



42  
207  
UNIVERSIDAD NACIONAL AUTÓNOMA  
DE MÉXICO

FACULTAD DE CIENCIAS

ESTUDIO COMPARATIVO ENTRE LOS ANALISIS DE  
CORRESPONDENCIAS SIMETRICO Y  
NO - SIMETRICO

TESIS PROFESIONAL  
QUE PARA OBTENER EL TITULO DE  
A C T U A R I O  
P R E S E N T A I  
RICARDO MEDINA ALVAREZ

DIRECTOR: DR. RUBEN HERNANDEZ CID

MEXICO, D. F.

TESIS CON  
FALLA DE ORIGEN

1991



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.

# CONTENIDO

Capítulo 1 Introducción	1
1.1 Relaciones entre variables cualitativas	1
1.2 Acerca de los temas tratados	2
Capítulo 2 Tablas de contingencia	4
2.1 Matriz asociada a una tabla de contingencia	4
2.2 Matriz de correspondencias y vectores suma	5
2.3 Matrices perfil de renglones y columnas	6
2.4 Ejemplo ilustrativo	7
Capítulo 3 Identificación de subespacios óptimos	11
3.1 Planteamiento general para identificar subespacios óptimos	12
3.2 Técnica para la solución de identificar subespacios óptimos	16

Capítulo 4 Análisis de Correspondencias Simétrico	21
4.1 Análisis en $R^J$	21
4.2 Análisis en $R^I$	23
4.3 Simetría	25
4.4 Resultados y Comentarios	27
Capítulo 5 Análisis de Correspondencias No-simétrico	33
5.1 Análisis en $R^I$	34
5.2 Análisis en $R^J$	35
5.3 Índice $\tau$ de Goodman-Kruskal	37
Capítulo 6 Aplicación	39
6.1 Ejemplo de datos médicos	40
6.2 Aplicación simétrica	41
6.3 Aplicación no-simétrica	43
6.4 Comentarios de resultados	45
Comentarios generales	47

Apéndice A Ejemplos de obtención de subespacios y asignación de pesos	49
A.1 Ejemplos en $\mathbb{R}^2$ y $\mathbb{R}^3$	49
A.2 Asignación de pesos a ejes y puntos	56
Apéndice B Propuesta de un programa en computadora	60
B.1 Instalación	61
B.1 ANACORR	61
Referencias	79

# 1. INTRODUCCION

## 1.1 Relaciones entre variables cualitativas.

En los últimos años, el estudio de relaciones entre variables cualitativas se ha realizado, en una importante parte, mediante el Análisis de Correspondencias.

El caso principal de una matriz de datos conveniente para el Análisis de Correspondencias, es una tabla de contingencia, la cual expresa la asociación de las observaciones entre dos variables cualitativas.

H. O. Hartley (1935) publicó un artículo en el cual da una formulación algebraica de la asociación entre los renglones y las columnas de una tabla de contingencia, a este artículo se le atribuye el origen matemático del Análisis de Correspondencias. Después, Louis Guttman (1941) trató el caso general de más de dos variables cualitativas, esto da como resultado lo que ahora se conoce como Análisis de Correspondencias Múltiple. El francés Jean Paul Benzécri (1977),

trabajando en un contexto lingüístico, fue el creador de la forma geométrica del Análisis de Correspondencias, donde el término francés *correspondence* fue usado para denotar "sistemas de asociaciones" entre elementos de dos conjuntos de datos ( rengiones y columnas).

Todos los métodos desarrollados para este tipo de análisis están basados en la hipótesis implícita de considerar simétricas las relaciones entre las variables. Esta suposición aunada al uso e interpretación de la ji-cuadrada de Pearson, medida de asociación (simétrica) en cual se apoya el Análisis de Correspondencias, podrían ser inconvenientes en casos en que la relación entre las variables no sea simétrica.

El análisis de relaciones no-simétricas para variables cuantitativas ya ha sido tratado en diferentes contextos como en Regresión Multivariada y un primer acercamiento al tratamiento no-simétrico de variables cualitativas fue propuesto por los italianos Luigi D'Ambra y Natale Lauro (1982).

## 1.2 Acerca de los temas tratados.

Este trabajo presenta métodos para realizar los Análisis de Correspondencias Simétrico (clásico) y No-simétrico, usando la técnica numérica de Descomposición en Valores Singulares.

Como en este trabajo, se desarrolla el Análisis de Correspondencias basado en tablas de contingencia, en el Capítulo 2 se tratan algunos conceptos relacionados con éstas, los cuales serán utilizados en capítulos posteriores, así

como un ejemplo de una tabla de contingencia la cual ilustra cada uno de los conceptos.

Los objetivos del Capítulo 3 son: primero, tratar el problema de identificar subespacios óptimos, el cual se considera como el problema principal del Análisis de Correspondencias; segundo, mostrar la teoría de la técnica de Descomposición en Valores Singulares (DVS) y tercero, mostrar cómo la técnica de DVS resuelve el problema de identificar subespacios óptimos.

Estos dos primeros capítulos son considerados de apoyo, pero básicos, para los capítulos principales, 4 y 5; donde se desarrolla la teoría relacionada con cada análisis.

Un ejemplo que ilustra el uso de ambos análisis es tratado en el Capítulo 6.

En el Apéndice A se muestran algunos ejemplos de obtención de subespacios de menor dimensión (en dimensiones 2 y 3) y además se hace notar cuál es la importancia de asignar pesos a puntos y ejes. Este capítulo podría servir como introducción para el Capítulo 3 y es muy recomendable para quienes empiezan a trabajar con Análisis de Correspondencias, aunque se requieren conocimientos de Álgebra Lineal y Estadística.

Dado que las técnicas que se usan en el Análisis de Correspondencias son un tanto laboriosas, el objetivo principal del Apéndice B, es el de explicar el funcionamiento del programa de cómputo ANACORR, el cual es una propuesta de un programa para poder obtener los resultados de ambos análisis.



## 2. TABLAS DE CONTINGENCIA

Las tablas de contingencia expresan la asociación de las observaciones entre variables cualitativas. Los siguientes conceptos son necesarios para desarrollar la teoría de los Análisis de Correspondencias. Las tablas de contingencias que se usarán en este trabajo, son las de dos variables cualitativas.

### 2.1 Matriz asociada a una tabla de contingencia

Sean  $I$  y  $J$ , dos variables cualitativas con  $I$  y  $J$  categorías respectivamente. Una tabla de contingencia asociada a estas variables se puede representar por medio de la matriz  $N_{I \times J}$ , donde la celda  $n_{ij}$  representa el número de observaciones clasificadas en la categoría  $i$  y en la categoría  $j$ , de las variables  $I$  y  $J$ .

Esta matriz debe cumplir

la celda  $n_{ij} \geq 0$ ; para  $i = 1 \dots I$  y  $j = 1 \dots J$  y

$$\sum_{j=1}^J n_{ij} > 0 \text{ para } i = 1 \dots I \quad \text{y} \quad \sum_{i=1}^I n_{ij} > 0 \text{ para } j = 1 \dots J;$$

## 2.2 Matriz de correspondencias y vectores suma

Sea  $n = \sum_{j=1}^J \sum_{i=1}^I n_{ij}$ , la suma de todas las celdas de  $N$ .

La *matriz de correspondencias*,  $P_{i \cdot}$ , se define como  $(1/n)N$ , por lo que el término general de  $P$  es  $p_{ij} = n_{ij}/n$ . Es claro que la suma de elementos de  $P$  es uno y se puede considerar como la matriz de densidad de probabilidad conjunta sobre las celdas de la matriz  $N$ .

El vector de las sumas de celdas por renglón de  $P$ , se denota como  $r = [p_{1 \cdot}, \dots, p_{h \cdot}, \dots, p_{I \cdot}]^T$ , y la  $h$ -ésima componente del vector es

$$p_{h \cdot} = \sum_{j=1}^J p_{ij} = \sum_{j=1}^J n_{ij}/n = (n_{h \cdot})/n$$

donde  $n_{h \cdot}$  es la suma de elementos del renglón de  $h$  de la matriz  $N$ .

La matriz diagonal, asociada con este vector se denota como

$$D_{r \cdot} = \text{diag}(p_{r \cdot})$$

Análogamente, el vector de las sumas de celdas por columna de  $P$  se denota como  $c = [p_{\cdot 1}, \dots, p_{\cdot h}, \dots, p_{\cdot J}]^T$  y la  $h$ -ésima componente del vector es

$$p_{\cdot h} = \sum_{i=1}^I p_{ij} = \sum_{i=1}^I n_{ij}/n = (n_{\cdot h})/n$$

donde  $n_{\cdot h}$  es la suma de elementos de la columna  $h$  de la matriz  $N$ .

La matriz diagonal, asociada con este vector se denota como

$$D_{c_{.j}} = \text{diag}(p_{.j}).$$

Los vectores  $r$  y  $c$  pueden ser considerados como densidades marginales sobre las celdas de la matriz  $N$ .

### 2.3 Matrices perfil de renglones y columnas

La matriz *perfil* de renglones de  $P$  se define como

$R_{r_{i.}} = D r^{-1} P$ , y representa la frecuencia relativa con respecto a la variable  $I$  en cada celda.

Notar que el término general,  $r_{ij}$ , se escribe como

$$r_{ij} = p_{ij} / p_{i.} = (n_{ij} / n) / (n_{i.} / n) = n_{ij} / n_{i.};$$

esto es, los renglones de  $P$  (o también de  $N$ ) divididos por su respectiva suma.

Análogamente, la matriz *perfil* de columnas de  $C$  se define como

$C_{c_{.j}} = P D c^{-1}$ , y representa a cada celda como la frecuencia relativa con respecto a la variable  $J$ , el término general;  $c_{ij} = p_{ij} / p_{.j} = n_{ij} / n_{.j}$ , esto es, los columnas de  $P$  (o de  $N$ ) divididos por su respectiva suma.

Estas matrices juegan un papel muy importante en el Análisis de Correspondencias y en el siguiente ejemplo se tratará de mostrar su importancia.

## 2.4 Ejemplo Ilustrativo

Una muestra de 193 empleados de una empresa son clasificados con respecto a: su hábito de fumar en las categorías NADA, POCO (1 a 5 cigarros al día), REGULAR (6 a 18 cigarros) y DEMASIADO (más de 18 cigarros); y el puesto que ocupan dentro de la empresa con las categorías GERENTE, SUBGERENTE, JEFE, ANALISTA y SECRETARIA. La tabla de contingencia asociada a la clasificación, es la siguiente

	NADA	POCO	REGULAR	DEMASIADO	Total
GERENTE	4	2	3	2	11
SUBGERENTE	4	3	7	4	18
JEFE	25	10	12	4	51
ANALISTA	18	24	33	13	88
SECRETARIA	10	6	7	2	25
Total	61	45	62	25	193

Si las variables cualitativas  $I$  y  $J$  son el puesto y el hábito de fumar, respectivamente, entonces los números de categorías asociado con estas variables son  $I=5$  y  $J=4$ .

La matriz de Correspondencias  $P_{5 \times 4}$  está dada por

$$P = \begin{bmatrix} 0.021 & 0.010 & 0.015 & 0.010 \\ 0.021 & 0.015 & 0.036 & 0.021 \\ 0.130 & 0.052 & 0.062 & 0.021 \\ 0.093 & 0.124 & 0.171 & 0.067 \\ 0.052 & 0.036 & 0.036 & 0.010 \end{bmatrix}$$

Como a ésta se le puede considerar como la matriz de densidad de probabilidad conjunta estimada sobre las celdas de la tabla de contingencia, se pueden dar varias relaciones como

$$\text{prob. (puesto=ANALISTA y hábito=REGULAR)} = 0.171$$

Los vectores de sumas de celdas por renglón y por columna están dados por

$$r = [0.057, 0.093, 0.264, 0.456, 0.130]^T$$

$$c = [0.316, 0.233, 0.321, 0.130]^T$$

Como ya se mencionó,  $r$  y  $c$  pueden ser considerados como densidades marginales estimadas sobre las celdas de la matriz  $N$  y se pueden dar los siguientes ejemplos

$$\text{prob. (puesto = SECRETARIA)} = 0.130 \text{ y}$$

$$\text{prob. (hábito = NADA)} = 0.316$$

La matriz perfil de renglones  $R_{5 \times 4}$  está dada por

$$R = \begin{bmatrix} 0.364 & 0.182 & 0.273 & 0.182 \\ 0.222 & 0.167 & 0.389 & 0.222 \\ 0.490 & 0.196 & 0.235 & 0.078 \\ 0.205 & 0.273 & 0.375 & 0.148 \\ 0.400 & 0.240 & 0.280 & 0.080 \end{bmatrix}$$

Lo que se logra al construir esta matriz es representar al total de observaciones en cada celda como frecuencias relativas con respecto a la variable *puesto*, es decir, esta matriz expresa las probabilidades condicionales de la forma  $P(\text{hábito dado puesto})$ . Por ejemplo, la frecuencia con respecto al total de observaciones en la celda GERENTE-NADA es 0.021 al igual que en la celda JEFE-DEMASIADO, pero si se observan las mismas celdas en la matriz perfil R, los valores son 0.364 para la primera y 0.078 para la segunda, esto significa que mientras que en la categoría GERENTE, la categoría NADA juega un papel importante (de 11 gerentes, 4 no fuman); la categoría DEMASIADO no es tan relevante para la de JEFE (de 51 jefes, solo 4 fuman demasiado).

