

APLICACIÓN DE LAS MEDIDAS DE CONCORDANCIA AL ANÁLISIS FACTORIAL DE CORRESPONDENCIAS

Ricardo García Ródenas,
M^a Luz López García,
Antonio Martínez Plaza,
Doroteo Verástegui Rayo

Antonio Martínez Plaza está en la Escuela Universitaria Politécnica de Albacete, Universidad de Castilla-La Mancha. Avda. España, s/n. 02071 Albacete. Ricardo García Ródenas, M^a Luz López García y Doroteo Verástegui Rayo están en Escuela Universitaria Politécnica de Almadén. Univ. de Castilla-La Mancha. Plaza Manuel Meca, 1. 13400 Almadén (Ciudad Real).

RESUMEN

En este trabajo formulamos tres problemas de optimización cuyas funciones objetivos se basan en las medidas de concordancia kappa de Cohen, alpha de Aicken y los modelos log-lineales. La solución óptimas de estos modelos nos permiten estimar que «caracteres» en Análisis de Correspondencias (ACF) son propios de una o varias «poblaciones». Esto complementa los objetivos básicos del ACF.

1. INTRODUCCIÓN

EL Análisis Factorial de Correspondencias (ACF), ver Cuadras (1991), es apropiado para representar geoméricamente las tablas de contingencias. Supongamos que los datos de interés corresponden a dos criterios de clasificación, a los que llamaremos «caracteres» y «poblaciones» las cuales se disponen como en la Tabla 1 de contingencia:

Uno de los objetivos del ACF es obtener una representación geométrica de las poblaciones H_1, \dots, H_k en relación de la distribución de frecuencias relativas de los «caracteres» A_1, \dots, A_m .

Este problema de representación de las poblaciones en dimensión reducida determinada por las coordenadas con referencia a los caracteres A_j se puede interpretar como un problema de representación de datos mediante análisis de componentes principales.

TABLA 1.

Poblaciones	Caracteres				
	A_1	A_2	...	A_m	Total
H_1	p_{11}	p_{12}	...	p_{1m}	$p_{1\cdot}$
H_2	p_{21}	p_{22}	...	p_{2m}	$p_{2\cdot}$
			...		
H_k	p_{k1}	p_{k2}	...	p_{km}	$p_{k\cdot}$
Total	$p_{\cdot 1}$	$p_{\cdot 2}$...	$p_{\cdot m}$	1

El objetivo de este trabajo es presentar varios criterios que permiten establecer la relación existente entre una población o un subconjunto de poblaciones afines y un subconjunto de caracteres característicos de ellos.

Los coeficientes de concordancia son especialmente útiles para medir el grado de asociación entre la relación de los dos conjuntos de atributos. La medida más popular para evaluar la concordancia entre varias variables cualitativas es el coeficiente kappa de Cohen, [Cohen (1960)], quien lo formuló para dos variables binarias. Este estadístico se desarrolló originariamente para medir la concordancia entre observadores, aplicándose particularmente a problemas de fiabilidad entre diagnósticos. En este contexto, se han presentado sucesivas generalizaciones [Cohen (1968), Fleiss (1971), Landis et al (1977), Fleiss et al. (1979), Chmura (1980)]. Todas estas generalizaciones se basan en que todas las variables cualitativas deben medirse en el mismo sistema de categorías, cuestión inherente al propio objetivo de medir la concordancia entre observadores. Su aplicación se ha extendido lejos de este problema específico, de hecho se emplea como una medida *similitud* para datos categóricos [Fleiss (1981)].

Varios autores [Feiss et al. (1973)] interpretan este coeficiente como una medida de correlación entre las categorías de dos variables nominales. En este contexto planteamos su aplicación.

Otro método para estudiar la concordancia entre ítems cualitativos lo proporciona los modelos log-lineales. Tanner (1985) muestra como se pueden analizar la concordancia entre variables cualitativas de manera análoga al análisis de la asociación en tablas de contingencias. Becker et al. (1991) y Graheam (1995) aplican estos modelos para describir la concordancia en determinadas situaciones.

