

ANÁLISIS DE TABLAS DE DISIMILITUDES REPRESENTACIÓN GEOMÉTRICA DE LA POBLACIÓN

Gérard Defives

A.D.E.M.

Grupo de Análisis de Datos y
Estadísticas Multidimensional
Escuela Básica de Ingeniería
Universidad de Los Andes

RESUMEN. Dada una población, cuyos individuos no se pueden describir mediante medidas cuantitativas, queda posible definir similitudes o disimilitudes entre los objetos a comparar, y luego se pueden obtener árboles de clasificación jerárquica, árbol máximo (o mínimo), grafo G3, filtrante, etc. En este trabajo se propone una representación geométrica de esta población. Ésta se basa en la construcción de un producto escalar entre los objetos del estudio, lo cual, mediante la diagonalización de la matriz simétrica, semi-definida positiva, de dicho producto, desemboca sobre una representación de la población parecida a las representaciones que se obtienen en Análisis de Componentes Principales.

0 INTRODUCCIÓN

En muchas situaciones, se desea comparar objetos mediante el cálculo de índices de similitud o de disimilitud, lo cual conduce generalmente a aplicar métodos de clasificación jerárquica (Benzécri, 1973 y Roux, 1985) o derivados de la teoría de grafos (G3 y filtrantes), (Degenne y Vergès, 1973). Varios métodos han sido propuestos para obtener una traducción geométrica de

una tabla de similitudes o disimilitudes. Cuando los objetos a comparar son modalidades de variables cualitativas, el análisis de Correspondencias puede aplicarse (Benzécri, 1973 y Di Giacomo, 1981). Nótese que en el ejemplo tratado por Di Giacomo, el índice de similitud no es otra cosa que la comparación de los perfiles en un tabla de contingencia. En este caso resulta fácil relacionar Análisis de Correspondencias y Clasificación. No es el caso si sólo se dispone de las similitudes o disimilitudes.

Un método consiste en considerar la matriz de similitudes \mathbf{S} ($n \times n$) como la propia matriz del producto escalar entre los objetos. Este método lo expone Escoufier (Escoufier, 1979). El paso de disimilitudes a similitudes no presenta dificultades. También se puede operar en base a la fórmula de W.S. Torgerson (Rao, 1975). Estos métodos se justifican principalmente cuando los objetos a comparar presentan el carácter de variables. En efecto, los términos de la diagonal de \mathbf{S} son todos la similitud máxima s_m . Dividir por s_m , conduce a una matriz análoga a la de las variables en Análisis de Componentes Principales (ACP).

En lo que sigue, se propone un método mejor adaptado a la representación de los individuos. En un artículo posterior, se describirá, de nuevo, este método propuesto y se ilustrará con un ejemplo.

1 ANÁLISIS DE UNA TABLA DE DISTANCIAS

Sean $\{\mathbf{A}_i\}_{i=1\dots n}$, n puntos ponderados por pesos respectivos p_i , tales que $\sum_i p_i = 1$, para los cuales sólo se conocen las distancias entre puntos $d(\mathbf{A}_i, \mathbf{A}_j) = d(i, j)$, distancia euclídea. Suponiendo la nube centrada en 0 (cero), se quiere reconstruir la nube de puntos en un espacio q -dimensional.

Al igual que en el ACP, se procede a partir de la matriz del producto escalar correspondiente a la distancia \mathbf{d} . Al producto escalar está asociada la norma $\|\mathbf{X}\|^2 = (\mathbf{X}|\mathbf{X})$, y la distancia se define por $d(A,B) = \|\mathbf{OA} - \mathbf{OB}\|$.

