

PRÁCTICA 3: DIAGRAMA DE VORONÓI Y CLUSTERING

Autor
Compañera de código

Martín Fernández de Diego
Belén Sánchez Centeno

1. INTRODUCCIÓN

A partir de un sistema X de 1000 elementos construido como muestra aleatoria entorno a unos centros definidos, se pide determinar una clasificación por vecindades —o clusters— a través de tres algoritmos distintos: KMeans, DBSCAN con métrica euclidiana y DBSCAN con métrica manhattan. También se pide representar el diagrama de Voronói correspondiente, entre otros.

2. MATERIAL USADO

2.1. Apartado i)

El **clustering por KMeans** toma k centroides arbitrarios y clasifica cada punto de cierto sistema X de forma que pertenezca al cluster asociado al centroide más cercano. Después, toma el punto medio de cada cluster como centroide de la siguiente iteración. El algoritmo cicla hasta que la distancia entre los nuevos puntos medios y los centroides de la anterior iteración sea inferior a cierto épsilon.

Se han utilizado las funciones:

- `plot_silhouette_kmeans(sistema)`
- `plot_clusters_kmeans(sistema)`

Para encontrar el **coeficiente de Silhouette**¹ con el algoritmo KMeans, en la primera función, el sistema X es entrenado para $k \in \{2, 3, \dots, 15\}$ vecindades con la instrucción `KMeans(num_vecindades, estado_aleatorio)`. A continuación, se calcula el valor con la librería `sklearn` y la función `silhouette_score(sistema, vecindades)` que recibe X y las vecindades devueltas por KMeans. Por último, se muestra una **primera gráfica** con los coeficientes —eje de ordenadas— correspondientes a cada número de vecindades —eje de abscisas— y se decide computacionalmente cuál es el número de vecindades que corresponde al mayor coeficiente de Silhouette.

La segunda función es la encargada de, dado el mejor k estimado por Silhouette, recalcular KMeans y, con ayuda de la plantilla ofrecida por el profesor *Robert Monjo*, **representar los clusters** por colores sobre una **segunda gráfica**.

¹El coeficiente de Silhouette es un indicador de ayuda a la decisión para saber cuán de buena ha sido una clasificación.

Adicionalmente, se añade la posición y nombre de los centros de los clusters para calcular y superponer el **teselado de Voronói**² con `voronoi = Voronoi(centros)` y `voronoi_plot_2d(voronoi)`.

2.2. Apartado ii)

El **clustering por DBSCAN** toma un punto arbitrario del sistema X y lo clasifica en función del número de puntos que haya en la ϵ -bola centrada en él respecto a n_0 . Si es mayor, es un centro, se crea una clase y se marcan sus vecinos —como miembros— para ser visitados. Si es menor, se clasifica como ruido a la espera de, quizás, revisitarlo. Antes de volver a tomar un punto arbitrario, se realizan las visitas pendientes y se marcan como puntos fronterizos si el número de puntos en su ϵ -bola es inferior a n_0 o, de nuevo centros, si es superior. Termina cuando no se pueden crear nuevas clases.

Se han utilizado las funciones:

- `plot_silhouette_dbscan(sistema, métrica)`
- `plot_clusters_dbscan(sistema, métrica)`

De manera semianáloga al apartado anterior, se halla el **coeficiente de Silhouette**, esta vez con el algoritmo DBSCAN, en la primera función. El sistema X es entrenado para un umbral $\epsilon \in (0, 1, 0, 4)$ con la instrucción `DBSCAN(epsilon, mínimo, métrica)`. A continuación, se obtiene el coeficiente de la misma forma que antes. Por último, se muestra una **primera gráfica** con los coeficientes —eje de ordenadas— correspondientes a cada valor de umbral —eje de abscisas— y se decide computacionalmente cuál es el umbral correspondiente al mayor coeficiente de Silhouette.

La segunda función, volviendo a hacer uso del material cedido en las plantillas, muestra una **segunda gráfica** con la **separación por clusters** en colores que arroja DBSCAN. Los puntos negros son los puntos de ruido del algoritmo.

