

MEMORIA PRÁCTICA 2 - GCOM: Diagrama de Voronói y Clustering

Nombre: Gabriel Alba Serrano
Subgrupo: U1

1. Introducción

La clasificación de grupos ó clustering de un sistema consiste en agrupar objetos por similitud, en grupos o conjuntos de manera que los miembros del mismo grupo tengan características similares. El algoritmo apropiado para dicha clasificación depende del conjunto de datos que se analiza y el uso que se le dará a los resultados.

La idea de un grupo de datos similares resulta incompleta y subjetiva, y por ello existen miles de algoritmos. Alguno de los modelos más conocidos se basan en la conectividad, centroides, distribución, densidad, etc. En esta práctica vamos a estudiar dos algoritmos, uno basado en agrupamiento por centroides (KMeans) y otro basado en agrupamiento por densidad (DBSCAN).

La principal característica del algoritmo KMeans es que clasifica en clusters de aproximadamente tamaño similar (con un número similar de elementos), debido a que siempre asignará cada elemento al centroide más cercano, según la correspondiente métrica. Esto a menudo provoca cortes incorrectos en los bordes de los grupos, ya que el algoritmo optimiza los centroides, no las fronteras. Los clusters en este caso se denominan Diagramas de Voronói.

Por otro lado la principal característica del DBSCAN es que los clusters están definidos como áreas de densidad más alta que en el resto del conjunto de datos. Los objetos en áreas esparcidas son conocidos como ruido o puntos frontera. Este algoritmo no proporciona un buen clustering en un sistema que no tiene bajadas de densidad de sus elementos, es decir, un sistema con una densidad similar en todas sus regiones.

2. Material usado

2.1. Método

Utilizando el software de python3 y su librería sklearn, hemos aplicado los algoritmos de clustering KMeans y DBSCAN a un sistema. Además hemos graficado los diagramas de Voronói de dicho sistema de acuerdo con la clasificación más óptima del algoritmo KMeans, es decir aquella con el coeficiente de Silhouette más próximo al valor de 1.

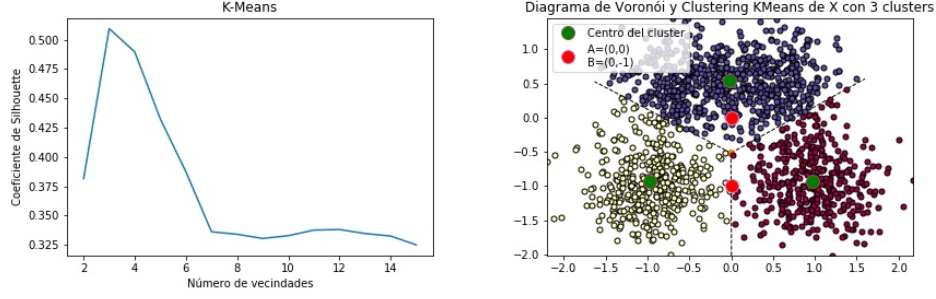
2.2. Datos

- Los archivos de texto Grados.en.la.facultad.matematicas.txt y Personas.en.la.facultad.matematicas.txt, siendo el primero un sistema de 1500 elementos en el que cada elemento tiene dos estados: X_1 ='nivel de estrés' y X_2 ='afición al rock'. Vamos a aplicar los distintos algoritmos de clustering a dicho sistema que denominaremos como X.
- Las plantillas 1 y 2 de Python que aplican un ejemplo de KMeans con métrica euclidiana y otro ejemplo de DBSCAN con métrica de Manhattan, respectivamente.

