Creado por: Isabel Maniega Ejercicio 5 In [1]: # Importing the libraries import numpy as np import matplotlib.pyplot as plt import pandas as pd import seaborn as sns In [2]: # Importing the dataset dataset = pd.read csv('Mall Customers.csv') dataset Genre Age Annual Income (k\$) Spending Score (1-100) CustomerID Out[2]: 0 1 Male 19 15 1 Male 21 15 81 2 3 Female 20 16 6 3 4 Female 23 16 77 17 4 5 Female 31 40 Female 120 195 196 35 79 126 196 197 Female 45 28 197 198 Male 32 126 74 198 199 137 Male 32 18 199 200 137 Male 30 83 200 rows × 5 columns K-Means In [3]: X = dataset.iloc[:, [3, 4]].valuesIn [4]: # Using the elbow method to find the optimal number of clusters from sklearn.cluster import KMeans wcss = []for i in range(1, 11): kmeans = KMeans(n_clusters = i, init = 'k-means++', random_state = 42) kmeans.fit(X)wcss.append(kmeans.inertia_) plt.plot(range(1, 11), wcss) plt.title('The Elbow Method') plt.xlabel('Number of clusters') plt.ylabel('WCSS') plt.show() The Elbow Method 250000 200000 150000 100000 50000 8 10 Number of clusters · Otra forma... In [5]: from yellowbrick.cluster import KElbowVisualizer from sklearn.cluster import KMeans In [6]: model = KMeans() visualizer = KElbowVisualizer(model, k=(1, 10)).fit(X)visualizer.show() Distortion Score Elbow for KMeans Clustering 0.06 elbow at k = 4, score = 73679.789250000 0.05 200000 distortion score 150000 100000 0.02 50000 0.01 1 2 3 4 5 6 7 8 9 Out[6]: <AxesSubplot: title={'center': 'Distortion Score Elbow for KMeans Clustering'}, xlabel='k', ylabel='distortion score'> Usaremos 5 Clusters In [7]: # Fitting K-Means to the dataset kmeans = KMeans(n_clusters = 5, init = 'k-means++', random_state = 42) y kmeans = kmeans.fit predict(X) y_kmeans 2, 3, 2, 3, 2, 3, 2, 3, 2, 3, 2, 3, 2, 3, 2, 3, 2, 3, 2, 3, 2, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 4, 1, 4, 0, 4, 1, 4, 1, 4, 0, 4, 1, 4, 1, 4, 1, 4, 1, 4, 0, 4, 1, 4], dtype=int32) In [8]: # Visualising the clusters $plt.scatter(X[y_kmeans == 0, 0], X[y_kmeans == 0, 1], s = 100, c = 'red', label = 'Cluster 1')$ plt.scatter(X[y_kmeans == 1, 0], X[y_kmeans == 1, 1], s = 100, c = 'blue', label = 'Cluster 2')
plt.scatter(X[y_kmeans == 2, 0], X[y_kmeans == 2, 1], s = 100, c = 'green', label = 'Cluster 3')
plt.scatter(X[y_kmeans == 3, 0], X[y_kmeans == 3, 1], s = 100, c = 'cyan', label = 'Cluster 4') $plt.scatter(X[y_kmeans == 4, 0], X[y_kmeans == 4, 1], s = 100, c = 'magenta', label = 'Cluster 5')$ plt.scatter(kmeans.cluster centers [:, 0], kmeans.cluster centers [:, 1], s = 300, c = 'yellow', label = 'Centr plt.title('Clusters of customers') plt.xlabel('Annual Income (k\$)') plt.ylabel('Spending Score (1-100)') plt.legend() plt.show() Clusters of customers 100 80 Cluster 1 60 Cluster 2 Cluster 3

Spending Score (1-100) Cluster 4 Cluster 5 40 Centroids 20 0 20 40 100 120 Annual Income (k\$) • Otra forma... In [9]: sns.scatterplot(data=dataset, x="Annual Income (k\$)", y="Spending Score (1-100)", hue=kmeans.labels_) plt.scatter(kmeans.cluster_centers_[:,0], kmeans.cluster_centers_[:,1], marker="X", c="r", s=80, label="centroids") plt.legend() plt.show() 100 80 Spending Score (1-100) 60 1 2 3 40 centroids 20 0 20 40 80 100 120 Annual Income (k\$) **Hierarchical Clustering** In [10]: # Using the dendrogram to find the optimal number of clusters import scipy.cluster.hierarchy as sch dendrogram = sch.dendrogram(sch.linkage(X, method = 'ward')) plt.title('Dendrogram') plt.xlabel('Customers') plt.ylabel('Euclidean distances') plt.show() Dendrogram 400 350 300 Euclidean distances 250 200 150 100 50 0 Customers Se observan 5 grupos

140

140

In [11]: # Fitting Hierarchical Clustering to the dataset from sklearn.cluster import AgglomerativeClustering hc = AgglomerativeClustering(n_clusters = 5, affinity = 'euclidean', linkage = 'ward') y_hc = hc.fit_predict(X) y_hc Out[11]: array([4, 3, 4, 4, 3, 4, 4, 3, 4, 4, 4, 4, 4, 4, 4, 4, 4, 4, 4, 4, 3, 4, 3, 4, 3, 4, 3, 4, 3, 4, 3, 4, 3, 4, 3, 4, 3, 4, 3, 4, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 2, 1, 2, 1, 2, 0, 2, 0, 2, 1, 2, 0, 2, 0, 2, 0, 2, 0, 2, 1, 2, 0, 2, 1, 2, 0, 2, 0, 2, 0, 2, 0, 2, 0, 2, 0, 2, 1, 2, 0, 2, 0, 2, 0, 2, 0, 2, 0, 2, 0, 2, 0, 2, 0, 2, 0, 2, 0, 2, 0, 2, 0, 2, 0, 2, 0, 2, 0, 2, 0, 2, 0, 2, 0, 2, 0, 2]) In [12]: # Visualising the clusters $plt.scatter(X[y_hc == 0, 0], X[y_hc == 0, 1], s = 100, c = 'red', label = 'Cluster 1')$ $plt.scatter(X[y_hc == 3, 0], X[y_hc == 3, 1], s = 100, c = 'cyan', label = 'Cluster 4')$ $plt.scatter(X[y_hc == 4, 0], X[y_hc == 4, 1], s = 100, c = 'magenta', label = 'Cluster 5')$ plt.title('Clusters of customers') plt.xlabel('Annual Income (k\$)') plt.ylabel('Spending Score (1-100)') plt.legend() plt.show() Clusters of customers 100 80 Spending Score (1-100) Cluster 1 60 Cluster 2 Cluster 3 Cluster 4 40 Cluster 5 20 0 20 40 100 120 140 Annual Income (k\$) · Otra forma... In [15]: sns.scatterplot(data=dataset, x="Annual Income (k\$)", y="Spending Score (1-100)", hue=hc.labels_) plt.legend() plt.show() 100 Spending Score (1-100) 60 2 40 20 0 0 40 20 60 80 100 120 140

Annual Income (k\$)

Creado por:

Isabel Maniega