

03061  
2es.



Universidad Nacional Autónoma de México

Unidad Académica de los Ciclos  
Profesionales y de Posgrado  
del C.C.H.

Instituto de Investigaciones  
en Matemáticas Aplicadas y en  
Sistemas

INFERENCIA ESTADISTICA EN EL ANALISIS DE  
CORRESPONDENCIAS

TESIS CON  
FALLA DE ORIGEN

T E S I S  
Que para obtener el grado de:  
MAESTRO EN ESTADISTICA E  
INVESTIGACION DE OPERACIONES  
Presenta el Actuario  
Eduardo Castaño Tostado

México, D. F.

Febrero, 1987



Universidad Nacional  
Autónoma de México



## **UNAM – Dirección General de Bibliotecas Tesis Digitales Restricciones de uso**

### **DERECHOS RESERVADOS © PROHIBIDA SU REPRODUCCIÓN TOTAL O PARCIAL**

Todo el material contenido en esta tesis está protegido por la Ley Federal del Derecho de Autor (LFDA) de los Estados Unidos Mexicanos (México).

El uso de imágenes, fragmentos de videos, y demás material que sea objeto de protección de los derechos de autor, será exclusivamente para fines educativos e informativos y deberá citar la fuente donde la obtuvo mencionando el autor o autores. Cualquier uso distinto como el lucro, reproducción, edición o modificación, será perseguido y sancionado por el respectivo titular de los Derechos de Autor.

## Índice-

1. Introducción.....	1
2. Análisis de Correspondencias, una introducción.....	4
3. Inferencia enfoque Asintótico.....	17
4. Inferencia enfoque de Remuestreo.....	29
5. Aplicaciones.....	50
6. Conclusiones y Discusión.....	59
Bibliografía.....	63
Apéndices.....	66

## 1. Introducción.

El Análisis de Correspondencias (A.C.) es una técnica encuadrada en el enfoque del Análisis de Datos siendo su propósito fundamental la descripción gráfica, a través de elementos algebraicos y geométricos, de las relaciones más relevantes en una tabla de contingencias. Lo anterior es la idea que tradicionalmente se ha tenido en el A.C., con lo que el carácter exploratorio del análisis tiene preponderancia sobre los aspectos inferenciales. En este trabajo se sigue la interpretación del Análisis de Correspondencias debida a la escuela de Estadística en Francia; sin embargo, el objetivo es presentar elementos inferenciales para enriquecer y validar las conclusiones del Análisis de Correspondencias.

Trabajos sobre inferencia en el Análisis de Correspondencias ya existen en la literatura, aunque sufren de algunos inconvenientes. Un primer grupo de artículos que se enfoca al tratamiento de los valores singulares, parte de la suposición de independencia entre los criterios de clasificación, supuesto que en la práctica es poco frecuente, ya que precisamente el A.C. tiene por objeto mostrar la estructura de dependencia de lo

tabla; trabajos relevantes en este sentido son los de Lebart (1976) y Corsten (1976). Un segundo grupo de resultados, debidos a O'Neill (1978, 1981), si bien llega a distribuciones asintóticas de los valores singulares, estas dependen fuertemente del cumplimiento de condiciones que involucran la estructura teórica de la tabla; estas condiciones son muy complejas y resulta completamente impráctico verificarlas, por lo que se utilizan las distribuciones sin tomarlas en cuenta. Una tercera alternativa es desarrollada por Goodman (1985), mediante un enfoque general de ajuste de modelos a una tabla de contingencias con dos criterios de clasificación, en términos de validar un sistema de pesos o calificaciones sobre los renglones y/o las columnas de la tabla bajo estudio. El Análisis de Correspondencias es un caso especial del tipo de modelos manejados por este autor; sin embargo las herramientas de análisis aún no han sido totalmente desarrolladas, siendo además un camino menos directo y más complejo que las propuestas anteriores.

Los tres grupos de trabajos mencionados anteriormente versan sobre inferencia en valores singulares del A.C.; Greenacre (1984), utilizando técnicas de remuestreo analiza la

estabilidad de la representación gráfica del A.C., pero que desde la perspectiva de este trabajo, pueden mejorarse substancialmente. Es propósito de esta tesis presentar propuestas de análisis de tipo inferencial sobre aspectos relevantes en la interpretación de resultados del A.C.; para ello se combinan ideas ya desarrolladas en los trabajos mencionados anteriormente con las propiedades de mayor alcance de los métodos de muestreo.

La estructura de la presentación es como sigue; en el capítulo dos se presentan los aspectos algebraicos y geométricos del Análisis de Correspondencias. En el capítulo tres se muestran los resultados referentes a los trabajos sobre inferencia a través de teoría asintótica, es decir los desarrollados por Lebart (1976), Corsten (1976), O'Neill (1978, 1981) y Goodson (1985). En el cuarto capítulo se hacen nuevas propuestas inferenciales, no sólo sobre los valores singulares sino en variables involucradas en la representación gráfica. En el capítulo cinco se presenta una aplicación del A.C. al análisis sensorial en tecnología de alimentos. Por último se dan conclusiones y se discute lo desarrollado en todo el trabajo.

## 2. Análisis de Correspondencias, una introducción.

El Análisis de Correspondencias es resultado de desarrollos geométricos y algebraicos. Existen distintos enfoques del A.C. en su utilización; uno de estos enfoques debido a la escuela francesa, y que se sigue en este trabajo, se debe primordialmente a los trabajos de Jean Paul Benzécri. Esta técnica se sitúa en el análisis multivariado de datos y su objetivo principal es la representación gráfica de las relaciones más relevantes que se hayan en forma subyacente en una tabla de contingencia de dos o más criterios de clasificación.

Es propósito de esta sección el presentar los detalles más importantes tanto algebraicos como geométricos del A.C. y para ello se sigue la notación de Greenacre (1984). La presentación por sencillez se hará para el caso de dos criterios de clasificación en la tabla de contingencias, aunque la generalización al caso multicriterio es inmediata, ya que a fin de cuentas las conclusiones susceptibles de obtenerse mediante el A.C. son sólo sobre las interacciones de orden uno entre los criterios de clasificación de la tabla bajo estudio (Greenacre, 1984).

## 2.1 Notación.

Sea  $N$  una tabla de contingencias con  $pq$  celdas y con frecuencias  $n_{ij}$ ,  $i=1, \dots, p+1$ ;  $j=1, \dots, q+1$  ( $p > q$ ), y denote por  $n$  a

$$n = \sum_{i=1}^p \sum_{j=1}^q n_{ij} \quad (2.1)$$

Llane  $\hat{P}$  a la matriz de frecuencias relativas

$$\hat{P} = (1/n)N \quad (2.2)$$

y a

$$C = \hat{P} \mathbf{1} \quad \text{y} \quad \zeta = \hat{P}' \mathbf{1}$$

donde  $\mathbf{1}$  es un vector de unos con las dimensiones adecuadas al caso,  $C$  representa el vector de sumas por renglón de  $\hat{P}$  y  $\zeta$  el vector de sumas por columna de la misma matriz. Denote a los elementos de estos dos vectores como  $p_{i.}$ ,  $i=1, \dots, p+1$ ;  $p_{.j}$ ,  $j=1, \dots, q+1$  respectivamente.

Se construyen las siguientes matrices:

$$R = D_r' \hat{P}, \quad C = D_c' \hat{P} \quad (2.3)$$

donde

$$D_r = \text{diag}(C), \quad D_c = \text{diag}(\zeta) \quad (2.4)$$

Los renglones de la matriz  $R$ , denotados por  $(C_{i.}, i=1, \dots, p+1)$  son llamados los perfiles renglón de la matriz  $\hat{P}$ ; análogamente, los renglones de la matriz  $C$  ( $\zeta_j, j=1, \dots, q$ ) son nombrados perfiles columna. Así, el  $i$ -ésimo perfil renglón es un vector cuyas entradas representan la distribución discreta marginal en el  $i$ -ésimo

renglón ( $i=1, \dots, p+1$ ) de la matriz  $P$  análogamente en el caso de los perfiles columna. En vista de lo anterior la estructura relativa de la tabla representada en dos formas por estos conjuntos de perfiles es invariante ante el tamaño total de la tabla  $n$ .

Ya que los perfiles renglón como los perfiles columna pueden considerarse como vectores de  $q+1$  y  $p+1$  entradas respectivamente, el objetivo del A.C. es buscar un subespacio de rango menor en los espacios respectivos, en los que las proyecciones de los correspondientes perfiles representen las asociaciones más relevantes que se dan en los espacio  $q+1$ -variado y  $p+1$ -variado respectivamente (figura 2.1). Posteriormente a la obtención de estos dos subespacios el A.C. obtiene una representación gráfica conjunta de ambos subespacios.

Si se denota por  $d_n(\underline{x}, \underline{y})$  a la distancia de  $\underline{x}$  a  $\underline{y}$  bajo la métrica  $n$ , un criterio para encontrar el subespacio ( $S^*$ ) en el espacio de perfiles renglón es

$$\min_{S^*} \sum_{i=1}^I w_i d_n(\underline{c}_i, \underline{s}_i) \quad \underline{s}_i \in S^* \quad (2.5)$$

donde  $w_i$  es un peso diferencial por perfil y  $\underline{s}_i$  la proyección ortogonal de  $\underline{c}_i$  en  $S^*$ .

En el A.C. se escoge  $w_i = p_i$ , y a la métrica  $n$  se le asigna la

matriz  $D_c^j$  es decir que los perfiles se ponderan en función de su frecuencia relativa total y las  $q+i$  coordenadas de cada uno de estos por su frecuencia relativa acumulada por estas en todos los perfiles. Así, (2.5) queda expresado como

$$\min \|R-S\|_{r,c} = \min \sum_i p_i (r_i - s_i)' D_c^j (r_i - s_i) \quad (2.6)$$

Análogamente para los perfiles columna, es decir, los renglones de la matriz  $C$ , se tiene que el criterio utilizado por el A.C. para encontrar el subespacio de mejor ajuste  $T$  es

$$\min \|C-T\|_{r,c} = \min \sum_j p_j (c_j - t_j)' D_r^j (c_j - t_j) \quad (2.7)$$

con  $t \in T$ .

Los centroides de los perfiles renglón y de los perfiles columna son respectivamente  $g = R1_c'$  y  $z = C1_r'$ . Se puede demostrar que estos están contenidos en los planos  $S$  y  $T$  por lo que se parte sin pérdida de generalidad de los perfiles centrados (Greenacre, 1984)

$$R - 1_c g', C - 1_r z'. \quad (2.8)$$

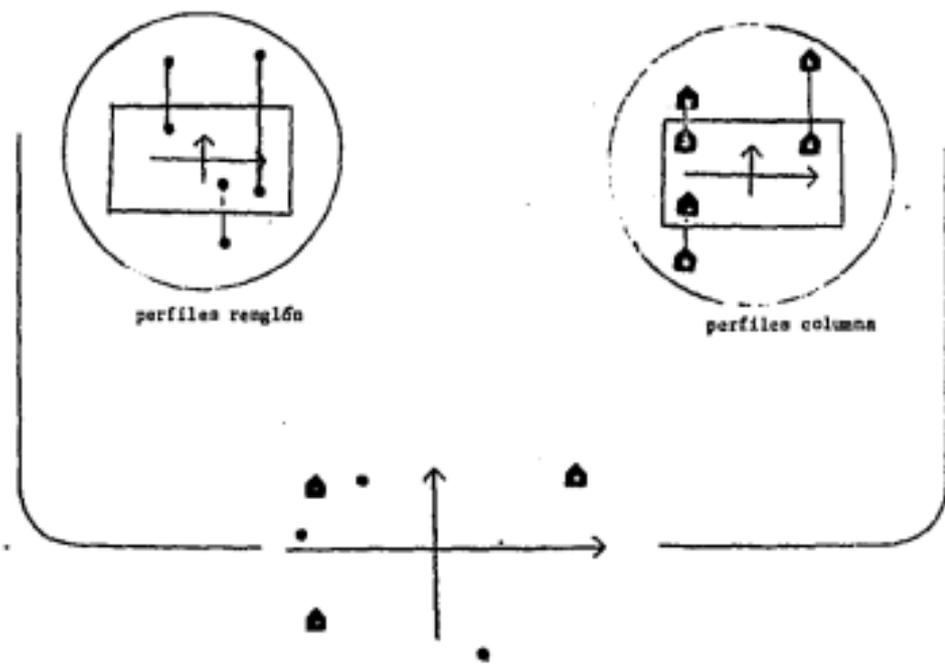
## 2.2 Descomposición en Valor Singular y el S.G.

Suponga que  $R - 1_c g'$  es de rango  $K$  ( $K \leq \min(p, q)$ ); se pueden encontrar  $K$  vectores ortogonales tales que (Green y Carroll, 1976):

$$R - 1_c g' = FB' \quad (2.9)$$

donde  $B = (b_1, b_2, \dots, b_k)$  son los ejes principales de los perfiles

Figura 2.1 Objetivo del Análisis de Correspondencias.



renglón y F las coordenadas de los perfiles renglón centrados, respecto a B. Análogamente

$$C - \underline{1}c' = GA' \quad (2.10)$$

con  $A = (a_1, a_2, \dots, a_K)$  ejes principales de los perfiles columna y G las coordenadas respectivas. Ahora,

$$\begin{aligned} D_r (R - \underline{1}c') &= \hat{P} - c c' \\ D_c (C - \underline{1}c') &= \hat{P}' - c c' \end{aligned} \quad (2.11)$$

con lo que encontrar A y B son problemas interrelacionados.

