

Diagrama de Voronói y clustering

Sergio González Montero (U2)

5 de marzo de 2024

1. Objetivo

A partir del sistema proporcionado por los archivos txt se pretende clasificar su contenido en *clusters* o celdas de Voronói utilizando el coeficiente de Silhouette, estimando con él un número óptimo de vecindades $k \in \{1, 2, \dots, 15\}$. También se explorará el número óptimo a través del umbral de distancia $\epsilon \in (0, 1, 0, 4)$. Además, se hará una predicción de pertenencia a grupos de ciertos puntos concretos especificados en apartado de resultados.

2. Material y datos

Se usaron los siguientes apartados de las notas de Robert Monjo de la asignatura de Geometría Computacional, a día 04/03/2024: Definición 2.2.5, Diagramas de Voronói; Algoritmo 2.2.2, Clustering o clasificación por k-medias / k-medoids; Definición 2.2.6, Coeficiente de Silhouette y Algoritmo 2.2.3, Algoritmo DBSCAN.

Se usaron los archivos *Personas_de_villa_laminera.txt* y *Franjas_de_edad.txt* que contienen los elementos que conforman el sistema de datos a analizar.

En cuanto al código, se usaron como base las plantillas *GCOM2024-practica2_plantilla1* y *GCOM2024-practica2_plantilla2*. En él, se han utilizado las librerías siguientes:

1. sklearn: aporta algoritmos de agrupamiento como KMeans o DBSCAN, valoración de dichos agrupamientos con Silhouette
2. numpy: operaciones matemáticas y cálculos numéricos
3. matplotlib.pyplot: para visualización de gráficos
4. scipy.spatial: añade el cálculo y representación del diagrama de Voronói así como las métricas euclidiana y de Manhattan

3. Resultados

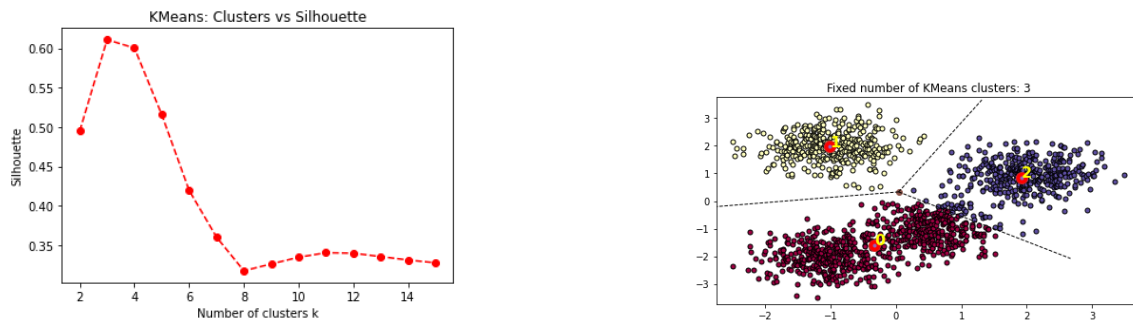


Figura 1: Optimum clusters: 3, Silhouette 0.611

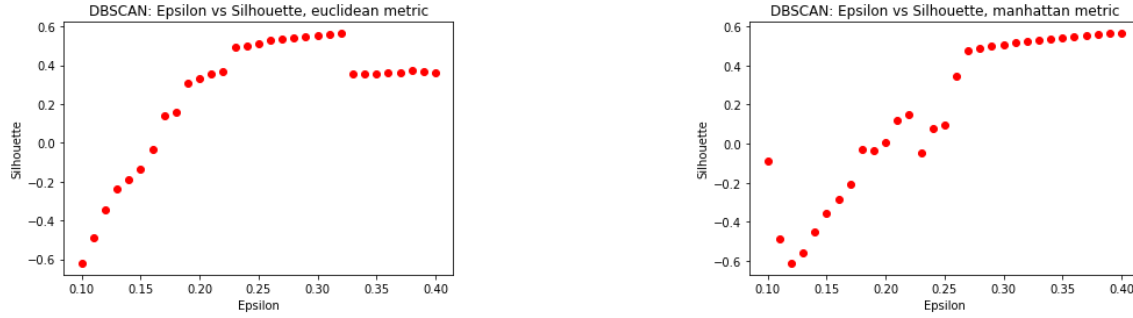


Figura 2: **Optimum euclidean epsilon: 0.32, Silhouette: 0.56357**
Optimum manhattan epsilon: 0.4, Silhouette: 0.56424