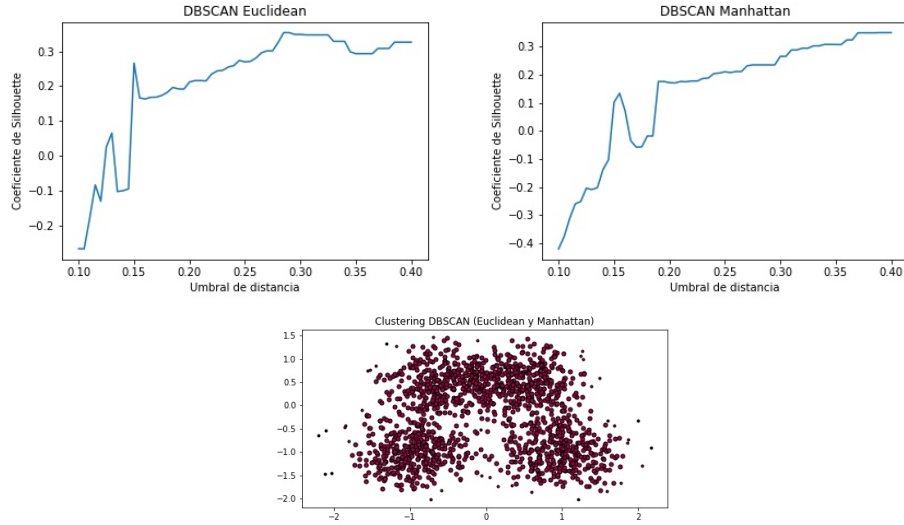
3. Resultados

1. Hemos calculado el coeficiente de Silhouette para cada $k \in \{2, 3, \dots, 15\}$, siendo k el número de vecindades (clusters) que se utiliza para aplicar el algoritmo KMeans. Para $k = 3$, el coeficiente

de Silhouette vale 0.510 que es el valor más próximo a 1 de todos los calculados, por lo tanto la clasificación KMeans del sistema X con 3 clusters es la más óptima.



- Además hemos calculado el coeficiente de Silhouette al aplicar el algoritmo de DBSCAN, primero con métrica euclidiana y luego con métrica de Manhattan. Ahora variando el umbral de distancia $\epsilon \in (0.1, 0.4)$ y fijando el número mínimo de elementos $n_0 = 10$. Con métrica 'Euclidean' el umbral de distancia óptimo es $\epsilon = 0.285$ y con métrica 'Manhattan' $\epsilon = 0.390$, es decir, para estos valores, sus respectivos coeficientes de Silhouette son los más altos ($s \approx 0.350$, con ambas métricas). Cabe destacar que para ambas métricas y con el umbral de distancia óptimo hemos obtenido que la clasificación de los elementos del sistema X a través del algoritmo DBSCAN ha sido de en un único cluster.



- Utilizando la función `kmeans.predict` hemos obtenido que los puntos $a := (0,0)$ y $b := (0,-1)$ pertenecen a los clusters 2 y 0, respectivamente, generados al aplicar el algoritmo KMeans con 3 clusters ó vecindades.

4. Conclusión

Observamos que el algoritmo KMeans clasifica los elementos de X en 3 clusters, frente al algoritmo DBSCAN que con ambas métricas (Euclidean y Manhattan) clasifica los elementos en un único cluster, es decir, no hace ninguna clasificación de los elementos del sistema, esto se debe a que el sistema X no sufre bajadas de densidad respecto a sus elementos. También hay que destacar que el valor óptimo del coeficiente de Silhouette de KMeans es mayor que el valor óptimo de DBSCAN. Por lo tanto la clasificación del algoritmo KMeans es más apropiada que la del DBSCAN para este sistema.