La descomposición en valor singular (véase apéndice A), da la pauta para encontrar estos dos conjuntos de vectores. Esta herramienta algebraica provee las matrices L, M y  $D_{\hat{u}}$  ( $K \times K$ ) tales que

$$\hat{P} - c c' = L D_{\hat{u}} M' \quad (2.12)$$

sujeito a que

$$L' D_r' L = M' D_c' M = I (K \times K) \quad (2.13)$$

Las columnas de la matriz L son K vectores ortogonales bajo la métrica  $D_r'$ , y constituyen una base ortonormal para los renglones de  $\hat{P} - c c'$ , mientras que las columnas de M son K vectores ortogonales bajo la métrica  $D_c'$ , siendo una base ortonormal para las columnas de la matriz  $\hat{P}' - c c'$ . Por último,  $D_{\hat{u}}$  es una matriz diagonal con elementos  $(\hat{u}_1, \hat{u}_2, \dots, \hat{u}_K > 0)$  llamados valores singulares. En vista de (2.12)-(2.13) se tiene que

$$\hat{P} - \underline{1}\underline{1}' = \sum_{k=1}^K \hat{u}_k \hat{u}_k' \quad (2.14)$$

Dado el carácter aditivo en (2.14) se puede aproximar la matriz (2.11) por  $Kt \leq K$  de los vectores de  $M$  y  $L$ . Si  $Kt=2$  se tiene una matriz de rango 2 que aproxima a la matriz de rango  $K$ . Esta forma de aproximar a (2.12) es óptima bajo el criterio

$$\min \|A - X\| = \min \sum_{i=1}^n p_{i1} (\beta_i - \hat{\beta}_i)' D_c^{-1} (\beta_i - \hat{\beta}_i) \quad (2.15)$$

con  $A = \hat{P} - \underline{1}\underline{1}'$  (Green y Carroll, 1976).

De (2.12) y (2.11) se tiene que las columnas de  $M$  representan una base para los renglones de la matriz

$$R = \underline{1}\underline{1}'$$

si se denota por  $\hat{F}$  las coordenadas respectivas, es decir,

$$\hat{F}M = (R - \underline{1}\underline{1}') \quad (2.16)$$

por (2.13) se tiene que

$$\hat{P} = (R - \underline{1}\underline{1}') D_c^{-1} M \quad (2.17)$$

Análogamente, si  $G$  denota las coordenadas de los renglones de

$$(C - \underline{1}\underline{1}')$$

respecto a  $L$ ,

$$\hat{G} = (C - \underline{1}\underline{1}') D_r^{-1} L \quad (2.18)$$

Si  $Kt=2$  se tendrá

$$\hat{F}2 = (R - \underline{1}\underline{1}') D_c^{-1} M2 \quad (2.19)$$

$$\hat{G}2 = (C - \underline{1}\underline{1}') D_r^{-1} L2 \quad (2.20)$$

serán las coordenadas de los perfiles renglón y de los perfiles columna en el subespacio de rango dos que mejor aproxima las relaciones en los espacios q y p variadas respectivamente, bajo el criterio (2.15), donde M2 y L2 representan los dos primeros ejes principales respectivos. Con estos resultados se tendrán las representaciones gráficas de los perfiles renglón y de perfiles columna. Debe renarcarse que estas representaciones son en dos espacios diferentes; el objetivo final del A.C. es la representación simultánea de ambos conjuntos, que formalmente no es posible realizar, pero que, a pesar de lo anterior, en la siguiente sección se muestran las expresiones que permitirán establecer tal representación conjunta ad hoc.

### 2.3 Fórmulas de Transición.

Las expresiones (2.17)-(2.18) están relacionadas mediante las llamadas fórmulas de transición. Para mostrar lo anterior, primero de (2.17) y (2.11) se tiene que

$$\hat{F} = D_r^{-1} (\hat{P} - \underline{c} \underline{c}') D_c^{-1} M. \quad (2.21)$$

Ahora de (2.12),

$$(\hat{P} - \underline{c} \underline{c}') D_c^{-1} M = L D \hat{U}$$

con lo que

$$\hat{F} = \hat{D}^{-1} L D \hat{U} \quad (2.22)$$

En forma análoga,

$$\hat{G} = \hat{D}^{-1} M D \hat{U} \quad (2.23)$$

De (2.21), (2.23) y (2.12) se tiene que

$$\begin{aligned} \hat{F} &= R \hat{G} \hat{D}^{-1} \\ \hat{G} &= C \hat{F} \hat{D}^{-1} \end{aligned} \quad (2.24)$$

llamadas las fórmulas de transición; estas expresiones significan que las coordenadas de un conjunto de perfiles se puede expresar en función de las coordenadas del otro conjunto de perfiles. Para explicar las implicaciones de estas expresiones, considere la coordenada del  $i$ -ésimo perfil renglón respecto al  $k$ -ésimo eje principal respectivo denotada por  $\hat{f}_{ik}$ ; de acuerdo con (2.24),

$$\hat{f}_{ik} = \sum_{j=1}^n r_{ij} \hat{g}_{jk} / \hat{u}_k \quad (2.25)$$

donde  $r_{ij}$  es la entrada  $j$  del perfil  $i$ ,  $\hat{g}_{jk}$  la coordenada del perfil columna  $j$  en su eje principal  $k$  y  $\hat{u}_k$  el valor singular respectivo. En vista de (2.25), se tiene que  $\hat{f}_{ik}$  es una combinación de las coordenadas de los perfiles columna con respecto a su eje  $k$ . Geométricamente, un perfil renglón tenderá a una posición que corresponderá a las coordenadas de los perfiles columna más importantes respecto a aquel. Las fórmulas de transición permiten así la superposición de gráficas de los dos tipos de perfiles

para su interpretación conjunta.

## 2.4 Interpretación de los Valores Singulares

La variación estimada de los perfiles renglón en el eje  $k$  es

$$\sum_1 p_{i.} (\hat{f}_{ik} - \bar{f}_{.k})^2 \quad k=1, \dots, K \quad (2.26)$$

$$\bar{f}_{.k} = \sum_1 p_{i.} \hat{f}_{ik} \quad k=1, \dots, K \quad (2.27)$$

Puede mostrarse que

$$(\bar{f}_{.1}, \bar{f}_{.2}, \dots, \bar{f}_{.k}) = 0$$

con lo que (2.26) es

$$\sum_1 p_{i.} \hat{f}_{ik}^2 \quad (2.28)$$

La expresión (2.28) puede verse como

$$(\hat{f}_{1k}, \hat{f}_{2k}, \dots, \hat{f}_{pk}) \begin{bmatrix} p_1 & 0 & 0 & \dots & 0 \\ 0 & p_2 & 0 & \dots & 0 \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & 0 & \dots & p_k \end{bmatrix} \begin{bmatrix} \hat{f}_{1k} \\ \hat{f}_{2k} \\ \vdots \\ \hat{f}_{pk} \end{bmatrix} \quad (2.29)$$

el producto del primer renglón de  $\hat{F}$  por sí mismo ponderado por la matriz  $D$ , es decir, como parte del producto de matrices  $\hat{F}' D \hat{F}$ .

De (2.22), la expresión (2.28) desarrollandola se llega a

$$\hat{F}' D \hat{F} = D_{\hat{U}} L' D_{\hat{R}} D_{\hat{R}}' L D_{\hat{U}} = D_{\hat{U}}^2$$

y

$$\text{traza}(\hat{F}' D \hat{F}) = \text{traza}(D_{\hat{U}}^2) \quad (2.30)$$

El susando correspondiente a  $\hat{U}_k$  es precisamente (2.28), con lo que  $\hat{U}_k^2$  es el estimador de la varianza de los perfiles renglón respecto al eje principal  $g_k, k=1, \dots, K$ . Si se recuerda los valores singulares están ordenados, por lo que en vista de la expresión (2.30),  $g_1$  será el eje de mayor variación en las coordenadas de los perfiles renglón y así hasta el eje  $g_K$ .

La expresión (2.6), pero ahora considerando la variación de los perfiles renglón respecto a su centroide es

$$\sum_i p_i \cdot (c_i - \bar{c}) \cdot D_c^T (c_i - \bar{c}) \quad (2.31)$$

que es igual a

$$\text{traza}[D_r (R - I\bar{c})^T D_c^T (R - I\bar{c})^T] \quad (2.32)$$

Ahora, por (2.17), (2.30) y (2.11) se tiene que

$$\begin{aligned} \text{traza}[D_0^2] &= \text{traza}[\hat{P}^T D_r \hat{P}] = \text{traza}[D_r \hat{P} \hat{P}^T] \\ &= \text{traza}[D_r D_r^T (\hat{P} - I\bar{c})^T D_c^T M M^T D_c (\hat{P} - I\bar{c})^T D_r^T] \quad (2.33) \\ &= \text{traza}[D_r (R - I\bar{c})^T D_c^T M M^T D_c^T (R - I\bar{c})^T] \end{aligned}$$

pero

$$D_c^T M M^T D_c^T = D_c^T$$

ya que

$$M^T D_c^T M = I$$

por lo que

$$\text{traza}[D_0^2] = \text{traza}[D_r (R - I\bar{c})^T D_c^T (R - I\bar{c})^T] \quad (2.34)$$

En vista de (2.34) la variación total original es recuperada

integrante en la descomposición en valor singular, en la variación de las coordenadas F.

Análogamente respecto a las coordenadas de los perfiles columna

$$\sum u_k^2 = \text{traza}(\hat{\theta}^T D_c \hat{\theta}) = \sum_j p_j (\epsilon_j - c)^2 / D_j^2 (\epsilon_j - c) \quad (2.35)$$

Por último, el A.C. provee de tres medidas para cuantificar la calidad de la representación gráfica; estas tres medidas son:

i) la calidad global de la representación en 1 dimensiones, representada por

$$(\hat{u}_1^2 + \dots + \hat{u}_l^2) / \sum_k \hat{u}_k^2 \quad 1 \leq l \leq K.$$

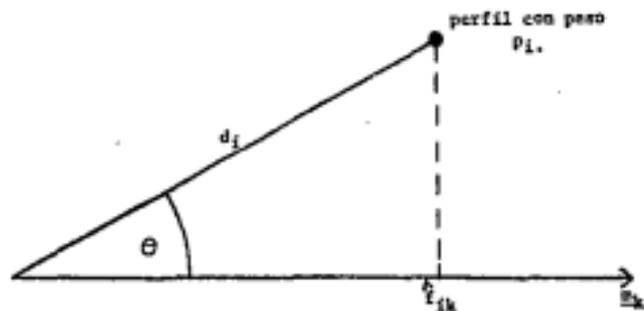
Usualmente se buscará que con  $l=1$  ó  $l=2$  se tenga la representación gráfica. Se puede interpretar como la proporción de la variación global que representan 1 dimensiones.

ii) La contribución relativa, que permite saber el ángulo de un perfil con respecto a uno de los ejes principales; de la figura 2.2, en términos de la definición del coseno de un ángulo se tiene que

$$(p_{i, ik}^{\wedge 2}) / (p_{i, d}^{\wedge 2}) = (f_{ik} / d_i)^2 = \cos^2 \theta$$

es decir, si  $\cos^2 \theta$  es cercano a 1 el ángulo del perfil  $i$  será pequeño con respecto al eje en cuestión.

Figura 2.2 Contribuciones relativa (correlación) y absoluta.



$$\text{Correlación} = (\hat{y}_{ik} / d_i)^2$$

$$\text{Contribución Absoluta} = (p_i \cdot \hat{y}_{ik})^2 / d_k^2$$

iii) La contribución absoluta, que cuantifica la aportación de un perfil en la variación total de un eje principal

$$(p_{i, k}^2) / \sum_k^2$$

En la práctica la utilidad de estas medidas es muy importante; respecto a la calidad global, lo que generalmente se desea es que con  $l=1$  ó  $l=2$  se tenga aglutinada la mayor parte de la variación original en la nube de perfiles multivariada ya que la aproximación de las relaciones a través de la gráfica unidimensional o bidimensional será mejor.

En el caso en que la calidad global no sea del todo convincente, es decir que oscile entre el 10% y 20%, las otras dos medidas mencionadas arriba entran en juego; la contribución absoluta identifica a aquellos perfiles que pesaron más en la orientación de cada uno de los ejes principales, mientras que la contribución relativa medirá la correlación de cada perfil con respecto a cada uno de los ejes principales, dando así una idea de que tan bien o que tan mal estuvieron representados. Los puntos que aparecen cercanos al origen de la gráfica estarán mal representados, con excepción de aquellos que coincidan con su centroide respectivo.

Así, mientras la calidad global es suficiente para cuando agrupa la mayor parte de la varianza original, las otras dos medidas calificarán en forma individual a cada perfil en términos de su aporte en la redistribución de la variabilidad en cada uno de los ejes generados por la DVB así como su correlación con éstos.

### 3. Inferencia en enfoque étnico.

Es propósito de este capítulo el presentar los resultados fundamentales hasta ahora encontrados en la literatura respecto a la inferencia estadística en el A.C. Estos resultados son presentados cronológicamente comenzando con los de Lebart (1976) y Corsten (1976) que operan bajo la suposición de independencia entre los criterios de clasificación de la tabla; posteriormente son mostrados para el caso de dependencia las distribuciones sobre los valores singulares, debidas a O'Neill (1978, 1981). Como tercer punto se establecen los elementos principales del enfoque de ajustes de modelos a tablas de contingencias debido a Goodman (1985). Finalmente se hacen comentarios generales sobre los desarrollos presentados y a partir de ellos se mencionan las necesidades a cubrir, que se atacarán con las herramientas propuestas en el capítulo cuatro.

