mathbf{a}\|^2 \right)$$

con

$$\mathcal{B}_{r,m} = \left\{ (D, \mathbf{a}) : \mathbf{a} \in \mathbb{R}^r, D = (A_o - B_o) \in \mathbb{R}^{r \times m}, A_o A_o^t = I_r = B_o B_o^t \right\}$$

que es similar al problema de PCA que trata Maronna (2005).

Notar que

- Si se estandarizan los vectores  $\mathbf{x}$  e  $\mathbf{y}$ ,

$$\mathbf{z} = \begin{pmatrix} \Sigma_{\mathbf{xx}}^{-1/2} (\mathbf{x} - E(\mathbf{x})) \\ \Sigma_{\mathbf{yy}}^{-1/2} (\mathbf{y} - E(\mathbf{y})) \end{pmatrix},$$

la matriz de covarianzas resultante para  $\mathbf{z}$  es

$$M = \begin{pmatrix} I & \Sigma_{\mathbf{xx}}^{-1/2} \Sigma_{\mathbf{xy}} \Sigma_{\mathbf{yy}}^{-1/2} \\ \Sigma_{\mathbf{yy}}^{-1/2} \Sigma_{\mathbf{yx}} \Sigma_{\mathbf{xx}}^{-1/2} & I \end{pmatrix}.$$

Notar que

- Si se estandarizan los vectores  $\mathbf{x}$  e  $\mathbf{y}$ ,

$$\mathbf{z} = \begin{pmatrix} \Sigma_{\mathbf{xx}}^{-1/2} (\mathbf{x} - E(\mathbf{x})) \\ \Sigma_{\mathbf{yy}}^{-1/2} (\mathbf{y} - E(\mathbf{y})) \end{pmatrix},$$

la matriz de covarianzas resultante para  $\mathbf{z}$  es

$$M = \begin{pmatrix} I & \Sigma_{\mathbf{xx}}^{-1/2} \Sigma_{\mathbf{xy}} \Sigma_{\mathbf{yy}}^{-1/2} \\ \Sigma_{\mathbf{yy}}^{-1/2} \Sigma_{\mathbf{yx}} \Sigma_{\mathbf{xx}}^{-1/2} & I \end{pmatrix}.$$

- Si  $M \begin{pmatrix} \mathbf{v}^t & \mathbf{w}^t \end{pmatrix}^t = \lambda \begin{pmatrix} \mathbf{v}^t & \mathbf{w}^t \end{pmatrix}^t$ ,  $0 < \lambda \neq 1$ , entonces

$$\begin{aligned} \Sigma_{\mathbf{xx}}^{-1/2} \Sigma_{\mathbf{xy}} \Sigma_{\mathbf{yy}}^{-1/2} \Sigma_{\mathbf{yx}} \left( \Sigma_{\mathbf{xx}}^{-1/2} \mathbf{v} \right) &= (\lambda - 1)^2 \left( \Sigma_{\mathbf{xx}}^{-1/2} \mathbf{v} \right) \\ \Sigma_{\mathbf{yy}}^{-1/2} \Sigma_{\mathbf{yx}} \Sigma_{\mathbf{xx}}^{-1/2} \Sigma_{\mathbf{xy}} \left( \Sigma_{\mathbf{yy}}^{-1/2} \mathbf{w} \right) &= (\lambda - 1)^2 \left( \Sigma_{\mathbf{yy}}^{-1/2} \mathbf{w} \right), \end{aligned}$$

- Luego, si  $(\mathbf{v}^t \ \mathbf{w}^t)^t$  es un autovector de  $M$ , entonces  $\Sigma_{\mathbf{xx}}^{-1/2}\mathbf{v}$  y  $\Sigma_{\mathbf{yy}}^{-1/2}\mathbf{w}$  son vectores canónicos relacionados con los grupos de variables  $\mathbf{x}$  e  $\mathbf{y}$ , respectivamente.

- Luego, si  $(\mathbf{v}^t \ \mathbf{w}^t)^t$  es un autovector de  $M$ , entonces  $\Sigma_{\mathbf{xx}}^{-1/2}\mathbf{v}$  y  $\Sigma_{\mathbf{yy}}^{-1/2}\mathbf{w}$  son vectores canónicos relacionados con los grupos de variables  $\mathbf{x}$  e  $\mathbf{y}$ , respectivamente.
- Y, si  $\lambda_j$  es un autovalor de  $M$  asociado con el autovector  $(\mathbf{v}_j^t \ \mathbf{w}_j^t)^t$ , entonces  $|\lambda_j - 1|$  es una correlación canónica relacionada con los grupos de variables  $\mathbf{x}$  e  $\mathbf{y}$ .

- Luego, si  $(\mathbf{v}^t \ \mathbf{w}^t)^t$  es un autovector de  $M$ , entonces  $\Sigma_{\mathbf{xx}}^{-1/2}\mathbf{v}$  y  $\Sigma_{\mathbf{yy}}^{-1/2}\mathbf{w}$  son vectores canónicos relacionados con los grupos de variables  $\mathbf{x}$  e  $\mathbf{y}$ , respectivamente.
- Y, si  $\lambda_j$  es un autovalor de  $M$  asociado con el autovector  $(\mathbf{v}_j^t \ \mathbf{w}_j^t)^t$ , entonces  $|\lambda_j - 1|$  es una correlación canónica relacionada con los grupos de variables  $\mathbf{x}$  e  $\mathbf{y}$ .

## Luego

ESTA CONEXIÓN NOS PERMITE CALCULAR LOS VECTORES CANÓNICOS A TRAVÉS DE UN ALGORITMO PARA COMPONENTES PRINCIPALES.

