

# Definición general del análisis de correspondencia y su interpretación

**Definición general:** Técnica descriptiva que resume gráficamente la información contenida en una tabla de contingencia.

## Interpretación:

- 1 Una manera de representar la relación entre dos variables categóricas en un espacio de dimensión menor, análogo a componentes principales, pero definiendo la distancia entre los puntos de manera coherente con la interpretación de los datos y en lugar de usar la distancia Euclidiana usamos la **distancia ji-cuadrada**
  - Desde este enfoque, el análisis de correspondencia es el equivalente de componentes principales para datos cualitativos.
- 2 Es un procedimiento objetivo de asignar valores numéricos a variables cualitativas, lo cual tiene mas relación con MDS.

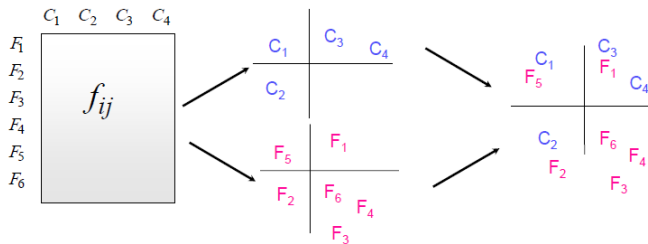
# Tabla de contingencia

Fue desarrollado por Benzecri (1973) para descubrir y entender la estructura y el modelo de los datos en tablas de contingencias

Tabla de contingencia con  $I$  filas y  $J$  columnas

X2	$1$	.....	$j$	.....	$J$	<i>Total</i>
X1						
$1$	$x_{11}$		$x_{1j}$		$x_{1J}$	$x_{1\cdot}$
$2$	$x_{21}$					
$\cdot$	$\cdot$					
$\cdot$	$\cdot$					
$i$	$\cdot$		$x_{ij}$			
$\cdot$	$\cdot$					
$I$	$x_{I1}$				$x_{IJ}$	$x_{I\cdot}$
<i>Total</i>	$x_{\cdot 1}$		$x_{\cdot j}$		$x_{\cdot J}$	$x_{\cdot \cdot}$

# Análisis de Correspondencias



**Objetivo:** Representar las filas y las columnas de una tabla de contingencia de dos vías como puntos en un espacio vectorial de baja dimensión, de forma que los correspondientes espacios se puedan superponer para obtener una representación conjunta

# Análisis de correspondencia

- Si *puntos filas* están muy cercanos indica que las filas tienen un perfil similar a través de las columnas (su distribución está condicionada a las columnas)
- Si *puntos columna* están muy cercanos indica que las columnas tienen un perfil similar a través de las filas (su distribución está condicionada a las filas)
- Puntos filas que son muy cercanos a los puntos columnas representan combinaciones que ocurren mas frecuentemente de las que podría esperarse de un modelo de independencia.
- Esto es, un modelo en el cual las  $I$  categorías de la variable  $\mathbf{X}_1$  (filas) están incorrelacionadas con las  $J$  categorías de la variable  $\mathbf{X}_2$  (columnas).

# Análisis de correspondencia

Los resultados de un análisis de correspondencia incluyen obtener la “mejor” representación de los datos en dos dimensiones, a través de las coordenadas de los puntos graficados y una medida llamada *inercia* que cuantifica la cantidad de información retenida en cada dimensión.

# Análisis de correspondencia. Desarrollo algebraico

- Sea  $\mathbf{X}$ , una tabla a dos vías  $I \times J$ , cuyos elementos  $x_{ij}$  representan las frecuencias de aparición (o conteos) de la  $i$ -ésima categoría de  $\mathbf{X}_1$  y la  $j$ -ésima categoría de  $\mathbf{X}_2$
- Si  $n$  es el total de frecuencias en la tabla de datos  $\mathbf{X}$ , primero se construye una matriz  $\mathbf{F} = \{f_{ij}\}$  de frecuencias relativas, dividiendo cada elemento de  $\mathbf{X}$  por  $n$ . Entonces

$$f_{ij} = \frac{x_{ij}}{n}, \quad i = 1, 2, \dots, I \quad j = 1, 2, \dots, J, \quad \mathbf{F} = \frac{1}{n} \mathbf{X}$$

La matriz  $\mathbf{F}$  es llamada la matriz de correspondencia y satisface

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

# Análisis de correspondencia. Desarrollo algebraico

- El análisis de la matriz  $F$  se puede hacer por filas o columnas, cualquier análisis de esta matriz debe ser equivalente al aplicado a su transpuesta.
- Se presentará el análisis por filas de la matriz  $F$  que será simétrico al análisis por columnas.