La matriz perfil de columnas  $C_{3 \times 4}$  se interpreta de forma análoga y está dada por

$$C = \begin{bmatrix} 0.066 & 0.044 & 0.048 & 0.080 \\ 0.066 & 0.067 & 0.113 & 0.160 \\ 0.410 & 0.222 & 0.194 & 0.160 \\ 0.295 & 0.533 & 0.532 & 0.520 \\ 0.164 & 0.133 & 0.113 & 0.080 \end{bmatrix}$$

Esta matriz representa las frecuencias de observaciones en cada celda como frecuencias relativas con respecto a la variable *hábito de fumar*, esto es, son probabilidades condicionales de la forma  $P(\text{puesto} \text{ dado } \text{hábito})$ . Por ejemplo, las celdas GERENTE-NADA y JEFE-DEMASIADO tienen una frecuencia con respecto al total de observaciones, como ya se ha visto, de 0.021, pero en la matriz perfil de columnas, estas celdas tienen valores de 0.066 y 0.161 respectivamente, lo que significa, que aunque el papel que juegan las categorías GERENTE y JEFE, en las categorías NADA y DEMASIADO respectivamente, no es tan fuerte, es un poco más importante el papel que juega la categoría JEFE en la categoría DEMASIADO.

Gracias a esta interpretación el Análisis de Correspondencias considera como puntos a las categorías de la variable  $I$  ( $J$ ), en un espacio de dimensión  $J$  ( $I$ ), donde los  $J$  ( $I$ ) ejes están dados por las  $J$  ( $I$ ) categorías de variable  $J$  ( $I$ ) y la componente  $j$  ( $i$ ) del punto  $i$  ( $j$ ) es  $r_{ij}$  ( $c_{ij}$ ); y trata de representarlos en un subespacio "conveniente" de menor dimensión, con el objeto de apreciar las *correspondencias* entre las categorías de ambas variables.

### **3. IDENTIFICACION DE SUBESPACIOS OPTIMOS**

Como se ha mencionado anteriormente, el objetivo principal del Análisis de Correspondencias es representar "lo mejor posible" las categorías de las variables cualitativas como puntos en un subespacio "conveniente".

En este capítulo se explicará, en primer término, qué se quiere decir con "lo mejor posible" y "conveniente", y después se propondrá una técnica para resolver el problema.



### 3.1 Planteamiento general para identificar subespacios óptimos

Supóngase que se tiene una nube de puntos  $U = \{u_1 \dots u_T\}$  en un espacio  $F$  de dimensión  $m$  y se requiere representar los puntos de esta nube en un subespacio  $S$  de dimensión  $k' < m$ , el objetivo es encontrar el subespacio  $S^*$  el cual esté "más cerca" de la nube de puntos.

¿Cómo definir distancia entre una nube de puntos a un subespacio dado?

Dada la distancia entre puntos, entonces, se puede definir la distancia entre un punto  $u$  y un subespacio  $S$  como la distancia más corta entre el punto  $u$  y todos los puntos contenidos en el subespacio, esto es

$$d(u, S) = \min ( d(u, y) ) \text{ para todo } y \in S.$$

Ahora, la distancia de la nube de puntos  $U = \{u_1 \dots u_T\}$  al subespacio  $S$ , se define como la suma de distancias al cuadrado (por conveniencias geométricas, como en Regresión y Análisis de Varianza), de los puntos  $u_i$  al subespacio  $S$ .

Si  $\Phi$  es la función que describe la distancia entre la nube  $U$  y el subespacio  $S$ ,  $\Phi$  está dada con la siguiente relación

$$\Phi(S; u_1 \dots u_T) = \sum_{i=1}^T w_i (u_i - \hat{a}_i)^T D_p (u_i - \hat{a}_i) = \sum_{i=1}^T w_i d(u_i, S) D_p$$

donde;

$\hat{a}_1, \dots, \hat{a}_T \in S$  son los puntos más cercanos a los puntos  $u_1 \dots u_T \in P$  respectivamente,

$w_i$  es el peso (masa) asignado al punto  $u_i$ .

$D_p$  es la matriz diagonal asociada a los pesos de las dimensiones (métrica).

Esta función depende del subespacio  $S$  y el objetivo es encontrar el subespacio óptimo  $S^*$  que minimice la función  $\Phi$ .

Si se elige al subespacio  $S$  de dimensión cero, éste sólo tendría al punto  $\bar{u}$  y la función objetivo es

$$\Phi(\bar{u}; u_1 \dots u_I) = \sum_{i=1}^I w_i (u_i - \bar{u})^T D_p (u_i - \bar{u})$$

El *centroide* de la nube  $I$ , es el punto que minimiza la función  $\Phi$ .

**Demostración.**

El centroide está dado por

$$\bar{u} = \frac{1}{I} \sum_{i=1}^I w_i u_i$$

su  $j$ -ésima coordenada es;

$$\bar{u}_j = \frac{1}{I} \sum_{i=1}^I w_i u_{ij}$$

y se tiene que demostrar que ésta, es la  $j$ -ésima coordenada del vector  $\bar{u}$  la cual minimiza la función  $\Phi$ .

Por otra parte, la función  $\Phi$  se puede reescribir como

$$\Phi(\bar{u}) = \sum_{i=1}^I w_i (q_1 (u_{i1} - s_1)^2 + \dots + q_j (u_{ij} - s_j)^2 + \dots + q_J (u_{iJ} - s_J)^2)$$

Entonces, la derivada de  $\Phi$  con respecto con  $s_j$  es

$$D\Phi(s_j) = -2 \sum_{i=1}^I w_i q_j (u_{ij} - s_j) = -2q_j \sum_{i=1}^I w_i (u_{ij} - s_j)$$

Ahora, igualando la derivada a cero se tiene que resolver la siguiente ecuación

$$\sum_{i=1}^I w_i (u_{ij} - s_j) = 0 \Rightarrow \sum_{i=1}^I w_i u_{ij} = \sum_{i=1}^I s_j w_i = s_j \sum_{i=1}^I w_i \Rightarrow$$

$$s_j = \frac{\sum_{i=1}^I w_i u_{ij}}{\sum_{i=1}^I w_i}, \text{ el cual es punto crítico de } \#.$$

Entonces, sólo resta verificar que este punto minimiza  $\#$  y por medio de la segunda derivada se tiene que

$$D^2 \#(s_j) = -2q_j \sum_{i=1}^I w_i (-1) = 2q_j > 0.$$

Por lo que el centroide minimiza  $\#$ . ■

Con base en este resultado se deduce que el subespacio óptimo  $S^*$ , debe contener al *centroide*.

**Demostración.** Sea  $S'$  un subespacio que no contiene al centroide  $\bar{u}$ , se demostrará que  $S'$  no es óptimo.

Sea  $u_i' \in S'$  es el punto más cercano a  $u_i$ , entonces la distancia de la nube de puntos  $I = \{u_1, \dots, u_I\}$ , al subespacio  $S'$  está dada por la función  $\#$

$$\#(S'; u_1, \dots, u_I) = \sum_{i=1}^I w_i (u_i - u_i')^T D_p (u_i - u_i').$$

Sin pérdida de generalidad se puede asumir que  $\sum_{i=1}^I w_i = 1$ . Sean  $u' \in S'$  el punto que está más cercano a  $\bar{u}$ ,  $t = \bar{u} - u'$  y, finalmente,  $u_i = u_i' + t$ . La función  $\#$  puede escribirse como

$$\begin{aligned}
 \#(S'; u_1 \dots u_I) &= \sum_{i=1}^I w_i (u_i - \hat{u}_i + \hat{u}_i - u_i')^T D_p (u_i - \hat{u}_i + \hat{u}_i - u_i') = \\
 &= \sum_{i=1}^I w_i (u_i - \hat{u}_i)^T D_p (u_i - \hat{u}_i) + \\
 &+ \sum_{i=1}^I w_i (\hat{u}_i - u_i')^T D_p (\hat{u}_i - u_i') + \\
 &+ 2 \sum_{i=1}^I w_i (u_i - \hat{u}_i)^T D_p (\hat{u}_i - u_i').
 \end{aligned}$$

El segundo sumando de esta expresión es, simplemente,  $t^T D_p t$  (cuadrado de la distancia entre  $\bar{u}$  y  $u'$ ); el tercero es 0, ya que  $D_p (\hat{u}_i - u_i') = D_p t$  y

$$\sum_{i=1}^I w_i (u_i - \hat{u}_i) = \sum_{i=1}^I w_i u_i - \sum_{i=1}^I w_i \hat{u}_i = \bar{u} - \sum_{i=1}^I w_i (u_i' + t) = \bar{u} - (u' + t) = 0$$

Entonces la función  $\#$  se puede escribir como

$$\#(S'; u_1 \dots u_I) = \sum_{i=1}^I w_i (u_i - \hat{u}_i)^T D_p (u_i - \hat{u}_i) + t^T D_p t.$$

Entonces el subespacio definido por todos los puntos de  $S'$  mas el vector  $t$  (los puntos  $\hat{u}_i$ ) es "mejor" subespacio que  $S'$ .

Por este resultado, los puntos  $\hat{u}_i$  se pueden escribir como

$$\hat{u}_i = \bar{u} + \sum_{k=1}^k f_{ik} m_k$$

donde  $m_1, \dots, m_k$  son los vectores base del subespacio  $S^*$ , y  $f_{ik}$  son las coordenadas con respecto a esta base. La función  $\Phi$  se puede escribir, entonces, como

$$\begin{aligned}\Phi(S; u_1, \dots, u_I) &= \sum_{i=1}^I w_i (u_i - \bar{u} - \sum_{k=1}^k f_{ik} m_k)^T D p (u_i - \bar{u} - \sum_{k=1}^k f_{ik} m_k) \\ &= \sum_{i=1}^I w_i \left[ (u_i - \bar{u}) - \sum_{k=1}^k f_{ik} m_k \right]^T D p \left[ (u_i - \bar{u}) - \sum_{k=1}^k f_{ik} m_k \right]\end{aligned}$$

Las variables de esta función objetivo son  $k^*$  ejes principales,  $m_1, \dots, m_{k^*}$ , implicando un total de  $Jk^*$  variables escalares.

### 3.2 Técnica para la solución de identificar subespacios óptimos

En esta sección se presentan algunos resultados relacionados a la técnica numérica de Descomposición en Valores Singulares y cómo puede ésta contribuir al Análisis de Correspondencias. Para una consulta más detallada se recomienda el texto de Micheal J. Greenacre (1984).

Teorema de DVS Ordinario. Cualquier matriz  $A_{I \times J}$  real, puede ser expresada como

$$A_{I \times J} = U_{I \times K} D_{K \times K} V_{K \times J}^T$$

donde;

$D_{K \times K}$  es una matriz diagonal de números positivos  $\alpha_1, \dots, \alpha_K$ ,

$K$  es el rango de  $A$  ( $K \leq \min(I, J)$ ),

$U^T U = V^T V = I$  (los vectores columna de  $U$  y  $V$  son ortonormales en el espacio euclidiano).

Una forma equivalente de expresar A es

$$A = \sum_{n=1}^K \alpha_n (u^n) (v^n)^T$$

donde  $u^1 \dots u^K$  y  $v^1 \dots v^K$  son las columnas de U y V.

Los valores  $\alpha_1 \dots \alpha_K$  son llamados valores singulares de A, los vectores  $u^1 \dots u^K$  son llamados vectores singulares izquierdos de A y los vectores  $v^1 \dots v^K$  son llamados vectores singulares derechos de A.

Los vectores singulares izquierdos,  $u^1 \dots u^K$ , forman una base ortonormal para las columnas de A y las coordenadas, con respecto a esta base, están dadas por los renglones de  $G = VD\alpha$ . Esto es, si la  $i$ -ésima columna de A es  $a^i = (a_{i1} \dots a_{iJ})^T$  y el  $i$ -ésimo renglón de G es  $g_i = (g_{i1} \dots g_{iK})$ , entonces

$$a^i = g_{i1}u^1 + g_{i2}u^2 + \dots + g_{iK}u^K.$$

De forma análoga, los vectores singulares derechos,  $v^1 \dots v^K$ , forman una base ortonormal para los renglones de A y las coordenadas, con respecto a esta base, están dadas por los renglones de  $F = UD\alpha$ . Esto es;

$$a_i = f_{i1}v^1 + f_{i2}v^2 + \dots + f_{iK}v^K.$$

Un caso especial de DVS es la eigendescomposición de una matriz simétrica  $B_{J \times J}$  de rango  $K \leq J$ , como sigue

$$B_{J \times J} = V_{J \times K} D_{K \times K} V_{K \times J}^T$$

En este caso los vectores singulares derechos e izquierdos son los mismos y conocidos como eigenvectores de B mientras que los valores singulares son conocidos como eigenvalores. La prueba de la existencia de la descomposición de la matriz A se facilita gracias a la existencia de la eigendescomposición de la

matriz cuadrada y simétrica B, tal que  $B = ATA = (UD\alpha V^T)^T (UD\alpha V^T) = VD\alpha(U^T U)D\alpha V^T = VD\alpha D\alpha V^T = VD\alpha^2 V^T$ .