# Estimadores propuestos - SM

Para acotar el efecto de "*casos*" atípicos, se considera una M-escala robusta  $\sigma$  y se define los **SM-vectores canónicos estandarizados robustos** como

$$(A_{SMo}, B_{SMo}, \mathbf{a}_{SM}) = \arg \min_{(A_o, B_o, \mathbf{a}) \in \mathcal{B}_{r,m}} \sigma(A_o, B_o, \mathbf{a}) \quad (6)$$

con  $\sigma = \sigma(A_o, B_o, \mathbf{a})$  definida por la ecuación

$$\frac{1}{n} \sum_{i=1}^n \rho \left( \frac{\|A_o \tilde{\mathbf{x}}_i - B_o \tilde{\mathbf{y}}_i - \mathbf{a}\|^2}{\sigma} \right) = \delta, \quad (7)$$

donde  $\rho : [0, \infty) \rightarrow [0, 1]$ , es no decreciente y derivable,  $\rho(0) = 0$ ,  $\lim_{x \rightarrow \infty} \rho(x) = 1$  y  $0 < \delta < 1$ .  $\tilde{\mathbf{x}} = \left( \hat{\Sigma}_{\mathbf{xx}}^{(R)} \right)^{-1/2} \mathbf{x}$ ,  $\tilde{\mathbf{y}} = \left( \hat{\Sigma}_{\mathbf{yy}}^{(R)} \right)^{-1/2} \mathbf{y}$ .  
Y, los **SM-vectores canónicos robustos** se definen de la siguiente manera,

$$A_{SM} = A_{SMo} \left( \hat{\Sigma}_{\mathbf{xx}}^{(R)} \right)^{-1/2}, \quad B_{SM} = B_{SMo} \left( \hat{\Sigma}_{\mathbf{yy}}^{(R)} \right)^{-1/2}.$$

# Estimador propuesto - SM

Consistencia Fisher

Si tomamos funcionales de dispersión  $S_x$  y  $S_y$  que son consistentes Fisher para  $\Sigma_{xx}$  y  $\Sigma_{yy}$  respectivamente, Luego, si se considera  $\tilde{\mathbf{x}} = S_x^{-1/2}(\mathbf{x} - \boldsymbol{\mu}_x)$  e  $\tilde{\mathbf{y}} = S_y^{-1/2}(\mathbf{y} - \boldsymbol{\mu}_y)$  y buscamos las soluciones  $\sigma$  de

$$E \left[ \rho \left( \frac{\|A\tilde{\mathbf{x}} - B\tilde{\mathbf{y}} - \mathbf{a}\|^2}{\sigma} \right) \right] = \delta.$$

Si tomamos la descomposición espectral de  $M$  dada por  $M = \sum_{i=1}^{p+q} \gamma_i \mathbf{t}_i \mathbf{t}_i^t$ , con  $\gamma_1 > \dots > \gamma_r > \gamma_{r+1} \geq \dots \geq \gamma_{p+q-r} > \gamma_{p+q-r+1} > \dots > \gamma_{p+q}$ ,  $\mathbf{t}_i^t \mathbf{t}_j = \delta_{ij}$ ,  $1 \leq i, j \leq p+q$ , con

$$\delta_{ij} = \begin{cases} 1 & i = j \\ 0 & i \neq j \end{cases}$$

# Estimador propuesto - SM

Consistencia Fisher

Como  $\mathbf{t}_i = (\mathbf{v}_i^t, \mathbf{w}_i^t)^t$ ,  $\mathbf{v}_i \in \mathbb{R}^p$ ,  $\mathbf{w}_i \in \mathbb{R}^q$ ,  $i = 1, \dots, p + q$ , si llamamos

$$A_o = \left( \frac{\mathbf{v}_{p+q-r+1}}{\|\mathbf{v}_{p+q-r+1}\|}, \dots, \frac{\mathbf{v}_{p+q}}{\|\mathbf{v}_{p+q}\|} \right)^t \in \mathbb{R}^{r \times p},$$

$$B_o = \left( \frac{\mathbf{w}_{p+q-r+1}}{\|\mathbf{w}_{p+q-r+1}\|}, \dots, \frac{\mathbf{w}_{p+q}}{\|\mathbf{w}_{p+q}\|} \right)^t \in \mathbb{R}^{r \times q},$$

$$\mathbf{a}_o = A_o \Sigma_{xx}^{-1/2} \boldsymbol{\mu}_x - B_o \Sigma_{yy}^{-1/2} \boldsymbol{\mu}_y.$$

En el siguiente teorema se establece la consistencia Fisher del estimador SM para vectores canónicos para familias elípticas con función densidad

$$f(\mathbf{z}, \boldsymbol{\mu}_0, \Sigma_0) = \frac{1}{\sqrt{|\Sigma_0|}} f_0((\mathbf{z} - \boldsymbol{\mu}_0)^t \Sigma_0^{-1} (\mathbf{z} - \boldsymbol{\mu}_0)),$$

donde  $f_0 : \mathbb{R}^+ \rightarrow \mathbb{R}^+$  es una función no creciente.

Y si además se verifica que:

- $f_0$  y  $\rho$  cumplen que existe un punto  $d$  en un intervalo no degenerado  $I$  de la intersección de los dominios tal que

$$\rho(u) < \rho(d) < \rho(v) \quad \text{y} \quad f_0(u) > f_0(d) > f_0(v)$$

para todo  $u$  y  $v$  en  $I$ , con  $u < d < v$ .

- Teorema 1.** Sea  $\mathbf{Z}$  un vector aleatorio con distribución elíptica. Luego, los estimadores SM son estimadores funcionales consistentes Fisher, esto es

$$(A_o, B_o, \mathbf{a}_o) = \arg \min_{\mathbf{a} \in \mathbb{R}^r, AA^t = I, r = BB^t} \sigma(A, B, \mathbf{a}).$$

- Corolario 1:** Sea  $\mathbf{Z}$  un vector aleatorio con distribución elíptica. Si  $\rho$  es diferenciable, con  $\rho' = \psi$ , luego, para alguna constante  $c > 0$ ,

$$E_F \left[ \psi \left( \frac{\|A_o \Sigma_{\mathbf{xx}}^{-1/2} \mathbf{x} - B_o \Sigma_{\mathbf{yy}}^{-1/2} \mathbf{y} - \mathbf{a}_o\|^2}{\sigma(A_o, B_o, \mathbf{a}_o)} \right) (\mathbf{z} - \boldsymbol{\mu})(\mathbf{z} - \boldsymbol{\mu})^t \right] = c \Sigma.$$

Los Teoremas que se enuncian a continuación establecen la consistencia de los estimadores SM de vectores canónicos para familias elípticas.