#### 3.1 Inferencia en caso de Independencia.

Bajo independencia entre los dos criterios de clasificación, Lebart (1976) mostró que la distribución de los valores singulares estimados mediante la descomposición en valor

singular en el A.C. puede ser aproximada por la distribución de los raíces correspondientes a una matriz Wishart central, con matriz de covarianzas la identidad de orden  $p-1$  y grados de libertad  $q-1$ . Trabajando con el vector cuyos elementos son

$$\sqrt{n}(\hat{n}_{1j} - n_{1j} / n_{1.}, n_{.j} / n_{.j}) \quad (3.1)$$

y utilizando la aproximación a la distribución Normal, además de aplicar transformaciones ortogonales a (3.1), llega a la aproximación deseada, es decir que si

$$\hat{u}_1 < \dots < \hat{u}_p$$

son los valores singulares estimados, la distribución asintótica conjunta de estos es

$$f(\hat{u}_1, \dots, \hat{u}_p) = c \prod_{1 < j < k} (n(\hat{u}_j^2 - \hat{u}_k^2)) \prod_{i=1}^p (n\hat{u}_i^2)^{1/2(q-p+1)} \exp\left(-\frac{1}{2} \sum_{i=1}^p n\hat{u}_i^2\right)$$

con

$$(3.2)$$

$$c = \frac{\pi^{p/2} (-1/2)^{pq}}{2} \prod_{i=1}^p \Gamma\left(\frac{1}{2}(p+1-i)\right) \Gamma\left(\frac{1}{2}(q+1-i)\right)^{-1}$$

Para detalles en la obtención de (3.2) véase Anderson (1958).

El resultado de Lebart también fue obtenido por Corsten (1976) según Greenacre (1984).

Por otra parte, Lebart (1976) mostró que el porcentaje de variación explicada hasta el  $l$ -ésimo eje principal, definido por

$$\sum_{i=1}^l \hat{u}_i^2 / \sum_{i=1}^p \hat{u}_i^2 \quad (3.3)$$

es independiente de

$$\sum_{i=1}^p \hat{Q}_i^2. \quad (3.4)$$

Este resultado es importante para el A.C. ya que aún en el caso de independencia, la representación gráfica producida podrá manejarse. En este sentido Lebart para ciertas dimensiones de una tabla de contingencias, generó por simulación tablas de significancia tanto para los porcentajes como para los 5 valores más grandes.

De esta manera, bajo independencia el problema inferencial sobre los valores singulares está resuelto. Si bien esto es un avance importante, en la práctica generalmente el caso de interés es cuando se presenta dependencia entre los criterios de clasificación, por lo que las conclusiones estadísticas de las herramientas propuestas son de un uso restringido constituyendo sólo puntos de referencia.

### 3.2 Inferencia en el caso de dependencia.

O'Neill (1978) obtuvo las distribuciones asintóticas de los valores singulares estimados en el caso de dependencia; sus resultados se presentan a continuación.

Denote por  $\mu$  el vector de  $p$  entradas cuyos elementos son

$\hat{u}_1, \hat{u}_2, \dots, \hat{u}_p$ . Particione a  $\hat{\mu}$  en dos subvectores  $\hat{\mu}_1$  de orden  $k=p-s$  y  $\hat{\mu}_2$  de orden  $s$ , donde  $\hat{\mu}_1$  corresponde a los valores estimados cuyos valores teóricos son distintos cero.

Lancaster (1958) mostró que bajo un esquema de correlación las probabilidades teóricas de la tabla  $P_{ij}, i=1, \dots, p, j=1, \dots, q$  son tales que

$$P_{ij} = P_{i.} \cdot P_{.j} \cdot \left( 1 + \sum_{k=1}^K u_k x_{ki} y_{kj} \right)$$

donde

$$\sum_i x_{wi} x_{vi} P_{i.} = \delta_{wv} \quad (3.5)$$

$$\sum_j y_{wj} y_{vj} P_{.j} = \delta_{wv}$$

$$\sum_i \sum_j x_{wi} y_{vj} P_{ij} = \delta_{wv}$$

A la vista de esta parametrización llamada la fórmula de reconstitución, los resultados de O'Neill (1978, 1981) son los siguientes:

i)  $n \hat{\mu}_2' I \hat{\mu}_2^k$  se distribuye asintóticamente como lo hacen las raíces de una matriz Wishart  $W_{q-p+s} (1, s)$  si y solo si para cada  $i, j, k, l > p-s$

$$\sum_{w=1}^u \left( \sum_{i=1}^k x_{wi} x_{ki} P_{i.} \right) \left( \sum_{j=1}^l y_{wj} y_{lj} P_{.j} \right) = 0. \quad (3.6)$$

ii)  $n(\hat{\mu}_1 - \mu_1) \sim AN(0, \Delta)$  (3.7)

donde  $\Delta$  depende en su cálculo explícito de la estructura teórica de la tabla (expresiones 10 y 11 de O'Neill, 1981).

Así en el caso de dependencia con valores singulares teóricos de multiplicidad 1, estas distribuciones serán de utilidad si se conoce la parametrización específica que sustentate la tabla bajo estudio.

### 3.3 Modelos de Asociación y Modelos de Correlación.

Goodman(1985) muestra otras posibilidades de análisis en una tabla de contingencias de dos criterios de clasificación. Su propuesta se orienta a la validación de un sistema de calificaciones (scoring system), para los renglones de la tabla y/o para las columnas de la misma. Los modelos propuestos son distintas parametrizaciones de las probabilidades teóricas por celda; estas parametrizaciones son divididas por el autor en dos clases:

- . modelos de correlación
- . modelos de asociación

A continuación se presentan los modelos de correlación ya que sólo en este caso se encontraron elementos importantes a reportarse desde la perspectiva del A.C.

### 3.3.1 Modelos de Correlación.

Este tipo de modelos parte de la parametrización llamada RC de correlación

$$P_{ij} = P_{i.} P_{.j} (1 + u x_i y_j) \quad (3.8)$$

donde  $x_i, y_j$  son calificaciones de renglón y de columna respectivamente tales que

$$\begin{aligned} \sum_i x_i P_{i.} &= \sum_j y_j P_{.j} = 0 \\ \sum_i x_i^2 P_{i.} &= \sum_j y_j^2 P_{.j} = 1 \end{aligned} \quad (3.9)$$

Dadas estas condiciones

$$u = \frac{\sum_i \sum_j x_i y_j P_{ij}}{\sum_i \sum_j x_i^2 P_{i.} \sum_j \sum_l y_l^2 P_{.l}} \quad (3.10)$$

con lo que este parámetro representa la correlación entre las  $x$ 's y  $y$ 's. De este hecho se desprende el nombre de estos modelos. Por otra parte

$$(P_{ij} - P_{i.} P_{.j}) / (P_{i.} P_{.j}) = u x_i y_j \quad (3.11)$$

es decir que  $u$  muestra la importancia de  $x_i y_j$  en la explicación de la dependencia; por último

$$\begin{aligned} \sum_i \sum_j (P_{ij} - P_{i.} P_{.j})^2 / (P_{i.} P_{.j}) &= u^2 \\ \sum_j y_j (P_{.j} / P_{i.}) &= u x_i \\ \sum_i x_i (P_{i.} / P_{.j}) &= u y_j \end{aligned} \quad (3.12)$$

Del modelo RC general de correlación, cuando los renglones y columnas tienen un orden específico se tiene

$$\begin{array}{ll}
 \text{Modelo U} & x_i - x_{i+1} = A_i, y_j - y_{j+1} = A_j \\
 \text{Modelo R} & y_j - y_{j+1} = A_j \\
 \text{Modelo C} & x_i - x_{i+1} = A_i
 \end{array} \quad (3.13)$$

Siendo el objetivo el seleccionar uno de los modelos propuestos en el sentido de encontrar un sistema de calificaciones  $\{x_i\}$   $\{y_j\}$  que mejor ajuste a los datos en la tabla de contingencias, en el ajuste de este tipo de modelos se utilizan la estadística ji-cuadrada de bondad de ajuste y el cociente de verosimilitudes para datos categóricos en el caso asintótico.

Por otra parte, para establecer la relación de este tipo de modelos con el A.C. se define el modelo RC de correlación saturado como

$$\begin{aligned}
 P_{ij} &= P_{i.} P_{.j} (1 + \sum_k u_k x_{ik} y_{jk}) \quad (3.14) \\
 \sum_i x_{ik} P_{i.} &= \sum_j y_{jk} P_{.j} = 0 \\
 \sum_i x_{ik}^2 P_{i.} &= \sum_j y_{jk}^2 P_{.j} = 1 \\
 \sum_i x_{ik} x_{ik'} P_{i.} &= \sum_j y_{jk} y_{jk'} P_{.j} = 0 \quad k \neq k'
 \end{aligned}$$

con lo que

$$u_k = \sum_i \sum_j x_{ik} y_{jk} P_{ij} \quad k=1, \dots, K. \quad (3.15)$$

La parametrización que maneja el A.C. es

$$P_{ij} = \sum_k P_{ik} P_{jk} (1 + \frac{x_{ik}^* y_{jk}^*}{u_k})$$

con

(3.16)

$$x_{ik}^* = \frac{x_{ik}}{u_k}$$

$$y_{jk}^* = \frac{y_{jk}}{u_k}$$

Así las restricciones sobre las calificaciones se modifican a

$$\sum_i x_{ik}^2 = \frac{u_k^2}{P_{ik}}$$

(3.17)

$$\sum_j y_{jk}^2 = \frac{u_k^2}{P_{jk}}$$

Con esto, el Análisis de Correspondencias es un caso especial del modelo RC de correlación saturado. En vista de lo anterior si las herramientas utilizadas para el caso del modelo RC no saturado con  $K=1$  son generalizadas para  $K \geq 1$  ( $1 \leq K \leq K$ ), el problema inferencial sobre los valores singulares del A.C. podrá ser atacado mediante este camino. Sin embargo, sobre la generalización Goodman sólo la señala como una posibilidad citando referencias para el caso; por otra parte la comparación entre los resultados del A.C. y los arrojados por el ajuste de sus modelos de correlación se hace en una tabla en la que el modelo no saturado (3.B) es suficiente estadísticamente, quedando por ver un caso en que al menos sean necesarios dos valores singulares en la explicación aceptable de la tabla. Como ya se mencionó, el objetivo del ajuste de los modelos de Goodman es la

validación de un sistema de calificaciones, por lo que al considerar los modelos no saturados con una sola dimensión sólo se supone que hay un único sistema; sin embargo en tablas en las que no baste una sola dimensión para su ajuste satisfactorio será necesario considerar la posibilidad de varios sistemas de calificaciones producidos no necesariamente por el mismo modelo. Para poder determinar que las últimas  $s$  raíces son iguales a cero estadísticamente, habría que ajustar el modelo saturado con  $p$  y con  $p-s$ , para posteriormente por medio del cociente de razón de verosimilitudes, saber cual es el adecuado de los dos. Sin embargo en términos de las restricciones en (3.14) a cumplir por las calificaciones y por las raíces, esto resulta ser muy complejo por el momento.

Una línea interesante de investigación consiste en la concreción de la generalización de las herramientas manejadas por Goodman, para enfrentar la situación señalada además de enfrentar el problema inferencial propio del A.C. en lo referente a la significancia de los valores singulares.

### 3.4 Comentarios.

Como se ha visto en lo expuesto en este tercer capítulo, el

problema de la inferencia en el A.C. ha sido visualizado en dos perspectivas. Se tiene en primer lugar la obtención de las distribuciones asintóticas de los valores singulares como lo hacen los trabajos de Corsten(1976), Lebart(1976) y O'Neill(1978,1981). La segunda perspectiva ataca el problema inferencial respecto a estos valores singulares desde un enfoque de mayor generalidad y por ende de mayor estructura, es decir, el trabajo de Goodman(1985).

Respecto al primer enfoque, especialmente respecto a los resultados de Lebart y Corsten, el hecho de que se parta de la suposición de independencia entre los criterios de clasificación limita fuertemente su uso en la práctica, en donde el caso de dependencia es el de interés en términos generales; los trabajos de O'Neill si bien suplen esta carencia al proveer de distribuciones en el caso de dependencia, su uso estará condicionado al cumplimiento de condiciones que en la práctica no es posible verificar. En cuanto al tratamiento de Goodman, se tiene el inconveniente de que aún no se han generalizado las herramientas de ajuste, con lo que no se ha tenido la posibilidad de comparar resultados en situaciones más generales como las que enfrenta el A.C., por lo que en el análisis de tablas que

requieran para su representación de más de una dimensión no se conoce con certeza el comportamiento de los modelos propuestos por este autor.

En el siguiente capítulo se presenta una propuesta de análisis inferencial sobre los valores singulares a partir de las ideas de O'Neill (1978) y combinándolas con elementos de técnicas de muestreo. Se decidió trabajar en este sentido porque se lo consideró más directo y menos complicado que el enfoque de Goodman. Sin embargo los comentarios sobre las posibilidades de generalización de este último debe considerarse como una línea futura de investigación.

Por otra parte, debe resaltarse el hecho de que los resultados presentados sólo han versado sobre los valores singulares. El A.C. si bien considera como parte importante el reconocimiento del número de dimensiones estadísticamente significativas, contiene otros elementos en los que sería interesante contar con herramientas de tipo inferencial. En el capítulo cuatro se presenta una estadística que permitirá realizar inferencias sobre los posibles agrupamientos de los perfiles de una misma característica, ya una vez proyectados sobre sus planos de mejor

ajuste. Nuevamente esto se hace a través de remuestros y resultados de teoría asintótica.