Como

$$\begin{aligned}\|\mathbf{X} - \mathbf{Y}\|^2 &= \|\mathbf{X}\|^2 + \|\mathbf{Y}\|^2 - 2(\mathbf{X}|\mathbf{Y}) \\ (\mathbf{X}|\mathbf{Y}) &= \frac{1}{2}[\|\mathbf{X}\|^2 + \|\mathbf{Y}\|^2 - \|\mathbf{X} - \mathbf{Y}\|^2]\end{aligned}$$

resulta que los productos escalares entre individuos se expresan en función de las distancias por:

$$(\mathbf{A}_i|\mathbf{A}_j) = \frac{1}{2}[\mathbf{d}^2(\mathbf{o},i) + \mathbf{d}^2(\mathbf{o},j) - \mathbf{d}^2(i,j)] = \mathbf{W}_{ij}$$

Pero en la situación presente, se conocen los $\mathbf{d}(i,j)$ distancias entre los individuos \mathbf{A}_i y \mathbf{A}_j pero no las distancias $\mathbf{d}(\mathbf{o},i)$ y $\mathbf{d}(\mathbf{o},j)$, distancias entre estos individuos y el centro de gravedad $\mathbf{0}$ (cero).

Para obtener estas distancias $\mathbf{d}(\mathbf{o},i)$, se puede proceder de la siguiente manera:

Sea \mathbf{I}_i la inercia de la nube de puntos con respecto al punto \mathbf{A}_i .

$$\mathbf{I}_i = \sum_{j=1}^n p_j \mathbf{d}^2(i,j) = \mathbf{d}^2(\mathbf{o},i) + \mathbf{I} \quad (\text{teorema de Huygens}) \quad (1)$$

donde \mathbf{I} representa la inercia de la nube de puntos con respecto a su centro de gravedad $\mathbf{0}$ (cero).

$$\sum_{i=1}^n p_i \mathbf{I}_i = \sum_{i=1}^n p_i \mathbf{d}^2(\mathbf{o},i) + \sum_{i=1}^n p_i \mathbf{I}$$

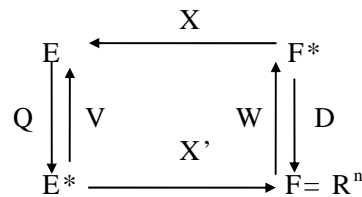
pero $I = \sum_{i=1}^n p_i d^2(o, i)$ y $\sum_{i=1}^n p_i = 1$, luego $\sum_{i=1}^n p_i I_i = 2I$

$$I = \frac{1}{2} \sum_{i=1}^n p_i I_i \quad (2)$$

se introduce en (1):

$$d^2(o, i) = I_i - I = I_i - \frac{1}{2} \sum_{j=1}^n p_j I_j$$

Una vez calculados los $d(o, i)$, se puede calcular la tabla de los productos escalares $W = [w_{ij}]$ y hacer jugar a esta matriz W el papel de la matriz $W = X'QX$ del esquema de dualidad:



Si d es una distancia euclídea, la matriz W obtenida es simétrica, semidefinida positiva, y verifica, si D denota la matriz $(n \times n)$ diagonal de los pesos, y $\underline{1}$ el vector $(n \times 1)$ cuyos términos son todos 1 (unos) : $WD\underline{1} = 0$.

En efecto, la simetría resulta de la simetría de d : $d(i, j) = d(j, i)$, y W sería definida positiva si no fuera por el centrado de la nube de puntos, lo que trae como consecuencia $WD\underline{1} = 0$.

$$2 \sum_{i=1}^n w_{ij} p_i = \sum_{i=1}^n [p_i d^2(o, i) + p_i d^2(o, j) - p_i d^2(i, j)] = I + d^2(o, j) - I_j = 0$$

En el caso presente, no se dispone de la matriz X_r que describe los individuos en un espacio vectorial $E = R^r$, ni X' que describe las variables en F . Sin embargo, se puede

obtener un mapa que represente la población $\{A_i\}$ y que distorsione, lo menos posible, las distancias entre los individuos A_i .