- **Teorema 2.** *Sea  $\mathbf{Z}$  un vector aleatorio con densidad elíptica. Si  $m - r + 1 \leq n(1 - \delta)$ , el problema de encontrar  $(A, B, \mathbf{a}) \in B_{r,m}$  que minimicen  $\sigma(A, B, \mathbf{a})$  sujeto a (7) tiene al menos una solución con probabilidad 1.*
- **Teorema 3.** *Sea  $\mathbf{Z}$  un vector aleatorio con densidad elíptica. Sea  $(A^{(n)}, B^{(n)}, \mathbf{a}^{(n)})$  una solución del problema planteado en el Teorema 2, entonces el estimador SM definido en (6) es un funcional consistente, esto es,*

$$\lim_{n \rightarrow \infty} A^{(n)} = A_0 \quad , \quad \lim_{n \rightarrow \infty} B^{(n)} = B_0 \quad y \quad \lim_{n \rightarrow \infty} \mathbf{a}^{(n)} = \mathbf{a}_0,$$

*casí seguramente.*

# Estimador propuesto - SM

## Función de Influencia

Sea  $F$  una familia de distribuciones elípticas, si  $\rho$  es dos veces diferenciable y,  $S_x$  y  $S_y$  son funcionales consistentes Fisher para  $\Sigma_{xx}$  y  $\Sigma_{yy}$  respectivamente.

Se muestra a continuación las gráficas de la norma de la función de influencia para el primer vector canónico  $\mathbf{v}_1$ , para el caso en el que  $F$  sea  $N_{p+q}(\mathbf{0}, \Sigma)$  con

$$\Sigma = \begin{pmatrix} \Sigma_{xx} & \Sigma_{xy} \\ \Sigma_{yx} & \Sigma_{yy} \end{pmatrix},$$

tal que  $\Sigma_{xx} = I_2$ ,  $\Sigma_{yy} = I_2$  y  $\Sigma_{xy} = \begin{pmatrix} 0,9 & 0 \\ 0 & 0,5 \end{pmatrix}$ .

# Estimador propuesto - SM

## Función de Influencia

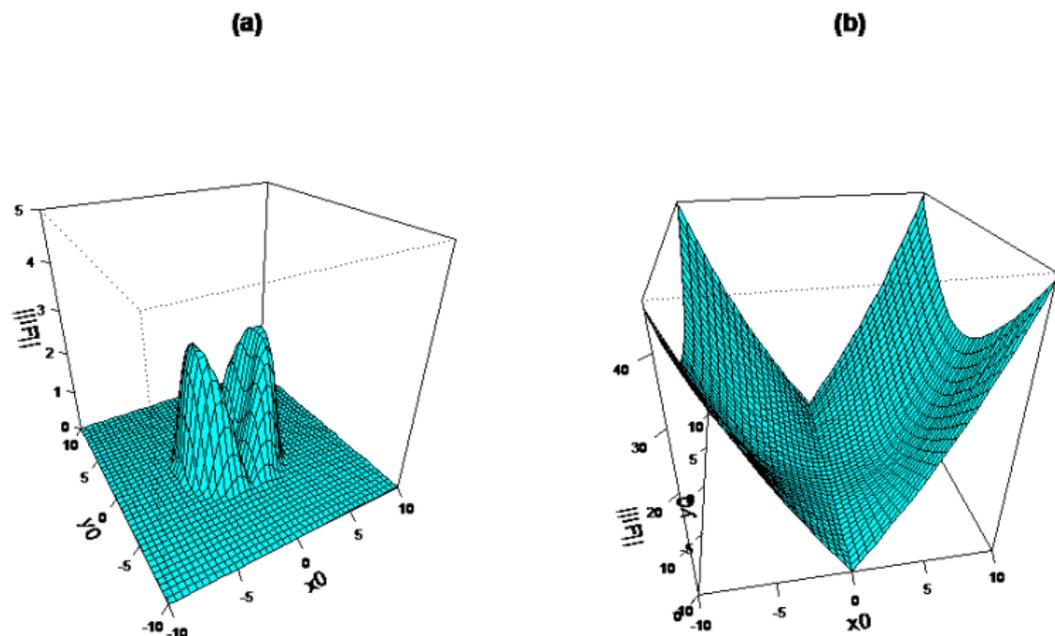


Figure: Norma de la Función de Influencia para  $(x_0, y_0, 0, 0)$  del primer vector canónico de (a) Estimador SM y (b) Estimador Clásico.

# Estimador propuesto - SM

## Función de Influencia

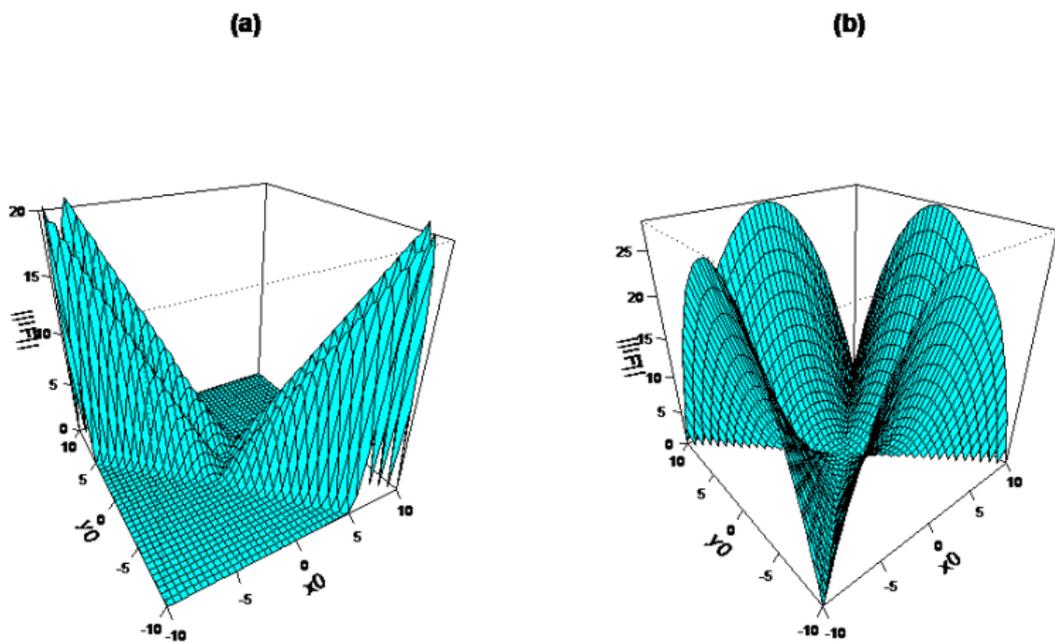


Figure: Norma de la Función de Influencia para  $(x_0, 0, y_0, 0)$  del primer vector canónico de (a) Estimador SM y (b) Estimador Clásico.

- Definimos Modelo de contaminaciones independientes cuando la contaminación puede afectar a las coordenadas de un vector en forma independiente (ver Alqallaf et al (2009) y Maronna and Yohai (2008)).