Figura 3: **Clusters según métrica**

En cuanto a las predicciones, siendo los puntos $a = [0.5, 0]$, $b = [0, -3]$, y los pares (etiqueta, centroide) $(0, [-0.32559490084985876, -1.5872946175637381])$, en color rojo, $(1, [-1.0077659574468085, 1.9790159574468091])$, en color amarillo y $(2, [1.922177033492823, 0.8609330143540673])$, en color verde, el punto a pertenece al cluster verde con la etiqueta 2, tanto por la métrica euclidiana como por la de manhattan, mientras que para el punto b se da una predicción de pertenencia al cluster rojo etiquetado como 0. Ambos resultados son consistentes con el dado por `kmeans.predict()`.

4. Conclusión

Como se ha podido comprobar la solución no es única pues, si nos guiamos únicamente por el coeficiente de Silhouette, podríamos decir que el algoritmo DBSCAN es más preciso con la métrica manhattan que con la euclidiana. Sin embargo, dependerá de la naturaleza de los datos, como su densidad o distribución, la elección de una u otra métrica para el cálculo de ϵ . De igual manera sucede con KMeans que, por defecto, usa la euclidiana y, aun así, difiere del resultado dado por DBSCAN con la misma métrica. Por tanto, se deduce que el coeficiente de Silhouette es una buena forma de deducir el número de clusters aunque en algunos casos puntuales haya que contrastar con otros métodos, como el índice de Calinski-Harabasz o la varianza explicada .

5. Código

Programa 1: practica2.py

```

1  # Sergio Gonzalez Montero
2  # Victor Martin Martin
3
4  import numpy as np
5
6  from sklearn.cluster import KMeans
7  from sklearn.cluster import DBSCAN
8  from sklearn import metrics
9  import matplotlib.pyplot as plt
10 from scipy.spatial import Voronoi, voronoi_plot_2d
11 from scipy.spatial.distance import euclidean, cityblock
12
13 print("-----Dataset-----\n")
14 # Definimos el sistema A
15 archivo1 = "Personas_de_villa_laminera.txt" # Asignar archivos
16 archivo2 = "Franjas_de_edad.txt"
17 X = np.loadtxt(archivo1, skiprows=1) # Carga datos de un de un txt
18 Y = np.loadtxt(archivo2, skiprows=1)
19 labels_true = Y[:,0]
20
21 header = open(archivo1).readline()
22 print(header)
23 print(X)
24
25 # Pinta el dataset
26 plt.plot(X[:,0], X[:,1], 'ro', markersize=1)
27 plt.title("Stress vs Sweets")
28 plt.xlabel("Stress")
29 plt.ylabel("Sweets")
30 plt.show()
31
32 print("\n-----KMeans-----\n")
33 # Calculo del numero de vecindades optimo
34 n_clusters = []
35 sil_k = []
36 for k in range(2, 16):
37     n_cluster = k
38     n_clusters.append(k)
39
40     # Usamos la inicializacion aleatoria "random_state=0"
41     kmeans = KMeans(n_clusters=k, random_state=0).fit(X)
42     labels = kmeans.labels_