5. Anexo con el script/código utilizado

```
1 """
2 GCOM - PRACTICA 2: DIAGRAMA DE VORONOI Y CLUSTERING
3 NOMBRE: Gabriel Alba Serrano
4 SUBGRUPO: U1
5 """
6
7 import numpy as np
8 from sklearn.cluster import KMeans
9 from sklearn.cluster import DBSCAN
10 from sklearn import metrics
11 from scipy.spatial import ConvexHull, convex_hull_plot_2d
12 from scipy.spatial import Voronoi, voronoi_plot_2d
13 import matplotlib.pyplot as plt
14
15 # Aqu tenemos definido el sistema X de 1500 elementos (personas) con dos
16 # estados: X1 = "nivel de estr s" y X2 = "afici n al rock"
17 archivo1 = "C:/Users/Gabriel/Documents/CC-MAT/Geometr a Computacional/Pr ctica 2/
18     Personas_en_la_facultad_matematicas.txt"
19 archivo2 = "C:/Users/Gabriel/Documents/CC-MAT/Geometr a Computacional/Pr ctica 2/
20     Grados_en_la_facultad_matematicas.txt"
21 X = np.loadtxt(archivo1)
22 Y = np.loadtxt(archivo2)
23 labels_true = Y[:,0]
24
25 #Si quisieramos estandarizar los valores del sistema, har amos:
26 #from sklearn.preprocessing import StandardScaler
27 #X = StandardScaler().fit_transform(X)
28
29 #Envolvente convexa, envoltura convexa o c psula convexa, de X
30 hull = ConvexHull(X)
31 convex_hull_plot_2d(hull)
32 plt.title('Envolvente convexa del sistema X')
33 plt.plot(X[:,0],X[:,1], 'ro', markersize=1)
34 plt.show()
35
36 """ APARTADO 1 """
37
38 """ Aplicamos al sistema X el algoritmo KMeans para cada n mero de vecindades
39 k = {2,3,...,15}, y comprobamos qu coeficiente de Silhouette es mayor seg n k,
40 y de esta forma obtenemos el n mero ptimo de celdas de Voronoi. """
41
42 neighborhoods = list(range(2,16))
43 silhouette = []
44 for k in neighborhoods:
45     #Usamos la inicializaci n aleatoria "random_state=0"
46     kmeans = KMeans(n_clusters = k, random_state=0).fit(X)
47     labels = kmeans.labels_
48     silhouette.append( metrics.silhouette_score(X, labels) )
49
50 # Gr fica del coeficiente de Silhouette respecto al n mero de vecindades
51 plt.plot(neighborhoods,silhouette)
52 plt.xlabel('N mero de vecindades')
53 plt.ylabel('Coeficiente de Silhouette')
54 plt.title('K-Means')
55 plt.savefig('Silhouette KMeans.jpg')
56 plt.show()
57
58 """ Definimos n_clusters := k , tal que max(silhouette) = silhouette[k]
59 De esta forma, n_clusters es el n mero de vecindades que mejor empareja los
60 elementos de cada cluster (vecindad), es decir es el n mero de vecindades m s
61 ptimo . """
62 n_clusters = neighborhoods[ silhouette.index(max(silhouette))]
63 print('El n mero ptimo de vecindades es: k = ' + str(n_clusters) + '\n')
64
65 """ Definimos la etiqueta que indica a qu cluster pertenece cada elemento y
66 los centros de cada cluster, cuando el n mero de vecindades es n_clusters """
```

```

67 kmeans = KMeans(n_clusters = n_clusters, random_state=0).fit(X)
68 labels = kmeans.labels_
69 centers = kmeans.cluster_centers_
70
71 """ Representamos la clasificaci n del sistema X, por el algoritmo KMeans
72 y su respectivo diagrama de Voronoi. """
73
74 unique_labels = set(labels)
75 colors = [plt.cm.Spectral(each)
76            for each in np.linspace(0, 1, len(unique_labels))]
77
78 """ Definimos las celdas de Voronoi seg n los centros de los clusters,y
79 graficamos el diagrama de Voronoi. """
80 vor = Voronoi(centers)
81 voronoi_plot_2d(vor)
82
83 """ Graficamos cada elemento del sistema X, seg n al cluster al que pertenezcan. """
84 for lab, col in zip(unique_labels, colors):
85     if lab == -1:
86         # Black used for noise.
87         col = [0, 0, 0, 1]
88
89     class_member_mask = (labels == lab)
90
91     xy = X[class_member_mask]
92     plt.xlim(min(X[:,0]), max(X[:,0]))
93     plt.ylim(min(X[:,1]), max(X[:,1]))
94     plt.plot(xy[:, 0], xy[:, 1], 'o', markerfacecolor=tuple(col),
95             markeredgecolor='k', markersize=5)
96
97 """ Graficamos los centros de cada cluster. """
98 plt.plot(centers[:,0], centers[:,1], 'o', markersize=12, markerfacecolor="green",
99         label='Centro del cluster')
100
101 """ Graficamos los puntos a=(0,0) y b=(0,-1), que los necesitaremos para el
102 apartado 3. """
103 problem = np.array([[0, 0], [0, -1]])
104 clases_pred = kmeans.predict(problem)
105 plt.plot(problem[:,0],problem[:,1], 'o', markersize=12, markerfacecolor="red", label='A
106         =(0,0) \nB=(0,-1)')
107
108 plt.title('Diagrama de Voron i y Clustering KMeans de X con ' + str(n_clusters) + '
109         clusters')
110 plt.legend()
111 plt.savefig('Diagrama de Voron i y Clustering KMeans con 3 clusters.jpg')
112 plt.show()
113
114 """ APARTADO 2 """
115
116 """ Aplicamos al sistema X el algoritmo DBSCAN para epsilon (umbral de distancia)
117 perteneciente al intervalo (0.1 , 0.4), y comprobamos qu coeficiente de Silhouette
118 es mayor seg n epsilon, y de esta forma obtenemos el n mero ptimo de celdas de
119 Voronoi. """
120
121 # Los posibles valores de epsilon van a ser los siguientes puntos:
122 epsilon = np.arange(0.1, 0.4, 0.005)
123
124 """ M trica 'Euclidean' """
125 s_euclidean = []
126
127 for e in epsilon:
128     db = DBSCAN(eps=e, min_samples=10, metric='euclidean').fit(X)
129     labels = db.labels_
130     s = metrics.silhouette_score(X, labels)
131     s_euclidean.append(s)
132
133 # Graficamos
134 plt.plot(epsilon, s_euclidean)

