

Análisis de Correspondencias



José A. Perusquía Cortés

Análisis Multivariado Semestre 2024-I



¿De qué va?

- Una técnica multivariada para analizar las asociaciones entre un conjunto de variables categóricas de forma gráfica (reducción de la dimensión).
- Es una técnica meramente descriptiva conocida desde Hirschfeld (1935) y redescubierta e impulsada por Jean-Paul Benzécri en Francia en los años 60's.
- Técnica similar a PCA pero para datos categóricos.

- Punto en un espacio multidimensional
- Un peso (o masa) asignado a cada punto
- Un centroide
- Una función de distancia entre puntos: **chi-squared distance**
 - Para dos renglones i, i'

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

- Para dos columnas j, j'

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

Ejemplo 1 : Estado de salud

- Estado de salud por grupo de edades

Grupo\Salud	Muy Bueno	Bueno	Regular	Malo	Muy Malo	Totales Ren.
16-24	243	789	167	18	6	1223
25-34	220	809	164	35	6	1234
35-44	147	658	181	41	8	1035
45-54	90	469	236	50	16	861
55-64	53	414	306	106	30	909
65-74	44	267	284	98	20	713
75+	20	136	157	66	17	396
Totales Col.	817	3542	1495	414	103	6371

Ejemplo 1 : Estado de salud

- Tabla de frecuencias por renglón

Grupo\Salud	Muy Bueno	Bueno	Regular	Malo	Muy Malo	Peso Renglón
16-24	0.199	0.645	0.137	0.015	0.005	0.192
25-34	0.178	0.656	0.133	0.028	0.005	0.194
35-44	0.142	0.636	0.175	0.040	0.008	0.162
45-54	0.105	0.545	0.274	0.058	0.019	0.135
55-64	0.058	0.455	0.337	0.117	0.033	0.143
65-74	0.062	0.374	0.398	0.137	0.028	0.112
75+	0.051	0.343	0.396	0.167	0.043	0.062
Peso Columna	0.128	0.556	0.235	0.065	0.016	1.000

Ejemplo 1 : Estado de salud

- Puntos en un espacio multidimensional: **perfiles por renglón**

Grupo\Salud	Muy Bueno	Bueno	Regular	Malo	Muy Malo	Peso Renglón
16-24	0.199	0.645	0.137	0.015	0.005	0.192
25-34	0.178	0.656	0.133	0.028	0.005	0.194
35-44	0.142	0.636	0.175	0.040	0.008	0.162
45-54	0.105	0.545	0.274	0.058	0.019	0.135
55-64	0.058	0.455	0.337	0.117	0.033	0.143
65-74	0.062	0.374	0.398	0.137	0.028	0.112
75+	0.051	0.343	0.396	0.167	0.043	0.062
Peso Columna	0.128	0.556	0.235	0.065	0.016	1.000

Ejemplo 1 : Estado de salud

- Pesos (masas) de cada perfil: **peso renglón**

Grupo\Salud	Muy Bueno	Bueno	Regular	Malo	Muy Malo	Peso Renglón
16-24	0.199	0.645	0.137	0.015	0.005	0.192
25-34	0.178	0.656	0.133	0.028	0.005	0.194
35-44	0.142	0.636	0.175	0.040	0.008	0.162
45-54	0.105	0.545	0.274	0.058	0.019	0.135
55-64	0.058	0.455	0.337	0.117	0.033	0.143
65-74	0.062	0.374	0.398	0.137	0.028	0.112
75+	0.051	0.343	0.396	0.167	0.043	0.062
Peso Columna	0.128	0.556	0.235	0.065	0.016	1.000

Ejemplo 1 : Estado de salud

- El centroide: peso columna

Grupo\Salud	Muy Bueno	Bueno	Regular	Malo	Muy Malo	Peso Renglón
16-24	0.199	0.645	0.137	0.015	0.005	0.192
25-34	0.178	0.656	0.133	0.028	0.005	0.194
35-44	0.142	0.636	0.175	0.040	0.008	0.162
45-54	0.105	0.545	0.274	0.058	0.019	0.135
55-64	0.058	0.455	0.337	0.117	0.033	0.143
65-74	0.062	0.374	0.398	0.137	0.028	0.112
75+	0.051	0.343	0.396	0.167	0.043	0.062
Peso Columna	0.128	0.556	0.235	0.065	0.016	1.000

1. Definir matrices diagonales \mathbf{D}_r y \mathbf{D}_c con las masas por renglón y columna.

2. Obtener la descomposición GSVD de $\mathbf{R} - \mathbf{1}\mathbf{c}^T$, i.e.,

$$\mathbf{R} - \mathbf{1}\mathbf{c} = \mathbf{N}\mathbf{\Lambda}\mathbf{M}^T$$

$$\mathbf{N}^T\mathbf{D}_r\mathbf{N} = \mathbf{M}^T\mathbf{D}_c^{-1}\mathbf{M} = \mathbf{I}$$

