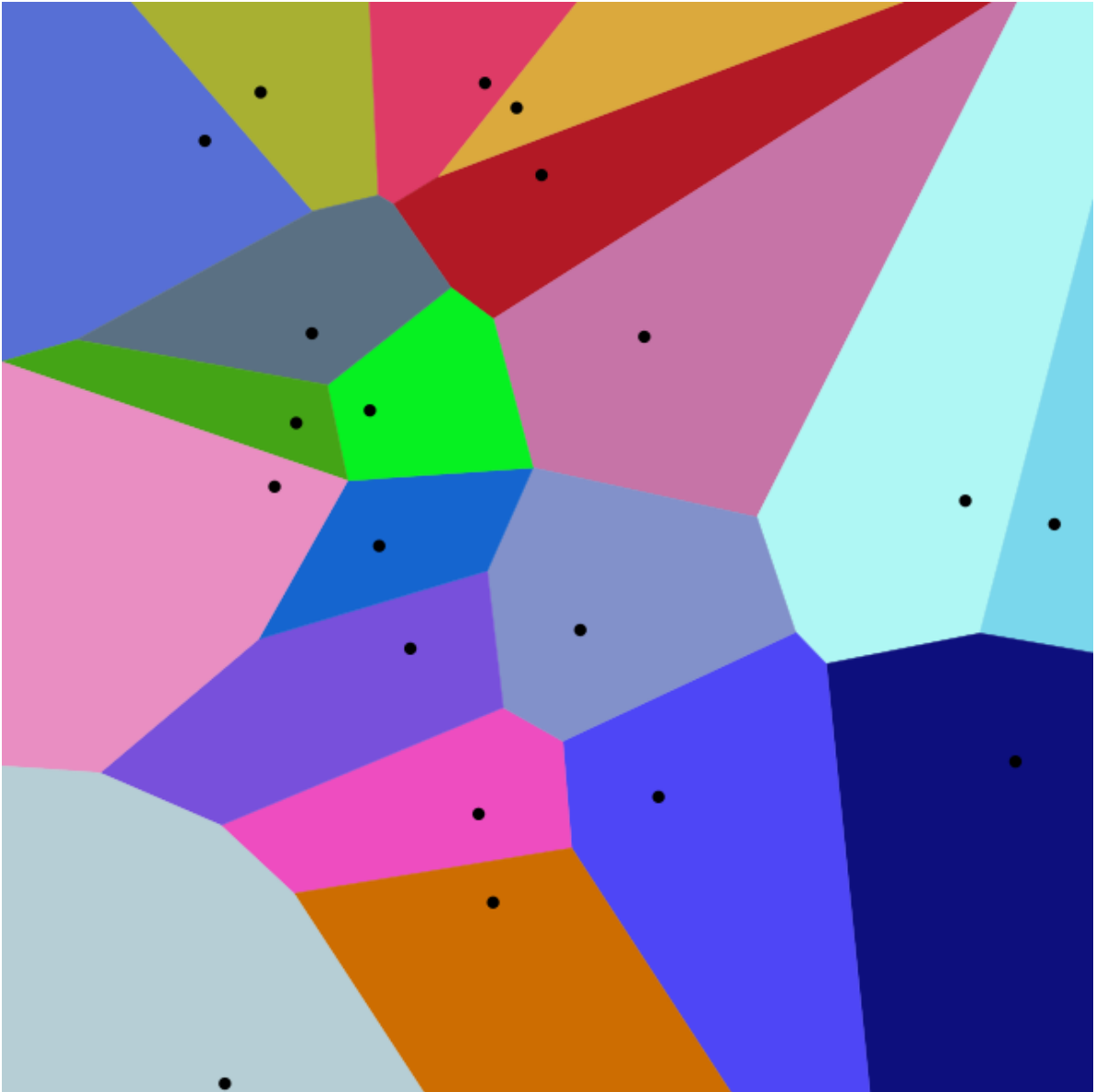


GEOMETRÍA COMPUTACIONAL

PRÁCTICA 3



DAVID SEIJAS PÉREZ

1. Introducción

En esta tercera práctica queremos aprender a clasificar un sistema con un número determinado de elementos, con dos estados cada uno de ellos, a partir de un número determinado de cluster o vecindades de Voronoi. Además, utilizaremos el coeficiente de Silhouette para determinar el número óptimo de estas vecindades.

2. Material Usado

Para el desarrollo de esta práctica he utilizado diversas librerías dadas por Python. Primero, como siempre, hemos usado *matplotlib.pyplot* para la representación de gráfica y, esta vez, además hemos utilizado *Voronoi* y *voronoi_plot_2d* para representar el diagrama de Voronoi. Hemos utilizado también las librerías *KMeans* y *DBSCAN* para usar los algoritmos con el mismo nombre. Por último, hemos usado *metrics* y *make_blobs* para, principalmente, calcular el coeficiente *Silhouette* y el sistema X, respectivamente.