```

```

134 plt.xlabel('Umbral de distancia')
135 plt.ylabel('Coeficiente de Silhouette')
136 plt.title('DBSCAN Euclidean')
137 plt.savefig('Silhouette DBSCAN Euclidean.jpg')
138 plt.show()
139
140 # Definimos el epsilon ptimo para la m trica Euclidean
141 epsilon_euclidean = format(epsilon[ s_euclidean.index(max(s_euclidean)) ] , '.3f')
142 print('El umbral de distancia ptimo para la m trica Euclidean es: = ' +
      epsilon_euclidean + '\n')
143
144
145 ### M trica 'Manhattan'
146 s_man = []
147
148 for e in epsilon:
149     db = DBSCAN(eps=e, min_samples=10, metric='manhattan').fit(X)
150     labels = db.labels_
151     s = metrics.silhouette_score(X, labels)
152     s_man.append(s)
153
154 # Graficamos
155 plt.plot(epsilon, s_man)
156 plt.xlabel('Umbral de distancia')
157 plt.ylabel('Coeficiente de Silhouette')
158 plt.title('DBSCAN Manhattan')
159 plt.savefig('Silhouette DBSCAN Manhattan.jpg')
160 plt.show()
161
162 # Definimos el epsilon ptimo para la m trica Manhattan
163 epsilon_man = format( epsilon[ s_man.index(max(s_man)) ] , '.3f' )
164 print('El umbral de distancia ptimo para la m trica Manhattan es: = ' +
      epsilon_man + '\n')
165
166 """ Finalmente en este apartado vamos a comparar gr ficamente el resultado de
167 aplicar el algoritmo DBSCAN con m trica 'Euclidean' y 'Manhattan' al conjunto X,
168 es decir, el clustering, frente al clustering obtenido al aplicar KMeans """
169
170 # Gr fica DBSCAN Euclidean
171
172 db = DBSCAN(eps=float(epsilon_euclidean), min_samples=10, metric='euclidean').fit(X)
173 labels = db.labels_
174 core_samples_mask = np.zeros_like(db.labels_, dtype=bool)
175 core_samples_mask[db.core_sample_indices_] = True
176
177 unique_labels = set(labels)
178 colors = [plt.cm.Spectral(each)
179           for each in np.linspace(0, 1, len(unique_labels))]
180
181 n_clusters_ = len(set(labels)) - (1 if -1 in labels else 0)
182 n_noise_ = list(labels).count(-1)
183
184 plt.figure(figsize=(8,4))
185 for k, col in zip(unique_labels, colors):
186     if k == -1:
187         # Black used for noise.
188         col = [0, 0, 0, 1]
189
190     class_member_mask = (labels == k)
191
192     xy = X[class_member_mask & core_samples_mask]
193     plt.plot(xy[:, 0], xy[:, 1], 'o', markerfacecolor=tuple(col),
194             markeredgecolor='k', markersize=5)
195
196     xy = X[class_member_mask & ~core_samples_mask]
197     plt.plot(xy[:, 0], xy[:, 1], 'o', markerfacecolor=tuple(col),
198             markeredgecolor='k', markersize=3)
199
200 plt.title('DBSCAN Euclidean del sistema X con ' + str(n_clusters_) + ' clusters')
201 plt.show()
202