- \mathbf{R} es la matriz de perfiles por renglón
- $\mathbf{c} = \mathbf{D}_c\mathbf{1}$ es el centroide

3. Las primeras dos coordenadas se encuentran con $\mathbf{N}_{(2)}\mathbf{\Lambda}_{(2)}$

Ejemplo 1 : Estado de salud

- Las matrices diagonales son:

$$\mathbf{D}_r = \text{diag}(.192,.194,.162,.135,.143,.112,.062)$$

$$\mathbf{D}_c = \text{diag}(.128,.556,.235,.065,.016)$$

- El centroide es:

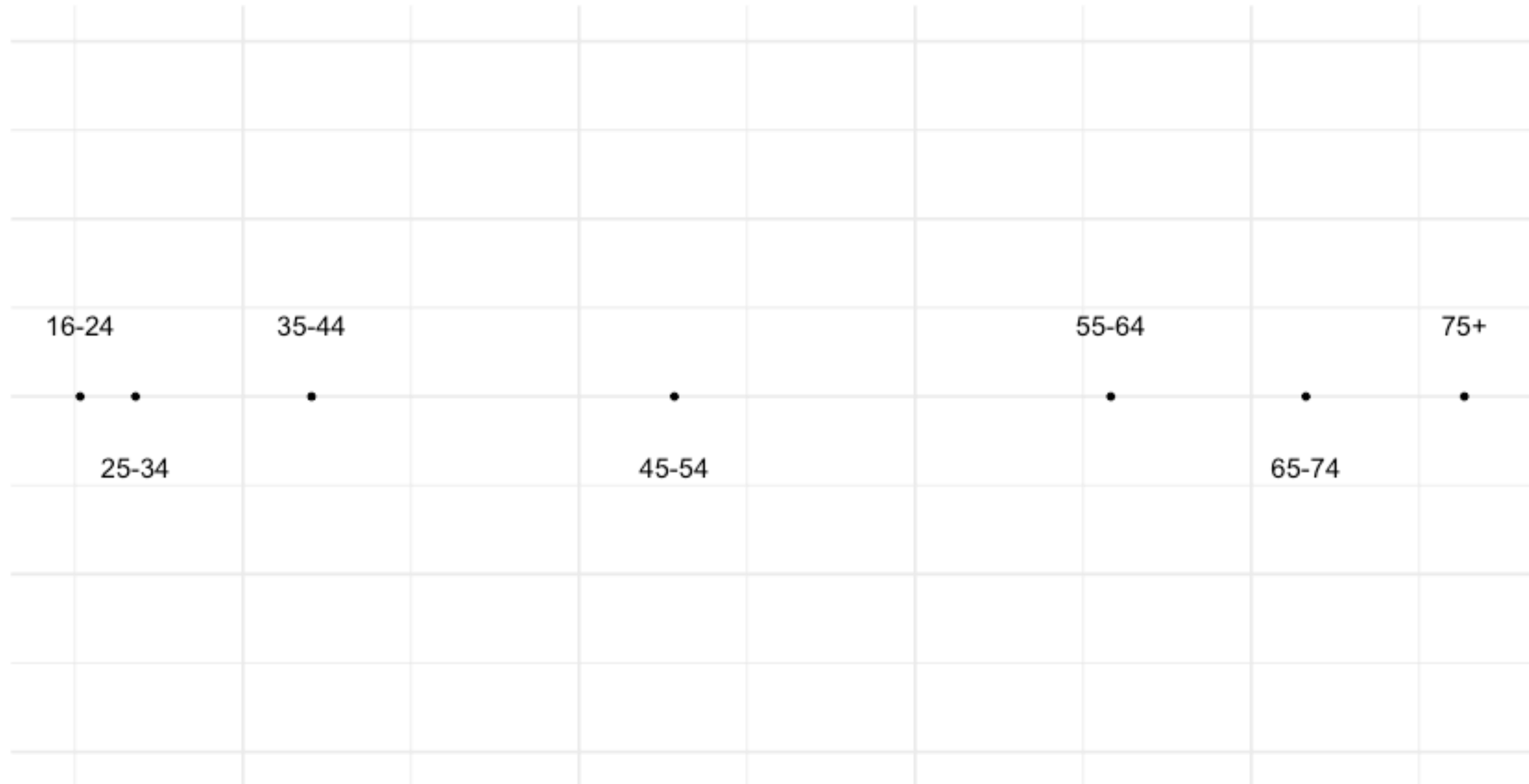
$$\mathbf{c} = (.128,.556,.235,.065,.016)^T$$