Aickin (1990) propone una medida de concordancia, alfa, parecida al coeficiente kappa de Cohen pero la concordancia esperada estimada mediante máxima verosimilitud.

2. MODELOS E ÍNDICES PARA MEDIR LA CONCORDANCIA

Nos planteamos que cada población posee un subconjunto de características propias. Éstos induce una partición dentro de $\{A_1, \dots, A_m\}$ cuando $k \geq m$ como muestra la figura 1.

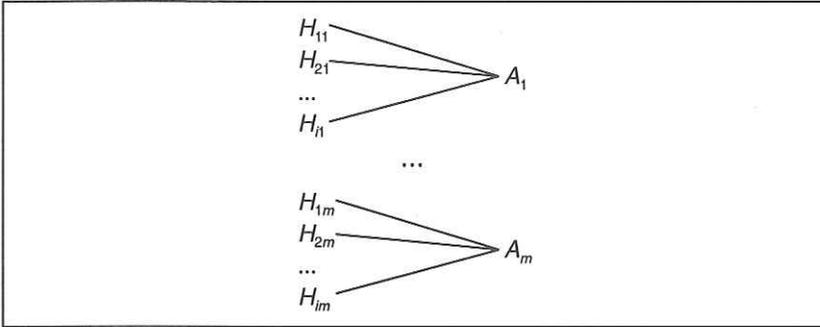


FIGURA 1

Este problema también tiene una formulación dual si el interés es encontrar que grupos presentan ciertas características (figura 2).

Ambos problemas se pueden plantear en un marco más general. Suponemos que existen $r \subseteq \{m, k\}$ grupos de poblaciones significativas que se asocian con r subconjuntos de características, tal como muestra la figura 3. En nuestro problema suponemos que r es conocido. En este sentido puede ser orientativo el número de factores significativos que obtiene el ACF.

Esta asociación induce un reagrupamiento en la tabla de contingencia que clasifica a nuestros datos. En la tabla 2 se observa un ejemplo.

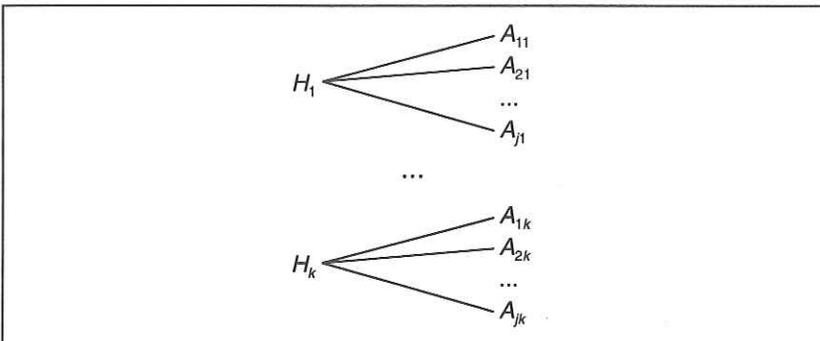


FIGURA 2

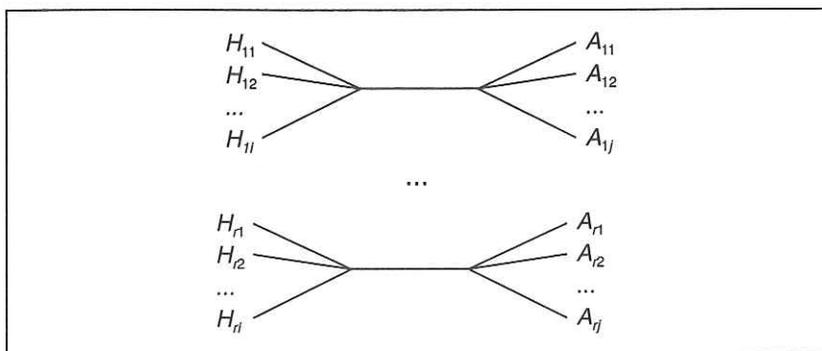


FIGURA 3

TABLA 2
Ejemplo de reagrupamiento.