Antes de mencionar cual es la importancia de DVS para este propósito, se generalizará el resultado con lo que se conoce como DVS Generalizado.

Teorema de Descomposición en Valores Singulares Generalizado.

Si  $\Omega_{kl}$  y  $\phi_{kl}$  son matrices simétricas definidas positivas, entonces cualquier matriz  $A_{kl}$  de rango K puede ser expresada como

$$A_{kl} = N_{k \times K} D\alpha_{K \times K} M^T_{K \times l}$$

donde los vectores columna de N y M son ortonormales con respecto a  $\Omega$  y  $\phi$ , es decir

$$N^T \Omega N = M^T \phi M = I$$

La importancia de DVS para este propósito es que si se eligen los valores singulares de la forma  $\alpha_1 \geq \alpha_2 \geq \dots \geq \alpha_K$  con sus respectivos vectores singulares y si se definen  $\Omega = D_w$  (matriz de pesos asignadas a cada punto) y  $\phi = D_p$  (matriz de pesos asignados a las dimensiones) entonces la matriz

$$A_{\{k^*\}} = N_{\{k^*\}} D\alpha_{\{k^*\}} M^T_{\{k^*\}} = \sum_{s=1}^{k^*} \alpha_s n_s m_s^T$$

minimiza la siguiente función

$$\|A - X\|_{D_p, D_w}^2 = \sum_{i=1}^I w_i (a_i - x_i)^T D_p (a_i - x_i)$$

entre todas las matrices X de rango menor o igual a K.

Recordando la función a minimizar, para encontrar el subespacio  $S^*$

$$e(S; u_1 \dots u_r) = \sum_{i=1}^r w_i [(u_i - \bar{u}) - \sum_{k=1}^k f_{ik} m_k]^T D q [(u_i - \bar{u}) - \sum_{k=1}^k f_{ik} m_k]$$

Entonces, si cada punto de la nube  $U = \{u_1 \dots u_r\}$  es asociado con un renglón de una matriz  $U_{kj}$  y, si  $A$  es definida como la matriz centrada de renglones de  $U$ , esto es, el  $i$ -ésimo renglón de  $A$  es  $u_i - \bar{u}$ , entonces DVS da la solución requerida, ya que los vectores  $m^1 \dots m^k$  (columnas de  $M$ ) definen una base ortonormal (ejes principales) para el subespacio óptimo y las coordenadas de los vectores  $u_i - \bar{u}$  con respecto a esta base están dadas en los renglones de la matriz  $F$ , donde  $F_{[k]} = N_{[k]}^* D \alpha_{[k]}^*$ .

El total de variación de la matriz  $A$  es cuantificado por la norma cuadrada  $\|A\|_{Dp, Dw}^2$ , la cual se puede expresar como el cuadrado de los valores singulares, esto da una idea de qué tan bien es representada la matriz sobre los ejes principales y se tiene que

$$\|A\|_{Dp, Dw}^2 = \sum_{i=1}^I w_i a_i^T D p a_i = \sum_{i=1}^K \alpha_i^2 \quad (\text{total de inercia}),$$

similarmente la variación de la aproximación  $A_{[k]}^*$  es

$$\|A_{[k]}^*\|_{Dp, Dw}^2 = \sum_{i=1}^k \alpha_i^2,$$

la variación no explicada es

$$\|A - A_{[k]}^*\|_{Dp, Dw}^2 = \sum_{i=k+1}^K \alpha_i^2,$$

el cociente  $r_k = \|A_{[k]}^*\|_{Dp, Dw}^2 / \|A\|_{Dp, Dw}^2$  es usado para cuantificar la calidad de la matriz  $A_{[k]}^*$  en la dimensión  $k$ .



Para obtener DVS generalizado es recomendable partir del DVS ordinario como sigue

Sea la matriz  $A$  a la cual se desea aplicar DVS generalizado, entonces

$A = N D\alpha M^T$  donde  $N^T D_w N = M^T D_q M = I$  y se desea conocer  $N$ ,  $D\alpha$ ,

$M$ .

Se definen una matriz auxiliar  $B$ , tal que

$$B = D_w^{\frac{1}{2}} A D_q^{\frac{1}{2}}$$

Ahora se aplica DVS ordinario a  $B$ , esto es

$B = U D\phi V^T$  donde  $U^T U = V^T V = I$  y se encuentran  $U$ ,  $D\phi$  y  $V$ .

Entonces,  $N = D_w^{-\frac{1}{2}} U$ ,  $M = D_q^{-\frac{1}{2}} V$  y  $D\phi = D\alpha$ , ya que;

$$B = D_w^{\frac{1}{2}} A D_q^{\frac{1}{2}} = D_w^{\frac{1}{2}} (N D\alpha M^T) D_q^{\frac{1}{2}} = (D_w^{\frac{1}{2}} N) D\alpha (M^T D_q^{\frac{1}{2}})$$

Se puede poner  $U = D_w^{\frac{1}{2}} N$  y  $V = D_q^{\frac{1}{2}} M$ , ya que

$$U^T U = (N^T D_w^{\frac{1}{2}}) (D_w^{\frac{1}{2}} N) = N^T D_w N = I \text{ y}$$

$$V^T V = (M^T D_q^{\frac{1}{2}}) (D_q^{\frac{1}{2}} M) = M^T D_q M = I, \text{ cumple con las restricciones}$$

de ortonormalización, entonces

$$N = D_w^{-\frac{1}{2}} U, M = D_q^{-\frac{1}{2}} V \text{ y } D\phi = D\alpha.$$

## 4. ANALISIS DE CORRESPONDENCIAS SIMETRICO

El Análisis de Correspondencias puede ser presentado como una técnica, para representar los renglones y columnas de una matriz de datos, como puntos en un espacio de menor dimensión, el cual tiene ciertas características.

### 4.1 Análisis en $R^J$

Considérese la nube de puntos  $N_i$  (I puntos en un espacio de dimensión  $J$ ) cuyas coordenadas, con respecto a la base estándar, están dadas por los renglones de la matriz perfil de renglones  $R$ . Es decir, el  $i$ -ésimo punto de la nube es

$$r_i = [r_{i1}, r_{i2}, \dots, r_{iJ}]^T =$$

$\{n_{i1}/n_i, \dots, n_{iJ}/n_i\}$ , esto es, el renglón  $i$  de la matriz asociada a la tabla de contingencia, en proporción con el total del renglón.

Ahora, antes de aplicar DVS para encontrar el subespacio de menor dimensión más cercano a  $N_i$ , se deben definir los pesos asociados a cada punto y los pesos asociados a las dimensiones, es decir, la matriz de pesos y la métrica.

Como matriz de pesos a puntos se elige la matriz  $D_r$ , de suma de elementos por renglón.

Como métrica se elige la matriz  $D_c^{-1}$ , inversa de la matriz de suma de elementos por columna (métrica ji-cuadrada).

Notar que, considerando la tabla de contingencia  $N$ , la proporción del renglón  $i$  (asociado con el punto  $i$ ) con respecto al total, está dado por

$$(n_{i1} + n_{i2} + \dots + n_{iU})/n = n_{i1}/n + n_{i2}/n + \dots + n_{iU}/n = p_{i1} + p_{i2} + \dots + p_{iU} = p_{i\cdot}$$

que es, justamente, el elemento  $i$  de la diagonal de la matriz  $D_r$  (matriz asignada a los pesos de puntos).

Ahora, el peso de la columna  $j$  (asociado con el eje  $j$ ) está dado por

$$(n_{1j} + n_{2j} + \dots + n_{ij})/n = n_{1j}/n + n_{2j}/n + \dots + n_{ij}/n = p_{1j} + p_{2j} + \dots + p_{ij} = p_{\cdot j}$$

que es el elemento  $j$  de la diagonal de la matriz  $D_c$  (inversa de la matriz asociada a la métrica).

Una vez definida la matriz de pesos a puntos se puede ver que el centroide de la nube  $N_i$  es el vector  $c$ , de suma de columnas de  $P$ , de la siguiente forma

$$\text{centroide de } N_i = \sum_{i=1}^I (p_{i\cdot}) r_i = \sum_{i=1}^I (p_{i\cdot}) [r_{i1}, \dots, r_{iJ}]^T =$$

$$\sum_{i=1}^I (p_{i\cdot}) [p_{i1}/p_{i\cdot}, \dots, p_{iJ}/p_{i\cdot}]^T = \sum_{i=1}^I [p_{i1}, \dots, p_{iJ}]^T =$$

$$[\sum_{i=1}^I p_{i1}, \dots, \sum_{i=1}^I p_{iJ}]^T = [p_{\cdot 1}, \dots, p_{\cdot J}]^T = c$$

Entonces la matriz centrada de  $R$ ,  $R^*$ , es aquella cuyo  $i$ -ésimo renglón es  $r_i - c$  ( $R^* = R - 1c^T$ , donde  $1$  es un vector de unos de dimensión  $I$ ).

El subespacio de dimensión  $k^*$ , el cual está más cerca a la nube de puntos  $N_i$  está definido por los primeros  $k^*$  vectores singulares derechos de  $R^*$ , es decir, los vectores singulares derechos definen los ejes principales del nube  $N_i$ .

Para conocer los ejes principales y las coordenadas de la nube con respecto a estos ejes; se tiene que obtener DVS generalizado de la matriz  $R^*$ , donde los vectores singulares izquierdos y derechos están ortonormalizados con respecto a  $D_r$  y  $D_c^{-1}$  respectivamente.

Esto es,  $R^*$  se puede descomponer como (DVS generalizado)

$$R^* = L D_\mu M^T \quad \text{donde} \quad L^T D_r L = M^T D_c^{-1} M = I$$

entonces, las columnas de  $M$  (vectores singulares derechos) definen los ejes principales y los renglones de  $L D_\mu = R^* D_c^{-1} M$  definen las coordenadas de la nube  $N_i$  con respecto a estos ejes.

Los ejes principales del mejor subespacio de dimensión  $k^*$  son las primeras  $k^*$  columnas de la matriz  $M$ .

#### 4.2 Análisis en $R^I$ (Análisis dual)

En la sección anterior, se ha tratado el problema de encontrar el mejor subespacio de menor dimensión para la nube de puntos  $N_i$ , asociada a la matriz perfil de renglones  $R$  (a la nube  $N_i$  suele llamársele nube de renglones). Con un tratamiento similar y *simétrico* se trabaja con la nube de columnas  $N^J$  ( $J$  puntos en un espacio de dimensión  $I$ ) cuyas coordenadas, con respecto a la base

estandard, están dadas por las columnas de la matriz perfil de columnas C (renglones de  $C^T$ ). Es decir el j-ésimo punto de la nube es

$$c_j = [c_{1j}, c_{2j}, \dots, c_{ij}]^T =$$

$[n_{1j}/n, \dots, n_{ij}/n, \dots]^T$ , es decir, la columna j de la matriz asociada a la tabla de contingencia, en proporción con el total de esa columna.

Como matriz de pesos a los puntos se elige la matriz  $D_c$ , de suma de columnas.

Notar que, el peso de la columna j (asociada con el punto j) con respecto al total, está dado por

$$(n_{1j} + n_{2j} + \dots + n_{ij}) / n = n_{1j}/n + n_{2j}/n + \dots + n_{ij}/n = p_{1j} + p_{2j} + \dots + p_{ij} = p_{.j}$$

que es, justamente, el elemento j de la diagonal de la matriz  $D_c$  (matriz asignada a los pesos de puntos).

Como, el peso del renglón i (asociado con el eje i), considerando a N, está dado por

$$(n_{i1} + n_{i2} + \dots + n_{ij}) / n = n_{i1}/n + n_{i2}/n + \dots + n_{ij}/n = p_{i1} + p_{i2} + \dots + p_{ij} = p_{i.}$$

que es el elemento i de la diagonal de la matriz  $D_r$ , entonces, como métrica (peso a las dimensiones) se elige a la matriz  $D_r^{-1}$  (métrica ji-cuadrada).

El centroide de la nube  $N^j$  es el vector  $r$ , de suma de renglones de P;

$$\text{centroide de } N^j = \sum_{j=1}^J (p_{.j}) r_j = \sum_{j=1}^J (p_{.j}) [c_{1j} \ c_{2j} \ \dots \ c_{ij}]^T =$$

$$\sum_{j=1}^J (p_{.j}) [p_{1j}/p_{.j}, \dots, p_{ij}/p_{.j}]^T = \sum_{j=1}^J [p_{1j}, \dots, p_{ij}]^T =$$

$$\left[ \sum_{j=1}^J p_{1j}, \dots, \sum_{j=1}^J p_{ij} \right]^T = [p_{1.}, \dots, p_{i.}]^T = r$$

La matriz centrada de  $C^T$ ,  $C^T$ , es aquella cuyo  $j$ -ésimo renglón es  $c_j - r$  ( $C^T = C^T - r1^T$ ,  $1$  es un vector de unos de dimensión  $J$ ).

Para encontrar los  $k^*$  ejes principales, los cuales definen a la nube de puntos  $N^J$ , se tiene que aplicar la técnica de DVS generalizado a la matriz  $C^T$ , donde los vectores singulares izquierdos y derechos son ortonormales con respecto a  $D_c$  y  $D_r^{-1}$ , en la cual los vectores singulares derechos definen a los ejes principales de  $N^J$ .