- Definimos Modelo de contaminaciones independientes cuando la contaminación puede afectar a las coordenadas de un vector en forma independiente (ver Alqallaf et al (2009) y Maronna and Yohai (2008)).
- Para el caso de PCA, dada la muestra  $\mathbf{z}_1, \dots, \mathbf{z}_n$  en  $\mathbb{R}^m$  (centrada), para encontrar vectores y componentes principales se debe encontrar  $\mathbf{v}_1, \dots, \mathbf{v}_q$  que forman una base del subespacio  $V$  de dimensión  $q < m$ , de forma de minimizar la siguiente expresión,

$$\begin{aligned} \sum_{i=1}^n \|\mathbf{z}_i - \hat{\mathbf{z}}_i\|^2 &= \sum_{i=1}^n \left\| \mathbf{z}_i - \sum_{j=1}^q (\mathbf{z}_i' \mathbf{v}_j) \mathbf{v}_j \right\|^2 \\ &= \sum_{i=1}^n \|\mathbf{z}_i - B\mathbf{a}_i\|^2. \end{aligned}$$

# Estimador propuesto - SMI

Para el caso de CCA, sean  $\mathbf{z}_i = (\mathbf{x}'_i, \mathbf{y}'_i)'$ ,  $i = 1, \dots, n$ , con covarianzas muestrales  $\hat{\Sigma}_{\mathbf{xx}}$  y  $\hat{\Sigma}_{\mathbf{yy}}$  respectivamente.

- Podemos escribir  $\hat{\Sigma}_{\mathbf{xx}}^{-1/2} \mathbf{x}_i = \sum_{k=1}^p \frac{(\mathbf{x}'_i \mathbf{t}_k)}{\sqrt{\lambda_k}} \mathbf{t}_k$ ,  $\lambda_1 > \dots > \lambda_p$  y  
 $\hat{\Sigma}_{\mathbf{yy}}^{-1/2} \mathbf{y}_i = \sum_{k=1}^q \frac{(\mathbf{y}'_i \mathbf{v}_k)}{\sqrt{\gamma_k}} \mathbf{v}_k$ ,  $\gamma_1 > \dots > \gamma_q$ .

# Estimador propuesto - SMI

Para el caso de CCA, sean  $\mathbf{z}_i = (\mathbf{x}'_i, \mathbf{y}'_i)'$ ,  $i = 1, \dots, n$ , con covarianzas muestrales  $\hat{\Sigma}_{\mathbf{xx}}$  y  $\hat{\Sigma}_{\mathbf{yy}}$  respectivamente.

- Podemos escribir  $\hat{\Sigma}_{\mathbf{xx}}^{-1/2} \mathbf{x}_i = \sum_{k=1}^p \frac{(\mathbf{x}'_i \mathbf{t}_k)}{\sqrt{\lambda_k}} \mathbf{t}_k$ ,  $\lambda_1 > \dots > \lambda_p$  y  $\hat{\Sigma}_{\mathbf{yy}}^{-1/2} \mathbf{y}_i = \sum_{k=1}^q \frac{(\mathbf{y}'_i \mathbf{v}_k)}{\sqrt{\gamma_k}} \mathbf{v}_k$ ,  $\gamma_1 > \dots > \gamma_q$ .
- Entonces cuando miramos la ecuación que queremos minimizar resulta

$$\begin{aligned} & \sum_{i=1}^n \left\| C \hat{\Sigma}_{\mathbf{xx}}^{-1/2} \mathbf{x}_i - B \hat{\Sigma}_{\mathbf{yy}}^{-1/2} \mathbf{y}_i - \mathbf{a} \right\|^2 \\ &= \sum_{i=1}^n \sum_{j=1}^r \left( C_j \sum_{k=1}^p \frac{(\mathbf{x}'_i \mathbf{t}_k)}{\sqrt{\lambda_k}} \mathbf{t}_k - B_j \sum_{k=1}^q \frac{(\mathbf{y}'_i \mathbf{v}_k)}{\sqrt{\gamma_k}} \mathbf{v}_k - a_j \right)^2 \\ &\approx \sum_{i=1}^n \sum_{j=1}^r \left( (\mathbf{C}_j, \mathbf{B}_j) \begin{pmatrix} T & 0 \\ 0 & V \end{pmatrix} (\mathbf{a}'_i, \mathbf{b}'_i)' - a_j \right)^2. \end{aligned}$$

- Para reducir la influencia de observaciones atípicas, proponemos usar un M-estimador de escala

$$(C^*, B^*, \boldsymbol{\mu}^*) = \arg \min_{\boldsymbol{\mu}, C, B} \sum_{j=1}^r \hat{\sigma}_j^2(\boldsymbol{\mu}, C, B),$$

- Para reducir la influencia de observaciones atípicas, proponemos usar un M-estimador de escala

$$(C^*, B^*, \boldsymbol{\mu}^*) = \arg \min_{\boldsymbol{\mu}, C, B} \sum_{j=1}^r \hat{\sigma}_j^2(\boldsymbol{\mu}, C, B),$$

- donde  $\hat{\sigma}_j = \hat{\sigma}_j(\boldsymbol{\mu}, C, B)$  satisface

$$g(\boldsymbol{\mu}_j, \mathbf{C}_j, \mathbf{B}_j) = \frac{1}{n} \sum_{i=1}^n \rho \left( \frac{(\mathbf{C}_j, \mathbf{B}_j) \begin{pmatrix} T & 0 \\ 0 & V \end{pmatrix} (\mathbf{a}'_i, \mathbf{b}'_i)' - \boldsymbol{\mu}_j}{\hat{\sigma}_j} \right) = \delta,$$

con  $\boldsymbol{\mu} \in \mathbb{R}^r$ ,  $C \in \mathbb{R}^{r \times p}$ ,  $B \in \mathbb{R}^{r \times q}$ , para  $r \leq \min(p, q)$ .

# Estudio de simulación

Medidas para evaluar el rendimiento de los estadísticos

Para realizar el estudio de simulación, se consideraron muestras aleatorias  $\mathbf{z}_j = \begin{pmatrix} \mathbf{x}_j \\ \mathbf{y}_j \end{pmatrix} \sim N_{p+q}(\mathbf{0}, \Sigma)$ , donde  $\Sigma$  que tiene todos los elementos diagonales iguales a 1 y los elementos fuera de la diagonal son todos iguales a  $h \in \{0.5, 0.6, 0.7, 0.8, 0.9\}$  (estructura de covarianzas utilizadas por Danilov et al (2012)).