#### 4. Inferencia estadística de resuestros.

Como se observó en el capítulo tres de este trabajo, los métodos de inferencia presentados sufren de inconvenientes ya que o se parte de supuestos poco realistas o aún no se han desarrollado a plenitud las posibilidades de su uso específico en el A.C.. Otro aspecto a resenar es que los esfuerzos han sido dirigidos a realizar inferencias sólo sobre la dimensionalidad de la representación, es decir, sobre la determinación estadística de cuales valores singulares se pueden considerar igual a cero, sin tomar en cuenta características de la representación gráfica. En este capítulo se persiguen dos objetivos; el primero respecto a proponer un algoritmo para probar dimensionalidad del A.C., mientras que el segundo consiste en realizar inferencias sobre la representación gráfica producida.

Así, en primer lugar se comenta sobre los fundamentos generales del método bootstrap desarrollados por Efron (1977), como las condiciones específicas a cumplir en los casos trabajados más adelante. En segundo lugar, aprovechando algunos resultados de O'Neill (1978) y combinándolos con las propiedades de la metodología de muestreo bootstrap, se presenta un algoritmo

estadístico para probar dimensionalidad sobre los valores singulares. En el mismo sentido también se reportan los resultados debidos a Beran (1985), respecto al uso del método bootstrap en inferencia respecto a porcentajes de variación recuperada con un número de valores singulares seleccionado de antemano.

En tercer lugar, como se mencionó en el capítulo tres, es necesario desarrollar elementos inferenciales en otros aspectos de relevancia práctica en el A.C., en este sentido se presenta una estadística que permitirá probar estadísticamente si dos perfiles, ya una vez proyectados en el plano de mejor ajuste, se pueden considerar iguales o no. En el desarrollo de esta estadística se utilizan resultados de teoría asintótica y elementos de muestreo tipo bootstrap.

#### 4.1 El Método Bootstrap.

Suponga que se tiene una muestra aleatoria de tamaño  $n$  de una función de distribución de probabilidad  $F$  no especificada,

$$\{X_i = x_i, X_i \sim F \text{ iid}, i=1, \dots, n\} \quad (4.1)$$

Dada una función  $R(X, F)$  el problema es estimar la distribución muestral de  $R$  en base a la observación de una muestra  $x$ . El método bootstrap en este caso es el siguiente (Efron, 1977):

i) Construya la distribución de probabilidad muestral  $F_n$  asignando  $(1/n)$  de probabilidad a cada  $x_i, i=1, \dots, n$ .

ii) Dada  $F_n$ , genere una muestra de tamaño  $n$

$$(X_1^*, X_2^*, \dots, X_n^*) \sim F_n, \text{ i.i.d.}, i=1, \dots, n$$

llamada la muestra bootstrap.

iii) Aproxime la distribución de  $R$  por la distribución bootstrap de  $R(X_n^*, F_n)$ .

Para el cálculo de la distribución de  $R(X_n^*, F_n)$  se tienen tres alternativas.

- i) Cálculo teórico directo.
- ii) Aproximación Monte Carlo de la distribución bootstrap.
- iii) Métodos de expansión por series de Taylor.

En el trabajo desarrollado en esta sección se usó ii), es decir, usando masivamente la computadora.

Por otra parte, un aspecto importante a señalar es que este método de resuestreo es aplicable también en el caso en que se conozca la familia paramétrica a la que pertenece  $F$  sólo desconociendo los valores de los parámetros que la indexan (Efron, 1977, observación K). Este es el caso de las aplicaciones de este método en los desarrollos presentados más adelante, ya que se está trabajando con tablas de contingencias que son gobernadas por la distribución Multinomial; así al considerar  $F$  se debe

pensar en un miembro de la familia Multinomial, es decir, si  $X = \text{vec}(N)$  vector de  $L$  - entradas entonces

$$P(X=x) = (n! / \prod_{i=1}^L x_i!) \prod_{i=1}^L p_i^{x_i}$$

donde

$$L = (p+1)(q+1), n = \sum_{i=1}^L x_i$$

$$x = (n_1, \dots, n_L), n_i = \sum_{j=1}^q (X_{ij} = 1) \quad (4.2)$$

$$p_i = \text{Prob}\{X_{ij} = 1\} \quad i=1, \dots, L; j=1, \dots, q,$$

denotada por  $M(n, (p_i, i=1, \dots, L))$ .

En este caso se tendrán que estimar  $(p_i)$ , con lo que

$$F_n = M(n, (\hat{p}_i, i=1, \dots, L))$$

con

$$\hat{p}_i = n_i / n \quad i=1, \dots, L$$

los estimadores máximo verosímil correspondientes (Lindgren, 1976).

En otro sentido, para usar el método bootstrap se requiere del cumplimiento de ciertas condiciones de regularidad de la función  $R(\dots)$ . Así, Efron (1977, observación 6), en el caso de que el espacio muestral sea finito con lo que la función  $F$  puede representarse como una tabla de probabilidades cuyos elementos son los correspondientes al vector  $L$ -dimensional  $P = \text{vec}(P)$ , define

$$Q(\hat{P}_V, P_V) = R(\hat{X}_V, F_n) \quad (4.3)$$

Ahora dado que

$$P_v^A \mid P_v \sim M(n, P_v) \quad (4.4)$$

y

$$P_v^B \mid P_v \sim M(n, P_v)$$

para muestras grandes  $\hat{P}_v^A$  estará cerca de  $P_v$  con lo que la distribución de (4.3) será razonablemente aproximada por aquella generada por la muestra bootstrap. En cuanto a la validez asintótica de esta aproximación se tienen que cumplir con las condiciones de regularidad (4.5):

- i.  $Q(P_v, P_v) = 0$
- ii.  $u(P_v, \hat{P}_v) = (\partial / \partial P_{v,1} [Q(P_v, \hat{P}_v)])$   
 exista continuamente para  $(P_v, \hat{P}_v)$  en una  
 vecindad abierta de  $(P_v, P_v)$ . (4.5)
- iii.  $u = u(P_v, P_v) \neq 0$
- iv.  $Q(\dots)$  no dependa de  $n$ .

Bajo estas condiciones y aplicando el teorema de Taylor (Serfling, 1980), y el hecho de que  $P_v^B$  y  $\hat{P}_v^A$  converge con probabilidad uno a  $P_v$  se tiene que

$$Q(\hat{P}_v^A, P_v) = (P_v^A - P_v) (u + o_n) \quad (4.6)$$

y

$$Q(P_v^B, \hat{P}_v^A) = (P_v^B - P_v) (u + o_n),$$

donde  $o_n$  convergen a cero con probabilidad uno. Ahora dado (4.4) y el

hecho de que  $\hat{P}_V^A$  converge a  $P_V$  con probabilidad uno,

$$n^{(1/2)} (\hat{P}_V^A - P_V) | P_V \sim AN(Q, \Omega_P) \quad (4.7)$$

y

$$n^{(1/2)} (P_V^A - P_V) | P_V \sim AN(Q, \Omega_P)$$

con  $\Omega_P$  matriz de covarianzas con elementos  $p_{v,i} (\delta_{in} - p_{v,n}^2)$ ,  $\delta_{in} = 1(0)$  si  $1 \leq i, n \leq m$ . De (4.6) y (4.7)

$$\sqrt{n} Q(P_V^A, P_V) | \hat{P}_V^A \quad (4.8)$$

es asintóticamente equivalente a

$$\sqrt{n} Q(\hat{P}_V^A, P_V) | P_V \quad (4.9)$$

ambas siendo asintóticamente Normales

$$N(Q, \Omega_P - \Omega_{P_V}). \quad (4.10)$$

Así, bajo el cumplimiento de las condiciones (4.5), el método bootstrap podrá utilizarse. En las secciones siguientes se hacen dos aplicaciones del enfoque bootstrap en el A.C. y en cada caso se hacen los señalamientos necesarios para satisfacer (4.5).

#### 4.2 Valores Singulares y el Método Bootstrap.

En el capítulo 3 se presentaron los resultados de la teoría asintótica en el caso de interés práctico, es decir, cuando existe dependencia entre los criterios de clasificación y se suponen distintos los valores singulares no cero. Sin embargo, se hizo

énfasis que el uso de estos resultados depende del conocimiento de la estructura teórica de la tabla considerada, lo cual es irreal en las aplicaciones.

Aun así, los resultados asintóticos pueden combinarse con el enfoque bootstrap, para dar respuestas útiles en la práctica.

Como se mencionó en el capítulo tres, los estimadores de los valores singulares centrados (3.7) son variables aleatorias cuya distribución asintótica es Normal en el caso de que los valores teóricos correspondientes sean distintos de cero. Basándose en este resultado, al generar

$$P_b^* \quad b=1, \dots, B \quad (4.11)$$

matrices de frecuencias relativas bootstrap, y que a partir de estas se obtenga

$$\{ \hat{u}_k^b, k=1, \dots, p \} \quad b=1, \dots, B \quad (4.12)$$

si para  $k$  fija, la distribución empírica corresponde estadísticamente a la distribución Normal se tendrá un elemento para afirmar que el correspondiente  $u_k$  es distinto de cero.

Este procedimiento será aplicable si las funciones  $Q_k$  cumplen con las condiciones de regularidad (4.5) presentadas en la introducción de este capítulo. Así, en forma trivial se tiene que

$$Q_k(u_k, u_k) \neq 0 \quad k=1, \dots, K$$

cumpliéndose la primera condición de regularidad, respecto a que las funciones  $D_k$  sean continuamente diferenciables respecto a  $P_k$ , resulta del hecho de que los valores singulares pueden verse como raíces de polinomios, que por teoría en esta materia (Lancaster, 1969), resultan cumplir con la condición requerida. El cumplimiento de la tercera condición se da ya que de otra manera los estimadores contrados de los valores singulares no serían asintóticamente Normales. Por último estandarizando por  $n^{1/2}$ , se cumple la cuarta condición.

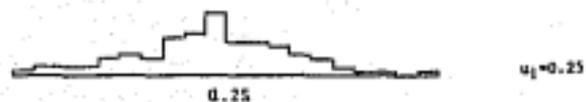
De esta manera la estadística a aplicarle el remuestreo será

$$\left( \hat{u}_k - u_k \right) \quad k=1, \dots, p. \quad (4.13)$$

Este planteamiento fue originalmente sugerido por la observación de los histogramas generados para cada valor singular a partir de un experimento (O'Neill, 1978, ejemplo 2) con 191 ensayos Monte Carlo (figura 4.1), los histogramas correspondientes a valores teóricos cero muestran un sesgo mientras que el correspondiente al valor teórico distinto de cero da una apariencia simétrica. Esta situación combinada con los resultados asintóticos ha hecho que específicamente se proponga el siguiente algoritmo:

- 1) A partir de la tabla de datos originales generarse B repeticiones bootstrap de la variable aleatoria

Figura 4.1 Histogramas con 191 ensayos Monte Carlo sobre Valores Singulares.



Multinomial tomando a P como el verdadero valor del parámetro, generándose las matrices (4.11).

ii) Aplique el A.C. a cada una de las tablas produciendo (4.12).

iii) En orden ascendente, a partir del último valor singular estimado distinto de cero numéricamente hablando, aplique pruebas de bondad de ajuste para determinar si la distribución bootstrap es Normal. Deténgase en el momento en que rechace la normalidad para algún valor singular, y a partir de éste hacia arriba considere que todos son distintos de cero estadísticamente.

Este algoritmo fue probado con el siguiente ejemplo:

Para una tabla 5x5 (p=q=5), y un tamaño total de 4096 individuos a contabilizar, se parte de las fórmulas de reconstitución (3.5) de la sección 3.2,

$$P_{ij} = p_{i.} p_{.j} \left( 1 + \sum_{k=1}^K u_{ik} x_{jk} y_{jk} \right)$$

con las condiciones (4.14)-(4.15) a continuación:

$$P_{ij} = (1/16, 4/16, 6/16, 4/16, 1/16) = P_{.j}, \quad (4.14)$$

$i, j = 1, \dots, 5.$

$$\begin{aligned}
 \bar{x}^{(1)} &= (6, 0, -2, 0, 6) / \sqrt{5} & \bar{y}^{(1)} &= (0, 1, 0, -1, 0) / \sqrt{2} \\
 \bar{x}^{(2)} &= (-2, -1, 0, 1, 2) & \bar{y}^{(2)} &= (1, 0, 0, 0, -1) / \sqrt{5} \quad (4.15) \\
 \bar{x}^{(3)} &= (2, 1, 0, 1, -2) & \bar{y}^{(3)} &= (-4, 1, 0, 1, -4) / \sqrt{10}
 \end{aligned}$$

y

$$u_1 = 0.25, u_2 = u_3 = u_4 = 0.$$

es decir con  $K=1$ .

Esta tabla fue definida por O'Neill (1978) en un experimento Monte Carlo para corroborar sus resultados asintóticos con los de muestra finita; debido a tiempo de máquina sólo generó 191 ensayos Monte Carlo. Partiendo de las condiciones (4.14)-(4.15), se generaron 191 ensayos Monte Carlo con  $B=50$  repeticiones bootstrap; esto fue realizado mediante programación en la máquina Burroughs B7800 de la U.N.A.M. en el compilador fortran 4 utilizando el generador de números pseudo aleatorios propio de esta máquina. A cada una de las matrices  $P_{ij}^k$  se realizó el A.C. correspondiente obteniéndose (4.12) y para cada  $k$  se aplicó la prueba Anderson-Darling caso tres cuya descripción es:

Sea  $X_1, \dots, X_n$  una muestra aleatoria de la cual se requiere saber si estadísticamente proviene de una distribución Normal con parámetros desconocidos  $\mu, \sigma^2$ . Calcule la media y la varianza muestral. Para cada  $X_i$

calcule

$$Z_i = F(X_i | \bar{X}, S^2) \quad i=1, \dots, n.$$

donde F es la dist.acumulativa Normal.