2.3. Apartado iii)

A partir del cálculo de los centros de las vecindades y para cierta definición de distancia, se puede estimar a qué cluster pertenece un punto dado.

Se han utilizado las funciones:

²El diagrama de Voronói de un conjunto de marcas es la división del espacio en regiones de forma que a cada punto se le asigna la región formada por los puntos que son más cercanos a ella que a ninguna otra marca.

- `distancia_euclidea(punto, centro)`
- `distancia_manhattan(punto, centro)`

Se calcula la distancia —euclidiana o manhattan— del punto a cada centroide del sistema. El centroide más cercano identificará el cluster al que pertenece el nuevo punto.

Dados dos puntos del plano $P = (p_1, p_2)$ y $Q = (q_1, q_2)$, la fórmula de la **distancia euclidiana** implementada es

$$d_e(P, Q) = \sqrt{(q_1 - p_1)^2 + (q_2 - p_2)^2}$$

y la fórmula de la **distancia manhattan** implementada es

$$d_m(P, Q) = |q_1 - p_1| + |q_2 - p_2|.$$

3. RESULTADOS

3.1. Apartado i)

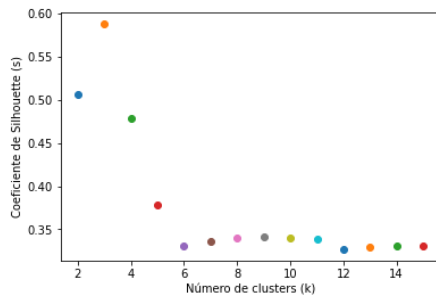


Fig. 1. Coeficientes de Silhouette por KMeans

El **clustering por KMeans** decide una clasificación en 3 clusters, en base al coeficiente de Silhouette más alto.

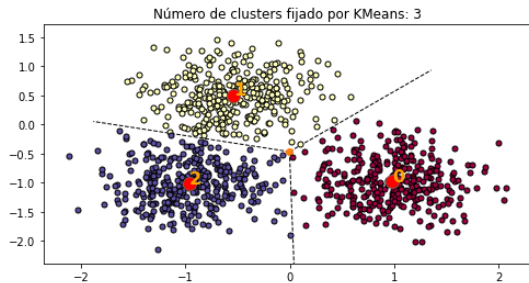


Fig. 2. Clusters por KMeans con teselado de Voronoi

3.2. Apartado ii)

3.2.1. Métrica euclidiana

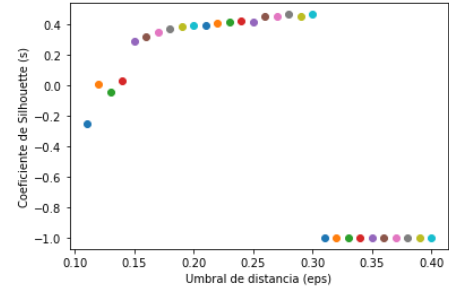


Fig. 3. Coeficientes de Silhouette por DBSCAN

El **clustering por DBSCAN** decide un umbral de 0.28 que se traduce en 2 clusters.

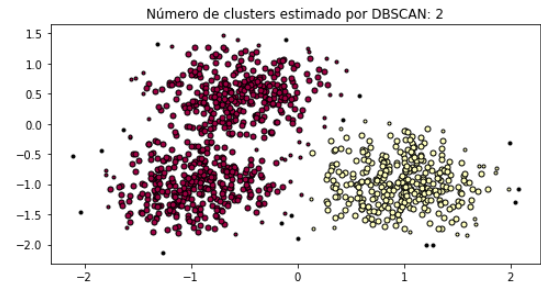


Fig. 4. Clusters por DBSCAN

3.2.2. Métrica manhattan

Por restricciones de formato y dada su similaridad con la euclidiana, no se mostrará el caso de la métrica manhattan.

3.3. Apartado iii)

El punto $(0, 0)$ se encuentra en el cluster 1 y el punto $(0, -1)$ en el cluster 2 según las distancias euclídea y manhattan. La comprobación corrobora este resultado.