Para evaluar el rendimiento de los estimadores, se consideró la siguiente medida propuesta por Branco et al (2005):

**Error cuadrático medio para vectores canónicos.**

$$MSE(\hat{\theta}_k) = \frac{1}{m} \sum_{j=1}^m \cos^{-1} \left( \frac{|\theta_{k,j}^t \hat{\theta}_{k,j}|}{\|\theta_{k,j}\| \|\hat{\theta}_{k,j}\|} \right),$$

siendo  $m$  la cantidad de replicaciones,  $\hat{\theta}_k$  el estimador del  $k$ -ésimo vector canónico,  $\theta_{k,j}$  el  $k$ -ésimo vector canónico

Se consideraron dos estructuras de contaminación:

- 1 Contaminación por celdas (independientes): celdas de la matriz de datos son elegidas al azar con probabilidad  $\varepsilon$  y su contenido reemplazado por un elemento proveniente de  $N_{p+q}(k\mathbf{v}_0, 0.5\Sigma)$ ,  $k \in \{1, \dots, 12\}$  y  $\mathbf{v}_0$  pertenece al subespacio de los autovectores de  $\Sigma$  asociados al menor autovalor.

Se consideraron dos estructuras de contaminación:

- 1 Contaminación por celdas (independientes): celdas de la matriz de datos son elegidas al azar con probabilidad  $\varepsilon$  y su contenido reemplazado por un elemento proveniente de  $N_{p+q}(k\mathbf{v}_0, 0.5\Sigma)$ ,  $k \in \{1, \dots, 12\}$  y  $\mathbf{v}_0$  pertenece al subespacio de los autovectores de  $\Sigma$  asociados al menor autovalor.
- 2 Contaminación por individuos (casos): El contenido de las primeras  $n\varepsilon$  filas de la matriz de datos son reemplazadas con elemento proveniente de  $N_{p+q}(k\mathbf{v}_0, 0.5\Sigma)$ ,  $k \in \{1, \dots, 12\}$  y  $\mathbf{v}_0$  pertenece al subespacio de los autovectores de  $\Sigma$  asociados al menor autovalor.

Se consideraron dos estructuras de contaminación:

- 1 Contaminación por celdas (independientes): celdas de la matriz de datos son elegidas al azar con probabilidad  $\varepsilon$  y su contenido reemplazado por un elemento proveniente de  $N_{p+q}(k\mathbf{v}_0, 0.5\Sigma)$ ,  $k \in \{1, \dots, 12\}$  y  $\mathbf{v}_0$  pertenece al subespacio de los autovectores de  $\Sigma$  asociados al menor autovalor.
  - 2 Contaminación por individuos (casos): El contenido de las primeras  $n\varepsilon$  filas de la matriz de datos son reemplazadas con elemento proveniente de  $N_{p+q}(k\mathbf{v}_0, 0.5\Sigma)$ ,  $k \in \{1, \dots, 12\}$  y  $\mathbf{v}_0$  pertenece al subespacio de los autovectores de  $\Sigma$  asociados al menor autovalor.
- **Tamaño de la muestra: 100.**

Se consideraron dos estructuras de contaminación:

- 1 Contaminación por celdas (independientes): celdas de la matriz de datos son elegidas al azar con probabilidad  $\varepsilon$  y su contenido reemplazado por un elemento proveniente de  $N_{p+q}(k\mathbf{v}_0, 0.5\Sigma)$ ,  $k \in \{1, \dots, 12\}$  y  $\mathbf{v}_0$  pertenece al subespacio de los autovectores de  $\Sigma$  asociados al menor autovalor.
  - 2 Contaminación por individuos (casos): El contenido de las primeras  $n\varepsilon$  filas de la matriz de datos son reemplazadas con elemento proveniente de  $N_{p+q}(k\mathbf{v}_0, 0.5\Sigma)$ ,  $k \in \{1, \dots, 12\}$  y  $\mathbf{v}_0$  pertenece al subespacio de los autovectores de  $\Sigma$  asociados al menor autovalor.
- Tamaño de la muestra: 100.
  - **Número de replicaciones: 150.**

Se consideraron dos estructuras de contaminación:

- 1 Contaminación por celdas (independientes): celdas de la matriz de datos son elegidas al azar con probabilidad  $\varepsilon$  y su contenido reemplazado por un elemento proveniente de  $N_{p+q}(k\mathbf{v}_0, 0.5\Sigma)$ ,  $k \in \{1, \dots, 12\}$  y  $\mathbf{v}_0$  pertenece al subespacio de los autovectores de  $\Sigma$  asociados al menor autovalor.
  - 2 Contaminación por individuos (casos): El contenido de las primeras  $n\varepsilon$  filas de la matriz de datos son reemplazadas con elemento proveniente de  $N_{p+q}(k\mathbf{v}_0, 0.5\Sigma)$ ,  $k \in \{1, \dots, 12\}$  y  $\mathbf{v}_0$  pertenece al subespacio de los autovectores de  $\Sigma$  asociados al menor autovalor.
- Tamaño de la muestra: 100.
  - Número de replicaciones: 150.
  - $p = q = 5$ .

Se consideraron dos estructuras de contaminación:

- 1 Contaminación por celdas (independientes): celdas de la matriz de datos son elegidas al azar con probabilidad  $\varepsilon$  y su contenido reemplazado por un elemento proveniente de  $N_{p+q}(k\mathbf{v}_0, 0.5\Sigma)$ ,  $k \in \{1, \dots, 12\}$  y  $\mathbf{v}_0$  pertenece al subespacio de los autovectores de  $\Sigma$  asociados al menor autovalor.
  - 2 Contaminación por individuos (casos): El contenido de las primeras  $n\varepsilon$  filas de la matriz de datos son reemplazadas con elemento proveniente de  $N_{p+q}(k\mathbf{v}_0, 0.5\Sigma)$ ,  $k \in \{1, \dots, 12\}$  y  $\mathbf{v}_0$  pertenece al subespacio de los autovectores de  $\Sigma$  asociados al menor autovalor.
- Tamaño de la muestra: 100.
  - Número de replicaciones: 150.
  - $p = q = 5$ .
  - $\varepsilon = 0.2$ .