Además, he cogido las plantillas *GCOM2022-practica3_plantilla1* y *GCOM2022-practica3_plantilla2* para reusar diversas partes del código como la representación gráfica de los clusters y los cálculos hechos en ambos algoritmos.

3. Metodología

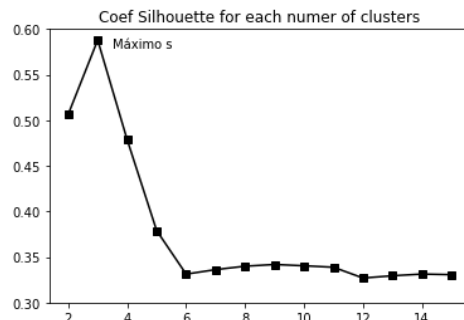
Para el primer apartado, he hecho la función *mostrarKMEANS*, reutilizando código de la plantilla 1, que muestra los clusters calculados por el algoritmo **KMeans** del sistema para un k dado. Esta función es utilizada solo para representar gráficamente el sistema junto al diagrama de Voronoi al final del apartado. He realizado, también, la función *apartado1* en la que realizamos el algoritmo KMeans para las distintas ks especificadas y calculamos el coeficiente Silhouette asociada a cada una. Posteriormente, muestro en una gráfica la relación de este coeficiente según el k correspondiente y, así, determinamos el mejor número de vecindades, asociado al máximo coeficiente. Para este k vuelvo a realizar el algoritmo para poder mostrar en una gráfica el sistema clusterizado con los datos apropiados.

En el segundo apartado, la función *mostrarDBSCAN* muestra, como en la plantilla 2, los clusters obtenidos según el algoritmo **DBSCAN**. En la función *apartado2* he desarrollado el algoritmo DBSCAN, como se aporta en la plantilla2, para una distancia dada por parámetros y para todos los valores de ϵ desde 0.1 a 0.4 con paso 0.005. Mi idea ha sido recorrer todos los ϵ aplicando el algoritmo para ver con cual obtenemos el mejor coeficiente de Silhouette y, así, tener el mejor número de cluster para cada distancia. Una vez hallado este umbral de distancia, he vuelto a aplicar el algoritmo para poder mostrar por pantalla la gráfica con los datos óptimos. Además, con una pequeña modificación de la función, no incluida en el código, he estudiado como se comportaba este algoritmo para distintos valores de ϵ con cada una de las distancias.

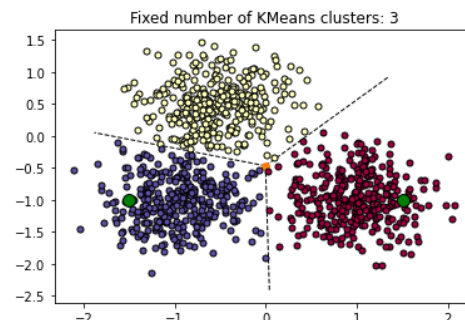
En ambos apartados, en las funciones de mostrar resultados, he representado los puntos a y b del apartado 3 resaltados para ver gráficamente la pertenencia de cada uno a su cluster.

4. Resultados

1. En el primer apartado he obtenido que el número óptimo de vecindades es $k = 3$ para las cuales obtenemos el coeficiente de Silhouette $\bar{s} = 0.588$. A continuación muestro las gráficas obtenidas.



(a) Relación entre \bar{s} y k



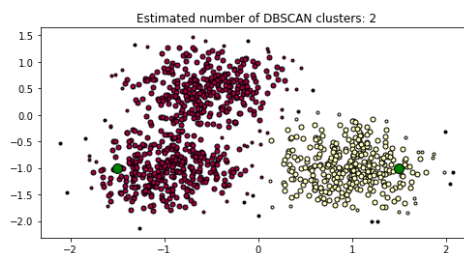
(b) Diagrama Voronoi con a y b

2. En el segundo apartado he obtenido que el número estimado de clusters es 2 para ambas distancias en los mejores casos. Para la distancia Euclídea (a) he obtenido:

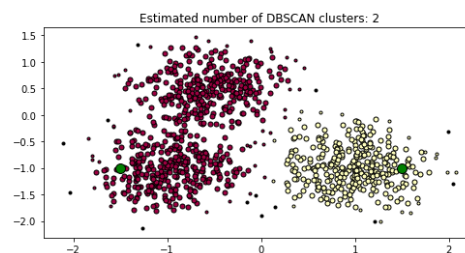
- Umbral de distancia con el que se alcanza el mejor \bar{s} : $\epsilon = 0.28$
- Número estimado de elementos ruidosos: 17
- Coeficiente de Silhouette: $\bar{s} = 0.467$

y para la distancia Manhattan (b):