- La matriz $\mathbf{R} - \mathbf{1c}^T$ está dada por:

0.071	0.089	-0.098	-0.050	-0.011
0.050	0.100	-0.102	-0.037	-0.011
0.014	0.080	-0.065	-0.025	-0.008
-0.023	-0.011	0.039	-0.007	0.003
-0.070	-0.101	0.102	0.052	0.017
-0.066	-0.182	0.163	0.072	0.012
-0.077	-0.213	0.161	0.102	0.027

Ejemplo 1 : Estado de salud

- La proyección en la primera coordenada



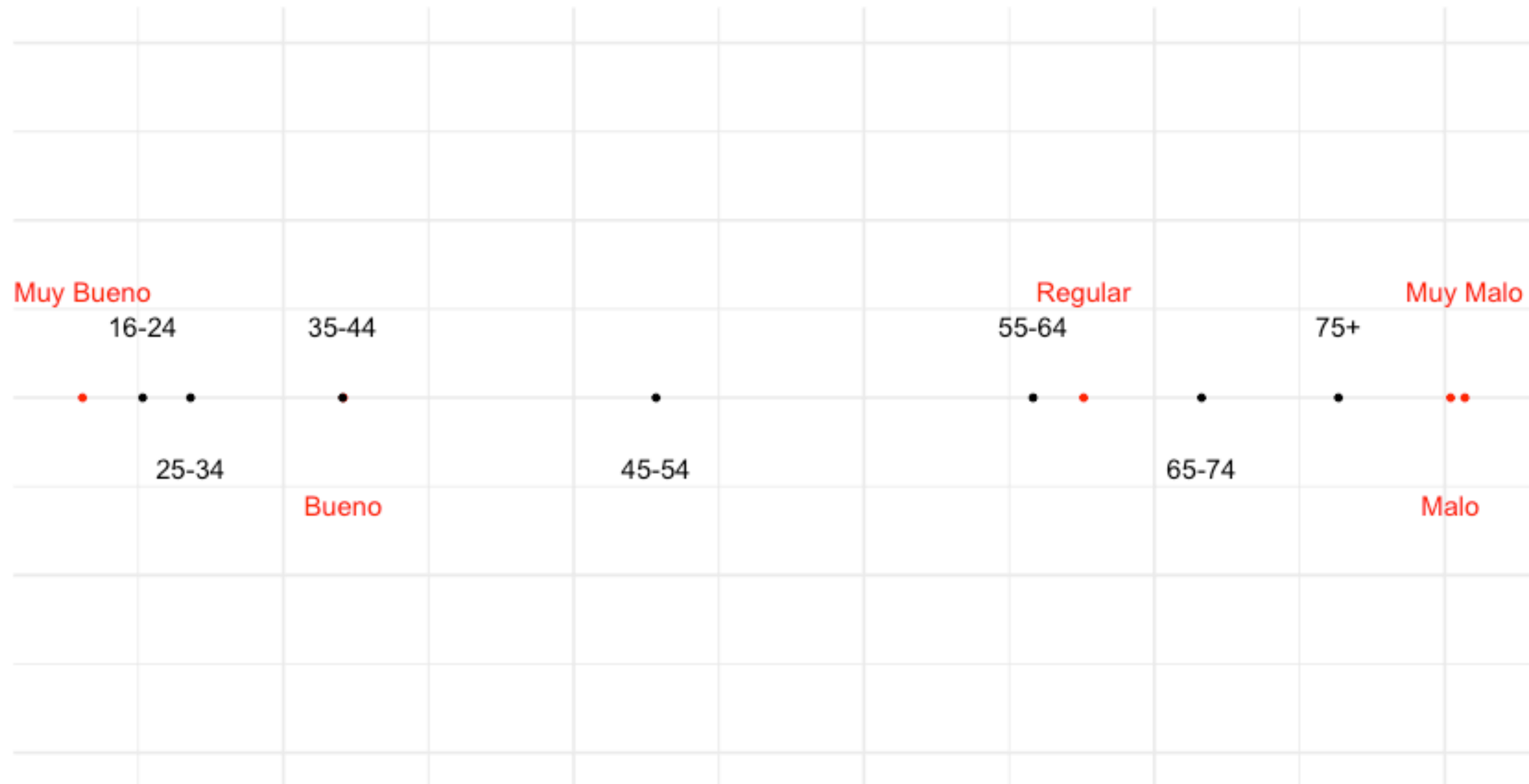
Ejemplo 1 : Estado de salud

- Usamos la transpuesta de la tabla de contingencia y repetimos.