## Proyección de las filas de la matriz $F$

- Las  $I$  filas de  $F$  pueden tomarse como  $I$  puntos en el espacio  $\mathbb{R}^J$
- Se busca una representación de estos  $I$  puntos en un espacio de dimensión menor que permita apreciar sus distancias relativas
- El objetivo es el mismo que en componentes principales pero tomando en cuenta la peculiaridad de los datos

# Análisis de correspondencia. Desarrollo algebraico

La peculiaridad de los datos proviene del hecho de que la frecuencia relativa de cada fila es distinta, lo que implica que:

- ① Todas las filas (puntos en  $\mathbb{R}^J$ ) no tienen el mismo peso, ya que algunas contienen mas datos que otras. Al representar el conjunto de las filas (puntos) debemos dar más peso aquellas filas que contienen mas datos.
- ② La distancia euclidiana entre los puntos no es una buena medida de su proximidad, entonces debe modificarse.

Denotamos por  $r_i$  a la **frecuencia relativa** o **masa** de la fila  $i$  de  $\mathbf{F}$ , como

$$r_i = \sum_{j=1}^J f_{ij} = \sum_{j=1}^J \frac{x_{ij}}{n}, i = 1, 2, \dots, I, \quad \mathbf{r}_{(I \times 1)} = \mathbf{F} \mathbf{1}_{(J \times 1)}$$



# Análisis de correspondencia. Desarrollo algebraico

- Debemos dar a cada fila un peso proporcional a su **frecuencia relativa o masa**.
- Los elementos del vector  $\mathbf{r}$  se toman como sus pesos, ya que son números positivos que suman 1
- Denotamos por  $\mathbf{R}$  a la matriz de frecuencias relativas condicionadas al total de la fila, llamada *matriz de perfiles*, que se obtiene como

$$\mathbf{R} = \mathbf{D}_r^{-1} \mathbf{F}, \text{ donde } \mathbf{D}_r = \text{diag}(r_1, r_2, \dots, r_I)$$

Entonces se transforma la matriz original  $\mathbf{F}$  en otra matriz cuyas casillas por fila suman uno.

- Cada fila de  $\mathbf{R}$ ,  $\mathbf{R}_i$  puede considerarse un punto o un vector en el espacio  $\mathbb{R}^J$ . Como la suma de los componentes de  $\mathbf{R}_i$  es uno, todos los puntos están en un espacio de dimensión  $J - 1$ .

# Análisis de correspondencia. Desarrollo algebraico

- La idea es proyectar estos puntos en un espacio de dimensión menor de tal forma que las filas que tengan la misma estructura estén lo más próximas y las que tengan una estructura muy diferente estén alejadas.
- Por tanto debemos definir una medida de distancia entre dos filas  $\mathbf{R}_a$  y  $\mathbf{R}_b$ . Una posibilidad es utilizar la distancia euclidiana
- La distancia euclidiana, tiene el inconveniente de tratar igual a todos los componentes de estos vectores.
- Sin embargo en este caso, el atributo  $j$  tiene diferente peso en cada fila, y por tanto la distancia euclidiana no será adecuada
- Para obtener comparaciones razonables debemos tener en cuenta la frecuencia relativa de aparición del atributo que estudiamos

# Análisis de correspondencia. Desarrollo algebraico

- Una manera intuitiva de construir las comparaciones, es ponderar las diferencias en frecuencia relativa entre dos atributos inversamente proporcional a la frecuencia de este atributo.
- Es decir, en lugar de sumar los términos

$$(R_{aj} - R_{bj})^2 = (f_{aj}/r_a - f_{bj}/r_b)^2$$

se considera la suma de los términos

$$(R_{aj} - R_{bj})^2 / c_j$$

donde

$$c_j = \sum_{i=1}^I f_{ij} = \sum_{i=1}^I \frac{x_{ij}}{n}, j = 1, 2, \dots, J, \quad \underset{(J \times 1)}{\mathbf{c}} = \underset{(J \times I)}{\mathbf{F}'} \underset{(I \times 1)}{\mathbf{1}}$$

$c_j$  representa la *frecuencia relativa o masa* de la columna  $j$  de  $\mathbf{F}$