Población	Caracteres					Total
	A_1	A_2	A_3	A_4	A_5	
H_1	p_{11}	p_{12}	p_{13}	p_{14}	p_{15}	$p_{1.}$
H_2	p_{21}	p_{22}	p_{23}	p_{24}	p_{25}	$p_{2.}$
H_3	p_{31}	p_{32}	p_{33}	p_{34}	p_{35}	$p_{3.}$
H_4	p_{41}	p_{42}	p_{43}	p_{44}	p_{45}	$p_{4.}$
Total	$p_{.1}$	$p_{.2}$	$p_{.3}$	$p_{.4}$	$p_{.5}$	1
Población	Caracteres				Total	
	$1 = A_2 + A_5$	$2 = A_1$	$3 = A_4 + A_3$			
$1 = H_1 + H_2$	$p_{11} + p_{15} + p_{22} + p_{25}$	$p_{11} + p_{21}$	$p_{14} + p_{13} + p_{24} + p_{23}$		$p_{1.} + p_{2.}$	
$2 = H_3$	$p_{32} + p_{35}$	p_{31}	$p_{34} + p_{33}$		$p_{3.}$	
$3 = H_4$	$p_{42} + p_{45}$	p_{41}	$p_{44} + p_{43}$		$p_{4.}$	
Total	$p_{.2} + p_{.5}$	$p_{.1}$	$p_{.4} + p_{.3}$		1	

2.1. Coeficiente de concordancia kappa de Cohen

Suponemos que cada uno de los sujetos de cierta población reciben dos tasas de modo independiente y medidas en una misma escala con r categorías. La distribución de probabilidad conjunta de ambas tasas se observa en la tabla 3.

Denotamos la probabilidad de la celda (i, j) , por p_{ij} , para $i, j = 1, 2, \dots, r$; y empleamos el símbolo "+" para denotar el sumatorio respecto el índice omitido. Los valores p_{i+} , $i = 1, \dots, r$ denotan la r probabilidades marginales de las filas y p_{+j} , $j = 1, \dots, r$ denotan las r probabilidades marginales de las columnas y $p_{++} = 1$. Las probabilidades espe-

TABLA 13
Distribución de probabilidad conjunta para dos tasas.

Tasa A	Tasa B				
	1	2	...	r	Total
1	p_{11}	p_{12}	...	p_{1r}	$p_{1.}$
2	p_{21}	p_{22}	...	p_{2r}	$p_{2.}$
.			...		
r	p_{r1}	p_{r2}	...	p_{rr}	$p_{r.}$
Total	$p_{.1}$	$p_{.2}$...	$p_{.r}$	1

radas bajo el supuesto de independencia entre las tasas las escribiremos como $\pi_{ij} = p_i + p_j$ para $i, j = 1, \dots, r$. La suma de las probabilidades observadas y esperadas en la diagonal principal las denotamos por

$$p_o = \sum_{i=1}^r p_{ii} \quad \text{y} \quad p_e = \sum_{i=1}^r \pi_{ii} \quad \text{respectivamente.}$$

La fórmula del coeficiente kappa para una tabla de contingencia cuadrada $r \times r$ es

$$\kappa = \frac{p_o - p_e}{1 - p_e} \tag{1}$$

Para estimar dicho coeficiente se toma la versión muestral de su definición reemplazando las correspondientes probabilidades por las estimadas en la muestra, que las representaremos por \hat{p}_{ij} , y $\hat{\pi}_{ij}$ para $i, j = 1, \dots, r$.

El coeficiente de concordancia kappa es una medida de la concordancia en la diagonal principal de tablas de contingencia cuadradas.

Nos planteamos estimar la asignación que maximiza dicho coeficiente. El anterior problema se formula matemáticamente mediante un problema de optimización.