Obtenga las estadísticas de orden de las  $Z_i$ 's para calcular

$$A_2 = -n - (1/n) \sum_{i=1}^n (2i-1) [\log(Z(i)) + \log(1-Z(n+1-i))].$$

Si  $A_2$  es menor o igual que el cuantil  $\hat{A}_2, \alpha$  entonces no se rechazará que la muestra provenga de la distribución Normal.

Si se recuerda, K denota al número de raíces mayores a cero; denótase por K a la variable aleatoria que contabiliza el número de raíces ordenadas distintas de cero detectadas por el algoritmo propuesto. Los resultados de la simulación se muestran en el cuadro 4.1. En esta simulación  $K=1$ , y de acuerdo a los resultados

$$P(K \geq 3) = \begin{cases} 0.95 & \text{significancia } 0.05 \\ 1.0 & \text{significancia } 0.01 \end{cases}$$

con lo que K resulta ser una cota superior de  $K_j$ ; sin embargo, via el algoritmo no se tomó como cero a algún valor singular que en realidad fuera positivo. Estos resultados aunque no son del todo

Cuadro 4.1. Probabilidades estimadas del número de raíces distintas de cero por medio del experimento Montecarlo con  $B=50$  repeticiones bootstrap en cada uno de los 191 ensayos.

k	$\hat{P}\{R=k\}$	
	$\alpha = 0.05$	$\alpha = 0.01$
0	0	0
1	.01	0
2	.04	0
3	.73	.49
4	.20	.51

satisfactorios, si sugieren que el algoritmo trabaja en forma conservadora evitando errores graves consistentes en no considerar a ejes principales que tengan algo que aportar en la explicación de la tabla. Tratando de explotar al máximo la información generada por la simulación, para determinar las causas en el comportamiento del algoritmo, se calcularon en forma marginal, es decir cada una de las raíces por separado sin tomar en cuenta a las restantes, los porcentajes en que se rechazó la hipótesis de que la raíz sea distinta de cero cuando lo es ( $\alpha$ ), y la proporción de veces en que se rechazó la hipótesis de que una raíz es cero cuando no lo es ( $\beta$ ). Estos porcentajes fueron calculados en base a dos niveles de significancia (0.05 y 0.01), mostrándose los resultados en la cuadro 4.2. A la vista de los resultados en este cuadro el algoritmo aplicado en forma marginal muestra tener un nivel de significancia adecuado aunque una potencia pobre para detectar ceros en el caso de raíces intermedias. Debe mencionarse que se aplicaron otras pruebas de bondad de ajuste al caso Normal y los resultados fueron menos satisfactorios que los arrojados por la prueba de Anderson-Darling.

Entre las posibles causas de los malos resultados en  $u_2$  y

Cuadro 4.2. Errores tipo 1 y tipo 2, Exp. Montecarlo 191 ensayos.

Valor Singular	$\alpha = 0.05$	$\alpha = 0.01$
$u_1 = 0.25$	$\beta = 0.03$	$\beta = 0.01$
$u_2 = 0$	$\beta = 0.92$	$\beta = 0.96$
$u_3 = 0$	$\beta = 0.94$	$\beta = 0.98$
$u_4 = 0$	$\beta = 0.19$	$\beta = 0.51$

$u_3$ , siguiendo las ideas de Efron y Tibshirani (1986), se puede argumentar que el número de repeticiones  $B=50$  es insuficiente para probar una hipótesis. En su trabajo, éstos autores afirman que en los casos en los que no sólo se quiere tener una idea de la precisión en la estimación de un parámetro, sino que se requiere construir por ejemplo un intervalo de confianza, se tendrá que incrementar sustancialmente a  $B$  para tener resultados confiables. En el presente trabajo para investigar la influencia del tamaño de  $B$ , se trabajaron dos casos; el primero fue un ensayo MonteCarlo de los generados por (4.14)-(4.15) y el segundo partiendo de la misma estructura pero considerando que el primer valor singular agrupaba el 90% de la variación total y el segundo el 10% restante, es decir,  $u_1 = 0.2236068$ ,  $u_2 = 0.0707107$ ,  $u_3 = u_4 = 0$ .

En ambos ejemplos se generaron 200, 500 y 1000 repeticiones bootstrap, construyéndose en cada caso histogramas (figuras 4.2 y 4.3), y se aplicó el algoritmo propuesto inicialmente, cuyos resultados se muestran en la cuadro 4.3.

En el primer ejemplo tanto  $B=500$  como  $B=1000$  dan resultados aceptables, si bien  $\alpha=0.01$  se presenta la tendencia a dar una cota superior. En el segundo ejemplo sólo con 1000 repeticiones se tienen resultados satisfactorios tanto en  $\alpha=0.01$  y

0.05. Observando los histogramas, se evidencia que los histogramas tienden a sesgarse cuando el valor teórico es cero, esto es marcado en los casos con  $B=1000$ . Aparentemente en estos ejemplos el número  $B$  sí tiene influencia.

En vista de lo anterior, será recomendable para aplicar el algoritmo realizar  $B \geq 1000$  repeticiones bootstrap.

Esta afirmación deberá ser corroborada en base a un experimento Monte Carlo con un número considerable de ensayos, experimento que no fue realizado por limitaciones en tiempo. Sin embargo los resultados apoyan los comentarios de Efron y Tibshirani en el sentido de mayor precisión a mayor  $B$ .

Una última propuesta del uso del enfoque Bootstrap se debe a Baran (1985); él propone trabajar con

$$r_k = \frac{\hat{0}_k^2}{\hat{0}_k^2} \frac{B}{k} \hat{0}_k^2 \quad (4.16)$$

es decir con el porcentaje de variación explicado por el eje  $k$ . Esta estadística si bien no es utilizable en términos de saber estadísticamente si un valor singular es cero, permite establecer cuantiles de la proporción  $r_k$  en base a la aproximación de su distribución bootstrap obtenida por simulación.

**Duadro 4.3. Efecto del tamaño de B en el algoritmo de dimensionalidad.**

**Ejemplo 1.**  $u_1 = 0.25, u_2 = u_3 = u_4 = 0$

B	0.05	0.01
200	para en $u_2$	para en $u_3$
500	para en $u_1$	para en $u_3$
1000	para en $u_1$	para en $u_2$

**Ejemplo 2.**  $u_1 = 0.2236068, u_2 = 0.0707107, u_3 = u_4 = 0$

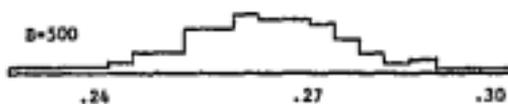
B	0.05	0.01
200	para en $u_2$	para en $u_3$
500	para en $u_1$	para en $u_2$
1000	para en $u_2$	para en $u_2$

Figura 4.2 Ensayo Montecarlo 1. Sensibilidad ante cambios en  $\beta$ .

Valor Teórico  $u_1=0.25$



sesgo 0.41  
curtosis 2.89  
 $\hat{\lambda}_2$  0.84



sesgo 0.19  
curtosis 2.93  
 $\hat{\lambda}_2$  0.52



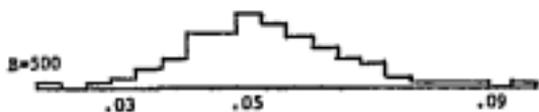
sesgo 0.11  
curtosis 2.90  
 $\hat{\lambda}_2$  0.37

Figura 4.2 (continuación). Ensayo MonteCarlo 1

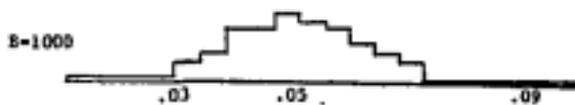
valor teórico  $\mu_2=0$



B=200  
sesgo 0.44  
curtosis 2.77  
 $\hat{A}_2$  0.98



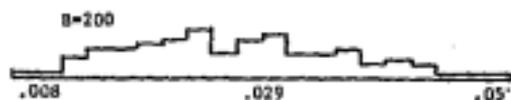
B=500  
sesgo 0.33  
curtosis 3.21  
 $\hat{A}_2$  0.90



B=1000  
sesgo 0.28  
curtosis 3.23  
 $\hat{A}_2$  0.86

Figura 4.1 (continuación)

valor teórico  $u_3=0.0$



sesgo 0.21  
curtosis 2.30  
 $\hat{\lambda}_2$  0.78



sesgo 0.25  
curtosis 2.65  
 $\hat{\lambda}_2$  0.97



sesgo 0.24  
curtosis 2.68  
 $\hat{\lambda}_2$  1.38

Figura 4.2 (continuación)

valor teórico  $u_4=0.0$

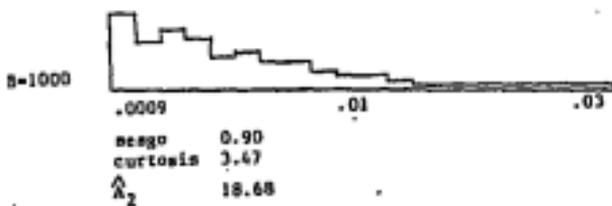
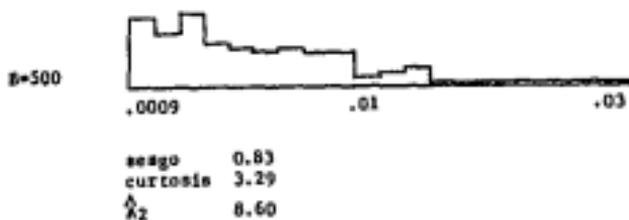
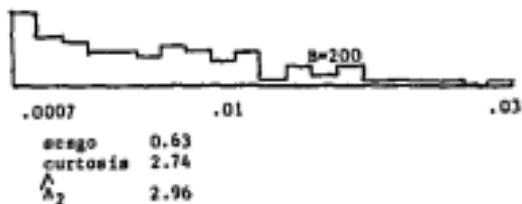
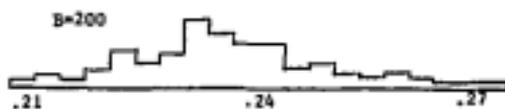


Figura 4.3 Ensayo MonteCarlo 2. Sensibilidad ante cambios en B.

valor teórico  $\mu_1=0.2236068$



sesgo	0.36
curtosis	2.91
$\hat{\Delta}_2$	0.67



sesgo	0.23
curtosis	2.81
$\hat{\Delta}_2$	0.72



sesgo	0.17
curtosis	2.83
$\hat{\Delta}_2$	0.65

Figura 4.3 (continuación).

valor teórico  $\lambda_2=0.0707107$

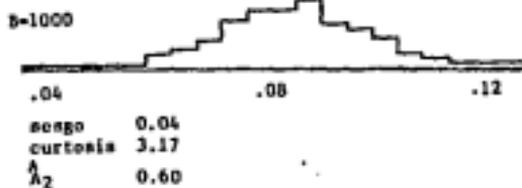
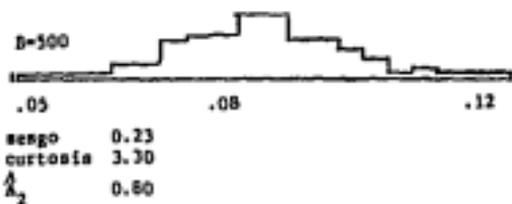
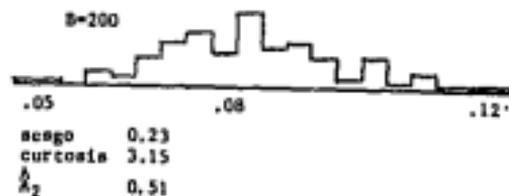


Figura 4.3 (continuación)

valor singular  $u_3=0.0$

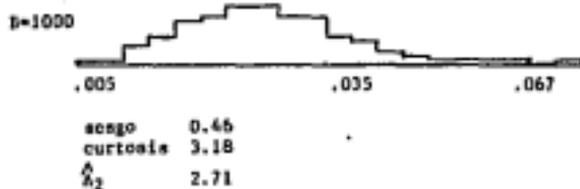
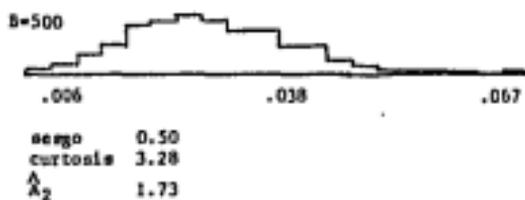
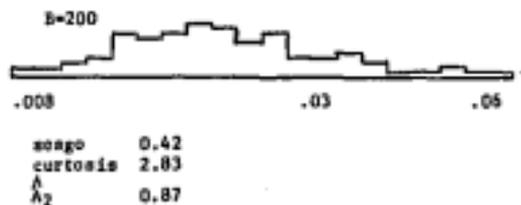
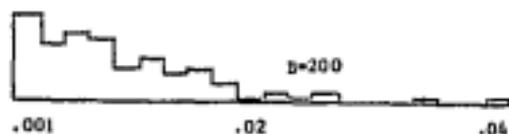
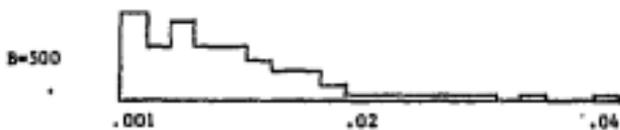