# Análisis de correspondencia. Desarrollo algebraico

La expresión de la distancia entre dos filas  $R_a$  y  $R_b$  de  $R$  está dada por

$$D^2(R_a, R_b) = \sum_{j=1}^J \left( \frac{f_{aj}}{r_a} - \frac{f_{bj}}{r_b} \right)^2 \frac{1}{c_j} = \sum_{j=1}^J \frac{(R_{aj} - R_{bj})^2}{c_j}$$

que en forma matricial se puede escribir como

$$D^2(R_a, R_b) = (R_a - R_b)' D_c^{-1} (R_a - R_b)$$

donde

$$D_c = \text{diag}(c_1, c_2, \dots, c_j)$$

- La distancia  $D^2(R_a, R_b)$  es la distancia  $\chi^2$  y equivale a la distancia euclidiana entre los vectores transformados

$$y_i = D_c^{-1/2} R_i$$

# Análisis de correspondencia. Desarrollo algebraico

- Por lo anterior, el problema se simplifica definiendo una matriz de datos transformados, sobre la que tiene sentido considerar la distancia euclidiana entre filas.
- Definiendo

$$\mathbf{Y} = \mathbf{R}\mathbf{D}_c^{-1/2} = \mathbf{D}_r^{-1}\mathbf{F}\mathbf{D}_c^{-1/2}$$

Obtenemos una matriz  $\mathbf{Y}$  cuyas entradas tienen la forma

$$y_{ij} = \left\{ \frac{f_{ij}}{r_i c_j^{1/2}} \right\}$$

Los elementos de  $\mathbf{Y}$  ya no suman **uno**, ni por renglón ni por columna

- Representan las frecuencias relativas condicionadas por fila, pero estandarizadas por su variabilidad, que depende de la raíz cuadrada de la frecuencia relativa de la columna

# Análisis de correspondencia. Desarrollo algebraico

- De esta manera, las entradas de  $\mathbf{Y}$  son directamente comparables entre si.
- La matriz de datos transformada  $\mathbf{Y}$  se puede tratar como una matriz de datos estándar, donde las filas representan a las observaciones y las columnas a las variables.
- La pregunta es como proyectarla de manera que se preserven las distancias relativas entre las filas, es decir, las filas con estructura similar aparezcan próximas en la proyección.
- Lo anterior implica encontrar un vector  $\mathbf{a}$  de norma la unidad, esto es,  $\mathbf{a}'\mathbf{a} = 1$  tal que el vector de puntos proyectados sobre  $\mathbf{a}$ ,

$$\mathbf{y}_p(\mathbf{a}) = \mathbf{Y}\mathbf{a}$$

Tenga variabilidad máxima

# Análisis de correspondencia. Desarrollo algebraico

- El vector  $\mathbf{a}$  se encuentra maximizando

$$Var(\mathbf{y}_p(\mathbf{a})) = \mathbf{y}_p(\mathbf{a})' \mathbf{y}_p(\mathbf{a}) = \mathbf{a}' \mathbf{Y}' \mathbf{Y} \mathbf{a}$$

Con la condición de que  $\mathbf{a}' \mathbf{a} = 1$ .

- El problema se convierte en un problema de componentes principales: el vector  $\mathbf{a}$  es vector propio de la matriz  $\mathbf{Y}' \mathbf{Y}$ .
- Sin embargo, este tratamiento de la matriz  $\mathbf{Y}$  como una matriz de variables continuas no es del todo correcto.
- En la matriz  $\mathbf{Y}$ , las filas tienen una frecuencia relativa distinta  $r_i$  y por tanto deben tener distinto peso.
- Las filas con mayor frecuencia relativa deben tener mas peso en la representación que aquellas con frecuencia relativa muy baja

# Análisis de correspondencia. Desarrollo algebraico

- De esta manera, las filas con gran número de individuos deben estar bien representadas, aunque esto sea a costa de representar peor las filas con pocos elementos.
- En consecuencia, se dará a cada fila de  $\mathbf{Y}$  un peso proporcional al número de datos que contiene, es decir  $\mathbf{D}_r^{1/2} \mathbf{Y}$
- Lo anterior implica maximizar ahora la suma de cuadrados ponderada:

$$\mathbf{m} = \mathbf{a}' \mathbf{Y}' \mathbf{D}_r \mathbf{Y} \mathbf{a} \quad \mathbf{a}' \mathbf{a} = 1$$

O equivalentemente

$$\mathbf{m} = \mathbf{a}' \mathbf{D}_c^{-1/2} \mathbf{F}' \mathbf{D}_r^{-1} \mathbf{F} \mathbf{D}_c^{-1/2} \mathbf{a}$$



# Análisis de correspondencia. Desarrollo algebraico

- Alternativamente, podemos construir una matriz de datos  $\mathbf{Z}$  definida por

$$\mathbf{Z} = \mathbf{D}_r^{-1/2} \mathbf{F} \mathbf{D}_c^{-1/2}$$

cuyos componentes son

$$z_{ij} = \left\{ \frac{f_{ij}}{\sqrt{r_i c_j}} \right\}$$

Lo cual estandariza las frecuencias relativas en cada casilla mediante el producto raíces cuadradas de las frecuencias relativas totales de la fila y la columna.

