

00321



UNIVERSIDAD NACIONAL AUTÓNOMA DE MÉXICO

76

FACULTAD DE CIENCIAS

Autorizo a la Dirección General de Bibliotecas de UNAM a difundir en formato electrónico e impreso el contenido de mi trabajo de investigación

NOMBRE: Pineda Santa Cruz Ivan Alejandro

FECHA: 18/ Mayo / 2003

FIRMA:

“ANÁLISIS DE CORRESPONDENCIAS PARA EL ESTUDIO DE ASOCIACIÓN ENTRE VARIABLES CATEGÓRICAS: UN ENFOQUE APLICADO”

TESIS CON FALLA DE ORIGEN

T E S I S

QUE PARA OBTENER EL TÍTULO DE ACTUARIO

PRESENTA:
IVAN ALEJANDRO PINEDA SANTA CRUZ



DIRECTOR: M. EN C. INGENIERO RAFAEL MENDOZA RÍOS



MÉXICO, D.F.

FACULTAD DE CIENCIAS
SECCION ESCOLAR

2003



Universidad Nacional
Autónoma de México

Dirección General de Bibliotecas de la UNAM

Biblioteca Central



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 esta 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.

PAGINACION DISCONTINUA



UNIVERSIDAD NACIONAL
DE INGENIERÍA

DRA. MARÍA DE LOURDES ESTEVA PERALTA
Jefa de la División de Estudios Profesionales de la
Facultad de Ciencias
Presente

Comunicamos a usted que hemos revisado el trabajo escrito: Análisis de correspondencias para el estudio de asociación entre variables categóricas: un enfoque aplicado.

realizado por Iván Alejandro Pineda Santa Cruz

con número de cuenta 09220414-2, quién cubrió los créditos de la carrera de actuaría.

Dicho trabajo cuenta con nuestro voto aprobatorio.

Atentamente

Director de Tesis	
Propietario	M. en C. Inocencio Rafael Madrid Ríos.
Propietario	Dr. Rubén Hernández Cid.
Propietario	Mat. Margarita Elvira Chávez Cano.
Suplente	M. en A.P. María Del Pilar Alonso Reyes.
Suplente	Act. Jaime Vázquez Alanilla.

[Handwritten signatures and initials]

Consejo Departamental de matemáticas

M. en C. José Antonio Flores Díaz

FACULTAD DE CIENCIAS

CONSEJO DEPARTAMENTAL

MATEMÁTICAS

Doy gracias:

Al corazón del cielo y al corazón de la tierra,
por el infinito amor con que me crearon.

Ome teotl

A mis padres: Constantino y Esther, por su amor
que transformado en paciencia y esmero
lograron pulirme como una piedra preciosa.

A mis hermanos: Hugo, Ricardo, Lilia, Edgar,
Oswaldo y Yolotzin.

A todos mis maestros, y en especial a:

M. en C. Inocencio Rafael Madrid Ríos.

Dr. Rubén Hernández Cid.

Mat. Margarita Elvira Chávez Cano.

M. en A.P. María Del Pilar Alonso Reyes.

Act. Jaime Vázquez Alamilla.

Por ser amigos y guías en esta aventura que se llama vida.

Porque ellos todo lo sacan de su corazón; obran con deleite,
todo lo obran con calma, con cuidado, crean, inventan
hábilmente lo que piensan, arreglan las cosas, dialogan con
su corazón, encuentran las cosas en su mente...

A todos mis amigos que han contribuido
con sus trazos a dibujar la historia de mi
vida.

A todas las personas que me ayudaron en la
realización de este trabajo.

Recuerda:

"Eres águila, eres ocelote, es tu don, es tu merecimiento...
que no desfallezcan, ni tu venerable mano, ni tu venerable
pie..."

Huehuetlahtolli

"Sólo quien se tome tiempo alcanzará la sabiduría".
Mexicah Tlamatiliztli.

"Y todos somos estrellas cuando descubrimos nuestra esencia".

Celsa Xolalohco

Capítulo I

Elementos básicos

1.1	Tabla de contingencia de dos dimensiones.	1
1.2	Conceptos geométricos y algebraicos obtenidos de la tabla de contingencia, necesarios para el desarrollo del análisis de correspondencias simple (AFC-Simple).	3
1.2.1	Peso o relevancia del renglón i , dado por los datos.	4
1.2.2	Peso o relevancia de la columna j , dado por los datos.	5
1.2.3	Perfiles por renglón y por columna.	5
1.2.4	Centroide.	6
1.2.5	Matriz perfil por renglones y por columnas.	7
1.2.6	Distancia J_i cuadrada.	9

Capítulo II

Análisis factorial de correspondencias simple.
(AFC-Simple)

2.1	Definición de asociación o dependencia para una tabla de contingencia de $I \times J$.	14
2.2	Prueba de independencia J_i cuadrada para una tabla de contingencia de $I \times J$.	15
2.3	Representación gráfica de la asociación entre factores con el AFC-Simple.	19
2.3.1	Formulación gráfica de ajuste de un subespacio vectorial para dos factores de la tabla de contingencia de $I \times J$.	19
2.3.2	Presentación algebraica del análisis de correspondencias.	23
2.3.2.1	Análisis en R^2 . Cálculo de factores	26
2.3.2.2.	Análisis en R^2 . Cálculo de factores	30
2.3.2.3	Relación del análisis de R^2 con el análisis en R^2 . Fórmulas de transición.	34
2.4	Calidad del análisis: contribuciones.	36
2.4.1	Inercia: Descomposición sobre los ejes factoriales. ¿Qué significado geométrico y estadístico tiene la inercia?	37
2.4.2	Inercia relativa y contribución absoluta.	38
2.4.3	Correlaciones al cuadrado: cosenos cuadrados.	41
2.4.4	Calidad de representación de un perfil.	43
2.5	Tratamiento de los elementos suplementarios de una tabla de contingencia.	43
2.6	Descripción de la gráfica del análisis de correspondencias para los datos de la tabla 2.1	45
	Resumen del capítulo	47

Capítulo III
Análisis factorial de correspondencias múltiples
(AFC-M)

		Pág
3.1	El análisis de correlación canónico (ACC)	50
3.2	Consideraciones para el análisis de correspondencias múltiples.	52
3.2.1	Matriz indicadora Z	53
3.2.2	Tabla de Burt asociada con Z	53
3.3	Análisis de correspondencias para una tabla de dos factores considerando diversas matrices	56
3.4	Generalización del análisis de correspondencias utilizando la matriz Z, para más de dos factores.	62
3.5	Resultados del análisis de correspondencias múltiples utilizando la matriz Z.	63
3.6	Dimensionalidad de la configuración de la p categorías en R^n	65
3.7	La mejor representación simultánea para los individuos de la matriz Z.	66
3.8	Ejemplo introductorio. Una forma alternativa de trabajar con más de dos factores.	66
	Utilización de análisis de correspondencias múltiples	72

Conclusiones

Conclusiones	87
Bibliografía	89

Anexos

Anexo 1	Resultados de álgebra lineal y álgebra de matrices	92
Anexo 2	Espacio métrico: determinación del cálculo de la distancia Ji cuadrada.	97
Anexo 3	Maximización de una forma cuadrática bajo una restricción cuadrática	99
Anexo 4	Análisis de correspondencias como un caso particular del análisis de correlación canónico.	101
Anexo 5	Salida del programa Xlstat V4 para los datos de la tabla 1 capítulo II	103
Anexo 6	Base de datos y resultados del AFC-M. Para el ejemplo de tres factores: sexo, sentencia y cargo. Capítulo III	106
Anexo 7	Base de datos para las flores Iris. Resultados del AFC-M.	108

Prefacio

Dada la demanda de profesionistas calificados en el manejo práctico de los métodos estadísticos es indudable la importancia que tiene la incorporación a los cursos que se imparten en la Facultad de Ciencias de la UNAM, la enseñanza de otras técnicas estadísticas como el análisis factorial de correspondencias, debido a su diversidad de aplicaciones prácticas.

El análisis de correspondencias es una técnica estadística que permite describir de manera gráfica la asociación de un conjunto de datos cualitativos.

Antes de que apareciera el análisis de correspondencias ya habian otros métodos estadísticos que permitían estudiar esta asociación por ejemplo la prueba χ^2 de Pearson o los modelos log lineales; sin embargo esta técnica empezó a generar interés debido a la relativa facilidad con que se puede interpretar una gráfica en lugar de una estadística de prueba, sobre todo para personas cuya formación no es matemática o estadística. Sin embargo, esta facilidad puede ser engañosa pues pueden darse por ciertas, conclusiones erróneas cuando existe falta de entendimiento de la técnica en sus supuestos, sus alcances y en la calidad de sus resultados.

La motivación para escribir este trabajo es mostrar los aspectos que deben tomarse en cuenta para una adecuada interpretación del análisis de correspondencias tanto en sus aspectos geométricos, algebraicos y prácticos.

TESIS CON
FALLA DE ORIGEN

Introducción

La estadística puede definirse como el desarrollo y aplicación de técnicas para coleccionar, analizar e interpretar datos cuantitativos o cualitativos. De tal manera que cuando se tiene la información o los datos acerca de un determinado fenómeno que se pretenda estudiar los métodos estadísticos pueden ser utilizados para tratar de describirlos o explicarlos.

Los métodos estadísticos se han desarrollado para estudiar la realidad tomando en cuenta la incertidumbre que un fenómeno aleatorio tiene. Una corriente de pensamiento que ha desarrollado algunos métodos para analizar datos, es la llamada escuela francesa cuyo enfoque se basa en la geometría y el álgebra lineal lo que permite realizar el estudio global de un gran número de variables, con un mínimo de supuestos. El análisis de correspondencias pertenece a esta escuela de pensamiento y permite estudiar la dependencia o asociación por ejemplo entre dos variables cualitativas o factores (variables que no toman valores numéricos), que generalmente son presentadas en forma de tablas de contingencia. En dichas tablas, los renglones y las columnas representan dos particiones de la misma población. Mediante el análisis de correspondencias es posible estudiar la "proximidad" entre las categorías de cada factor, y también se puede estudiar la proximidad entre las categorías de dos o más factores.

H.O.Hartley publicó en 1935 un artículo en el cual da una formulación algebraica de la correlación entre los renglones y columnas de una tabla de contingencia, de tal manera que se le puede atribuir el origen matemático del análisis de correspondencias. Louis Guttman (1941) trató el caso general de más de dos variables cualitativas desarrollando lo que hoy se conoce como análisis de correspondencias múltiples. La forma algebraica geométrica del análisis de correspondencias fue desarrollada por el francés Jean-Paul-Benzecri (1976) dentro de la llamada escuela francesa del análisis de datos. El término francés *correspondence*, fue usado para denotar *sistemas de asociaciones* entre los elementos de dos conjuntos: los renglones y las columnas de una tabla de contingencia.

El análisis de correspondencias puede ser utilizado para estudiar la asociación de más de dos factores generalizándose las ideas desarrolladas para el caso de dos factores además la interpretación de los factores sigue teniendo las mismas reglas que en el análisis de correspondencias simple. Sin embargo, existen algunos inconvenientes en el momento de interpretar los resultados ya que como se aumenta la dimensión del espacio para que las categorías de los factores puedan ser representados en un espacio de menor dimensión se tiene una pérdida de información que se traduce en que los porcentajes de inercia en los ejes son bajos y la descripción de la gráfica no es del todo satisfactoria, es decir, la parte de la inercia explicada pierde aquí su interés ya que los valores propios no representan más que una pequeña parte de la inercia total.

Las formas de realizar el análisis de correspondencias son las siguientes: utilizando la información por individuo que se encuentra en una matriz indicadora denominada por Z o bien utilizando la matriz de Burt; la cual recoge la clasificación de la muestra de acuerdo a cada factor y es en sí una generalización de la tabla de contingencia. Cuando se trabaja con la matriz indicadora Z , se puede hacer uso del análisis de correlación canónico.

En el capítulo I se introducen conceptos básicos que son requeridos para la comprensión del análisis de correspondencias.

En el capítulo II. Análisis Factorial de Correspondencias simples. Se estudiarán los elementos que permiten comprender el tratamiento geométrico y algebraico de los datos, y se desarrollará el análisis de correspondencias para una tabla de contingencia de $I \times J$.

En el capítulo III. Análisis de Correspondencias Múltiples. Se presenta el análisis de correspondencias múltiples, como una generalización del análisis de correspondencias simple, además presentarán las formas equivalentes entre el análisis de correspondencias llevado a cabo con los datos de la matriz indicadora Z o bien mediante el uso de la matriz de Burt. Así mismo, se presentarán ejemplos que ayuden a comprender la interpretación del análisis.

Finalmente se presentarán las conclusiones que se obtuvieron con el estudio del análisis de correspondencias aplicado a los ejemplos propuestos.

Para la lectura de esta tesis, se requieren conocimientos de álgebra lineal y de geometría analítica por lo que se han incluido varios anexos y bibliografía recomendada de modo que el lector interesado puede consultarlos al final del presente trabajo.

El paquete estadístico computacional que se utilizó fue Xlstat V.4.

Capítulo I

Elementos básicos.

En este capítulo, se mencionan los elementos que se necesitarán para la comprensión del análisis factorial de correspondencias simple (AFC-simple).

Conceptos tales como espacio vectorial, distancia entre los elementos que lo conforman y la representación gráfica de los elementos dentro del espacio, permiten el manejo de la información contenida en la tabla de contingencia de la cual hace uso el AFC.

1.1 Tabla de contingencia de dos dimensiones

Las poblaciones se describen mediante medidas numéricas llamadas parámetros. Desafortunadamente la mayoría de las veces la obtención de estos parámetros resulta demasiado costosa y es de difícil obtención. Uno de los objetivos de la estadística es inferir acerca de alguna o varias características de la población con base en la información contenida en una muestra haciendo uso de estimadores para esos parámetros. Un estimador puede ser entendido como una fórmula para calcular un valor estimado o representativo de un parámetro de la población basada en las mediciones contenidas en una muestra.

Existe toda una metodología para obtener muestras representativas de una población, para los propósitos del presente trabajo supondremos que se cuenta con una muestra de tamaño T de una distribución multinomial.

En forma general podemos decir que una tabla de contingencia de $I \times J$ es un arreglo en el cual una muestra de T observaciones es clasificada con respecto a dos factores o variables cualitativas.

Una variable cualitativa es una transformación que va de un conjunto de unidades de estudio a un conjunto de valores.

X : Unidades de estudio \rightarrow *Valores*

De tal manera que si V es un subconjunto del conjunto de valores, entonces

i) Si V sólo tiene estructura de orden, la variable X es ordinal, es decir todas aquellas variables que establecen cierta jerarquía. Por ejemplo: grado de estudios, nivel económico, preferencias, etc.

ii) Si V no tiene ninguna estructura en particular, la variable X es nominal, es decir son aquellas a las que se les asigna un nombre o etiqueta. Por ejemplo: estado civil, religión que profesa una persona, tipo de enfermedad, sexo, etc.

A estos dos tipos de variables comúnmente se les conoce como variables cualitativas,

1

TESIS CON
FALLA DE ORIGEN

categorías o factores. En lo que sigue se hará referencia a las variables cualitativas como factores.

		Factor 2	
		j	
Factor 1	i	k_{ij}	$k_{i.} = \sum_j k_{ij}$
		$k_{.j} = \sum_i k_{ij}$	$T = \sum_{ij} k_{ij}$

Fig. 1 Tabla de contingencia de dos factores

En la figura 1 se muestra esquemáticamente la forma de la tabla de contingencia en la que el primer factor tiene I categorías y el segundo factor tiene J categorías que corresponden a los renglones y columnas respectivamente.

La frecuencia observada en la i-ésima categoría del factor I y la j-ésima categoría del factor J se define como k_{ij} . El término k_{ij} debe cumplir con lo siguiente: $k_{ij} \geq 0$ para $i=1,2,3,\dots,I$; $j=1,2,3,\dots,J$. Y, además para los renglones y columnas en la tabla se debe cumplir que $\sum_{j=1}^J k_{ij} > 0$ para $i=1,2,3,\dots,I$ y $\sum_{i=1}^I k_{ij} > 0$ para $j=1,2,3,\dots,J$. Lo que significa que debe de existir para cada columna o renglón al menos algún valor distinto de cero.

Si suponemos que la muestra de tamaño T se distribuye como una variable aleatoria multinomial. Es decir se tiene la probabilidad

$$P(e_{11} = k_{11}, e_{12} = k_{12}, \dots, e_{IJ} = k_{IJ}) = \frac{T! p_{11}^{k_{11}} p_{12}^{k_{12}} \dots p_{IJ}^{k_{IJ}}}{T! p_{1.}^{k_{1.}} p_{.j}^{k_{.j}}}$$

con

$$\sum_{i=1}^I \sum_{j=1}^J k_{ij} = T$$

$$\sum_{i=1}^I \sum_{j=1}^J p_{ij} = 1$$

Donde p_{ij} es la probabilidad de que de los T elementos de la población k_{ij} estén en la categoría i del factor I y la categoría j del factor J (celda (i,j)) de la tabla de contingencia. De tal manera que se tiene que estimar el valor de p_{ij} .

Mediante el método de máxima verosimilitud se sabe que un estimador de p_{ij} es $\hat{p}_{ij} = \frac{k_{ij}}{T}$ además por este mismo método un estimador para $p_{i.}$ es $\hat{p}_{i.} = \frac{k_{i.}}{T}$ y para $p_{.j}$ es $\hat{p}_{.j} = \frac{k_{.j}}{T}$ (Christensen, pág 41-43).

Como se verá más adelante estos estimadores son los que definen los pesos o masas de las

categorías de los factores y además definen a la matriz de correspondencias.

Dentro de la tabla de contingencia se tienen los siguientes conceptos: *Renglón marginal* o renglón de totales por columna: es aquel que se obtiene al sumar los elementos de cada una de las j columnas de la tabla de contingencia y cuyo j -ésimo término es $k_{.j} = \sum_{i=1}^I k_{ij}$

Columna marginal; de manera análoga se tiene la columna de totales por renglón la cual se obtiene al sumar los elementos de cada uno de los i renglones de la tabla de contingencia y cuyo i -ésimo término es $k_{i.} = \sum_{j=1}^J k_{ij}$. La suma de todos los elementos de la tabla de contingencia corresponde al tamaño de la muestra llamado también total o gran total que denotaremos como $T = \sum_{j=1}^J \sum_{i=1}^I k_{ij}$. Véase Fig 1.

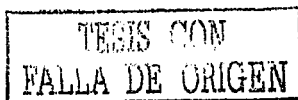
1.2 Conceptos geométricos y algebraicos obtenidos de la tabla de contingencia, necesarios para el desarrollo del AFC-simple.

A continuación se expone el desarrollo algebraico del AFC en el cual la información de las categorías de la tabla de contingencia es tratada como vectores en diferentes espacios ; con esto se obtiene la ventaja de la aplicación del algebra lineal y la geometría analítica.

Con lo mencionado hasta el momento podemos relacionar la tabla de contingencia con una matriz, que denominaremos como $K_{I \times J}$ que tiene la siguiente forma:

$$\begin{array}{c}
 \text{Factor 2} \\
 \begin{array}{cccccc}
 & 1 & 2 & \dots & j & \dots & J \\
 \text{Factor 1 } i & \left[\begin{array}{cccccc}
 k_{11} & k_{12} & \dots & k_{1j} & \dots & k_{1J} \\
 \dots & \dots & \dots & \dots & \dots & \dots \\
 k_{i1} & k_{i2} & \dots & k_{ij} & \dots & k_{iJ} \\
 \dots & \dots & \dots & \dots & \dots & \dots \\
 k_{I1} & k_{I2} & \dots & k_{Ij} & \dots & k_{IJ}
 \end{array} \right] & = & K_{I \times J} & (1)
 \end{array}
 \end{array}$$

En este trabajo se usarán las siguientes notaciones: la notación punto, en la cual los puntos indican las sumas sobre los índices particulares donde aparecen puntos; por ejemplo $k_{.j}$ indicará la suma de los elementos que forman el renglón j de la tabla de contingencia, ($k_{.j} = k_{1j} + k_{2j} + k_{3j} + \dots + k_{Ij} = \sum_{i=1}^I k_{ij}$), mientras que $k_{i.}$ indicará la suma de los elementos que forman la columna i de la tabla de contingencia ($k_{i.} = k_{i1} + k_{i2} + \dots + k_{iJ} = \sum_{j=1}^J k_{ij}$); notación suma y la notación desarrollada
 $k_{i.}$ Notación punto.
 $k_{i1} + k_{i2} + \dots + k_{iJ}$ Notación desarrollada.
 $\sum_{j=1}^J k_{ij}$ Notación suma.



De esta manera se puede pensar que la tabla de contingencia puede ser vista como una matriz que está formada por I vectores renglón, cada uno con J componentes; o, formada por J vectores columna, cada uno con I componentes. Es decir en una matriz se encuentran implícitos dos espacios vectoriales: R^I para los renglones y R^J para las columnas.

La ventaja de trabajar con matrices, es desde luego el uso matemático de los espacios vectoriales el cual permite desarrollar un proceso de análisis de datos para estudiar la asociación entre variables categóricas con los datos presentados en la tabla de contingencia.

Los términos de la matriz $K_{I \times J}$ representan las frecuencias absolutas que aparecen en la tabla de contingencia.

Se puede obtener la matriz de frecuencias relativas (proporciones) con respecto al total de observaciones T multiplicando a la matriz K por el recíproco de la suma total T (recordemos que el estimador máximo verosímil de p_{ij} es $\hat{p}_{ij} = \frac{k_{ij}}{T}$ al cual de ahora en adelante lo denotaremos por f_{ij}) La matriz F representa la probabilidad conjunta estimada para cada celda, es decir la probabilidad de que la información esté clasificada en la categoría i del factor 1 y la categoría j del factor 2

$$F_{I \times J} = \frac{1}{T} K_{I \times J}$$

$$F_{I \times J} = \begin{bmatrix} f_{11} & f_{12} & \dots & f_{1j} & \dots & f_{1J} \\ f_{21} & f_{22} & \dots & f_{2j} & \dots & f_{2J} \\ \dots & \dots & \dots & \dots & \dots & \dots \\ f_{i1} & f_{i2} & \dots & f_{ij} & \dots & f_{iJ} \\ \dots & \dots & \dots & \dots & \dots & \dots \\ f_{I1} & f_{I2} & \dots & f_{Ij} & \dots & f_{IJ} \end{bmatrix} \quad \text{con } f_{ij} = \hat{p}_{ij} = \frac{k_{ij}}{T} \quad (\text{II})$$

1.2.1 Peso o relevancia del renglón i , dado por los datos.

El peso del renglón i puede entenderse como la importancia que tiene el i -ésimo elemento del primer factor dentro del conjunto de los renglones que es determinado en términos de proporción con respecto a los demás renglones, y se calcula como el cociente entre la suma de los elementos del i -ésimo renglón de la matriz K ($k_{i\cdot}$) y el total (T),

$$\frac{k_{i\cdot}}{T} = \frac{\sum_{j=1}^J k_{ij}}{T}$$

Con la matriz F el peso del renglón i es $f_{i\cdot} = \sum_{j=1}^J f_{ij}$.

Así se tienen las siguientes equivalencias:

$$f_{i\cdot} = \sum_{j=1}^J f_{ij} = \frac{k_{i\cdot}}{T} = \frac{\sum_{j=1}^J k_{ij}}{T}$$

Estos i términos forman un vector, el vector de pesos por renglones:

$$f' = \begin{bmatrix} f_{1\cdot} \\ f_{2\cdot} \\ \dots \\ f_{i\cdot} \end{bmatrix} \quad \text{Nótese que } \sum_{i=1}^I f_{i\cdot} = 1 \quad (III)$$

1.2.2 Peso o relevancia de la columna j , dado por los datos.

De manera análoga, el peso de la columna j puede entenderse como la importancia que tiene el j -ésimo elemento del segundo factor dentro del conjunto de las columnas el cual es determinado por la tabla de contingencia en términos de proporción con respecto a las demás columnas, y se calcula como el cociente entre el total de la j -ésima columna de la matriz K y la suma total.

$$\frac{k_{\cdot j}}{T} = \frac{\sum_{i=1}^I k_{ij}}{T}$$

Con la matriz F el peso de la columna j es $f_{\cdot j} = \sum_{i=1}^I f_{ij}$.

Teniendo las siguientes equivalencias:

$$f_{\cdot j} = \sum_{i=1}^I f_{ij} = \frac{k_{\cdot j}}{T} = \frac{\sum_{i=1}^I k_{ij}}{T}$$

Estos j términos forman el vector $c_{j\cdot}$; el vector, el vector de pesos por columnas que denotaremos por:

$$c = [f_{\cdot 1} \quad f_{\cdot 2} \quad \dots \quad f_{\cdot p}] \quad (IV)$$

Observación: $\sum_{j=1}^J \sum_{i=1}^I \frac{k_{ij}}{T} = \sum_{i=1}^I \sum_{j=1}^J f_{ij} = \sum_{i=1}^I f_{i\cdot} = \sum_{j=1}^J f_{\cdot j} = 1$

1.2.3 Perfiles por renglón y por columna.

A los elementos del i -ésimo renglón de la matriz K dividido entre su total ($k_{i\cdot}$) se le conoce como perfil renglón.

$$\text{Esto es: } f_{ij} = \left[\frac{k_{ij}}{k_{i\cdot}}, \frac{k_{i2}}{k_{i\cdot}}, \dots, \frac{k_{iJ}}{k_{i\cdot}} \right] \quad (V)$$

Observación: La notación f_{ij} quiere decir que el i -ésimo perfil tiene J elementos, es decir es un vector del espacio vectorial R^J .

A la j -ésima columna de la matriz K dividida entre el total de la columna se le conoce como perfil columna.

$$\text{Esto es: } f_j = \begin{bmatrix} \frac{k_{1j}}{k_{.j}} \\ \frac{k_{2j}}{k_{.j}} \\ \dots \\ \frac{k_{ij}}{k_{.j}} \end{bmatrix} \quad (\text{VI})$$

Observación: La notación f_j quiere decir que el j -ésimo perfil tiene I elementos, es decir es un vector del espacio vectorial R^I .

1.2.4 Centroide.

Hasta ahora hemos caracterizado algunos elementos que se pueden obtener de la tabla de contingencia en cada uno de los dos espacios generados por los dos factores que intervienen. En la estadística descriptiva, por ejemplo, se habla de medidas de tendencia central dentro de las cuales el promedio es una de las más utilizadas y que en términos generales pone de manifiesto alrededor de que valor oscilan los valores de una muestra, es decir, es un centro de gravedad; nótese que cuando se obtiene este promedio cada elemento considerado tiene igual importancia dentro de la muestra es decir tienen masa uno. Hasta ahora, hemos trabajado con vectores, por lo que podemos preguntarnos ¿existirá un vector promedio alrededor del cual se encuentran los vectores renglón y columna de la matriz K , y que además tome en cuenta su peso o masa?

Por ejemplo tomemos el caso de los renglones o sea el espacio R^I . Como hemos visto tenemos

$$I \text{ vectores } \begin{bmatrix} f_1 \\ f_2 \\ \dots \\ f_I \end{bmatrix} \text{ y cada uno de estos vectores tiene por peso } \begin{pmatrix} f_{1.} \\ f_{2.} \\ \dots \\ f_{I.} \end{pmatrix} \text{ el centro de gravedad o}$$

centroide o vector promedio es un vector $\in R^I$ al cual lo denotaremos por G_{R^I} y es tal que:

$$G_{R^I} = \frac{\sum_{i=1}^I f_i f_{i.}}{\sum_{i=1}^I f_{i.}} = \sum_{i=1}^I f_i f_{i.} = (f_{.1}, f_{.2}, \dots, f_{.I}) \quad (\text{VII})$$

El centroide para el espacio de los renglones representa el vector promedio de los I perfiles por renglón ponderados.

De manera análoga se tiene que el centroide para el espacio de las columnas es:

$$G_{R'} = \frac{\sum_{j=1}^J f_{ij} f_j}{\sum_{j=1}^J f_j} = \sum_{j=1}^J f_{ij} f_j = (f_{i1}, f_{i2}, \dots, f_{iJ})' \quad (VIII)$$

El cual representa el vector promedio de los perfiles por columna ponderados

1.2.5 Matriz perfil por renglones y columnas.

Al hacer para cada uno de los renglones o columnas la división respectiva para obtener los perfiles resulta tedioso. Una solución es el manejo matricial de la información.

Denotaremos por $R_{I \times J}$ a la matriz perfil por renglones

En términos algebraicos tenemos que:

$$R_{I \times J} = D_r^{-1} F_{I \times J} = \begin{bmatrix} \frac{f_{11}}{f_{1.}} & \frac{f_{12}}{f_{1.}} & \dots & \frac{f_{1j}}{f_{1.}} & \dots & \frac{f_{1J}}{f_{1.}} \\ \frac{f_{21}}{f_{2.}} & \frac{f_{22}}{f_{2.}} & \dots & \frac{f_{2j}}{f_{2.}} & \dots & \frac{f_{2J}}{f_{2.}} \\ \dots & \dots & \dots & \dots & \dots & \dots \\ \frac{f_{i1}}{f_{i.}} & \frac{f_{i2}}{f_{i.}} & \dots & \frac{f_{ij}}{f_{i.}} & \dots & \frac{f_{iJ}}{f_{i.}} \\ \dots & \dots & \dots & \dots & \dots & \dots \\ \frac{f_{J1}}{f_{J.}} & \frac{f_{J2}}{f_{J.}} & \dots & \frac{f_{Jj}}{f_{J.}} & \dots & \frac{f_{JJ}}{f_{J.}} \end{bmatrix} = \begin{bmatrix} f_j \\ f_j \\ \dots \\ f_j \\ \dots \\ f_j \end{bmatrix} \quad (IX)$$

y cuyo término general es $\frac{k_{ij}}{k_{i.}} = \frac{f_{ij}}{f_{i.}}$

donde D_r^{-1} es la matriz diagonal asociada al vector de pesos por renglón es decir

$$D_r^{-1} = \begin{bmatrix} \frac{1}{f_{1.}} & 0 & 0 & 0 \\ 0 & \frac{1}{f_{2.}} & 0 & 0 \\ 0 & 0 & \dots & 0 \\ 0 & 0 & 0 & \frac{1}{f_{i.}} \end{bmatrix}$$

“Lo que se logra al construir la matriz $R_{I \times J}$ es expresar las probabilidades condicionales de la forma:

$\Pr(\text{categoría del factor (2)} | \text{categoría del factor(1)})$ ”.

De manera análoga denotaremos por $C_{I \times J}$ la matriz perfil por columnas.

$$C_{I \times J} = F_{I \times J} D_c^{-1}$$

$$= \begin{bmatrix} \frac{f_{11}}{f_{.1}} & \frac{f_{12}}{f_{.2}} & \dots & \frac{f_{1j}}{f_{.j}} & \dots & \frac{f_{1I}}{f_{.I}} \\ \frac{f_{21}}{f_{.1}} & \frac{f_{22}}{f_{.2}} & \dots & \frac{f_{2j}}{f_{.j}} & \dots & \frac{f_{2I}}{f_{.I}} \\ \dots & \dots & \dots & \dots & \dots & \dots \\ \frac{f_{j1}}{f_{.1}} & \frac{f_{j2}}{f_{.2}} & \dots & \frac{f_{jj}}{f_{.j}} & \dots & \frac{f_{jI}}{f_{.I}} \\ \dots & \dots & \dots & \dots & \dots & \dots \\ \frac{f_{I1}}{f_{.1}} & \frac{f_{I2}}{f_{.2}} & \dots & \frac{f_{Ij}}{f_{.j}} & \dots & \frac{f_{II}}{f_{.I}} \end{bmatrix} = [f_{.1} \ f_{.2} \ \dots \ f_{.j} \ \dots \ f_{.I}] \quad (X)$$

cuyo término general es $\frac{k_{ij}}{k_{.j}} = \frac{f_{ij}}{f_{.j}}$

donde $D_{.j}^{-1}$ es la matriz diagonal asociada al vector de sumas por columna de F es decir.

$$D_{.j}^{-1} = \begin{bmatrix} \frac{1}{f_{.1}} & 0 & 0 & 0 \\ 0 & \frac{1}{f_{.2}} & 0 & 0 \\ 0 & 0 & \dots & 0 \\ 0 & 0 & 0 & \frac{1}{f_{.j}} \end{bmatrix}$$

De manera análoga las celdas que forman la matriz C representan la probabilidad condicional de la forma:

Pr(categoría del factor(1) | categoría del factor(2))"

Estas probabilidades condicionales nos van a permitir presentar de una manera intuitiva el concepto de asociación, y nos permitirá introducir el análisis de correspondencias.

Ya que la matriz de datos puede verse como un arreglo de vectores por renglones o un arreglo de vectores por columnas, entonces quedan determinados dos espacios vectoriales R^j y R^I . Dentro de esos espacios vectoriales los renglones en R^j y las columnas en R^I , permiten determinar la matriz perfil por renglones y la matriz perfil por columnas respectivamente. Cada uno de estos renglones o columnas tienen asignado un peso o masa en función de la frecuencia de los datos.