- Umbral de distancia con el que se alcanza el mejor \bar{s} : $\epsilon = 0.365$
- Número estimado de elementos ruidosos: 13
- Coeficiente de Silhouette: $\bar{s} = 0.546$



(a) Métrica Euclídea



(b) Métrica Manhattan

Además, he analizado el comportamiento del algoritmo en las dos distancias según el parámetro ϵ . Para la distancia euclídea, con valores de $\epsilon > 3$ he obtenido un solo cluster con más elementos ruidosos y peores valores de \bar{s} cuanto mayor es el valor. Decrementando ϵ se mantienen los dos clusters hasta 0.15, a partir de donde empiezan a aumentar mucho los clusters y llegamos, incluso, a obtener valores $\bar{s} < 0$. El comportamiento para la distancia Manhattan es parecido, aunque se comporta mejor con valores más grande ϵ , mantiene los dos clusters para valores más grandes de este parámetro y con mejores \bar{s} . Sin embargo, para valores más pequeños se comporta

peor, a partir de $\epsilon < 2$ empiezan a aumentar mucho los clusters obteniendo valores más pequeños del coeficiente de Silhouette, hasta aproximadamente $\bar{s} = -0.43$

3. Para el apartado 3, en la función *apartado1* observamos que el punto *a* pertenece al cluster 2 y el *b* al 0, según la notación del algoritmo KMeans. La pertenencia de estos puntos a las vecindades con los dos algoritmos se muestran en los gráficos.

5. Conclusión

Me parecen fascinantes los métodos de clasificación vistos y cómo pueden variar las clasificaciones modificando ligeramente los datos. Por ejemplo, en el apartado 2 para algunos valores de ϵ se pueden obtener resultados realmente malos, incluso con valores del coeficiente de Silhouette negativo y una cantidad de clusters absurda para el sistema que estamos estudiando. Esto nos hace ver la importancia de determinar los valores óptimos para los cuales los algoritmos de clasificación se comportan como esperamos. Además, con el apartado 2 vemos también como se consiguen mejores resultados con la distancia Manhattan, lo cual me ha sorprendido bastante pues pensaba que sería al revés, mejor con la euclídea.

6. Anexo: Código

```
1  '''
2  DAVID SEIJAS PEREZ
3  Practica 3
4  '''
5
6  import numpy as np
7  from sklearn.cluster import DBSCAN
8  from sklearn.cluster import KMeans
9  from sklearn import metrics
10 from sklearn.datasets import make_blobs
11 from scipy.spatial import Voronoi, voronoi_plot_2d
12 import matplotlib.pyplot as plt
13
14
15 '''
16 Funcion que muestra los clusters calculados con KMEANS
17 '''
18 def mostrarKMEANS(X, labels, nClusters):
19     unique_labels = set(labels)
20     colors = [plt.cm.Spectral(each)
21               for each in np.linspace(0, 1, len(unique_labels))]
22
23     #plt.figure(figsize=(8,4))
24     for k, col in zip(unique_labels, colors):
25         if k == -1:
26             # Black used for noise.
27             col = [0, 0, 0, 1]
```

```

28
29     class_member_mask = (labels == k)
30
31     xy = X[class_member_mask]
32     plt.plot(xy[:, 0], xy[:, 1], 'o', markerfacecolor=tuple(col),
33              markeredgecolor='k', markersize=5)
34
35 plt.plot(problem[:,0],problem[:,1],'ko', markersize=10,
36          markerfacecolor="green")
37 plt.axis('tight')
38 plt.title('Fixed number of KMeans clusters: %d' % nClusters)
39 plt.show()
40
41 def apartado1():
42     ks = [i for i in range(2, 16)]
43     coefsSil = [0]*14
44
45     centers = [[-0.5, 0.5], [-1, -1], [1, -1]]
46     X, labels_true = make_blobs(n_samples=1000, centers=centers,
47                                cluster_std=0.4,
48                                random_state=0)
49
50     #Algoritmo KMEANS para las distintas vencidades k
51     for i in range(len(ks)):
52         kmeans = KMeans(n_clusters=ks[i], random_state=0).fit(X)
53         labels = kmeans.labels_
54         silhouette = metrics.silhouette_score(X, labels)
55         #mostrarKMEANS(X, labels, problem, ks[i])
56
57         '''
58         # Etiqueta de cada elemento (punto)
59         print(kmeans.labels_)
60         # ndice de los centros de vencindades o regiones de Voronoi
61         para cada elemento (punto)
62         print(kmeans.cluster_centers_)
63         #Coeficiente de Silhouette
64         '''
65         print("Silhouette Coefficient for k = " + str(ks[i])
66               + ": %0.3f" % silhouette)
67         coefsSil[i] = silhouette
68
69     plt.ylim(0.3, 0.6)
70     plt.plot(ks, coefsSil, 'ks-')
71     plt.title('Coef Silhouette for each numer of clusters')
72     plt.text(3.5, 0.58, "M ximo s")
73     plt.show()
74     #Con esto vemos que el mejor n mero de vecindades es k=3
75     k = 3
76
77     #Calculamos de nuevo todo para mostrar los clusters
78     kmeans = KMeans(n_clusters=k, random_state=0).fit(X)
79     labels = kmeans.labels_