Ejemplo 1 : Estado de salud

- Graficamos las dos variables al mismo tiempo



1. Calcular la **matriz de correspondencia** $\mathbf{P} = \frac{\mathbf{N}}{n}$
2. Definir matrices diagonales \mathbf{D}_r y \mathbf{D}_c con las sumas por renglón y columna.

3. Obtener la descomposición SVD de

$$\mathbf{D}_r^{-\frac{1}{2}} (\mathbf{P} - \mathbf{r}\mathbf{c}^T) \mathbf{D}_c^{-\frac{1}{2}} = \mathbf{U}\mathbf{\Lambda}\mathbf{V}^T$$

5. Obtener las **coordenadas estándar**

$$\mathbf{X} = \mathbf{D}_r^{-\frac{1}{2}} \mathbf{U} \qquad \mathbf{Y} = \mathbf{D}_c^{-\frac{1}{2}} \mathbf{V}$$

6. Obtener las **coordenadas principales**

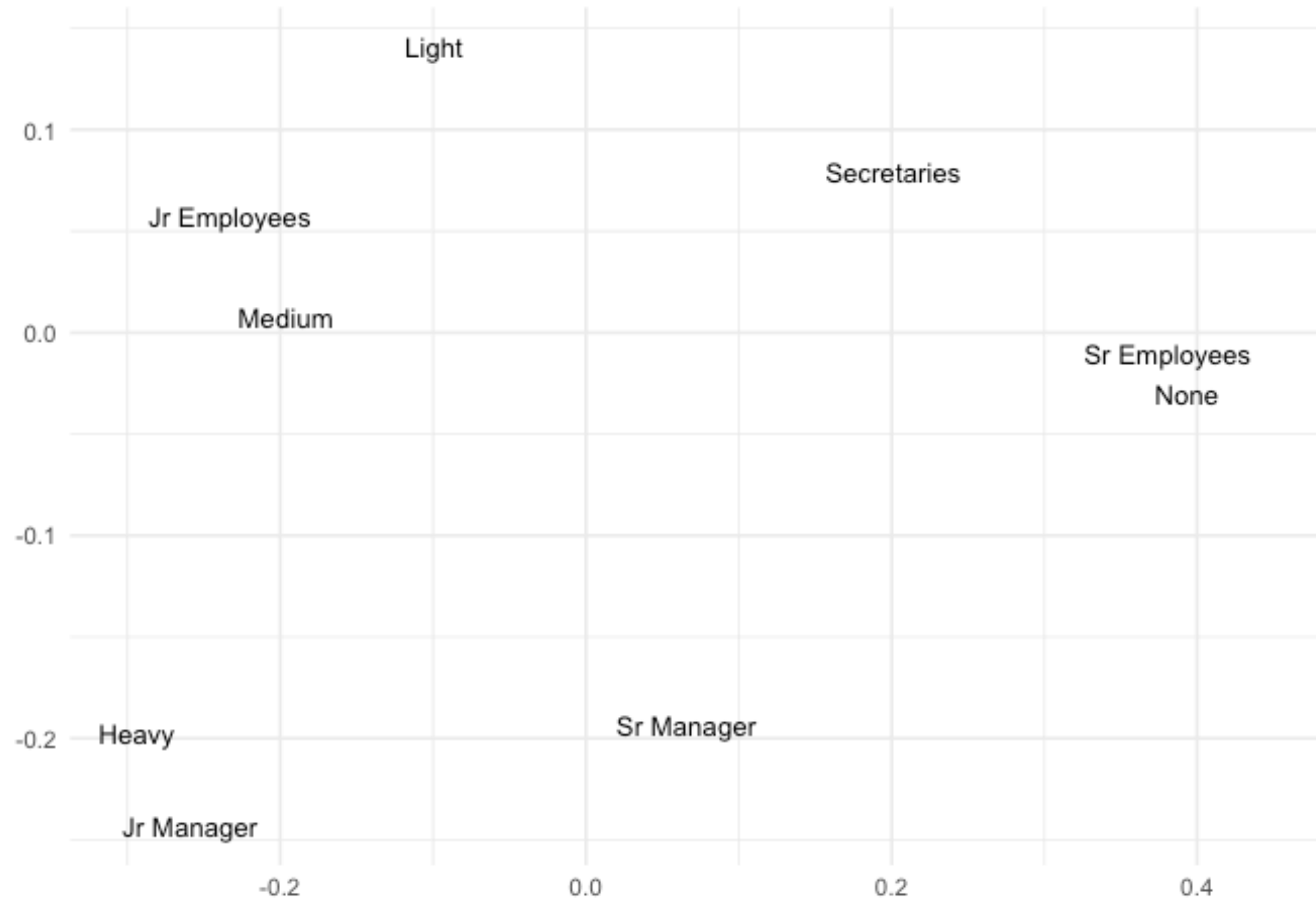
$$\mathbf{F} = \mathbf{X}\mathbf{\Lambda} \qquad \mathbf{G} = \mathbf{Y}\mathbf{\Lambda}$$