# Estudio de simulación

Los estimadores que se consideran en el estudio de simulación son:

- **Estimador clásico.**

Los estimadores que se consideran en el estudio de simulación son:

- **Estimador clásico.**
- **Estimador robusto de 2 pasos para la matriz de covarianzas (Agostinelli et al (2015)).**

Los estimadores que se consideran en el estudio de simulación son:

- **Estimador clásico.**
- **Estimador robusto de 2 pasos para la matriz de covarianzas (Agostinelli et al (2015)).**
- **Estimador SM (Adrover and Donato (2015)).** Para computar la escala se usó la función (Maronna (2005) que hace convergente al algoritmo iterativo)

$$\rho(t) = \min \left\{ 1, 1 - (1 - t)^3 \right\},$$

con  $\delta = 0,5$ .

Los estimadores que se consideran en el estudio de simulación son:

- **Estimador clásico.**
- **Estimador robusto de 2 pasos para la matriz de covarianzas (Agostinelli et al (2015)).**
- **Estimador SM (Adrover and Donato (2015)).** Para computar la escala se usó la función (Maronna (2005) que hace convergente al algoritmo iterativo)

$$\rho(t) = \min \left\{ 1, 1 - (1 - t)^3 \right\},$$

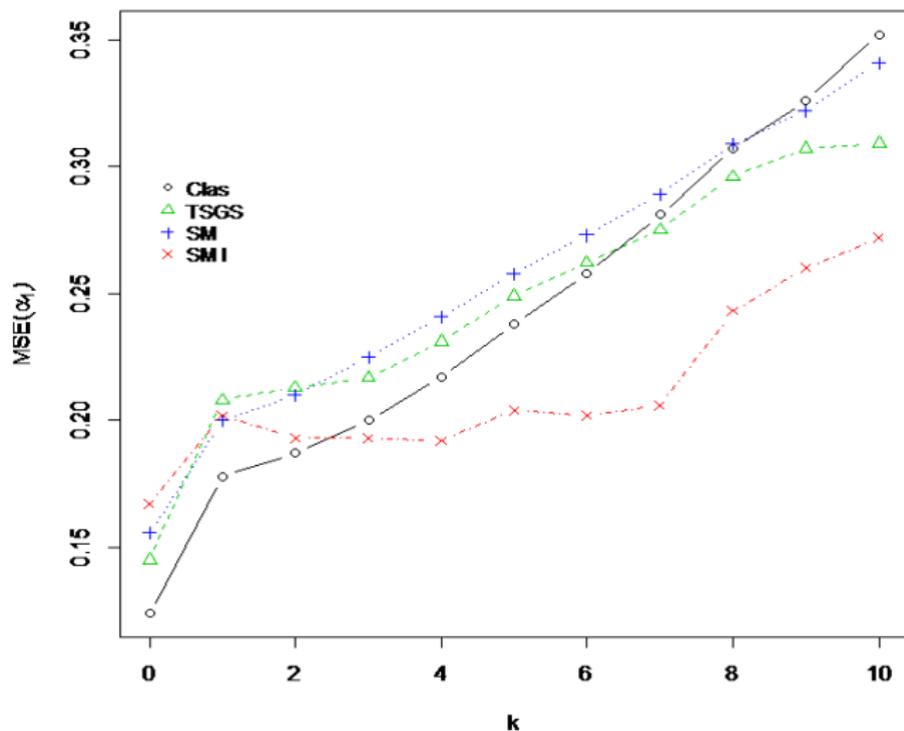
con  $\delta = 0,5$ .

- **Estimador SMI.** Para computar el algoritmo se utilizó la función bicuadrada de Tukey

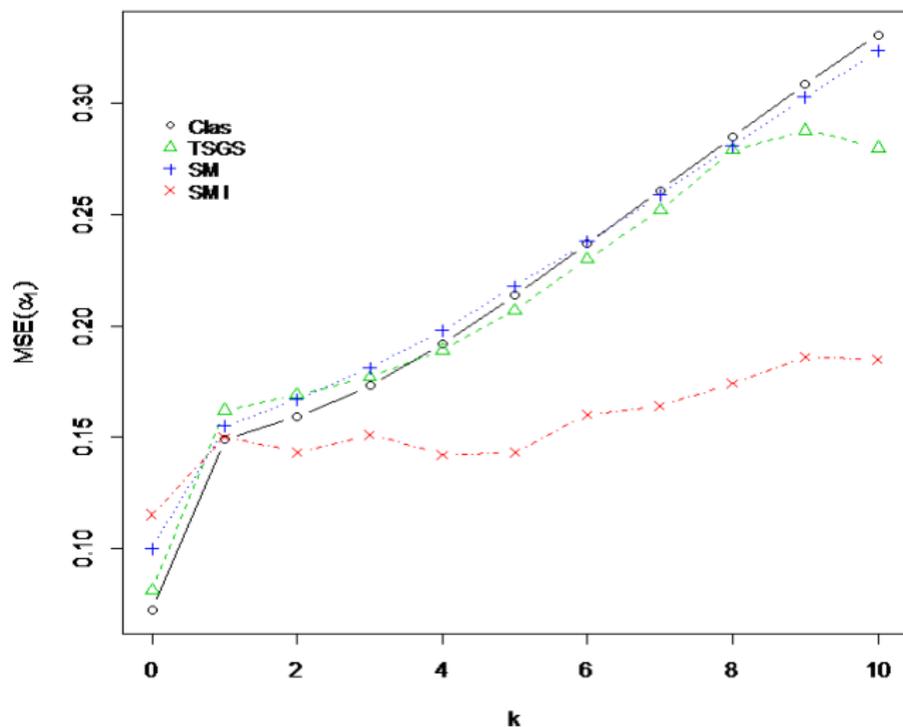
$$\rho(t) = \min \left\{ 1, 1 - \left( 1 - \left( \frac{t}{c} \right)^2 \right)^3 \right\},$$

con  $c = 1.54764$  y  $\delta = 0,5$ . (Boente and Salibian Barrera (2015)).