```

```

203
204 # Gráfica DBSCAN Manhattan
205
206 db = DBSCAN(eps=float(epsilon_man), min_samples=10, metric='manhattan').fit(X)
207 labels = db.labels_
208 core_samples_mask = np.zeros_like(db.labels_, dtype=bool)
209 core_samples_mask[db.core_sample_indices_] = True
210
211 unique_labels = set(labels)
212 colors = [plt.cm.Spectral(each)
213           for each in np.linspace(0, 1, len(unique_labels))]
214
215 n_clusters_ = len(set(labels)) - (1 if -1 in labels else 0)
216 n_noise_ = list(labels).count(-1)
217
218 plt.figure(figsize=(8,4))
219 for k, col in zip(unique_labels, colors):
220     if k == -1:
221         # Black used for noise.
222         col = [0, 0, 0, 1]
223
224     class_member_mask = (labels == k)
225
226     xy = X[class_member_mask & core_samples_mask]
227     plt.plot(xy[:, 0], xy[:, 1], 'o', markerfacecolor=tuple(col),
228             markeredgecolor='k', markersize=5)
229
230     xy = X[class_member_mask & ~core_samples_mask]
231     plt.plot(xy[:, 0], xy[:, 1], 'o', markerfacecolor=tuple(col),
232             markeredgecolor='k', markersize=3)
233
234 plt.title('DBSCAN Manhattan del sistema X con ' + str(n_clusters_) + ' clusters')
235 plt.show()
236
237 # Conclusión
238 print('Como conclusión observamos que el algoritmo KMeans clasifica los elementos de
239 X en 3 clusters, frente al algoritmo DBSCAN que con ambas métricas (Euclidean y
240 Manhattan) clasifica los elementos en un único cluster.')
241
242 print('Hay que destacar que la clasificación del algoritmo KMeans es más ptima que
243 la del DBSCAN porque su coeficiente de Silhouette es el que más se aproxima a 1.
244 \n')
245
246
247 """ APARTADO 3 """
248
249 problem = np.array([[0, 0], [0, -1]])
250 clases_pred = kmeans.predict(problem)
251 print('Según el clustering resultante de aplicar el algoritmo KMeans, los puntos a
252 =(0,0) y b=(0,-1) deberían pertenecer a los clusters ' + str(clases_pred[0]) + ' y
253 ' + str(clases_pred[1]) + ', respectivamente.')

```