Ejemplo 2 : Trabajadores y hábitos de fumar

- Encuesta a trabajadores de una empresa sobre sus hábitos de fumar

Staff\Nivel	None	Light	Medium	Heavy	Totales Ren.
Sr Managers	4	2	3	2	11
Jr Managers	4	3	7	4	18
Sr Employees	25	10	12	4	51
Jr Employees	18	24	33	13	88
Secretaries	10	6	7	2	25
Totales Col.	61	45	62	25	193

Ejemplo 2 : Trabajadores y hábitos de fumar



- Como PCA se busca explicar la mayor cantidad de varianza definida como:

$$\text{Inercia} = \sum_{i,j} \frac{\left(p_{ij} - r_i c_j\right)^2}{(r_i c_j)}$$

- Equivalentemente, $\text{Inercia} = \frac{\chi^2}{n}$, donde χ^2 es el estadístico de Pearson y n el total de observaciones
- Los valores singulares al cuadrado $\lambda_1^2, \lambda_2^2, \dots$ son las inercias principales y explican la inercia total.

Ejemplo 2 : Trabajadores y hábitos de fumar

- La inercia total:

$$\text{Inercia} = 0.08518986$$

- Las inercias principales:

$$\lambda_1^2 = 0.07475911; \quad \lambda_2^2 = 0.01001718; \quad \lambda_3^2 = 0.0004135741$$

- Porcentaje explicado acumulado:

$$0.8775588 \quad 0.9951453 \quad 1$$

- Las coordenadas de los renglones **F** y la de las columnas **G** están relacionadas

$$\mathbf{F} = \mathbf{RGA}^{-1} \qquad \mathbf{G} = \mathbf{CFA}^{-1}$$

- Nos da una forma de añadir perfiles suplementarios para columnas y renglones

Ejemplo 2 : Trabajadores y hábitos de fumar

- Imaginar que se tiene un promedio nacional de fumadores

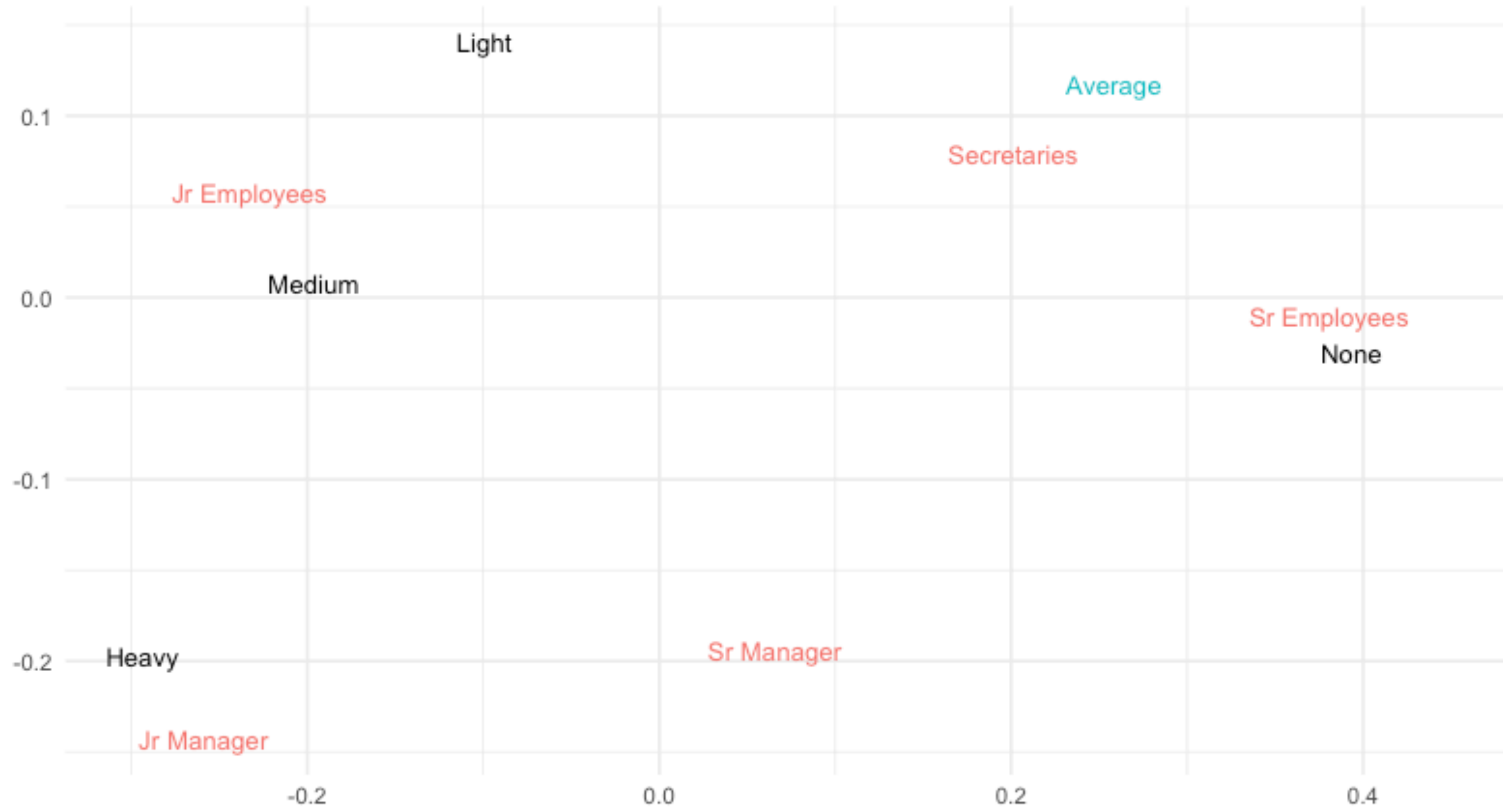
Staff\Nivel	None	Light	Medium	Heavy
Sr Managers	4	2	3	2
Jr Managers	4	3	7	4
Sr Employees	25	10	12	4
Jr Employees	18	24	33	13
Secretaries	10	6	7	2
Promedio	42%	29%	20%	9%

