

Taller 6 AI: Clustering

Facultad de Ingeniería
Departamento de Electrónica

Nota: fecha máxima de entrega del informe **viernes 13 de mayo de 2022 a las 11:59 p.m.**
Por cada minuto de retraso en la entrega se descontará una (1) décima.

Objetivo:

- Estudiar el modelo de conectividad y el modelo de centroide. Para ello, se utilizará el agrupamiento jerárquico tipo aglomerativo y el algoritmo k -means. Además, se determinará el número “óptimo” de grupos, se interpretarán los resultados obtenidos y su utilidad como estrategia en el deporte analizado.

1. Utilice el conjunto de datos *Basketball data set* disponible en: <https://sci2s.ugr.es/keel/dataset.php?cod=1293sub1> para estudiar el aprendizaje automático no supervisado de tipo agrupamiento usando un modelo de conectividad (agrupamiento jerárquico) y un modelo de centroide (algoritmo k -means) y así realizar un análisis exploratorio de los datos.

Con base en lo anterior, realice:

- a.) Estadísticas descriptivas de los datos:
 - i.) Describa estadísticamente el conjunto de observaciones.
 - ii.) Obtenga los histogramas de las variables de entrada y analice si las observaciones provienen de una población con una distribución normal (Gaussiana).
 - iii.) Obtenga los diagramas de dispersión.
 - iv.) Examine la dependencia entre las variables de entrada con base en el criterio que considere idóneo (i.e., matriz de covarianza, coeficiente de correlación de *Pearson*, etc.).
- b.) Agrupamiento jerárquico (modelo de conectividad):
 - i.) Utilice el algoritmo de agrupamiento jerárquico aglomerativo para agrupar los datos.
 - ii.) Grafique el dendrograma y estime el número “óptimo” de grupos a través de la técnica vista en clase basada en la distancia en el dendrograma.
 - iii.) Utilice la técnica del codo (*Elbow Method*) para tener un criterio adicional del número “óptimo” de grupos.
 - iv.) Seleccione el número “óptimo” de grupos (k). Justifique su respuesta.
- c.) Agrupamiento con k -means (modelo de centroide):

- i.) Utilice el algoritmo k -means para agrupar los datos, teniendo como referencia el k seleccionado en el ítem previo.
- ii.) Grafique los clústeres en 3D (escoja las variables que considere pertinentes).
- iii.) Grafique el coeficiente de silueta e interprete los resultados.
- d.) Concluya sobre los resultados obtenidos.