# Análisis de correspondencia. Desarrollo algebraico

- Considerando esta matriz  $\mathbf{Z}$ , el problema de encontrar el vector  $\mathbf{a}$  se transforma al problema de maximizar

$$\mathbf{m} = \mathbf{a}' \mathbf{Z}' \mathbf{Z} \mathbf{a} \quad \mathbf{a}' \mathbf{a} = 1$$

- Lo cual es el problema de componentes principales cuya solución está dada por los vectores y valores propios de  $\mathbf{Z}' \mathbf{Z}$  :

$$\mathbf{D}_c^{-1/2} \mathbf{F}' \mathbf{D}_r^{-1} \mathbf{F} \mathbf{D}_c^{-1/2} \mathbf{a} = \lambda \mathbf{a}$$

donde  $\mathbf{a}$  debe ser un vector propio de la matriz  $\mathbf{Z}' \mathbf{Z}$  y  $\lambda$  su valor propio asociado.

# Análisis de correspondencia. Desarrollo algebraico

Vamos a comprobar que la matriz  $\mathbf{Z}'\mathbf{Z}$  siempre tiene como valor propio máximo el 1 y como vector propio  $\mathbf{D}_c^{-1/2}\mathbf{1}$ .

Multiplicando  $\mathbf{D}_c^{-1/2}\mathbf{F}'\mathbf{D}_r^{-1}\mathbf{F}\mathbf{D}_c^{-1/2}\mathbf{a} = \lambda\mathbf{a}$  por  $\mathbf{D}_c^{-1/2}$  se obtiene

$$\mathbf{D}_c^{-1}\mathbf{F}'\mathbf{D}_r^{-1}\mathbf{F}\mathbf{D}_c^{-1/2}\mathbf{a} = \lambda(\mathbf{D}_c^{-1/2}\mathbf{a})$$

Las matrices

$$\mathbf{D}_r^{-1}\mathbf{F} \quad \text{y} \quad \mathbf{F}\mathbf{D}_c^{-1}$$

representan las matrices de perfiles por filas y por columnas y su suma por filas y columnas correspondiente es **uno**.

Por tanto

$$\mathbf{D}_r^{-1}\mathbf{F} \underset{J \times 1}{\mathbf{1}} = \underset{I \times 1}{\mathbf{1}} \quad \text{y} \quad \mathbf{D}_c^{-1}\mathbf{F}' \underset{I \times 1}{\mathbf{1}} = \underset{J \times 1}{\mathbf{1}}$$

que implica que

$$\mathbf{D}_c^{-1}\mathbf{F}'\mathbf{D}_r^{-1}\mathbf{F}\mathbf{1} = \mathbf{1}$$

Por tanto  $\mathbf{D}_c^{-1}\mathbf{F}'\mathbf{D}_r^{-1}\mathbf{F}$  tiene un valor propio de 1 asociado a un vector propio de  $\mathbf{1}$

# Análisis de correspondencia. Desarrollo algebraico

- Si esta expresión

$$D_c^{-1} F' D_r^{-1} F \mathbf{1} = \mathbf{1}$$

la multiplicamos por  $D_c^{1/2}$ , obtenemos

$$\begin{aligned} D_c^{1/2} D_c^{-1} F' D_r^{-1} F \mathbf{1} &= D_c^{1/2} \mathbf{1} \\ \Rightarrow D_c^{-1/2} F' D_r^{-1} F D_c^{-1/2} (D_c^{1/2} \mathbf{1}) &= (D_c^{1/2} \mathbf{1}) \end{aligned}$$

Es decir,  $\mathbf{Z}'\mathbf{Z}$  tiene un valor propio de 1 con vector propio  $D_c^{1/2} \mathbf{1}$ .

- Olvidando esta solución trivial que no proporciona información sobre la estructura de las filas, tomamos el mayor valor propio menor que 1 y su vector propio asociado  $\mathbf{a}$ .
- Entonces proyectando la matriz  $\mathbf{Y}$  sobre la dirección encontrada dada por el vector  $\mathbf{a}$ , obtenemos el vector proyectado:

$$\mathbf{y}_r(\mathbf{a}) = \mathbf{Y}\mathbf{a} = D_r^{-1} F D_c^{-1/2} \mathbf{a}$$

el vector  $\mathbf{y}_r(\mathbf{a})$  es la mejor representación de las filas de la tabla de contingencia en un espacio de **una** dimension.