Se trata entonces de hallar la matriz Y ($q \times n$), cuyas columnas representarán a los individuos en un espacio vectorial de dimensión q , referido a una base ortogonal, y tal que los productos escalares calculados entre estas representaciones ($Y'Y$) sean lo más cercano posible a la matriz W . Tomando en cuenta los pesos atribuidos a los individuos, esto consiste en maximizar:

$$RV(WD, Y'YD) = \frac{\text{Tr}(WDY'YD)}{\sqrt{\text{Tr}(WD)^2 \text{Tr}(Y'YD)^2}}$$

y la solución es la matriz Y , cuyas filas son los vectores propios de WD asociados a los q mayores valores propios (Escoufier, 1979). Se normalizarán estos vectores propios de tal forma que $YDY' = \Lambda$, la matriz diagonal ($q \times q$) de los q mayores valores propios de WD para que la inercia de la proyección de la nube corresponda a la inercia de la nube. (Y cuando $q = p =$ dimensión del más pequeño espacio vectorial que contiene la nube, la inercia inicial está totalmente reconstituida. Esta dimensión p es el número de valores propios de WD diferentes de cero (Defives, 1983).

2 ANÁLISIS DE UNA TABLA DE DISIMILITUDES

En la práctica, el proceso arrancará de una disimilitud, que no es una distancia. Si esta disimilitud no verifica la desigualdad del triángulo, puede resultar que algunos $d^2(o,i)$ sean negativos, o que la matriz W no sea positiva.

Por ejemplo, para tres puntos a,b,c, cuyas disimilitudes están dadas por la tabla siguiente:

d	a	b	c
a	0	1	2
b		0	4
c			0

provistos de pesos idénticos (1/3), se tienen las inercias:

$I_a = 5/3$; $I_b = 17/3$; $I_c = 20/3$, y la inercia con respecto al origen $\mathbf{0}$, colocado en el centro de gravedad:

$$I = \frac{1}{6}(I_a + I_b + I_c) = 7/3$$

Luego: $d^2(\mathbf{o},\mathbf{a}) = -2/3$; $d^2(\mathbf{o},\mathbf{b}) = 10/3$; $d^2(\mathbf{o},\mathbf{c}) = 13/3$.

El paliativo a esta situación consiste en cambiar la disimilitud \mathbf{d} por otro índice de disimilitud δ , mediante la agregación a todos los $\mathbf{d}(\mathbf{i},\mathbf{j})$, fuera de la diagonal, de una misma longitud \mathbf{l} , de tal forma que cualquiera sean $\mathbf{i}, \mathbf{j}, \mathbf{k}$

$$\mathbf{d}(\mathbf{i}, \mathbf{k}) + \mathbf{l} \leq \mathbf{d}(\mathbf{i}, \mathbf{j}) + \mathbf{l} + \mathbf{d}(\mathbf{j}, \mathbf{k}) + \mathbf{l}$$

o sea: $\mathbf{d}(\mathbf{i}, \mathbf{k}) \leq \mathbf{d}(\mathbf{i}, \mathbf{j}) + \mathbf{d}(\mathbf{j}, \mathbf{k}) + \mathbf{l}$

El uso de la disimilitud $\delta(\mathbf{i},\mathbf{j}) = \mathbf{d}(\mathbf{i},\mathbf{j}) + \mathbf{l}$ si $\mathbf{i} \neq \mathbf{j}$, $\delta(\mathbf{i},\mathbf{i}) = \mathbf{0}$, conserva la misma ordenanza, o sea: las clasificaciones resultantes de \mathbf{d} y δ son idénticas. Sólo los nudos están trasladados de una longitud \mathbf{l} . Pero δ es una distancia, luego la matriz $\hat{\mathbf{W}} = [\hat{\mathbf{W}}_{ij}]$, de los productos escalares asociados a δ es simétrica, semi-definida positiva, y $\hat{\mathbf{W}}\mathbf{D}\mathbf{1} = \mathbf{0}$.

