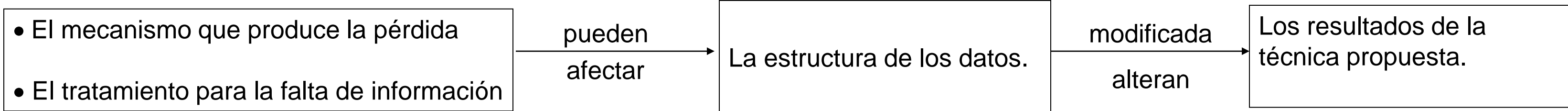


# INFLUENCIA DEL MECANISMO, DE LA ESTRUCTURA DE LOS DATOS Y DEL TRATAMIENTO DE LA INFORMACIÓN INCOMPLETA EN EL ANÁLISIS ESTADÍSTICO DE LA EPH

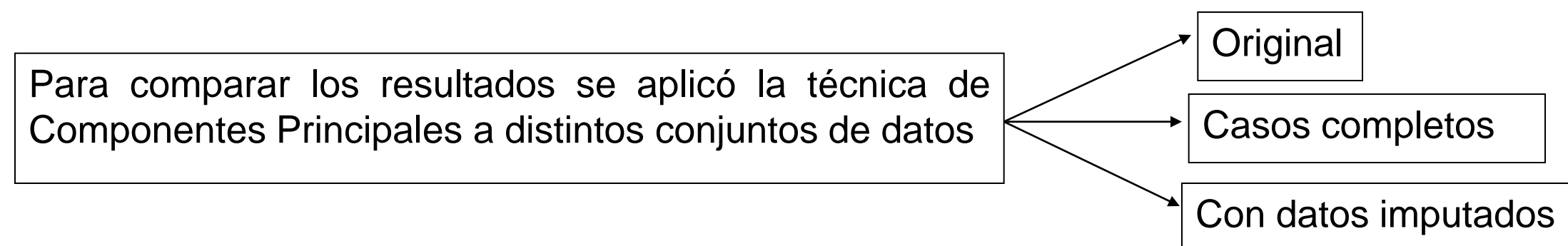
Badler, Clara; Alsina, Sara; Beltrán, Celina; Puigsubirá, Cristina; Vitelleschi, Ma. Susana

Instituto de Investigaciones Teóricas y Aplicadas de la Escuela de Estadística, Fac. Cs. Económicas y Estadística. Universidad Nacional de Rosario

## INTRODUCCIÓN

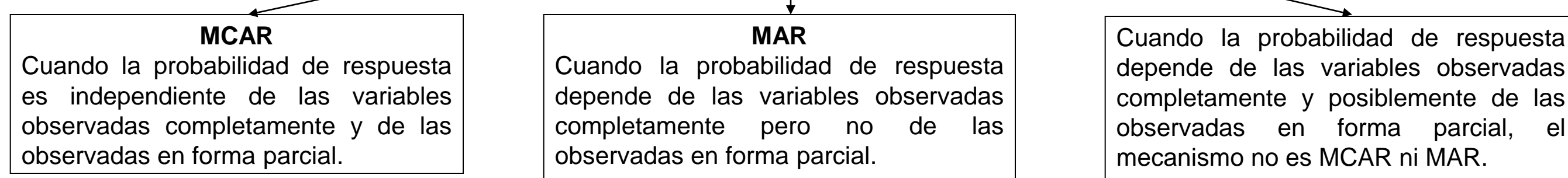


Los procedimientos de simulación de pérdidas y el uso posterior de técnicas de análisis multivariado permiten describir e interpretar un determinado fenómeno, considerando diferentes situaciones.