# Análisis de correspondencia. Desarrollo algebraico

- De igual manera, si extraemos el vector propio asociado al siguiente valor propio mayor de  $\mathbf{Z}'\mathbf{Z}$ , obtenemos una segunda coordenada y podemos representar las filas en un espacio de **dos** dimensiones.
- Las coordenadas de la representación de cada fila están dadas por las filas de la matriz

$$\mathbf{C}_r = \mathbf{Y}\mathbf{A}_2 = \mathbf{D}_r^{-1}\mathbf{F}\mathbf{D}_c^{-1/2}\mathbf{A}_2$$

Donde  $\mathbf{A}_2 = \begin{bmatrix} \mathbf{a}_1 : \mathbf{a}_2 \end{bmatrix}$  contiene en columnas los dos vectores propios de  $\mathbf{Z}'\mathbf{Z}$

- La matriz  $\mathbf{C}_r$  es de tamaño  $I \times 2$  y las dos coordenadas de cada fila proporcionan la mejor representación de las filas de la matriz  $\mathbf{F}$  en un espacio de dos dimensiones

# Análisis de correspondencia. Desarrollo algebraico

- El procedimiento se extiende sin dificultad para obtener representaciones en más dimensiones, calculando vectores propios adicionales de la matriz  $\mathbf{Z}'\mathbf{Z}$ .
- En resumen, el procedimiento que hemos presentado para buscar una buena representación de las filas de la tabla de contingencia es:
  - 1 Caracterizar las filas por sus frecuencias relativas condicionadas, y considerarlas como puntos en el espacio (Obtención de  $\mathbf{R}$ )
  - 2 Definir la distancia entre los puntos mediante la distancia  $\chi^2$ , que tiene en cuenta que cada coordenada de las filas tiene distinta precisión (Obtención de  $\mathbf{Y}$ )
  - 3 Proyectar los puntos de  $\mathbf{Y}$  sobre las direcciones de máxima variabilidad, teniendo en cuenta que cada fila tiene un peso distinto e igual a sus frecuencias relativas. (Obtención de  $\mathbf{Z}$ )

# Análisis de correspondencia. Desarrollo algebraico

- El procedimiento operativo para obtener la mejor representación bidimensional de las filas de la tabla de contingencia es:
  - 1 Calcular la matriz  $\mathbf{Z}'\mathbf{Z}$  y obtener sus valores y vectores propios.
  - 2 Tomar los dos vectores propios  $\mathbf{a}_1$  y  $\mathbf{a}_2$ , asociados a los valores propios mayores (menores a 1 ) de esta matriz.
  - 3 Calcular las proyecciones  $\mathbf{D}_r^{-1}\mathbf{F}\mathbf{D}_c^{-1/2}\mathbf{a}_i$ ,  $i = 1, 2$  y representarlas gráficamente en un espacio bidimensional.

## Ejemplo: Análisis de correspondencia por fila

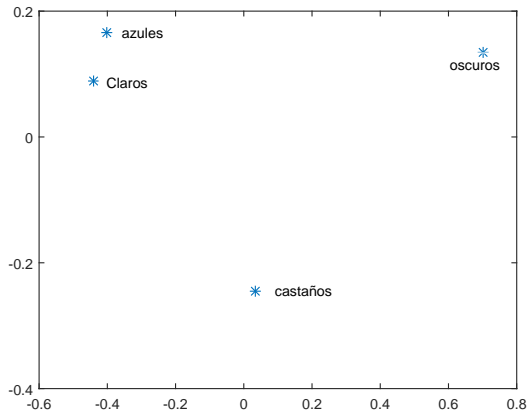
Aplicar un análisis de correspondencia por fila a la siguiente tabla de contingencia

	Color de cabello					
Color de ojos	Rubio	Pelirrojo	Castaño	Oscuro	Negro	Total
Claros	688	116	584	188	4	<b>1580</b>
Azules	326	38	241	110	3	<b>718</b>
Castaños	343	84	909	412	26	<b>1774</b>
Oscuros	98	48	403	681	85	<b>1315</b>
<b>total</b>	<b>1455</b>	<b>286</b>	<b>2137</b>	<b>1391</b>	<b>118</b>	<b>5387</b>

- La tabla presenta la clasificación de  $n = 5387$  estudiantes escoceses por el color de sus ojos, que tiene 4 categorías posibles ( $I = 4$ ), y por el color de su cabello, que tiene 5 categorías posibles ( $J = 5$ ).
- La tabla tiene interés histórico debido a que fue utilizada por Fisher (1940) para ilustrar un método de análisis de tablas de contingencias relacionado al análisis de correspondencia