4. CONCLUSIÓN

Los algoritmos de clusterización han arrojado resultados diferentes para un mismo sistema X. El funcionamiento de cada algoritmo está directamente relacionado. Mientras que KMeans categoriza en base al teselado de Voronoi una vez encontrados unos centroides suficientemente buenos, DBSCAN expande la región de una categoría sin importar su centroide mientras haya suficiente densidad de puntos.

Dado que el coeficiente de Silhouette debe interpretarse como ayuda a la decisión, se podría tomar un umbral menor a 0,15 para diferenciar más de 2 clusters con DBSCAN.

5. ANEXO CON EL SCRIPT Y CÓDIGO UTILIZADO

5.1. Código

```
1  """
2  PRÁCTICA 3: DIAGRAMA DE VORONOI Y CLUSTERING
3  Belén Sánchez Centeno
4  Martín Fernández de Diego
5  """
6
7  import numpy as np
8  import matplotlib.pyplot as plt
9
10 from sklearn.cluster import KMeans
11 from sklearn.cluster import DBSCAN
12 from sklearn import metrics
13 from sklearn.datasets import make_blobs
14 from scipy.spatial import ConvexHull, convex_hull_plot_2d
15 from scipy.spatial import Voronoi, voronoi_plot_2d
16
17 """
18 Dado un conjunto de puntos X
19 muestra la gráfica de los coeficientes de Silhouette para cada número de clusters
20 devuelve el número óptimo de clusters asociado al mayor coeficiente de Silhouette
21 """
22 def plot_silhouette_kmeans(X):
23     # Mostramos los coeficientes de Silhouette para cada k
24     # y obtenemos el k asociado al mayor coeficiente de todos
25     max_s = -1
26     for k in range(2,16):
27         kmeans = KMeans(n_clusters=k, random_state=0).fit(X)
28         labels = kmeans.labels_
29         silhouette = metrics.silhouette_score(X, labels)
30         # Decidimos computacionalmente el número óptimo de clusters
31         if max_s < silhouette:
32             max_s = silhouette
33             max_k = k
34         plt.plot(k, silhouette, 'o')
35     plt.xlabel("Número de clusters (k)")
36     plt.ylabel("Coeficiente de Silhouette (s)")
37     plt.show()
38
39     return max_k
40
41 """
42 Dado un conjunto de puntos X y el número de vecindades óptimo n_clusters
43 muestra los clusters estimados por el algoritmo KMeans
44 devuelve la instancia de KMeans para el apartado iii)
45 """
46 def plot_clusters_kmeans(X,n_clusters):
47     # Tomamos el número de vecindades óptimo devuelto por el coeficiente de Silhouette
48     # y volvemos a ejecutar KMeans para mostrar los clusters por colores
49     kmeans = KMeans(n_clusters=n_clusters, random_state=0).fit(X)
50     labels = kmeans.labels_
51     centers = kmeans.cluster_centers_
52     # Mantenemos la misma proporción al resto en el tamaño de la gráfica
53     fig = plt.figure(figsize=(8,4))
54     ax = fig.add_subplot(111)
55
56     # Mostramos el teselado de Voronoi
57     vor = Voronoi(centers)
58     voronoi_plot_2d(vor,ax=ax)
59
60     # Representamos el resultado con un plot
61     unique_labels = set(labels)
62     colors = [plt.cm.Spectral(each) for each in np.linspace(0, 1, len(unique_labels))]
63     for k, col in zip(unique_labels, colors):
64         if k == -1:
65             # Black used for noise.
66             col = [0, 0, 0, 1]
```

```

67         class_member_mask = (labels == k)
68         xy = X[class_member_mask]
69         plt.plot(xy[:, 0], xy[:, 1], 'o', markerfacecolor=tuple(col), markeredgecolor='k
70         ', markersize=5)
71     # Mostramos los centros
72     plt.plot(centers[:,0],centers[:,1],'o', markersize=12, markerfacecolor="red")
73     for i in range(len(centers)):
74         plt.text(centers[i,0],centers[i,1],str(i),color='orange',fontsize=16,fontweight=
75         'black')
76     # Configuramos los atributos de la gráfica con sus límites
77     plt.title('Número de clusters fijado por KMeans: %d' % n_clusters)
78     plt.xlim([np.min(X[:,0])-0.25,np.max(X[:,0])+0.25])
79     plt.ylim([np.min(X[:,1])-0.25,np.max(X[:,1])+0.25])
80     plt.show()
81     return kmeans
82 """
83 Dado un conjunto de puntos X, el tipo de métrica y la sensibilidad de búsqueda del
84 umbral de distancia
85 muestra la gráfica de los coeficientes de Silhouette para cada umbral de distancia
86 devuelve el umbral de distancia óptimo asociado al mayor coeficiente de Silhouette
87 """
88 def plot_silhouette_dbscan(X,metric,step=0.01):
89     # Mostramos los coeficientes de Silhouette para cada épsilon
90     # y obtenemos el épsilon asociado al mayor coeficiente de todos
91     max_s = -1
92     for epsilon in np.arange(0.11,0.4,step):
93         # Utilizamos el algoritmo de DBSCAN para mínimo 10 elementos
94         db = DBSCAN(eps=epsilon, min_samples=10, metric=metric).fit(X)
95         labels = db.labels_
96         # Aseguramos el valor de Silhouette si el número de clusters es 1
97         n_clusters_ = len(set(labels)) - (1 if -1 in labels else 0)
98         silhouette = metrics.silhouette_score(X, labels) if n_clusters_ != 1 else -1
99         # Decidimos computacionalmente el número óptimo de clusters
100        if max_s < silhouette:
101            max_s = silhouette
102            max_eps = epsilon
103        plt.plot(epsilon, silhouette, 'o')
104    plt.xlabel("Umbral de distancia (eps)")
105    plt.ylabel("Coeficiente de Silhouette (s)")
106    plt.show()
107    return max_eps
108 """
109 Dado un conjunto de puntos X, el tipo de métrica y el umbral de distancia
110 muestra los clusters estimados por el algoritmo DBSCAN con la métrica dada
111 """
112 def plot_clusters_dbscan(X,metric,epsilon):
113     # Tomamos el épsilon óptimo devuelto por el coeficiente de Silhouette
114     # y volvemos a ejecutar DBSCAN para mostrar los clusters por colores
115     db = DBSCAN(eps=epsilon, min_samples=10, metric=metric).fit(X)
116     core_samples_mask = np.zeros_like(db.labels_, dtype=bool)
117     core_samples_mask[db.core_sample_indices_] = True
118     labels = db.labels_
119     # Number of clusters in labels, ignoring noise if present.
120     n_clusters_ = len(set(labels)) - (1 if -1 in labels else 0)
121     n_noise_ = list(labels).count(-1)
122
123     print("Número óptimo de vecindades: ", n_clusters_)
124
125     unique_labels = set(labels)
126     colors = [plt.cm.Spectral(each) for each in np.linspace(0, 1, len(unique_labels))]
127     plt.figure(figsize=(8,4))
128     for k, col in zip(unique_labels, colors):
129         # Black used for noise.
130         if k == -1:
131             col = [0, 0, 0, 1]
132         class_member_mask = (labels == k)