## METODOLOGÍA

### MECANISMOS DE PÉRDIDA



Para la evaluación del supuesto MCAR se utiliza un test propuesto por Little, aplicado a un esquema monótono con tres variables, en el que se observa:

- $Y_1$  en forma completa ( $n$  unidades)
- $Y_2$  en  $n_2$  ( $<n$ ) unidades
- $Y_3$  en  $n_3$  ( $<n_2$ ) unidades

$$d^2 = \frac{SSB_1}{MST_1} + \frac{SSB_{2,1}}{MST_{2,1}} = \frac{(n-1)2F_1}{(n-3)+2F_1} + \frac{(n_2-1)F_{2,1}}{(n_2-2)+F_{2,1}} \sim \chi_3^2$$

- $SSB_1$ ,  $MST_1$  y  $F_1$  son la suma de cuadrados entre grupos, el cuadrado medio total y la estadística F del análisis de la variancia de  $Y_1$  sobre el esquema de pérdida
- $SSB_{2,1}$ ,  $MST_{2,1}$  y  $F_{2,1}$  son la suma de cuadrados entre grupos, el cuadrado medio total y la estadística F del análisis de covariancia de la variable  $Y_2$  sobre los grupos determinados por los esquemas donde  $Y_2$  es observada, ajustando por  $Y_1$

### ESTUDIO DE LA ESTRUCTURA DE LOS DATOS

- Se divide el conjunto de datos originales en grupos con respecto a una determinada variable ( $Y_j$ ). En este caso, la variable presenta pérdida de información simulada.
- Se utiliza como punto de corte para dicha división los percentiles de la variable  $Y_j$ .
- Se calcula la correlación entre las variables que presentan información incompleta en cada uno de los subconjuntos de datos resultantes de la división y la matriz de correlaciones para el conjunto completo de datos.
- Se analiza la tendencia que presentan las correlaciones en los distintos grupos.
- Se investiga de qué forma se ven afectados los resultados obtenidos de la aplicación de técnicas estadísticas basadas en las correlaciones.

### TRATAMIENTO DE LA INFORMACION FALTANTE

**CASOS COMPLETOS:** Consiste en utilizar sólo las unidades con información completa en todas las variables. Es adecuado bajo el supuesto MCAR.

**IMPUTACIÓN POR EL MÉTODO DE BUCK:** Asigna al valor perdido el obtenido como predicción de la regresión lineal sobre las variables observadas en esa unidad:

$$\hat{Y}_{ij,per} = \hat{\beta}' Y_{obs,j}$$

\* $Y_{ij,per}$  valor faltante de la variable  $Y_j$  en la  $i$ -ésima unidad  
 •  $Y_{obs,i}$  vector de los valores observados en dicha unidad  
 •  $\hat{\beta}$  vector de coeficientes estimados de regresión.  
 Realiza el supuesto MCAR, sin embargo los estimadores obtenidos pueden ser consistentes aún cuando el mecanismo sea MAR.

### COMPONENTES PRINCIPALES

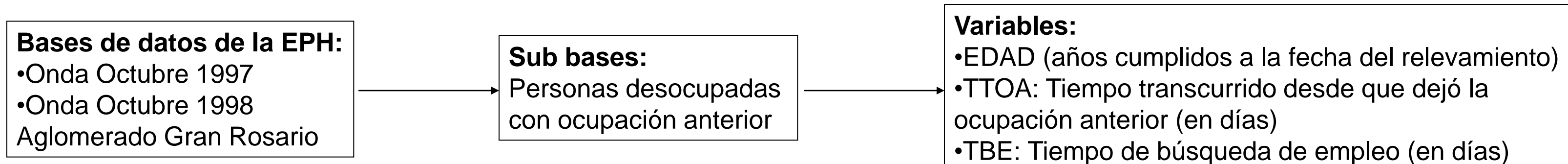
Transforma el conjunto de datos en uno de menor dimensión a través de nuevas variables que expliquen el máximo de la variación de las variables originales