# Ejemplo: Análisis de correspondencia por fila



# Proyección de las columnas de $F$

- El análisis anterior para las filas de  $F$  se puede aplicar también para las columnas. Las  $J$  columnas de  $F$  son puntos en  $\mathbb{R}^I$
- Se busca una representación de estos  $J$  puntos en un espacio de dimensión menor que  $I$  que permita apreciar sus distancias relativas.
- Para lo cual construimos las frecuencias relativas por columna

$$c_j = \sum_{i=1}^I f_{ij} = \sum_{i=1}^n \frac{x_{ij}}{n}, j = 1, 2, \dots, J \quad \underset{(J \times 1)}{\mathbf{c}} = \underset{(J \times I)}{\mathbf{F}}' \underset{I \times 1}{\mathbf{1}}$$

y

$$\mathbf{D}_c = \text{diag}(c_1, c_2, \dots, c_J)$$

- La mejor representación de los  $J$  puntos (columnas) en un espacio de dimensión menor, considerando la distancia  $\chi^2$  conducirá, por simetría a estudiar la matriz equivalente

$$\mathbf{Y} = \mathbf{D}_c^{-1} \mathbf{F}' \mathbf{D}_r^{-1/2}$$

# Proyección de las columnas de $F$

- Ahora la matriz que contiene las frecuencias relativas de las filas de  $F'$  es  $D_c$  y la matriz que contiene las frecuencias relativas de las columnas de  $F'$  es  $D_r$ .
- Intercambiando el papel de las matrices, las direcciones de proyección son los vectores propios de la matriz

$$ZZ' = D_r^{-1/2} F D_c^{-1} F' D_r^{-1/2}$$

- $Z'Z$  y  $ZZ'$  tienen los mismos valores propios no nulos, entonces  $ZZ'$  también tendrá un valor propio de 1 asociado a un vector propio que es una solución trivial.
- Denotando por  $b$  al vector propio asociado al mayor valor propio distinto de 1 de  $ZZ'$ , entonces la mejor representación de las columnas de la matriz en un espacio de **una dimensión** está dada por

$$y_c(b) = Y' b = D_c^{-1} F' D_r^{-1/2} b$$

# Proyección de las columnas de $F$

- Análogamente, la mejor representación en **dos** dimensiones de las columnas de la matriz  $Y$  están dadas por las coordenadas definidas por las filas de la matriz

$$C_c = Y' B_2 = D_c^{-1} F' D_r^{-1/2} B_2$$

donde  $B_2 = \begin{bmatrix} b_1 & b_2 \end{bmatrix}$  contiene en columnas los dos vectores propios asociados a los valores propios mayores de  $ZZ'$  menores a 1.

- La matriz  $C_c$  es de tamaño  $J \times 2$  y cada fila es la mejor representación de las columnas de la matriz  $F$  en un espacio de **dos dimensiones**

# Análisis conjunto

- Dada la simetría del problema conviene representar conjuntamente las filas y las columnas de la matriz  $F$ .
- Las matrices  $Z'Z$  y  $ZZ'$  tienen los mismos valores propios no nulos
- Además, los vectores propios de ambas matrices que corresponden al mismo valor propio están relacionados.
- Sea  $\mathbf{a}_i$  un vector propio de  $Z'Z$  asociado al valor propio  $\lambda_i$  :

$$Z'Z\mathbf{a}_i = \lambda_i\mathbf{a}_i$$

- Multiplicando por  $Z$

$$ZZ'(\mathbf{Za}_i) = \lambda_i(\mathbf{Za}_i)$$

- Entonces  $\mathbf{b}_i = \mathbf{Za}_i$  es un vector propio de  $ZZ'$  asociado a  $\lambda_i$

# Análisis conjunto

Una forma rápida para obtener los vectores propios es calcularlos directamente de la matriz de dimensión más pequeña,  $\mathbf{Z}'\mathbf{Z}$  ó  $\mathbf{Z}\mathbf{Z}'$  y obtener los otros vectores propios como  $\mathbf{Z}\mathbf{a}_i$  ó  $\mathbf{Z}'\mathbf{b}_i$ .