Se puede observar que de no verificarse la desigualdad triangular, la matriz W no resulta definida positiva. En el caso del ejemplo, usando

$W_{ij} = \frac{1}{2}[\mathbf{d}^2(\mathbf{o}, \mathbf{i}) + \mathbf{d}^2(\mathbf{o}, \mathbf{j}) - \mathbf{d}^2(\mathbf{i}, \mathbf{j})]$, se obtiene:

$$W = \frac{1}{6} \begin{bmatrix} -4 & 5 & -1 \\ 5 & 20 & -25 \\ -1 & -25 & 26 \end{bmatrix}$$

Si se agrega $\mathbf{l} = \mathbf{1}$ a todas las disimilitudes $\mathbf{d}(\mathbf{i}, \mathbf{j})$, se obtiene la distancia δ :

δ	a	b	c
a	0	2	3
b		0	5
c			0

Ahora, $\mathbf{I}_a = 13/3$; $\mathbf{I}_b = 29/3$; $\mathbf{I}_c = 43/3$; $\mathbf{I} = 38/9$

$\delta^2(\mathbf{o}, \mathbf{a}) = 1/9$; $\delta^2(\mathbf{o}, \mathbf{b}) = 49/9$; $\delta^2(\mathbf{o}, \mathbf{c}) = 64/9$

y el producto escalar está dado por la matriz:

$$\hat{W} = \frac{1}{6} \begin{bmatrix} 1 & 7 & -8 \\ 7 & 49 & -56 \\ -8 & -56 & 64 \end{bmatrix}$$

\hat{W} es semi-definida positiva; $\hat{W}\mathbf{D}\mathbf{1} = \mathbf{0}$.

Se verifica también que $\text{tr } \hat{W}\mathbf{D} = \mathbf{I} = 38/9$.

Esto es general, ya que se ha construido \hat{W} sobre δ que es una verdadera distancia.

Del punto de vista de las disimilitudes, agregar L a d para obtener δ , es perfectamente legítimo. Esto equivale a no tomar en consideración un elemento común a todos los individuos a comparar. La dificultad radica en hallar el más pequeño l tal que $\delta = d + l$ sea una distancia.

Para todo $i = 1, \dots, n$, y todo $k = i + 1, \dots, n$, y todo $j = 1, \dots, n$, se calcula $l(i, j, k) = d(i, k) - d(i, j) - d(j, k)$ y se toma el mayor de los $l(i, j, k)$ como l .

Otra estrategia consiste en despreocuparse de hallar una distancia, y buscar directamente la matriz del producto escalar. Para esto, se diagonaliza la matriz W , calculada con la disimilitud d . Sea k el más pequeño valor propio (negativo) de W , la matriz $M = W - kI$ es una matriz simétrica semidefinida positiva, porque si X es vector propio de W para algún valor propio h : $MX = WX - kX = (h - k)X$, luego X es vector propio de M para el valor propio $(h - k)$. Siendo k el más pequeño de los valores propios de w , $h - k \geq 0$.

Si m_{ij} denota el término de la i -ésima fila y j -ésima columna de M , la distancia entre individuos, inducida por la matriz M es:

$$\delta^2(i, j) = m_{ii} + m_{jj} - k - 2w_{ij}$$

$$\delta^2(i, j) = w_{ii} - k + w_{jj} - k - 2w_{ij} = d^2(i, j) - 2k$$

$\delta(i, j)$ es una distancia que produce la misma ordenanza que la disimilitud d porque si:

$$d(i, j) < d(i', j')$$

$$d^2(i,j) < d^2(i',j')$$

$$d^2(i,j) - 2k < d^2(i',j') - 2k$$

es decir: $\delta^2(i,j) < \delta^2(i',j')$

$$\delta^2(i,j) < \delta^2(i',j') \text{ ya que } \delta > 0$$

Pero M no verifica la relación $MD\underline{1} = 0$ ya que $WD\underline{1} - kD\underline{1} = -kD\underline{1}$.

Esto se puede corregir fácilmente, porque si M es semi-definida positiva, la matriz $\hat{M} = (I - \underline{1}\underline{1}'D)M(I - D\underline{1}\underline{1}')$ lo es también. Además $\hat{M}D\underline{1} = 0$ y

$$\hat{m}_{ii} + \hat{m}_{jj} - 2\hat{m}_{ij} = m_{ii} + m_{jj} - 2m_{ij}$$

luego, al usar \hat{M} en lugar de M , resultan las mismas distancias δ .