$$CP_j = a_{1j} Y_1 + a_{2j} Y_2 + \dots + a_{pj} Y_p \quad j=1,2,\dots$$

$$Var(CP_j) = a_j' \Sigma a_j \quad j=1,2,\dots$$

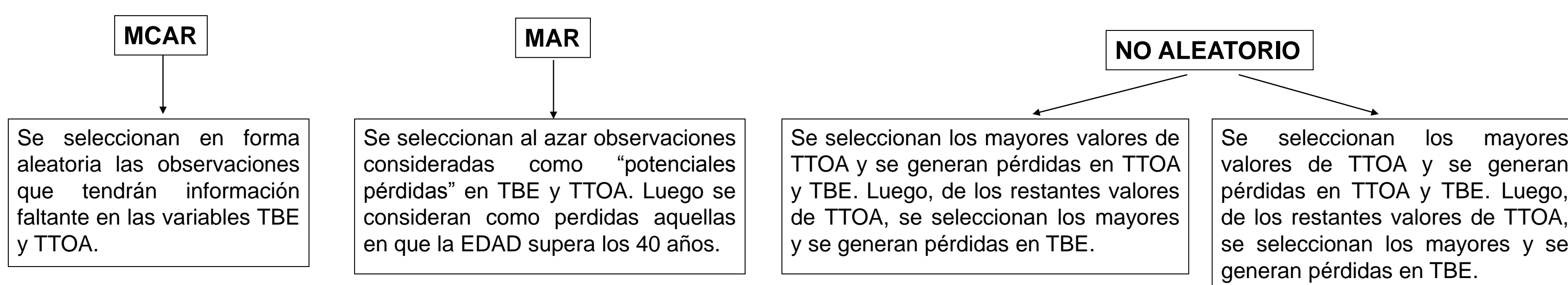
$\Sigma$  = matriz de variancias y covariancias de las variables  $Y$ .  
 $a_j$  = autovector normalizado asociado al  $j$ -ésimo autovalor de la matriz de correlaciones.

## RESULTADOS



### SIMULACIÓN DE PERDIDA DE INFORMACIÓN

Mediante el programa SAS se generan pérdidas para obtener un esquema monótono en el cual la variable EDAD se observa en forma completa, TTOA presenta 10% de pérdida y TBE un 20%. Las simulaciones se realizan bajo distintos mecanismos de pérdida:



Cuadro 1: Test para evaluar el supuesto MCAR

ONDA	MECANISMO DE PERDIDA	F <sub>1</sub>	F <sub>2,1</sub>	d <sup>2</sup> <sub>(GL=3)</sub>	PROB	DECISIÓN*
Octubre 1997 (N=171)	MCAR	1.1636	0.639	2.96	0.40	No se rechaza
	MAR	35.9433	0.031	50.975	4.95E <sup>-11</sup>	Se rechaza
	Valores máx. Valores mín.	0.7435 0.01	182.77 8.13	85.02 7.9	2.57E <sup>-18</sup> 0.04	Se rechaza Se rechaza
Octubre 1998 (N=142)	MCAR	0.0518	0.456	0.56	0.90	No se rechaza
	MAR	25.65	2.346	40.33	9.06E <sup>-9</sup>	Se rechaza
	Valores máx. Valores mín.	2.82 0.1035	400.28 0.759	102.09 0.97	5.52E <sup>-22</sup> 0.81	Se rechaza No se rechaza

\* $\alpha=0.05$

- Para la onda Octubre de 1997, en todo los mecanismos no MCAR, se rechaza la hipótesis
- Para la onda Octubre de 1998, en el mecanismo no aleatorio por el cual se pierde información en unidades con valores bajos de TTOA, no se rechaza la hipótesis.