El hecho de que se trabaje con vectores que pertenezcan a un espacio vectorial en el cual se puedan medir distancias entre vectores nos permite por ejemplo tratar de encontrar vectores parecidos en términos de sus proximidades, es decir podemos trabajar en un espacio métrico, por lo que tenemos una manera de medir distancias: métrica. Usualmente se ha trabajado en un espacio euclidiano en el cual se tiene la distancia usual:

Un espacio euclidiano por ejemplo de n dimensiones es el conjunto de todos los vectores con n componentes para los cuales las operaciones de suma vectorial y multiplicación por escalar son permitidas, y además, para cualquiera dos vectores en el espacio, existe un número no negativo

llamado la distancia Euclidiana (métrica). La distancia Euclidiana para \mathbf{a} y $\mathbf{b} \in \mathbb{R}^n$ está definida como:

$$d^2(\mathbf{a}, \mathbf{b}) = \|\mathbf{a} - \mathbf{b}\|^2 = (a_1 - b_1)^2 + (a_2 - b_2)^2 + \dots + (a_n - b_n)^2$$

De manera alternativa podemos definir $\|\mathbf{a} - \mathbf{b}\|^2$ en términos del producto escalar de $(\mathbf{a} - \mathbf{b})$ el cual es

$$d^2(\mathbf{a}, \mathbf{b}) = \|\mathbf{a} - \mathbf{b}\|^2 = \langle (\mathbf{a} - \mathbf{b}), (\mathbf{a} - \mathbf{b}) \rangle = [(\mathbf{a} - \mathbf{b})'(\mathbf{a} - \mathbf{b})] \text{ donde } (\mathbf{a} - \mathbf{b}) \text{ es el vector de diferencia.}$$

Nótese que esta métrica no hace diferencia con respecto a la importancia que tienen los vectores, es decir todos tienen la misma masa. Sin embargo, en el análisis de datos, los métodos deben ser adaptados a situaciones prácticas donde los individuos o las categorías no necesariamente tengan la misma "importancia". Así, por ejemplo si lo que se quiere es determinar la proximidad o cercanía entre elementos que se consideren tengan el mismo peso o sean igualmente importantes se podría pensar en utilizar la distancia euclidiana. Pero, también se da el caso en que el investigador pueda privilegiar a algunos elementos más que otros en función de su importancia con respecto a algún criterio o que la estructura interna de los datos así lo establezca, de tal manera que cada elemento tiene asociado una masa f_w tal que $f_w > 0$ y $\sum_w f_w = 1$. Así, si se quiere determinar la proximidad o cercanía entre estos elementos ponderados, entonces tendremos que utilizar alguna métrica que tome en cuenta ese peso o masa; una métrica que satisficiera lo anterior es la llamada: métrica Ji cuadrada, la cual define un espacio métrico ponderado.

1.2.6 Distancia Ji cuadrada

En un espacio ponderado en el que a los vectores de este espacio se les asigna una masa o un peso. Algunos autores denominan a este espacio como un espacio euclidiano ponderado en el cual se define el producto punto o producto escalar de la siguiente forma: para dos vectores \mathbf{y} , \mathbf{w} , con masa q_1 y q_2 respectivamente tenemos:

$$\langle \mathbf{y}, \mathbf{w} \rangle = \mathbf{y}' D_q \mathbf{w} = \sum_j q_j y_j w_j, \text{ Donde } D_q \text{ es la matriz diagonal asociada a los pesos de los puntos.}$$

De tal manera que la distancia al cuadrado de los dos puntos en este espacio euclidiano ponderado es:

$$d^2(\mathbf{y}, \mathbf{w}) = \langle \mathbf{y} - \mathbf{w}, \mathbf{y} - \mathbf{w} \rangle = (\mathbf{y} - \mathbf{w})' D_q (\mathbf{y} - \mathbf{w}) = \sum_j q_j (y_j - w_j)^2.$$

En particular, cuando se trabaja con una tabla de contingencia se tiene la llamada distancia Ji cuadrada, la cual está definida para dos perfiles por renglones i e i' como



$$d^2(i, i') = \sum_{j=1}^J \frac{1}{f_{j'}} \left(\frac{f_{ij}}{f_{i'}} - \frac{f_{ij'}}{f_{j'}} \right)^2 \quad \text{y para perfiles por columna } j \text{ y } j' \text{ como:}$$

$$d^2(j, j') = \sum_{i=1}^I \frac{1}{f_{i'}} \left(\frac{f_{ij}}{f_{j'}} - \frac{f_{ij'}}{f_{j'}} \right)^2$$

Esta distancia difiere de la distancia usual Euclidiana en que cada término al cuadrado está ponderado por el inverso del peso correspondiente a cada término del otro espacio.

Hasta el momento lo que se ha hecho es lo siguiente:

1. Se ha visto a la tabla de contingencia como una matriz, la cual define de manera natural el espacio de los renglones y el espacio de las columnas R^I y R^J respectivamente.

2. Se han obtenido los elementos por perfiles dependiendo en que espacio estemos trabajando. A cada uno de estos perfiles se les ha asignado una masa que está en relación a su importancia dentro del espacio definido por los elementos que conforman la tabla de contingencia.

3. Se ha encontrado el centroide de cada uno de los espacios definidos en la tabla de contingencia.

5. La masa que tienen los perfiles establecen la necesidad de trabajar en un espacio métrico ponderado, donde la métrica es la distancia ji-cuadrada.

Ejemplo 1 : Relación entre los espacios R^I y R^J con los perfiles por columna y por renglón

Tabla 1.1. Tabla de contingencia (K): número de muebles producidos con defectos clasificados de acuerdo a tipo de defecto y turno de producción

Turno de producción	Tipo de defecto				Total
	A	B	C	D	
1	15	21	45	13	94
2	23	31	34	5	93
3	33	17	49	20	119
Total	74	69	128	38	309

← Matriz de frecuencias absolutas K_{ij}

Tabla 1.2 Tabla de frecuencias relativas (F): porcentaje de muebles producidos con defectos clasificados de acuerdo a tipo de defecto y turno de producción

Turno de producción	Tipo de defecto				Total
	A	B	C	D	
1	0.049	0.038	0.146	0.042	0.304
2	0.034	0.1	0.11	0.016	0.311
3	0.107	0.025	0.139	0.035	0.305
Total	0.239	0.223	0.414	0.123	1

← Matriz de correspondencias

Cuadro 1: Relación entre los espacios R^I y R^J con los perfiles por columna y por renglón

Espacio R^I para turno de producción, hay 3 renglones son los "perfiles por renglón" (puntos en R^I)
 - Las coordenadas del renglón $i=1,2,3$ son $\frac{1}{i}$
 perfil para el renglón $i: (\frac{1}{i}, \frac{1}{i}, \dots, \frac{1}{i})$ son

Turno de producción	Perfiles por turno (renglones)			
	Tipo de defecto			
	A	B	C	D
1	0.15957	0.2234	0.47972	0.1383
2	0.27083	0.32992	0.35417	0.05208
3	0.27731	0.14266	0.41175	0.12827

- El peso del i -ésimo renglón es $F_i: \begin{bmatrix} 0.304 \\ 0.311 \\ 0.385 \end{bmatrix}$

- Los I perfiles por renglón corresponden a los I renglones de la matriz: $D^I \cdot F$
 - La distancia entre dos renglones i e i' es:

$$d^2(i, i') = \sum_{j=1}^4 F_j \left[\frac{1}{i} - \frac{1}{i'} \right]^2$$

$$d^2(1, 2) = \sum_{j=1}^4 F_j \left[\frac{1}{1} - \frac{1}{2} \right]^2 = 0.194$$

Adicionalmente a lo revisado se verá en el capítulo siguiente:

- El criterio de ajuste está determinado por una forma cuadrática en la cual intervienen los pesos de los perfiles los cuales aparecen en la matriz D^I .
- La matriz de distancias entre renglones está determinada por una forma cuadrática en la cual interviene D^I .

Espacio R^J para tipo de defecto, hay 4 columnas son los "perfiles por columna" (puntos en R^J)
 - Las coordenadas de la columna $j=1,2,3,4$, son $\frac{1}{j}$
 y los perfiles por columna $j: (\frac{1}{j}, \frac{1}{j}, \dots, \frac{1}{j})$ son

Turno de producción	Perfiles por tipo de defecto (columnas)			
	Tipo de defecto			
	A	B	C	D
1	1.2027	0.30435	0.35156	0.34211
2	0.35136	0.44928	0.25563	0.13158
3	0.44895	0.24638	0.24791	0.52522

- El peso de la j -ésima columna es $F_j: \begin{bmatrix} 0.239 \\ 0.223 \\ 0.414 \\ 0.123 \end{bmatrix}$

- Los J perfiles por columna corresponden a las J columnas de la matriz $F \cdot D^J$
 - La distancia entre dos columnas j y j' es:

$$d^2(j, j') = \sum_{i=1}^3 F_i \left[\frac{1}{j} - \frac{1}{j'} \right]^2$$

$$d^2(A, B) = \sum_{i=1}^3 F_i \left[\frac{1}{1} - \frac{1}{2} \right]^2 = 0.168$$

Adicionalmente a lo revisado se verá en el capítulo siguiente:

- El criterio de ajuste está determinado por una forma cuadrática en la cual intervienen los pesos de los perfiles los cuales aparecen en la matriz D^J .
- La matriz de distancias entre columnas está determinada por la forma cuadrática asociada con: D^J

En el siguiente capítulo mencionaremos lo que es asociación entre las categorías haciendo referencia a la prueba Ji-cuadrada de Pearson para independencia y mostraremos el desarrollo del AFC-simple para representar gráficamente de manera óptima la asociación entre los perfiles renglones y columnas a través de la proximidad entre ellos ya sea entre un mismo espacio o entre ambos espacios mediante un escalamiento.

Capítulo II

Análisis factorial de correspondencias simple.

El análisis factorial de correspondencias (AFC), desarrollado por un grupo de estadísticos franceses desde principios de 1960, se reporta que es teóricamente equivalente a casos especiales de algunas técnicas clásicas que han aparecido en diferentes contextos desde la mitad de los años 30's, como por ejemplo el análisis de correlación canónico y el análisis de componentes principales. El nombre de análisis factorial tiene dos acepciones la primera referida al estudio de la asociación entre factores o variables cualitativas y el segundo a la técnica del análisis multivariado llamada análisis de factores en la que se busca resumir en menos dimensiones la información contenida en una matriz de datos.

El AFC, es una herramienta estadística exploratoria para estudiar la asociación entre las categorías de factores presentadas en una tabla de contingencia, en este sentido puede ser utilizado como una técnica complementaria a otros métodos estadísticos como los modelos log-lineales.

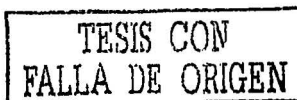
El AFC determina gráficamente la posición de una serie de unidades de una población bajo estudio, que pueden referirse a objetos o sujetos (segmentos del mercado, grupos de individuos o personas físicas, sectores, productos, etc.) con relación a una serie de factores de interés (atributos, características, escalas de valoración, etc.) clasificados por dos factores de interés y de acuerdo a la proximidad que tienen las categorías dentro de cada factor y entre factores, reporta o refleja en dos o en a lo mucho tres dimensiones la estructura de asociación tanto de las categorías dentro de cada factor como entre las categorías de los factores, determinando además un escalamiento de las categorías de los factores.

Para introducir las ideas referentes al AFC-S utilizaremos el siguiente ejemplo donde se estudia la asociación entre dos factores:

Supongamos que se quiere clasificar los defectos encontrados en los muebles producidos en una fábrica, de acuerdo a 1) el tipo de defecto y 2) el turno en que fue producido el mueble. De tal manera que un total de 309 muebles defectuosos fueron registrados y clasificados de acuerdo a uno de los siguientes tipos de defecto: A, B, C o D. Al mismo tiempo, cada mueble fue identificado de acuerdo al turno de producción en el cual fue fabricado. La información es presentada en la siguiente tabla de contingencia:

Tabla 2.1 Tabla de contingencia número de muebles con defectos
Por tipo de defecto y turno de producción

Turno de producción	Tipo de defecto				Total
	A	B	C	D	
1	15	21	45	13	94
2	26	31	34	5	96
3	33	17	49	20	119
Total	74	69	128	39	309



2.1 Definición de asociación o dependencia para una tabla de contingencia de $I \times J$.

De la Tabla 2.1, para el caso del espacio R^2 -que es el espacio de los renglones- se refiere al turno en que fueron fabricados los muebles, de tal forma que para cada turno se tiene un vector con cuatro entradas es decir nos referimos al espacio R^4 , para el caso del espacio R^2 -que es el espacio de las columnas- se refiere al tipo de defecto que tienen los muebles, de tal manera que para cada tipo de defecto se tiene un vector con tres entradas es decir nos referimos al espacio R^3 .

La probabilidad de que los elementos de una población pertenezcan a una categoría j del factor tipo de defecto, sea mayor o menor que en una categoría i del factor turno que en otra categoría i' del mismo factor se representa como:

$$\Pr(x \in j | x \in i) \neq \Pr(x \in j | x \in i') \quad \text{para alguna } i \neq i'$$

Si esta desigualdad se cumple para alguna categoría j del factor tipo de defecto, entonces existe una asociación entre los factores turno de producción y tipo de defecto. Esta ecuación indica que la probabilidad de que un mueble x pertenezca a la categoría j del factor tipo de defecto depende de la categoría del factor turno que se considere.

Como una explicación hipotética, en el caso de que se presentara la asociación pudiera ser que si pensamos que los turnos de fabricación corresponden al matutino, vespertino y nocturno; puede darse el caso por ejemplo que para trabajadores que no estén muy habituados a trabajar en la noche su capacidad productiva se ve disminuida por efecto del sueño, de tal manera que se esperaría un mayor número de muebles defectuosos en el turno nocturno presentándose de esta manera la posible asociación.

Para hablar del término de independencia se tiene lo siguiente:

$$\text{Si } \Pr(x \in j | x \in i) = \Pr(x \in j | x \in i') = \Pr(x \in j) \quad i=1,2,3 \text{ con } i \neq i'$$

Es decir, si estas igualdades se cumplen para toda categoría j del factor tipo de defecto $j=1,2,3,4$ entonces las variables no están asociadas. Esto quiere decir que la probabilidad de que un mueble x pertenezca a la categoría j del factor tipo de defecto es independiente de las categorías del factor turno de producción.

Existen métodos para probar la independencia entre las categorías de los factores, dentro de los cuales la prueba Ji-cuadrada es una de las más usadas. Aplicaremos esta prueba estadística a los datos de la tabla de contingencia 2.1

2.2 Prueba de independencia χ^2 para una tabla de contingencia de $I \times J$.

Datos: Una muestra aleatoria de tamaño T es obtenida. Estas observaciones son clasificadas de acuerdo a dos criterios (factores): Para el ejemplo de la tabla 2.1 tenemos que $T=309$.

k_{ij} representa la frecuencia de muebles en la categoría i del primer factor turno de producción y la categoría j del segundo factor tipo de defecto.

Supuestos:

1. La muestra T es aleatoria.
2. Cada observación es clasificada en exactamente una de las I categorías del primer factor y exactamente en una de las J categorías del segundo factor.

Hipótesis:

H_0 : La aparición del defecto j en los muebles producidos, es independiente del turno de producción i para todo $i = 1, 2, 3$ y $j = 1, 2, 3, 4$.

Es decir

$$H_0 : p_{ij} = p_{i.}p_{.j} \quad \forall i, j \quad \text{Vs} \quad H_1 : p_{ij} \neq p_{i.}p_{.j} \quad \text{para alguna } i, j$$

Estadístico de prueba:

Sea $E_{ij} = \frac{k_{ij}T}{T}$ el valor esperado para la celda ij suponiendo cierta la hipótesis H_0 de independencia. El estadístico está dado por: $\chi^2_c = \sum_{i=1}^I \sum_{j=1}^J \frac{(k_{ij} - E_{ij})^2}{E_{ij}}$.

Para los datos de la tabla 2.1 $\chi^2_c = 19.177$

Distribución nula del estadístico:

La distribución nula del estadístico está dada aproximadamente por la distribución Ji-cuadrada con $(I-1)(J-1)$ grados de libertad ($\chi^2_{((I-1)(J-1))}$). La aproximación a la Ji-cuadrada es satisfactoria si todas las E_{ij} son mayores o iguales que .5 y que al menos la mitad son más grandes que 1.

Para el ejemplo se tiene que $\chi^2_c \sim \chi^2_{(2)(3)} = \chi^2_{(6)}$.

Regla de decisión:

Rachazar H_0 si $\chi^2_c > \chi^2_{((1-\alpha), (I-1)(J-1))}$, donde $\chi^2_{((1-\alpha), (I-1)(J-1))}$ es el cuantil de orden $1-\alpha$ de una distribución χ^2 con $((I-1)(J-1))$ grados de libertad.

Por lo tanto para el ejemplo se tiene que para un nivel de significancia $\alpha = 0.05$:

$\chi_{(6)}^{(95)} = 12.591$, de tal manera que se rechaza la hipótesis H_0 ya que $\chi^2 = 19.177 > \chi_{(95,6)}^2 = 12.591$

Por lo que se concluye que existe asociación entre el turno de fabricación del mueble y el tipo de defecto que presenta.

Una pregunta que quizás interese a la persona encargada del control de calidad de la fábrica es el poder determinar cuales categorías de los factores son los que están asociados; ya que la prueba Ji-cuadrada únicamente establece si existe asociación o no, pero no especifica cuales son las categorías causantes de ésta.

Un camino para descubrir que factores causan la asociación entre los factores es haciendo uso de la expresión que nos proporcionaba la definición de asociación:

$$\Pr(x \in j | x \in i) \neq \Pr(x \in j | x \in i') \quad \text{para alguna } i \neq i' \quad (1)$$

la cual tiene dos posibilidades:

$$\Pr(x \in j | x \in i) > \Pr(x \in j | x \in i') \quad \text{para alguna } i \neq i' \quad (2)$$

ó

$$\Pr(x \in j | x \in i) < \Pr(x \in j | x \in i') \quad \text{para alguna } i \neq i' \quad (3)$$

De manera intuitiva observamos de 2 que es más probable que la categoría j del segundo factor ocurra con la categoría i del primer factor que con la categoría i' .

De manera análoga de 3 observamos que es más probable que la categoría j del segundo factor ocurra con la categoría i' del primer factor que con la categoría i .

Ahora bien, si lográramos representar estas posibles asociaciones en un gráfico que muestre las similitudes entre las categorías, esperaríamos para 2 que la proximidad de la categoría j a la categoría i , sea menor que la proximidad de la categoría j a la categoría i' , dado que es más probable que lo primero ocurra.

Como se recordará la matriz perfil por renglones y por columnas definidas en el capítulo uno - tablas 2.2 y 2.3 -, contine las probabilidades condicionales

Tabla 2.2. Perfiles por turno (renglones)

Turno de producción	Tipo de defecto				Total
	A	B	C	D	
1	0.1887	0.2234	0.4782	0.1093	1
2	0.2703	0.3232	0.3517	0.0548	1
3	0.2731	0.1485	0.4178	0.1607	1

Tabla 2.3. Perfiles por tipo de defecto (columnas)

Turno de producción	Tipo de defecto			
	A	B	C	D
1	0.2027	0.30435	0.35158	0.34211
2	0.35135	0.44928	0.26563	0.13158
3	0.44565	0.24538	0.38281	0.52632
Total	1	1	1	1

Por ejemplo:

1) $\Pr(x \in C \mid x \in 3)$ es la $\Pr(\text{los muebles tengan el defecto } C \text{ cuando se fabricaron en el tercer turno}) = 0.41176$

2) $\Pr(x \in 3 \mid x \in C)$ es la $\Pr(\text{los muebles se hayan fabricado en el tercer turno cuando presentan el defecto } C) = 0.38281$

3) $\Pr(x \in A \mid x \in 2) = 0.27083$

4) $\Pr(x \in 2 \mid x \in A) = 0.35135$

Una interpretación gráfica de la información presentada en las tablas 2.2 y 2.3, es la siguiente:

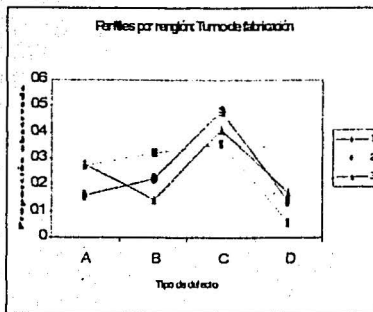


Fig 2.1

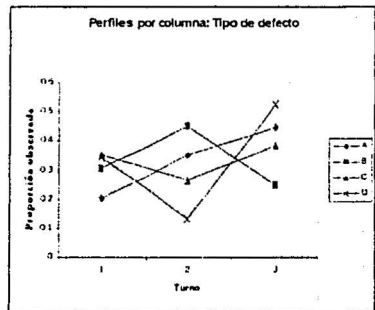


Fig 2.2

Así, en la figura 2.1 se tiene que el tipo de defecto A es más probable de ocurrir en las turnos 2 y 3 que en el primer turno. Para el defecto B se observa que en el segundo turno es más probable que este ocurra que en el tercer turno. Para los tipos de defectos C y D es casi igual de probable que ocurran en todos los turnos.

Si nos fijamos en la figura 2.2 vemos que en el turno 1 es más probable la ocurrencia de los defectos B, C y D que el defecto A. Para el turno 2, es más probable que ocurra el defecto B que el defecto D y para el turno 3 es más probable que ocurra el defecto D que el defecto B.

Estas proporciones observadas nos dan una idea acerca de cómo es la asociación; para conseguir una mejor apreciación de ésta se construyen las tablas 2.4 y 2.5 que representan las desviaciones de los perfiles renglón y columna respectivamente con respecto al centroide del espacio vectorial respectivo.

Tabla 2.4. Desviación de los perfiles renglón (Turno de producción) con respecto al centroide.

Turno de producción	Tipo de defecto			
	A	B	C	D
1	-0.07991	0.0001	0.06448	0.01522
2	0.03135	0.06662	-0.06007	-0.07089
3	0.03783	-0.08044	-0.00247	0.04509

Tabla 2.5. Desviación de los perfiles columna (Tipo de defecto) con respecto al centroide.

Turno de producción	Tipo de defecto			
	A	B	C	D
1	-0.1015	0.00014	0.04736	0.0379
2	0.04057	0.1386	-0.04505	-0.1791
3	0.06083	-0.13874	-0.0023	0.1412

Una presentación gráfica de la información contenida en las tablas 2.4 y 2.5 se presenta en las figuras 2.3 y 2.4:

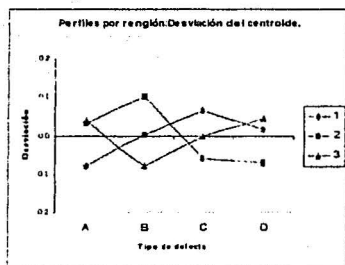


Fig 2.3

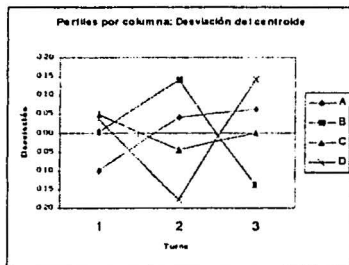


Fig 2.4

De la figura 2.3 interpretamos que el defecto A se encuentra más alejado en un sentido negativo del centroide para el turno 1 en el espacio de los renglones (R^4), lo cual nos habla de una asociación entre el turno 1 y el defecto A. Para el caso del defecto B, tenemos que el turno 2 y 3 se encuentran alejados del centroide en un sentido positivo y negativo respectivamente lo cual nos hace suponer una asociación entre el tipo de defecto B y los turnos 2 y 3. El defecto C por otra parte parece que tiene una asociación no tan marcada con respecto a los turnos 1 y 2, mientras que con el tercer turno no parece existir asociación. Por último, con respecto al defecto D se tiene que presenta una posible asociación con respecto al turno 2.

De la figura 2.4 se tiene que el turno 1 está más asociado con respecto al tipo de defecto A en un sentido negativo, mientras que existe una asociación en un sentido positivo entre el defecto B y el turno 2 y una asociación en un sentido negativo entre el defecto D y el turno 2. Se observa que, además existe una asociación en un sentido positivo entre el turno 3 y el defecto D y una asociación en un sentido negativo entre el turno 3 y el defecto B.

La forma de describir la asociación presentada arriba es un mejor acercamiento para conocerla, esto en el sentido de que únicamente se están describiendo los perfiles tal cual se obtienen de la tabla de contingencia, pero sin tomar en cuenta el peso que cada uno de ellos tiene en su respectivo

espacio, para lograr introducir este elemento en el análisis es necesario ponderar cada uno de los elementos. Recuérdese que precisamente en el capítulo I se hizo referencia a un espacio vectorial ponderado.

Así, el problema de descubrir la asociación de manera gráfica entre las categorías de los factores se puede tratar de la siguiente manera:

Para cada uno de los elementos que forman la matriz perfil por renglón y por columnas que son vectores en los espacios R' y R'' respectivamente, se pretende encontrar una representación gráfica en un espacio de menor dimensión prefiriéndose una línea recta, enseguida un plano o bien un hiperplano, en el cual se tomen en cuenta la ponderación de los elementos y que presente las proximidades entre elementos en cada espacio y de manera simultánea se consideren los dos espacios para tener una representación óptima. Categorías muy próximas entre sí, con un alejamiento considerable del centroide implicarán una asociación fuerte entre ellas.

Como se puede observar el problema arriba presentado equivale a encontrar vectores que generen el subespacio óptimo buscado.

2.3 Representación gráfica de la asociación entre factores con el AFC.

En el capítulo anterior se mencionó que los elementos de la tabla de contingencia ya sea por renglón o por columna podían ser vistos como vectores. Ahora bien, esos vectores pueden ser representados en un sistema de coordenadas, a través de unos vectores que forman una base para el espacio vectorial. Si pudiéramos ver esa representación veríamos esos vectores como una nube de puntos.

Dentro de esa nube se esperaría que aquellos vectores parecidos estarán más próximos y aquellos vectores que sean distintos estarán más alejados. Por lo que el problema de encontrar un subespacio óptimo para esos vectores, será el encontrar un nuevo sistema de referencia -base- que conserve la estructura de la posición de los vectores y que además sea posible representarlos visualmente. Para encontrar esos vectores que forman la base que buscamos existen algunos métodos de ajuste como el de mínimos cuadrados.

En las siguientes secciones se muestra de manera gráfica la idea de este ajuste para el AFC y también se muestra la parte algebraica.

2.3.1. Formulación gráfica de ajuste de un subespacio vectorial para dos factores de la tabla de contingencia de $I \times J$.

Se tienen dos nubes de puntos $N(I)$ y $N(J)$

Un elemento de la nube $N(I)$ es un par formado por un perfil renglón y su peso respectivo, es decir:

$N(I) = \{f_i, f_{i^*}\}$, donde f_i representa al perfil renglón i-ésimo

Tabla 2.6. Nube de puntos $N(I)$: Perfiles por renglón

Turno	A	B	C	D	Peso
1	0.160	0.223	0.479	0.138	0.304
2	0.271	0.323	0.354	0.052	0.311
3	0.277	0.143	0.412	0.168	0.385

Un elemento de la nube $N(J)$ es un par formado por un perfil columna y su masa respectiva, es decir:

$N(J) = \{f_j, f_{j^*}\}$, donde f_j representa al perfil columna j-ésimo.

Tabla 2.7. Nube de puntos $N(J)$: perfiles por columna.

Turno	A	B	C	D
1	0.203	0.304	0.352	0.342
2	0.351	0.449	0.266	0.132
3	0.446	0.246	0.383	0.526
Peso	0.239	0.223	0.414	0.123

Gráficamente la nube de puntos $N(J)$ se ve así:

Nube de puntos $N(J)$: Perfiles por columna

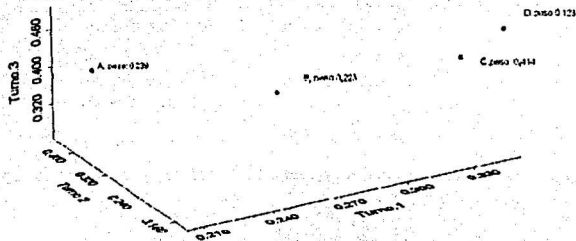


Fig. 2.5

Para estos dos conjuntos buscamos una representación que esté en un subespacio de la menor dimensión posible el cual puede ser R o R^2 y además que sea lo más confiable en el sentido de que se conserve la relación de asociación entre las categorías.

Para lograr esto se hará la proyección de los elementos de la nube de puntos en un espacio lineal S que puede ser una línea, un plano o un hiperplano, etc.; el cual contendrá al centro de gravedad G , si se consigue que la distancia de los puntos a su proyección sea mínima entonces el subespacio S será la aproximación a la representación deseada.

La razón por la que el centroide debe estar contenido en el subespacio, es porque éste es el punto más cercano a la nube de puntos en el sentido de minimiza la suma de cuadrados de las distancias entre los puntos de la nube y dicho punto. Por lo que la búsqueda de subespacios óptimos se restringirá a aquellos que contengan al centroide.

Para encontrar el subespacio de menor dimensión se tienen dos alternativas:

1. Se define una función llamada de cercanía entre la nube de puntos $N(I)$ o $N(J)$ con el subespacio S .

Por ejemplo para la nube de puntos $N(I)$

$$\Psi(S, N(I)) = \sum_{i=1}^I f_i \cdot d^2(f_i, S)$$

donde

$d^2(f_i, S)$ es la distancia mínima al cuadrado entre el perfil f_i y el subespacio S .

f_i es el peso o masa del perfil i -ésimo.

Esta función establece que el subespacio S es el que mejor ajusta a los perfiles en términos de minimizar la distancia que existe entre estos y el subespacio, por lo que el problema se reduce a: minimizar la suma de las distancias elevadas al cuadrado ponderadas de los elementos de la nube de puntos $N(I)$ al subespacio S .

2. La otra forma es maximizar la suma de distancias ponderadas del origen (el centroide) a las proyecciones de los puntos de la nube sobre el subespacio S .

Desde luego este subespacio S es generado por algunos vectores por lo que a fin de cuentas lo que importa es encontrar los vectores que lo generan.

Para presentar gráficamente esta idea y sin perder generalidad trabajemos con la nube de puntos por renglón $N(I)$ y además supongamos que los elementos de la nube tienen pesos unitarios. Tomemos un perfil cualquiera f_j de la nube, y proyectémoslo de manera ortogonal sobre el subespacio S generado por el vector u -hasta ahora desconocido- denotemos a esta proyección como $Pr_{oy_u} f_j$.

Observemos los segmentos $\overline{Gf_j}$, $\overline{GPr_{oy_u} f_j}$, $\overline{f_j Pr_{oy_u} f_j}$. Véase figura 2.6.

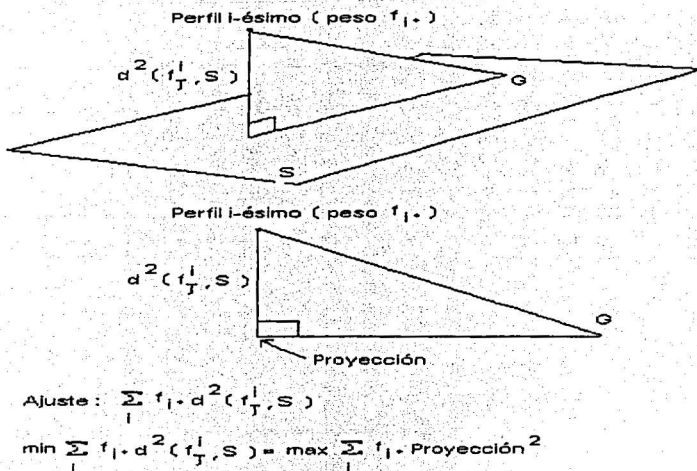


Fig 2.6

Lo que se quiere es minimizar la longitud o distancia del segmento $\overline{f_j \text{Proy}_u f_j}$, para todos los I elementos de la nube, de tal manera que la suma de las distancias de todos los puntos con respecto al subespacio generado por el vector u sea mínima. Tomemos la distancia al cuadrado $\overline{f_j \text{Proy}_u f_j}^2$. Aplicando el Teorema de Pitágoras tenemos que $\overline{f_j \text{Proy}_u f_j}^2 + \overline{G \text{Proy}_u f_j}^2 = \overline{G f_j}^2$, sumando sobre los I elementos $\sum_{j=1}^I \overline{f_j \text{Proy}_u f_j}^2 + \sum_{j=1}^I \overline{G \text{Proy}_u f_j}^2 = \sum_{j=1}^I \overline{G f_j}^2$, como el problema es minimizar $\sum_{j=1}^I \overline{f_j \text{Proy}_u f_j}^2$, esto implica que $\sum_{j=1}^I \overline{f_j \text{Proy}_u f_j}^2 = \sum_{j=1}^I \overline{G f_j}^2 - \sum_{j=1}^I \overline{G \text{Proy}_u f_j}^2$.

Nótese que $\sum_{j=1}^I \overline{f_j \text{Proy}_u f_j}^2 = \sum_{j=1}^I f_n d^2(f_j, S)$ para $f_n = 1 \quad i=1,2,3,\dots,I$ y para el subespacio generado por el vector u . Lo que implica que minimizar la suma de las distancias al cuadrado entre los perfiles y el subespacio generado por el vector u es equivalente a maximizar $\sum_{j=1}^I \overline{G \text{Proy}_u f_j}^2$, que es la suma de las proyecciones al cuadrado sobre el subespacio generado por el vector u (ver Fig. 2.6)

Para encontrar la solución al problema de encontrar el mejor subespacio que ajuste a los puntos algunos autores hacen uso de la técnica de descomposición en valores singulares (SVD) y otros lo plantean como un problema de optimización cuya solución corresponde a encontrar los eigenvalores y eigenvectores de la matriz que representa la suma ponderada de cuadrados de las proyecciones de los puntos de la nube. En este trabajo se presenta la segunda forma.