Variables. Para $i = 1, \dots, k; j = 1, \dots, m; s = 1, \dots, r$ definimos:

$$z_{is} = \begin{cases} 1 & \text{si la población } i \text{ la incluimos en el grupo } s \\ 0 & \text{en caso contrario} \end{cases}$$

$$y_{js} = \begin{cases} 1 & \text{si la población } j \text{ la incluimos en el grupo } s \\ 0 & \text{en caso contrario} \end{cases}$$

$$x_{js} = \begin{cases} 1 & \text{si la población } i \text{ posee la característica } j \\ 0 & \text{en caso contrario} \end{cases}$$

Restricciones. Las restricciones a que sometemos nuestras variables modeliza el tipo de asignación deseada. Consideramos que todas las características A_1, \dots, A_m corresponden alguna población H_1, \dots, H_k y que toda población posee alguna característica relevante. Esta suposición se modeliza mediante el siguiente conjunto de restricciones:

$$\sum_{s=1}^r y_{js} = 1 \quad j = 1, \dots, m$$

$$\sum_{s=1}^r z_{is} = 1 \quad i = 1, \dots, k$$

$$\sum_{j=1}^m y_{js} \geq 1 \quad s = 1, \dots, r$$

$$\sum_{i=1}^k z_{is} \geq 1 \quad s = 1, \dots, r$$

Las restricciones que relacionan las variables x_{ij} con las z_{is} e y_{js} son

$$x_{ij} = \sum_{s=1}^r z_{is} y_{js} \quad \begin{array}{l} i = 1, \dots, k \\ j = 1, \dots, m \end{array} \quad (2)$$

Función objetivo. Calculemos, en función de las variables de decisión, el valor de kappa en la tabla plegada $r \times r$. El porcentaje observado en la muestra de casos concordantes, una vez plegada la tabla, se expresa por:

$$\hat{p}_o = \sum_{i,j}^{k,m} \hat{p}_{ij} x_{ij} \quad (3)$$

El porcentaje esperado de casos concordantes en la muestra bajo el supuesto de independencia se calcula mediante la expresión:

$$\hat{p}_e = \sum_{s=1}^r \bar{p}_{s+} \cdot \bar{p}_{+s}$$

donde \bar{p}_{s+} y \bar{p}_{+s} , $s = 1, \dots, r$ son las probabilidades marginales de la tabla plegada en r categorías. Calculemos estos valores en función de las probabilidades marginales de la tabla original $k \times m$.

$$\bar{p}_{s+} = \sum_{i=1}^k \hat{p}_{i+} z_{is} \quad s = 1, \dots, r$$

$$\bar{p}_{+s} = \sum_{j=1}^m \hat{p}_{+j} y_{js} \quad s = 1, \dots, r$$

$$\begin{aligned} \hat{p}_e &= \sum_{s=1}^r \left(\sum_{i=1}^k \hat{p}_{i+} z_{is} \right) \left(\sum_{j=1}^m \hat{p}_{+j} y_{js} \right) = \sum_{s=1}^r \left(\sum_{i,j=1}^{k,m} \hat{p}_{i+} z_{is} \hat{p}_{+j} y_{js} \right) = \\ &= \sum_{i,j=1}^{k,m} \left(\sum_{s=1}^r \hat{p}_{i+} \hat{p}_{+j} z_{is} y_{js} \right) = \sum_{i,j=1}^{k,m} \hat{p}_{i+} \hat{p}_{+j} \left(\sum_{s=1}^r z_{is} y_{js} \right) \end{aligned} \quad (3)$$

Sustituyendo la relación (2) en (3), obtenemos:

$$\hat{p}_e = \sum_{i,j=1}^{k,m} \hat{p}_{i+} \hat{p}_{+j} x_{ij} \quad (4)$$

La versión muestral del coeficiente kappa en términos de nuestras variables de decisión x_{ij} se obtiene al sustituir (3) y (4) en (1)

$$\kappa = \frac{\sum_{i,j=1}^{k,m} (\hat{p}_{ij} - \hat{p}_{i+} \hat{p}_{+j}) x_{ij}}{1 - \sum_{i,j=1}^{k,m} \hat{p}_{i+} \hat{p}_{+j} x_{ij}}$$