### CAMBIOS DETECTADOS EN LA CORRELACIÓN TTOA\*TBE

Cuadro 2: Cambios producidos en la correlación TTOA\*TBE. Onda Octubre de 1997

Conjunto de datos	Correlación TTOA*TBE
Datos originales	0.43
Mec. No aleatorio	
Valores Máx. (Casos Completos)	0.85
Mec. No aleatorio	
Valores Máx. (Imputación Buck)	0.92

Cuadro 3: Cambios producidos en la correlación TTOA\*TBE. Onda Octubre de 1998

Conjunto de datos	Correlación TTOA*TBE
Datos originales	0.02
Mec. No aleatorio	
Valores Máx. (Casos Completos)	0.78
Mec. No aleatorio	
Valores Máx. (Imputación Buck)	0.93

La correlación entre las variables que presentan pérdidas se ve sobrestimada al trabajar con los casos completos como al imputar por el método de Buck.

### ESTRUCTURA DE LOS DATOS

Se divide la muestra en grupos definidos por los percentiles 20, 40, 60 y 80 de la variable TTOA

Gráfico 1: Correlación TTOA\*TBE según grupo de TTOA. Onda Octubre de 1997

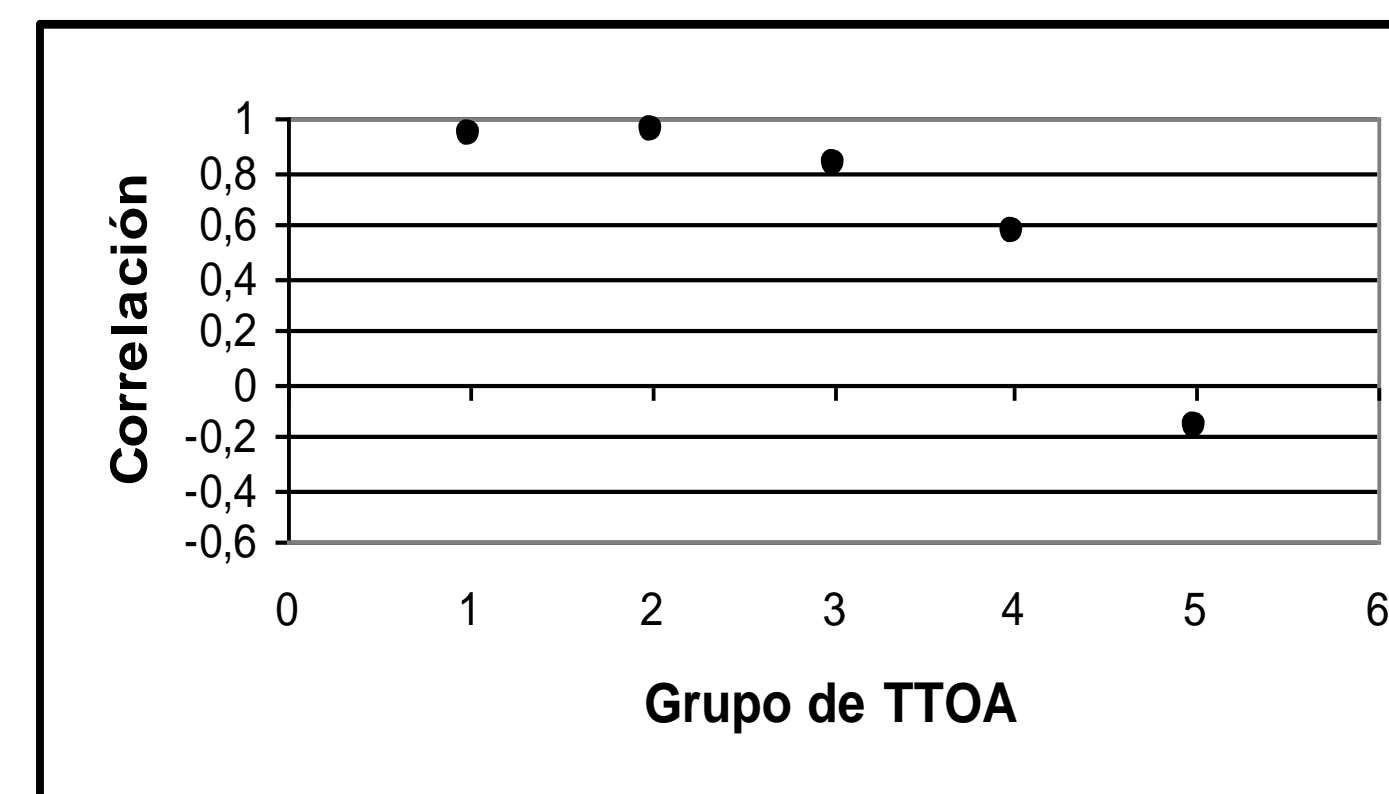
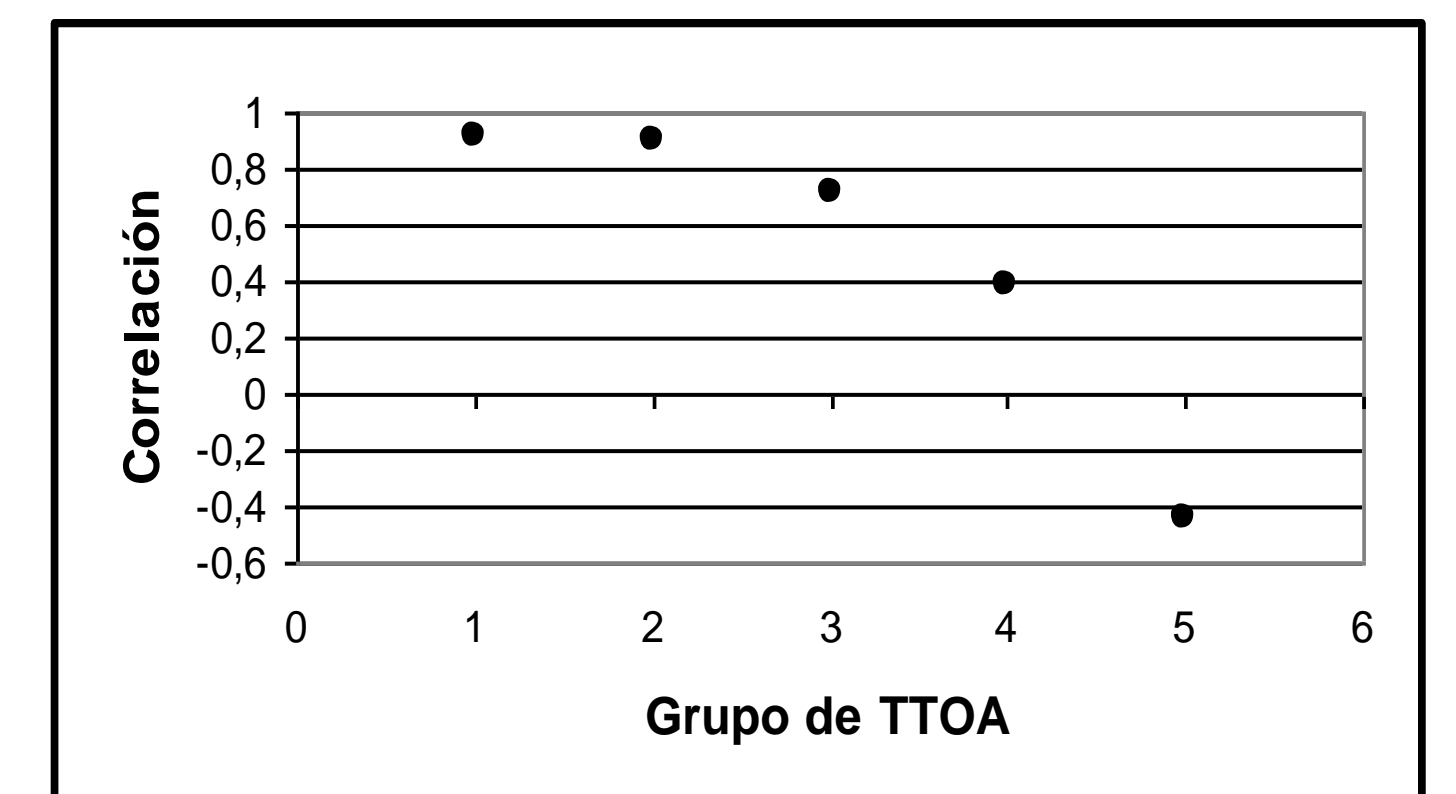


Gráfico 2: Correlación TTOA\*TBE según grupo de TTOA. Onda Octubre de 1998



La correlación entre las variables TTOA y TBE disminuye al aumentar el valor de la variable TTOA.