Es decir,  $C^T$  se puede descomponer como

$$C^T = U D_\mu V^T \quad \text{donde } U^T D_c U = V^T D_r^{-1} V = I$$

entonces, las columnas de  $V$  (vectores singulares derechos) definen los ejes principales y los renglones de  $U D_\mu = C^T D_r^{-1} V$  definen las coordenadas de la nube  $N^J$  con respecto a estos ejes.

Los ejes principales del mejor subespacio de dimensión  $k^*$  son las primeras  $k^*$  columnas de la matriz  $V$ .

### 4.3 Simetría

Observar la siguiente comparación entre los dos Análisis ( $R^J$  y  $R^I$ ):

ANÁLISIS EN  $R^J$ , NUBE  $N^J$  ( $I$  PUNTOS DE DIMENSION  $J$ ):

*Puntos :*

$I$  renglones, de la matriz perfil de renglones  $R$ .

*Matriz de pesos asociados a los puntos :*

la matriz diagonal  $D_r$ .

*Métrica :*

la matriz diagonal  $D_c^{-1}$ .

**Centroide :**  
vector de suma de columnas  $c$ .

**ANÁLISIS EN  $R^J$  NUBE  $N^J$  (J PUNTOS DE DIMENSION I):**

**Puntos :**  
J columnas, de la matriz perfil de columnas, C.

**Matriz de pesos asociados a los puntos :**  
la matriz diagonal  $D_c$ .

**Métrica :**  
la matriz diagonal  $D_r^1$ .

**Centroide :**  
vector de suma de columnas  $r$ .

Se puede apreciar que existe una "gran" simetría, impuesta por las definiciones de métricas y matriz de pesos, entre los análisis en  $R^J$  y  $R^I$ , además el trato para las variables es exactamente el mismo, esto es, los resultados obtenidos, al efectuar el análisis en  $R^I$  ( $R^J$ ) sobre una matriz N, asociada a una tabla de contingencia, son exactamente los mismos; a los resultados después de efectuar el análisis en  $R^J$  ( $R^I$ ) donde se ha trabajado con una matriz transpuesta a N.

Existen algunos casos en los cuales esta "suposición simétrica" resulta inconveniente, como cuando se requiera trabajar con variables cualitativas no-simétricas, es decir, que una preceda a otra o explique a otra, que tengan una dependencia de estructura. Algunos casos podrían ser: droga vs. reacción, diagnóstico vs. medicamento, nivel socio-económico vs. voto. Es por este tipo de casos que en el capítulo siguiente se propone un análisis, desde el punto de vista

de Análisis de Correspondencias, para estudiar este tipo de variables, al cual se le ha denominado Análisis de Correspondencias No-simétrico.

#### 4.4 Resultados y comentarios

1.- La variación global de cada nube de puntos es cuantificada por su *total de inercia*, esto es, la suma de distancias al cuadrado a sus respectivos centroides (considerando los pesos a los puntos y la métrica).

Total de inercia de la nube  $N_i$ .

$$in(N_i) = \sum_{i=1}^I (p_{i.}) (r_i - c)^T Dc^{-1} (r_i - c)$$

Total de inercia de la nube  $N^j$ .

$$in(N^j) = \sum_{j=1}^J (p_{.j}) (c^j - r)^T Dr^{-1} (c^j - r)$$

2.- El total de inercia en ambas nubes es igual al coeficiente de contingencia de significancia al cuadrado calculado sobre  $N$ , esto es, la estadística ji-cuadrada para "independencia" dividida por el total  $n$ .

$$in(N_i) = in(N^j) = X^2/n.$$

$$\text{donde } X^2 = \sum_{i=1}^I \sum_{j=1}^J (n_{ij} - e_{ij})^2 / e_{ij} \text{ y } e_{ij} = n_{i.} \cdot n_{.j} / n$$

**Demostración.**

$$in(N_i) = \sum_{i=1}^I (p_{i.}) (r_i - c)^T Dc^{-1} (r_i - c)$$



$$\begin{aligned}
&= \sum_{i=1}^I (p_{i\cdot}) \sum_{j=1}^J (1/p_{\cdot j}) (r_{ij} - p_{\cdot j})^2 \\
&= \sum_{i=1}^I \sum_{j=1}^J (p_{i\cdot}/p_{\cdot j}) (p_{ij}/p_{i\cdot} - p_{\cdot j})^2 \\
&= \sum_{i=1}^I \sum_{j=1}^J (p_{i\cdot}/p_{\cdot j}) (p_{ij} - p_{i\cdot}p_{\cdot j})^2 / p_{i\cdot}^2 \\
&= \sum_{i=1}^I \sum_{j=1}^J (p_{ij} - p_{i\cdot}p_{\cdot j})^2 / p_{i\cdot}p_{\cdot j}
\end{aligned}$$

Ahora el mismo procedimiento con el total de inercia de la nube  $N^j$ :

$$\begin{aligned}
\text{in}(N^j) &= \sum_{j=1}^J (p_{\cdot j}) (c_j - \bar{x})^T D_x^{-1} (c_j - \bar{x}) \\
&= \sum_{j=1}^J (p_{\cdot j}) \sum_{i=1}^I (1/p_{i\cdot}) (c_{ij} - p_{i\cdot})^2 \\
&= \sum_{i=1}^I \sum_{j=1}^J (p_{\cdot j}/p_{i\cdot}) (p_{ij}/p_{\cdot j} - p_{i\cdot})^2 \\
&= \sum_{i=1}^I \sum_{j=1}^J (p_{\cdot j}/p_{i\cdot}) (p_{ij} - p_{i\cdot}p_{\cdot j})^2 / p_{\cdot j}^2 \\
&= \sum_{i=1}^I \sum_{j=1}^J (p_{ij} - p_{i\cdot}p_{\cdot j})^2 / p_{i\cdot}p_{\cdot j}
\end{aligned}$$

Con esto se ha demostrado que  $\text{in}(N_i) = \text{in}(N^j)$ , ahora se demostrará que  $\text{in}(N^j) = X^2/n$ , para concluir la prueba.

$$\text{in}(N^j) = \sum_{i=1}^I \sum_{j=1}^J (p_{ij} - p_{i\cdot}p_{\cdot j})^2 / p_{i\cdot}p_{\cdot j}$$

$$\begin{aligned}
 &= \sum_{i=1}^I \sum_{j=1}^J (n_{ij}/n - n_{i.}n_{.j}/n^2)^2 / (n_{i.}n_{.j}/n^2) \\
 &= \sum_{i=1}^I \sum_{j=1}^J (n_{ij}/n - e_{ij}/n)^2 / (e_{ij}/n) \\
 &= \sum_{i=1}^I \sum_{j=1}^J ((n_{ij} - e_{ij})^2 / n^2) / (e_{ij}/n) \\
 &= \sum_{i=1}^I \sum_{j=1}^J (n_{ij} - e_{ij})^2 / (ne_{ij}) \\
 &= 1/n \sum_{i=1}^I \sum_{j=1}^J (n_{ij} - e_{ij})^2 / (e_{ij}) = \chi^2/n. \blacksquare
 \end{aligned}$$

3.- Los respectivos subespacios de dimensión  $k^*$  de las nubes  $N_i$  y  $N^j$ , están definidos por los  $k^*$  vectores singulares derechos e izquierdos respectivamente de  $P - rc^T$ , con las métricas  $Dr^{-1}$  y  $Dc^{-1}$ , correspondientes a los  $k^*$  valores singulares más grandes. En otras palabras, los vectores singulares derechos e izquierdos definen los ejes principales de la nube de  $N_i$  y de la nube  $N^j$  respectivamente.

Esto es, la matriz  $P - rc^T$  se puede descomponer como

$$P - rc^T = A D_{\mu} B^T \quad \text{donde } A^T Dr^{-1} A = B^T Dc^{-1} B = I$$

$\mu_1 \geq \dots \geq \mu_k > 0$ . Entonces las columnas de  $A$  y  $B$  definen los ejes principales de las nubes  $N_i$  y  $N^j$  respectivamente.

**Demostración.** Recordar que, DVS de  $R^* = R - 1c^T = Dr^{-1}P - 1c^T$ ;

$$Dr^{-1}P - 1c^T = L D_{\mu} M^T \quad \text{donde } L^T Dr L = M^T Dc^{-1} M = I$$

entonces, las columnas de  $M$  definen los ejes principales y los renglones de  $L D_{\mu}$  las coordenadas principales.

Si se multiplica por  $D_r$  por la izquierda se obtiene

$$P - r e^T = (D_r L) D_\mu M^T \text{ donde } (D_r L)^T D_r^{-1} (D_r L) = M^T D_c^{-1} M = I$$

$$\text{Notar que } (D_r L)^T D_r^{-1} (D_r L) = (L^T D_r) D_r^{-1} (D_r L) = L^T D_r L$$

Esto muestra que las columnas de  $M$  son iguales a las columnas de  $B$ .

Para mostrar que los ejes principales de la nube  $N^j$  se sigue un razonamiento análogo. ■

**Comentario.** Este resultado dice que no es necesario aplicar dos veces DVS, a  $N_i$  y  $N^j$ , para encontrar los ejes principales asociados a cada nube. Pero no dice nada, directamente, acerca de las coordenadas de las nubes con respecto a los ejes (coordenadas principales). En los siguientes resultado se verá como a partir de este, se encontrarán las coordenadas principales de las nubes.

4.- Las respectivas coordenadas de las nubes  $N_i$  y  $N^j$  con respecto a sus ejes principales (coordenadas principales) están relacionadas con los ejes principales de la otra nube. Esto es,

Coordenadas principales de la nube  $N_i$ :

$$F = R^T D_c^{-1} B, \text{ entonces } F = D_r^{-1} A D_\mu.$$

Coordenadas principales de la nube  $N^j$ :

$$G = C^T D_r^{-1} A, \text{ entonces } G = D_c^{-1} B D_\mu.$$

**Demostración.** Se puede escribir F como

$F = D r^{-1} (P - r c^T) D c^{-1} B$  y si al DVS,  $P - r c^T = A D_{\mu} B^T$  se multiplica por  $D c^{-1} B$  por la derecha se obtiene

$$(P - r c^T) D c^{-1} B = A D_{\mu} B^T D c^{-1} B = A D_{\mu} I.$$

entonces  $F = D r^{-1} A D_{\mu}$ .

La demostración para G es análoga. ■

**Comentario.** Las matrices F y G definen las coordenadas de las nubes  $N_i$  y  $N_j$  con respecto a todos los ejes principales. Las coordenadas de los puntos con respecto a un subespacio óptimo de dimensión  $k'$  están contenidas en los renglones de las primeras  $k'$  columnas de F y G.

5.- Los conjuntos de coordenadas de F y G están relacionadas por medio de las siguientes fórmulas.

Transición de F a G:

$$G = C^T F D_{\mu}^{-1}.$$

Transición de G a F:

$$F = R^{-1} G D_{\mu}^{-1}.$$

**Demostración.** Se tiene que  $F = D r^{-1} A D_{\mu}$  entonces  $F D_{\mu}^{-1} = D r^{-1} A$  y como,  $G = C^T D r^{-1} A$  entonces  $G = C^T F D_{\mu}^{-1}$ .

De igual forma,  $G = D c^{-1} B D_{\mu}$  entonces  $G D_{\mu}^{-1} = D c^{-1} B$  y como,  $F = R^{-1} D c^{-1} B$  entonces  $F = R^{-1} G D_{\mu}^{-1}$ . ■

**Comentario.** A partir de este resultado se deduce que no es necesario calcular las dos matrices de coordenadas principales, sino que a partir de una se puede calcular la otra.

## **5. ANALISIS DE CORRESPONDENCIAS NO-SIMETRICO**

La suposición de simetría impuesta sobre las variables, en el análisis anterior, puede ser inconveniente en casos en los que alguna de las variables tenga antecedentes lógicos de la otra, es decir, cuando haya una dependencia de estructura entre las variables (variables no-simétricas). Algunos casos que podrían presentar estas características podrían ser estudios con variables tales como: droga-reacción o sexo-deporte practicado. En estos casos resulta inconveniente dar un tratamiento simétrico a estas variables.

Los italianos Luigi D'Ambra y Natale Lauro (1982) desarrollan el Análisis en Componentes Principales sobre Subespacios Referencia, para el uso de variables no-simétricas. A partir de esta técnica, obtienen el Análisis de Correspondencias No-simétrico, para variables cualitativas no-simétricas.

En este trabajo es presentado el Análisis de Correspondencias No-simétrico usando la herramienta de Descomposición en Valores Singulares

El Análisis de Correspondencias No-simétrico hace un cierto tipo de discriminación para cada una de estas variables y su objetivo es evaluar la influencia de las  $J$  categorías, de la variable cualitativa  $J$ , a la cual se le dará el tratamiento de variable explicativa; sobre la distribución de la variable cualitativa  $I$ , a la cual se le dará el tratamiento de variable respuesta. En otras palabras, que tanta influencia hay de  $J$  sobre  $I$ .

Para tratar de establecer la dependencia de estructura entre la variable respuesta  $I$  y la variable explicativa  $J$ , se realizan dos análisis, en el espacio de columnas y en el espacio de renglones, obteniendo así los subespacios duales para cada nube de puntos.

### 5.1 Análisis en $R^I$

Considérese una nube de puntos  $N^J$  ( $J$  puntos en un espacio de dimensión  $I$ ) en la cual, se puede relacionar los  $J$  puntos con la variable explicativa  $J$  y las  $I$  dimensiones con la variable respuesta  $I$ .