```

```

134         xy = X[class_member_mask & core_samples_mask]
135         plt.plot(xy[:, 0], xy[:, 1], 'o', markerfacecolor=tuple(col), markeredgecolor='k
136         ', markersize=5)
137         xy = X[class_member_mask & ~core_samples_mask]
138         plt.plot(xy[:, 0], xy[:, 1], 'o', markerfacecolor=tuple(col), markeredgecolor='k
139         ', markersize=3)
140     plt.title('Número de clusters estimado por DBSCAN: %d' % n_clusters_)
141     plt.show()
142
143     def distancia_euclidea(punto,centro):
144         return np.sqrt((punto[0]-centro[0])**2 + (punto[1]-centro[1])**2)
145
146     def distancia_manhattan(punto,centro):
147         return abs(punto[0]-centro[0]) + abs(punto[1]-centro[1])
148
149     # FORMATO
150     class Formato:
151         BOLD = "\033[1m"
152         RESET = "\033[0m"
153
154     # Aquí tenemos definido el sistema X de 1000 elementos de dos estados
155     # construido a partir de una muestra aleatoria entorno a unos centros:
156     centers = [[-0.5, 0.5], [-1, -1], [1, -1]]
157     X, labels_true = make_blobs(n_samples=1000, centers=centers, cluster_std=0.4,
158                                random_state=0)
159
160     #Si quisieramos estandarizar los valores del sistema, haríamos:
161     #from sklearn.preprocessing import StandardScaler
162     #X = StandardScaler().fit_transform(X)
163
164     #Envolvente convexa, envoltura convexa o cápsula convexa
165     hull = ConvexHull(X)
166     convex_hull_plot_2d(hull)
167
168     plt.plot(X[:,0],X[:,1],'ro', markersize=1)
169     plt.show()
170
171     # APARTADO i)
172     print("\n" + Formato.BOLD + "Apartado i)" + Formato.RESET)
173
174     max_k = plot_silhouette_kmeans(X)
175     print("Número óptimo de vecindades: ", max_k)
176     kmeans = plot_clusters_kmeans(X,max_k)
177
178     # APARTADO ii)
179     print("\n" + Formato.BOLD + "Apartado ii)" + Formato.RESET)
180
181     # Euclidean
182     euclidean_metric = 'euclidean'
183     euclidean_max_eps = plot_silhouette_dbscan(X,euclidean_metric)
184     print("Umbral de distancia euclidiana óptimo: ",euclidean_max_eps)
185     plot_clusters_dbscan(X,euclidean_metric,euclidean_max_eps)
186
187     # Manhattan
188     manhattan_metric = 'manhattan'
189     manhattan_max_eps = plot_silhouette_dbscan(X,manhattan_metric)
190     print("Umbral de distancia manhattan óptimo: ",manhattan_max_eps)
191     plot_clusters_dbscan(X,manhattan_metric,manhattan_max_eps)
192
193     # APARTADO iii)
194     print("\n" + Formato.BOLD + "Apartado iii)" + Formato.RESET)
195
196     centers = kmeans.cluster_centers_
197
198     puntol = [0,0]
199     distancias_euclideas = [distancia_euclidea(puntol,centers[i]) for i in range(len(centers
200     ))]
201     cluster1 = distancias_euclideas.index(min(distancias_euclideas))
202     print("El punto ",puntol,"se encuentra en el cluster ",cluster1,"según la distancia
203     euclídea")
204     distancias_manhattan = [distancia_manhattan(puntol,centers[i]) for i in range(len(