Al perder los valores más altos de TTOA, se utiliza información de las unidades que presentan una correlación mayor que la de la muestra original.

Sobrestimación de la correlación.

### APLICACIÓN DE COMPONENTES PRINCIPALES

Se presentan los resultados para los mecanismos MCAR y NO ALEATORIO (Valores máx.), ya que en estas situaciones los mismos presentan diferencias.

#### MECANISMO MCAR

Cuadro 4: Coeficientes y porcentaje de variancia explicada de las dos primeras componentes, onda Octubre 1997.

Variable	COEFICIENTES					
	DATOS ORIGINALES		METODO DE BUCK		CASOS COMPLETOS	
	CP1	CP2	CP1	CP2	CP1	CP2
EDAD	0.37	0.92	0.36	0.93	0.35	0.93
TBE	0.67	-0.15	0.67	-0.17	0.67	-0.17
TTOA	0.64	-0.37	0.65	-0.34	0.65	-0.32
% de variancia explicada	81.2		80.2		80.1	

Cuadro 5: Coeficientes y porcentaje de variancia explicada de las dos primeras componentes, onda Octubre 1998.

Variable	COEFICIENTES					
	DATOS ORIGINALES		METODO DE BUCK		CASOS COMPLETOS	
	CP1	CP2	CP1	CP2	CP1	CP2
EDAD	0.69	-0.05	0.69	-0.03	0.70	-0.03
TBE	0.40	0.86	0.43	0.81	0.39	0.85
TTOA	0.60	-0.51	0.58	-0.58	0.60	-0.52
% de variancia explicada	73.3		72.5		72.4	

#### MECANISMO NO ALEATORIO (Valores máx.)

Cuadro 6: Coeficientes y porcentaje de variancia explicada de las dos primeras componentes, onda Octubre 1997.

Variable	COEFICIENTES					
	DATOS ORIGINALES		METODO DE BUCK		CASOS COMPLETOS	
	CP1	CP2	CP1	CP2	CP1	CP2
EDAD	0.37	0.92	0.24	0.97	0.26	0.97
TBE	0.67	-0.15	0.69	-0.15	0.69	-0.15
TTOA	0.64	-0.37	0.68	-0.19	0.68	-0.21
% de variancia explicada	81.2		95.0		97.5	

Cuadro 7: Coeficientes y porcentaje de variancia explicada de las dos primeras componentes, onda Octubre 1998.

Variable	COEFICIENTES					
	DATOS ORIGINALES		METODO DE BUCK		CASOS COMPLETOS	
	CP1	CP2	CP1	CP2	CP1	CP2
EDAD	0.69	0.05	-0.02	0.99	0.05	0.99
TBE	0.40	-0.86	0.71	0.02	0.71	-0.01
TTOA	0.60	0.51	0.71	0.01	0.71	-0.05
% de variancia explicada	73.3		97.8		92.7	

Se puede observar como los resultados de la aplicación de componentes principales se ven afectados cuando se ha modificado la estructura de las relaciones entre las variables.

## DISCUSIÓN

Al analizar bases de datos con información incompleta debe considerarse:

- la naturaleza de las pérdidas
- los supuestos sobre las características del mecanismo que las produjo
- el tratamiento aplicado para "completar" los valores perdidos
- la estructura de los datos originales y la resultante de la aplicación del tratamiento seleccionado.