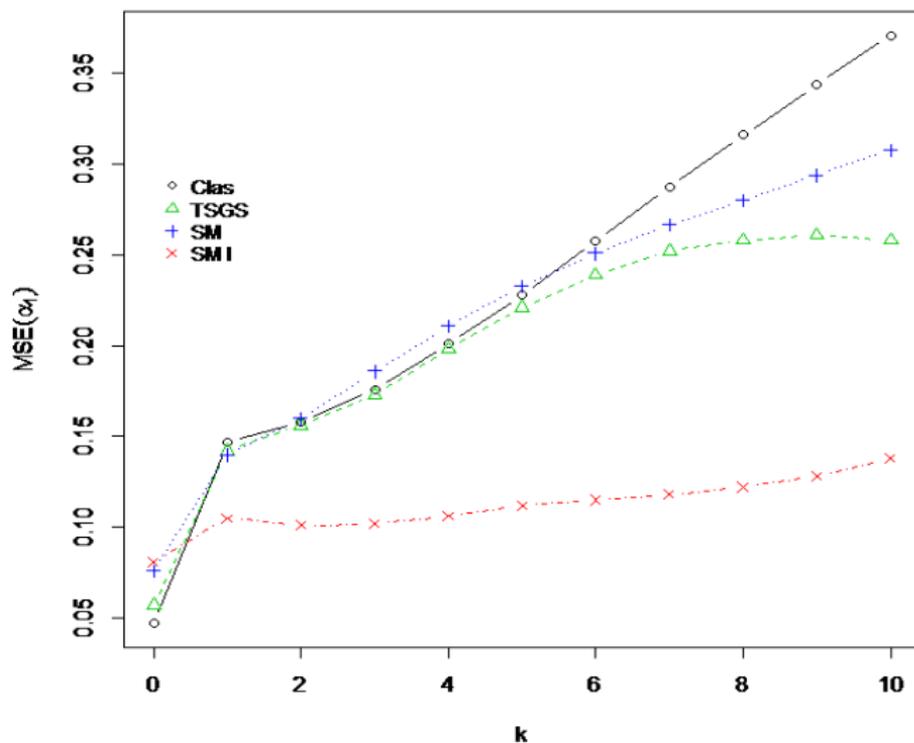
Contaminación por celdas ( $h=0.6$ )



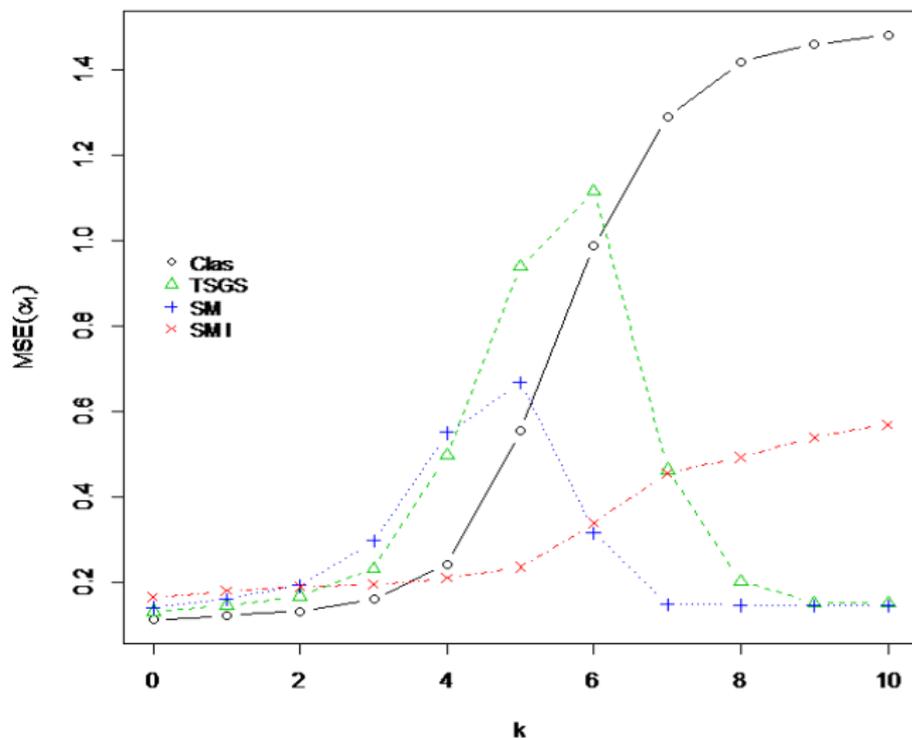
## Contaminación por celdas ( $h=0.8$ )



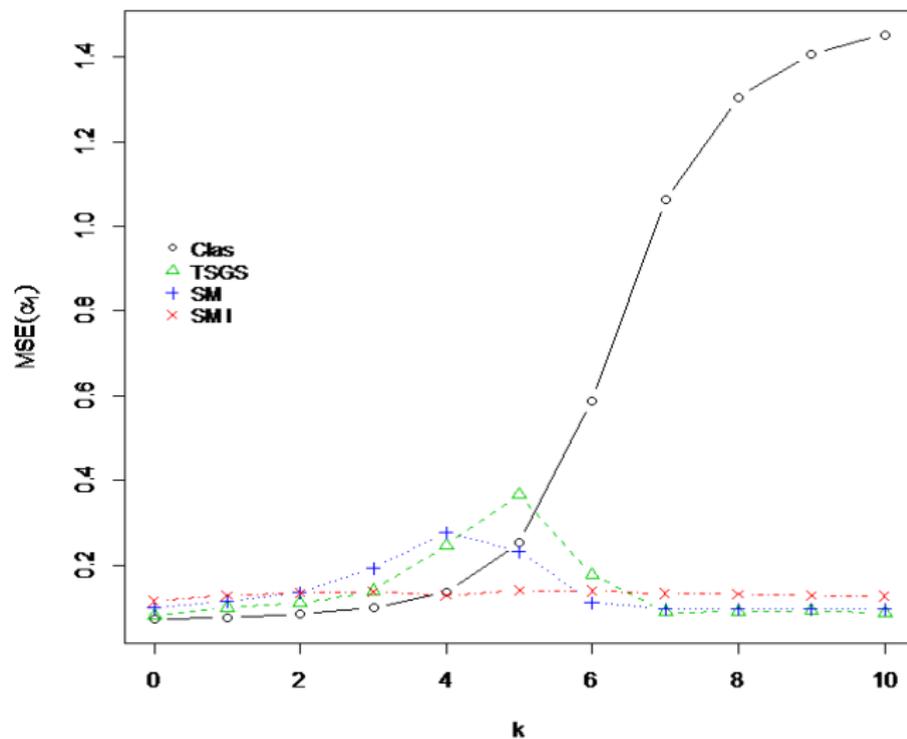
## Contaminación por celdas ( $h=0.9$ )