```

```

        centers))]
199 cluster1 = distancias_manhattan.index(min(distancias_manhattan))
200 print("El punto ",punto1,"se encuentra en el cluster ",cluster1,"según la distancia
    manhattan")
201 print("Comprobación: ",kmeans.predict([punto1])[0])
202
203 punto2 = [0,-1]
204 distancias_euclideas = [distancia_euclidea(punto2,centers[i]) for i in range(len(centers
    ))]
205 cluster2 = distancias_euclideas.index(min(distancias_euclideas))
206 print("El punto ",punto2,"se encuentra en el cluster ",cluster2,"según la distancia
    euclídea")
207 distancias_manhattan = [distancia_manhattan(punto2,centers[i]) for i in range(len(
    centers))]
208 cluster2 = distancias_manhattan.index(min(distancias_manhattan))
209 print("El punto ",punto2,"se encuentra en el cluster ",cluster2,"según la distancia
    manhattan")
210 print("Comprobación: ",kmeans.predict([punto2])[0])

```

5.2. Ejecución

```

1 Apartado i)
2 Número óptimo de vecindades: 3
3
4 Apartado ii)
5 Umbral de distancia euclidiana óptimo: 0.2799999999999999
6 Número óptimo de vecindades: 2
7 Umbral de distancia manhattan óptimo: 0.3599999999999999
8 Número óptimo de vecindades: 2
9
10 Apartado iii)
11 El punto [0, 0] se encuentra en el cluster 1 según la distancia euclídea
12 El punto [0, 0] se encuentra en el cluster 1 según la distancia manhattan
13 Comprobación: 1
14 El punto [0, -1] se encuentra en el cluster 2 según la distancia euclídea
15 El punto [0, -1] se encuentra en el cluster 2 según la distancia manhattan
16 Comprobación: 2

```