Figura 4.3 (continuación)

valor teórico  $u_q=0.0$



sesgo 1.12  
 curtosis 4.90  
 $\hat{\lambda}_2$  3.52



sesgo 1.12  
 curtosis 5.05  
 $\hat{\lambda}_2$  7.99



sesgo 1.02  
 curtosis 4.21  
 $\hat{\lambda}_2$  18.34

#### 4.2 Agrupamientos de Perfiles en el Plano de Mejor Ajuste.

Como se mencionó en el capítulo dos de este trabajo, el objetivo primordial del A.C. es la representación gráfica de las relaciones relevantes en una tabla de contingencias con dos criterios de clasificación. En cuanto al problema inferencial la mayoría de los esfuerzos se han centrado en los valores singulares; por otra parte, en la práctica la interpretación de la cercanía en el plano de mejor ajuste de dos perfiles renglón o columna es uno de los aspectos de mayor preponderancia en el A.C.. En este sentido se cree que sería útil contar con una herramienta inferencial que nos permita trabajar sobre los posibles agrupamientos de perfiles. Greenacre (1984), ha utilizado el enfoque bootstrap en lo que él llama la estabilidad de la representación en dos dimensiones. Específicamente, a partir de la tabla original genera (4.11), obteniendo en cada caso las coordenadas proyectadas sobre el plano de mejor ajuste de los perfiles originales, mediante la fórmula de transición

$$\frac{A}{F}_b = D^{-1} P \frac{A}{G}_b D^{-1} \quad b=1, \dots, B \quad (4.17)$$

donde  $\frac{A}{F}_b$  son las coordenadas de los perfiles renglón de la repetición  $b$ , pero en el espacio ortogonal correspondiente a las

coordenadas de los perfiles columna en su espacio ortonormal original  $\hat{D}$  y  $D_0$  la matriz diagonal de valores singulares mayores que cero en forma análoga para los perfiles columna se realiza el procedimiento mediante la respectiva fórmula de transición. De un análisis visual Greenacre obtiene conclusiones sobre la estabilidad. Para mayores detalles consúltese Greenacre (1984, cap. 8). Esto presenta problemas, ya que sus conclusiones dependen fuertemente de las intersecciones entre las nubes de perfiles proyectados, hecho que puede tener implicaciones de poca claridad.

Se denota a  $(\hat{f}_{11}, \hat{f}_{12})$  como las coordenadas del  $i$ -ésimo perfil respecto al plano de mejor ajuste. Considere el vector de diferencias por coordenadas en el plano de mejor ajuste de dimensión dos entre el perfil renglón  $i$  y el perfil renglón  $i'$ , es decir,

$$\hat{d}_{11, i'} = (\hat{f}_{11} - \hat{f}_{11, i'}, \hat{f}_{12} - \hat{f}_{12, i'}) \quad (4.18)$$

y a la estadística

$$\hat{d}_{11, i'}^T \hat{\Sigma}_{11, i'}^{-1} \hat{d}_{11, i'} \quad (4.19)$$

donde  $\hat{\Sigma}_{11, i'}$  es la matriz de covarianzas de (4.18). Esta estadística será utilizable en la práctica en la medida que se conozca esta matriz. El método bootstrap es utilizado para estimar



2x2, es decir dos perfiles renglón y dos perfiles columna. En este caso se tendrá una recta de mejor ajuste, en vista de (4.20)

$$\begin{pmatrix} \hat{f}_{11} \\ \hat{f}_{21} \end{pmatrix} \begin{pmatrix} 1/r_1 & 0 \\ 0 & 1/r_2 \end{pmatrix} \begin{pmatrix} p_{11} & p_{12} \\ p_{21} & p_{22} \end{pmatrix} \begin{pmatrix} \hat{q}_{11} / \hat{\Delta}_1 \\ \hat{q}_{21} / \hat{\Delta}_1 \end{pmatrix}$$

por lo que, por ejemplo

$$\partial \partial p_{11} \begin{pmatrix} \hat{f}_{11} \\ \hat{f}_{21} \end{pmatrix} = \begin{pmatrix} [p_{12} (\hat{q}_{11} - \hat{q}_{21})] / [r_1^2 \hat{\Delta}_1] \\ 0 \end{pmatrix}$$

con lo que para que sea distinta de cero se requiere de que

$$q_{11} \neq q_{21} \text{ y } p_{12} \neq 0.$$

Para tablas de mayores dimensiones sin embargo, se espera que las condiciones no sean tan fáciles de manejar.

Por último, la cuarta condición de (4.5) se cumple trivialmente.

Dado lo anterior se justifica el uso del enfoque bootstrap, ya que así como se tiene la equivalencia asintótica de la distribución de  $\hat{g}_{11}^A$  con la ley de  $g_{11}^b$ , sin embargo (Serfling, 1980), la convergencia en distribución implica convergencia en momentos de orden dos si  $\{g_{11}^b\}^2$  es uniformemente integrable.

Para ello será suficiente mostrar que  $H(p)$  es uniformemente integrable. Se tiene la siguiente demostración:

Existe una vecindad cerrada alrededor de  $P$  tal que  $H(\cdot)$  está bien

definida sobre ella. Sea  $\sqrt{P}$  tal vecindad.  $H(\cdot)$  es continua sobre  $\sqrt{P}$  y por lo tanto es acotada. Ahora,

$$P \xrightarrow{A} P \text{ c.p.1}$$

cuando el tamaño de muestra tiende a infinito y, por lo mismo,

$$P \xrightarrow{A} \hat{P} \text{ c.p.1}$$

cuando  $n$  tiende a infinito; entonces

$$P \xrightarrow{A} P \text{ c.p.1}$$

Si definimos

$$P^0 = \begin{cases} P & \text{si } P \in \sqrt{P} \\ P & \text{si } P \notin \sqrt{P} \end{cases}$$

entonces

$$P^0 \in \sqrt{P}$$

siempre, y

$$n \int (H(P_n) - H(P^0)) \xrightarrow{A} 0 \text{ c.p.1} \quad \int > 0$$

por lo que  $H(P_n)$  y  $H(P^0)$  son asintóticamente equivalentes. Por lo tanto  $H(P^0)$  converge en distribución a lo mismo que  $H(P_n)$ .

$H(P^0)$  es acotada, por lo tanto es uniformemente integrable (apéndice 2); también, son uniformemente integrable, por teorema A sección 1.4 de Serfling (1980), los momentos de  $H(P^0)$  convergen a los momentos con respecto a la distribución límite, por la equivalencia asintótica también convergerán a lo

siano los momentos de  $H(P_i)$ .

Por lo tanto,  $\{d_{ii}^2\}$  es uniformemente integrable.

Así el estimador bootstrap  $\hat{\Phi}_{ii}^{-1}$  converge a  $\Phi_{ii}^{-1}$ , la matriz de covarianzas de  $\hat{d}_{ii}$ .

Dado todo lo anterior la estadística propuesta es

$$D_{ii}^2 = \hat{d}_{ii}^{-1} \cdot \hat{\Phi}_{ii}^{-1} \cdot \hat{d}_{ii} \quad (4.22)$$

donde

$$\hat{\Phi}_{ii}^{-1} = (B-1)^{-1} \sum_{b=1}^B (d_{i,jb} - \bar{d}_{i,j}) (d_{i,jb} - \bar{d}_{i,j})'$$

En cuanto a la distribución de ésta se tienen los siguientes argumentos:

i)  $\text{vec}(\hat{P}-P)$  es asintóticamente Normal con media cero (Serfling, 1980 pág 109).

ii) Dada la justificación del uso del método bootstrap, se cumplen las condiciones para que  $\hat{d}_{ii}$  sea una variable aleatoria con distribución asintótica Normal (Serfling, 1980, pág. 122 teorema A).

iii) En base a ii) y Serfling (1980, pag. 130), asintóticamente

$$\hat{d}_{ii} \cdot \hat{\Phi}_{ii}^{-1} \cdot \hat{d}_{ii} \sim \chi^2_{(2, d_{ii})} \cdot \Phi_{ii}^{-1} \cdot d_{ii} \quad (4.23)$$

iv) Bajo la hipótesis  $H_0$  de que el perfil  $i$  es igual al perfil  $i'$  en el plano de mejor ajuste se espera que  $d_{ii} = 0$  con lo que

$$\frac{A_{d_{ii}}}{d_{ii}} \cdot \frac{I_{ii}^{-1}}{I_{ii}} \cdot \frac{A_{d_{ii}}}{d_{ii}} \sim \chi^2_{(2)} \quad (4.24)$$

Así, el perfil  $i$  y el perfil  $j$  son considerados estadísticamente iguales en el plano de mejor ajuste a un nivel  $\alpha$  si

$$\frac{A_{d_{ii}}}{d_{ii}} \cdot \frac{I_{ii}^{-1}}{I_{ii}} \cdot \frac{A_{d_{ii}}}{d_{ii}} < \chi^2_{(2), \alpha}$$

## 5. Aplicaciones.

En la presente sección se presentan dos aplicaciones en el área de tecnología de alimentos, específicamente en el campo de la evaluación sensorial. Se detallan tanto el uso descriptivo de la técnica como las herramientas de corte inferencial discutidas en la sección anterior. A continuación se presenta la descripción general de la evaluación sensorial de una fruta llamada Jiotilla. Interesa el estudio de esta fruta considerada exótica y perteneciente a la familia de las cactáceas, por sus propiedades alimenticias y porque representa una fuente de ingresos adicionales para las comunidades rurales que tienen acceso a ella. Se prepararon 13 tratamientos de esta fruta:

TRATAMIENTO	IDENTIFICADOR
1. En fresco (control)	FFRE
2. Confitada en estufa	CEST
3. Confitada en sol	CSOL
4. Confitada al aire	CAIR
5. Hervida	HEHV
6. Almibar	ALMI
7. Confitada por proceso lento	CPLE
8. Confitada por proceso rápido	CPRA
9. Glasada por proceso lento	GPLE
10. Glasada por proceso rápido	GPRA
11. Confitada y envasada en cloruro de polietileno	CECP
12. Confitada y envasada en colofón	CECE
13. Confitada y envasada en cloruro de polivinilo.	CECV

Estos tratamientos fueron evaluados por un panel de 15 jueces no entrenados, que calificaron según una escala

organoléptica (hedónica) de nueve categorías:

CATEGORIA	IDENTIFICADOR
Gusta extremadamente	GEX
Gusta mucho	GMO
Gusta moderadamente	GMD
Gusta ligeramente	GLI
Indiferente	IND
Disgusta ligeramente	DLI
Disgusta moderadamente	DMD
Disgusta mucho	DMO
Disgusta extremadamente	DEX

Esto fue realizado en olor, sabor, color, textura y aspecto de cada uno de los trece tratamientos aplicados a la fruta, se presentan a continuación los resultados de las primeras dos evaluaciones.

### 5.1 OLOB.

Los datos se muestran en el cuadro 5.1. Se evidencia la nula contabilización desde DMO hasta DEX por lo que no se consideran en el análisis. La representación producida por A.C. se muestra en sus dos primeros ejes en la figura 5.1.

En forma descriptiva se observa que entre estos dos primeros ejes se tiene un 73% de la variación en los datos; respecto a la calidad de cada perfil en la representación con los dos primeros ejes, los perfiles renglón, es decir, las categorías de la escala de preferencias, en el cuadro 5.4 las contribuciones absolutas y las correlaciones. Así, en general están bien representados en estos dos ejes, ya que las contribuciones relativas son mayores al

0.5 salvo el caso de DLI. En cuanto a la contribución a los ejes, en el primero de ellos GEX y GLI son los preponderantes sucediendo esto también en el segundo eje, es decir, que ambos ejes son representantes de las posiciones relativas al gusto extremo y al gusto ligero.

En cuanto a los perfiles columna, los tratamientos, se tiene los resultados en cuanto a contribuciones absolutas y correlaciones en el cuadro 5.5.

En vista de las correlaciones, FFRE, CAIR, MERM y ALMI son tratamientos mal representados; los demás en general no tienen mayores problemas. En cuanto a las contribuciones absolutas, CECP y GPLE son los que dan primordialmente la dirección al eje uno; respecto al segundo eje se unen al grupo anterior CEST y CECE.

En cuanto a la orientación de los trece tratamientos respecto a la escala de preferencias se tiene que GPLE descolla hacia el gusto extremo siguiendo en orden descendente CPLE, GPRA, FFRE, CPRA y CEST. Posteriormente se agrupan hacia un gusto moderado MERM, ALMI, CSOL y CAIR. Por último CECV, CECE y CECP son poco aceptados.

Utilizando las herramientas de corte inferencial comentadas en el

capítulo cuatro, en primer lugar los estimadores de los valores singulares junto con la estadística Anderson-Darling ( $A_2$ ) con 1000 repeticiones bootstrap son mostrados en el cuadro 5.2; vía el algoritmo de dimensionalidad de la sección 4.1, se tiene que los primeros cuatro valores singulares son estadísticamente distintos de cero, ya que tanto al 0.01 y al 0.05 se dio esta situación.

Con la misma información proporcionada por las 1000 repeticiones bootstrap, en base a (4.16) se construyeron intervalos al 95% para porcentajes de variación marginal y acumulada para cada uno de los valores singulares, mostrándose estos resultados en el cuadro 5.3, con lo que, la representación global en los dos primeros ejes oscilará entre 63 y 79 % lo que es aceptable en general y se observa que nuevamente hasta el valor singular 4 no incluye el valor cero por ciento.