- Encontramos su representación como:

$$f_{11}^* = \frac{(.42 * -0.39330845) + (.29 * 0.09945592) + (.2 * 0.19632096) + (.09 * 0.29377599)}{.2734211} = .258$$

$$f_{12}^* = \frac{(.42 * -0.030492071) + (.29 * 0.141064289) + (.2 * 0.007359109) + (.09 * -0.197765656)}{0.1000859} = .118$$

Ejemplo 2 : Trabajadores y hábitos de fumar



Ejemplo 2 : Trabajadores y hábitos de fumar

- De forma similar podemos añadir columnas

Staff\Nivel	None	Light	Medium	Heavy	Drinking	Not Drinking
Sr Managers	4	2	3	2	0	11
Jr Managers	4	3	7	4	1	17
Sr Employees	25	10	12	4	5	46
Jr Employees	18	24	33	13	10	78
Secretaries	10	6	7	2	7	18
Promedio	42%	29%	20%	9%		

- **Observación:** Ya **no** es una tabla de contingencia

Ejemplo 2 : Trabajadores y hábitos de fumar



Análisis de correspondencias múltiple

- Extensión del análisis de correspondencias simple que permite analizar la relación de K variables categóricas dependientes, cada una con J_k niveles tales que $\sum_k J_k = J$
- Es importante que las variables sean “homogéneas”, e.g. no mezclar variables de opinión con variables demográficas.
- A grandes rasgos se puede ver como un análisis de correspondencias simple en una matriz indicadora.

- Las variables pueden ser cuantitativas también, siempre que se agrupen.
- El análisis se puede hacer en la matriz indicadora \mathbf{X} o en la matriz de Burt $\mathbf{X}^T \mathbf{X}$ (puede ser computacionalmente más sencillo)
- Es necesario re-escalar los eigenvalores, e.g. Greenacre (1993) propuso:

$$\lambda_i^c = \begin{cases} \left[\left(\frac{K}{K-1} \right) \left(\lambda_i - \frac{1}{K} \right) \right]^2 & \text{si } \lambda_i > \frac{1}{K} \\ 0 & \text{si } \lambda_i \leq \frac{1}{K} \end{cases}$$