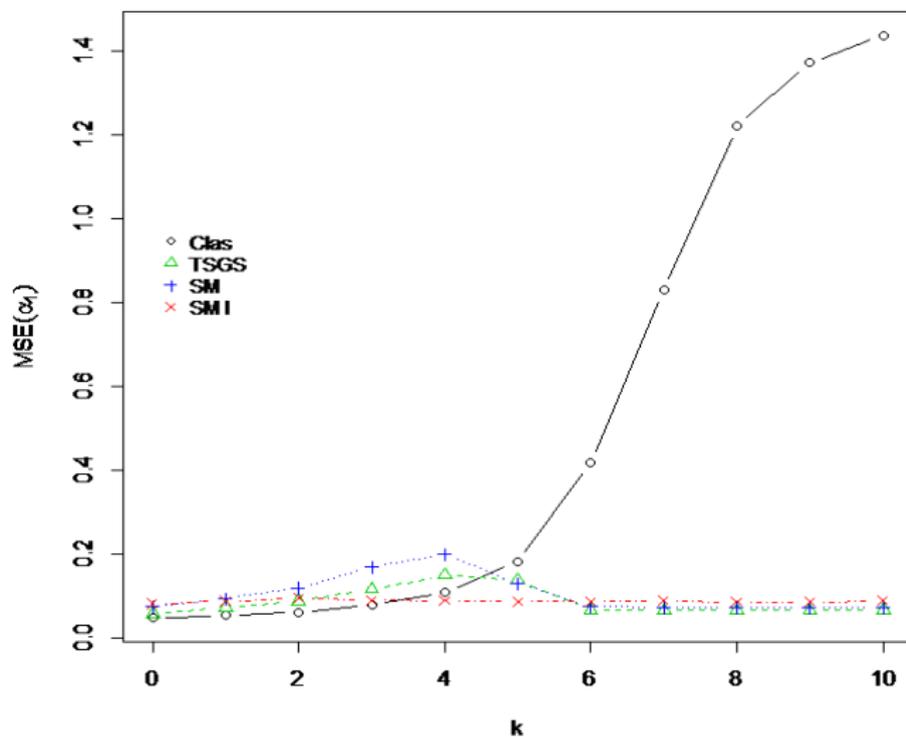
## Contaminación por individuos ( $h=0.6$ )



## Contaminación por individuos ( $h=0.8$ )



## Contaminación por individuos ( $h=0.9$ )



- 1 Se presentaron dos estimadores robustos para vectores canónicos, uno para el modelo de contaminaciones por celdas y otro para contaminaciones por individuo

- 1 Se presentaron dos estimadores robustos para vectores canónicos, uno para el modelo de contaminaciones por celdas y otro para contaminaciones por individuo
- 2 Se exploró la relación CCA - PCA que permitió desarrollar los algoritmos iterativos.

- 1 Se presentaron dos estimadores robustos para vectores canónicos, uno para el modelo de contaminaciones por celdas y otro para contaminaciones por individuo
- 2 Se exploró la relación CCA - PCA que permitió desarrollar los algoritmos iterativos.
- 3 Se comparó el rendimiento de los estimadores propuestos a través de un estudio de simulación.

- 1 Se presentaron dos estimadores robustos para vectores canónicos, uno para el modelo de contaminaciones por celdas y otro para contaminaciones por individuo
- 2 Se exploró la relación CCA - PCA que permitió desarrollar los algoritmos iterativos.
- 3 Se comparó el rendimiento de los estimadores propuestos a través de un estudio de simulación.
- 4 El estimador SMI muestra el mejor comportamiento global para los dos modelos de contaminación analizados.

- 1 Se presentaron dos estimadores robustos para vectores canónicos, uno para el modelo de contaminaciones por celdas y otro para contaminaciones por individuo
- 2 Se exploró la relación CCA - PCA que permitió desarrollar los algoritmos iterativos.
- 3 Se comparó el rendimiento de los estimadores propuestos a través de un estudio de simulación.
- 4 El estimador SMI muestra el mejor comportamiento global para los dos modelos de contaminación analizados.
- 5 Los estimadores SM y TSGS tienen mejor rendimiento en el modelo de contaminación por individuo.

- 1 Se presentaron dos estimadores robustos para vectores canónicos, uno para el modelo de contaminaciones por celdas y otro para contaminaciones por individuo
- 2 Se exploró la relación CCA - PCA que permitió desarrollar los algoritmos iterativos.
- 3 Se comparó el rendimiento de los estimadores propuestos a través de un estudio de simulación.
- 4 El estimador SMI muestra el mejor comportamiento global para los dos modelos de contaminación analizados.
- 5 Los estimadores SM y TSGS tienen mejor rendimiento en el modelo de contaminación por individuo.
- 6 **Muchas gracias!!!**

- Adrover, J. G., Donato, S. M. (2015) "A robust predictive approach for canonical correlation analysis". *Journal of Multivariate Analysis*. 133: 356-376.
- Agostinelli, C., Leung, A., Yohai, V. and Zamar, R. (2015) "Robust estimation of multivariate location and scatter in the presence of cellwise and casewise contamination" *Test*. 24(3): 441-461.
- Alqallaf, F., Van Aelst, S., Yohai, V. J. and Zamar, R. H. (2009). "Propagation of outliers in multivariate data". *The Annals of Statistics* 37(1): 311-331.
- Boente, G. and Salibian Barrera, M. (2015). "S-estimators for functional principal component analysis". *Journal of the American Statistical Association* 110(511): 1100-1111.
- Brillinger, David. (1975). *Time Series: Data Analysis and Theory*. Estados Unidos de América. Holt, Rinehart and Winston.

- Branco, J.A., Croux, C., Filzmoser, P., and Oliveira, M.R. (2005) "Robust Canonical Correlations: A Comparative Study". *Computational Statistics*, 20(2): 203-229.
- Danilov, M., Yohai, V. and Zamar, R. (2012) "Robust Estimation of Multivariate Location and Scatter in the presence of Missing Data". *Journal of the American Statistical Association*, 107(499): 1178-1186.
- Maronna, R.A. (2005) "Principal Components and Orthogonal Regression Based on Robust Scales". *Technometrics*, 47(3): 264-273.
- Maronna, R. A. and Yohai, V. J. (2008) "Robust Low-Rank Approximation of Data Matrices With Elementwise Contamination" *Technometrics*. 50(3): 295-304.
- Seber, G. A. F. (1984). *Multivariate Observations*. Estados Unidos de América. John Wiley & Sons, Inc.