De la misma forma que en el caso simétrico, las coordenadas de  $N^J$ , con respecto a la base estándar, están dadas por las columnas de la matriz perfil de columnas  $C$ . Entonces el  $j$ -ésimo punto de la nube es

$$c^j = [c_{1j}, c_{2j}, \dots, c_{ij}]^T = [n_{1j}/n_j, \dots, n_{ij}/n_j]^T$$

Es decir, la columna  $j$  de la matriz asociada a la tabla de contingencia, en proporción con el total de esa columna. Esta nube se eligió así ya que está relacionada con la variable explicativa.

La matriz de pesos a los puntos seguirá siendo  $D_c$ .

En el caso simétrico se eligió como métrica a la matriz  $D_r^{-1}$ , pero en el caso no-simétrico podría ser inconveniente, ya que cuando la frecuencia de la categoría  $i$ , de la variable respuesta  $l$ , es baja, entonces, el coeficiente  $1/p_{il}$  causa que esta categoría tenga un papel importante. Es por esto que la métrica a elegir es la matriz identidad  $(I_{\omega})$ , es decir, el peso a la dimensión  $i$  es 1, para que el papel que tenga cada categoría sea el mismo.

Notar que el centroide es el mismo, es decir,  $r$  (suma de renglones de  $P$ ), ya que la nube  $N^J$  es la misma y los pesos a los puntos también y por lo tanto la matriz centrada  $C^T$  sigue siendo igual.

Ahora, aplicando DVS:

$$C^T = U D_{\mu} V^T \quad \text{donde} \quad U^T D_c U = V^T V = I$$

Las columnas de  $V$  definen los ejes principales y los renglones de  $U D_{\mu} = C^T V$  son las coordenadas (coordenadas principales) con respecto a estos ejes.

## 5.2 Análisis en $R^J$

El Análisis en  $R^J$ , en el caso no-simétrico, es muy diferente al del caso simétrico, ya que se tiene que elegir una nube  $N_l$  ( $l$  puntos de dimensión  $J$ ) más representativa considerando que hay un tratamiento diferente para cada una de



las variables  $J$  e  $I$ . Observar que la variable  $I$  está relacionada con los puntos y que la variable  $J$  con las dimensiones.

Entonces, para cada punto de la nube  $N_i$ , se consideran las  $J$  distribuciones condicionales  $p_{i\cdot}/p_{i\cdot}$ , con referencia al vector de suma de renglones  $r$ .

Sea la matriz  $T_{Dc}$  con estas características, entonces el término general  $t_{ij}$  está dado por

$$t_{ij} = [p_{i\cdot}/p_{i\cdot} - p_{i\cdot}]$$

Es decir, las coordenadas del  $i$ -ésimo punto, con respecto a la base estandard, están dadas por el  $i$ -ésimo renglón de la matriz  $T$

$$\begin{aligned} t_i &= [p_{i1}/p_{i\cdot} - p_{i\cdot}, \dots, p_{iJ}/p_{i\cdot} - p_{i\cdot}] = \\ &= [(p_{i1} - p_{i\cdot}p_{1\cdot})/p_{i\cdot}, \dots, (p_{iJ} - p_{i\cdot}p_{J\cdot})/p_{i\cdot}] = \\ &= [(n_{i1} - n_{i\cdot}n_{1\cdot}/n)/n_{i\cdot}, \dots, (n_{iJ} - n_{i\cdot}n_{J\cdot}/n)/n_{i\cdot}] \end{aligned}$$

Como matriz de pesos a los puntos, relacionada con la variable  $I$ , se elige a la matriz identidad, ya que se quiere dar el mismo peso a cada punto. Como matriz de pesos a las dimensiones (métrica) se elige la matriz  $Dc$ , ya que el peso a cada dimensión está en proporción del total de cada una de las columnas de  $P$ .

Es fácil ver que el centroide es el vector  $0$ ;

centroide de  $N_i =$

$$\sum_{i=1}^I (1) [p_{i1}/p_{i\cdot} - p_{i\cdot}, \dots, p_{iJ}/p_{i\cdot} - p_{i\cdot}] / n =$$

$$[\sum_{i=1}^I p_{i1}/p_{i\cdot} - \sum_{i=1}^I p_{i\cdot}, \dots, \sum_{i=1}^I p_{iJ}/p_{i\cdot} - \sum_{i=1}^I p_{i\cdot}] / n =$$

$$[p_{\cdot 1}/p_{\cdot 1} - 1, \dots, p_{\cdot J}/p_{\cdot J} - 1] / n = 0.$$

Ahora para encontrar los ejes principales del mejor subespacio de dimensión  $k^*$ , se tiene que aplicar DVS de la siguiente manera:

$$T = M D_{\mu} L^T \quad \text{donde } M^T M = L^T D_c L = I$$

Las columnas de  $N$  definen los ejes principales y los renglones de  $M D_{\mu} = T D_c L$ , son las coordenadas con respecto a estos ejes.

### 5.3 Índice $\tau$ de Goodman-Kruskal

El análisis clásico de tablas de contingencia sugiere medidas de asociación no simétricas como el índice  $\tau$  de descomposición de predictibilidad debido a Goodman-Kruskal, el cual permite dar una interpretación probabilística.

El Análisis de Correspondencias No-simétrico se basa en el índice descomposición de predictibilidad de Goodman-Kruskal (1954), ya que ofrece algunas ventajas más, que si se usara el índice  $\phi^2$  de Pearson ( $\phi^2 = X^2/n$ ).

El índice  $\tau$  está dado por

$$\tau = (\Sigma (p_{1j}^2/p_{.j}) - \Sigma p_{1.}^2) / (1 - \Sigma p_{1.}^2).$$

Si se comparan ambos índices se puede observar que:

- $\tau$  contrariamente a  $\phi^2$ , tiene definido un límite superior;
- $\tau$  contrariamente a  $\phi^2$ , no varía si los elementos de la tabla de contingencia son multiplicados por una constante;
- $\tau$ , es menos sensible cuando la distribución marginal de la variable respuesta es muy asimétrica;

- $\tau$  contrariamente a  $\phi^2$ , crece cuando la varianza de la variable respuesta decrece;

- $\tau$  y  $\phi^2$ , decrecen cuando la varianza de la variable explicativa decrece.

## 6. APLICACION

En este capítulo, se mostrará un ejemplo de ambos tipos de análisis usando el programa ANACORR (Apéndice B). Este, es un ejemplo médico donde se intenta mostrar algunas de las ventajas y desventajas aunadas al uso de los dos análisis, aunque cada uno de ellos esté basado en supuestos diferentes.

En la primera parte del capítulo, se definirán las variables cualitativas a usar, así como las categorías de cada una de ellas, y a partir de éstas se construirá la tabla de contingencia. En la segunda y tercera partes, se le aplicarán a este caso los Análisis de Correspondencias Simétrico y No-simétrico, respectivamente, mostrando en ambas partes los resultados obtenidos en ANACORR. En la última parte, se tratan algunos comentarios de los resultados obtenidos.

## 6.1 Datos médicos

Se considerará un caso de estudio, donde los datos conciernen a la frecuencia, con la cual seis medicamentos fueron prescritos para siete enfermedades.

Las variables cualitativas, sus categorías y las claves de éstas, se describen en los siguientes cuadros.

MEDICAMENTO	
1. Penicilina	PENI
2. Tifomicina	TIFO
3. Tetraciclina	TETR
4. Eritromicina	ERIT
5. Tiofenicina	TIOF
6. Gentamicina	GENT

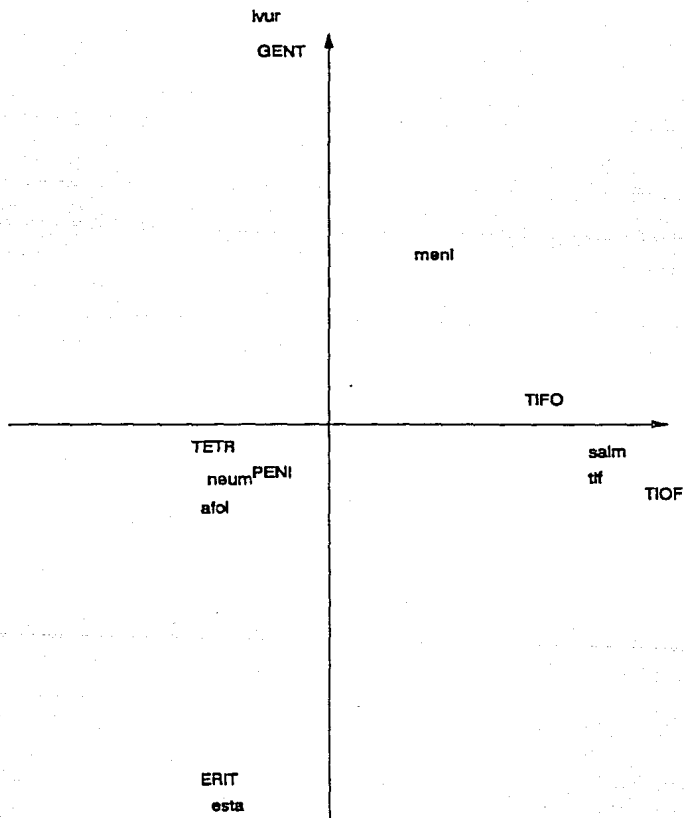
ENFERMEDAD	
1. Tifoidea	tif
2. Salmonelosis	salm
3. Afección Otorrino-laringea	afol
4. Neumocistosis	neum
5. Meningitis	meni
6. Infección vías urinarias	ivur
7. Estafilococcia	esta

69 pacientes que presentaron tener síntomas de alguna de estas enfermedades se les preguntó que medicamento se les prescribió, a partir de esta información (Hospital La Raza) se obtiene la siguiente tabla de contingencia.

MEDICAMENTOS	ENFERMEADES							Tot ren
	tif	salm	afol	neum	meni	ivur	esta	
PENI	0	0	8	7	2	4	3	24
TIFO	4	2	0	0	2	0	0	8
TETR	0	0	5	5	0	2	1	13
ERIT	0	0	3	2	0	0	3	8
TIOF	2	1	0	0	0	0	0	3
GENT	0	0	3	3	1	6	0	13
Tot col	6	3	19	17	5	12	7	69

Este, parece ser un caso en donde una variable cualitativa (ENFERMEDAD) precede a otra (MEDICAMENTO) y por lo tanto parece ser más conveniente realizar Análisis de Correspondencias No-simétrico.

## 6.2 Aplicación simétrica



Al realizar Análisis de Correspondencias Simétrico se obtuvieron los siguientes valores singulares

$$\alpha_1^2 = 0.8783$$

$$\alpha_2^2 = 0.2037$$

$$\alpha_3^2 = 0.0492$$

$$\alpha_4^2 = 0.0286$$

$$\alpha_5^2 = 0.00004$$

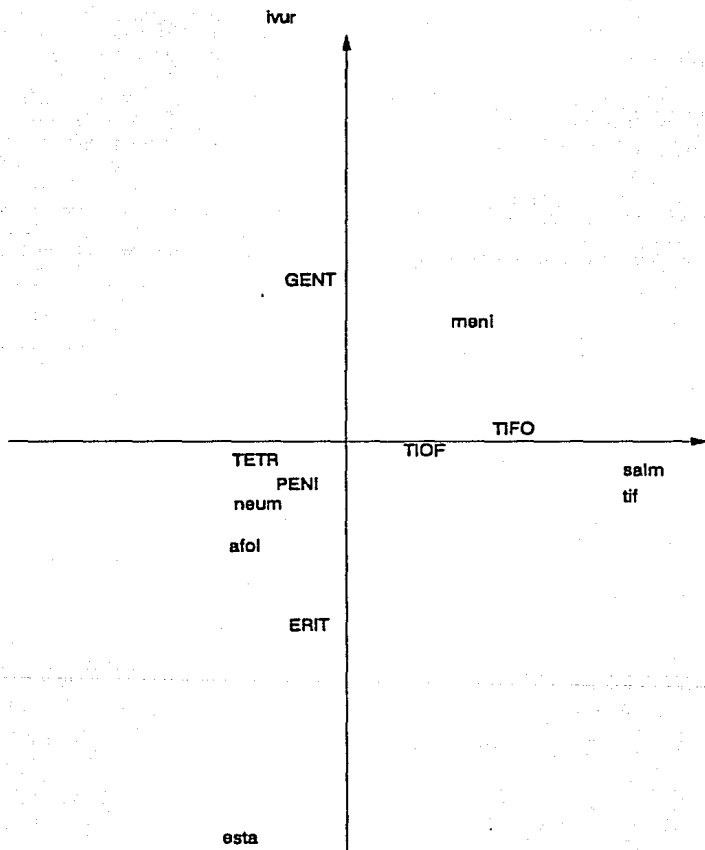
Por lo tanto, el total de inercia (ji-cuadrada de Pearson) obtenido en este caso es 1.159.

El "mejor" desplegado (utilizando los ejes de mayor inercia) de puntos simétrico está dado en la gráfica anterior.