Por otra parte se aplicó la estadística  $D^2$  para detectar agrupamientos tanto en el caso de los perfiles renglón como de los perfiles columna. Antes de mostrar algún resultado, debe mencionarse que sólo se reportan los correspondientes a los perfiles que tuvieron una correlación acumulada mayor al 0.5 con los dos primeros ejes principales; así bien lo anterior es

arbitrario, se tomo esta determinación ya que esta estadística producía agrupamientos demasiado extendidos al considerar a los perfiles mal representados. La determinación estadística de cual perfil está bien representado en un conjunto de ojeas es un asunto que no ha sido atacado hasta ahora y que se cree será necesario tratar en un futuro.

A un nivel de significancia global de 0.036 (cada comparación 0.001 ( $\frac{9}{2}$  comparaciones)), los tratamientos se agrupan de la siguiente manera:

1. GPLE, CPLE, GPRA, CPRA, CEST.
2. CECV, CECE, CCEP.
3. CSOL como puente entre 1 y 2.

Respecto a la escala de preferencias con un nivel global de 0.01 (cada comparación con 0.001, ( $\frac{5}{2}$  comparaciones)), se tiene que no hay agrupamiento alguno.

Conjuntando ambos grupos de resultados se tiene que:

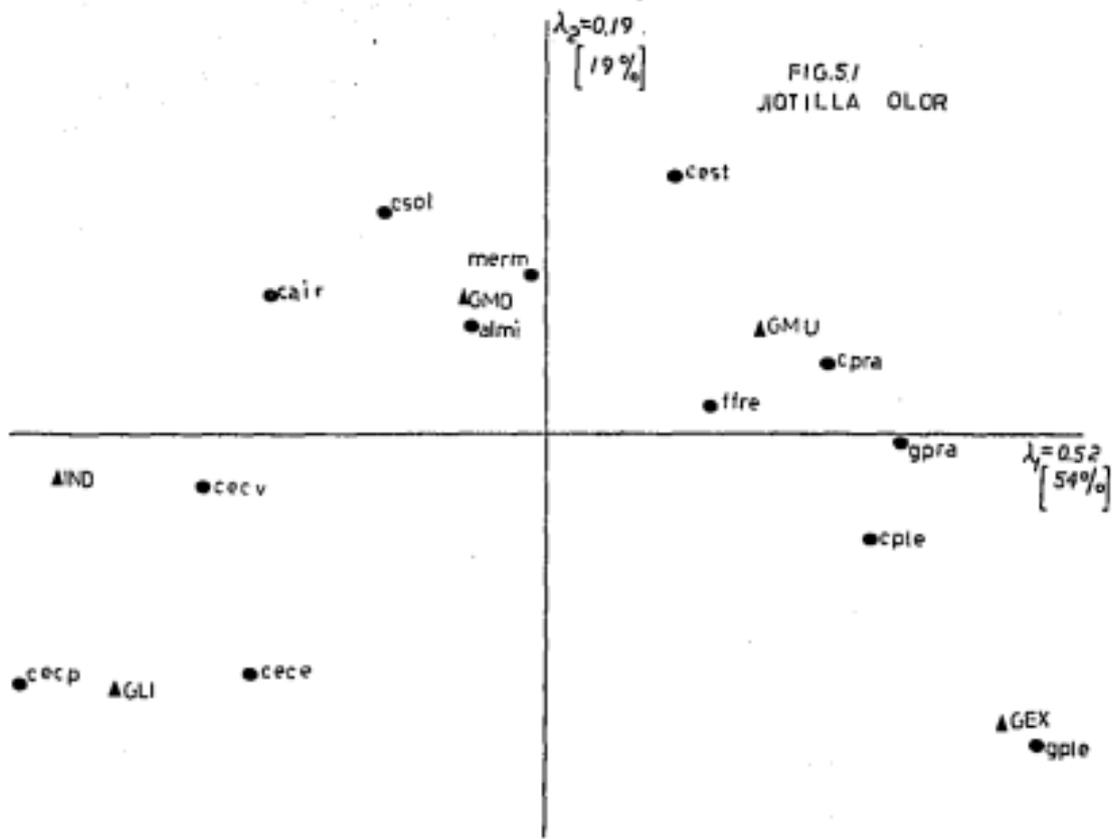
1. GPLE, CPLE, GPRA, CPRA y CEST se orientan hacia GEX y GBU.
2. CECV, CECE y CCEP se orientan hacia el GBI e IND.
3. CSOL se sitúa hacia el gusto moderado.

Los resultados generales se muestran en el cuadro 5.6.

De lo anterior, se tiene que los glaseados y los confitados tanto



FIG.5/1  
JOTILLA OLCR



**Cuadro 5.2 Estadística Anderson-Darling para los valores singulares, DLR.**

No.	valor	$\hat{A}_2$
1	0.72	0.32
2	0.46	0.31
3	0.40	0.34
4	0.26	0.19
5	0.14	57.93

**Cuadro 5.3. Intervalos de confianza para porcentajes marginales y acumulados con porcentaje marginal observado.**

	observado	marginal 11%	1e%	11%	acumulado 1e%
1	54.19	37.75	51.16	37.75	51.16
2	19.53	20.99	32.83	63.22	79.26
3	16.31	12.91	22.17	82.29	95.21
4	7.71	4.31	13.25	93.05	100.00
5	2.25	0	6.90		

**Cuadro 5.4. Contribuciones Absolutas y Correlaciones.  
Escala de preferencias, CLDR.**

	CONTRIBUCIONES ABSOLUTAS		CORRELACIONES	
	eje 1	eje2	eje1	eje2
DEX	0.314	0.381	0.648	0.283
SHU	0.149	0.101	0.604	0.148
GHO	0.024	0.181	0.153	0.417
GLI	0.279	0.276	0.621	0.221
IND	0.223	0.006	0.558	0.005
DLI	0.011	0.055	0.098	0.184

**Cuadro 5.5 Contribuciones Absolutas y Correlaciones.**  
**TRATAMIENTOS, OLDR.**

	CONTRIBUCIONES ABSOLUTAS		CORRELACIONES	
	#Jo 1	#Jo 2	#Jo 1	#Jo 2
FFRE	0.023	0.002	0.288	0.010
DEST	0.013	0.167	0.142	0.643
CSOL	0.023	0.118	0.220	0.413
CATR	0.066	0.046	0.343	0.007
MEMH	0	0.061	0.007	0.463
ALMI	0.005	0.026	0.108	0.190
CPLE	0.084	0.031	0.787	0.103
CPRA	0.043	0.010	0.789	0.045
GPLE	0.203	0.241	0.659	0.282
GPRA	0.103	0	0.762	0
CECP	0.240	0.150	0.776	0.175
CECE	0.076	0.141	0.337	0.227
CECV	0.100	0.007	0.848	0.020



por proceso lento como por proceso rápido son tratamientos de muy alta aceptación; en el mismo sentido se encuentra la jicotilla confitada por estufa. Siguiendo de este grupo de alta aceptación, se encuentra la fruta confitada al sol. Por último, los tratamientos envasados se ven devaluados en su olor.

## 5.2. Saboreo.

Los datos se muestran en el cuadro 5.7. La representación gráfica producida por el A.C. se muestra en la figura 5.2.

Descriptivamente, los dos primeros ejes contemplan el 71% de la variación en los datos; en cuanto a la representación de las categorías de la escala de preferencias se tiene en base al cuadro 5.11, los siguientes comentarios:

Se tiene que GMD y DEX tienen graves problemas en su representación. En cuanto a la contribución absoluta, en el primer eje sobresalen DLI, DMD y GMU con lo que este eje muestra la posición relativa entre el disgusto moderado y el gusta mucho; en el segundo eje nuevamente DMD pesa pero ahora superado por DLI. Respecto a los perfiles columna, es decir los tratamientos a la fruta se tiene en base al cuadro 5.12, lo siguiente:

Se tienen problemas en la representación de FFRE, CEST, CSOL, CAIR y

MERM (correlación < 0.5). En cuanto a las contribuciones absolutas CCEP y CECV son los de mayor influencia en ambos ejes, en cuanto a la orientación específica de los tratamientos a la fruta respecto a la escala de preferencias, se tiene un agrupamiento fuertemente aglutinado alrededor de GEX y GMU, estando constituido por CPLE, GPLE, GPRA y CEST siguiendo MERM, FFRE, ALMI y CSOL. Hacia GLI se tiene a CECE, CPRA, CAIR y CECV. Por último, se separa totalmente CCEP hacia DM0 y DMU.

En cuanto a la utilización de las herramientas inferenciales, se tiene en primer término que los seis primeros valores singulares pueden considerarse distintos de cero (véase cuadro 5.8).

Como en el caso de la evaluación del OLOR se calcularon intervalos al 95% de confianza para los porcentajes marginales y acumulados (cuadro 5.9). De esta manera, la representación en los dos primeros ejes en su calidad global es satisfactoria (del 58% al 75%).

Al utilizar la estadística propuesta en el capítulo anterior omitiendo a los no representados, los agrupamientos más relevantes en los tratamientos a una significancia global de 0.028 (cada comparación 0.001,  $\binom{8}{2}$  comparaciones), son:

1. CPLE, GPLE, GPRA, ALMI.

2. CPRA, CECV.

3. CAIR.

4. CCEP.

Respecto a la escala de preferencias a una significancia global de 0.021 (cada comparación 0.001,  $\binom{7}{2}$  comparaciones), se tiene:

1. GEX, GMU.

2. GLI, IND, DLI.

3. DMU, DMD.

Conjuntando ambos grupos de resultados se tiene que:

1. CPLE, GPLE, GPRA, ALMI se orientan hacia GEX-GMU.

2. CPRA, CECV orientados hacia el GLI-IND-DLI.

3. CAIR orientado hacia GLI-IND-DLI.

4. CCEP orientado hacia DMD-DMU.

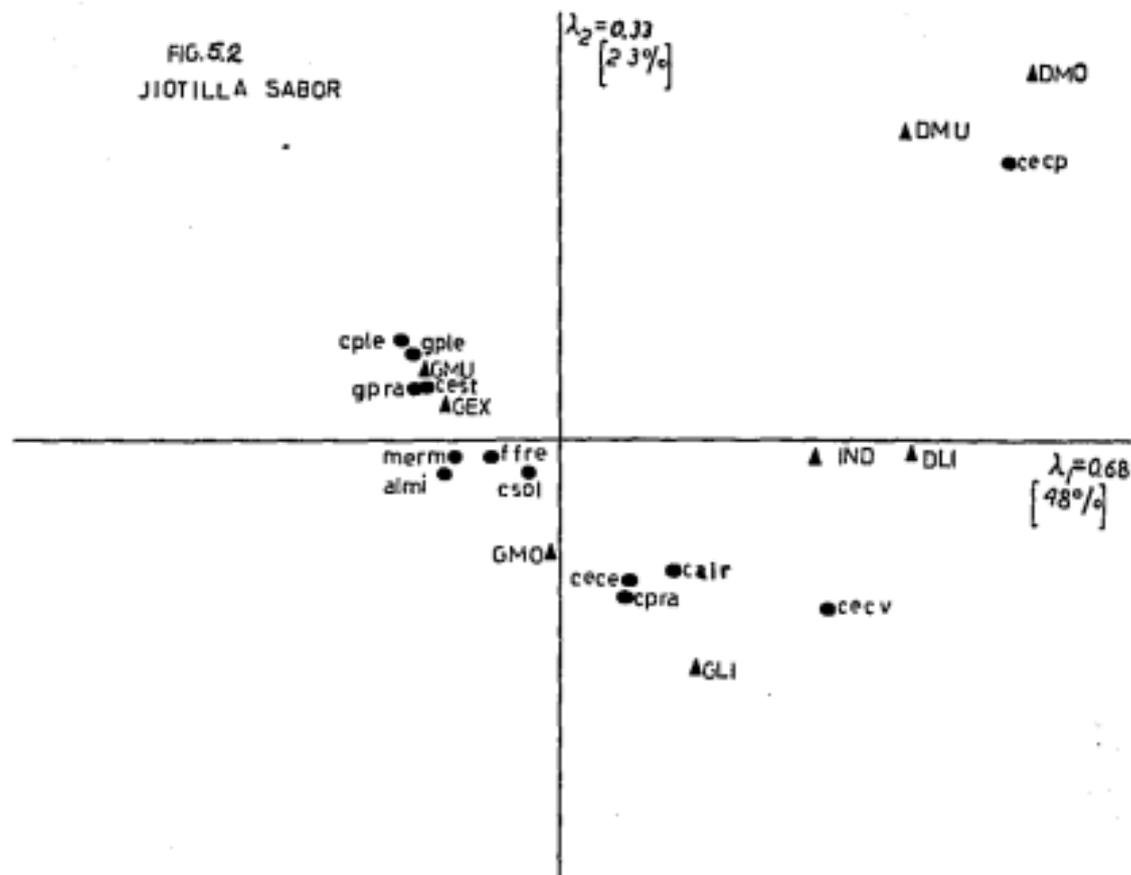
Los resultados generales se muestran en el cuadro 5.13.

El proceso lento es muy aceptado tanto en confitado como en glaseado, cosa que no sucede con el proceso rápido dado que sólo el glaseado es muy aceptado. Así, el confitado por proceso rápido, el confitado y envasado en cloruro de polivinilo y la jiotilla confitada al aire se muestran como tratamientos no aceptados. Por último, el confitado envasado en celofán se ve





FIG. 5.2  
JIOTILLA SABOR



Cuadro 5.8 Estadística Anderson-Darling para los valores singulares, SABOR.