```

```

43     silhouette = metrics.silhouette_score(X, labels)
44     sil_k.append(silhouette)
45
46 plt.plot(n_clusters, sil_k, 'ro--')
47 plt.title("KMeans: Clusters vs Silhouette")
48 plt.xlabel("Number of clusters k")
49 plt.ylabel("Silhouette")
50 plt.show()
51 max_s_index = sil_k.index(max(sil_k))
52 # Numero de clusters asociado al Silhouette maximo
53 k_optimo = n_clusters[max_s_index]
54 print(f"Optimum clusters: {k_optimo} \n\
55 Silhouette {round(max(sil_k), 3)}\n")
56
57 # Se vuelve a calcular para la representacion particular con k_optimo
58 kmeans = KMeans(n_clusters=k_optimo, random_state=0).fit(X)
59 labels = kmeans.labels_
60 centroids = kmeans.cluster_centers_
61 unique_labels = set(labels)
62 colors = [plt.cm.Spectral(each)
63            for each in np.linspace(0, 1, len(unique_labels))]
64
65 fig = plt.figure(figsize=(8,4))
66 ax = fig.add_subplot(111)
67
68 for k, col in zip(unique_labels, colors):
69     if k == -1:
70         # Black used for noise.
71         col = [0, 0, 0, 1]
72
73     class_member_mask = (labels == k)
74
75     xy = X[class_member_mask]
76     plt.plot(xy[:, 0], xy[:, 1], 'o', markerfacecolor=tuple(col),
77              markeredgecolor='k', markersize=5)
78
79 # Graficado de centroides
80 plt.plot(centroids[:,0],centroids[:,1], 'o',
81          markersize=12, markerfacecolor="red")
82 for i in range(len(centroids)):
83     plt.text(centroids[i,0],centroids[i,1],str(i),
84             color='yellow',fontsize=16,fontweight='black')
85 # Diagrama de Voronoi
86 vor = Voronoi(centroids)
87 voronoi_plot_2d(vor,ax=ax)
88 # Acomodamiento de los ejes respecto al dataset

```

```

89 plt.xlim([min(X[:,0])-0.25,max(X[:,0])+0.25])
90 plt.ylim([min(X[:,1])-0.25,max(X[:,1])+0.25])
91
92 plt.title('Fixed number of KMeans clusters: %d' % k_optimo)
93 plt.show()
94
95 print("\n-----DBSCAN-----\n")
96 e = []
97 sil_e = []
98 def dbscan_silhouette(metric):
99     """
100     metric : str, metrica a usar por DBSCAN
101     Grafica epsilon vs Silhouette y calcula el
102     epsilon optimo para cada metrica
103     """
104     for epsilon in np.arange(0.10,0.4,0.01):
105
106         # Utilizamos el algoritmo de DBSCAN para minimo 10 elementos
107         db = DBSCAN(eps=epsilon, min_samples=10, metric=metric).fit(X)
108         labels = db.labels_
109         # Number of clusters in labels, ignoring noise if present.
110         n_clusters_ = len(set(labels)) - (1 if -1 in labels else 0)
111         silhouette=metrics.silhouette_score(X,labels)if n_clusters_!=1 else -1
112         sil_e.append(silhouette)
113         e.append(round(epsilon, 2))
114         plt.plot(round(epsilon, 2), silhouette, 'ro--')
115         plt.title(f"DBSCAN: Epsilon vs Silhouette, \
116 {metric} metric")
117         plt.xlabel("Epsilon")
118         plt.ylabel("Silhouette")
119         plt.show()
120
121         max_s_index = sil_e.index(max(sil_e))
122         # Numero de clusters asociado al Silhouette maximo
123         e_optimo = e[max_s_index]
124         printed = print(f"Optimum {metric} epsilon: {e_optimo}\nSilhouette: \
125 {round(max(sil_e), 3)}")
126         return printed, e_optimo
127
128 dbscan_silhouette('euclidean')[0]
129 dbscan_silhouette('manhattan')[0]
130
131 def plot_dbscan(metric):
132     """
133     metric : str, metrica a usar por DBSCAN
134     Grafica los diagramas correspondientes segun metrica