- El porcentaje de inercia se puede calcular como

$$\frac{\lambda_i^c}{\sum_i \lambda_i^c}$$

- Greenacre propuso en su lugar estimarla a través de

$$\frac{\lambda_i^c}{\bar{\mathcal{J}}}$$

Donde

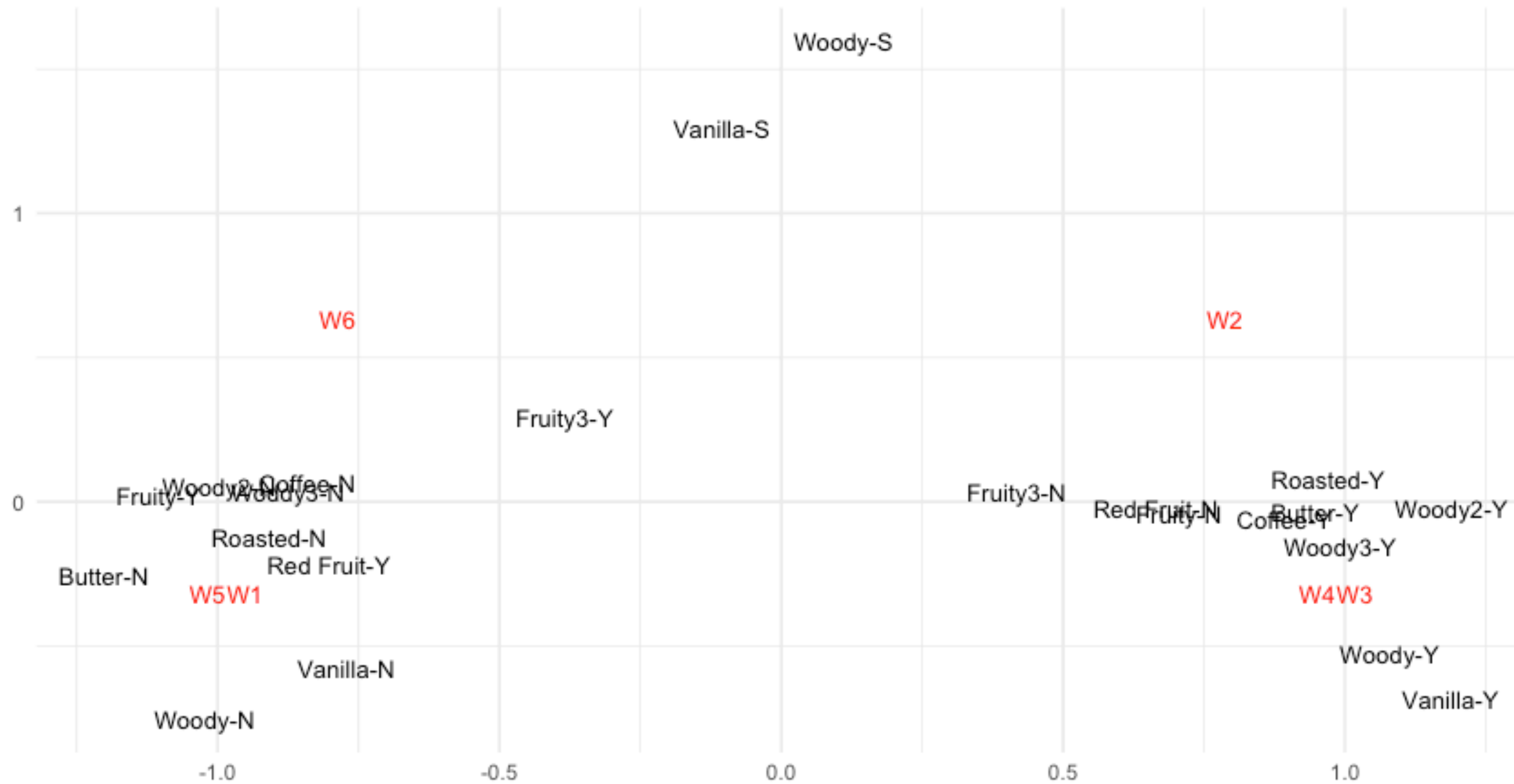
$$\bar{\mathcal{J}} = \frac{K}{K-1} \left(\sum_i \lambda_i^2 - \frac{J-K}{K^2} \right)$$

Ejemplo 3: Características de vinos

		Expert 1							Expert 2							Expert 3							
Wine	Oak								red														
	Type	fruity		woody			coffee		fruit		roasted		vanillin			woody		fruity		butter		woody	
W1	1	1	0	0	0	1	0	1	1	0	0	1	0	0	1	0	1	0	1	0	1	0	1
W2	2	0	1	0	1	0	1	0	0	1	1	0	0	1	0	1	0	0	1	1	0	1	0
W3	2	0	1	1	0	0	1	0	0	1	1	0	1	0	0	1	0	0	1	1	0	1	0
W4	2	0	1	1	0	0	1	0	0	1	1	0	1	0	0	1	0	1	0	1	0	1	0
W5	1	1	0	0	0	1	0	1	1	0	0	1	0	0	1	0	1	1	0	0	1	0	1
W6	1	1	0	0	1	0	0	1	1	0	0	1	0	1	0	0	1	1	0	0	1	0	1
W?	?	0	1	0	1	0	.5	.5	1	0	1	0	0	1	0	.5	.5	1	0	.5	.5	0	1