2.3.2 Presentación algebraica del análisis de correspondencias.

Algunos resultados importantes que se refieren al subespacio óptimo S son los siguientes:

1. El centroide de la nube de puntos $N(I)$, es el punto que minimiza la función que describe la distancia entre la nube de puntos $N(I)$ y el subespacio S.
2. El subespacio óptimo S debe contener al centroide.

Para garantizar que el centroide esté contenido en el subespacio S se lleva a cabo una traslación a un nuevo sistema de referencia en el cual el origen sea precisamente el centroide (el origen del sistema de referencia para el subespacio lineal). Para fines ilustrativos se trabajará en esta parte de la exposición con las matrices originales sin considerar por el momento esta transformación lineal, sin embargo al momento de considerar los ejemplos se tomará en cuenta.

Definición 1

Sea A una matriz simétrica, se dice que A es positiva definida si $\lambda > 0, \forall \lambda \neq 0$.

Sea $X_{I \times J}$ una matriz de datos, que se refiere a una tabla de contingencia de I categorías para el primer factor y J categorías para el segundo factor, sea $M_{J \times J}$ una matriz simétrica, positiva definida; que defina la distancia en R^J . Sea $N_{I \times J}$ una matriz diagonal cuyos elementos son los pesos de los I puntos (si se considera como matriz de distancias a la matriz diagonal de pesos por renglón entonces la matriz $N_{I \times I}$ que representará la distancia para el espacio R^I) Véase cuadro 1 Capítulo I.

Sea u un vector unitario en el espacio R^J que satisfaga la ecuación $u^t M u = 1$. Las proyecciones de los I puntos sobre el eje que define el vector u son los elementos del vector $\hat{v}_{I \times 1}$ tal que $\hat{v} = X M u$. La suma de las distancias del centroide a la proyección ponderadas al cuadrado es $\hat{v}^t N \hat{v}$, a este término se le conoce como ajuste ya que es precisamente lo que corresponde a $\sum_{i=1}^I \overline{G \text{Pr } o y_{u_i}}^2$ y de acuerdo a lo expresado con anterioridad en la medida que esta cantidad sea maximizada ocurrirá que la distancia entre los perfiles y el subespacio generado por el vector u es menor.

En términos matriciales esto queda expresado de la siguiente manera:

$$\hat{v}^t N \hat{v} = u^t M X^t N X M u \quad (4)$$

esta cantidad debe ser maximizada con la restricción:

$$u^t M u = 1 \quad (5)$$

Es decir el problema es maximizar la suma del cuadrado de las proyecciones ponderadas sobre el subespacio generado por el vector u con la métrica M y además con la restricción de que el

vector u tenga norma 1 para la misma métrica definida por M .

De la expresión anterior observamos que lo único que desconocemos es el vector u , por lo que el problema ahora es su determinación. Para resolver este problema veamos la siguiente proposición:

Proposición 1

Sea A y M matrices simétricas, además M es positiva definida, un forma cuadrática para la matriz A es $u^t A u$, maximizar esta forma cuadrática bajo la restricción $u^t M u = 1$. Se escribe:

$$\max u^t A u$$

$$\text{s.a } u^t M u = 1$$

y u_α es el eigenvector de la matriz $M^{-1} A$ que corresponde al eigenvalor más grande λ_α . ($\alpha = 1, 2, \dots$, rango de A), el cual maximiza la forma cuadrática $u_\alpha^t A u_\alpha$. Esta proposición además establece que el vector unitario u_α , para la distancia definida por M ($u_\alpha^t M u_\alpha = 1$) es M -ortogonal a u_β (es decir $u_\alpha^t M u_\beta = 0$ para $\alpha > \beta$) †.

Utilizando la proposición anterior para (4), se tiene

$$\max \hat{v}^t N \hat{v} = u^t M X^t N X M u$$

s.a

$$u^t M u = 1$$

$$\text{Sea } A = M X^t N X M$$

$$\Rightarrow \max \hat{v}^t M \hat{v} = u^t M X^t N X M u = u^t A u$$

s.a

$$u^t M u = 1$$

que es de la misma forma que la proposición 1 además tenemos que

1) $A_{J_{NM}} = M X^t N X M$ es simétrica.

Dem:

Sabemos que para dos matrices A y B $(A B)^t = B^t A^t$, y que $(A^t)^t = A$

$$\Rightarrow A^t = (M X^t N X M)^t$$

$$= (M (X^t N X M))^t = (X^t (N X M))^t M^t = (N (X M))^t X M^t = (X M)^t N^t X M^t$$

$$= M^t X^t N^t X M^t$$

pero sabemos por hipótesis que M es simétrica y que N es una matriz diagonal

$$\Rightarrow A' = M X' N X M$$

$\therefore A = A'$ es decir A es simétrica. †

2) Por hipótesis sabemos que la matriz de distancia definida por M es simétrica y positiva definida.

Así u_α es el eigenvector de la matriz $M^{-1} M X' N X M = X' N X M$ correspondiente al eigenvalor más grande λ_α ($\alpha = 1, 2, \dots, J$)

Al vector u_α se le conoce como el α eje principal. Este vector satisface la ecuación $X' N X M u_\alpha = \lambda_\alpha u_\alpha$.

El operador "proyección" sobre el eje u_α , definido como $\varphi_\alpha = M u_\alpha$, es llamado o se le conoce como "factor". Y satisface la ecuación:

$$M X' N X \varphi_\alpha = \lambda \varphi_\alpha,$$

ya que $X' N X M u_\alpha = \lambda u_\alpha$ y si premultiplicamos por M tenemos $M X' N X (M u_\alpha) = \lambda (M u_\alpha)$ que da la relación anterior.

Observación: Los factores definidos como φ_α son α ($\alpha = 1, \dots, \text{rango de } A$) para cada vector u_α definido por su respectivo eigenvalor. Además la norma de estos factores bajo la distancia definida por M^{-1} es uno:

$$\varphi_\alpha' M^{-1} \varphi_\alpha = u_\alpha' M M^{-1} M u_\alpha = u_\alpha' M u_\alpha = 1$$

Al resolver el problema de maximización hemos encontrado los vectores u que generan al mejor espacio vectorial que ajusta a los perfiles.

En el análisis de correspondencias se muestran los perfiles por renglón y los perfiles por columna de una tabla de contingencia como puntos de un espacio vectorial de menor dimensión. Ambas representaciones pueden ser presentadas en una sola gráfica que muestra los perfiles de ambos factores de manera simultánea.

Sin embargo, todavía tenemos un problema importante: ¿cuántos vectores de los u 's que fueron encontrados son necesarios para generar el subespacio adecuado S y dar una buena representación de los perfiles en la gráfica?. Dejaremos por un momento la respuesta a esta pregunta, y obtengamos la representación de los perfiles para el ejemplo de la tabla 2.1. Los resultados fueron obtenidos mediante el uso del paquete mathematica V 4.0.

2.3.2.1 Análisis en R' . Cálculo de factores .

Los elementos que definen a este espacio para los datos del ejemplo son:

1. 1- perfiles por renglón, los cuales están centrados en el origen. Cada uno de estos vectores o puntos se encuentran en los renglones de la matriz R^* referida a los renglones de la matriz perfil por renglón tomando su desviación con respecto al centroide, en el espacio R' .

$$R^* = R_{R'} - 1r$$

donde

$R = D_r^{-1} F$ es la matriz perfil por renglón

$$1^*r = \begin{pmatrix} 1 \\ 1 \\ 1 \\ 1 \end{pmatrix} \begin{pmatrix} 0.239 & 0.223 & 0.111 & 0.133 \end{pmatrix} = \begin{pmatrix} 0.239 & 0.223 & 0.114 & 0.123 \\ 0.239 & 0.223 & 0.114 & 0.123 \\ 0.239 & 0.223 & 0.114 & 0.123 \\ 0.239 & 0.223 & 0.114 & 0.123 \end{pmatrix}$$

así, los elementos de las celdas de la matriz R^* para el ejemplo de la tabla 2.1 están dados en la tabla 2.4.

Tabla 2.4. Perfiles Renglón Desviación del centroide

Turbo	Tipo de defecto			
	A	B	C	D
1	-0.07991	0.00010	0.06448	0.01532
2	0.03135	0.03652	-0.06007	-0.07099
3	0.03783	-0.06044	-0.00247	0.04609

Para fines del desarrollo en la explicación la matriz R^* que se refiere a la matriz perfil por renglones considerando su desviación con respecto al centroide, será considerada conceptualmente como la matriz $D_r^{-1} F$ que es la forma matricial para calcular la matriz perfil por renglones.

2. La distancia definida entre dos perfiles de R^* es

$$d^2(i, i') = \sum_{j=1}^J \frac{1}{f_j} \left(\frac{f_{ij}}{f_j} - \frac{f_{i'j}}{f_j} \right)^2$$

3. La distancia está caracterizada por la la matriz diagonal del recíproco de los pesos por columna.

$$D_c^{-1} = \begin{pmatrix} \frac{1}{0.239} & 0 & 0 & 0 \\ 0 & \frac{1}{0.223} & 0 & 0 \\ 0 & 0 & \frac{1}{0.414} & 0 \\ 0 & 0 & 0 & \frac{1}{0.123} \end{pmatrix}$$

4. Los pesos de los vectores para este espacio están dados por la la matriz diagonal de pesos por renglón

$$D_r = \begin{pmatrix} 0.304 & 0 & 0 \\ 0 & 0.311 & 0 \\ 0 & 0 & 0.385 \end{pmatrix}$$

5. La matriz transpuesta de correspondencias es:

$$F' = \begin{pmatrix} 0.048 & 0.084 & 0.106 \\ 0.067 & 0.100 & 0.055 \\ 0.145 & 0.110 & 0.158 \\ 0.042 & 0.016 & 0.064 \end{pmatrix}$$

Así, de acuerdo al problema de encontrar un subespacio de menor dimensión que R^4 tenemos que encontrar los vectores u_a que generarán este nuevo subespacio; estos vectores son unitarios para la distancia definida en R^4 la cual está dada por D_c^{-1} y además deberán satisfacer la siguiente restricción $u_a' D_c^{-1} u_a = 1$. Es decir estos vectores serán ortonormales para la distancia definida.

Una vez encontrados esos vectores u_a podremos proyectar los perfiles de la matriz R^* en cada eje generado por el vector u_a conservando la distancia de cada perfil afectada por su peso. La proyección está dada por:

$R^* D_c^{-1} u_a = \hat{v}_{r \times 1}$, lo que se traduce como que la matriz perfil por renglones es proyectada a sobre el vector u_a , lo cual da por resultado un vector en el cual cada entrada representa la proyección del elemento i -ésimo sobre el vector u_a .

La expresión anterior puede presentarse de manera matricial como $D_r^{-1} F D_c^{-1} u_a$.

Recordando que la cantidad que debe ser maximizada es la suma de cuadrados de las proyecciones de la nube de puntos al centroide (en la que se considera el peso de cada elemento dentro de la nube) esta suma de cuadrados está dada por la siguiente forma cuadrática $\hat{v}' D_r \hat{v}$.

Estableciéndose el siguiente problema de maximización:

$$\max \hat{v}^t D_r \hat{v} = u_a^t D_c^{-1} F^t D_r^{-1} F D_c^{-1} u_a$$

s. a.

$$u_a^t D_c^{-1} u_a = 1.$$

Sea

$$A = D_c^{-1} F^t D_r^{-1} F D_c^{-1}$$

$$M = D_c^{-1}$$

Por lo tanto de acuerdo a la proposición 1 se tiene que: u_a es el eigenvector de la matriz

$$M^{-1} A = (D_c^{-1})^{-1} D_c^{-1} F^t D_r^{-1} F D_c^{-1} = F^t D_r^{-1} F D_c^{-1} \text{ es decir } u_a \text{ cumple con}$$

$$F^t D_r^{-1} F D_c^{-1} u_a = \lambda_a u_a$$

Para el ejemplo tenemos que

$$S = F^t D_r^{-1} F D_c^{-1}$$

$$= \begin{bmatrix} 0.048 & 0.084 & 0.106 \\ 0.067 & 0.100 & 0.055 \\ 0.145 & 0.110 & 0.158 \\ 0.042 & 0.016 & 0.064 \end{bmatrix} \begin{bmatrix} -0.079 & 0.0001 & 0.064 & 0.015 \\ 0.031 & 0.0996 & -0.060 & -0.070 \\ 0.037 & -0.0804 & -0.602 & 0.045 \end{bmatrix} \begin{bmatrix} \frac{1}{0.239} & 0 & 0 & 0 \\ 0 & \frac{1}{0.223} & 0 & 0 \\ 0 & 0 & \frac{1}{0.414} & 0 \\ 0 & 0 & 0 & \frac{1}{0.123} \end{bmatrix}$$

$$= \begin{bmatrix} 0.0116 & -0.0009 & -0.0052 & -0.0033 \\ -0.0008 & 0.0249 & -0.0042 & -0.0291 \\ -0.0091 & -0.0079 & 0.0057 & 0.0128 \\ -0.0016 & -0.0160 & 0.0038 & 0.0196 \end{bmatrix}$$

Los eigenvalores de esta matriz determinan los vectores u_a que son conocidos como los α ejes principales. Al vector $\varphi_a = D_c^{-1} u_a$, se le conoce como primer factor, y corresponde al operador proyección. Además φ_a es un eigenvector de la matriz

$$D_c^{-1} F^t D_r^{-1} F, \text{ ya que si premultiplicamos (1) por } D_c^{-1}, \text{ tenemos que } D_c^{-1} F^t D_r^{-1} F (D_c^{-1} u_a) = \lambda_a (D_c^{-1} u_a).$$

Las proyecciones de los I puntos sobre el eje principal u_a quedan determinados como:

$$R D_c^{-1} u_a = D_r^{-1} F D_c^{-1} u_a = D_r^{-1} F \varphi_a, \text{ ya que } \varphi_a = D_c^{-1} u_a$$

En general la matriz S no es simétrica y por lo tanto no es posible aplicar el resultado de la

proposición 1. Para encontrar la solución a nuestro problema haremos lo siguiente:

Sea

$S = F^T D_r F D_c^{-1}$ y llamemos

$\hat{A} = F^T D_r^{-1} F$, la cual es simétrica.

Dem.

$$\hat{A}^t = (F^T D_r^{-1} F)^t = (D_r^{-1} F)^t (F^T)^t = F^T (D_r^{-1})^t F = F^T D_r^{-1} F = \hat{A}^t$$

Además la matriz D_c^{-1} es una matriz diagonal. Por lo anterior se puede escribir como $D_c^{-1} = D_c^{-\frac{1}{2}} D_c^{-\frac{1}{2}}$. Así, S puede ser reescrita como: $S = \hat{A} D_c^{-\frac{1}{2}} D_c^{-\frac{1}{2}}$.

Como nuestro problema era encontrar los vectores que cumplieran con lo siguiente: $S u = \lambda u$, entonces tenemos: $\hat{A} D_c^{-\frac{1}{2}} D_c^{-\frac{1}{2}} c u = \lambda u$.

Si premultiplicamos en los dos lados de la igualdad por $D_c^{-\frac{1}{2}}$, se tiene

$$D_c^{-\frac{1}{2}} \hat{A} D_c^{-\frac{1}{2}} D_c^{-\frac{1}{2}} u = \lambda D_c^{-\frac{1}{2}} u. \quad (6)$$

Haciendo $w = D_c^{-\frac{1}{2}} u$ en (6) tenemos lo siguiente: $D_c^{-\frac{1}{2}} \hat{A} D_c^{-\frac{1}{2}} w = \lambda w$.

La matriz $S_1 = D_c^{-\frac{1}{2}} \hat{A} D_c^{-\frac{1}{2}}$ es simétrica

Dem.

$$S_1^t = (D_c^{-\frac{1}{2}} \hat{A} D_c^{-\frac{1}{2}})^t = (\hat{A} D_c^{-\frac{1}{2}})^t (D_c^{-\frac{1}{2}})^t = (D_c^{-\frac{1}{2}})^t \hat{A}^t (D_c^{-\frac{1}{2}})^t = D_c^{-\frac{1}{2}} \hat{A} D_c^{-\frac{1}{2}} = S_1$$

además la matriz S_1 tiene los mismos eigenvalores λ que S.

La matriz S_1 para los datos del ejemplo de la tabla 2.1 es:

$$S_1 = \begin{pmatrix} 0.0116 & -0.0008 & -0.0069 & -0.0023 \\ -0.0008 & 0.0249 & -0.0058 & -0.0216 \\ -0.0069 & -0.0058 & 0.0057 & 0.0070 \\ -0.0023 & -0.0216 & 0.0070 & 0.0196 \end{pmatrix}$$

Los eigenvalores son:

$$\lambda_1 = 0.0463$$

$$\lambda_2 = 0.0157$$

$$\lambda_3 = -2.72 \times 10^{-11}$$

$$\lambda_4 = 9.7911 \times 10^{-19}$$

Los eigenvectores son:

$$W_1 = \begin{pmatrix} -0.0719, & -0.7200, & 0.2287, & 0.6511 \end{pmatrix}$$

$$W_2 = \begin{pmatrix} -0.8528, & 0.2447, & 0.4610, & 0.0144 \end{pmatrix}$$

$$W_3 = \begin{pmatrix} -0.4894, & -0.4725, & -0.6435, & -0.3505 \end{pmatrix}$$

$$W_4 = \begin{pmatrix} 0.1670, & -0.4453, & 0.5665, & -0.6729 \end{pmatrix}$$

Los eigenvectores de la matriz original se obtiene como: $u = D_c^{-\frac{1}{2}} w$, y los factores son obtenidos de la siguiente manera:

$$\varphi = D_c^{-1} u = D_c^{-1} w$$

los cuales son:

$$\text{La matriz de factores es: } \begin{pmatrix} -0.147087 & -1.74231 & -1.00004 & 0.341292 \\ -1.52358 & 0.517987 & -0.999958 & -0.942305 \\ 0.355445 & 0.716465 & -1.00007 & 0.9804 \\ 1.95722 & 0.0411931 & -0.999996 & -1.91915 \end{pmatrix}$$

Así, la proyección sobre los ejes principales son:

$$\text{La matriz de proyecciones es: } \begin{pmatrix} 0.069701 & 0.186109 & 2.7103 \times 10^{-11} & 5.55112 \times 10^{-17} \\ -0.319412 & -0.0433942 & 2.72575 \times 10^{-11} & 8.30367 \times 10^{-17} \\ 0.19986 & -0.107494 & 2.72267 \times 10^{-11} & -9.71445 \times 10^{-17} \end{pmatrix}$$

2.3.2.2 Análisis en R^l . Cálculo de factores.

Se hará el mismo análisis para encontrar los factores y las proyecciones para el espacio R^l cuyos elementos se encuentran en la matriz perfil por columnas.

Los elementos que definen a este espacio son:

1. J-perfiles por columna, los cuales están centrados en el origen. Cada uno de estos vectores o puntos se encuentran en los renglones de la matriz C^* referida a los renglones de la matriz perfil por columnas tomando su desviación con respecto al centroide en el espacio R^l .

Tabla 2.3. Perfiles por columna

Turno	Tipo de defecto			
	A	B	C	D
1	0.2027	0.30435	0.35156	0.34211
2	0.35135	0.44928	0.26563	0.13158
3	0.44595	0.24638	0.38281	0.52632
Total	1	1	1	1

$$C^* = C_{IKJ} - c1$$

donde $C_{IKJ} = F_{IKJ} D_c^{-1}$ es la matriz perfil por columnas.

$$c^*1 = \begin{pmatrix} 0.304 \\ 0.311 \\ 0.385 \end{pmatrix} \begin{pmatrix} 1 & 1 & 1 & 1 \end{pmatrix} = \begin{pmatrix} 0.304 & 0.304 & 0.304 & 0.304 \\ 0.311 & 0.311 & 0.311 & 0.311 \\ 0.385 & 0.385 & 0.385 & 0.385 \end{pmatrix}$$

De manera similar a lo establecido en la sección anterior, la matriz C^* será referida conceptualmente como la matriz perfil por columnas considerando su desviación con respecto al centroide: $F D_c^{-1}$. Otra forma equivalente de trabajar esta matriz es considerando su transpuesta $D_c^{-1} F^t$, y esto es precisamente lo que se hará en esta sección, de tal manera que C^* esta dada por la tabla 2.5.

Tabla 2.5. Perfiles columna: Desviación del centroide

Turno	Tipo de defecto			
	A	B	C	D
1	-0.10180	0.00014	0.04726	0.03730
2	0.04067	0.13860	-0.04526	-0.17910
3	0.06063	-0.13874	-0.00200	0.14120

2. La distancia definida entre dos puntos o columnas es:

$$d^2(j, j') = \sum_{r=1}^t \frac{1}{f_r} \left(\frac{f_{jr}}{f_r} - \frac{f_{j'r}}{f_r} \right)^2$$

3. La distancia está caracterizada por la forma cuadrática D_F^{-1} : la matriz diagonal de pesos por renglón.

$$D_F^{-1} = \begin{pmatrix} \frac{1}{0.304} & 0 & 0 \\ 0 & \frac{1}{0.311} & 0 \\ 0 & 0 & \frac{1}{0.385} \end{pmatrix}$$

4. La matriz F es

$$F = \begin{pmatrix} 0.0485 & 0.0679 & 0.14563 & 0.0420 \\ 0.0841 & 0.1003 & 0.1100 & 0.0161 \\ 0.1067 & 0.0550 & 0.1585 & 0.0647 \end{pmatrix}$$



De manera análoga al análisis en R^d , queremos encontrar vectores v_a ortonormales de tal manera que para la distancia definida por la matriz D_r^{-1} se tenga $v_a^t D_r^{-1} v_a = 1$.

Así, el vector de las "J" proyecciones sobre el eje definido por v es \hat{w} tal que

$$\hat{w} = D_c^{-1} F^t D_r^{-1} v, \text{ donde el factor de proyección es } \psi = D_r^{-1} v$$

De acuerdo a lo establecido en la interpretación gráfica queremos maximizar la suma ponderada de cuadrados de las proyecciones tomando en cuenta los pesos de cada uno de los elementos, los cuales se encuentran en la matriz diagonal D_c .

Así el problema queda representado como:

$$\text{Max } \hat{w}^t D_c \hat{w} = v^t D_r^{-1} F^t D_c^{-1} D_c D_c^{-1} F^t D_r^{-1} v = v^t D_r^{-1} F^t D_c^{-1} F^t D_r^{-1} v$$

s.a

$$v_a^t D_r^{-1} v_a = 1$$

$$\text{Sea } A = D_r^{-1} F^t D_c^{-1} F^t D_r^{-1}$$

$$M = D_r^{-1}$$

de manera semejante que al análisis en R^d , por la proposición 1 tenemos que:

v_a es el eigen vector de la matriz $(D_r^{-1})^{-1} D_r^{-1} F^t D_c^{-1} F^t D_r^{-1} = F^t D_c^{-1} F^t D_r^{-1}$ de manera que se satisfice la ecuación:

$$F^t D_c^{-1} F^t D_r^{-1} v_a = \lambda_a v_a \quad (7)$$

Para los datos del ejemplo tenemos:

$$S^* = F^t D_c^{-1} F^t D_r^{-1}$$

$$= \begin{pmatrix} 0.0485 & 0.0679 & 0.14561 & 0.0420 \\ 0.6841 & 0.1003 & 0.1100 & 0.0161 \\ 0.1067 & 0.0330 & 0.1583 & 0.0647 \end{pmatrix} \begin{pmatrix} -0.1015 & 0.01067 & 0.06083 \\ 0.00014 & 0.1386 & -0.13874 \\ 0.01734 & -0.04505 & -0.0023 \\ 0.0379 & -0.1791 & 0.1112 \end{pmatrix} \begin{pmatrix} \frac{1}{0.3104} & 0 & 0 \\ 0 & \frac{1}{0.311} & 0 \\ 0 & 0 & \frac{1}{0.385} \end{pmatrix}$$

$$= \begin{pmatrix} 0.0119 & -0.0090 & -0.0021 \\ -0.0086 & 0.0301 & -0.0174 \\ -0.0025 & -0.0221 & 0.0199 \end{pmatrix}$$

Los eigenvalores de esta matriz determinan a los vectores v_α que es conocido como el primer α eje principal, al operador $\psi_\alpha = D_r^{-1} v_\alpha$ se le conoce como el factor α -ésimo. Además ψ_α satisface la ecuación $D_r^{-1} F D_c^{-1} I^T \psi_\alpha = \lambda_\alpha \psi_\alpha$, es decir ψ_α (premultiplicando (7) por D_r^{-1}) es un eigenvector de la matriz $D_r^{-1} F D_c^{-1} I^T$.

Las coordenadas de los "J" puntos sobre el eje generado por v_α son: $D_c^{-1} I^T D_r^{-1} v = D_c^{-1} I^T \psi_\alpha$. Donde $\psi_\alpha = D_r^{-1} v_\alpha$

Infortunadamente, de manera general la matriz S' no es simétrica por lo que no se puede aplicar la proposición I. Para poder resolver el problema, haremos lo siguiente:

Sea

$$S' = F D_c^{-1} I^T D_r^{-1}$$

La matriz $\hat{A}' = F D_c^{-1} I^T$ es simétrica y además $D_r^{-1} = D_r^{-\frac{1}{2}} D_r^{-\frac{1}{2}}$. De tal manera que $S' = \hat{A}' D_r^{-\frac{1}{2}} D_r^{-\frac{1}{2}}$. Si premultiplicamos por $D_r^{-\frac{1}{2}}$ y hacemos $w = D_r^{-\frac{1}{2}} v$, entonces el problema original tiene la siguiente forma:

Encontrar los eigenvalores λ y eigenvectores w de:

$$D_r^{-\frac{1}{2}} \hat{A}' D_r^{-\frac{1}{2}} w = \lambda w$$

La matriz $S_2 = D_r^{-\frac{1}{2}} \hat{A}' D_r^{-\frac{1}{2}}$ es simétrica y tiene los mismos eigenvalores de S' . Con esta nueva matriz los factores se calculan como: $\psi_\alpha = D_r^{-\frac{1}{2}} w$.

$$S_2 = \begin{pmatrix} 0.0119 & -0.0091 & -0.0024 \\ -0.0085 & 0.0301 & -0.0194 \\ -0.0023 & -0.0199 & 0.0199 \end{pmatrix}$$

cuyos eigenvalores son:

$$\lambda_1 = 0.0463$$

$$\lambda_2 = 0.0157$$

$$\lambda_3 = 1.2845 \times 10^{-6}$$

asociados a estos eigenvalores, sus eigenvectores son:

$$w1 = (0.1688, -0.7933, 0.5849)$$

$$w2 = (0.8298, -0.2060, -0.5185)$$

$$w3 = (-0.5516, -0.5576, -0.6202)$$

ψ_a se le conoce como el factor correspondiente al eigenvalor λ_a .

Así, la matriz de factores es:

$$\text{La matriz de factores es: } \begin{pmatrix} 0.306265 & 1.50514 & -1.00052 \\ -1.42254 & -0.369434 & -1.00002 \\ 0.942701 & -0.835667 & -0.999574 \end{pmatrix}$$

Y, las proyecciones de los perfiles sobre los ejes principales

$$\text{La matriz de proyecciones es: } \begin{pmatrix} -0.0309726 & -0.212319 & -2.41289 \cdot 10^{-7} \\ -0.327257 & 0.352777 & -4.17077 \cdot 10^{-7} \\ 0.0770458 & 0.0501763 & 1.40693 \cdot 10^{-6} \\ 0.400143 & 0.00354634 & -7.14046 \cdot 10^{-6} \end{pmatrix}$$

2.3.2.3 Relación del análisis de R' con el análisis en R^J : fórmulas de transición.

$$\text{En } R^J \text{ tenemos que } F^T D_r^{-1} F D_c^{-1} u_a = \lambda_a u_a. \quad (8)$$

es decir u_a es un eigenvector de la matriz $F^T D_r^{-1} F D_c^{-1}$ correspondiente al eigenvalor λ_a .

Si premultiplicamos (8) por $F D_c^{-1}$ tenemos que

$(F D_c^{-1} F^T D_r^{-1} F D_c^{-1}) u_a = \lambda_a (F D_c^{-1} u_a)$, lo que significa que $F D_c^{-1} u_a$ es un eigenvector de la matriz $F D_c^{-1} F^T D_r^{-1} F D_c^{-1}$ para el mismo eigenvalor λ_a en R^J lo cual implica que v_a en R^J es proporcional a $F D_c^{-1} u_a$ en R^J .

Dándose las fórmulas de transición:

$$v_a = \frac{1}{\sqrt{\lambda_a}} F^T D_r^{-1} u_a$$

$$u_a = \frac{1}{\sqrt{\lambda_a}} F^T D_r^{-1} v_a$$

Nótese que

$$\psi_a = D_r^{-1} v_a = D_r^{-1} \frac{1}{\sqrt{\lambda_a}} F D_c^{-1} u_a = \frac{1}{\sqrt{\lambda_a}} D_r^{-1} F D_c^{-1} u_a = \frac{1}{\sqrt{\lambda_a}} D_r^{-1} F \varphi_a$$

$$\varphi_a = D_c^{-1} u_a = D_c^{-1} \frac{1}{\sqrt{\lambda_a}} F^T D_r^{-1} v_a = \frac{1}{\sqrt{\lambda_a}} D_c^{-1} F^T D_r^{-1} v_a = \frac{1}{\sqrt{\lambda_a}} D_c^{-1} F^T \psi_a$$

Así, por ejemplo una vez hecho el análisis para R^J y si queremos la representación de los datos en R^I se tiene:

$$\psi_1 = \frac{1}{\sqrt{0.11403}} \begin{pmatrix} \frac{1}{0.304} & 0 & 0 \\ 0 & \frac{1}{0.311} & 0 \\ 0 & 0 & \frac{1}{0.385} \end{pmatrix} \begin{pmatrix} 0.0485 & 0.0679 & 0.11563 & 0.0120 \\ 0.0641 & 0.1001 & 0.1100 & 0.0161 \\ 0.1067 & 0.0550 & 0.1585 & 0.0647 \end{pmatrix} \begin{pmatrix} -0.1470 \\ -1.5235 \\ 0.3554 \\ 1.8571 \end{pmatrix}$$

$$= \begin{pmatrix} 0.2926 \\ -1.4378 \\ 0.9287 \end{pmatrix}$$

$$\psi_2 = \frac{1}{\sqrt{0.01374}} \begin{pmatrix} \frac{1}{0.304} & 0 & 0 \\ 0 & \frac{1}{0.311} & 0 \\ 0 & 0 & \frac{1}{0.385} \end{pmatrix} \begin{pmatrix} 0.0185 & 0.0679 & -0.11563 & 0.0420 \\ 0.0641 & 0.1003 & 0.1100 & 0.0161 \\ 0.1067 & -0.0550 & 0.1585 & 0.0647 \end{pmatrix} \begin{pmatrix} -1.7423 \\ 0.5179 \\ 0.7164 \\ 0.0411 \end{pmatrix}$$

$$= \begin{pmatrix} 1.4832 \\ -0.3903 \\ -0.8566 \end{pmatrix}$$

Las relaciones anteriores muestran que las coordenadas de los puntos sobre un eje principal en el espacio R^I son proporcionales a las coordenadas de los puntos en el otro espacio R^J correspondientes al mismo eigenvalor.

Las relaciones anteriores pueden ser escritas en forma explícita es decir:

$$\psi_{ai} = \frac{1}{\sqrt{\lambda_a}} \sum_{j=1}^J \frac{f_{ij}}{\sqrt{\lambda_j}} \varphi_{aj}$$

.....

$$\varphi_{aj} = \frac{1}{\sqrt{\lambda_a}} \sum_{i=1}^I \frac{f_{ij}}{\sqrt{\lambda_j}} \psi_{ai}$$

Estas fórmulas de transición son las que permiten representar gráficamente de manera simultánea los perfiles de los dos espacios, sobre un subespacio (plano) generado por los dos primeros ejes principales. Y son las fórmulas que nos permitirán en el próximo capítulo establecer

las relaciones entre varios factores.

En la Fig 2.7 se muestran los perfiles para tipo de defecto y turno, representados en un subespacio de dimensión dos.

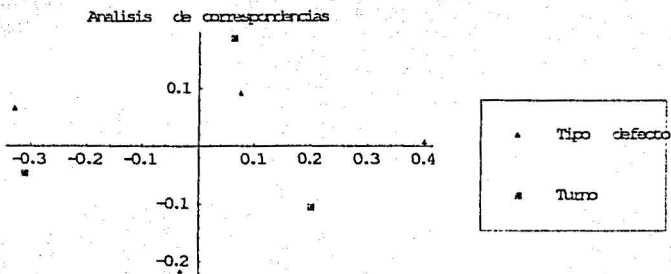


Fig. 2.7 Perfiles por tipo de defecto y turno en un espacio de dos dimensiones.

Sin embargo, existen algunas preguntas que se tienen que hacer:

- 1) ¿Porqué nada más se eligieron dos ejes?
- 2) ¿Qué calidad de representación tienen los perfiles en esta gráfica?
- 3) ¿Puede hacerse una interpretación en terminos de proximidad de manera simultánea entre los perfiles por renglón y los perfiles por columna?

Estas preguntas serán respondidas en las siguientes secciones.

2.4 Calidad del análisis: contribuciones.

En el análisis de correspondencias, cuando se obtienen las proyecciones sobre los nuevos ejes de referencia o ejes factoriales se presenta una gráfica de los perfiles, generalmente en un espacio de menor dimensión (dos o tres a lo mucho).