Alternativamente, se puede utilizar la descomposición en valores singulares de la matriz  $\mathbf{Z}$  ó  $\mathbf{Z}'$ . La descomposición aplicada a  $\mathbf{Z}$  es

$$\mathbf{Z} = \mathbf{B}_s \mathbf{\Lambda}_s \mathbf{A}_s' = \sum_{i=1}^s \lambda_i^{1/2} \mathbf{b}_i \mathbf{a}_i'$$

donde  $\mathbf{B}_s$  contiene en columnas los vectores propios de  $\mathbf{Z}\mathbf{Z}'$ ,  $\mathbf{A}_s$  los vectores propios de  $\mathbf{Z}'\mathbf{Z}$  y  $\mathbf{\Lambda}_s$  es diagonal y contiene los valores singulares,  $\lambda_i^{1/2}$ , ó raíces de los valores propios no nulos de  $\mathbf{Z}\mathbf{Z}'$  ó  $\mathbf{Z}'\mathbf{Z}$  y donde  $s = \min(I, J)$ .

# Análisis conjunto

- Entonces, la representación de las filas en el espacio conjunto se obtiene mediante

$$y_r(a) = Ya = D_r^{-1}FD_c^{-1/2}a$$

- Y la representación de las columnas en el mismo espacio se obtiene como

$$y_c(b) = Y'b = D_c^{-1}F'D_r^{-1/2}b$$

- La representación de la matriz  $Z$  con  $h$  dimensiones (normalmente  $h=2$ ) implica aproximar esta matriz mediante.

$$\hat{Z}_h = B_h\Lambda_hA'_h$$

- Esto es equivalente a obtener una aproximación a la tabla de contingencia observada  $F$  mediante la expresión

$$\hat{F}_h = D_r^{1/2}\hat{Z}_hD_c^{1/2}$$

# Análisis conjunto

- Si se desea eliminar el valor propio **uno** desde el principio, dado que no aporta información de interés, se puede reemplazar la matriz  $\mathbf{F}$  por  $\mathbf{F} - \hat{\mathbf{F}}_e$ , donde  $\hat{\mathbf{F}}_e$  es la matriz de frecuencias esperadas que viene dada por

$$\hat{\mathbf{F}}_e = \frac{1}{n} \mathbf{rc}'$$

- Se puede probar que la matriz  $\mathbf{F} - \hat{\mathbf{F}}_e$  tiene rango  $r - 1$ , y ya no tiene el valor propio igual a la unidad.
- La proporción de variabilidad explicada por cada dimensión se calcula como en componentes principales, descartando el valor propio igual a uno y tomando la proporción que representa cada valor propio con respecto a la suma de todos los valores propios distintos de 1.



# Análisis conjunto

El análisis de correspondencia de una tabla de contingencia de dimensiones  $I \times J$  se realiza en los siguientes pasos:

- 1 Se calcula la tabla de frecuencias relativas  $F$ .
- 2 Se calcula la tabla estandarizada  $Z$  de frecuencias relativas que tiene las mismas dimensiones de la tabla original,  $I \times J$ , y que se obtiene dividiendo cada celda de  $F$  por la raíz de los totales de su fila y columna

$$z_{ij} = \left\{ \frac{f_{ij}}{\sqrt{r_i c_j}} \right\}$$

- 3 Se calculan los  $h$  vectores propios (normalmente  $h=2$ ) asociados a los valores propios mayores (distintos de 1), de la matriz de menor dimension de  $ZZ'$  y  $Z'Z$ .

# Análisis conjunto

- Si obtenemos los vectores propios  $\mathbf{a}_i$  de  $\mathbf{Z}'\mathbf{Z}$ , los  $\mathbf{b}_i$  de  $\mathbf{Z}\mathbf{Z}'$  se obtienen por  $\mathbf{b}_i = \mathbf{Z}\mathbf{a}_i$ .
- Análogamente si se obtienen los  $\mathbf{b}_i$  de  $\mathbf{Z}\mathbf{Z}'$ ,  $\mathbf{a}_i = \mathbf{Z}'\mathbf{b}_i$ .
- Las  $I$  filas de la matriz se representarán como  $I$  puntos en  $\mathbb{R}^h$  y las coordenadas de cada fila vienen dadas por

$$\mathbf{C}_r = \mathbf{D}_r^{-1/2} \mathbf{Z} \mathbf{A}_2$$

donde  $\mathbf{A}_2$  tiene en columnas los dos vectores propios de  $\mathbf{Z}'\mathbf{Z}$ . Las  $J$  columnas se representarán como  $J$  puntos en  $\mathbb{R}^h$  y las coordenadas de cada columna son

$$\mathbf{C}_c = \mathbf{D}_c^{-1/2} \mathbf{Z}' \mathbf{B}_2$$

donde  $\mathbf{B}_2$  tiene en columnas los dos vectores propios de  $\mathbf{Z}\mathbf{Z}'$ .