Esto conduce al siguiente problema de maximización:

$$(P_1) \quad \text{Maximizar } \kappa = \frac{\sum_{i,j=1}^{k,m} (\hat{p}_{ij} - \hat{p}_{i+} \hat{p}_{+j}) x_{ij}}{1 - \sum_{i,j=1}^{k,m} \hat{p}_{i+} \hat{p}_{+j} x_{ij}}$$

$$\text{Sujeto a:} \quad x_{ij} = \sum_{s=1}^r z_{is} y_{js} \quad \begin{array}{l} i = 1, \dots, k \\ j = 1, \dots, m \end{array}$$

$$\sum_{s=1}^r y_{js} = 1 \quad j = 1, \dots, m$$

$$\sum_{s=1}^r z_{is} = 1 \quad i = 1, \dots, k$$

$$\sum_{j=1}^m y_{js} \geq 1 \quad s = 1, \dots, r$$

$$\sum_{i=1}^k z_{is} \geq 1 \quad s = 1, \dots, r$$

$$x_{ij}, z_{is}, y_{js} \in \{0, 1\}$$

2.2. Modelos lineales para el estudio de la concordancia

Se han aplicado los modelos log-lineales para el estudio de la concordancia para el caso de dos tasas. El modelo log-lineal general para este caso es

$$p_{ij} = \frac{\exp(\lambda_{i1} + \lambda_{j2} + \tau_{ij})}{\sum_a \sum_b \exp(\lambda_{a1} + \lambda_{b2} + \tau_{ab})}$$

donde los parámetros del modelo satisfacen

$$\sum_i \lambda_{i1} = \sum_j \lambda_{j2} = \sum_i \tau_{ij} = \sum_j \tau_{ij} = 0$$

Para describir la concordancia en escalas nominales, Tanner (1985) utiliza el modelo de casi-independencia con

$$\kappa_{ij} = \tau l \quad (i = j) \quad (5)$$

donde $l(\cdot)$ es la función indicadora que vale 1 cuando la tasas coinciden y 0 en otro caso. Para nuestro caso dicha función es desconocida y es la que queremos estimar. La relación existe entre la función indicadora y nuestras variables viene dada por $l(i = j) = x_{ij}$.

Se puede demostrar que la estimación máximo verosímil de los parámetros de nuestro modelo $\tau, \lambda_{i1}, \lambda_{j2}, l(\cdot)$ conduce al siguiente problema de optimización:

$$(P_2) \quad \text{Maximizar } \mu + \tau \hat{p}_0 + \sum_{i=1}^r \bar{p}_{i+} \lambda_{i1} + \sum_{j=1}^r \bar{p}_{+j} \lambda_{j2}$$

$$\text{Sujeto a:} \quad \sum_{i=1}^r \lambda_{i1} = 0$$

$$\sum_{j=1}^r \lambda_{j2} = 0$$

$$\sum_{i,j=1}^r \exp(\mu + \lambda_{i1} + \lambda_{j2} + x_{ij} \tau) = 1$$

$$\bar{p}_{s+} = \sum_{i=1}^k \hat{p}_{i+} z_{is} \quad s = 1, \dots, r$$

$$\bar{p}_{+s} = \sum_{j=1}^m \hat{p}_{+j} y_{js} \quad s = 1, \dots, r$$

$$\hat{p}_0 = \sum_{i,j=1}^r \hat{p}_{ij} x_{ij}$$

$$x_{ij} = \sum_{s=1}^r z_{is} y_{js} \quad \begin{array}{l} i = 1, \dots, k \\ j = 1, \dots, m \end{array}$$

$$\sum_{s=1}^r y_{js} = 1 \quad j = 1, \dots, m$$

$$\sum_{s=1}^r z_{is} = 1 \quad i = 1, \dots, k$$

$$\sum_{j=1}^m y_{js} \geq 1 \quad s = 1, \dots, r$$

$$\sum_{i=1}^n z_{is} \geq 1 \quad s = 1, \dots, r$$

$$x_{ij}, z_{is}, y_{js} \in \{0, 1\}$$

2.3. Coeficiente alfa de Aickin

Aickin (1990) propone una medida de concordancia, alfa, que está basado en el modelo descrito por Kraemer et. al. (1988). Éste conduce a una estimación similar al coeficiente kappa de Cohen, pero reemplazando la concordancia esperada por su estimación máximo verosímil p_e' .