Para poder interpretar está gráfica de una manera adecuada es necesario conocer con que certeza se están describiendo las relaciones entre las categorías analizadas, es decir, que tan buena es la representación de los perfiles en la gráfica.

Algunos elementos que nos ayudarán a conocer la calidad de la representación se presentan a continuación.

2.4.1 Inercia: descomposición sobre los ejes factoriales ¿Qué significado geométrico y estadístico tiene la inercia?

Regresando a las preguntas finales de la sección 2.3.2.3

Se define la inercia de un punto i con masa f_i con respecto al centroide G como la distancia al cuadrado del punto i al centroide G ponderada por su masa, es decir :

$$I_G(i, f_i) = f_i d^2(G, i).$$

Y, se define la inercia de una nube de puntos $N(I)$ o $N(J)$ como:

$$I_G N(I) = \sum_{i=1}^I f_i d^2(G, i)$$

$$I_G N(J) = \sum_{j=1}^J f_j d^2(G, j)$$

Otra manera de representar la inercia del punto j -ésimo es:

$$I(j) = f_j (\text{Proy}_u^2 M_1 + \text{Proy}_u^2 M_2 + \dots + \text{Proy}_u^2 M_J)$$

La inercia de la nube en la dirección del eje α es el eigenvalor λ_α correspondiente a este eje. La inercia es la suma extendida a varios puntos de la nube, de las cantidades:

$$\lambda_\alpha = \sum_{i=1}^I f_i \text{Proy}_\alpha^2(i)$$

de tal manera que λ_α aparece como la media ponderada sobre los elementos de la nube $N(I)$, con su sistema asociado de masas f_i $i = 1, \dots, I$, de Las proyecciones al cuadrado generadas por el α -ésimo factor.

Regresando a la idea que se desarrolló para encontrar, para la nube de puntos $N(I)$ o $N(J)$ el subespacio de menor dimensión que ajustara de manera óptima a los puntos, se obtuvo un subespacio de dos dimensiones (plano), el cual generado por dos vectores que son ortogonales y que generan a los elementos que son proyectados (perfiles), es decir se puede descomponer la inercia en dos distancias debido al teorema de pitágoras, por lo que tenemos lo siguiente:

La inercia total = Inercia correspondiente al eje 1 + Inercia correspondiente al eje 2.

Y, demás la inercia correspondiente al primer eje factorial está registrada con el primer eigenvalor λ_1 , y la inercia correspondiente al segundo eje está registrada por el eigenvalor λ_2 . De manera adicional se cumple la siguiente relación:

$$\lambda_1 \leq \lambda_2 \leq \lambda_3 \leq \lambda_4.$$

Como un comentario adicional, se tiene que el número de ejes factoriales de la nube de puntos $N(I)$ no puede exceder al $\min\{(J-1), (I-1)\}$.

Para el centroide G de la nube de puntos $N(I)$ siendo el origen del sistema de ejes coordenadas se tiene que todos sus componentes son ceros. En lenguaje probabilístico, la media ponderada del cuadrado de una función centrada en cero es llamada varianza. Ahora como $Proy_a$ es de media cero sobre I con el sistema de masas f_a asociado para cada perfil, la media ponderada por los masas o pesos f_a del cuadrado de la proyección $Proy_a^2$ es su varianza, así, diremos que el factor a tiene λ_a por varianza, de tal manera que el factor 1 tiene λ_1 de varianza, el segundo factor tiene λ_2 y así sucesivamente. Un aspecto importante para mencionar es el hecho de que en análisis de correspondencias λ_1 no puede ser mayor que 1.

2.4.2- Inercia relativa y contribución absoluta.

1) ¿Porqué nada más se eligieron dos ejes?

La inercia relativa y la contribución absoluta son algunas medidas que ayudan a saber que tan buena es la representación de los datos por mediante el análisis.

La inercia relativa es la proporción de la inercia total correspondiente a cada eje factorial. Es decir, que porcentaje de la inercia está siendo representado por cada eje factorial.

$$r_a = \frac{\lambda_a}{\sum_{j=1}^2 \lambda_j} = \frac{\lambda_a}{Tr_{inca}}$$

En nuestro ejemplo tenemos que

$$r_1 = \frac{\lambda_1}{\sum_{j=1}^2 \lambda_j} = \frac{\lambda_1}{Tr_{inca}} = \frac{0.0463}{0.0463+0.0157} = 74.6\% \text{ de la inercia es explicada por el primer eje.}$$

$$r_2 = \frac{\lambda_2}{\sum_{j=1}^2 \lambda_j} = \frac{\lambda_2}{Tr_{inca}} = \frac{0.0157}{0.0463+0.0157} = 25.3\% \text{ de la inercia es explicada por el segundo eje.}$$

Así, se eligieron nada más los dos primeros ejes para los resultados del ejemplo de la tabla 2.1 porque entre ambos recogen el 100% de la inercia.

Para interpretar la información que proporcionan los ejes principales calcularemos las contribuciones absolutas para cada elemento j de la nube de puntos $N(J)$ que representa los defectos de los muebles.

Estas contribuciones se representan por $((ca_a(j)))$ e indican la proporción de la varianza explicada por cada perfil del factor tipo de defecto en relación a cada eje principal.

Sabemos que $u_a^t D_c^{-1} u_a = 1$, para la distancia definida en R^J . Además sabemos que $\varphi_a = D_c^{-1} u_a$, lo cual implica que $u_a = D_c \varphi_a$, dado que D_c^{-1} es positiva definida. Así, tenemos que $u_a^t D_c^{-1} u_a = (D_c \varphi_a)^t D_c^{-1} D_c \varphi_a = 1$

$$= \varphi_a^t D_c^t D_c^{-1} D_c \varphi_a = \varphi_a^t D_c D_c^{-1} D_c \varphi_a = 1 \text{ (debido a que } D_c^{-1} \text{ es simétrica)}$$

$$= \varphi_a^t D_c \varphi_a = 1,$$

lo cual de manera explícita queda expresado como: $\sum_{j=1}^J f_j \varphi_{aj}^2 = 1$

Recordemos que la proyección del elemento j -ésimo en R^I es

$$\varphi_{aj} = \frac{1}{\sqrt{\lambda_a}} \sum_{i=1}^I \frac{f_{ij}}{f_j} \psi_{ai}, \text{ desarrollando tenemos que:}$$

$$\sqrt{\lambda_a} \varphi_{aj} = \sum_{i=1}^I \frac{f_{ij}}{f_j} \psi_{ai}, \text{ sea } \hat{\varphi}_{aj} = \sqrt{\lambda_a} \varphi_{aj},$$

$$\text{por lo tanto tenemos que: } \hat{\varphi}_{aj} = \sum_{i=1}^I \frac{f_{ij}}{f_j} \psi_{ai}.$$

Así, la varianza de el conjunto de puntos proyectado sobre el eje a con respecto a G es:

$$\sum_{j=1}^J f_j \hat{\varphi}_{aj}^2 = \sum_{j=1}^J f_j (\sqrt{\lambda_a} \varphi_{aj})^2 = \sum_{j=1}^J f_j \lambda_a \varphi_{aj}^2 = \lambda_a$$

donde f_j es la masa asignada a la j -ésima columna y $\hat{\varphi}_{aj}^2$ representa la distancia de la proyección del perfil j -ésimo con respecto al centroide definido por las columnas.

Dado que la varianza está representada por el eigenvalor λ_a , pero además este valor representa la inercia del conjunto de puntos en R^I con respecto al centro de gravedad. Es decir $\lambda_1 + \lambda_2 + \lambda_3 + \dots + \lambda_I$ es igual a la inercia de todos los puntos con respecto a los ejes principales y además esta suma representa la varianza que cada uno de los ejes recoge.

Además $\lambda_1 + \lambda_2 + \lambda_3 + \dots + \lambda_I = \sum_{i=1}^I \lambda_i$, que corresponde a la inercia total de ambas nubes es igual al coeficiente de contingencia en media cuadrática $\phi^2 = \frac{\chi^2}{n}$ la cual en principio no es acotada y su interpretación para valores cercanos a cero se puede afirmar que los factores no están asociadas entre sí.

El cociente $\frac{\lambda_a \varphi_{aj}^2 f_j}{\lambda_a} = f_j \varphi_{aj}^2 = ca_a(j)$ Representa la contribución absoluta del elemento j -ésimo al eje principal a .

J	Masa	Contribución primer eje	Contribución Segundo eje
A	0.2394	0.00022	0.01141
B	0.2233	0.02391	0.00095
C	0.4142	0.00245	0.00336
D	0.12297	0.01969	3.78×10^{-6}

Observación

$$\sum_{j=1}^J ca_{\alpha}(j) = \lambda_{\alpha}$$

Si ahora, dividimos la contribución de cada perfil entre su respectivo eigenvalor tendremos la siguiente tabla

<i>J</i>	<i>Masa</i>	<i>Contribución %</i> primer eje	<i>Contribución %</i> Segundo eje
<i>A</i>	0.2394	0.00496	0.7250
<i>B</i>	0.2233	0.51648	0.00095
<i>C</i>	0.4142	0.05311	0.21395
<i>D</i>	0.12297	0.4252	0.00024

Que representa la contribución en términos de porcentaje de cada perfil tipo de defecto a la inercia que recoge el eje respectivo. En este ejemplo vemos que para el primer eje la categoría Tipo de defecto B y D son las que aportan una mayor inercia; mientras que para el segundo eje las categorías que mayor cantidad de inercia aportan son el Tipo de defecto A y C.

$$\sum_{j=1}^J \frac{ca_{\alpha}(j)}{\lambda_{\alpha}} = 1$$

De manera análoga en el espacio de los renglones o R^J tenemos que

$$f_{i\alpha} \psi_{\alpha i}^2 = ca_{\alpha}(i)$$

	<i>Masa</i>	<i>Contribución</i> primer eje	<i>Contribución</i> Segundo eje
1	0.304	0.00120	0.01052
2	0.311	0.02977	0.00074
3	0.385	0.01537	0.00444

donde $\sum_{i=1}^I ca_{\alpha}(i) = \lambda_{\alpha}$

Si ahora, dividimos la contribución de cada perfil entre su respectivo eigen_valor tendremos la siguiente tabla

<i>I</i>	<i>Masa</i>	Contribución % primer eje	Contribución % Segundo eje
1	0.304	0.02603	0.66879
2	0.311	0.64297	0.04739
3	0.385	0.33212	0.28256

Que representa la contribución en términos de porcentaje de la contribución de cada perfil Turno a la inercia que recoge el eje respectivo. En este ejemplo vemos que para el primer eje la categoría "segundo" y "tercer" turno son las que aportan una mayor inercia; mientras que para el segundo eje se tienen las categorías con mayor cantidad de inercia son el "primer" y "tercer" turno.

2.4.3 Correlaciones al cuadrado: cosenos cuadrados. ($Cor^2(j)$)

2) ¿Qué calidad de representación tienen los perfiles en esta gráfica?

Explica la parte de la varianza de una variable explicada por un eje principal.

Para entender este punto supongamos que estamos en el espacio definido por las columnas, es decir, estamos en R^I . Como estamos trabajando en un espacio ponderado hemos visto que la distancia entre dos puntos se calcula de la siguiente manera

$$d^2(j, j') = \sum_{r=1}^I \frac{1}{f_r} \left(\frac{f_{jr}}{f_r} - \frac{f_{j'r}}{f_r} \right)^2$$

Pero veamos que pasa con un elemento especial del espacio R^I : el centroide. Por lo que para calcular la distancia entre un columna j y el centroide tendremos

$$d^2(j, G) = \sum_{r=1}^I \frac{1}{f_r} \left(\frac{f_{jr}}{f_r} - f_r \right)^2, \text{ dado que el centroide se puede calcular como: } G_j = \sum_{r=1}^I f_r \frac{f_{jr}}{f_r} = \sum_{r=1}^I f_{jr} = f_j$$

El cuadrado de la proyección de la variable j sobre el eje α es igual a:

$$d_{\alpha}^2(j, G) = (\sqrt{\lambda_{\alpha}} \varphi_{\alpha j})^2$$

Observación

$$\sum_u d_u^2(j, G) = d^2(j, G)$$

Así, tenemos que $\cos^2 \omega = \frac{d_u^2(j, G)}{d^2(j, G)} = Cor(j)$

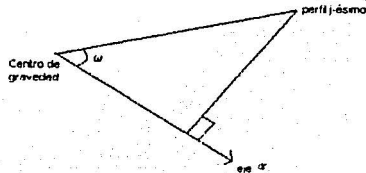


Fig 2.8 $\cos^2 \omega = \frac{d_u^2(j, G)}{d^2(j, G)} = Cor(j)$

En particular para el ejemplo tenemos para las columnas:

J	Dist.1erEje	Dist.2doEje	Dist.Total	cos ² 1er eje	cos ² 2do eje	C.R.*
A	0.00095	0.04766	0.04862	0.01972	0.98027	0.9999
B	0.10709	0.00426	0.11135	0.96173	0.03826	1
C	0.00393	0.00913	0.01406	0.42195	0.57801	1
D	0.16011	0.00003	0.16014	0.9999	0.00019	0.9999

*calidad de representación

Ahora, para los renglones:

I	Dist.1erEje	Dist.2doEje	Dist.Total	cos ² 1er eje	cos ² 2do eje	C.R.*
1	0.00396	0.03116	0.03500	0.10272	0.8972	1
2	0.09372	0.00239	0.09612	0.96173	0.0244	1
3	0.03994	0.01155	0.05149	0.77362	0.2243	1

*calidad de representación

2.4.4 Calidad de representación de un perfil.

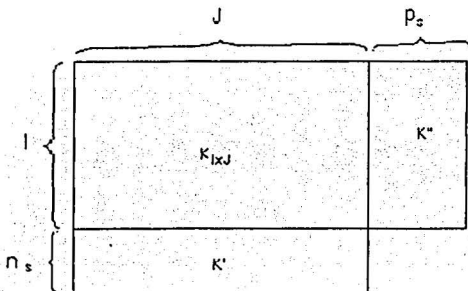
Explica la varianza explicada por los primeros α ejes principales

$$QLT = Cor(1) + Cor(2) + \dots + Cor(\alpha)$$

2.5. Tratamiento de los elementos suplementarios en la tabla de contingencia.

Para fines prácticos, dentro del análisis de correspondencias existe la posibilidad representar categorías suplementarias también denominadas ilustrativas, las cuales son otras categorías de la tabla de contingencia que se quieran añadir para proyectarlas en el subespacio previamente encontrado por el análisis de correspondencias y sin que la inserción de tales categorías interfiera en el resultado de la determinación del subespacio óptimo. De este modo se pueden crear varios escenarios, los cuales representan posibles hipótesis de investigación. Lo que se espera es que el nuevo elemento quede situado en la proximidad de los elementos que se le parecen.

Retornemos a la matriz de datos original $K_{I \times J}$, la cual quede ser incrementada por p_s categorías del factor que define las columnas o por n_s categorías del factor que define los renglones.



Estamos interesados en proyectar los perfiles de los p_s nuevos puntos con respecto a los p puntos analizados en el espacio R^L . Sea k''_{ij} la i -ésima coordenada de la j -ésima columna suplementaria. El perfil de esta columna suplementaria es el vector cuyo i -ésimo componente es $\frac{k''_{ij}}{k''_{.j}}$ es decir

$$j'' = \begin{bmatrix} \frac{k''_y}{k''_y} \\ \frac{k''_y}{k''_y} \\ \frac{k''_y}{k''_y} \\ \dots \\ \frac{k''_y}{k''_y} \end{bmatrix} \quad \text{con } k''_y = \sum_{j=1}^I k''_y$$

Dado que ya se obtuvieron las fórmulas de transición para las proyecciones sobre los ejes principales (4), lo que se hace es proyectar el punto j'' sobre el eje α , usando la fórmula:

$$\varphi_{\alpha j''} = \frac{1}{\sqrt{k''_\alpha}} \sum_{j=1}^I \left(\frac{k''_y}{k''_y} \right) \psi_{\alpha j}$$

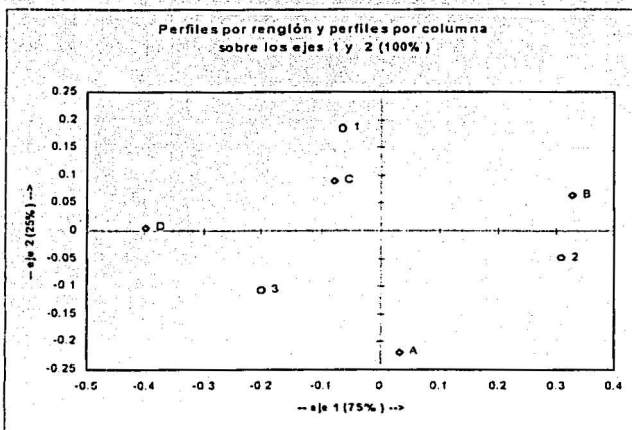
Para el caso de un categoría por renglón complementaria tenemos que los perfiles de las n_s categorías suplementarias con respecto a los I puntos analizados en el espacio R^J . Sea k'_j la j -ésima coordenada del i -ésimo renglón suplementario. El perfil de este renglón es el vector cuyo j -ésimo elemento es $\frac{k'_{ij}}{k'_{i\alpha}}$ es decir:

$$i' = \left[\frac{k'_{i1}}{k'_{i\alpha}} \quad \frac{k'_{i2}}{k'_{i\alpha}} \quad \frac{k'_{i3}}{k'_{i\alpha}} \quad \dots \quad \frac{k'_{iJ}}{k'_{i\alpha}} \right] \quad \text{donde } k'_{i\alpha} = \sum_{j=1}^J k'_{ij}$$

De nuevo lo que se hace es proyectar el punto i' sobre el eje α , usando la fórmula:

$$\psi_{\alpha i'} = \frac{1}{\sqrt{k'_{i\alpha}}} \sum_{j=1}^J \left(\frac{k'_{ij}}{k'_{i\alpha}} \right) \varphi_{\alpha j}$$

2.6 Descripción de la gráfica del Análisis de correspondencias para los datos de la tabla 2.1.



En base a las consideraciones hechas con las contribuciones observamos que con los dos primeros ejes principales se explica el 100% de la inercia, y además la calidad de representación también es del 100% para cada perfil, tanto por renglón o por columna. Se observa que en el primer eje el cual recoge el 75% de la inercia ésta es generada principalmente por el tipo de defecto D y B, en el caso del turno se tiene que el tercero y el segundo son los que más aportan a la inercia del primer eje. Con lo anterior se puede establecer que debido al alejamiento del origen (centroide) el tipo de defecto D se encuentra asociado mayoritariamente con el turno de fabricación 3, en contraparte con el tipo de defecto B que está más asociado con el turno de producción 2; es importante señalar que tomando en cuenta únicamente el primer eje no se muestra alguna relación con respecto a los tipos de defecto A y C y turno de producción 1.

Con el segundo eje, se puede ver que el tipo de defecto A es opuesto al B y C y además se aprecia que existe una asociación no tan marcada entre el turno 3 y el defecto A la cual es opuesta a la asociación entre el tipo de defecto C y el primer turno.

Observación:

En vista de lo anteriormente expuesto, es necesario poner atención en lo siguiente:

1. Inercia acumulada por los ejes que queremos retener.
2. La calidad de representación de los puntos.
3. Tener en cuenta que u y $-u$ son solución al problema de optimización, lo que ocasiona que se puedan obtener gráficas "aparentemente" diferentes. Las discrepancias son debidas a que en ocasiones lo que aparecía a la "derecha" ahora aparece a la izquierda y lo de "arriba" aparece ahora "abajo".

Resumen del capítulo

A manera de resumen se plantean en forma esquemática los pasos realizados hasta el momento para el desarrollo del análisis de correspondencias.

1. Se tienen los datos de un fenómeno o problema de investigación en una tabla de contingencia de dos factores con I y J categorías respectivamente.

2. El número de categorías en las variable determina los espacios correspondientes a los renglones (R^I) y a las columnas (R^J).

3. La tabla de contingencia puede ser vista como un matriz de $I \times J$.

4. Se busca en el análisis índices o factores a partir de los datos originales de tal manera que conserven la mayor cantidad de varianza original y además generen un espacio óptimo en el cual la distancia de los puntos a dicho espacio sea mínima.

5. Se encontraron para el espacio R^J esos factores o índices por medio de la relación $\varphi_{\alpha} = D_{\alpha}^{-1} u_{\alpha}$, donde el vector u_{α} maximiza la forma cuadrática asociada a suma ponderada de cuadrados de las proyecciones sobre el espacio generado (este vector se conoce como el α eje principal) con la restricción $u_{\alpha}' D_{\alpha}^{-1} u_{\alpha} = 1$ (norma uno). Y donde se utiliza la distancia Ji-cuadrada.

6. Encontrados los factores se procedió a encontrar las proyecciones sobre los ejes principales.

7. Para saber que tan buena es esta representación en los ejes se caculan las contribuciones respectivas de cada eje.

8. Todo este desarrollo se debe hacer con origen en el centroide debido a que el subespacio óptimo debe contener al centroide.

9. Todo esto es aplicable de la misma manera al espacio en R^I y para el cual se establece una espacio de dualidad con respecto al espacio R^J mediante las fórmulas de transición encontradas.

Capítulo III

Análisis de correspondencias múltiples.

El análisis de correspondencias simple, el cual ha sido descrito en el capítulo II, puede ser generalizado para poder estudiar la asociación de más de dos factores de manera simultánea, sin embargo como se aumenta la dimensión del espacio, para que las categorías de los factores puedan ser representados en un espacio de menor dimensión existe una pérdida de información en los factores por lo que los porcentajes de inercia en los ejes son bajos y la descripción de la gráfica no es del todo satisfactoria. En este caso la interpretación de los factores sigue las mismas reglas que en el análisis de correspondencias simple, la única peculiaridad reside en que la parte de la inercia explicada pierde aquí su interés ya que los valores propios no representan más que una pequeña parte de la inercia total.

Existen dos formas de presentar el análisis de correspondencias, una utilizando la información por individuos que se encuentra en una matriz indicadora Z o bien utilizando la matriz de Burt; la cual recoge la clasificación de la muestra de acuerdo a cada factor y es en si una generalización de la tabla de contingencia. Ambas formas son equivalentes en el sentido de la descripción entre las categorías y difieren únicamente en los valores de las inercias que cada uno arroja.

Cuando se trabaja con la matriz indicadora Z , lo que se hace es una extensión del análisis de correlación canónico. El propósito de presentar al Análisis de Correlación Canónico es el de ejemplificar una forma alternativa de trabajar con el AFC, recordemos que en el capítulo anterior el desarrollo del AFC se hizo con base a una tabla de contingencia, pero también se puede comenzar el análisis a través de los datos obtenidos de manera directa es decir, la información codificada. Esta información se presenta en una matriz indicadora Z en la que aparece un 1 para la categoría del factor que se haya elegido. El tener esta matriz, en la cual se le asigna un valor numérico a una característica cualitativa nos permitirá el uso de otras técnicas como el análisis de correlación canónico y además nos permitirá estudiar el problema de la generalización del análisis de correspondencias para el caso en que se tengan más de dos factores.

Es importante mencionar que la forma en que se quiera empezar con el análisis depende de las características particulares del problema que esté siendo atacado, así por ejemplo, si se quiere hacer un estudio de discriminación con respecto a un grupo de variables que puedan describir de manera adecuada el perfil de un trabajador (Pérez Soto Claudia: Análisis de correspondencias múltiples como una técnica para el estudio de datos cualitativos, 1990) lo acertado sería utilizar la matriz indicadora Z , para hacer un análisis por individuos y por variables. Pero en dado caso de que se quiera estudiar la asociación entre las categorías sería recomendable el uso de la matriz de Burt.

3.1 El análisis de correlación canónico (ACC).

El Análisis de correlación canónico (ACC) es una técnica multivariada que investiga la relación entre dos conjuntos de variables que pueden ser un conjunto de variables de respuesta y un conjunto de variables explicativas, ambos conjuntos de variables continuas.

Desde el punto de vista del análisis de datos el ACC consiste en estudiar las relaciones lineales existentes entre dos grupos de variables cuantitativas observadas sobre un mismo conjunto de individuos.

Para la utilización del análisis de correlación canónico con propósitos descriptivos no se requieren supuestos distribucionales. En estos casos, las variables de predicción y las variables de respuesta pueden ser medidas en una escala nominal u ordinal. Sin embargo para probar la significancia de la relación entre las variables canónicas, los datos deben presentar los requerimientos de normalidad multivariada y homogeneidad de varianzas; es decir deben ser cuantitativas.

Como el ACC pretende poner en evidencia las relaciones lineales que existen entre los dos conjuntos de variables, la pregunta de interés es ¿cuál es el criterio para decir que las relaciones lineales entre el conjunto de variables explicativas y de respuesta son próximas?

Así, en el ACC se tienen los siguientes elementos:

$\{x^1, x^2, x^3, \dots, x^p\}$ y $\{y^1, y^2, y^3, \dots, y^q\}$ dos grupos de variables cuantitativas medidas sobre un conjunto I de individuos.

Estos grupos de variables definen matrices así, por ejemplo la matriz $X = [x^1 \ x^2 \ x^3 \ \dots \ x^p]$ está asociada a la aplicación X y es la matriz de datos con respecto al primer paquete de variables, de manera análoga $Y = [y^1 \ y^2 \ y^3 \ \dots \ y^q]$ es la matriz de datos para el segundo grupo de variables.

Como un caso particular se tienen el ACC clásico en el cual se considera que cada individuo tiene asociado un peso p_i , positivo y además $\sum_{i=1}^n p_i = 1$. Para cada paquete de variables se considera que cada variable está centrada, por lo que el centro de gravedad de la nube de individuos coincide con el origen.

El espacio de variables tiene asociada la métrica de pesos M

$$M = D_p = \begin{pmatrix} p_1 & 0 & 0 \\ 0 & \dots & 0 \\ 0 & 0 & p_n \end{pmatrix}$$

para este caso las matrices $V_{XX} = X^T D_p X$, $V_{YY} = Y^T D_p Y$, $V_{XY} = X^T D_p Y$, representan matrices de varianzas y covarianzas asociadas respectivamente a los paquetes de variables X y Y

entre ambos grupos.

Así, los factores canónicos pueden ser obtenidos a partir de las ecuaciones:

$$V_{X'X}^{-1} V_{X'Y} V_{Y'Y}^{-1} V_{Y'X} a = \lambda a$$

$$V_{Y'Y}^{-1} V_{Y'X} V_{X'X}^{-1} V_{X'Y} b = \lambda b$$

$$\text{con } \|a\|_{V_{X'X}} = 1 \text{ y } \|b\|_{V_{Y'Y}} = 1$$

de tal manera que las variables canónicas ξ^i y η^i verifican que $\xi^i = X a$, y $\eta^i = Y b$, y además si $\lambda_i \neq 0$, las ecuaciones de transición pueden ser expresadas como:

$$a_i = \frac{1}{\sqrt{\lambda_i}} V_{X'X}^{-1} V_{X'Y} b_i$$

$$b_i = \frac{1}{\sqrt{\lambda_i}} V_{Y'Y}^{-1} V_{Y'X} a_i$$

Además, como las variables canónicas son centradas y reducidas se tiene que $\text{corr}^2(\xi^i, \eta^i) = \lambda_i$, de lo que se deduce que el i -ésimo valor propio λ_i es el cuadrado del i -ésimo coeficiente de correlación canónica.

La interpretación geométrica del análisis de correlación canónico es la siguiente:

De manera general se tiene

$$a = \frac{1}{\sqrt{\lambda_i}} V_{X'X}^{-1} V_{X'Y} b$$

$$b = \frac{1}{\sqrt{\lambda_i}} V_{Y'Y}^{-1} V_{Y'X} a$$

Si se premultiplican los elementos de cada ecuación por X y Y respectivamente, se tiene:

$$Xa = \frac{1}{\sqrt{\lambda_i}} X (X' X)^{-1} X' Y b \quad \text{donde } V_{X'X}^{-1} = (X' X)^{-1} \text{ y } V_{X'Y} = X' Y$$

$$Yb = \frac{1}{\sqrt{\lambda_i}} Y (Y' Y)^{-1} Y' X a \quad \text{donde } V_{Y'Y}^{-1} = (Y' Y)^{-1} \text{ y } V_{Y'X} = Y' X$$

Llamemos φ_X y φ_Y los subespacios lineales de R^n , que son generados por las columnas de X y Y respectivamente. Las combinaciones lineales de a y b definen los puntos de φ_X y de φ_Y respectivamente, cuyas coordenadas son los componentes de los vectores Xa y Yb respectivamente. Por lo que las matrices

$$P_X = X (X' X)^{-1} X'$$

$$P_Y = Y (Y' Y)^{-1} Y'$$

son los operadores de proyección sobre φ_X y φ_Y respectivamente. Es decir cada vector es



colineal a la proyección del otro.

Como los vectores Xa y Yb son unitarios, se tiene que:

$$\sqrt{\lambda} = \cos \omega = \cos(\angle Xa, Yb)$$

Así, se tiene que la primera raíz canónica λ es el coseno al cuadrado de el ángulo más pequeño entre los subespacios ρ_X y ρ_Y .

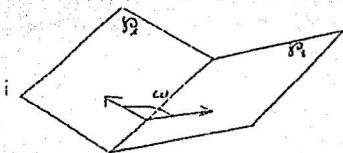


Fig 3.1. Interpretación geométrica del ACC

Hasta ahora se han mostrado en que consiste en ACC de manera general. Las ideas aquí presentadas nos van a permitir entender parte de la generalización de la técnica del análisis de correspondencias múltiples.

Para introducir la generalización del análisis de correspondencias múltiples (AFC-múltiples) presentamos algunas consideraciones

3.2 Consideraciones para el análisis de correspondencias múltiple.

- Se tienen Q factores. Para un factor q en particular se tiene un conjunto de p_q categorías.

Así entonces, el número total de categorías que se tienen es $p = \sum_{q=1}^Q p_q$

- El tamaño de muestra que es considerada para los Q factores es n .

Lo anterior queda ejemplificado a través de lo siguiente: cuando en algunos estudios estadísticos se tienen una gran cantidad de preguntas para las cuales las respuestas son categorías y además se tiene que para cada pregunta las categorías son mutuamente excluyentes; entonces para cada pregunta o factor que tiene k categorías se tienen a lo más k grupos en las cuales la muestra de tamaño n puede ser repartida. Es decir si preguntamos a 10 personas acerca de su religión (suponiendo que sólo consideramos católica, protestante, musulmana) entonces a lo más se pueden formar con la muestra los tres grupos referentes a cada religión o pudiera darse el caso de que los 10 fueran católicos o los 10 fueran protestantes o que los 10 fueran musulmanes.

3.2.1 Matriz indicadora $Z=[Z_1, Z_2, \dots, Z_q, \dots, Z_Q]$.

Como se ha dicho si tenemos Q factores, entonces la persona únicamente puede escoger una sola categoría para cada factor q $q=1,2,\dots,Q$. Esta información se recoge en la matriz Z , es decir, se denota por Z a la matriz con I renglones (número de individuos) y p (número de categorías) columnas que describe para cada individuo i de la muestra de tamaño I su elección o clasificación de acuerdo a las categorías de los factores.

La matriz Z es la yuxtaposición de las Q submatrices tales que

$$Z=[Z_1, Z_2, \dots, Z_q, \dots, Z_Q]$$

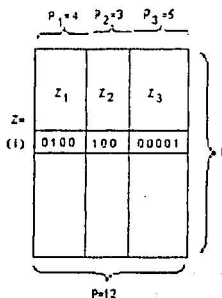


Fig 3.2. Matriz indicadora

La submatriz Z_q (la cual tiene I renglones y P_q columnas) es tal que su i -ésimo renglón tiene $P_q - 1$ veces el valor cero y una vez el valor uno, en la columna correspondiente a la categoría del factor q en el cual fue clasificado el sujeto i .

Esta matriz representa la información por individuo, pero si quisieramos trabajar con una matriz con datos condensados entonces podemos hacer uso de la llamada tabla de contingencia asociada a la matriz de Burt.

3.2.2 Tabla de Burt asociada con Z .

La tabla de Burt condensa la información generada por la matriz indicadora Z , y puede ser vista como bloques de matrices, donde los bloques son tablas de contingencia tomadas de dos en dos. Los bloques que contienen a los elementos de la diagonal son matrices diagonales debido a que corresponden a la tabla de contingencia de una variable en ella misma, por lo que los elementos de la diagonal son las frecuencias de dicha variable.

A la matriz $B=Z'Z$ se le conoce como tabla de contingencia de Burt asociada con Z y tiene la siguiente forma:

$$B = \begin{bmatrix} Z'_1 Z_1 & \dots & Z'_1 Z_q & \dots & Z'_1 Z_Q \\ \dots & \dots & \dots & \dots & \dots \\ Z'_q Z_1 & \dots & Z'_q Z_q & \dots & Z'_q Z_Q \\ \dots & \dots & \dots & \dots & \dots \\ Z'_Q Z_1 & \dots & Z'_Q Z_q & \dots & Z'_Q Z_Q \end{bmatrix}$$

Nótese que la matriz B está formada por Q^2 bloques. La q -ésima submatriz $Z'_q Z_q$ es una matriz diagonal cuyos elementos son las frecuencias para cada categoría del factor q . La submatriz $Z'_q Z_{q'}$ es la tabla de contingencia que clasifica la muestra de tamaño n en los factores q y q' .