### 6.3 Aplicación no-simétrica

Considerando a la variable MEDICAMENTO como respuesta y como explicativa a la variable ENFERMEDAD, se efectúa Análisis de Correspondencias No-simétrico, con el objetivo de poder prescribir alguno de los medicamentos (categorías de la variable MEDICAMENTO) cuando un individuo presente tener síntomas de alguna de las enfermedades (categorías de la variable ENFERMEDAD). Se obtiene para valor de  $\tau = 0.180$  (Índice Goodman-Kruskal) y el desplegado, no-simétrico, obtenido es el siguiente





## 6.4 Comentarios de las aplicaciones

Si se observan ambas gráficas se puede notar que en la gráfica simétrica, la cercanía entre algunos puntos es más marcada que en la gráfica no-simétrica, es decir, las "parejas" entre categorías de ambas variables son más notorias. Por ejemplo, las categorías *ERIT-esta* y *GENT-ivur* aparecen claramente como "parejas" en la gráfica simétrica, no así en la no-simétrica.

Con referencia en la tabla de contingencia es evidente que, por ejemplo, *ERIT* es intermedia con respecto a las enfermedades para las cuales el medicamento es usado (*afoI* 3, *neum* 2, *esta* 3), y como el Análisis de Correspondencias No-simétrico asigna un mismo peso a las categorías *afoI*, *neum* y *esta* (y a las otras categorías de la variable ENFERMEDAD), en la gráfica se aprecia que *ERIT* está intermedia a estas categorías. A diferencia el Análisis de Correspondencias Simétrico considera los pesos de las categorías y como el menor de estos pesos es el de la categoría *esta*, en la gráfica simétrica forman una notoria "pareja". En otras palabras, la interpretación que se le podría dar en el caso no-simétrico es: de 69 personas a 8 se les prescribió eritromicina de las cuales 3 presentaron síntomas de afección otorrino-laríngea; 2, de neumocistosis y 3, de estafilococcia, entonces eritromicina podría estar explicado por estas tres enfermedades. En el caso simétrico la interpretación podría ser: a 8 se les prescribió eritromicina; de las cuales 3 tenían afección otorrino-laríngea (de 12 totales), 2 neumocistosis (de 17) y 3 estafilococcia (de 7); y a 7 se les encontraron síntomas de estafilococcia; de las cuales a 3 les recetaron penicilina (de 24), a 1 tetraciclina (de 13) y a 3 eritromicina (de 8); entonces se puede apreciar,

apreciar, considerando los totales marginales, que hay más dependencia entre eritromicina y estafilococcia que, por ejemplo, entre eritromicina y neumocistosis, o entre penicilina y estafilococcia.

## **Comentarios generales**

En este trabajo, se han desarrollado los Análisis de Correspondencias Simétrico (clásico) y No-simétrico, basados en la poderosa técnica numérica de Descomposición en Valores Singulares (DVS).

El Análisis de Correspondencias Simétrico, parte de dos variables cualitativas no independientes (sentido de probabilidad) y trata de encontrar la dependencia entre estas variables; no así el Análisis de Correspondencias No-simétrico, el cual parte de dos variables no-simétricas (variable explicativa y otra de respuesta) y trata de mostrar la dependencia de estructura entre ambas variables.

Considerando lo anterior puede observarse que los análisis están basados sobre hipótesis diferentes, es decir se puede elegir uno u otro dependiendo de lo que se quiera investigar. En este trabajo se ha querido mostrar la diferencia de los métodos para desarrollar cada análisis, estos métodos, basados en DVS, están apoyados teóricamente en Análisis en Componentes Principales, para el caso simétrico, y en Análisis en Componentes Principales sobre un Subespacio Referencia (N. Lauro, D'Ambra), para el caso no-simétrico.

Por último, dado que realizar los análisis manualmente resultaría casi imposible, muy probablemente erróneo (considerando errores de cálculo) y no se conoce algún programa que realice Análisis de Correspondencias No-simétrico, se

ha propuesto el programa en computadora ANACORR el cual realiza ambos análisis, en una forma muy sencilla.

## **APENDICE A**

### **EJEMPLOS DE OBTENCION DE SUBESPACIOS Y ASIGNACION DE PESOS**

#### **A.1 Ejemplos de obtención de subespacios de dimensión menor para puntos en $R^2$ y $R^3$**

Considérese que se está interesado en las medidas de altura y de peso de personas. Así, los vectores de datos, de orden 2, para tres personas (personas A, B, C) son:

$$\mathbf{x}_A = [175 \ 70]^T \quad \mathbf{x}_B = [170 \ 68]^T \quad \mathbf{x}_C = [150 \ 60]^T.$$

Estos vectores pueden ser vistos como puntos en el plano, (espacio de dimensión 2) los valores de cada vector son las coordenadas del punto.

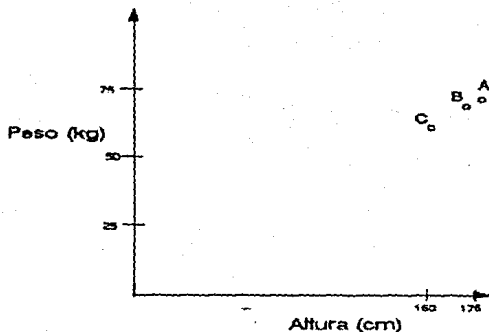


Fig 1.1

Se puede observar que los tres puntos;  $x_A$ ,  $x_B$  y  $x_C$ ; están en una línea recta que, además, pasa por el origen. Esto significa que los tres vectores de datos pueden ser expresados como múltiplos de un solo vector, por ejemplo el vector  $b = [10 \ 4]^T$ :

$$x_A = 17.5b \quad x_B = 17b \quad x_C = 15b$$

Así estos tres puntos se representan en un subespacio de dimensión 1, definido por el vector base  $b$  y con coordenadas, con respecto a  $b$ , de 17.5, 17 y 15 respectivamente.

$$\mathbf{b} = 10\mathbf{e}_1 + 4\mathbf{e}_2$$

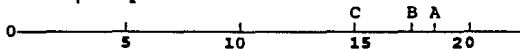


FIG 1.3

Esta nueva dimensión, la cual es una combinación de las dimensiones originales de altura y de peso, puede ser interpretada como una dimensión de "tamaño" y la fig 1.3 muestra dónde se encuentra cada persona sobre esta dimensión, así la persona B es más pequeña que A y C es mucho más pequeña que B.

Es claro que hay una infinidad de caminos para la elección de un vector base de la dimensión de "tamaño". Por ejemplo,  $\mathbf{c} = [5 \ 2]^T$  es otro vector base y las coordenadas de los tres puntos con respecto a  $\mathbf{c}$  son 35, 34 y 30 respectivamente.

Similarmente la base estándar en el espacio de 2 dimensiones es solo una, de un infinito de posibles bases.

Considérese otro ejemplo en el espacio de dos dimensiones, donde los tres vectores

$$\mathbf{y}_1 = [1500 \ 2000]^T \quad \mathbf{y}_2 = [2000 \ 1000]^T \quad \mathbf{y}_3 = [2200 \ 600]^T$$

representan las exportaciones e importaciones de cierto producto en tres años consecutivos.



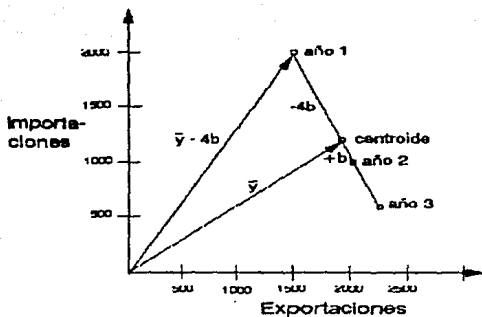


Fig 1.4

Estos tres puntos están sobre una línea recta, pero, ésta no pasa por el origen y los vectores no pueden ser expresados, directamente, como múltiplo de un vector base.

Entonces se usará el hecho de que cualquier vector puede ser expresado como la suma de un vector fijo desde el origen mas un múltiplo de un vector a lo largo de la recta.

El vector fijo a elegir es el centroide (vector de medias)

$$\bar{y} = [(1500+2000+2200)/3, (2000+1000+600)/3]^T = [1900 \ 1200]^T$$

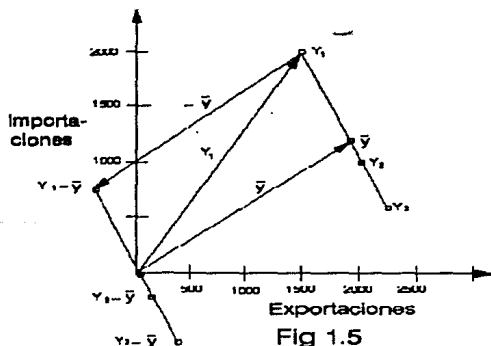
El vector  $b$  a lo largo de la recta puede ser elegido como  $b = [100 \ -200]^T$  y así los tres puntos  $y_1$ ,  $y_2$  e  $y_3$  pueden expresarse como

$$y_1 = \bar{y} - 4b \quad y_2 = \bar{y} + b \quad y_3 = \bar{y} + 3b$$

Alternativamente, se expresan los vectores como desviaciones desde su centroide

$$z_1 = y_1 - \bar{y} \quad z_2 = y_2 - \bar{y} \quad z_3 = y_3 - \bar{y}$$

Este "centrado" da como resultado que el origen, sea el centroide de los vectores  $z_1$ ,  $z_2$  y  $z_3$ ; y éstos sean múltiplos del vector base  $b$ .



La representación en un espacio de una dimensión, de los tres puntos está dada en la fig. 1.6., y las coordenadas de los puntos centrados son -4, 1, 3 respectivamente.

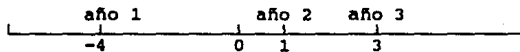


FIG 1.6

Esta dimensión podría interpretarse como una medida de "crecimiento" de la manufactura local del producto, con movimiento a la derecha indicando un "crecimiento" positivo. El vector base  $\mathbf{b}$  ha sido etiquetado como una unidad de "crecimiento" positivo: un incremento en exportaciones de 100 acompañado de un decremento en importaciones de 200. Para una interpretación más real de la situación se necesita información adicional; como, por ejemplo, el consumo local del producto año con año, ya que, el decrecimiento en importaciones se debe a un bajo consumo.

Supóngase que se tiene un consumo local por año, de 2500, 2700, 3200 respectivamente, los vectores nuevos son de dimensión 3:

$$\mathbf{y}_1 = [1500 \ 2000 \ 2500]^T \quad \mathbf{y}_2 = [2000 \ 1000 \ 2700]^T \quad \mathbf{y}_3 = [2200 \ 600 \ 3200]^T$$

El centroide de estos puntos con la nueva dimensión es

$$\bar{\mathbf{y}} = [1900 \ 1200 \ 2800]$$

Como  $\mathbf{b} = [100 \ -200]^T = 100\mathbf{e}_1 - 200\mathbf{e}_2$  donde  $\mathbf{e}_1$  y  $\mathbf{e}_2$  son los vectores unitarios de los ejes de exportaciones e importaciones entonces,

$$y_1 = [1500 \ 2000 \ 2500]^T = \bar{y} - 4b - 300e_3$$

donde  $e_3$  es el vector unitario sobre el nuevo eje de consumo. Representando los puntos centrados en un espacio de dos dimensiones con ejes  $b$  y  $e_3$  da como resultado la fig 1.8. y las coordenadas de los puntos con respecto a estos ejes son  $[-4 \ -300]$ ,  $[1 \ -100]$  y  $[3 \ 400]$ .

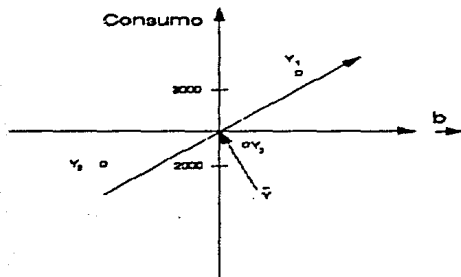


Fig 1.8

Estos puntos están muy próximos a una recta, definida por un vector  $c$ . Si se obtienen las desviaciones de los tres puntos con respecto a este vector, se podrían reducir los puntos originales en un espacio de dimensión 3 a un espacio de dimensión 1, desplegado a lo largo de un solo eje definido por el vector  $c$ . Este

vector es ahora la combinación del eje de consumo y el eje definido por  $\mathbf{b}$ , en otras palabras, combinación de las tres dimensiones originales.

Este ejemplo, en tres dimensiones, es un simple ejemplo de reducción de dimensionalidad de un conjunto, o *nube*, de puntos en un espacio multidimensional. En aplicaciones actuales, usualmente se trabaja con puntos de en espacios de dimensión mucho más alta.

La pregunta obligada es ¿cómo encontrar un subespacio de menor dimensión el cual esté "más cerca" a la nube de puntos?. Lo cual es tratada en el Capítulo 2.

## A.2 Asignación de pesos a ejes y puntos

Como se sabe, se puede definir distancia entre dos puntos a partir de la definición de producto escalar de vectores.