# La distribución Ji cuadrada

- La prueba de independencia entre las variables fila y columna en una tabla de contingencia  $I \times J$  se realiza con el estadístico:

$$\chi^2 = \sum \frac{(\text{fr. observadas} - \text{fr. esperadas})^2}{\text{fr. esperadas}}$$

Que bajo la hipótesis de independencia sigue, asintóticamente, una distribución  $\chi^2$  con  $(I - 1) \times (J - 1)$  grados de libertad

- El estadístico  $\chi^2$  para probar la independencia se puede escribir como

$$\chi^2 = \sum_{i=1}^I \sum_{j=1}^J \frac{(nf_{ij} - nr_i c_j)^2}{nr_i c_j}$$

O también como

$$\chi^2 = n \sum_{i=1}^I r_i \sum_{j=1}^J \left( \frac{f_{ij}}{r_i} - c_j \right)^2 \frac{1}{c_j}$$

# La distribución Ji cuadrada

- Se puede probar que esta expresión es equivalente a calcular las distancias entre los vectores de la matriz de frecuencias relativas por filas,  $\mathbf{R}$  definida como  $\mathbf{R} = \mathbf{D}_r^{-1}\mathbf{F}$  si medimos la distancia con la métrica  $\chi^2$
- También se puede probar que el valor medio o centroide de las filas de  $\mathbf{R}$  está dado por el vector cuyos componentes son las frecuencias relativas de las columnas, es decir por el vector  $\mathbf{c}$ .
- De igual forma, el centroide de las columnas de  $\mathbf{R}$  está dado por el vector  $\mathbf{r}$  de frecuencias relativas de las filas

# La distribución Ji cuadrada

- La distancia de cualquier vector fila  $\mathbf{R}_i$  a su media  $\mathbf{c}$ , con la métrica  $\chi^2$  será

$$(\mathbf{R}_i - \mathbf{c})' \mathbf{D}_c^{-1} (\mathbf{R}_i - \mathbf{c})$$

donde la matriz  $\mathbf{D}_c^{-1}$  se obtiene como en  $\mathbf{D}^2(\mathbf{R}_a, \mathbf{R}_b)$ .

- La suma de todas estas distancias, ponderadas por su importancia se conoce como **inercia total de la tabla** y se define como

$$I_T = \sum_{i=1}^I r_i (\mathbf{R}_i - \mathbf{c})' \mathbf{D}_c^{-1} (\mathbf{R}_i - \mathbf{c})$$

Lo cual se puede escribir como

$$I_T = \sum_{i=1}^I r_i \sum_{j=1}^J \left( \frac{f_{ij}}{r_i} - c_j \right)^2 / c_j$$

# La distribución Ji cuadrada

- Si se compara con

$$\chi^2 = n \sum_{i=1}^I r_i \sum_{j=1}^J \left( \frac{f_{ij}}{r_i} - c_j \right)^2 \frac{1}{c_j}$$

la **inercia total** es igual a  $\frac{\chi^2}{n}$

- También se puede demostrar que la **inercia total** es la suma de los valores propios al cuadrado de la matriz  $\mathbf{Z}'\mathbf{Z}$  eliminando el **uno**, es decir

$$I_T = \sum_{i=1}^I r_i (\mathbf{R}_i - \mathbf{c})' \mathbf{D}_c^{-1} (\mathbf{R}_i - \mathbf{c}) = \sum_{k=1}^{J-1} \lambda_k^2$$

- Por tanto, el análisis de las filas (o de las columnas debido a que el problema es simétrico) se puede ver como una descomposición de los componentes del estadístico  $\chi^2$  en sus fuentes de variación.

# La distribución Ji cuadrada

- Esto implica también que existe una inercia para cada dimensión de la representación, definida por los valores propios  $\lambda_i$
- Debido a que la inercia es una medida de la variación total de los datos de la tabla, ¿como se interpretaría un valor grande para la proporción

$$(\lambda_1^2 + \lambda_2^2) / \sum_{k=1}^{J-1} \lambda_k^2$$

?

- Geométricamente se dice que la asociación en los datos es bien representada por los puntos en un plano, y esta aproximación representa casi toda la variación en los datos más allá de lo que se podría representar con una solución de rango 1.
- Algebraicamente se puede decir que la aproximación

$$\mathbf{Z} \approx \mathbf{B}_2 \mathbf{\Lambda}_2 \mathbf{A}'_2 = \lambda_1^{1/2} \mathbf{b}_1 \mathbf{a}'_1 + \lambda_2^{1/2} \mathbf{b}_2 \mathbf{a}'_2$$

es **muy buena**