La matriz diagonal D es una matriz cuadrada que tiene los mismos elementos en la diagonal que B ; esos elementos en la diagonal son las frecuencias para cada una de las categorías. La matriz D también tiene Q^2 bloques. Y además solamente las Q submatrices diagonales son diferentes a la matriz cero. De tal manera que el q -ésimo bloque $D_q = Z'_q Z_q$ es una matriz diagonal cuyos elementos son las frecuencias correspondientes a las categorías del factor q .

Ejemplo 1 : Supongamos que tenemos 10 personas a las cuales las clasificamos de acuerdo a los siguientes factores: Z_1 = sexo (masculino , femenino), Z_2 = estado civil (soltero, casado, viudo) y Z_3 = partido político de preferencia (PRI,PAN,PRD). La matriz Z que representa la información es la siguiente:

$$Z = \begin{array}{c} \begin{matrix} Z_1 & Z_2 & Z_3 \end{matrix} \\ \left[\begin{array}{ccc|c} 10 & 100 & 100 & 1 \\ 01 & 100 & 010 & 2 \\ 01 & 010 & 100 & 3 \\ 10 & 100 & 100 & 4 \\ 10 & 010 & 001 & 5 \\ 10 & 001 & 001 & 6 \\ 01 & 001 & 010 & 7 \\ 01 & 010 & 100 & 8 \\ 10 & 100 & 010 & 9 \\ 01 & 100 & 010 & 10 \end{array} \right] \end{array} \left. \vphantom{\begin{matrix} Z_1 & Z_2 & Z_3 \\ 1 \\ 2 \\ 3 \\ 4 \\ 5 \\ 6 \\ 7 \\ 8 \\ 9 \\ 10 \end{matrix}} \right\} \text{individuo}$$

En esta tabla se lee por ejemplo que para el individuo 4: es hombre, soltero y tiene preferencia por el PRI. Para el décimo individuo: es mujer, soltera y con preferencia por el PAN. Se observa además que de los 10 individuos hay 5 hombres y 5 mujeres; 5 solteros, 3 casados y 2 viudos; 4 con preferencias por el PRI, 4 con preferencias por el PAN y 2 que tienen preferencias por el PRD.

La matriz de Burt se obtiene de la siguiente manera:

$$Z' = \begin{bmatrix} Z'_1 \\ Z'_2 \\ Z'_3 \end{bmatrix}$$

$$B = \begin{bmatrix} 1 & 0 & 0 & 1 & 1 & 1 & 0 & 0 & 1 & 0 \\ 0 & 1 & 1 & 0 & 0 & 0 & 1 & 1 & 0 & 1 \\ 1 & 1 & 0 & 1 & 0 & 0 & 0 & 0 & 1 & 1 \\ 0 & 0 & 1 & 0 & 1 & 0 & 0 & 1 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 1 & 1 & 0 & 0 & 0 \\ 1 & 0 & 1 & 1 & 0 & 0 & 0 & 1 & 0 & 0 \\ 0 & 1 & 0 & 0 & 0 & 0 & 1 & 0 & 1 & 1 \\ 0 & 0 & 0 & 0 & 1 & 1 & 0 & 0 & 0 & 0 \end{bmatrix} \begin{bmatrix} 10 & 100 & 100 \\ 01 & 100 & 010 \\ 01 & 010 & 100 \\ 10 & 100 & 100 \\ 10 & 010 & 001 \\ 10 & 001 & 001 \\ 01 & 001 & 010 \\ 01 & 010 & 100 \\ 10 & 100 & 010 \\ 01 & 100 & 010 \end{bmatrix}$$

Z' Z

Haciendo las operaciones se tiene que

$$B = \begin{bmatrix} 5 & 0 & 3 & 1 & 1 & 2 & 1 & 1 \\ 0 & 5 & 2 & 2 & 1 & 2 & 3 & 0 \\ 3 & 2 & 5 & 0 & 0 & 2 & 3 & 0 \\ 1 & 2 & 0 & 3 & 0 & 2 & 0 & 1 \\ 1 & 1 & 0 & 0 & 2 & 0 & 1 & 1 \\ 2 & 2 & 1 & 2 & 0 & 1 & 0 & 0 \\ 1 & 3 & 3 & 0 & 1 & 0 & 1 & 0 \\ 2 & 0 & 0 & 1 & 0 & 0 & 0 & 2 \end{bmatrix}$$

Tabla 1. Matriz de Burt.

	Hombre	Mujer	Soltero	Casado	Viudo	PRI	PAN	PRD
Hombre	5	0	3	1	1	2	1	2
Mujer	0	5	2	2	1	2	3	0
Soltero	3	2	5	0	0	2	3	0
Casado	1	2	1	3	0	2	0	1
Viudo	1	1	0	0	2	0	1	1
PRI	2	2	2	2	0	3	3	1
PAN	1	3	3	0	1	0	4	0
PRD	2	0	0	1	1	0	0	2

Nótese que la matriz B está formada por Q^2 bloques, es decir el cruce del factor sexo consigo mismo y con los otros dos, el cruce del factor estado civil consigo mismo y con los otros dos y el

cruce del factor preferencia política, igual, consigo mismo y los otros dos, en total se tiene 9 bloques. La matriz $Z_1^1 Z_1$ (en rojo) es una matriz diagonal que representa la frecuencia de cada categoría del factor sexo (recuérdese que se tenían 5 hombres y 5 mujeres); la matriz $Z_2^1 Z_2$ (en azul) es una matriz diagonal que representa la frecuencia de cada categoría del factor estado civil (recuérdese que se tenían 5 solteros, 3 casados y 2 viudos), la matriz $Z_3^1 Z_3$ (en verde) es una matriz diagonal que representa la frecuencia de cada categoría del factor preferencia por el partido político (recuérdese que se tenían 4 personas con preferencias por el PRI, 4 con preferencias por el PAN y 2 que tienen preferencias por el PRD). El bloque $Z_1^1 Z_2^1$ es la tabla de contingencia que clasifica a la muestra de 10 personas con respecto a los factores sexo y estado civil; de tal manera que se lee en la tabla que de los 5 hombres en la muestra 3 eran solteros, 1 era casado y 1 viudo. De manera semejante se lee para el caso de las mujeres.

El bloque $Z_1^1 Z_3^1$ es la tabla de contingencia que clasifica a la muestra de 10 personas con respecto a los factores sexo y preferencia política; de tal manera que se lee en la tabla que de los 5 hombres en la muestra 2 preferían al PRI, 1 al PAN y 2 al PRD. De manera semejante se lee para el caso de las mujeres.

El bloque $Z_2^1 Z_3^1$ es la tabla de contingencia que clasifica a la muestra de 10 personas con respecto a los factores estado civil y preferencia política; de tal manera que se lee en la tabla que de las 5 personas solteras en la muestra 2 preferían al PRI, 3 al PAN y ninguna al PRD. De las 3 personas casadas 2 preferían al PRI, ninguna al PAN y 1 al PRD; y, de las dos personas viudas ninguna prefería al PRI, una al PAN y una al PRD.

Hasta aquí, se ha mencionado la forma en que se obtiene la información, resulta importante el hecho de que es posible hacer el análisis con los datos de la matriz indicadora (datos crudos, sin procesar) y haciendo uso de la matriz de Burt (datos procesados).

En lo que sigue se expondrán las equivalencias para el caso de dos factores del análisis de correspondencias utilizando las matrices que hemos descrito hasta el momento:

- Tabla de contingencia entre dos factores.
- Matriz Indicadora Z.
- Matriz de Burt.

Esto para mostrar las formas usuales en que se reporta el análisis de correspondencias para el caso de más de dos factores.

3.3 Análisis de correspondencias para una tabla de dos factores considerando diversas matrices.

La matriz de respuestas contenidas en la matriz Z es $Z = [Z_1 | Z_2]$

Es equivalente, desde el punto de vista de describir la relación entre categorías, llevar a cabo alguno de los siguientes análisis:

1. Un análisis de correspondencias de la matriz $Z_{1 \times p}$
2. Un análisis de correspondencias de la matriz de Burt $B_{p \times p}$
3. Un análisis de correspondencias de la tabla de contingencia $K_{1 \times 2} = Z_1' Z_2$
4. Un análisis canónico de los dos conjuntos de bloques Z_1 y Z_2 .

Para mostrar lo anterior a excepción del punto 4 sea lo siguiente:

Equivalencia entre el análisis 1 y 2.

El α -ésimo factor ϕ_α extraído del análisis 1 es el eigenvector de la siguiente ecuación:

$$\frac{1}{Q} D^{-1} Z' Z \phi_\alpha = \mu_\alpha \phi_\alpha$$

Lo anterior se muestra de lo siguiente; recordando la notación utilizada para el análisis de correspondencias utilizado en el capítulo II en el cual se estableció que el factor ϕ era el eigenvector de la matriz $A = D_c^{-1} F' D_r^{-1} F$

Haciendo

$$F = \frac{1}{\sqrt{Q}} Z$$

$$D_c = \frac{1}{\sqrt{Q}} D$$

$$D_r = \frac{1}{n} I_n \text{ con } I_n \text{ la matriz identidad.}$$

sustituyendo se tiene que:

$$A = D_c^{-1} F' D_r^{-1} F = \left(\frac{1}{\sqrt{Q}} D \right)^{-1} \frac{1}{\sqrt{Q}} Z' \left(\frac{1}{n} I_n \right)^{-1} \frac{1}{\sqrt{Q}} Z = \frac{1}{Q} D^{-1} Z' Z.$$

Así, los factores son :

$$\frac{1}{Q} D^{-1} Z' Z \phi_\alpha = \mu_\alpha \phi_\alpha \quad (3.1)$$

Lo anterior se ejemplifica con los datos del ejemplo I:

Se tiene la matriz Z que es la siguiente:

$Z =$

	Sexo		Estado civil			Preferencia política			Total
	Hombre	Mujer	Soltero	Viudo	Casado	PRI	PAN	PRD	
	1	0	1	0	0	1	0	0	3
	0	1	1	0	0	0	1	0	3
	0	1	0	1	0	1	0	0	3
	1	0	1	0	0	1	0	0	3
	1	0	0	1	0	0	0	1	3
	1	0	0	0	1	0	0	1	3
	0	1	0	0	1	0	1	0	3
	0	1	0	1	0	1	0	0	3
	1	0	1	0	0	0	1	0	3
	0	1	1	0	0	0	1	0	3
Total	5	5	5	3	2	4	4	2	30

entonces la matriz F está dada la expresión siguiente donde Q es el número de factores.

$$F = \frac{1}{1 \times 3} Z = \frac{1}{30} Z \quad \text{= matriz de correspondencias}$$

$$D_c = \frac{1}{10 \times 3} D = \begin{pmatrix} 5 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 5 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 5 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 3 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 2 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 4 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 4 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 2 \end{pmatrix}$$

la matriz diagonal de pesos por columnas

$$D_r = \frac{1}{10} / 10 = \begin{pmatrix} \frac{1}{10} & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & \frac{1}{10} & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & \frac{1}{10} & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & \frac{1}{10} & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & \frac{1}{10} & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & \frac{1}{10} & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & \frac{1}{10} & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & \frac{1}{10} & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & \frac{1}{10} & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & \frac{1}{10} \end{pmatrix}$$

Ahora, consideremos la matriz B tal que la matriz $B=Z'Z$. B es simétrica. Además los renglones y las columnas marginales son los elementos de la matriz QD .

En el análisis de la matriz B

Sea la matriz

$$F = \frac{1}{nQ^2} B$$

las matrices

$$D_c = D_r = \frac{1}{nQ} J$$

Así la matriz A ser diagonalizada es:

$$A = D_c^{-1} F' D_r^{-1} F = \left(\frac{1}{nQ} D\right)^{-1} \frac{1}{nQ^2} B' \left(\frac{1}{nQ} D\right)^{-1} \frac{1}{nQ^2} B = \frac{1}{Q^2} D^{-1} B' D^{-1} B$$

si premultiplicamos $\frac{1}{Q} D^{-1} Z' Z \phi_a = \mu_a \phi_a$ por $\frac{1}{Q} D^{-1} B$ obtenemos lo siguiente:

$$\left(\frac{1}{Q} D^{-1} B\right) \frac{1}{Q} D^{-1} Z' Z \phi_a = \left(\frac{1}{Q} D^{-1} B\right) \mu_a \phi_a$$

\Leftrightarrow

$$\frac{1}{Q^2} D^{-1} B D^{-1} B \phi_a = \left(\frac{1}{Q} D^{-1} B\right) \mu_a \phi_a, \text{ pero } \left(\frac{1}{Q} D^{-1} B\right) \text{ es un eigenvector para } (3.1)$$

por lo que

$$\frac{1}{Q^2} D^{-1} B D^{-1} B \phi_a = \mu_a^2 \phi_a$$

por lo que los factores son idénticos para ambos análisis.

La matriz de Burt para los datos del ejemplo 1 son:

	Hombre	Mujer	Soltero	Casado	Viudo	PRÍ	PAN	PRD	Total
Hombre	5	0	3	1	1	2	1	2	15
Mujer	0	5	2	2	1	2	3	0	15
Soltero	3	2	5	0	0	2	3	0	15
Casado	1	2	0	3	0	2	0	1	9
Viudo	1	1	0	0	2	0	1	1	6
PRÍ	2	2	2	2	0	4	0	0	12
PAN	1	3	3	0	1	0	4	0	12
PRD	2	0	0	1	1	0	0	2	6
Total	15	15	15	9	6	12	12	6	90

por lo que las matrices son

$$F = \frac{1}{10 \times 3^2} B = \frac{1}{90} B$$

← matriz de correspondencias

$$D_c = D_l = \frac{1}{10 \times 3} D = \frac{1}{30} D = \begin{pmatrix} \frac{5}{30} & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & \frac{5}{30} & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & \frac{5}{30} & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & \frac{1}{30} & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & \frac{2}{30} & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & \frac{1}{30} & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & \frac{1}{30} & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & \frac{1}{30} & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & \frac{2}{30} \end{pmatrix}$$

Nótese que la matriz anterior representa la matriz diagonal de pesos por renglón y a la matriz diagonal de pesos por columna. Y se obtiene si a la columna y renglón marginal se divide entre la suma total. Igual que se hacía en el análisis de correspondencias simple expuesto en el capítulo II.

Las equivalencias anteriores fueron presentadas utilizando 3 factores como caso particular cuando se tienen dos factores está la siguiente equivalencia.

Ahora, la equivalencia entre el análisis I y el análisis 3

Se mostrará que para cada par de factores $(\varphi_\alpha, \psi_\alpha)$ relativos al mismo eigenvalor λ_α el cual es extraído del análisis de la tabla de contingencia $K=Z_1'Z_2$, hay un factor ϕ_α del análisis Z (o B) tal que:

$$\phi_\alpha = \begin{bmatrix} \varphi_\alpha \\ \psi_\alpha \end{bmatrix}$$

Recordemos que la matriz $D_1 = Z_1'Z_1$ y la matriz $D_2 = Z_2'Z_2$ son las matrices diagonales que representan el renglón y columna marginales de la tabla de contingencia.

$$D = \begin{bmatrix} D_1 & 0 \\ 0 & D_2 \end{bmatrix}$$

Recordemos que las formulas de transición fueron establecidas de la siguiente manera:

$$\varphi_\alpha = \frac{1}{\sqrt{\lambda_\alpha}} D_1^{-1} Z_1' Z_2 \psi$$

(3.2)

$$\psi_\alpha = \frac{1}{\sqrt{\lambda_\alpha}} D_2^{-1} Z_2' Z_1 \varphi_\alpha$$

De (3.2) tenemos que $\varphi_\alpha \sqrt{\lambda_\alpha} = D_1^{-1} Z_1' Z_2 \psi_\alpha$ y sumando a ambos lados φ_α y agrupando tenemos lo siguiente:

$(1 + \sqrt{\lambda_\alpha}) \varphi_\alpha = D_1^{-1} (D_1 \varphi_\alpha + Z_1' Z_2 \psi_\alpha)$, haciendo lo mismo con la otra ecuación se tiene lo siguiente:

$$(1 + \sqrt{\lambda_\alpha}) \psi_\alpha = D_2^{-1} (D_2 \psi_\alpha + Z_2' Z_1 \varphi_\alpha)$$

Este sistema de ecuaciones en forma matricial se reescribe de la siguiente manera

$$\begin{bmatrix} D_1 & 0 \\ 0 & D_2 \end{bmatrix}^{-1} \begin{bmatrix} D_1 & Z_1' Z_2 \\ Z_2' Z_1 & D_2 \end{bmatrix} \begin{bmatrix} \varphi_\alpha \\ \psi_\alpha \end{bmatrix} = (1 + \sqrt{\lambda_\alpha}) \begin{bmatrix} \varphi_\alpha \\ \psi_\alpha \end{bmatrix}$$

esto a su vez representa lo siguiente:

$$D^{-1} Z' Z \begin{bmatrix} \varphi_\alpha \\ \psi_\alpha \end{bmatrix} = (1 + \sqrt{\lambda_\alpha}) \begin{bmatrix} \varphi_\alpha \\ \psi_\alpha \end{bmatrix}$$

\Leftrightarrow

$$D^{-1} Z' Z \phi_\alpha = (1 + \sqrt{\lambda_\alpha}) \phi_\alpha$$

Como estamos tratando con dos factores, y como se mostró en la ecuación (3.1) para el análisis de la matriz Z, entonces haciendo Q=2 se tiene finalmente:

$$\frac{1}{Q} D^{-1} Z' Z \phi_\alpha = \left(\frac{1 + \sqrt{\lambda_\alpha}}{2} \right) \phi_\alpha$$

lo cual tiene los mismos factores que el análisis 1, sin embargo el eigenvalor ahora es:

$$\mu_\alpha = \frac{1 + \sqrt{\lambda_\alpha}}{2}$$

Así, si λ_α es el eigenvalor más grande obtenido del análisis de la tabla de contingencia K entonces la ecuación anterior muestra el α -ésimo eigenvalor más grande del análisis de la matriz Z.

Por transitividad los tres análisis son equivalentes, por lo tanto hacer un análisis de correspondencias de la tabla de contingencias es equivalente a realizar un análisis de correspondencias de la matriz indicadora y es equivalente también a efectuar un análisis de correspondencias de la tabla de Burt.

En la sección siguiente se hará una generalización para resolver el problema de encontrar los componentes y los ejes para el caso de más de dos factores, la idea que se pretende aprovechar es la generalización del análisis de correlación canónico a través del cual se concluye que los eigenvalores se obtienen de la utilización de la matriz de Burt B y como se mostró en el apartado

anterior son proporcionales a los que se obtienen utilizando la matriz Z.

3.4 Generalización del análisis de correspondencias utilizando la matriz Z, para más de dos factores.

De la figura 1 se observa que $Z = [Z_1, Z_2, \dots, Z_Q]$ tiene P columnas que es el número de categorías totales y cada columna corresponde a un punto en R^n . Consideremos el espacio R^n . Cada submatriz Z_q genera un subespacio lineal ζ_q con p_q dimensiones. El rango de la matriz Z es a lo más igual a $P - (Q - 1)$ ya que todos esos subespacios tienen en común al menos el primer bisector (el vector cuyos componentes son iguales a uno).

Sea φ_q el vector cuyas p_q componentes son las coordenadas de un punto m_q de ζ_q en la base definida por las columnas de Z_q .

Las coordenadas de m_q en R^n son los componentes de $m_q = Z_q \varphi_q$

La distancia al cuadrado de este punto m_q al origen es :

$$\varphi_q' Z_q' Z_q \varphi_q = \varphi_q' D_q \varphi_q$$

Así, el análisis de correspondencias de la tabla de contingencia para dos preguntas q y q' se reduce a estudiar las posiciones relativas de los subespacios ζ_q y $\zeta_{q'}$. Y esto precisamente corresponde al análisis canónico de la matriz $[Z_q | Z_{q'}]$.

Por lo que las ecuaciones de transición quedan de la siguiente forma:

$$\varphi_q = \frac{1}{\sqrt{\lambda}} (Z_q' Z_q)^{-1} Z_q' Z_{q'} \varphi_{q'} = \frac{1}{\sqrt{\lambda}} D_q^{-1} Z_q' Z_{q'} \varphi_{q'}$$

$$\varphi_{q'} = \frac{1}{\sqrt{\lambda}} (Z_{q'}' Z_{q'})^{-1} Z_{q'}' Z_q \varphi_q = \frac{1}{\sqrt{\lambda}} D_{q'}^{-1} Z_{q'}' Z_q \varphi_q$$

De manera análoga a lo hecho para el análisis canónico si premultiplicamos por Z_q y $Z_{q'}$ respectivamente tenemos lo siguiente:

$$Z_q \varphi_q = \frac{1}{\sqrt{\lambda}} Z_q (Z_q' Z_q)^{-1} Z_q' Z_{q'} \varphi_{q'}$$

$$Z_{q'} \varphi_{q'} = \frac{1}{\sqrt{\lambda}} Z_{q'} (Z_{q'}' Z_{q'})^{-1} Z_{q'}' Z_q \varphi_q$$

Donde se obtiene que los operadores de proyección sobre ζ_q y $\zeta_{q'}$ son:

$$P_q = Z_q (Z_q' Z_q)^{-1} Z_q'$$

$$P_{q'} = Z_{q'} (Z_{q'}' Z_{q'})^{-1} Z_{q'}'$$

lo que finalmente proporciona las coordenadas de los puntos m_q y $m_{q'}$ para ambos

subespacios:

$$m_q = \frac{1}{\sqrt{k}} P_q m_{q'}$$

$$m_{q'} = \frac{1}{\sqrt{k}} P_{q'} m_q$$

De manera análoga a lo establecido en el análisis canónico la proyección ortogonal de m_q sobre $\varphi_{q'}$ es colineal a $m_{q'}$ y viceversa.

Una generalización para el problema del análisis canónico cuando se tienen más de dos factores.

Sean $\varphi_1, \varphi_2, \dots, \varphi_Q$, respectivamente, los vectores de los componentes de los Q puntos, m_1, m_2, \dots, m_Q en las bases Z_1, Z_2, \dots, Z_Q , y sea $m = m_1 + m_2 + \dots + m_Q$

La cantidad a ser maximizada es

$$\|m\|^2 = \sum \{ \varphi_q' Z_q' Z_q \varphi_q \mid q \in Q \text{ y } q' \in Q \}$$

con la restricción

$$\sum \{ \varphi_q' D \varphi_q \mid q \in Q \} = Qn$$

Si ϕ es el vector con p componentes definido por

$$\phi' = \{ \varphi_1, \varphi_2, \dots, \varphi_Q \}$$

El problema llega a ser el

$$\max \phi' B \phi$$

$$\text{s.a } \phi' D \phi = Qn$$

Los factores requeridos ϕ son los eigenvectores de la matriz $D^{-1}B$, relativos a los eigenvalores más grandes. Ellos son proporcionales a aquellos que son extraídos del análisis de correspondencias de la matriz Z y coinciden, con los extraídos del análisis de la matriz B.

3.5 Resultados del análisis de correspondencias múltiples utilizando la matriz Z.

Recordemos que los ϕ factores extraídos del análisis de la matriz Z son tales que:

$$\frac{1}{Q} D^{-1} B \phi = \mu \phi,$$

lo cual reorganizando los términos de esta ecuación para mostrar los componentes φ_q de ϕ relativos al factor q, de acuerdo a los bloques de las matrices D y B, se tiene:

$$\frac{1}{Q} \sum_{q=1}^Q D_q^{-1} Z_q' Z_q \varphi_q = \mu \varphi_q$$

Los Q subconjuntos de puntos correspondientes a las p_q categorías del factor q tienen el mismo centro de gravedad, el cual es también el centro de gravedad del conjunto entero de puntos. J_q denota el subconjunto de los p valores del índice j correspondientes al factor q (J_q tiene p_q elementos)

Las coordenadas del subconjunto de puntos relativos al factor q son las columnas de $Z_q D_q^{-1} = Z_q (Z_q' Z_q)^{-1}$ y los elementos de la diagonal de $\frac{1}{T} D_q$ son las masas relativas de los p_q puntos del subconjunto q.

La i-ésima coordenada del centro de gravedad G_q es

$$g_{qi} = \sum_{j \in J_q} \frac{d_{ij} z_{ij}}{T} = \frac{1}{T} \quad \text{ya que } \sum_{j \in J_q} z_{ij} = 1$$

De tal manera que g_{qi} es independiente de q, es decir $g_{qi} = g_i$.

Los φ_q factores correspondientes a los factores no triviales están centrados porque esos factores corresponden a un análisis del conjunto de puntos después de la traslación del origen a G.

Algunos resultados del análisis de correspondencias múltiples.

1. La suma de los pesos de las columnas $Z_q = \frac{1}{Q}$.

Es decir cada factor recibe el mismo peso, el cual es distribuido sobre todas las categorías de acuerdo a la frecuencia de cada categoría del factor

Por ejemplo:

$$\text{El peso para la submatriz } Z_1 = \frac{1}{3} = 0.16 + .016 = 0.33333$$

$$\text{El peso para la submatriz } Z_2 = \frac{1}{3} = 0.16 + 0.1 + .06 = 0.33333$$

$$\text{El peso para la submatriz } Z_3 = \frac{1}{3} = 0.13 + .013 + .06 = 0.33333$$

2. El centroide de los perfiles columna de Z_q es el centro de la gráfica que es el de todos los perfiles columna, de tal manera que la nube de los perfiles para cada factor está balanceada en el origen.

3. La inercia total de los perfiles columna y renglón es:

$$\ln(P) = \frac{1}{Q-1}$$

4. La inercia de los perfiles columna de Z_q es:

$$\ln(J_q) = \frac{c_q - 1}{Q}$$

De esto se deduce que la inercia con la que contribuye un factor se incrementa linealmente con el número de categorías del factor.

5. La inercia de una categoría en particular (j) es:

$$\ln(j) = \frac{1}{Q} - c_j^2$$

De esto se deduce que la inercia con la que contribuye una categoría incrementa en la medida que decrece la frecuencia de esta categoría, con un límite superior de $\frac{1}{Q}$

6. El número de dimensiones no triviales con inercia positiva es a lo más J-Q, es decir el número de dimensiones de los perfiles por renglón y por columnas es a lo más J-Q.

7. Los perfiles renglón están en el centroide de las perfiles por columna, al reescalar por el inversa de la raíz cuadrada de las inercias principales a través de los ejes respectivos.

Al igual que en el capítulo II nuestro problema consiste en encontrar una representación gráfica que nos permita apreciar en un espacio de menor dimensión, los perfiles de las categorías para poder estudiar posibles asociaciones entre los factores, por tal razón surge de nuevo la pregunta de ¿cuántas gráficas puedo obtener y cuál es la calidad de cada una?

3.6 Dimensionalidad de la configuración de las p categorías en R^n

Las coordenadas de las categorías en R^n son las columnas de ZD^{-1} . Ellas generan un subespacio cuya dimensión es el rango de ZD^{-1} y por lo tanto el rango de $Z = [Z_1, Z_2, \dots, Z_q, \dots, Z_Q]$

Todos los p_q subespacios generados por las columnas de las Z_q matrices tienen en común el primer bisector Δ , por lo tanto el rango máximo de Z es

$$p_1 + (p_2 - 1) + \dots + (p_Q - 1) = p - Q + 1$$

Por lo tanto el rango máximo de $D^{-1}Z'Z$, la matriz a ser diagonalizada es $p-Q+1$. pero el análisis con respecto al origen O, el primer bisector Δ es el eigenvector correspondiente al eigenvalor 1, en el análisis con respecto al centro de gravedad G, (p-Q) se encuentran (p-Q) eigenvalores diferentes de cero.

Así, escogiendo una base en el subespacio del conjunto de puntos, el problema puede ser reducido a encontrar eigenvalores y eigenvectores para una matriz de rango (p-Q).

3.7 La mejor representación simultánea para los individuos de la matriz Z

Nosotros buscamos las abscisas de los I individuos y las p categorías en el mismo eje tales que:

1. La abscisa de un individuo particular i es la media aritmética de su respuesta.
2. La abscisa de una categoría j es la media aritmética de las abscisas de los individuos que escogieron esto.

Como se ha mencionado las formulas de transición obtenidas para el análisis de la matriz Z son:

$$\Psi = \frac{1}{\sqrt{p}} \frac{1}{Q} Z \phi$$

$$\phi = \frac{1}{\sqrt{I}} \frac{1}{Q} D^{-1} Z' \Psi$$

donde Ψ_i es la abscisa del individuo i , y ϕ_j de la categoría j .

3.8 Ejemplo introductorio. Una forma alternativa de trabajar con más de dos factores.

El siguiente ejemplo aparece en el libro Applied Multivariate Data Analysis, Jobson, y establece lo siguiente:

Los datos que se presentan en la tabla 3.1 fueron obtenidos de los archivos de un estudiante de la carrera de servicios de advertencias legales para pobres y establece las frecuencias observadas para el número de personas a las que se les imputaba por un delito (cargo) y el resultado de la querrela es decir si fueron sentenciados o absueltos por el delito por el que se les acusaba. Esta información estaba disponible tanto para hombres como para mujeres. El objetivo era examinar la asociación entre el cargo por el que se acusa a una persona y el resultado penal para hombres y mujeres. Véase Tabla 3.1

Tabla 3.1 Tabla de contingencia. Frecuencias observadas para personas acusadas por un cargo penal y el veredicto del juez, tanto para hombres como para mujeres

Veredicto del juez	Sexo	Cargo penal				Total
		Daños impulsivos	Robo < \$1000	Desorden público	Poseción de narcóticos	
Absuelto	Hombre	8	11	5	7	43
	Mujer	5	15	3	1	30
Sentenciado	Hombre	105	32	11	23	208
	Mujer	32	57	8	2	122
Total		150	115	25	33	403

Una forma de averiguar esa relación sería el estudio de manera simultánea de los tres factores,

lo cual puede ser hecho con el análisis de correspondencias múltiples. Otra forma es la siguiente, se pueden concatenar algunas variables, por ejemplo el veredicto del juez y el sexo de la persona, de tal manera que se tiene la siguiente tabla de contingencia Tabla 3.2.

Tabla 3.2. Tabla de contingencia. Frecuencias observadas para personas acusadas por un cargo penal y concatenación de sexo y veredicto del juez.

Sexo/Veredicto del juez	Cargo Penal					Total
	Daños impulsivos	Robo <\$1000	Desorden público	Poseción de narcóticos	Otros	
Hombre/Absuelto	8	11	5	7	12	43
Mujer/Absuelta	5	15	3	1	6	30
Hombre/Condenado	105	32	11	23	37	208
Mujer/Condenada	32	57	6	2	25	122
Total	150	115	25	33	80	403

De la tabla anterior se observa que para el factor Sexo/veredicto del juez y para el factor cargo penal se tienen cuatro y cinco categorías respectivamente. Los resultados de la prueba Ji cuadrada de independencia son:

$\chi^2 = 70.9167$ que se distribuye $\chi^2_{(12)}$ y el cuantil de orden 0.95 de una distribución Ji-cuadrada con 12 grados de libertad es 21.0261 por lo que se rechaza la hipótesis de independencia entre los renglones y las columnas de la tabla de contingencia; además el nivel de significancia descriptivo $p\text{-value} = \Pr(\chi^2 > 70.9167) \cong 0.0000$. Por lo que se concluye que existe asociación entre las categorías.

Como se realizó en el capítulo II, un estudio gráfico exploratorio podría ser de utilidad. Así se tienen las siguientes gráficas con respecto a los perfiles con desviación del origen.

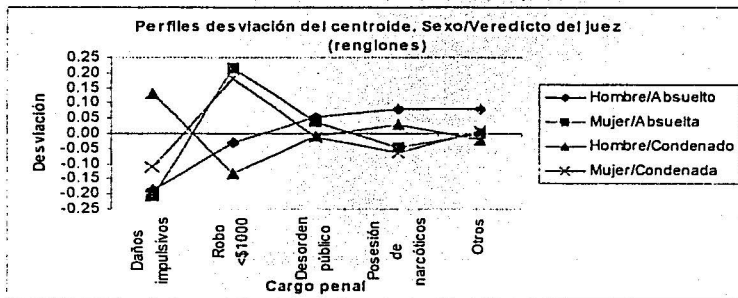


Fig 3.1

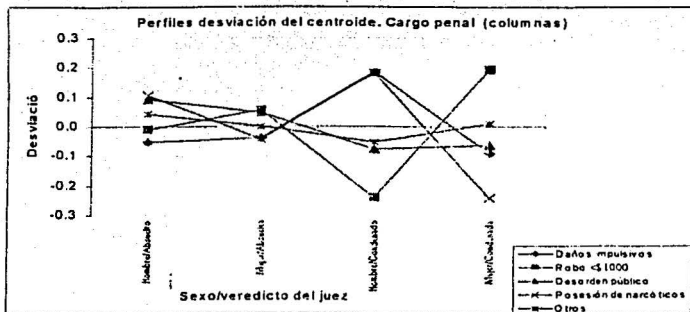


Fig 3.2

En la Fig 3.1 se observa que para el caso en que los hombres son condenados la desviación mayor se encuentra en la categoría de daños impulsivos, en contraparte con robo < \$1000. Para el caso de las mujeres que si son sentenciadas se observa que la desviación mayor aparece en la categoría de robo < \$1000 en contraparte con posesión de narcóticos.

De la Fig. 3.2 se observa que para el caso de posesión de narcóticos y daños impulsivos los hombres si son condenados, en contraparte a con robo < \$1000. De manera análoga para robo abajo < \$1000 las mujeres son mayoritariamente sentenciadas mientras que ocurre lo contrario para las mujeres son condenadas cuando se les enfrenta por "otros" tipos de cargos.

Estas dos gráficas dan una idea general de cómo es la relación entre las categorías de los factores.