El producto escalar entre los vectores  $\mathbf{a} = [a_1, a_2, \dots, a_J]^T$  y  $\mathbf{b} = [b_1, b_2, \dots, b_J]^T$ , se denota como

$$\langle \mathbf{a}, \mathbf{b} \rangle = \sum_{j=1}^J a_j b_j = \mathbf{a}^T \mathbf{b}$$

Entonces la distancia entre  $\mathbf{a}$  y  $\mathbf{b}$  es

$$d(\mathbf{a}, \mathbf{b}) = \left( \sum_{j=1}^J (a_j - b_j)^2 \right)^{1/2} = \langle \mathbf{a} - \mathbf{b}, \mathbf{a} - \mathbf{b} \rangle^{1/2}$$

Estos resultados dependen, absolutamente, de un sistema perpendicular de coordenadas, ya que la base estándar de vectores  $\mathbf{e}_1, \mathbf{e}_2, \dots, \mathbf{e}_J$  son *ortogonales* ( $\langle \mathbf{e}_i, \mathbf{e}_j \rangle = 0$  si  $i \neq j$ ), y además estos vectores tienen longitud de una unidad,

*normalizados*, entonces la base estandard, elemental para estos resultados es una base *ortonormal*.

En la definición anterior de distancia se presupone que las dimensiones asociadas a cada dimensión son de la misma naturaleza. Se trata de longitudes medidas de acuerdo a una misma cantidad.

En casi todas las situaciones que se podrían presentar, a cada eje (dimensión) se le asocia una característica que implícitamente tiene unidad diferente.

Recordando el ejemplo de las medidas de altura y peso sobre diferentes personas, es claro que se está trabando con unidades completamente diferentes (cm, para altura y kg, para peso). La altura siempre tendrá un valor mucho más alto al del peso, así la medida de altura contribuirá más, relativamente, al resultado de la distancia; pero si la unidad de altura es m, entonces, el peso es el que contribuirá más en el resultado de la distancia.

Claramente, no es conveniente que las distancias dependan directamente de la elección de la escala de las medidas.

Es necesario, entonces, "ponderar" cada eje con respecto a las diferencias de las variables y así trabajar en un Espacio Euclídiano con Peso. Comúnmente, se dividen las medidas por sus respectivas desviaciones estándares. Por ejemplo, suponer que la desviación estandard de la altura y peso de una muestra de gente es 30 cm y 10 kg respectivamente. Sean  $x$  e  $y$  los vectores que representan las medidas de dos personas, los vectores estandarizados son

$\mathbf{x}_a = [x_1/30 \ x_2/10]^T$  y  $\mathbf{y}_a = [y_1/30 \ y_2/10]^T$  y su producto escalar  $\langle \mathbf{x}_a, \mathbf{y}_a \rangle = x_1 y_1 / 30^2 + x_2 y_2 / 10^2$ .

Este producto escalar puede ser escrito como

$$\langle \mathbf{x}_a, \mathbf{y}_a \rangle = x_1 y_1 / 30^2 + x_2 y_2 / 10^2 = \mathbf{x}^T \mathbf{D}_a^{-1} \mathbf{y}^T$$

donde

$$\mathbf{D}_a^{-1} = \begin{bmatrix} 1/30^2 & 0 \\ 0 & 1/10^2 \end{bmatrix}$$

es la matriz diagonal de las inversas de las varianzas, asociada con el Espacio Euclidiano Pesado.

Si se está trabajando con puntos en un espacio de dimensión  $J$ , la matriz asociada (comunmente llamada métrica), será de  $J \times J$ , simétrica, y definida positiva.

Las métricas más utilizadas son las diagonales, ya que son de fácil interpretación.

En general, se define el producto escalar como

$\mathbf{x}^T \mathbf{D}_p \mathbf{y} = \sum_j p_j x_j y_j$ , donde  $p_1, \dots, p_j$  son números reales positivos elementos de la diagonal de  $\mathbf{D}_p$  ( $\mathbf{D}_p = \text{diag}(p_j)$ ) que significan; el peso asociado a cada dimensión del espacio.

Entonces la distancia al cuadrado se expresa como

$$d^2(\mathbf{x}, \mathbf{y}) = (\mathbf{x} - \mathbf{y})^T \mathbf{D}_p (\mathbf{x} - \mathbf{y}) = \sum_j p_j (x_j - y_j)^2.$$

Una base,  $\mathbf{b}_1, \dots, \mathbf{b}_j$ , ortonormal para este espacio (ortonormal con respecto a la métrica  $\mathbf{D}_p$ ) debe satisfacer

$\mathbf{b}_i^T \mathbf{D} \mathbf{p} \mathbf{b}_j = 0$ , si  $i \neq j$  (ortogonal) y

$\mathbf{b}_i^T \mathbf{D} \mathbf{p} \mathbf{b}_i = 1$  (normal).

Lo anterior referente a pesos asignados a ejes, y tratando lo referente a pesos asignados a puntos, se pretende que los métodos usados puedan adaptarse a situaciones prácticas, entonces se considera, que las observaciones no necesariamente tengan la misma importancia. En la estadística clásica cada observación es, generalmente, considerada como aportadora de la misma cantidad de información. En la realidad, se puede privilegiar a algunas observaciones más que otras en función de su importancia con respecto a cierto criterio. De tal forma, a cada observación se le asociará un peso  $w_i$  tal que;  $w_i > 0$ . Es claro, que si todas las observaciones se consideren con la misma importancia, entonces  $w_i = 1/(\text{total de obs.})$ .

Generalmente estos pesos de puntos serán almacenados en una matriz diagonal (matriz de pesos asignados a l puntos)  $\mathbf{D} \mathbf{w}_{1 \times l}$ .

Así, el centroide del conjunto de puntos,  $\mathbf{l} = \{\mathbf{x}_1 \ \mathbf{x}_2 \ \dots \ \mathbf{x}_l\}$ , con pesos  $w_1 \ w_2 \ \dots \ w_l$  se define como

$$\bar{\mathbf{x}} = \frac{\sum_{i=1}^l w_i \mathbf{x}_i}{\sum_{i=1}^l w_i}.$$

Normalmente, se supondrá que  $\sum_{i=1}^l w_i = 1$ , a menos que se especifique lo contrario, en este caso la expresión para el centroide es

$$\bar{\mathbf{x}} = \sum_{i=1}^l w_i \mathbf{x}_i$$



## **APENDICE B**

### **PROPUESTA DE UN PROGRAMA EN COMPUTADORA**

Para realizar los Análisis de Correspondencias Simétrico y No-simétrico, es decir, obtener resultados y un desplegado de puntos, se propone el programa ANACORR, el cual realiza los análisis, en un ambiente de microcomputadora amigable y de fácil operación. ANACORR está programado en el lenguaje de Turbo Pascal (Versión 5.0), apoyado en el manejador de bases de datos Btrieve (Versión 5.1) y basado en la teoría presentada en los capítulos anteriores.

ANACORR está basado en un ambiente de menús en los cuales, para ejecutar alguna opción se puede lograr de dos maneras:

a) presionar las teclas de flechas, en la dirección de la opción deseada hasta llegar a ésta y después presionar la tecla <ENTER>.

b) Presionar la tecla de la letra de la opción que tiene el color diferente.

## B.1 Instalación

Es recomendable que el programa ANACORR se pueda ejecutar en un disco duro, es por esto que cuenta con una herramienta para instalarlo del disco flexible al disco duro, para lo cual se deben de seguir los siguientes pasos:

1) colocarse en el lugar, del disco duro, donde será instalado el programa. Por ejemplo si se quiere trabajar en el directorio PRUEBA, se debe estar en C:\PRUEBA>.

2) cambiarse al disco flexible (generalmente A: o B:) y elegir el subdirectorio ANACORR (A:>CD ANACORR),

3) teclear el comando, A:\ANACORR\INSTALA C:

La herramienta INSTALA copiará todo lo necesario para trabajar con ANACORR sobre el directorio C:\PRUEBA.

## B.2 ANACORR

Para entrar a ANACORR se deben seguir dos pasos:

1) Cambiarse al subdirectorio ANACORR

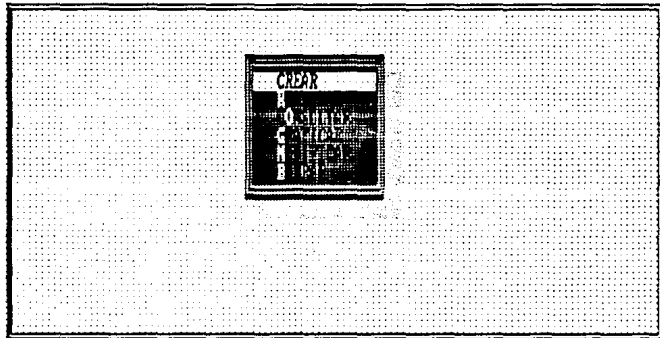
CD ANACORR

2) Teclear ANACORR

\ANACORR>ANACORR

Una vez realizado esto, ya se está dentro del programa y éste desplegará la siguiente pantalla, correspondiente al MENU PRINCIPAL.

ANALISIS DE CORRESPONDENCIAS U. N. A. M.	NOMBRE PRINCIPAL
---	------------------



Crear estudio

A continuación se explicará a detalle como pueden ser utilizadas las opciones de este menú.

**CREAR.** Esta opción permite alimentarle al programa la información necesaria de un estudio determinado, para después realizar, alguno de los análisis. La opción está formada por tres pantallas en las cuales se puede apreciar la información requerida por ANACORR.

ANÁLISIS DE CORRESPONDENCIAS  
U N A M

CREAR UN ESTUDIO

## CARACTERÍSTICAS DEL ESTUDIO

NOMBRE DEL ESTUDIO : ESTUDIO DE CUCOS Y FLECO

ARCHIVO DEL ESTUDIO : ESTUD

NOMBRE DE LA 1er VARIABLE : CUANTOS

NOMBRE DE LA 2da VARIABLE : ESTUD FLECO

NO. DE CATEGORIAS DE 1er VAR : 2

NO. DE CATEGORIAS DE 2da VAR : 2

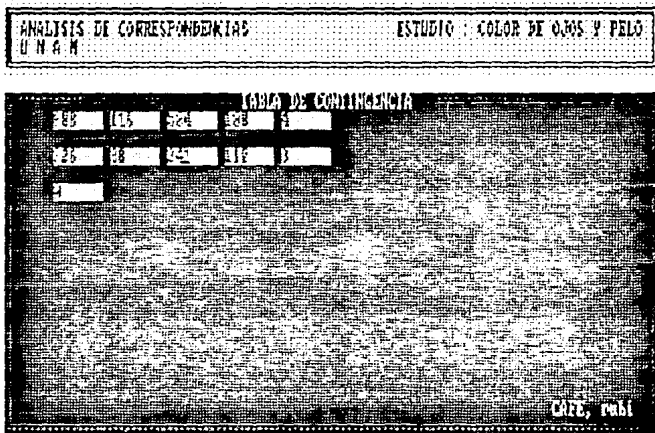
DEFINICION CORRECTA (S/N) :

ANÁLISIS DE CORRESPONDENCIAS  
U N A M

ESTUDIO : COLOR DE OJOS Y PELO

NOMBRE DE CATEGORÍAS	
COLOR OJOS	COLOR PELO
1 CLAR	1 rubi
2 AZUL	2 rojo
3 CAFE	3 cast
4 OSC	4 canb
	5 negro

DEFINICION CORRECTA (S/D)



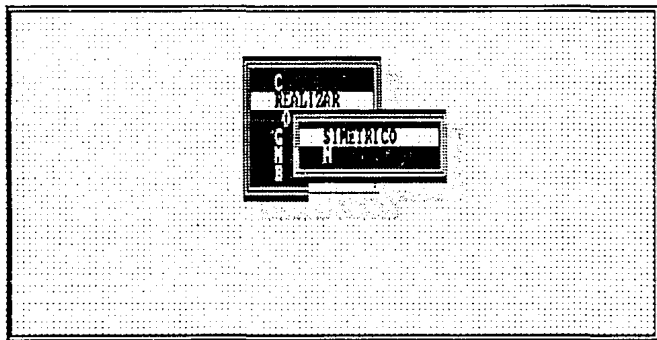
En la primer pantalla se requieren datos generales del estudio, tales como: nombre del estudio, nombre del archivo en el cual se almacenará la información, el programa no permitirá que el nombre de este archivo ya se le haya asignado a otro; nombre y número de categorías de las variables, para el caso no-simétrico se considerará la primer variable como respuesta y la segunda como explicativa. La información que se requiere en la segunda pantalla es, básicamente, el nombre o una clave para identificar las categorías de cada variable. La última pantalla consta de la tabla de contingencia.

Al final de cada pantalla, se pregunta si la información que se ha tecleado es correcta, si ésta no lo es se puede volver a teclear desde el principio de la

pantalla, además con la tecla <ESC> en cualquier momento se puede cancelar la creación del estudio.

**REALIZAR.** Una vez creado un estudio, el siguiente paso sería aplicarle alguno de los análisis, éste es el objetivo de esta opción.

Al ejecutar la opción de REALIZAR, se permitirá elegir el tipo de análisis (SIMETRICO o NO-SIMETRICO).



Realizar Análisis de Correspondencias Simétrico

Ya elegido el tipo de análisis se mostrará la lista de todos los estudios creados a los cuales se les pueda realizar este tipo de análisis.