```

```

135     """
136     db = DBSCAN(eps=dbscan_silhouette(metric)[1],
137                 min_samples=10, metric=metric).fit(X)
138     core_samples_mask = np.zeros_like(db.labels_, dtype=bool)
139     core_samples_mask[db.core_sample_indices_] = True
140     labels = db.labels_
141     # Number of clusters in labels, ignoring noise if present.
142     n_clusters_ = len(set(labels)) - (1 if -1 in labels else 0)
143     n_noise_ = list(labels).count(-1)
144
145     unique_labels = set(labels)
146     colors = [plt.cm.Spectral(each)
147               for each in np.linspace(0, 1, len(unique_labels))]
148
149     plt.figure(figsize=(8,4))
150     for k, col in zip(unique_labels, colors):
151         if k == -1:
152             # Black used for noise.
153             col = [0, 0, 0, 1]
154
155         class_member_mask = (labels == k)
156
157         xy = X[class_member_mask & core_samples_mask]
158         plt.plot(xy[:, 0], xy[:, 1], 'o', markerfacecolor=tuple(col),
159                 markeredgecolor='k', markersize=5)
160
161         xy = X[class_member_mask & ~core_samples_mask]
162         plt.plot(xy[:, 0], xy[:, 1], 'o', markerfacecolor=tuple(col),
163                 markeredgecolor='k', markersize=3)
164
165     # Graficado de centroides
166     plt.plot(centroids[:,0],centroids[:,1], 'o',
167             markersize=12, markerfacecolor="red")
168     for i in range(len(centroids)):
169         plt.text(centroids[i,0],centroids[i,1],str(i),
170                 color='yellow',fontsize=16,fontweight='black')
171
172     # Diagrama de Voronoi
173     vor = Voronoi(centroids)
174     voronoi_plot_2d(vor,ax=ax)
175     # Acomodamiento de los ejes respecto al dataset
176     plt.xlim([min(X[:,0])-0.25,max(X[:,0])+0.25])
177     plt.ylim([min(X[:,1])-0.25,max(X[:,1])+0.25])
178     plt.title(f'Estimated number of DBSCAN clusters\
179 ({metric}): %d' % n_clusters_)
180     plt.show()

```

```
181 plot_dbscan('euclidean')
182 print('\n')
183 plot_dbscan('manhattan')
184
185 print("\n-----PREDICTION-----\n")
186 a, b = [1/2,0], [0,-3]
187 print(f"Points: a = {a}, b = {b}")
188 print("(Label, centroid)")
189 for i in range(len(centroids)):
190     print(f"({i}, {list(centroids[i])})")
191
192 a_d2 = [euclidean(a,centroid) for centroid in centroids]
193 b_d2 = [euclidean(b,centroid) for centroid in centroids]
194 a_dmanhattan = [cityblock(a,centroid) for centroid in centroids]
195 b_dmanhattan = [cityblock(b,centroid) for centroid in centroids]
196 cluster_a_e = a_d2.index(min(a_d2))
197 cluster_b_e = b_d2.index(min(b_d2))
198 cluster_a_m = a_dmanhattan.index(min(a_dmanhattan))
199 cluster_b_m = b_dmanhattan.index(min(b_dmanhattan))
200
201 print(f"\nPoint {a} belongs to cluster {cluster_a_e},\
202       green, by euclidean metric")
203 print(f"Point {b} belongs to cluster {cluster_b_e},\
204       red, by euclidean metric")
205 print(f"Point {a} belongs to {cluster_a_m},\
206       green, by manhattan metric")
207 print(f"Point {b} belongs to {cluster_b_m},\
208       red, by manhattan metric")
209
210 print("\nPrediction by kmeans.predict()")
211 a_predict = kmeans.predict([a])[0]
212 b_predict = kmeans.predict([b])[0]
213 print(f"Point a belongs to cluster {a_predict},\
214       green, by kmeans.predict()")
215 print(f"Point b belongs to cluster {b_predict},\
216       red, by kmeans.predict()")
```