Haciendo uso del análisis de correspondencias (AFC) se tiene lo siguiente:

Tabla 3.3 Valores propios y porcentaje de la inercia.

Valor	Valores propios		
	1	2	3
% de la inercia	0.1419	0.0329	0.0012
% acumulado	81%	19%	1%

De la tabla 3.3 se tiene que los tres primeros ejes factoriales recogen el 100% de la inercia total, y si nos fijamos en los dos primeros ejes estos recogen el 99%. Por lo que de acuerdo a lo establecido en el capítulo II, la gráfica con los dos primeros ejes explica el 99% de la dispersión de los datos, la cual está repartida en el primer eje con un 81% y con el segundo eje con un 19%.

La gráfica del AFC es la siguiente:

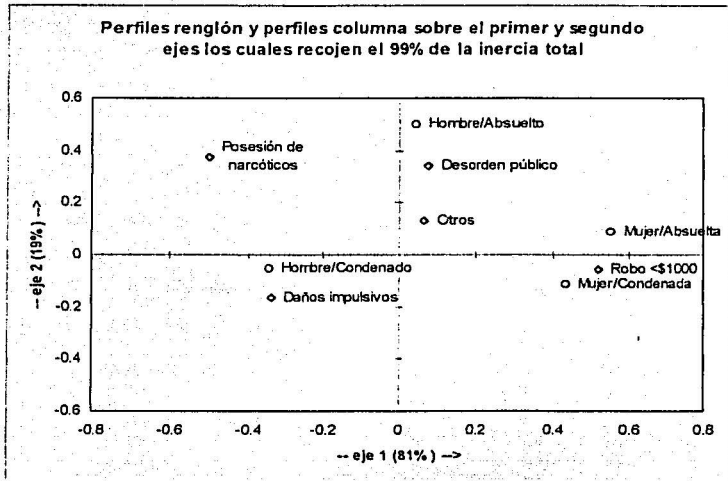


Fig 3.3

Recordemos que la interpretación de la gráfica se hace a partir de la posición de los perfiles en los ejes y la interpretación de los ejes tiene que ver con qué tan dispersos están los perfiles sobre éste.

Por lo que la interpretación de la gráfica puede ser la siguiente:

1) Tomando en cuenta el primer eje:

a) Considerando los perfiles por renglón se observan dos grupos; por una parte los hombres que son condenados y por otro lado las mujeres para las cuales no existe gran diferencia entre ser condenadas o ser absueltas.

b) En cuanto a los perfiles por columna se observa también que existen dos grupos que contrastan de manera fuerte, por un lado los cargos de posesión de narcóticos y daños impulsivos y por el otro el robo < \$1000.

2) Tomando en cuenta el segundo eje:

a') Se observa para los hombres que son condenados, las mujeres que son condenadas y las

mujeres que son absueltas, se encuentran muy cerca del centroide (origen), mientras que los hombres que son absueltos presentan un alejamiento considerable del centroide.

b') Se observa que robo abajo <\$1000 y otros se encuentran muy cerca del centroide, mientras que posesión de narcóticos y desorden público .

3) Tomando en cuenta los perfiles por renglón y por columnas de manera simultánea, se observa una relación entre los hombres que son condenados y los cargos de daños impulsivos y posesión de narcóticos en el sentido de alejamiento del origen, esta relación está en contraste con las mujeres que son condenadas por robo < \$1000 y que también son sentenciadas por el mismo cargo. Con respecto al segundo eje parece que los hombres que son sentenciados están relacionados con posesión de narcóticos y daños públicos en contraste con daños impulsivos.

Lo anterior sugiere lo siguiente: los hombres tienden a tener una fuerte asociación en las tasas de condena por daños impulsivos, mientras que las mujeres están relacionadas con el robo <\$1000.

Esta interpretación debe de ser sostenida por las contribuciones y la calidad de la representación de los perfiles.

Tabla 3.4 Componentes de los perfiles renglón sobre los ejes factoriales

Perfiles renglón	Peso	Inercia	% inercia	Coordenadas en los ejes		
				1	2	3
Hombre/Absuelto	0.1067	0.0273	0.1554	0.0443	0.5035	-0.0278
Mujer/Absuelta	0.0744	0.0243	0.1383	0.5541	0.0692	0.1067
Hombre/Condenad	0.5161	0.0627	0.3561	-0.3447	-0.0503	0.0053
Mujer/Condenada	0.3027	0.0616	0.3503	0.4359	-0.1137	-0.0262

De la tabla 3.4 se observa que los perfiles Hombre/condenado y mujer/condenada, tienen mayor peso y además esos perfiles aportan la mayor cantidad de inercia.

Tabla 3.5 Componentes de los perfiles columna sobre los ejes factoriales

Perfiles columna	Peso	Inercia	% inercia	Coordenadas en los ejes		
				1	2	3
Daños impulsivos	0.3722	0.0528	0.2589	-0.3364	-0.1634	0.0082
Robo <\$1000	0.2854	0.0787	0.4473	0.5219	-0.0582	0.0041
Desorden público	0.0620	0.0083	0.0471	0.0751	0.3420	0.1053
Posesión de narcót	0.0819	0.0317	0.1802	-0.4982	0.3728	-0.0138
Otros	0.1985	0.0047	0.0265	0.0663	0.1293	-0.0485

Para la tabla 3.5 se observa que los perfiles que más peso tienen son Daños impulsivos, Robo<\$1000 y otros. Aquí mientras Daños impulsivos , Robo<\$1000 aportan una gran cantidad de inercia, el perfil Otros no aporta mucho.

Ahora analicemos las contribuciones que cada perfil aportan a cada eje, recuérdese que el primer eje aportaba la mayor cantidad de inercia y es interesante saber como está distribuida esa inercia entre los perfiles.

Tabla 3.6 Contribuciones (en porcentaje) de los perfiles renglón a la inercia asociada a cada eje factorial

Perfiles renglón	Ejes		
	1	2	3
Hombre/Absuelto	0.0015	0.8232	0.0686
Mujer/Absuelta	0.1610	0.0180	0.7465
Hombre/Condenada	0.4322	0.0397	0.0120
Mujer/Condenada	0.4053	0.1191	0.1729

Tabla 3.7 Contribuciones (en porcentaje) de los perfiles columna a la inercia asociada a cada eje factorial

Perfiles columna	Ejes		
	1	2	3
Daños impulsivos	0.3004	0.3024	0.0208
Robo <\$1000	0.5478	0.0294	0.0041
Desorden público	0.0025	0.2209	0.5727
Posesión de narcót	0.1432	0.3464	0.0130
Otros	0.0061	0.1010	0.3894

De la tabla 3.6 para el primer eje principal se tiene que los perfiles Hombre/condenado y Mujeres/condenada son los que más aportan al primer eje (un 43.22% y un 40.53% respectivamente), y además son los que determinan su orientación. Para el segundo eje se tiene Hombre/ absuelto (82.32%) es el perfil que más contribuye y por lo mismo el que determina la orientación del segundo eje (Véase Fig 3.3).

De manera análoga para la tabla 3.7 en el primer eje se tiene que Daños impulsivos, Robo< \$1000 contribuyen con mayor cantidad de inercia (30.04% y 54.78% respectivamente), mientras que para el segundo eje se tiene que Daños impulsivos, Posesión de narcóticos y Desorden público (30.24%, 22.09% y 34.64% respectivamente) son los que contribuyen a la orientación del eje. (Véase Fig 3.3).

Ahora bien, la calidad de representación de los perfiles en la gráfica es determinada por los $\cos^2\theta$ o las correlaciones al cuadrado, recuérdese que si $\cos^2\theta$ es grande , entonces θ es pequeño y puede decirse que el perfil va en dirección del eje o que está correlacionado con el eje factorial, estos elementos se presentan en las siguientes tablas.

Tabla 3.8 Cosenos cuadrados de los ángulos de los perfiles renglón con respecto a los ejes factoriales

Perfiles renglón	Ejes		
	1	2	3
Hombre/Absuelto	0.0077	0.9893	0.0030
Mujer/Absuelta	0.9389	0.0243	0.0368
Hombre/Condenada	0.9790	0.0208	0.0002
Mujer/Condenada	0.9331	0.0635	0.0034

Tabla 3.9 Cosenos cuadrados de los ángulos de los perfiles columna con respecto a los ejes factoriales

Perfiles columna	Ejes		
	1	2	3
Daños impulsivos	0.8106	0.1889	0.0005
Robo <\$1000	0.9877	0.0123	0.0001
Desorden público	0.0421	0.8750	0.0829
Poseción de narcót	0.6407	0.3588	0.0005
Otros	0.1873	0.7124	0.1003

A partir de la tabla 3.8 se puede decir que el primer eje factorial guarda una fuerte correlación con respecto a los perfiles Mujer/condenada, Hombre/condenado y Mujer/absuelta; mientras que el segundo eje se encuentra fuertemente correlacionado con el perfil Hombre/absuelto; esto corrobora lo que ya se había dicho en cuanto a la dirección de los ejes

Similarmente, para la tabla 3.9 se observa que Daños impulsivos, Robo < \$1000 están altamente correlacionados con el primer eje factorial, mientras que Desorden público se encuentra muy correlacionado con el segundo eje.

Además, como se recordará en el capítulo II, la calidad de representación de cualquier perfil en la gráfica es la suma de las correlaciones al cuadrado, por lo que se puede decir que la calidad de representación de cualquier perfil tanto para los renglones como para las columnas es muy buena con los dos primeros ejes.

Este ejemplo muestra una forma particular de hacer un estudio de la relación concatenando dos factores, sin embargo el problema que ahora surge es el siguiente, ¿por qué no concatenar, por ejemplo, el sexo de la persona y el cargo penal y estudiar la asociación con el factor veredicto del juez, o concatenar el veredicto del juez con el cargo penal y estudiar la asociación con respecto al sexo de la persona?. Este tipo de estudios nos llevaría al estudio gráfico de la asociación condicional y parcial, pero que pasaría con la asociación total, en la que se utilicen de manera simultánea los tres factores.

Una forma alternativa para llevar a cabo el estudio de la asociación en el que se involucren los tres factores de manera simultánea es mediante el uso del análisis de correspondencias múltiples.

Del cual haremos su introducción a continuación:

Utilización del análisis de correspondencias múltiples (AFCM)

El objetivo es determinar si existe asociación o no entre las categorías de los tres factores, sexo, resultado de la sentencia y el cargo por el que se acusa.

Para iniciar la interpretación del análisis de correspondencias múltiple se hacen las siguientes interrogantes que se tratarán de contestar mediante el análisis:

1. ¿Quiénes son más sentenciados los hombres o las mujeres?
2. Para cada tipo de delito, ¿son igualmente sentenciados hombres y mujeres?
3. ¿Qué delito es más sentenciado?

El procedimiento para la interpretación del Análisis Factorial de Correspondencias Múltiples (AFCM), principalmente considera:

1. El número de ejes que se utilicen para la construcción del gráfico.

De acuerdo a lo establecido hasta el momento se tiene que el número total de ejes que recuperan el 100% de la estructura de la información es 6, por lo que debemos escoger entre estos seis ejes aquellos que recuperen la mayor cantidad de inercia. De manera análoga a lo establecido en el análisis de correspondencias simple se escogen aquellos ejes que conserven la mayor cantidad de inercia y además se considera que las categorías de los factores estén bien representadas en la gráfica. Así, para hacer la interpretación utilizaremos el subespacio generado por un solo eje (una línea recta) tomando en cuenta la proporción de inercia que se representa en los eigenvalores.

Valeurs pr.	1	2	3	4	5	6
Valeur	0.4626	0.3658	0.3333	0.3333	0.2745	0.2105
% de vari.	23%	19%	17%	17%	14%	11%
% cumulé	23%	42%	59%	76%	89%	100%

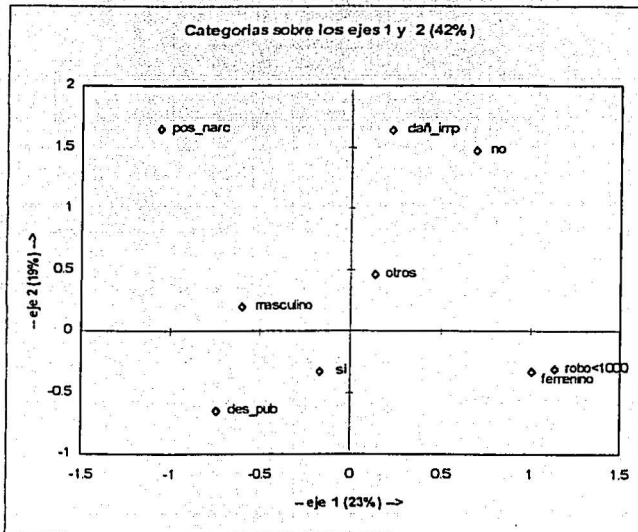
El subespacio generado por el primer eje recupera el 23% de la inercia total, el segundo el 19% y así sucesivamente. Considerar cualquiera de estos seis posibles subespacios no ayuda a establecer la relación de asociación que pudiera existir debido a la poca información que recogen los ejes.

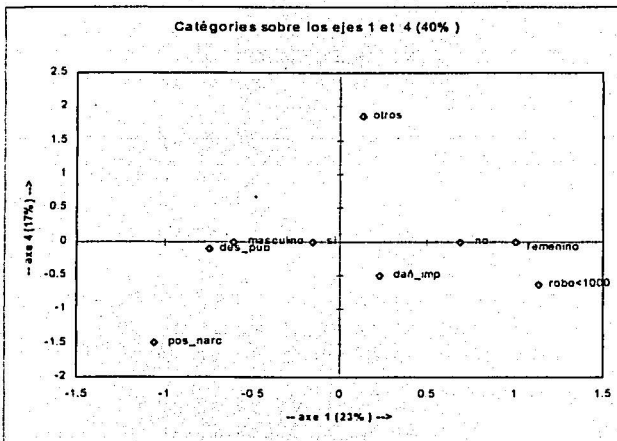
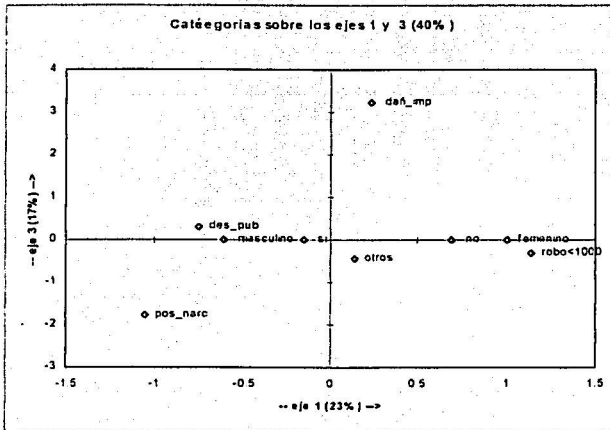
Si consideramos como el mejor subespacio el formado ahora por los pares de ejes, tendríamos los siguientes:

Gráfica	% de la inercia
eje 1 Vs eje 2	42*
eje 1 Vs eje 3	40*
eje 1 Vs eje 4	40*
eje 1 Vs eje 5	37
eje 1 Vs eje 6	34
eje 2 Vs eje 3	36
eje 2 Vs eje 4	36

eje 2 Vs eje 5	33
eje 2 Vs eje 6	30
eje 3 Vs eje 4	34
eje 3 Vs eje 5	31
eje 3 Vs eje 6	28
eje 4 Vs eje 5	31
eje 4 Vs eje 6	28
eje 5 Vs eje 6	25

De tal manera que si queremos conservar la mayor cantidad de inercia escogeríamos la interpretación de las gráficas: eje1 Vs eje 2, eje1 Vs eje 3, eje1 Vs eje 4.





Primera conclusión a partir de los gráficos.

1. Las mujeres están asociadas con los delitos de robo<1000. Este delito no es penalizado.

2. Los hombres están asociado con los delitos de posesión de narcótico y con daños impulsivos. Estos delitos son penalizados.

Si empleamos las restricciones debidas a la calidad de la representación de las categorías debemos tener en cuenta lo siguiente: la representación gráfica debe ser nítida, como una foto: para ello se requiere buena cantidad de inercia explicada, buena calidad de representación de la foto, buena correlación de los puntos con los ejes; se requieren que en el peor de los casos se pueda consultar simultáneamente todas las fotos, aunque ello propicie una actitud anti_parsimoniosa dentro del análisis de datos.

En base a lo anterior y consultando la salida del paquete Xlstat V4. para los datos de este problema tenemos lo siguiente Véase anexo 6.

a) Ninguna de las gráficas con dos ejes explica más del 42%. Por lo que la cantidad de información que se pierde es de más de la mitad.

b) La calidad de representación utilizando los dos primeros ejes (aquellos que nos proporcionan una mayor cantidad de inercia) es a lo más del 62% por lo que no podemos hablar de una buena representación con estos dos ejes. Si tomamos el primer y tercer ejes (40% de la inercia total) no hay una mejora, únicamente la categoría de daños impulsivos aumenta a un 68% de calidad en la representación. Utilizando el primer y cuarto ejes (40 % de la inercia total) tampoco hay un aumento en la calidad de la representación únicamente la categoría cargo "otros" aumenta a un 86%. En base a estos resultados podemos concluir que con una gráfica formada por pares de ejes la calidad de representación no es adecuada y por lo tanto no podemos decir todavía nada de manera segura.

Es importante notar en este momento que la representación de la asociación será más sencilla gráficamente cuando en realidad ésta se encuentre bien estructurada, cuando no es así entonces se tiene que hacer una intensa búsqueda de con cuales y cuántos ejes se puede representar la asociación y que las categorías de los factores quedan bien representadas, considerando al mismo tiempo que la inercia recuperada sea alta.

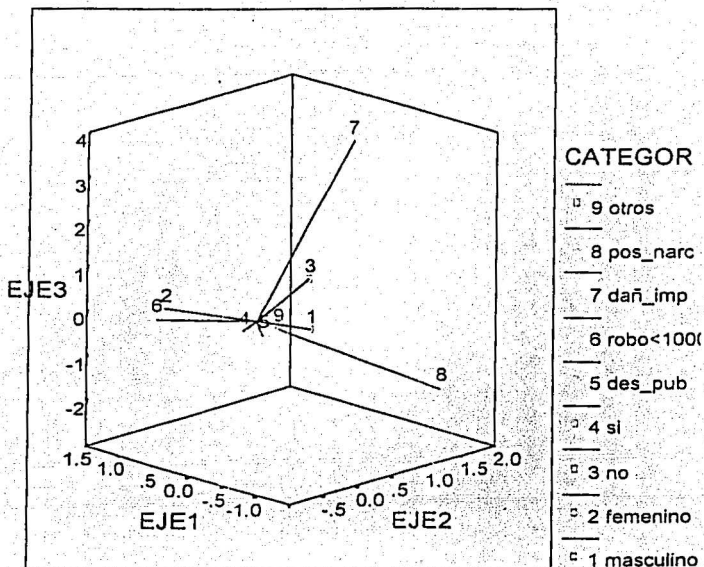
Una alternativa sería estudiar la representación con los tres primeros ejes que corresponde a 59% de inercia recuperada y en cuanto a la calidad por arriba de 59% para todas las categorías, con excepción de el cargo de otros, como puede verse en el apéndice relativo a la salida del paquete Xlstat V4.

Ahora bien, si aumentamos una dimensión a la gráfica tenemos las siguientes posibles gráficas.

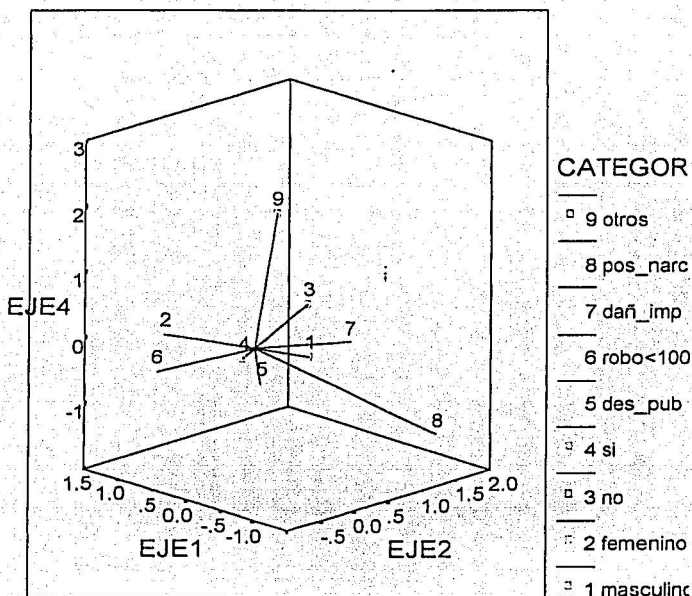
Gráfica	% de inercia
eje1 Vs eje2 Vs eje 3	59*
eje1 Vs eje2 Vs eje 4	59*
eje1 Vs eje2 Vs eje 5	56
eje1 Vs eje2 Vs eje 6	53
eje1 Vs eje3 Vs eje 4	57*
eje1 Vs eje3 Vs eje 5	54
eje1 Vs eje3 Vs eje 6	51
eje1 Vs eje4 Vs eje 5	54
eje1 Vs eje4 Vs eje 6	51
eje2 Vs eje3 Vs eje 4	53
eje2 Vs eje3 Vs eje 5	50
eje2 Vs eje3 Vs eje 6	47
eje2 Vs eje4 Vs eje 5	50
eje2 Vs eje4 Vs eje 6	47
eje2 Vs eje5 Vs eje 6	44
eje3 Vs eje4 Vs eje 5	48
eje3 Vs eje4 Vs eje 6	45
eje3 Vs eje5 Vs eje 6	42

De manera análoga las gráficas que más inercia aportan son: eje1 Vs eje2 Vs eje3, eje1 Vs eje2 Vs eje4 y eje1 Vs eje3 Vs eje 4. Las cuales se presentan a continuación:

TESIS CON
FALLA DE ORIGEN

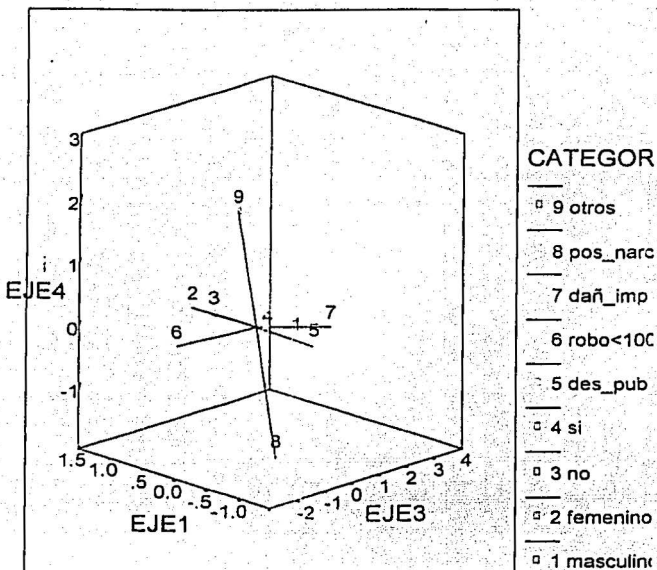


De esta gráfica que representa el 59% de la inercia total se puede notar que los hombres (1) están relacionados con daños impulsivos (7) y con posesión de narcóticos (8) y además esta asociación apunta en la dirección de no penalizados (3). A diferencia de las mujeres (2) que están asociadas con robo<1000 (6) y esta asociación está dada también con el hecho de que si son sentenciadas (4). Nótese que la calidad de representación es del 59% para todas las categorías a excepción del cargo otros.



De esta gráfica que también recoge el 59% de la inercia total se puede notar que los hombres (1) están relacionados con posesión de narcóticos (8) y no son sentenciados (3), nótese que aunque también parece que está asociado con daños impulsivos (7) la calidad de representación de esta categoría es del 20% lo cual es muy bajo. Por otro lado, las mujeres (2) están asociadas con robo<1000 (6) y además si son sentenciadas (4).

TESIS CON
 FALLA DE ORIGEN



De esta gráfica lo que más se puede decir es que los daños impulsivos (7) son opuestos a robo<1000 (6) y además la relación con el sexo sigue siendo la misma que en las gráficas anteriores, sólo que ahora posesión de narcóticos (8) no tiene una buena calidad de representación 58% y esta gráfica no es informativa con respecto a la asociación con el sexo ya que para hombres (1) y mujeres (2) la calidad de representación es apenas del 11%.

Del ejemplo anterior se aprecia la dificultad para observar la estructura de la asociación de los datos, desde luego que no es sencilla la interpretación de una gráfica de mayor dimensión y que además tiene una calidad de representación modesta o baja. Este problema muy probablemente se deba a los datos. En estos casos se sugiere hacer uso de otras técnicas que pueden complementar las interpretaciones que se han hecho, por ejemplo, se puede hacer uso de los modelos log-lineales.

Bajo las consideraciones anteriores se tiene que: el delito más frecuente en las mujeres es el (robo<1000) y además es sancionado, a diferencia del daño impulsivo asociado con el hombre y que no es sancionado.

Ahora bien, veamos un ejemplo que se ha manejado mucho en la literatura estadística y que se refiere a las flores de Iris, trabajada por Fisher. Se reporta (Dallas E. Johnson) que aplicando la técnica del análisis discriminante a la base de datos compuesta por 50 observaciones para el tipo de iris setosa, 50 para el tipo de iris versicolor, 50 para el tipo de iris virginica; con las siguientes variables:

1. Sepallen : longitud del sépalo.
2. Sepalwid: ancho del sépalo.
3. Petallen: longitud del pétalo.
4. Petalwid: ancho del pétalo.

Se obtiene tres grupos diferenciados para cada especie de flor de iris. En el análisis reportado el porcentaje de flores que fueron clasificadas correctamente fue es del 98%, y lo que se concluye es.

a) Las variables explicativas individualmente diferencian en promedio a las especies, por lo que son importantes para la clasificación de estas.

b) La variable que más aporta a la discriminación es el largo del pétalo.

c) Con las variables explicativas estudiadas la especie setosa se identifica claramente, pero las especies versicolor y virginica podrían confundirse.

Ante este conocimiento a priori, seleccionamos dos de las tres variedades de iris: la setosa y la virginica ya que son los grupos que más se diferencian entre si. Lo que nos interesa conocer es si existe alguna característica que determine a que variedad pertenece, es decir, si por ejemplo medidas pequeñas en las variables están relacionadas con alguna especie en particular.

Para llevar a cabo lo anterior, dado que la escala de medida es de razón, se propusieron las siguientes categorías, que se formaron de la siguiente manera: 1) se obtuvieron los valores mínimo y máximo de las variables estudiadas. 2) como se deseaba que fueran tres categorías: corto, mediano, largo, entonces se tomo $\Delta E = \frac{\max - \min}{3}$ como la longitud del intervalo, y así las categorías fueron:

corto: $(\min, \min + \Delta E)$

mediano: $(\min + \Delta E, \min + 2\Delta E)$

largo: $(\min + 2\Delta E, \max)$

Con estas categorías tenemos una medida cualitativa que nos relaciona las longitudes de las variables.

La nueva base de datos se presenta en el anexo 7, y a ésta se le aplicó un análisis de correspondencias múltiples, para las variables: variedad de flor, sepalen, sepalwid, petallen, petalwid. Por lo que en total se tienen catorce categorías y cinco factores.

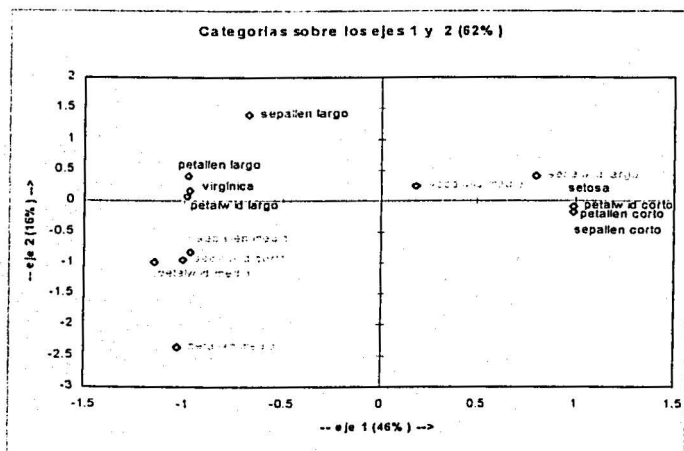
Los resultados se muestran a continuación:

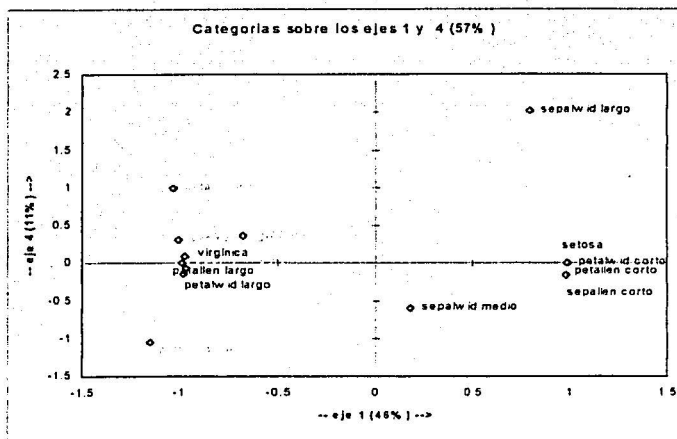
1. El número de ejes que se tienen es: $14-5=9$.
2. Dentro de esos 9 ejes debemos escoger aquellos que conserven la mayor cantidad de inercia

Valores pr.	1	2	3	4	5	6	7	8	9
Valor	0.6223	0.2163	0.2227	0.2026	0.1343	0.1053	0.0275	0.0000	0.0000
% de in.	48%	16%	12%	11%	7%	6%	2%	0%	0%
% acumulé	48%	64%	76%	87%	94%	100%	100%	100%	100%

de la tabla anterior se observa que el primer eje recoge el 48% de la inercia total, el segundo el 16% y el tercero y cuarto el 12% y 11% respectivamente. Lo anterior sugiere que si se está captando la estructura de asociación en los datos y además de una manera fuerte, con los tres primeros ejes se recoge casi el 75% de la inercia total.

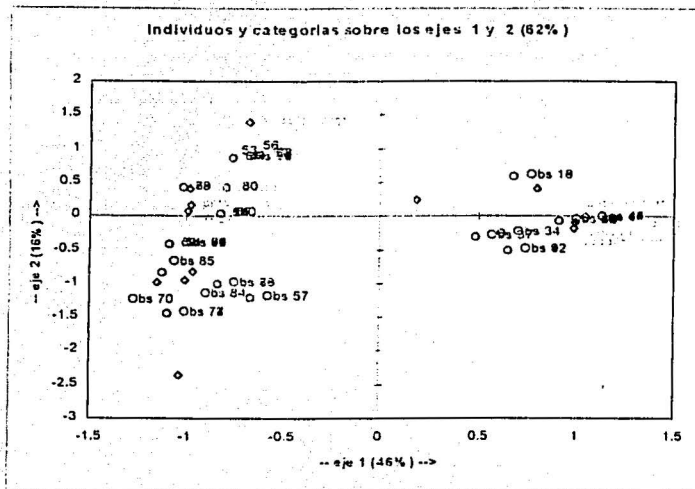
3. La gráfica de correspondencias con los dos primeros ejes es:





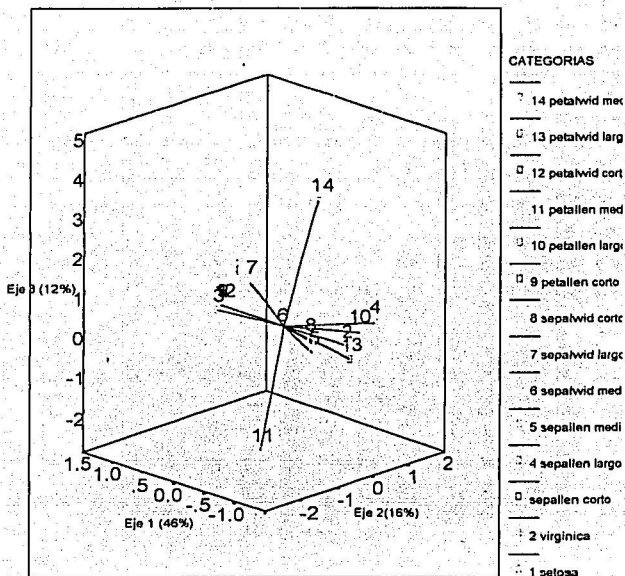
De las gráficas anteriores se concluye que las variedades de iris son opuestas totalmente, esto muestra como se sabía una clara diferenciación entre estas variedades, ahora bien, otro punto interesante era determinar que características se relacionaban con estas variedades. Así, de las gráficas se concluye que para el caso de la iris virginica la asociación está mayoritariamente dada por longitudes grandes tanto del sépalo como del pétalo y una anchura en el pétalo grande. Pero para el caso de la iris setosa se tiene que está asociada mayoritariamente a longitudes pequeñas tanto del pétalo, el sépalo y el ancho del pétalo. Hay que hacer notar que en las gráficas anteriores las categorías en rojo representan aquellas cuya calidad de representación era baja (abajo del 50%).

**TESIS CON
FALLA DE ORIGEN**



Si quisiéramos hacer un análisis sobre los individuos entonces vemos que estos se encuentran también en dos grupos bastante definidos y además como era de esperarse, las observaciones del 1 al 50 se encuentran relacionadas con la variedad de iris setosa, y los datos del 50 al 100 se encuentran relacionados con la variedad virginica, además las calidades de representación de los individuos es bastante buena considerando los dos primeros ejes. Véase datos anexo 7.