ANÁLISIS DE CORRESPONDENCIAS U N A M		REALIZAR UN ESTUDIO SIMÉTRICO			
Nombre de Archivo	Nombre de Estudio	Nombre 1er Var	Nombre 2da Var	Categ. 1er.	Categ. 2da
CIUDADO	CIUDADO EN EMPRESA	PUESTO	TIPO RUMADOR	5	4
CIUDAD	CIUDAD DE AÑOS Y FELO	COLOS OJOS	CARRE FELO	4	3
ISRAEL	FRECUENCIA ISRAELI	FRECUENCIA	RESIDENCIA	7	5

F2-Realizar estudio ESC-Salir

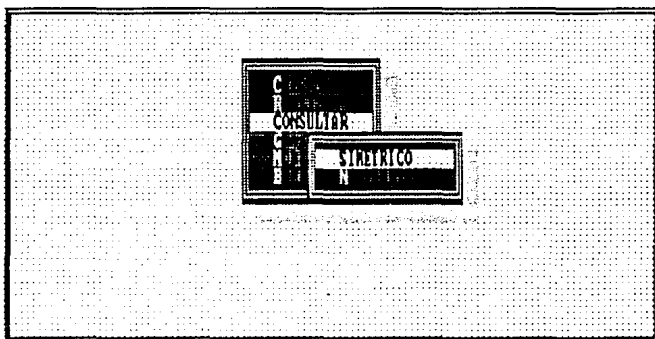
Cuando se haya seleccionado el estudio, para llevar a cabo el análisis solo se tiene que presionar la tecla <F2>.

**CONSULTAR.** Permite consultar los resultados de un estudio al cual se le aplicó el análisis requerido.

Al ejecutar esta opción se tendrá que elegir el tipo de análisis que se desea consultar.



ANÁLISIS DE CORRESPONDENCIAS MENU PRINCIPAL  
 U N A M



Consultar un estudio con Análisis de Correspondencias SIMETRICO

Cuando se elija el análisis requerido se mostrarán todos los estudios a los cuales se les ha realizado éste análisis.

ANÁLISIS DE CORRESPONDENCIAS  
U N A M

CONSULTAR UN ESTUDIO SIMETRICO

Nombre de Archivo	Nombre de Estudio	Nombre 1er Var	Nombre 2da Var	Categ. 1er	Categ. 2da
CICARRO	CICARRO EN EMPRESA	PUESTO	TIPO SUJECION	5	4
CALOR	CALOR DE AÍOS Y FELD	COLON UICE	CALOR FELD	4	3
ISRAEL	PREOCUPACION ISRAELI	PREOCUPACION	RESIDENCIA	7	5

F2=Ver estudio    F3=Imprimir estudio    ESC=Salir

Para consultar el estudio en la pantalla se tiene que seleccionar éste y después presionar la tecla <F2>, para imprimir los resultados, solo se tiene que presionar la tecla <F3>.

Los comandos que se pueden realizar para consultar los resultados en la pantalla, pueden observarse al presionar la tecla de ayuda <F1>.

ANÁLISIS DE CORRESPONDENCIAS D N A M		CONSULTAR UN ESTUDIO SIMETRICO	
COLOR, SIM			
Comandos			
E	Línea arriba	(Up)	
8	Línea abajo	(Down)	
1	Página arriba	(PgUp)	
2	Página abajo	(PgDn)	
M	Scroll derecha	(^ -)	
	Scroll izquierda	(^ _)	
	Inicio de archivo	(^PgUp)	(Home)
	Fin de archivo	(^PgDn)	(End)
	Salir	(Esc)	

TABLA DE CONT				
688	116	384	183	4
126	22	24	118	2
140	64	123	413	26
98	48	133	681	85
= Líneas: 1 - 18			Totales: 95 líneas 2486 bytes	
= Ayuda =				

GRAFICAR. Con esta opción se puede visualizar la representación en el plano de las categorías, con la posibilidad de elegir los ejes principales por medio de los cuales se desea hacer dicha representación.

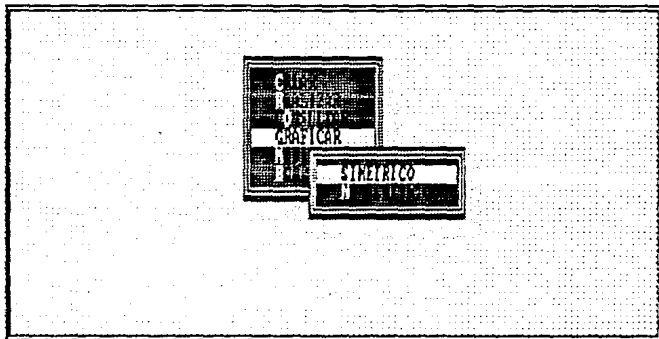
Al igual que en la opción de CONSULTAR, después de haber elegido el tipo de análisis requerido se mostrarán todos los estudios a los cuales se les ha realizado este análisis. Esto lo muestran las siguientes pantallas.

ANÁLISIS DE CORRESPONDENCIA U N A M	CONSULTAR UN ESTUDIO SIMETRICO																				
COLOR.SIM																					
E 8 1 2 M M	<table border="1" style="margin: auto;"> <tr><th colspan="2">Comandos</th></tr> <tr><td>Línea arriba</td><td>(Up)</td></tr> <tr><td>Línea abajo</td><td>(Down)</td></tr> <tr><td>Página arriba</td><td>(PgUp)</td></tr> <tr><td>Página abajo</td><td>(PgDn)</td></tr> <tr><td>Scroll derecha</td><td>(←)</td></tr> <tr><td>Scroll izquierda</td><td>(→)</td></tr> <tr><td>Inicio de archivo</td><td>(Home)</td></tr> <tr><td>Fin de archivo</td><td>(End)</td></tr> <tr><td>Salir</td><td>(Esc)</td></tr> </table>	Comandos		Línea arriba	(Up)	Línea abajo	(Down)	Página arriba	(PgUp)	Página abajo	(PgDn)	Scroll derecha	(←)	Scroll izquierda	(→)	Inicio de archivo	(Home)	Fin de archivo	(End)	Salir	(Esc)
Comandos																					
Línea arriba	(Up)																				
Línea abajo	(Down)																				
Página arriba	(PgUp)																				
Página abajo	(PgDn)																				
Scroll derecha	(←)																				
Scroll izquierda	(→)																				
Inicio de archivo	(Home)																				
Fin de archivo	(End)																				
Salir	(Esc)																				
TABLA DE CONT																					
688	116	584	183	4																	
726	38	241	110	3																	
343	84	929	412	16																	
48	48	483	681	85																	
= Líneas: 1 - 18		Totales: 95 líneas 2486 Bytes																			
Esc-Ataque		Esc-Salir																			

GRAFICAR. Con esta opción se puede visualizar la representación en el plano de las categorías, con la posibilidad de elegir los ejes principales por medio de los cuales se desea hacer dicha representación.

Al igual que en la opción de CONSULTAR, después de haber elegido el tipo de análisis requerido se mostrarán todos los estudios a los cuales se les ha realizado este análisis. Esto lo muestran las siguientes pantallas.

ANÁLISIS DE CORRESPONDENCIAS      NOME PRINCIPAL  
U N A H

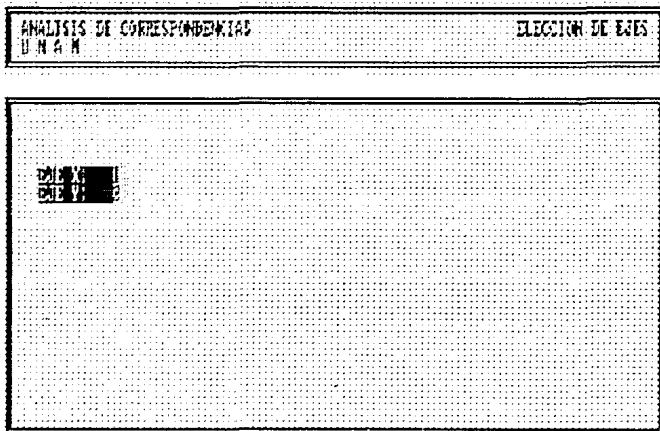


Gráficar un estudio con análisis de Correspondencias SÍMTRICO

ANÁLISIS DE CORRESPONDENCIAS U N A M		GRÁFICAS UN ESTUDIO SIMÉTRICO			
Nombre de Archivo	Nombre de Estudio	Nombre 1er Var	Nombre 2da Var	Categ. 1er	Categ. 2da
CICARRO	CICARRO EN DEPRESA	PUESTO	TIPO FUMADOR	5	4
COLOS	CARNE DE ANOS Y FLEA	COLOS CUBA	CARNE FLEA	4	3
ISRAEL	FEDERACION ISRAELI	FEDERACION	RESISTENCIA	7	5

F2=Vérbica    F3=Impresión    ESC=Salir

Como se puede observar en esta pantalla la representación se puede ver en la pantalla (tecla <F2>) o se puede enviar a la impresora (tecla <F3>). Cuando ya se ha elegido el estudio a graficar y en caso que el número de ejes principales sea mayor que dos, se tendrá que elegir en cual de los ejes se realizará la representación (recordar que la mejor representación se obtiene con los dos primeros ejes).



**MODIFICAR.** Cuando a un estudio creado es necesario hacerle algún cambio se tiene que recurrir a la opción de MODIFICAR, la cual, cuando se haya seleccionado, mostrará todos los estudios creados con la posibilidad de hacerle algún cambio presionando la tecla <F2>.

ANÁLISIS DE CORRESPONDENCIAS U N A M	MODIFICAR UN ESTUDIO
---	----------------------

Nombre de Archivo	Nombre de Estudio	Nombre 1er Var	Nombre 2da Var	Categ. 1er	Categ. 2da
CIGARRA	CIGARRA EN EMPRESA	PUESTO	TIPO FEMARDE	5	4
COLOM	COLOM DE OROS V PULO	COLOM OROS	COLOM PULO	4	5
ISRAEL	PERCEPCION ISRAEL	PERCEPCION	RESISTENCIA	7	5

F2=Modificar estudio. ESC=Salir

La opción de MODIFICAR consta de las mismas pantallas que la opción de CREAR con la diferencia de que éstas tienen la información alimentada cuando se creó el estudio.



ANÁLISIS DE CORRESPONDENCIAS  
U N A M

MODIFICACION DEL ESTUDIO

## CONDICIONES DEL ESTUDIO

NOMBRE DEL ESTUDIO : ENZEL DE GUS Y FIERA

ARCHIVO DEL ESTUDIO : ARCHIVO

NOMBRE DE LA 1er VARIABLE : ENZEL DE GUS

NOMBRE DE LA 2da VARIABLE : ENZEL DE FIERA

NO. DE CATEGORIAS DE 1er VAR : 2

NO. DE CATEGORIAS DE 2da VAR : 2

DEFINICION CORRECTA (S/D) :

ANÁLISIS DE CORRESPONDENCIAS  
U N A M

ESTUDIO : COLOR DE OJOS Y PELO

NOMBRE DE CATEGORÍAS	
COLOR OJOS	COLOR PELO
1 CLAR	
2 AZUL	
3 NEGRO	

ANÁLISIS DE CORRESPONDENCIAS					ESTUDIO : COLOR DE OJOS Y PELLO				
U N A M									
<b>TABLA DE CONTINGENCIA</b>									
308	115	324	428	1					
226	88	141	316	1					
643	54	163	112	26					
58	49	110	391	35					
DEFINICION CORRECTA (S/D) :									

Es importante mencionar, que cuando se haya realizado algún análisis a un estudio que se modifica éste se considerará como un estudio al cual no se le ha realizado ningún tipo de análisis.

**BORRAR.** Esta opción se usa cuando ya no se desea considerar algún estudio creado y se decide eliminarlo.

ANÁLISIS DE CORRESPONDENCIAS  
U. N. A. M.

FORRAR EN ESTUDIO

Nombre de Archivo	Nombre de Estudio	Nombre 1er Var	Nombre 2da Var	Categ. 1er	Categ. 2da
CICARRO	CICARRO EN EMPRESA	FUJATO	TIPO INICIALES	5	4
ISRAEL	FRENTEPORACION ISRAELI	FRENTEPORACION	MINIENCIA	7	3
GUOTIPELO	COLOS DE OROS Y FELD	COLOS GJCS	COLOS FELD	4	5

12-ForrarEstudio ESC-3aliv

ESTA TESIS NO DEBE  
SALIR DE LA BIBLIOTECA

## Referencias

D'Ambra Luigi, Lauro Natale (1982). Principal components analysis onto reference subspaces.

D'Ambra Luigi, Lauro Natale (1983). Non-symmetrical correspondence analysis.

D'Ambra Luigi, Lauro Natale (1984). Non-symmetrical analysis of three-way contingency tables.

Greenacre Micheal J. (1984). "Theory and applications of correspondence analysis". Academic Press, London.

Hernández Cid Rubén (1988). Notas del curso de Análisis de datos.

Lang Serge (1976). "Algebra lineal". Fondo Educativo Interamericano, México.

Leach Chris (1982). "Fundamentos de estadística. Enfoque no paramétrico para ciencias sociales". LIMUSA, México.

Mendenhall William, Scheaffer Richard L., Wackerly Dennis D. (1986). "Estadística matemática con aplicaciones". Grupo Editorial Iberoamerica, México.

**Paquetes usados para el desarrollo del programa ANACORR:**

Btrieve (Ver 5.1).

Turbo Pascal (Ver 5.0).

Turbo Pascal, Graphix Toolbox (Ver 4.0).

Turbo Pascal, Numerical Toolbox (Ver 4.0).

Turbo Professional (Ver 4.0).

Xtrieve Plus (Ver 3.0).