Demostración:

$$M = (I - \underline{1}\underline{1}'D)M(I - D\underline{1}\underline{1}') = M - \underline{1}\underline{1}'DM - MD\underline{1}\underline{1}' + \underline{1}\underline{1}'DM D\underline{1}\underline{1}'$$

$$\hat{m}_{ij} = m_{ij} - \sum_{k=1}^n p_k m_{kj} - \sum_{k=1}^n p_k m_{ik} + \sum_{k=1}^n \sum_{h=1}^n p_k p_h m_{kh}$$

de ahí sale: $\hat{m}_{ii} + \hat{m}_{jj} - 2\hat{m}_{ij} = m_{ii} + m_{jj} - 2m_{ij}$

Sustituyendo y observando que

$$\sum_{k=1}^n p_k m_{ki} = \sum_{k=1}^n p_k m_{ik}$$

por la simetría de \mathbf{M} :

$$\sum_{j=1}^n \mathbf{p}_j \hat{\mathbf{m}}_{ij} = \sum_{j=1}^n \mathbf{p}_j \mathbf{m}_{ij} - \sum_{j=1}^n \sum_{k=1}^n \mathbf{p}_j \mathbf{p}_k \mathbf{m}_{kj} - \sum_{j=1}^n \mathbf{p}_j \sum_{k=1}^n \mathbf{p}_k \mathbf{m}_{ik} + \sum_{k=1}^n \sum_{h=1}^n \mathbf{p}_k \mathbf{p}_h \mathbf{m}_{kh} = \mathbf{0}$$

lo que comprueba que $\hat{\mathbf{M}} \mathbf{D} \mathbf{1} = \mathbf{0}$.

Por otra parte, \mathbf{M} siendo simétrica, $\hat{\mathbf{M}}$ también porque $(\mathbf{I} - \mathbf{1} \mathbf{1}' \mathbf{D})' = \mathbf{I} - \mathbf{D} \mathbf{1} \mathbf{1}'$ y $\mathbf{M} = \mathbf{P} \Delta \mathbf{P}^{-1} = \mathbf{P} \Delta \mathbf{P}'$ conduce a $\hat{\mathbf{M}} = (\mathbf{I} - \mathbf{1} \mathbf{1}' \mathbf{D}) \mathbf{P} \Delta \mathbf{P}' (\mathbf{I} - \mathbf{D} \mathbf{1} \mathbf{1}')$, o sea

$$\mathbf{M} = \mathbf{Q} \Delta \mathbf{Q}' \text{ con } \mathbf{Q} = (\mathbf{I} - \mathbf{1} \mathbf{1}' \mathbf{D}) \mathbf{P} \text{ y } \mathbf{Q}^{-1} = \mathbf{Q}'.$$

\mathbf{M} y $\hat{\mathbf{M}}$ tienen los mismos valores propios, luego si \mathbf{M} es semi-definida positiva, $\hat{\mathbf{M}}$ también lo es.

El primer método modifica las disimilitudes \mathbf{d} agregándoles $\mathbf{1}$ para obtener la distancia δ , y el segundo método equivale a modificar \mathbf{d}^2 restándole $2 \mathbf{k}$, (\mathbf{k} siendo el más pequeño valor propio de \mathbf{W}), para obtener δ^2 .

Ambos métodos desembocan sobre una matriz $\hat{\mathbf{W}}$ o $\hat{\mathbf{M}}$, simétricas, semi-definidas positivas, que verifican $\hat{\mathbf{W}} \mathbf{D} \mathbf{1} = \mathbf{0}$, $\hat{\mathbf{M}} \mathbf{D} \mathbf{1} = \mathbf{0}$, y pueden jugar el papel de la matriz de productos escalares entre individuos, cuya diagonalización conduce a la matriz \mathbf{Y} , que representa los individuos en un espacio vectorial referido a una base ortogonal.

¿A cuál método dar la preferencia?