	$\hat{\lambda}$	$\hat{A}_n (B=1000)$
1	0.82	2
2	0.57	1.17
3	0.42	0.74
4	0.31	0.50
5	0.26	0.64
6	0.20	0.40
7	0.14	0.22
8	0.03	0.88
		28.74

Cuadro 5.9. Porcentajes marginales y acumulados. SABOR.

k	marginal			acumulado	
	observado	11%	1%	11%	1%
1	48.88	36.22	48.98	36.22	48.98
2	23.01	18.42	29.02	57.76	75.28
3	12.22	9.90	18.83	73.51	86.94
4	7.19	6.06	13.09	84.88	94.40
5	5.03	3.13	9.05	92.44	98.57
6	2.16	1.08	5.56	97.08	99.85
7	1.44	0.12	2.65	99.36	100.00
8	0.07	0.00	0.61	100.00	

**Cuadro 5.10 Contribuciones absolutas y Correlaciones.  
Escala de preferencias, SABOR.**

	CONTRIBUCIONES ABSOLUTAS		CORRELACIONES	
	eje1	eje2	eje1	eje2
GEX	0.110	0.024	0.485	0.050
GMU	0.179	0.039	0.623	0.149
GMO	0	0.147	0.001	0.464
GLI	0.054	0.330	0.214	0.666
IND	0.103	0	0.646	0.001
DLI	0.277	0	0.760	0
DMO	0.219	0.280	0.579	0.356
DMU	0.058	0.099	0.418	0.340
DEX	0	0.011	0.001	0.038

**Cuadro 3.11 Contribuciones Absolutas y Correlaciones.  
Tratamientos, SAGOR.**

	CONTRIBUCIONES ABSOLUTAS		CORRELACIONES	
	eje1	eje2	eje1	eje2
FFRE	0.012	0.001	0.302	0.010
CEST	0.048	0.020	0.412	0.082
CSOL	0.003	0.03	0.018	0.008
CAIR	0.029	0.084	0.231	0.349
MERM	0.030	0.001	0.454	0.008
ALMI	0.035	0.002	0.702	0.023
CPLE	0.063	0.052	0.364	0.226
CPRA	0.008	0.126	0.084	0.610
GPLE	0.059	0.042	0.585	0.201
GPRA	0.048	0.020	0.674	0.136
CECP	0.489	0.401	0.714	0.281
CECI	0.008	0.103	0.062	0.381
CECV	0.167	0.144	0.527	0.219



#### 6. Conclusiones y Discusión.

En el desarrollo de este trabajo se realizaron principalmente cuatro tareas; la primera consistió en la exposición de los detalles del A.C. desde un punto de vista descriptivo. La segunda tarea estuvo conformada por la revisión de herramientas inferenciales ya existentes en la literatura. Como tercera parte, en base a las dos primeras se propusieron procedimientos para enfrentar distintos aspectos inferenciales en el A.C. Por último, la cuarta tarea consistió de la aplicación tanto los aspectos descriptivos como los inferenciales propuestos, a un ejemplo de análisis sensorial proveniente del área de tecnología de alimentos.

La evaluación de estas tareas es como sigue: de la primera se obtuvieron los conocimientos básicos necesarios para la correcta aplicación o interpretación de los resultados descriptivos arrojados por el A.C.; respecto a la revisión bibliográfica, se puede decir que en general las herramientas estudiadas tienen inconvenientes. En el caso de los trabajos de Lebart (1976) y Corsten (1976), el hecho de partir de la independencia entre los criterios de clasificación de la tabla, limita su uso práctico ya que la dependencia es, en términos generales, el caso común y de

interés. Referente a las distribuciones asintóticas desarrolladas por O'Neill (1978, 1981), como ya fue señalado en su momento, requiere del cumplimiento de condiciones que en la práctica no es posible verificar. En cuanto a la metodología de ajuste de modelos propuesta por Goodman (1985), si bien su objetivo primordial no es del A.C., al menos en el caso de contar con una tabla de contingencias de rango 1, provee los elementos suficientes para realizar inferencias en forma indirecta sobre la dimensionalidad; sin embargo su desarrollo para casos más generales está pendiente.

Respecto a la tercera tarea, se puede decir que las herramientas propuestas trabajan, a la luz de los resultados, en forma adecuada; en el caso del algoritmo para probar sobre dimensionalidad, la mezcla de una parte de la teoría asintótica con el método bootstrap, parece ser satisfactorio tanto en su concepto (utilizar teoría que sí es utilizable directamente y cubrir sus inconvenientes con remuestreo), como en los resultados que arroja. Por otra parte, en cuanto a la estadística  $D^2$  para detectar agrupamientos de perfiles, ésta resulta ser uno de los primeros intentos para realizar inferencias en otro aspecto que

no sean los valores singulares, sus resultados son plausibles.

En cuanto a la evaluación de la cuarta tarea, se puede decir que es la primera vez que en México se realizan este tipo de aplicaciones y que de la retroalimentación con los tecnólogos en alimentos, los resultados ofrecidos por el A.C. son fácilmente entendidos y resultan ser útiles para sus fines y objetivos de análisis. El A.C. es sólo una posibilidad de análisis entre otras posibles; así las posibilidades de aplicación son amplias.

En cuanto a líneas futuras de investigación o aspectos que quedan pendientes se tiene:

i) Diseñar y realizar un experimento Monte Carlo que considere un número mayor o igual a 1000 de repeticiones bootstrap, para tener elementos de juicio más fuertes sobre el funcionamiento del algoritmo para probar sobre la significancia de los valores singulares. Nótese que este experimento tiene una magnitud de mucha importancia.

ii) Diseñar más herramientas inferenciales en el A.C. (en este caso se pueden considerar para ello a las contribuciones absolutas y las correlaciones, en términos de la determinación estadística de la orientación de un eje, así como la calidad de la representación de un perfil en particular).

iii) Los modelos propuestos por Goodman (1985) deben ser generalizados en forma explícita para poder ver tanto el comportamiento de estos modelos en el caso de contar con más de un sistema de calificaciones como sus implicaciones para el A.C.

iv) Profundizar en la validez la aplicación del método bootstrap en el A.C.

v) Sobre la posibilidad de que la combinación de resultados ya desarrollados con el enfoque bootstrap pueda ser explotada en otras áreas, en el caso de dimensionalidad se piensa en el tema de colinealidad en análisis de regresión.

### Bibliografía.

- Anderson, T.W. (1958). "An Introduction to Multivariate Statistical Theory". John Wiley and Sons.
- Beran, R. y Brivastava, M.S. (1985). "Bootstrap tests and confidence regions for functions of a covariance matrix. *Annals of Statistics*, vol. 13 no. 95-115.
- Corsten, L.C.A. (1976). "Matrix Approximation as key to application of multivariate methods". Proceedings of the 9th International Conference of the Biometric Society, pp 61-77, Raleigh, North Carolina, USA.
- Efron, B. (1977). "Bootstrap Methods: Another look at the Jackknife". *Annals of Statistics*, vol. 7 no. 1, 1-26.
- Efron, B. y Tibshirani, R. (1986). "Bootstrap Methods for Standard Errors, Confidence Intervals and other measures of Statistical Accuracy". *Statistical Science*, vol. 1 no. 1, 54-77.
- Green E.P. y Carroll, J.D. (1976). "Mathematical Tools for Applied Multivariate Analysis". Academic Press.
- Greenacre, M.J. (1984). "Theory and Applications of Correspondence Analysis". Academic Press.
- Goodman, L.A. (1965). "The Analysis of cross-classified data having

ordered and/or unordered categories: Association models, correlation models and asymmetry models for contingency tables with or without missing entries". *The Annals of Statistics*, vol. 13, no. 1 pp 10-69.

Kirch, A. (1973). "Introduction to Statistics with Exerices". Holt, Rinehart and Winston, Inc.

Lancaster, H. O. (1958). "The Structure of Bivariate Distributions". *Annals of Mathematical Statistics*, 29, 719-736.

Lancaster, P. (1969). "Theory of Matrices". Academic Press.

Lebart, L. (1976). "The significance of Eigenvalues Issued From Correspondence Analysis of Contingency Tables", en *COMPSTAT 1976*, Physica Verlag Wien, editado por J. Gordesch y P. Naeye, pp 38-45.

Lindgren, B. W. (1976). "Statistical Theory". Collier MacMillan.

O'Neill, M. E. (1981). "Asymptotic Distributions of the Canonical Correlations from Contingency Tables". *Australian Journal of Statistics*, 20(1), 75-82.

O'Neill, M. E. (1978). "Distributional Expansions for Canonical Correlations from Contingency Tables", *Journal of the Royal Statistical Society Series B*, 40 no. 3 303-312.

Rubinstein, R. Y. (1981). "Simulation and the Monte Carlo Method". John Wiley and Sons.

Serfling, R. J. (1980). "Approximation Theorems of Mathematical Statistics". John Wiley and Sons.

Tabet, N. (1973). Programa de Análisis de Correspondencias, parte de tesis Doctoral Universidad de París VI.

### Apéndice B. Descomposición en Valor Singular.

Sea  $A$  una matriz  $p \times q$  real; entonces  $A$  puede expresarse como el producto de tres matrices

$$A = U D_a V^T \quad (a.1)$$

donde  $D_a$  es una matriz diagonal  $K \times K$  con números positivos  $a_1, a_2, \dots, a_K$ ,  $K$  el rango de  $A$ . Las matrices  $U$   $p \times K$  y  $V$   $q \times K$  tales que  $U^T U = V^T V = I$ , es decir ortogonales por columna.

En vista de esta descomposición, la matriz  $A$  puede expresarse como

$$A = \sum_{k=1}^K a_k u_k v_k^T \quad (a.2)$$

donde  $u_k, v_k$  son las columnas de  $U$  y  $V$ . Los valores  $a_1, a_2, \dots, a_K$  son llamados los valores singulares, mientras que  $(u_k)$  y  $(v_k)$  los vectores propios izquierdos y derechos.

Los vectores izquierdos constituyen una base ortonormal para las columnas de  $A$  mientras que los vectores propios derechos conforman una base ortonormal de los renglones de  $A$ .

La existencia de esta descomposición se puede demostrar a través de la descomposición propia (eigenestructura) de  $A^T A$  y de  $AA^T$ .

La unicidad de la DVS será un hecho si todos los valores singulares son distintos entre sí; en el caso de que existan espates las porciones correspondientes en  $U$  y  $V$  no serán únicas. Debe señalarse que  $D_a$  sí que será única. En cuanto al

problema de los empates hay que recordar sin embargo que los subespacios generados por vectores propios asociados a un mismo valor singular son iguales.

Dado (a.2), la matriz A puede aproximarse por matrices de rango menor a K, es decir,

$$A = \sum_{k=1}^K a_k u_k x_k^* \quad 1 \leq k \leq K \quad (a.3)$$

Esta forma de aproximar es óptima bajo el criterio

$$\begin{aligned} \min_X \|A-X\| &= \min_X \sum_i \sum_j (a_{ij} - x_{ij})^2 \\ &= \min \text{traza}((A-X)(A-X)^*) \end{aligned} \quad (a.4)$$

Por último, en el caso de que se desea realizar una DVB generalizada, es decir, expresar A como

$$A = N D M^* \quad (a.5)$$

sujeto a  $N^* \mathcal{T} N = M^* \mathcal{X} M = I$ , con  $\mathcal{T}$  y  $\mathcal{X}$  no necesariamente matrices identidad, resólvase la DVB simple para  $A$  y con esto

$$N = \mathcal{T}^{1/2} U \quad M = \mathcal{X}^{1/2} V. \quad (a.6)$$

En este caso la aproximación de A por matrices de rango menor será óptima bajo

$$\begin{aligned} \min_X \text{traza}(\mathcal{T}(A-X)\mathcal{X}(A-X)^*) \\ = \min_X \sum_i w_i (a_i - x_i)^2 \mathcal{X} (a_i - x_i) \end{aligned} \quad (a.7)$$

si  $\mathcal{T} = \text{diag}(w_i)$ .

**Definición 1.**

1. Si  $(Y_n)$  es tal que  $P(|Y_n| < k) = 1$  para toda  $n$ , entonces  $(Y_n)$  es uniformemente integrable, es decir,

$$\lim_{c \rightarrow \infty} \sup_n E(|Y_n| I(|Y_n| > c)) = 0$$

donde

Para toda  $c > k$   $E(|Y_n| > c) = 0$  para toda  $n$  por lo tanto,

$$\sup_n E(|Y_n| > c) = 0 \text{ para toda } c > k.$$

De esto se tiene el resultado.

2.  $(Y_n)$  y  $(X_n)$  son asintóticamente equivalentes y acotadas.  $(Y_n)$  es uniformemente integrable y  $X_n \xrightarrow{d} F$  entonces,

$$Y_n \xrightarrow{d} F \text{ y } E(X_n) \rightarrow \int x dF.$$

donde

$Y_n - X_n \rightarrow 0$  c.p. i; por lo tanto (Slutsky)  $Y_n \xrightarrow{d} F$  en distribución.  $Y_n$  es uniformemente integrable implica que

$$E(Y_n) \rightarrow \int x dF$$

por teorema A sec 1.4 de Serfling (1980).

$$E(X_n) = E(X_n - Y_n) + E(Y_n)$$

Por el teorema de convergencia dominada (o acotada)

$$E(X_n - Y_n) \rightarrow 0$$

y por lo tanto

$$E\left[\frac{1}{n}\right] \longrightarrow \int x dF.$$