A pesar de que con estas gráficas se tiene una corroboración acerca del comportamiento que tienen los datos obtenido a través del análisis discriminante, si seguimos con la idea de recuperar la mayor cantidad de información posible entonces agreguemos una dimensión más.



Se observa que efectivamente la iris setosa (1) es opuesta a la iris virginica (2) ambas con una calidad de representación del 99%. Además se aprecia que la iris virginica está asociada a petallen largo, (10) petalwid largo(13) y sepallen largo (4) tal como se había concluido en las gráficas anteriores. Por otro lado la iris setosa está asociada a petalwid corto (12), petallen medio (11), sepallen corto (3). Lo cual concuerda con lo establecido en las conclusiones de las gráficas con dos ejes. Sin embargo, con esta tres dimensiones se recupera el 75% de la inercia total, lo cual es muy bueno y se garantiza una calidad de representación buena para las categorías.

Así, podemos concluir con base al análisis que medidas cortas en el largo y ancho del pétalo y del sépalo están asociadas fuertemente con la iris setosa, mientras que medidas grandes de estas variables están asociadas a las iris virginica.

Como conclusión para este capítulo tenemos

1) En el capítulo anterior se mostró que existía una relación entre la cantidad total de inercia y

el coeficiente de contingencia en media cuadrática ϕ^2 que nos permitía interpretar para cada gráfica el grado de asociación entre las categorías; en el AFC-multiple no existe esta relación, por lo que la interpretación de la asociación se vuelve más heurística al momento de interpretar la gráfica, por lo que esta interpretación debe estar basada en buena medida con el conocimiento que se tenga de la posible asociación.

2) A veces es posible reducir la dimensión de la tabla de contingencia para más de dos factores y convertirla en una tabla de contingencia de dos dimensiones por medio de la concatenación de algunas categorías o a través de procedimientos acorde a los objetivos del estudio y que buscan generar nuevas variables para trabajar con ellas.

3) A veces resulta difícil determinar la estructura de asociación entre las categorías de los factores estudiados a través del análisis de correspondencias múltiple, sin embargo ante estas problemáticas se recomienda recurrir a otras técnicas estadísticas que ayuden a interpretar los resultados.

4) Cuando la estructura de asociación en los datos es clara, entonces el análisis de correspondencias resulta muy útil, ya que presenta en una gráfica esa relación tomando en cuenta los elementos que permiten determinar la calidad de las gráficas.

Conclusiones.

La posible asociación entre las categorías de los factores puede ser estudiada a través de un AFC-S o un AFC-M según sea el caso, mediante representaciones gráficas que pueden ser presentadas como una línea recta, un plano, un plano en tres dimensiones o un hiperplano. Como una similitud podemos decir que estas representaciones son semejantes a un conjunto de "fotografías", en las que con un grado de nitidez aceptable en los elementos que la conforman se puede estudiar la posible asociación entre las categorías de los factores de la tabla de contingencia ya sea entre las categorías de los factores o bien por categorías entre factores; dicha asociación puede establecerse en la medida que:

i) Los perfiles de las categorías se alejen del origen (centroide) y muestren un proximidad entre sí, si esto ocurre se podrá decir que existe una asociación fuerte.

ii) La cantidad de inercia que sea recogida por los ejes que definen la(s) representación(es) gráfica(s) sea alta. A medida que se aumenta el número de dimensiones en la gráfica aumenta la cantidad de inercia representada.

iii) La calidad de representación de los perfiles dentro de la(s) representación(es) gráfica(s) sea buena, lo que da la confianza para poder hablar de la posible asociación.

Para el caso del AFC-S, se observó que cuando existen pocas categorías la cantidad de inercia que recogen las gráficas es alta, sin embargo en la práctica puede suceder que el número de categorías de los factores en la tabla de contingencia sean muchos por lo que también el número de ejes lo serán, decidir cuáles representaciones se conservarán para el análisis requerirá un conocimiento muy particular acerca de la asociación que se quiere estudiar ya que la interpretación va más allá de las consideraciones geométricas por que puede darse el caso de que las gráficas seleccionadas conserven una cantidad de inercia alta pero tengan una calidad de representación baja o al contrario la cantidad de representación sea alta pero la cantidad de inercia sea baja. Así, se sugiere que se conserven los ejes que cumplan con: cantidad de inercia alta y calidad de representación alta; considerándose las posibles combinaciones de gráficas de dos dimensiones o las posibles combinaciones de gráficas en tres dimensiones etc.

Para el caso del AFC-M con la introducción de la matriz indicadora Z, el análisis presenta el problema de que la cantidad de inercia explicada disminuye, por lo que hacer una interpretación semejante al AFC-S resulta complicado debido a que la inercia se distribuye en un mayor número de ejes lo que ocasiona que la interpretación esté basada en gráficas con mayor dimensión o bien la interpretación de gráficas de dos dimensiones que recogen poca inercia; esto evidencia la necesidad de contar con un conocimiento general del contexto del problema para el que se utiliza el AFC.

Por otra parte puede darse el caso en que los datos presenten una estructura interna de asociación marcada por ejemplo en los datos para el estudio de las "iris de fisher". En el cual sin bien en el AFC-M no se recogía una cantidad de inercia muy alta con los tres primeros ejes, la interpretación esta de acuerdo a lo reportado en la literatura estadística.

Ante estos problemas de interpretación del análisis para tablas de contingencia con muchas categorías o con más de dos factores se han propuesto para el primer caso reducir el número de categorías colapsando éstas de tal manera que se reduzca el número de ejes en la gráfica y para el segundo caso reducir la dimensión de una tabla de contingencia a una de dos dimensiones mediante la concatenación de las categorías de un factor sobre otro, de tal manera que se pueda aplicar un AFC-S. Esto desde luego lleva el riesgo en la metodología que se utiliza para abordar el problema bajo estudio, ya que colapsar categorías o concatenar factores puede ser incongruente con la investigación en que se utilice el AFC.

Ante esta situación se recomienda un estudio a cerca de la metodología del uso del AFC-M y la utilización de técnicas complementarias como pueden ser los modelos log-lineales.

Bibliografía

1. J.P. Benzecri. *Correspondence Analysis Handbook*. Marcel Dekker, Inc, New York, 1992.
2. Ludovic Lebart, Alain Marineau, kenneth M Warwick. *Multivariate Descriptive Analysis (Correspondence Analysis and Related Techniques For Large Matrices)*, John Wiley & Sons , USA, 1984.
3. Michael J. Greenacre. *Theory and Applications of Correspondence Analysis*, Academic Press, London, 1984.
4. J.M. Comejo. *Técnicas de Investigación social. El Análisis de Correspondencias (Teoría y Práctica)*. Biblioteca Universitaria de Ciencias Sociales, Serie Medium, Barcelona España, 1988.
5. Brian.Sidney Everit. *The Analysis of Contingence Tables*. London, Chapman and Hall, 1977.
6. Francisco Casanova del Ángel. *These de doctorat de 3e Cycle. Mathematiques; Traitement des donnees statistiques*. Paris, 1980.
7. Rubén Hernández. *Apuntes: Análisis de datos*.
8. K.V. Mardia. *Multivariate Analysis* . Academic Press, USA, 1979.
9. Alvin C. Rencher. *Metods of Multivariate Analysis*. Wiley Series in Probability and Mathematical Statistics, USA, 1995.
10. Bryan F.J. Manly. *Multivariate Statistical Methods. A Primer*. Chapman and Hall, Great Britain, 1990.
11. Eduardo Solar González, *Apuntes de Álgebra Lineal*, UNAM, 1985
12. Brian,S. Everitt and Graham Dunn. *Applied Multivariate Data Analysis*. Ed. Edward Arnold A. division of Hodder &Stoughton, Londres 1991.
13. Journal of Econometrics 22 (1983) 139-167. North-Holland Publishing Company *Analyzing rectangular tables by joint and constrained multidimensional scaling*. Willenm.J.Heiser and Jaqueline Meulman.
14. Alan,Agresti. *Analysis of Ordinal Categorical Data*. Wiley Series in Probability and Mathematical Statistics, USA, 1984
15. *Interpreting Multivariate Data*. Edited by Vic Barnett.Wiley Series in probability and mathematical statistics. 1981.

16. *A users guide of principal Components Analysis*. Wiley Series in Probability and Mathematical Statistics, USA
18. W.J Conover. *Practical Nonparametric Statistics*. John Wiley & Sons, Inc. 3rd. Edition, USA, 1999.
19. Visauta Vinacua Bienvenido. *Análisis Estadístico con SPSS para Windows. Estadística Multivariante*. Mc Graw Hill, España, 1999.
20. J.D. Jobson. *Applied Multivariate Data Analysis*. Vol II. ,Springler Verlag.
21. Luis Joaristi Olariaga, Luis Lizasoain Hernández, Análisis de correspondencias. Cuadernos de estadística, Madrid, España, Ed. La muralla, 1999.
- 22 Dallas E. Johnson, *Métodos multivariados aplicados al análisis de datos*, México, Ed. International Thomson, 2000.
23. Lara Pérez Soto, Claudia. *Análisis de Correspondencias Múltiples como una técnica para el estudio de datos cualitativos*. Tesis de licenciatura, Facultad de Ciencias, UNAM, 1990.
24. Humberto Soto de la Rosa, *Algunas Técnicas para el análisis de datos categóricos*, Tesis de licenciatura, Instituto Tecnológico Autónomo de México. 1999.
25. Nuria Guerra Vargas, *Análisis de correspondencias simples y múltiples (tres enfoques teóricos, aplicaciones y programación)*, Tesis de licenciatura, Instituto Tecnológico Autónomo de México. 1996.

Anexos

Anexo 1. Resultados de álgebra lineal y álgebra de matrices

1. Espacio Vectorial

Sea V un conjunto no vacío y sea $(K, +, \cdot)$ un campo. Se dice que V es un espacio vectorial sobre K si están definidas dos leyes de composición llamadas adición y multiplicación por un escalar, tal que:

a) La adición asigna a cada pareja ordenada (\vec{u}, \vec{v}) de elementos de V un único elemento $(\vec{u} + \vec{v}) \in V$, llamado la suma de V .

$$b) \forall \vec{u}, \vec{v}, \vec{w} \in V \text{ se tiene que } \vec{u} + (\vec{v} + \vec{w}) = (\vec{u} + \vec{v}) + \vec{w}$$

$$c) \exists \vec{0} \in V \text{ tal que } \vec{0} + \vec{v} = \vec{v}, \forall \vec{v} \in V$$

$$d) \forall \vec{v} \in V \exists -\vec{v} \in V \text{ tal que } -\vec{v} + \vec{v} = \vec{0}$$

$$e) \forall \vec{u}, \vec{v} \in V \vec{u} + \vec{v} = \vec{v} + \vec{u}$$

f) La multiplicación por un escalar asigna a cada pareja ordenada (α, \vec{v}) de elementos $\alpha \in K$ y $\vec{v} \in V$ un único elemento $\alpha\vec{v} \in V$, llamado producto de α por \vec{v} .

$$g) \forall \alpha \in K; \vec{u}, \vec{v} \in V \text{ se tiene que } \alpha(\vec{v} + \vec{u}) = \alpha\vec{u} + \alpha\vec{v}$$

$$h) \forall \alpha, \beta \in K; \vec{v} \in V \text{ se tiene que } (\alpha + \beta)\vec{v} = \alpha\vec{v} + \beta\vec{v}$$

$$i) \forall \alpha, \beta \in K; \vec{v} \in V \text{ se tiene que } \alpha(\beta\vec{v}) = (\alpha\beta)\vec{v}$$

$$j) \text{ Si } 1 \text{ es la unidad en } K \text{ se tiene que } \forall \vec{v} \in V, 1\vec{v} = \vec{v}$$

A los elementos de V se les llama vectores y a los elementos de K se les llama escalares.

2. Combinación lineal

Un vector \vec{w} es una combinación lineal de los vectores $\vec{v}_1, \vec{v}_2, \dots, \vec{v}_n$ si puede ser expresado en la forma:

$$\vec{w} = \alpha_1 \vec{v}_1 + \alpha_2 \vec{v}_2 + \dots + \alpha_n \vec{v}_n$$

donde $\alpha_1, \alpha_2, \dots, \alpha_n$ son escalares.

3 Vectores linealmente independientes y vectores linealmente dependientes

Sea $S = \{\vec{v}_1, \vec{v}_2, \dots, \vec{v}_n\}$ un conjunto de vectores:

i) S es linealmente dependiente si existen escalares $\alpha_1, \alpha_2, \dots, \alpha_n$, no todos iguales a cero, tales

que

$$\vec{0} = \alpha_1 \vec{v}_1 + \alpha_2 \vec{v}_2 + \dots + \alpha_n \vec{v}_n$$

ii) S es linealmente independiente si

$$\vec{0} = \alpha_1 \vec{v}_1 + \alpha_2 \vec{v}_2 + \dots + \alpha_n \vec{v}_n$$

$$\text{con } \alpha_1 = \alpha_2 = \dots = \alpha_n = 0.$$

4. Conjunto generador.

Sea V un espacio vectorial sobre K , y sea $G = \{\vec{v}_1, \vec{v}_2, \dots, \vec{v}_m\}$ un conjunto de vectores de V .

Se dice que G es un generador de V si para todo vector $\vec{x} \in V$ existen escalares $\alpha_1, \alpha_2, \dots, \alpha_m$ tales que $\vec{x} = \alpha_1 \vec{v}_1 + \alpha_2 \vec{v}_2 + \dots + \alpha_m \vec{v}_m$. es decir todos los vectores de V son combinaciones lineales de los elementos

de G .

5. Base de un espacio vectorial V

Se llama base de un espacio vectorial V a un conjunto generador de V que es linealmente independiente.

5. Dimensión de un espacio vectorial

Sea V un espacio vectorial sobre K . Si $B = \{\vec{v}_1, \vec{v}_2, \dots, \vec{v}_n\}$ es una base de V se dice que V es de dimensión n , lo cual se denota con

$$\dim V = n$$

En particular, si $V = \{\vec{0}\}$,

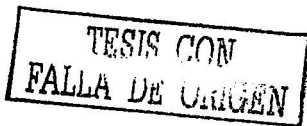
$$\dim V = 0$$

6. Transformaciones :

Entenderemos por transformaciones lineales a las funciones vectoriales de variable vectorial, es decir funciones del tipo $\vec{w} = f(\vec{v})$. En terminos formales tenemos: Si V y W son espacios vectoriales una función $T : V \rightarrow W$ recibe el nombre de transformación. Los espacios V y W se llaman, respectivamente dominio y codominio de la transformación.

Operaciones elementales de matrices.

7. Para multiplicar dos matrices, el número de columnas en la primera debe ser igual al número de renglones en la segunda. De ese modo, la multiplicación de matrices está dada de la siguiente manera:



$$A_{m \times n} B_{n \times p} = C_{m \times p}$$

8. Definición de Producto:

Sean $A = [a_{ij}]$ y $B = [b_{ij}]$ dos matrices con elementos en \mathbb{C} (números complejos), de $m \times n$ y $n \times q$ respectivamente. El producto AB es una matriz $P = [p_{ij}]$, de $m \times q$, definida por:

$$p_{ij} = \sum_{k=1}^n a_{ik} b_{kj}; \text{ para } i = 1, 2, 3, \dots, m \text{ y } j = 1, 2, 3, \dots, q.$$

9. Cuando todos los productos de matrices están definidos, se cumple la ley asociativa para la multiplicación de las mismas. Por Tanto $A(BC) = (AB)C$.

10. La multiplicación de matrices no siempre es conmutativa. Es decir AB no siempre es igual a BA .

11. Definición (Suma)

Sean $A = [a_{ij}]$ y $B = [b_{ij}]$ dos matrices de $m \times n$ con elementos en \mathbb{C} (números complejos). La suma de $A + B$ es una matriz $S = [s_{ij}]$, de $m \times n$, definida por

$$s_{ij} = a_{ij} + b_{ij}, \text{ para } i = 1, 2, 3, \dots, m \text{ y } j = 1, 2, 3, \dots, n.$$

12. Teorema.

Si A, B son matrices de $m \times n$ cuyos elementos son números complejos, y sea $\alpha \in \mathbb{R}$ entonces,

- i). Existe una matriz $-A$ de $m \times n$ tal que $A + (-A) = 0$
- ii). $\alpha(A + B) = \alpha A + \alpha B$

Álgebra de matrices.

13. Si

$$A = \begin{bmatrix} a_{11} & a_{12} & \dots & a_{1p} \\ a_{21} & a_{22} & \dots & a_{2p} \\ \dots & \dots & \dots & \dots \\ a_{n1} & a_{n2} & \dots & a_{np} \end{bmatrix}$$

entonces A^t se conoce como la transpuesta de la matriz A . Ésta se define por

$$A^t = \begin{bmatrix} a_{11} & a_{12} & \dots & a_{1n} \\ a_{21} & a_{22} & \dots & a_{2n} \\ \dots & \dots & \dots & \dots \\ a_{1p} & a_{2p} & \dots & a_{np} \end{bmatrix}$$

14. $(AB)' = B'A'$

15. Se dice que una matriz A es simétrica si $A = A'$

16. La traza de una matriz cuadrada $A_{p \times p}$, se define como la suma de sus elementos en la diagonal. Se escribe:

$$\text{tr}(A) = \sum_{i=1}^p a_{ii}$$

17. Se dice que una matriz cuadrada es una matriz diagonal si todos sus elementos que no están en la diagonal son ceros.

18. El rango de una matriz se define como el número máximo de renglones (columnas) en A que son linealmente independientes. De modo equivalente, el rango de A es la dimensión del subespacio vectorial generado por los renglones (columnas) de la matriz A .

19. Proposiciones:

a) El Rango de $A_{m \times n} \leq \min(m, n)$

b) Rango de $A = \text{rango de } A'$

20. Dos vectores X y Y , son ortogonales si $XY' = 0$. Además se dice que dos vectores ortogonales son ortonormales si $X'X = 1$ y $Y'Y = 1$.

21. Se dice que un vector X está normalizado si $X'X = 1$

Formas cuadráticas

22. Una forma cuadrática en p variables x_1, x_2, \dots, x_p , es una función de la forma $\sum_{i=1}^p \sum_{j=1}^p a_{ij}x_i x_j$. Una forma cuadrática de este tipo siempre se puede escribir como:

$$X'AX, \text{ para alguna matriz simétrica } A_{p \times p}, \text{ en donde } X = \begin{bmatrix} x_1 \\ x_2 \\ \dots \\ x_p \end{bmatrix}$$

23. se dice que una matriz simétrica A es

a) Positiva definida si $X'AX > 0$, para toda $X \neq 0$.

b) Positiva semidefinida si $X'AX \geq 0$, para toda X y si $X'AX = 0$ para algún X diferente de cero.

Eigenvalores y Eigenvectores.

24. Los eigenvalores de una matriz simétrica $A_{p \times p}$, son las raíces de una ecuación polinomial dada por

$$|A - \lambda I| = 0,$$

la cual es una ecuación polinomial en λ , de p -ésimo grado

25. A cada eigenvalor de A le corresponde un vector diferente de cero llamado eigenvalor, que satisface

$$Ac_i = \lambda_i c_i, \text{ para } i = 1, 2, 3, \dots, p$$

26. Si A es una matriz simétrica de números reales, entonces sus eigenvalores y eigenvectores también consistirán en números reales.

27. Los eigenvectores no son únicos, de modo que a menudo se normalizan de tal manera que $c_i^T c_i = 1$.

28. La traza de una matriz simétrica es igual a la suma de sus eigenvalores, es decir, $tr(A) = \sum_{i=1}^p \lambda_i$.

Ángulos

29. El ángulo θ entre X y Y se expresa por:

$$\cos \theta = \frac{X^T Y}{\|X\| \|Y\|}$$

Anexo 2. Espacio métrico: determinación del cálculo de la distancia Ji cuadrada

Una distancia o una métrica en un conjunto S es una función $d: S \times S \rightarrow R$ que satisface las siguientes propiedades:

- $d(x,y) \geq 0$ para toda $x,y \in S$
- $d(x,y) = 0 \Leftrightarrow x=y$
- $d(x,y) = d(y,x)$ para toda $x,y \in S$
- $d(x,y) \leq d(x,z) + d(z,y)$ para toda $x,y,z \in S$

Un espacio métrico (S,d) es un conjunto S con una métrica d en S.

Para ilustrar un poco como se ha llegado a la fórmula de distancia para dos renglones tomemos la matriz perfil por renglones $R_{n \times p} = D_n^{-1} F$, cuya forma general es:

$$\begin{bmatrix} \frac{f_{11}}{f_{1\cdot}} & \frac{f_{12}}{f_{1\cdot}} & \frac{f_{13}}{f_{1\cdot}} & \dots & \frac{f_{1j}}{f_{1\cdot}} \\ \frac{f_{21}}{f_{2\cdot}} & \frac{f_{22}}{f_{2\cdot}} & \frac{f_{23}}{f_{2\cdot}} & \dots & \frac{f_{2j}}{f_{2\cdot}} \\ \dots & \dots & \dots & \dots & \dots \\ \frac{f_{n1}}{f_{n\cdot}} & \frac{f_{n2}}{f_{n\cdot}} & \frac{f_{nj}}{f_{n\cdot}} & \dots & \frac{f_{nj}}{f_{n\cdot}} \end{bmatrix}$$

Ahora calculemos la distancia entre dos renglones de acuerdo a la métrica definida por Dc^{-1} (matriz diagonal inversa de los pesos de las columnas). Los renglones se encuentran en un espacio R^j ponderado por los pesos asociados a la matriz Dc^{-1} , para cada elemento del renglón, es decir

$$F_j = \left(\frac{f_{11}}{f_{1\cdot}}, \frac{f_{12}}{f_{1\cdot}}, \frac{f_{13}}{f_{1\cdot}}, \dots, \frac{f_{1j}}{f_{1\cdot}} \right)$$

$$F'_j = \left(\frac{f_{i1}}{f_{i\cdot}}, \frac{f_{i2}}{f_{i\cdot}}, \frac{f_{i3}}{f_{i\cdot}}, \dots, \frac{f_{ij}}{f_{i\cdot}} \right)$$

$$\begin{aligned} \text{lo cual implica que } (F_j - F'_j) &= \left(\frac{f_{11}}{f_{1\cdot}} - \frac{f_{i1}}{f_{i\cdot}}, \frac{f_{12}}{f_{1\cdot}} - \frac{f_{i2}}{f_{i\cdot}}, \frac{f_{13}}{f_{1\cdot}} - \frac{f_{i3}}{f_{i\cdot}}, \dots, \frac{f_{1j}}{f_{1\cdot}} - \frac{f_{ij}}{f_{i\cdot}} \right) \\ &= \left(\frac{f_{11}}{f_{1\cdot}} - \frac{f_{i1}}{f_{i\cdot}}, \frac{f_{12}}{f_{1\cdot}} - \frac{f_{i2}}{f_{i\cdot}}, \frac{f_{13}}{f_{1\cdot}} - \frac{f_{i3}}{f_{i\cdot}}, \dots, \frac{f_{1j}}{f_{1\cdot}} - \frac{f_{ij}}{f_{i\cdot}} \right) \end{aligned}$$

De acuerdo a la definición de producto punto en este espacio ponderado tenemos.

$$d^2(i, i') = \langle i - i', i - i' \rangle$$

$$= \left(\frac{f_{1i}}{f_{.i}} - \frac{f_{1i'}}{f_{.i'}}, \frac{f_{2i}}{f_{.i}} - \frac{f_{2i'}}{f_{.i'}}, \dots, \frac{f_{ji}}{f_{.i}} - \frac{f_{ji'}}{f_{.i'}} \right) Dc^{-1} \left(\frac{f_{1i}}{f_{.i}} - \frac{f_{1i'}}{f_{.i'}}, \frac{f_{2i}}{f_{.i}} - \frac{f_{2i'}}{f_{.i'}}, \dots, \frac{f_{ji}}{f_{.i}} - \frac{f_{ji'}}{f_{.i'}} \right)'$$

donde

$$Dc^{-1} = \begin{bmatrix} \frac{1}{f_{.1}} & 0 & 0 & 0 \\ 0 & \frac{1}{f_{.2}} & 0 & 0 \\ 0 & 0 & \dots & 0 \\ 0 & 0 & 0 & \frac{1}{f_{.j}} \end{bmatrix}$$

y por lo tanto

$$d^2(i, i') = \frac{1}{f_{.1}} \left[\frac{f_{1i}}{f_{.i}} - \frac{f_{1i'}}{f_{.i'}} \right]^2 + \frac{1}{f_{.2}} \left[\frac{f_{2i}}{f_{.i}} - \frac{f_{2i'}}{f_{.i'}} \right]^2 + \dots + \frac{1}{f_{.j}} \left[\frac{f_{ji}}{f_{.i}} - \frac{f_{ji'}}{f_{.i'}} \right]^2$$

$$= \sum_{j=1}^J \frac{1}{f_{.j}} \left[\frac{f_{ji}}{f_{.i}} - \frac{f_{ji'}}{f_{.i'}} \right]^2$$

Que es la misma expresión que se planteó en el capítulo I. Ponderar las diferencias con respecto a los elementos de los espacios respectivos implica introducir las matrices diagonales de pesos $Dc_{j \cdot}^{-1}$ para el caso de $R^j : Dr_{j \cdot}^{-1}$ y algunos autores llaman a estas matrices "métricas" o la distancia. De tal manera que se habla de la distancia definida por la matriz Dc^{-1} , para el espacio R^j de los renglones.

Habría que probar que la distancia así definida en el espacio ponderado es efectivamente una distancia. Por lo pronto dejaremos eso de lado y supondremos que la distancia está bien definida.

Una de las razones para escoger la distancia Ji-cuadrada, es que en esta se verifica la propiedad de equivalencia distribucional, que se expresa para el caso de perfiles por renglón y columna como:

"Si dos renglones de la tabla de contingencia original, son proporcionales y, se reemplazan por un sólo renglón el cual es la suma de los dos renglones, entonces las distancias entre las columnas no cambian. Esta propiedad es importante porque garantiza invarianza en resultados sin importar como fueron recogidos.

Anexo 3

Maximización de una forma cuadrática bajo un restricción cuadrática.

Encontremos un vector u tal que maximice la cantidad $u' Au$ con la restricción $u' Mu = 1$, donde A y M son matrices simétricas, y adicionalmente M es positiva definida.

Es decir:

$$\text{Max } u' Au$$

$$\text{s.a. } u' Mu = 1$$

Desarrollo:

La forma cuadrática $u' Au$ puede ser escrita como:

$$u' Au = \sum_i \sum_j a_{ij} u_i u_j \quad (1)$$

Derivando esta cantidad para para las p componentes del vector u sucesivamente, se observa que el vector de las derivadas parciales de $u' Au$ se escriben en forma matricial como:

$$\frac{\partial(u' Au)}{\partial u} = 2Au \quad (2)$$

y

$$\frac{\partial(u' Mu)}{\partial u} = 2Mu \quad (3)$$

El lagrangiano es:

$$\mathcal{L} = u' Au - \lambda(u' Mu - 1) \quad (4)$$

donde λ es un multiplicador de Lagrange

Para encontrar un máximo se requiere igualar las derivadas de la ecuación de Lagrange a cero:

$$\frac{\partial \mathcal{L}}{\partial u} = 2Au - 2\lambda Mu = 0 \quad (5)$$

De la ecuación anterior se deduce que:

$$Au = \lambda Mu \quad (6)$$

Cuando premultiplicamos a ambos lados de la ecuación por u' , tomando en cuenta que el hecho de que $u' Mu = 1$ se tiene que

$$u' Au = \lambda u' Mu = \lambda, \quad (7)$$

el valor del parámetro λ es máximo.

Como la matriz M es positiva definida, y por lo tanto no singular, la relación (6) se escribe como:

$$M^{-1} Au = \lambda u \quad (8)$$

u es el eigenvector de la matriz $M^{-1}A$ que corresponde al eigenvalor más grande λ .

De aquí en adelante llamaremos u_1 al vector u que corresponde al valor más grande λ_1 tal que la ecuación (6) es cierta. Encontraremos un vector u_2 tal que es M -ortogonal a u_1 , es decir $u_1' Mu_2 = 0$; y tenga norma 1, es decir $u_2' Mu_2 = 1$.

Igualando a cero las derivadas parciales de la ecuación de Lagrange

$$L = u_1' Au_2 - \lambda_2 (u_2' Mu_2 - 1) - \mu_2 u_1' Mu_2 \quad (9)$$

donde λ_2 y μ_2 son los dos multiplicadores de Lagrange.

La condición para un máximo o un mínimo es escrita para u_2 como:

$$\frac{\partial L}{\partial u_2} = 2Au_2 - 2\lambda_2 Mu_2 - \mu_2 Mu_1 = 0 \quad (10)$$

Si multiplicamos los elementos de la ecuación (10) por u_1' , se tiene que $\mu_2 = 0$ porque

$$u_1' Au_2 = u_1' Mu_2 = 0 \Rightarrow \mu_2 u_1' Mu_1 = 0 \Rightarrow \mu_2 = 0$$

Así, de acuerdo a la restricción de ortogonalidad se tiene que

$$2Au_2 - 2\lambda_2 Mu_2 = 0 \quad (11)$$

es decir:

$$Au_2 = \lambda_2 Mu_2 \quad (12)$$

Como M es no singular, u_2 es el segundo eigenvector de $M^{-1}A$, relativo al segundo eigenvalor más grande λ_2 , si este es único.

En forma general podemos escribir que para encontrar el α -ésimo vector se resuelve el sistema:

$$Au_\alpha = \lambda_\alpha Mu_\alpha$$

y si M es no singular, $M^{-1}Au_\alpha = \lambda_\alpha u_\alpha$, α no puede ser más grande que el orden de la matriz A .

Anexo 4. Análisis de correspondencias como un caso particular del análisis de correlación canónico.

Los datos con los que trabaja el AFC, cuando hace uso de la información por individuos es $X=(x^1 \ x^2 \ \dots \ x^I)$ y $Y=(y^1 \ y^2 \ \dots \ y^J)$ donde x^i y y^j representan las variables indicadoras asociadas respectivamente a las categorías de los factores X y Y .

Esta forma de presentar el problema, permite estudiar la relación entre las características cualitativas a través del análisis de la dependencia entre dos grupos de variables cuantitativas muy particulares (las variables indicadoras).

De acuerdo a lo planteado el análisis canónico de las matrices X y Y consiste en encontrar parejas de variables canónicas (ξ^i, η^j) lo más correlacionadamente posible. De tal manera que las variables canónicas ξ^i, η^j transforman una característica cualitativa en una característica cuantitativa. A esta operación de cuantificar se conoce en estadística como "escalamiento de una variable cualitativa".

En esas condiciones, el AFC de la tabla original permite obtener los mejores escalamientos (en el sentido de máxima correlación) de las variables cualitativas X y Y .

Se identifican las matrices V_{XX} y V_{YY} por:

$$V_{XX} = \frac{1}{n} \begin{pmatrix} n_{1.} & 0 & 0 & 0 & 0 \\ 0 & \dots & 0 & 0 & 0 \\ 0 & 0 & n_{2.} & 0 & 0 \\ 0 & 0 & 0 & \dots & 0 \\ 0 & 0 & 0 & 0 & n_{I.} \end{pmatrix} = Dp_X$$

$$V_{YY} = \frac{1}{n} \begin{pmatrix} n_{.1} & 0 & 0 & 0 & 0 \\ 0 & \dots & 0 & 0 & 0 \\ 0 & 0 & n_{.j} & 0 & 0 \\ 0 & 0 & 0 & \dots & 0 \\ 0 & 0 & 0 & 0 & n_{.J} \end{pmatrix} = Dp_Y$$

Estas matrices tienen en sus diagonales los pesos relativos de cada renglón y de cada columna respectivamente.

Y , además la matriz V_{XY} está dada por:

$$V_{XY} = \frac{1}{n} X^T Y = \frac{1}{n} \begin{pmatrix} n_{11} & \dots & n_{1j} & \dots & n_{1J} \\ \dots & \dots & \dots & \dots & \dots \\ n_{i1} & \dots & n_{ij} & \dots & n_{iJ} \\ \dots & \dots & \dots & \dots & \dots \\ n_{J1} & \dots & n_{Jj} & \dots & n_{JJ} \end{pmatrix} = \frac{1}{n} C = C^*$$

de acuerdo al problema planteado para el análisis canónico clásico, los factores canónicos son solución de las ecuaciones:

$$D_{p_x}^{-1} C^* D_{p_x}^{-1} C^{*'} a_k = \lambda_k a_k$$

$$D_{p_y}^{-1} C^{*'} D_{p_y}^{-1} C^* b_k = \lambda_k b_k$$

Nótese que la matriz $D_{p_x}^{-1} C^*$ contiene los perfiles por renglón para cada categoría i de X . Así mismo, la matriz $D_{p_y}^{-1} C^{*'}$ contiene los perfiles por columna para cada categoría j de Y .

Los factores del AFC son entonces los vectores propios del producto de dos matrices de perfiles:

$$P_r P_c' a_k = \lambda_k a_k$$

$$P_c' P_r b_k = \lambda_k b_k$$

Anexo 5. Salida del programa Xlstat V 4.0 para los datos de la tabla 1. capítulo II..

Interpretación de la salida del paquete:

XLSTAT - Analyse Factorielle des Correspondances

Plage de données :

Test de l'hypothèse d'indépendance entre les lignes et les colonnes :

Distance du Chi2 sur les données = 19.1780 ~ Nombre de degrés de liberté = 6

Probabilité associée : 0.0039

Chi2 limite pour l'intervalle de confiance choisi = 12.5916

Sur la base de ce test on doit rejeter l'hypothèse d'indépendance entre lignes et colonnes.