En el ejemplo, la disimilitud \mathbf{d} conducía a una matriz \mathbf{W} , no positiva, cuyos valores propios son 8,34 ; 0 ; -1.34. Luego la matriz \mathbf{M} es:

$$M = \frac{1}{6} \begin{bmatrix} 4,05 & 5 & -1 \\ 5 & 28,05 & -25 \\ -1 & -25 & 34,05 \end{bmatrix}$$

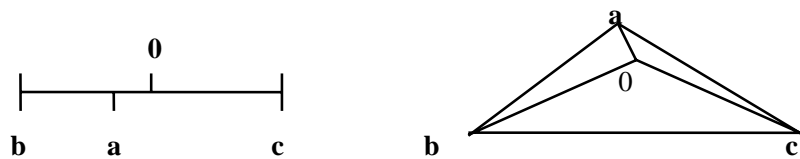
Resultan las distancias:

δ^2	a	b	c
a	0	3,68	6,68
b		0	18,68
c			0

δ	a	b	c
a	0	1,91	2,59
b		0	4,32
c			0

En este caso, las distancias δ obtenidas por el segundo método son más cercanas a las disimilitudes originales d que por el primer método. No se puede afirmar que será general. Restar el doble del menor valor propio de W de los cuadrados de las disimilitudes d , afecta menos a las disimilitudes altas que a las bajas, y esto es difícil de justificar. Sin embargo, como no afecta la ordenanza, no tiene mucha importancia visto lo que la selección de un índice de disimilitud contiene ya de arbitrario.

Las configuraciones geométricas resultantes de ambos métodos se dan a continuación:



3 CONCLUSIÓN

A partir de una tabla de disimilitudes $\mathbf{d}(\mathbf{i},\mathbf{j})$ se pueden calcular las $\mathbf{d}(\mathbf{o},\mathbf{j})$ que faltan para construir la matriz \mathbf{W} cuyos términos son:

$$w_{ij} = \frac{1}{2} [d^2(\mathbf{o},\mathbf{i}) + d^2(\mathbf{o},\mathbf{j}) - d^2(\mathbf{i},\mathbf{j})]$$

Si \mathbf{d} resulta ser una distancia euclídea, \mathbf{W} es un producto escalar. De lo contrario, la construcción de la distancia δ por la agregación a \mathbf{d} de la menor longitud \mathbf{L} que permita verificar la desigualdad triangular, o restar de la diagonal de \mathbf{W} el menor de sus valores propios negativos, permite definir sobre la población un producto escalar cuya matriz es $\hat{\mathbf{W}}$ en el primer caso, $\hat{\mathbf{M}}$ en el segundo.

Los vectores propios \mathbf{Y} de $\hat{\mathbf{W}}\mathbf{D}$ o de $\hat{\mathbf{M}}\mathbf{D}$ según la táctica escogida, normados de tal forma que $\mathbf{Y}\mathbf{Y}' = \mathbf{\Lambda}$ (la diagonal de los \mathbf{q} mayores valores propios) dan las \mathbf{q} coordenadas de cada individuo en una representación \mathbf{q} -dimensional de la población. Los criterios de apreciación de la calidad de la representación serán los mismos que en el ACP, en particular podrá considerarse la fracción de la inercia total reconstituida.

BIBLIOGRAFÍA

- Benzécri J.P. (1973): **L'Analyse des données**. París, Dunod.
- Defives G. (1983): "Inmersión de un conjunto de puntos en un espacio vectorial". **Sistemas Mérida**: ULA. No. 2. 6-8.

- Degenne A. y Vergès P. (1973): "Introduction á l'analyse des similitudes". **Revue Française de Sociologie**. XIV. 471-512.
- Di Giacomo J.P. (1981): "Aspects méthodologiques de l'analyse des représentations sociales". **Cahiers de Psychologie Cognitive**. 1-397-422.
- Escoufier Y. (1979): **Cours d'Analyse des données**. Cap. 2. U.S.T.L. Montpellier.
- Rao C.R. (1965): "The use and interpretation of P.C.A. in applied research". **Sankhya The Indian Journal of Statistics**, A 26. 329-358.
- Roux M. (1985): **Algorithmes de classification**. París: Masson.