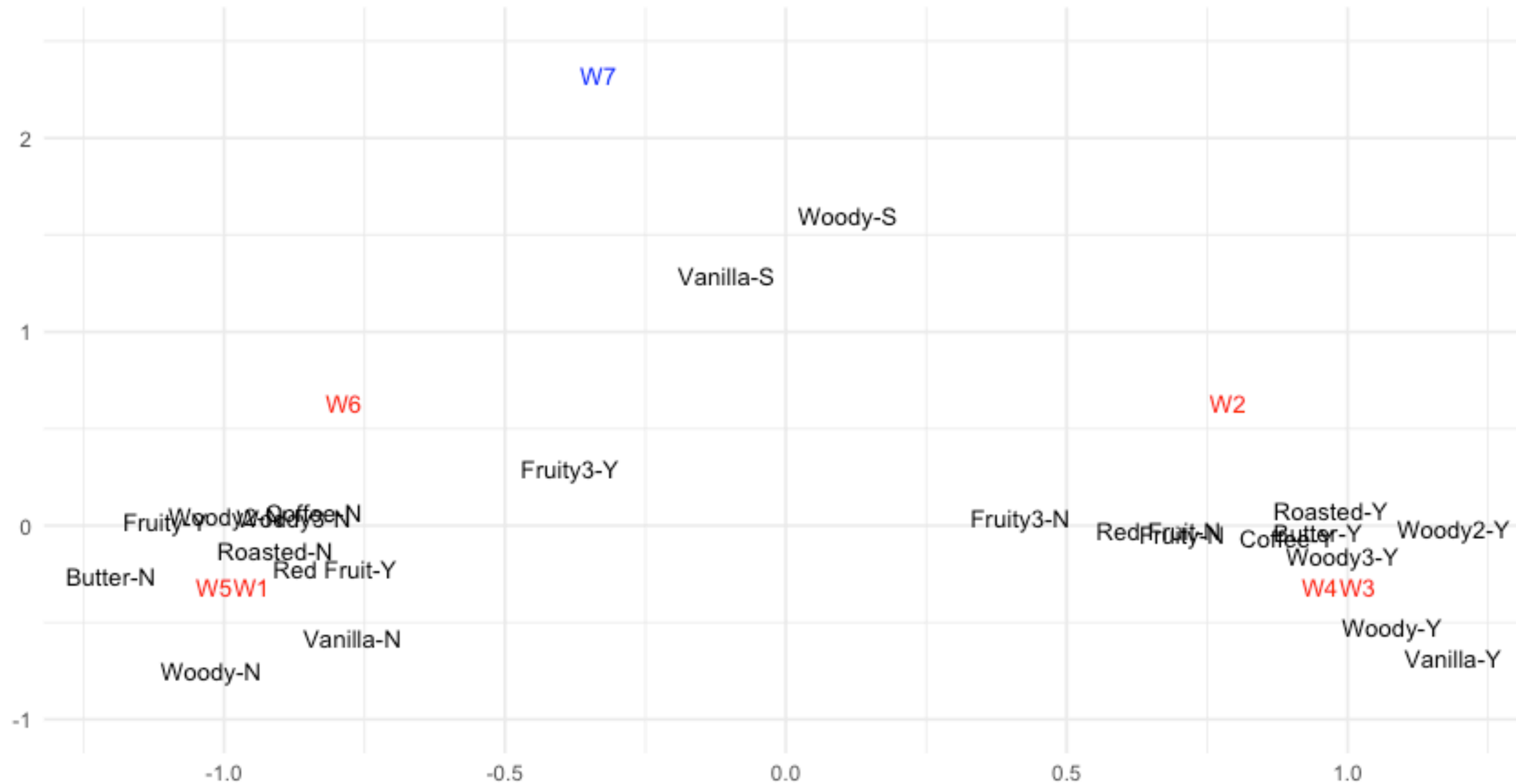
Ejemplo 3: Características de vinos

- Coordenadas (modificadas)



Ejemplo 3: Características de vinos

- Añadiendo el vino del que no se conoce el tipo de barrica



Ejemplo 3: Características de vinos

- Inercias principales

$$\lambda_1^2 = 0.8532; \quad \lambda_2^2 = 0.2; \quad \lambda_3^2 = .1151; \quad \lambda_4^2 = 0.0317$$

- Contribución a la inercia total

$$0.7110 \quad 0.8776 \quad .9736 \quad 1$$

Ejemplo 3: Características de vinos

- Inercias corregidas

$$\lambda_1^2 = 0.7004; \quad \lambda_2^2 = 0.0123; \quad \lambda_3^2 = .0003; \quad \lambda_4^2 = 0$$

- Contribución a la inercia total

$$0.9519 \quad 0.9687 \quad .9691 \quad .9691$$