$$\alpha = \frac{p_0 - p_e'}{1 - p_e'}$$

Aickin (1990) demostró que si consideramos el modelo de Tanner de casi-independencia, [(5)], para la estimación de p_e' la expresión anterior adopta la forma:

$$\alpha = \frac{\exp(\tau) - 1}{\exp(\tau)} p_0$$

La estimación de las equivalencias que maximiza el coeficiente alfa de Aickin conduce el siguiente problema de optimización:

$$(P_3) \quad \text{Maximizar } z = \frac{\exp(\tau) - 1}{\exp(\tau)} \hat{p}_0$$

Sujeto a:

$$\sum_{i=1}^r \lambda_{i1} = 0$$

$$\sum_{j=1}^r \lambda_{j2} = 0$$

$$\sum_{i,j=1}^r \exp(\mu + \lambda_{i1} + \lambda_{j2} + x_{ij} \tau) = 1$$

$$\bar{p}_{s+} = \sum_{i=1}^k \hat{p}_{i+} z_{is} \quad s = 1, \dots, r$$

$$\bar{p}_{+s} = \sum_{j=1}^m \hat{p}_{+j} y_{js} \quad s = 1, \dots, r$$

$$x_{ij} = \sum_{s=1}^r z_{is} y_{js} \quad \begin{array}{l} i = 1, \dots, k \\ j = 1, \dots, m \end{array}$$

$$\hat{p}_0 = \sum_{i,j=1}^r \hat{p}_{ij} x_{ij} = 1$$

$$\sum_{s=1}^r y_{js} = 1 \quad j = 1, \dots, m$$

$$\sum_{s=1}^r z_{is} = 1 \quad i = 1, \dots, k$$

$$\sum_{j=1}^m y_{js} \geq 1 \quad s = 1, \dots, r$$

$$\sum_{i=1}^k z_{is} \geq 1 \quad s = 1, \dots, r$$

$$x_{ij}, z_{is}, y_{js} \in \{0, 1\}$$

3. MÉTODOS COMPUTACIONALES

El conjunto de restricciones (2) no son lineales. Este conjunto de restricciones son equivalentes al siguiente conjunto de restricciones lineales:

$$\begin{array}{r}
 i = 1, \dots, k \\
 2y_{js} \leq z_{is} + x_{ij} \quad j = 1, \dots, m \\
 s = 1, \dots, r
 \end{array}$$

Si reemplazamos el conjunto de restricciones (2) por el conjunto de restricciones (7) el problema de optimización P1 se convierte en un problema de *programación lineal fraccional entera* [ver Bazaraa (1979)]. Existen varios algoritmos para resolverlo. Hemos adaptado el algoritmo descrito en García et al. (1995) basado en *el método simplex* [ver Bazaraa (1989)].

Planteamos el siguiente problema de *programación lineal* [ver Bazaraa (1989)]

$$(P_4) \quad \text{Maximizar } Z(x) = \sum_{i,j=1}^r [\hat{p}_{ij} - (1 - \lambda) \hat{p}_{i+} \hat{p}_{+j}] x_{ij} - \lambda$$

Sujeto al mismo conjunto de restricciones de (P₁) donde λ es una constante positiva prefijada de antemano.

Algoritmo. El siguiente algoritmo converge a una solución óptima del problema P₁.

- **Paso 1.** Resolver el problema P₄ mediante el algoritmo del método simplex para un valor $\lambda = \lambda_0$ tal que $0 < \lambda_0 < 1$. Sea z_0^* su valor óptimo y x_0^* una solución óptima.
- **Paso 2.** Sea λ_n el valor de λ para la iteración n-ésima. z_n^* el valor óptimo y x_n^* una solución óptima para el problema P₄ en dicha iteración.

