

Ejemplo: Datos de Arrestos en los E.U.A. (1973)
Estadísticas (por cada 100,000 habitantes)
de asaltos, asesinatos y violaciones.
La cuarta columna es el porcentaje de la
población que vive en áreas urbanas.

La función "Usarrest.R" estandariza los datos, calcula una matriz de distancias ID y con esta obtiene un agrupamiento el cual depende del método jerárquico aglomerativo que se le indique a la función de R "hclust()", posteriormente se le puede aplicar la función "viz_dend()" ⁽¹⁾, al objeto que contiene el agrupamiento, con el fin de producir un dendograma. La Figura en el archivo "Figure D1.pdf" contiene dendogramas producidos con los métodos jerárquicos aglomerativos: "average Linkage" (1), "complete Linkage" (2) y "single Linkage" (3).

En algunos libros de Análisis Multivariado se

(1) Paquete "factoextra" de R

propone juzgar la "calidad" de un agrupamiento calculando la correlación entre las distancias D y las "distancias cofiléticas" en el agrupamiento. Para dos individuos x_i y x_j su distancia cofilética se define como la altura (en el dendograma) del nodo en donde x_i y x_j fueron "clasificados en el mismo grupo". Si T denota la matriz de distancias cofiléticas para un agrupamiento, se define el "coeficiente de correlación cofilética" c como ⁽¹⁾

$$c = \frac{\sum_{i < j} (d_{ij} - \bar{d})(T_{ij} - \bar{T})}{\sqrt{\sum_{i < j} (d_{ij} - \bar{d})^2} \sqrt{\sum_{i < j} (T_{ij} - \bar{T})^2}}$$

donde \bar{d} es el promedio $\bar{d} = \frac{\sum_{i < j} d_{ij}}{\frac{n(n+1)}{2}}$

y \bar{T} es el promedio $\bar{T} = \frac{\sum_{i < j} T_{ij}}{\frac{n(n+1)}{2}}$

Si c tiene valores "grandes" (cercaños a 1)

(1) Sokal y Rohlf (1962) The comparison of dendograms by objective methods. Taxon, 11, 33-40

entonces se considerará que el agrupamiento propuesto por el dendograma es una "buena descripción de los datos". En la práctica un valor arriba de 0.75 se considera "grande". La función "cophenD.R" calcula los coeficientes de correlación copenético para los distintos métodos de agrupamiento jerárquico aglomerativo que tiene la función "hclust()". Al parecer el agrupamiento "más adecuado" para los datos de arrestos es el que se obtiene usando average linkage.

Se puede obtener un agrupamiento al cortar el dendograma a una altura determinada, o bien cuando se identifica la altura en que se forman un número de grupos K pre-determinado. Si se usa una altura ó un valor de K predeterminados, eso depende de la aplicación con la que se está trabajando.

No obstante se puede encontrar un valor "óptimo" de K (cuando no se tiene idea de este), usando una medida de "qué tan adecuado" resulta considerar a K como el número de grupos.

Esta medida se conoce como el "ancho promedio de la silueta"⁽²⁾ del agrupamiento

(2) Alan Julian Izenman "Modern Multivariate Statistical Techniques, Regression Classification and Manifold Learning". Springer Verlag páginas 426 - 429.