} Prueba Ji cuadrada de indepe

Valeurs propres et pourcentage d'inertie correspondant :

Valeurs pr.	1	2
Valeur	0.0463	0.0157
% d'inertie	75%	25%
% cumulé	75%	100%

} Valores propios, que definen a los vectores propios que generan la base para el nuevo sistema de referencia: los ejes factoriales, solución al problema de maximización.

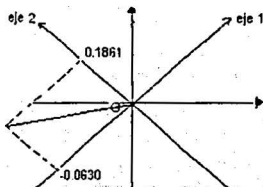
$$\text{Traza} = 0.0463 + 0.0157 = 0.062$$

Composantes des profils lignes sur les axes factoriels :

	Poids	Inertie	Inertie normée	1	2
1	0.3042	0.0117	0.1892	-0.0630	0.1861
2	0.3107	0.0305	0.4912	0.3094	-0.0490
3	0.3851	0.0198	0.3196	-0.1999	-0.1075

Perfil renglón Peso 0.062 Inertie/baza proyección eje 1 proyección eje 2

} Componentes de los perfil por renglón: turno



$$\begin{aligned} \text{Inercia (Turno 1)} &= 0.3042 \cdot (-0.0630)^2 + 0.3042 \cdot (0.1861)^2 \\ &= 0.0117427 \\ &\quad \uparrow \\ &\text{primer turno} \end{aligned}$$

**TESIS CON
FALLA DE ORIGEN**

Composantes des profils colonnes sur les axes factoriels :

	Poids	Inertie	Inertie normée	1	2
A	0.2395	0.0117	0.1883	0.0317	-0.2186
B	0.2233	0.0250	0.4023	0.3280	0.0650
C	0.4142	0.0058	0.0929	-0.0755	0.0898
D	0.1230	0.0196	0.3165	-0.3996	0.0052

Componentes de los perfiles por columna: Tipo de defecto

▼ Perfil columna ▼ Peso ▼ n h62 ▼ Inertie/traza ▼ proyección eje 1 ▼ proyección eje 2

Contributions des lignes aux inerties associées aux axes factoriels

	1	2
1	0.0280	0.6698
2	0.6419	0.0474
3	0.3321	0.2829

Con tribuciones relativas de los perfiles renglón a la inercia que recoge cada eje

$$= \frac{\text{Contribución absoluta}}{\lambda_k} = \frac{0.3042 * (0.083)^2}{0.0463} = 0.028$$

primer turno

Contributions des colonnes aux inerties associées aux axes factoriels

	1	2
A	0.0052	0.7274
B	0.5186	0.0599
C	0.0523	0.2124
D	0.4240	0.0002

Contribuciones relativas de los perfiles columna a la inercia que recoge cada eje

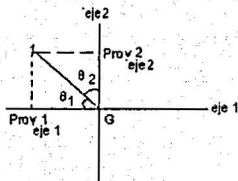
$$= \frac{\text{Contribución absoluta}}{\lambda_k} = \frac{0.2395 * (0.2186)^2}{0.0157} = 0.72$$

Defecto A

Cosinus carrés des angles des vecteurs lignes avec les axes

	1	2
1	0.1027	0.8973
2	0.9755	0.0245
3	0.7756	0.2244

Cosenos al cuadrado: correlaciones al cuadrado de los perfiles renglón con respecto a cada eje



$$d^2(G, \text{Turno 1}) = \text{Proy}_{\text{eje 1}}^2 + \text{Proy}_{\text{eje 2}}^2$$

$$\cos^2 \theta_1 = \frac{d_{\text{eje 1}}^2(G, \text{Turno 1})}{d^2(G, \text{Turno 1})} = \frac{-0.0630^2}{-0.063^2 + 0.1861^2} = 0.102817$$

▲
correlación con el
primer eje primer turno

$$\cos^2 \theta_2 = \frac{d_{\text{eje 2}}^2(G, \text{Turno 1})}{d^2(G, \text{Turno 1})} = \frac{0.1861^2}{-0.063^2 + 0.1861^2} = 0.897182$$

▲
correlación con el
segundo eje primer turno

Cosinus carrés des angles
des vecteurs-colonnes avec les axes

	1	2
A	0.0205	0.9795
B	0.9622	0.0378
C	0.4202	0.5798
D	0.9998	0.0002

} Cosenos al cuadrado: correlaciones al cuadrado de los perfiles por columna con respecto a cada eje

TESIS CON
FALLA DE ORIC

Anexo 6. Base de datos y resultados del AFCM. Para el ejemplo de tres factores : Sexo, sentencia, y cargo. Capítulo III

Sexo	sentenci	Cargo	frecuencia
HOMBRE	SI	DAN_IMP	8
HOMBRE	SI	ROBO<10000	11
HOMBRE	SI	DES_PUB	5
HOMBRE	SI	NARCOT	7
HOMBRE	SI	OTROS	12
MUJER	SI	DAN_IMP	5
MUJER	SI	ROBO<10000	15
MUJER	SI	DES_PUB	3
MUJER	SI	NARCOT	1
MUJER	SI	OTROS	6
HOMBRE	NO	DAN_IMP	105
HOMBRE	NO	ROBO<10000	32
HOMBRE	NO	DES_PUB	11
HOMBRE	NO	NARCOT	23
HOMBRE	NO	OTROS	37
MUJER	NO	DAN_IMP	32
MUJER	NO	ROBO<10000	57
MUJER	NO	DES_PUB	6
MUJER	NO	NARCOT	2
MUJER	NO	OTROS	25

XLSTAT - Analyse Factorielle des Correspondances Multiples / Début le 03/02/2003 à 09:43:02 a.m.
 Plage de données : classeur = Libro2 / feuille = Hoja1 / plage = Hoja1!\$A\$2:\$C\$21
 Variables supplémentaires : 0
 Observations supplémentaires : 0

Tableau de Burt :

	HOMBRE	MUJER	SI	NO	DAN_IMP	ROBO<10000	DES_PUB	NARCOT	OTROS
HOMBRE	251	0	43	208	113	43	13	30	49
MUJER	0	152	30	122	37	72	9	3	31
SI	43	30	73	0	13	26	3	8	18
NO	208	122	0	330	137	89	17	25	62
DAN_IMP	113	37	13	137	150	0	0	0	0
ROBO<10000	43	72	26	89	0	115	0	0	0
DES_PUB	16	9	8	17	0	0	25	0	0
NARCOT	30	3	8	25	0	0	0	33	0
OTROS	49	31	18	62	0	0	0	0	60

Valeurs propres et coordonnées des catégories sur les axes correspondants :

Valeurs pr.	1	2	3	4	5	6
Valeur	0.4626	0.3858	0.3333	0.3333	0.2745	0.2105
% de vari.	23%	19%	17%	17%	14%	11%
% cumulé	23%	42%	59%	76%	89%	100%
Catégories	1	2	3	4	5	6
HOMBRE	-0.6075	0.1969	0.0000	0.0000	0.1378	-0.4228
MUJER	1.0032	-0.3251	0.0000	0.0000	-0.2275	0.6982
SI	0.6917	1.4704	0.0000	0.0000	1.3466	0.2577
NO	-0.1530	-0.3253	0.0000	0.0000	-0.2979	-0.0570
DAN_IMP	-0.7480	-0.6449	0.3174	-0.1023	0.6097	0.4777
ROBO<1000	1.1321	-0.3125	-0.3063	-0.6238	0.0963	-0.7956
DES_PUB	0.2312	1.6407	3.2181	-0.4996	-1.3283	-0.0664
NARCOT	-1.0558	1.6479	-1.7898	-1.4328	-1.1625	0.8184
OTROS	0.1385	0.4660	-0.4306	1.3605	-0.3870	-0.0689

**TESIS CON
FALLA DE ORIGEN**

Contribucións, contribucións totales e custos de custos

Categorías de Áreas	Efecto F	Puntos de vista F	Contribucións en millóns de euros						Contribucións en millóns de euros						
			1	2	3	4	5	6	1	2	3	4	5	6	
H.A.M.P.R.	281	20 76	333 79	42 52	0 00	0 00	39 83	305 29	0 00	0 00	0 00	0 00	0 00	0 00	0 00
M.I.P.R.	142	11 57	51 53	89 28	0 00	0 00	41 77	100 00	0 00	0 00	0 00	0 00	0 00	0 00	0 00
Suma de valores 1	423	32 33	385 32	131 80	0 00	0 00	81 60	405 29	0 00	0 00	0 00	0 00	0 00	0 00	0 00
M	73	8 04	125 99	461 84	0 00	0 00	833 91	8 34	0 11	0 48	0 00	0 00	0 40	0 40	0 20
MCO	28	27 20	27 84	100 83	0 00	0 00	175 83	8 45	0 11	0 48	0 00	0 00	0 20	0 20	0 20
Suma de valores 2	101	35 24	153 83	562 67	0 00	0 00	1009 74	16 79	0 22	0 96	0 00	0 00	0 60	0 60	0 40
CO.M.P.	180	12 41	132 42	376 31	75 58	95	338 05	270 84	0 23	0 23	0 08	0 01	0 23	0 23	0 11
A.M.O.M.I.E.S.I.	113	2 01	131 22	46 52	53 84	223 74	6 47	576 31	0 01	0 04	0 04	0 16	0 01	0 01	0 28
DE.B.P.R.	24	2 07	4 82	343 71	1284 52	19 21	267 87	0 87	0 03	0 18	0 46	0 02	0 12	0 12	0 20
M.A.M.C.O.T.	31	2 73	132 83	307 14	518 81	87 88	270 83	175 03	0 10	0 24	0 28	0 20	0 12	0 12	0 28
CO.M.P.	63	8 42	5 63	76 35	161 46	1284 52	72 78	3 01	0 00	0 08	0 08	0 08	0 04	0 04	0 04
Suma de valores 3	613	31 13	676 31	1070 42	2015 00	2015 00	468 37	1500 32	0 34	0 36	0 36	0 36	0 36	0 36	0 36

Detalle de custos de custos

Categorías de Áreas	Distribución F	1	2	3	4	5	6
H.A.M.P.R.	1 01						
M.I.P.R.	1 06						
M	4 52						
MCO	0 22						
CO.M.P.	1 99						
A.M.O.M.I.E.S.I.	2 74						
DE.B.P.R.	19 12						
M.A.M.C.O.T.	11 21						
CO.M.P.	0 24						

Notas: Los valores positivos representan los aumentos.

Contribucións de custos totales de custos

Categorías de Áreas	Puntos de vista F	Distribución F	Contribucións en millóns de euros						Contribucións en millóns de euros					
			1	2	3	4	5	6	1	2	3	4	5	6
Cha 1	1 40	17 17	-0 02	1 38	0 51	-0 17	3 77	0 64	0 48	1 58	0 20	0 02	12 84	0 48
Cha 2	2 76	27 98	1 99	2 41	-0 58	-1 18	3 24	2 31	2 10	3 74	0 28	1 08	10 08	8 32
Cha 3	1 24	33 14	0 38	2 87	-1 18	-0 86	6 22	-0 28	0 88	10 13	13 00	0 31	0 04	3 11
Cha 4	1 74	38 12	-1 28	4 71	-2 70	-2 28	0 54	1 26	0 85	14 28	5 44	3 87	0 27	1 80
Cha 5	2 38	39 88	0 28	3 77	-0 06	3 72	2 42	-0 98	0 19	10 12	0 35	10 21	0 39	0 41
Cha 6	1 24	13 10	1 08	4 63	-0 41	-0 13	2 46	0 58	0 58	0 22	0 13	0 01	0 27	0 38
Cha 7	3 72	43 38	8 27	1 73	-0 88	-1 38	3 00	0 43	15 40	1 83	0 26	1 46	8 11	0 24
Cha 8	0 74	21 28	1 64	2 88	3 22	-0 80	-0 23	1 12	1 43	4 21	7 71	1 08	0 06	0 46
Cha 9	0 20	5 78	0 31	1 93	-1 02	-0 88	-0 03	1 26	0 85	1 46	0 78	0 45	0 00	1 46
Cha 10	1 46	20 42	2 20	2 12	-0 89	2 49	1 14	1 68	2 40	2 80	0 28	0 18	1 18	2 84
Cha 11	26 00	87 87	-7 58	-4 28	1 88	0 41	2 40	-0 02	30 79	11 82	2 83	6 77	7 77	0 02
Cha 12	7 84	38 53	1 03	-1 34	-1 00	-2 04	-0 23	-0 24	0 97	1 18	0 74	3 58	0 00	32 84
Cha 13	2 71	58 47	0 69	2 48	-1 18	0 95	-3 14	-1 32	0 42	4 86	20 27	0 68	0 42	2 04
Cha 14	8 71	82 30	-0 27	3 91	-0 30	-1 12	-1 04	1 19	8 78	8 84	17 88	12 72	14 73	1 84
Cha 15	8 16	90 88	-1 88	1 12	-1 81	5 83	-2 12	-0 02	1 84	0 78	1 10	31 78	4 28	6 03
Cha 16	7 84	37 91	0 28	-3 03	1 04	-0 33	0 30	0 43	0 04	8 94	0 80	0 08	0 08	24 63
Cha 17	14 14	83 16	7 33	-3 80	-1 54	-2 72	-2 09	-0 85	20 68	9 78	1 35	5 80	3 68	0 85
Cha 18	14 46	33 88	1 80	1 03	-1 38	-0 71	-2 96	1 02	0 81	1 28	18 42	0 37	7 85	1 22
Cha 19	0 50	6 72	0 14	0 78	-1 48	-1 22	-1 82	1 93	0 00	3 27	1 88	1 11	2 32	2 85
Cha 20	8 20	68 20	2 42	-0 40	-1 24	1 37	-2 40	2 18	3 16	0 19	1 42	21 47	7 09	0 08

Custos totales de custos de custos

Categorías de Áreas	1	2	3	4	5	6
Cha 1	0 38	3 13	0 21	0 20	0 78	2 02
Cha 2	0 14	0 24	0 01	0 00	0 43	1 18
Cha 3	0 03	0 47	0 81	0 21	0 00	0 00
Cha 4	0 04	0 96	0 19	0 14	0 01	0 04
Cha 5	0 70	3 43	0 02	0 38	0 16	0 01
Cha 6	0 08	0 03	0 01	0 00	0 48	0 41
Cha 7	0 48	0 07	0 21	0 04	0 21	0 03
Cha 8	0 13	0 31	0 49	0 01	0 00	0 08
Cha 9	0 02	0 38	0 18	0 13	0 00	0 26
Cha 10	0 24	0 22	0 02	0 34	0 08	0 12
Cha 11	0 04	0 21	0 04	0 00	3 10	0 00
Cha 12	0 03	0 08	0 03	0 12	0 00	0 17
Cha 13	0 01	0 12	0 06	0 02	0 17	0 33
Cha 14	0 28	0 17	0 28	0 14	0 18	0 32
Cha 15	0 08	0 02	0 04	0 31	0 07	0 10
Cha 16	0 00	0 41	0 03	0 00	0 00	0 08
Cha 17	0 08	0 18	0 07	0 04	0 28	0 01
Cha 18	0 08	0 06	0 37	0 09	0 28	0 00
Cha 19	0 08	0 07	0 24	0 19	0 28	0 08
Cha 20	0 17	0 20	0 03	0 16	0 19	0 03

TESIS CON FALLA DE ORIGEN

Anexo 7. Base de datos para las flores de Iris. Resultados del AFCM.

Capitulo III

Especimen	Estado de floración (centímetros)			Categorías		
	Longitud	Diámetro	Petalos	Superior	Medio	Posterior
1. iris-001	51	26	14	02	1	1
2. iris-002	49	3	14	02	1	1
3. iris-003	47	32	13	03	1	1
4. iris-004	48	31	15	03	1	1
5. iris-005	5	26	14	03	1	1
6. iris-006	34	29	17	04	1	1
7. iris-007	48	24	14	03	1	1
8. iris-008	6	24	15	02	1	1
9. iris-009	44	26	14	02	1	1
10. iris-010	48	31	15	01	1	1
11. iris-011	39	27	15	02	1	1
12. iris-012	46	34	18	02	1	1
13. iris-013	48	3	14	01	1	1
14. iris-014	43	3	11	01	1	1
15. iris-015	50	4	12	03	1	1
16. iris-016	57	44	18	04	1	1
17. iris-017	54	30	13	04	1	1
18. iris-018	51	25	14	02	1	1
19. iris-019	37	38	17	03	1	1
20. iris-020	51	30	15	03	1	1
21. iris-021	54	34	17	04	1	1
22. iris-022	51	37	18	04	1	1
23. iris-023	44	30	1	05	1	1
24. iris-024	51	23	17	05	1	1
25. iris-025	46	34	18	02	1	1
26. iris-026	5	34	16	04	1	1
27. iris-027	5	34	16	04	1	1
28. iris-028	5	34	16	04	1	1
29. iris-029	5	34	16	04	1	1
30. iris-030	5	34	16	04	1	1
31. iris-031	48	31	18	02	1	1
32. iris-032	54	14	15	04	1	1
33. iris-033	53	41	18	01	1	1
34. iris-034	55	43	14	02	2	1
35. iris-035	48	31	15	02	1	1
36. iris-036	8	34	12	02	1	1
37. iris-037	56	25	13	02	2	1
38. iris-038	48	36	14	01	1	1
39. iris-039	44	3	12	02	1	1
40. iris-040	5	34	15	03	1	1
41. iris-041	5	34	15	03	1	1
42. iris-042	48	23	12	03	1	1
43. iris-043	44	23	12	02	1	1
44. iris-044	5	38	16	00	1	1
45. iris-045	51	36	18	04	1	1
46. iris-046	48	32	14	02	1	1
47. iris-047	53	37	18	02	1	1
48. iris-048	5	33	14	02	1	1
49. iris-049	63	33	9	25	2	2
50. iris-050	36	27	11	19	2	2
51. iris-051	71	3	38	21	2	2
52. iris-052	63	29	36	16	2	2
53. iris-053	66	3	36	23	2	2
54. iris-054	78	3	40	11	2	2
55. iris-055	48	23	10	17	2	2
56. iris-056	73	29	43	14	2	2
57. iris-057	67	23	38	11	2	2
58. iris-058	72	24	41	29	2	2
59. iris-059	63	12	51	2	2	2
60. iris-060	64	37	53	10	2	2
61. iris-061	66	3	50	21	2	2
62. iris-062	57	25	8	2	2	2
63. iris-063	56	28	81	24	2	2
64. iris-064	64	32	53	22	2	2
65. iris-065	6	3	39	14	2	2
66. iris-066	77	24	47	22	2	2
67. iris-067	77	26	50	21	2	2
68. iris-068	4	22	5	13	2	2
69. iris-069	69	32	47	23	2	2
70. iris-070	56	34	49	1	2	2
71. iris-071	77	28	47	2	2	2
72. iris-072	83	27	49	18	2	2
73. iris-073	67	23	37	11	2	2
74. iris-074	72	22	36	14	2	2
75. iris-075	62	26	48	14	2	2
76. iris-076	81	3	49	14	2	2
77. iris-077	64	18	44	7	2	2
78. iris-078	72	1	50	1	2	2
79. iris-079	74	26	61	12	2	2
80. iris-080	18	38	44	2	2	2
81. iris-081	64	28	56	22	2	2
82. iris-082	63	24	51	18	2	2
83. iris-083	61	28	50	14	2	2
84. iris-084	77	3	61	33	2	2
85. iris-085	82	24	58	24	2	2
86. iris-086	64	31	55	18	2	2
87. iris-087	6	3	48	16	2	2
88. iris-088	68	31	54	21	2	2
89. iris-089	67	31	56	24	2	2
90. iris-090	69	31	51	23	2	2
91. iris-091	56	27	41	14	2	2
92. iris-092	68	22	58	23	2	2
93. iris-093	67	33	47	25	2	2
94. iris-094	67	33	47	25	2	2
95. iris-095	63	28	5	18	2	2
96. iris-096	46	3	32	3	2	2
97. iris-097	42	34	34	24	2	2
100. iris-098	38	3	31	16	2	2

Categorías
1: corto
2: medio
3: largo

TESIS CON FALLA DE ORIGEN

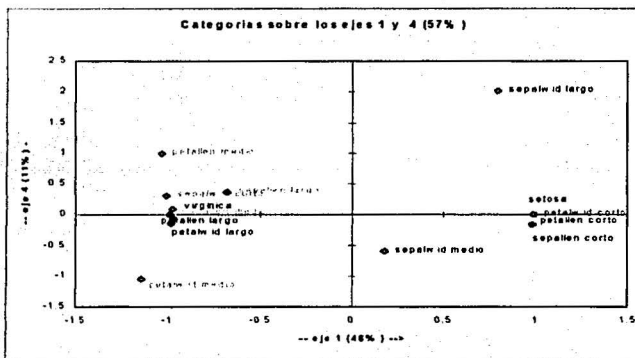
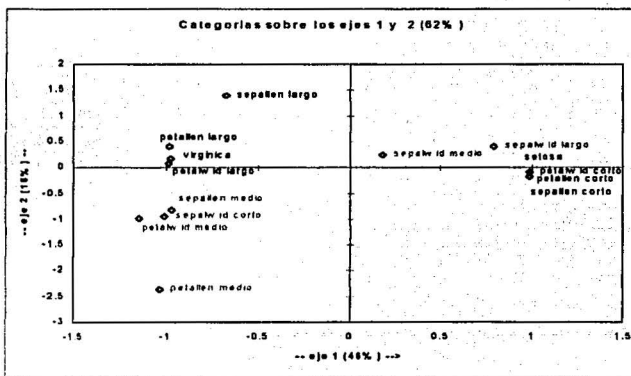
LISTAS CARTAS PARA LOS OBSERVADORES

Observación	1	2	3	4	5	6	7	8	9
Obs 1	0.07	0.00	0.01	0.12	0.00	0.00	0.03	0.00	0.01
Obs 2	0.07	0.03	0.01	0.12	0.00	0.00	0.03	0.00	0.01
Obs 3	0.07	0.00	0.01	0.12	0.00	0.00	0.03	0.00	0.01
Obs 4	0.07	0.00	0.01	0.12	0.00	0.00	0.03	0.00	0.01
Obs 5	0.07	0.00	0.01	0.12	0.00	0.00	0.03	0.00	0.01
Obs 6	0.06	0.00	0.04	0.35	0.01	0.02	0.03	0.00	0.00
Obs 7	0.07	0.00	0.01	0.12	0.00	0.00	0.03	0.00	0.01
Obs 8	0.07	0.00	0.04	0.35	0.00	0.00	0.03	0.00	0.01
Obs 9	0.28	0.18	0.03	0.03	0.03	0.44	0.00	0.00	0.00
Obs 10	0.07	0.00	0.01	0.12	0.00	0.00	0.03	0.00	0.01
Obs 11	0.06	0.04	0.05	0.01	0.02	0.00	0.03	0.00	0.01
Obs 12	0.07	0.00	0.01	0.12	0.00	0.00	0.03	0.00	0.01
Obs 13	0.07	0.00	0.01	0.12	0.00	0.00	0.03	0.00	0.01
Obs 14	0.07	0.00	0.01	0.12	0.00	0.00	0.03	0.00	0.01
Obs 15	0.20	0.14	0.04	0.48	0.72	0.00	0.13	0.00	0.00
Obs 16	0.20	0.14	0.04	0.48	0.72	0.00	0.13	0.00	0.00
Obs 17	0.06	0.00	0.04	0.35	0.01	0.02	0.03	0.00	0.00
Obs 18	0.07	0.00	0.01	0.12	0.00	0.00	0.03	0.00	0.01
Obs 19	0.07	0.14	0.04	0.48	0.72	0.00	0.13	0.00	0.00
Obs 20	0.06	0.00	0.04	0.35	0.01	0.02	0.03	0.00	0.00
Obs 21	0.07	0.00	0.01	0.12	0.00	0.00	0.03	0.00	0.01
Obs 22	0.06	0.00	0.04	0.35	0.01	0.02	0.03	0.00	0.00
Obs 23	0.07	0.00	0.01	0.12	0.00	0.00	0.03	0.00	0.01
Obs 24	0.07	0.00	0.01	0.12	0.00	0.00	0.03	0.00	0.01
Obs 25	0.07	0.00	0.01	0.12	0.00	0.00	0.03	0.00	0.01
Obs 26	0.07	0.00	0.01	0.12	0.00	0.00	0.03	0.00	0.01
Obs 27	0.07	0.00	0.01	0.12	0.00	0.00	0.03	0.00	0.01
Obs 28	0.07	0.00	0.01	0.12	0.00	0.00	0.03	0.00	0.01
Obs 29	0.07	0.00	0.01	0.12	0.00	0.00	0.03	0.00	0.01
Obs 30	0.07	0.00	0.01	0.12	0.00	0.00	0.03	0.00	0.01
Obs 31	0.07	0.00	0.01	0.12	0.00	0.00	0.03	0.00	0.01
Obs 32	0.07	0.00	0.01	0.12	0.00	0.00	0.03	0.00	0.01
Obs 33	0.06	0.00	0.04	0.35	0.01	0.02	0.03	0.00	0.00
Obs 34	0.17	0.03	0.04	0.34	0.98	0.21	0.11	0.00	0.00
Obs 35	0.07	0.00	0.01	0.12	0.00	0.00	0.03	0.00	0.01
Obs 36	0.07	0.00	0.01	0.12	0.00	0.00	0.03	0.00	0.01
Obs 37	0.15	0.08	0.00	0.07	0.07	0.14	0.45	0.00	0.00
Obs 38	0.07	0.00	0.01	0.12	0.00	0.00	0.03	0.00	0.01
Obs 39	0.07	0.00	0.01	0.12	0.00	0.00	0.03	0.00	0.01
Obs 40	0.07	0.00	0.01	0.12	0.00	0.00	0.03	0.00	0.01
Obs 41	0.07	0.00	0.01	0.12	0.00	0.00	0.03	0.00	0.01
Obs 42	0.18	0.03	0.00	0.03	0.00	0.18	0.00	0.00	0.00
Obs 43	0.07	0.00	0.01	0.12	0.00	0.00	0.03	0.00	0.01
Obs 44	0.07	0.00	0.01	0.12	0.00	0.00	0.03	0.00	0.01
Obs 45	0.06	0.00	0.04	0.35	0.01	0.02	0.03	0.00	0.00
Obs 46	0.07	0.00	0.01	0.12	0.00	0.00	0.03	0.00	0.01
Obs 47	0.06	0.00	0.04	0.35	0.01	0.02	0.03	0.00	0.00
Obs 48	0.06	0.00	0.01	0.12	0.00	0.00	0.03	0.00	0.01
Obs 49	0.06	0.00	0.04	0.35	0.01	0.02	0.03	0.00	0.00
Obs 50	0.07	0.00	0.01	0.12	0.00	0.00	0.03	0.00	0.01
Obs 51	0.53	0.00	0.03	0.08	0.13	0.23	0.00	0.00	0.00
Obs 52	0.06	0.00	0.01	0.00	0.22	0.31	0.30	0.00	0.00
Obs 53	0.40	0.52	0.03	0.01	0.04	0.30	0.00	0.00	0.00
Obs 54	0.04	0.09	0.01	0.00	0.22	0.03	0.00	0.00	0.00
Obs 55	0.33	0.09	0.03	0.08	0.13	0.23	0.00	0.00	0.00
Obs 56	0.40	0.52	0.03	0.01	0.04	0.30	0.00	0.00	0.00
Obs 57	0.10	0.33	0.15	0.07	0.20	0.37	0.04	0.00	0.00
Obs 58	0.32	0.09	0.00	0.04	0.30	0.38	0.00	0.00	0.00
Obs 59	0.32	0.09	0.00	0.04	0.30	0.38	0.00	0.00	0.00
Obs 60	0.40	0.52	0.03	0.01	0.04	0.30	0.00	0.00	0.00
Obs 61	0.53	0.00	0.03	0.08	0.13	0.23	0.00	0.00	0.00
Obs 62	0.06	0.00	0.01	0.00	0.22	0.31	0.00	0.00	0.00
Obs 63	0.40	0.52	0.03	0.01	0.04	0.30	0.00	0.00	0.00
Obs 64	0.04	0.09	0.01	0.00	0.22	0.03	0.00	0.00	0.00
Obs 65	0.04	0.09	0.01	0.00	0.22	0.03	0.00	0.00	0.00
Obs 66	0.53	0.00	0.03	0.08	0.13	0.23	0.00	0.00	0.00
Obs 67	0.53	0.00	0.03	0.08	0.13	0.23	0.00	0.00	0.00
Obs 68	0.18	0.35	0.01	0.44	0.00	0.02	0.02	0.00	0.00
Obs 69	0.52	0.09	0.00	0.04	0.00	0.35	0.00	0.00	0.00
Obs 70	0.11	0.40	0.12	0.03	0.04	0.00	0.00	0.00	0.00
Obs 71	0.40	0.52	0.03	0.01	0.04	0.00	0.00	0.00	0.00
Obs 72	0.25	0.45	0.13	0.07	0.06	0.00	0.00	0.00	0.00
Obs 73	0.08	0.32	0.04	0.00	0.00	0.35	0.00	0.00	0.00
Obs 74	0.25	0.45	0.13	0.07	0.06	0.00	0.00	0.00	0.00
Obs 75	0.40	0.52	0.03	0.01	0.04	0.00	0.00	0.00	0.00
Obs 76	0.04	0.09	0.01	0.00	0.04	0.00	0.00	0.00	0.00
Obs 77	0.25	0.45	0.13	0.07	0.06	0.00	0.00	0.00	0.00
Obs 78	0.17	0.24	0.09	0.01	0.18	0.11	0.00	0.00	0.00
Obs 79	0.64	0.09	0.41	0.30	0.22	0.03	0.00	0.00	0.00
Obs 80	0.11	0.03	0.48	0.07	0.31	0.02	0.00	0.00	0.00
Obs 81	0.09	0.64	0.01	0.00	0.03	0.73	0.00	0.00	0.00
Obs 82	0.18	0.35	0.01	0.44	0.00	0.02	0.00	0.00	0.00
Obs 83	0.06	0.09	0.01	0.00	0.22	0.33	0.00	0.00	0.00
Obs 84	0.20	0.11	0.62	0.03	0.04	0.00	0.00	0.00	0.00
Obs 85	0.20	0.11	0.62	0.03	0.04	0.00	0.00	0.00	0.00
Obs 86	0.40	0.52	0.03	0.01	0.04	0.00	0.00	0.00	0.00
Obs 87	0.53	0.00	0.03	0.08	0.13	0.23	0.00	0.00	0.00
Obs 88	0.53	0.00	0.03	0.08	0.13	0.23	0.00	0.00	0.00
Obs 89	0.17	0.24	0.09	0.01	0.18	0.11	0.00	0.00	0.00
Obs 90	0.40	0.52	0.03	0.01	0.04	0.00	0.00	0.00	0.00
Obs 91	0.40	0.52	0.03	0.01	0.04	0.00	0.00	0.00	0.00
Obs 92	0.40	0.52	0.03	0.01	0.04	0.00	0.00	0.00	0.00
Obs 93	0.54	0.00	0.01	0.00	0.22	0.03	0.00	0.00	0.00
Obs 94	0.40	0.52	0.03	0.01	0.04	0.00	0.00	0.00	0.00
Obs 95	0.40	0.52	0.03	0.01	0.04	0.00	0.00	0.00	0.00
Obs 96	0.40	0.52	0.03	0.01	0.04	0.00	0.00	0.00	0.00
Obs 97	0.64	0.09	0.01	0.00	0.22	0.03	0.00	0.00	0.00
Obs 98	0.53	0.00	0.03	0.08	0.13	0.23	0.00	0.00	0.00
Obs 99	0.53	0.00	0.03	0.08	0.13	0.23	0.00	0.00	0.00
Obs 100	0.53	0.00	0.03	0.08	0.13	0.23	0.00	0.00	0.00

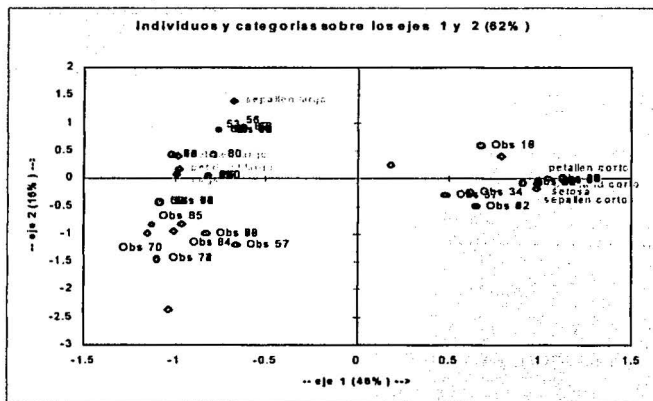
TESIS ()
FALLA DE () EN

Tabla 1. Matriz de Burt.

	Hombre	Mujer	Soltero	Casado	Viudo	PRI	PAN	PRD
Hombre	5	0	3	1	1	2	1	2
Mujer	0	5	2	2	1	2	3	0
Soltero	3	2	5	0	0	2	3	0
Casado	1	2	0	3	0	2	0	1
Viudo	1	1	0	0	2	0	1	1
PRI	2	2	2	2	0	4	0	0
PAN	1	3	3	0	1	0	4	0
PRD	2	0	0	1	1	0	0	2



TESIS CON
FALLA DE ORIGEN



25. Ronald Christensen, *Log-Linear Models*, Springer-Verlag, New York, USA, 1990.