Si $z_n^* = 0$ entonces x_n^* es una solución óptima del problema P₁.

En caso contrario definimos

$$\lambda_{n+1} = \lambda_n + \frac{z_n^*}{1 - \sum_{i,j=1}^r \hat{p}_{i+} \hat{p}_{+j} (x_n^*)_{ij}}$$

y volvemos a resolver el problema P₄ para dicho valor de λ .

Los problemas P₂ y P₃ son difíciles de resolver. Su resolución pasa por aplicar una generalización del método de los multiplicadores de Lagrange descrito en Clarke (1990).

4. CONCLUSIONES

El ACF es capaz de representar geométricamente las poblaciones y caracteres simultáneamente. Dicha representación permite identificar gráficamente ciertas asociaciones entre ambos conjuntos de categorías.

En este trabajo hemos formulado tres modelos que nos permiten estimar, maximizando las medidas de concordancia kappa y alpha, y mediante la estimación máximo verosímil en los modelos log-lineales de casi-independencia, las anteriores relaciones.

Estas medidas tienen una distribución muestral conocida y ésta puede ser aplicada para determinar la significación y cuantificar la asociación establecida.

5. REFERENCIAS

- AICKIN, M. (1990): «Maximun likelihood estimation of agreement in the constant predictive probability model, and its relation to Cohen's kappa». *Biometrics*, 46, 293-302.
- BAZARAA, M. and JARVIS, J. (1989): *Programación Lineal y Flujo en Redes*. Limusa, México.
- BAZARAA, M. and SHETTY (1979): *Nonlinear Programming Theory and Algorithms*, capítulo 11. John Wiley & Sons.
- BECKER, M. and AGRESTI, A. (1991): «Log-linear modelling of pairwise interobserver agreement on a categorical scale». *Statistics in Medicine*, 11, 101-114.
- CLARKE, F. H. (1990): *Optimization and Nonsmooth Analysis*. Capítulo 6. Wiley, New York.
- COHEN, J. (1960): «A coefficient of agreement for nominal scales». *Educ. Psychol. Meas.* 20, 37-46.
- CHMURA, H. (1980): «Extension of the Kappa Coefficient». *Biometrics* 36, 207-216.
- COHEN, J. (1968): «Weighted kappa: nominal scale agreement with provision for scaled disagreement or partial credit». *Psychol. Bull.* 70, 213-220.
- CUADRAS, M. C. (1991): *Métodos de Análisis Multivariante*, capítulo 14, PPU, Barcelona.
- FLEISS, J. L. (1981): *Statistical Methods for Rates and Proportions*, capítulo 13, 2nd edn; Wiley, New York.
- FLEISS, J. L. (1971): «Measuring nominal scale agreement among many raters». *Psychol. Bull.* 76, 378-383.
- FLEISS, J. L. and COHEN, J. (1973): «The Equivalence of Weighted Kappa and the Interclass Correlation Coefficient as Measures of Reliability». *Educ. Psychol. Meas.* 33, 613-619.
- FLEISS, J. L. and CUZICK, J. (1979): «The reliability of dichotomous judgments: Unequal numbers of judges per subject». *Appl. Psychol. Meas.* 3, 537-542.
- GARCÍA, R., LÓPEZ, M. L. y VERÁSTEGUI, D. (1995): en prensa.
- GRAHEAM, P. (1995): «Modelling covariate effects in observer agreement studies: the case of nominal scale agreement». *Statistics in Medicine*, 14, 299-310.
- LANDIS, J. R. and KOCH, G. G. (1977): «A one-way components of variance model for categorical data». *Biometrics* 33, 671-679.
- KRAEMER, H. C. and BLOCH, D. A. (1988): «Kappa coefficients in epidemiology: an appraisal of a reappraisal», *Journal of Clinical Epidemiology*, 41, 959-968.
- TANNER, M.A. and YOUNG, M.A. (1985): «Modeling agreement among raters», *Journal of the American Statistical Association*, 80, 175-180.