```

```

79
80     #Mostramos Diagrama de Voronoi junto con los clusters
81     vor = Voronoi(kmeans.cluster_centers_)
82     voronoi_plot_2d(vor)
83     mostrarKMEANS(X, labels, k)
84     #Al quitar el plt.figure de mostrarClusters no nos elimina la
      figura de voronoi creada antes
85
86     #Calculamos a qu clusters pertenecen los puntos a y b (Apartado
      3)
87     clases_pred = kmeans.predict(problem)
88     print("Clasificaci n de los puntos a=(0,0) y b=(0,-1) para Kmeans
      :")
89     print(clases_pred)
90     print("\n-----\n")
91
92
93
94     '''
95     Apartado 2
96     '''
97
98     '''
99     Funcion que muestra los cluster calculados con DBSCAN
100    '''
101    def mostrarDBSCAN(X, labels, core_samples_mask, nClusters):
102        unique_labels = set(labels)
103        colors = [plt.cm.Spectral(each)
104                  for each in np.linspace(0, 1, len(unique_labels))]
105
106        plt.figure(figsize=(8,4))
107        for k, col in zip(unique_labels, colors):
108            if k == -1:
109                # Black used for noise.
110                col = [0, 0, 0, 1]
111
112            class_member_mask = (labels == k)
113
114            xy = X[class_member_mask & core_samples_mask]
115            plt.plot(xy[:, 0], xy[:, 1], 'o', markerfacecolor=tuple(col),
116                    markeredgecolor='k', markersize=5)
117
118            xy = X[class_member_mask & ~core_samples_mask]
119            plt.plot(xy[:, 0], xy[:, 1], 'o', markerfacecolor=tuple(col),
120                    markeredgecolor='k', markersize=3)
121
122            plt.plot(problem[:,0],problem[:,1], 'ko', markersize=10,
123                    markerfacecolor="green")
124            plt.title('Estimated number of DBSCAN clusters: %d' % nClusters)
125            plt.show()
126
127    def apartado2(distancia):
128        epsilons = [0.1 + 0.005*i for i in range(60)]

```

```

129     silMax = 0
130     epsilonMax = 0.1
131
132     centers = [[-0.5, 0.5], [-1, -1], [1, -1]]
133     X, labels_true = make_blobs(n_samples=1000, centers=centers,
134                                cluster_std=0.4,
135                                random_state=0)
136
137     #Algoritmo DBSCAN
138     for i in range(len(epsilons)):
139         db = DBSCAN(eps=epsilons[i], min_samples=10, metric=distancia)
140             .fit(X)
141         core_samples_mask = np.zeros_like(db.labels_, dtype=bool)
142         core_samples_mask[db.core_sample_indices_] = True
143         labels = db.labels_
144
145         # Number of clusters in labels, ignoring noise if present.
146         n_clusters_ = len(set(labels)) - (1 if -1 in labels else 0)
147         n_noise_ = list(labels).count(-1)
148
149         silhouette = metrics.silhouette_score(X, labels)
150         if(silhouette > silMax):
151             silMax = silhouette
152             epsilonMax = epsilons[i]
153
154     db = DBSCAN(eps=epsilonMax, min_samples=10, metric=distancia).fit(
155         X)
156     core_samples_mask = np.zeros_like(db.labels_, dtype=bool)
157     core_samples_mask[db.core_sample_indices_] = True
158     labels = db.labels_
159     n_clusters_ = len(set(labels)) - (1 if -1 in labels else 0)
160     n_noise_ = list(labels).count(-1)
161
162     print(distancia + " DISTANCE")
163     print('Best epsilon: %0.3f' % epsilonMax)
164     print('Estimated number of clusters: %d' % n_clusters_)
165     print('Estimated number of noise points: %d' % n_noise_)
166     print("Adjusted Rand Index: %0.3f"
167           % metrics.adjusted_rand_score(labels_true, labels))
168     print("Silhouette Coefficient: %0.3f"
169           % silMax)
170     print("\n-----\n")
171
172     mostrarDBSCAN(X, labels, core_samples_mask, n_clusters_)
173
174     problem = np.array([[ -1.5, -1], [1.5, -1]])
175     apartado1()
176     apartado2("euclidean")
177     apartado2("manhattan")

```