

Practica 1: Código Huffman y Primer Teorema de Shannon

Pablo Jiménez Poyatos : U1

6 de marzo de 2024

1. Introducción

En el ámbito del análisis de datos, dos técnicas cruciales son el clustering y el diagrama de Voronói. Utilizando datos de la población "Villa Laminera", aplicaremos estas técnicas con los algoritmos KMeans y DBSCAN para agrupar a los individuos en clusters óptimos. Evaluaremos la calidad de los clusters utilizando el coeficiente de Silhouette y compararemos los resultados obtenidos con diferentes métricas y algoritmos.

2. Material usado

En esta investigación, he utilizado las plantillas proporcionadas en el campus virtual, “GCOM2024- practica2_plantilla1.py” y GCOM2024- practica2_plantilla2.py” como punto de partida para desarrollar mi código en Python. El archivo de datos, “Personas_de_villa_laminera.txt”, proporciona información crucial sobre las características de los individuos, con variables como el nivel de estrés y la afinidad por los dulces.

Para la manipulación y análisis de los datos, se importaron varias bibliotecas de Python. Esto incluyó **numpy** para operaciones con matrices y arrays, **matplotlib.pyplot** para la creación de gráficos y visualizaciones, y por último, las funciones y clases específicas de **scikit-learn** para implementar los algoritmos de clustering, como KMeans y DBSCAN, así como métricas de evaluación como el coeficiente de Silhouette. Además, se utilizó la biblioteca **scipy.spatial** para trabajar con el diagrama de Voronoi.

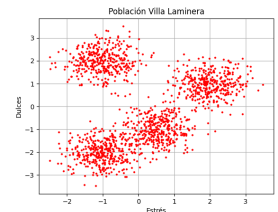


Figura 1: Grafico dispersion datos

El código **comienza** con el pre-procesamiento de los datos, donde se cargan y se muestran en un grafico de dispersión (Figura 1). En el **primer apartado**, se utiliza el algoritmo KMeans para determinar el número óptimo de clusters en el conjunto de datos. Se calcula el coeficiente de Silhouette para diferentes números de clusters (k) y se grafican en función de k para visualizar cómo varía la calidad de la agrupación. Luego, se identifica el valor máximo del coeficiente de Silhouette y se determina el número óptimo de clusters correspondiente. Con esta información, se entrena un modelo KMeans y se clasifican los datos utilizando un diagrama de Voronoi.

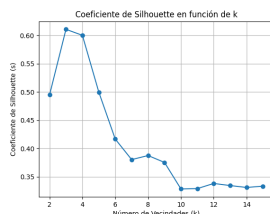


Figura 2: Relación del coeficiente de Silhouette en funcion de k

En el **segundo apartado**, se utiliza el algoritmo DBSCAN para estimar el número óptimo de clusters de manera automática. Se define un rango de valores para el parámetro eps y se establece un valor fijo para el número mínimo de puntos requeridos para formar un cluster. Luego, se itera sobre diferentes métricas de distancia y se calcula el coeficiente de Silhouette para cada combinación de parámetros. Finalmente, se grafican los resultados y se comparan con los del apartado anterior.

En el **tercer apartado**, se realizan predicciones utilizando el modelo KMeans entrenado en el primer apartado sobre dos puntos específicos en el espacio de características. Se utiliza el método predict del modelo KMeans para asignar cada punto al cluster más cercano según la estructura definida por el modelo.

3. Resultados

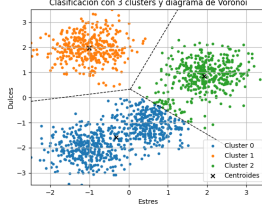


Figura 3: Diagrama de Voronoi usando KMEANS

Para el sistema A, se obtuvo el coeficiente de Silhouette (s) utilizando el algoritmo KMeans para diferentes números de vecindades k en el rango $[2,15]$. La Figura 2 muestra el valor de s en función de k , lo que permite determinar el número óptimo de vecindades. Basándonos en esta gráfica, se selecciona el valor de k que maximiza el coeficiente de Silhouette. En cuyo caso es $k = 3$ con coeficiente de Silhouette $s = 0,6108$.

A continuación, en la Figura 3, se muestra el diagrama de Voronoi que representa las regiones de influencia de cada cluster en el espacio de características. Esta visualización proporciona una representación clara de cómo los datos se agrupan en clusters y cómo se distribuyen las regiones de influencia en función de los centroides.

En el caso del algoritmo DBSCAN, se calculó el coeficiente de Silhouette para el mismo sistema A utilizando las métricas euclidean y 'manhattan'. En la Figura 4, se muestra el valor de s en función del umbral de distancia (eps) para ambas métricas. En él podemos observar que el mayor coeficiente de Silhouette para DBSCAN con la métrica euclidean es $s = 0,5636$ que le corresponde con umbral de distancia óptimo $Eps = 0,32$ y se estima que haya 3 clusters. Por otra parte, para la métrica manhattan, se estima el mismo número de clusters 3, y parecido coeficiente de Silhouette $s = 0,5642$ pero diferente umbral de distancia, $Eps = 0,4$. A continuación veremos las pequeñas diferencias entre algoritmos.

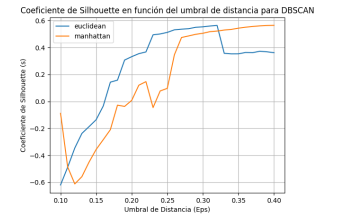
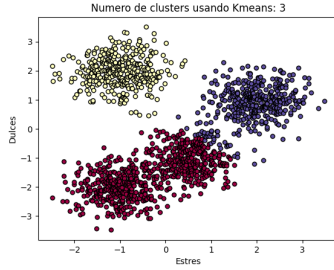
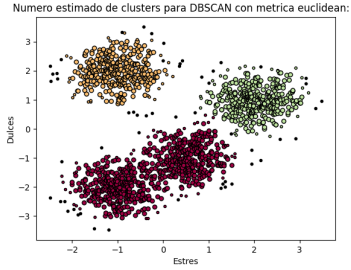


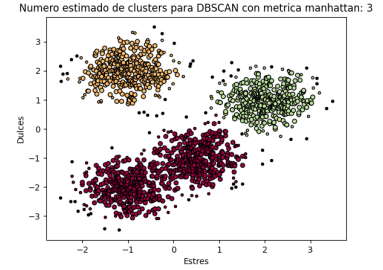
Figura 4: Relacion s en función de eps



(a) KMEANS



(b) DBSCAN euclidean



(c) DBSCAN manhattan

Figura 5: Comparación de las nubes de puntos dependiendo del algoritmo utilizado

Se puede observar que la única diferencia es, a parte de la distribución de los *core points* y *noise points* (se representa con el tamaño de los puntos), que los puntos que se encuentran en la frontera entre los clusters de color rojo y (verde o azul), usando KMEANS se quedan de color azul porque la distancia al centroide azul es más pequeña, pero usando DBSCAN se incluyen en el rojo porque hay más densidad de puntos en su vecindad.

Por último, para determinar la franja de edad a la que pertenecen las personas con coordenadas $a = (\frac{1}{2}, 0)$ y $b = (0, -3)$, vamos a graficarlos (Figura 6). En él, se puede observar como la persona a pertenecería al cluster número 2 (color verde) y la persona b, al cluster numero 0 (color azul). Para comprobarlo, utilizamos la función `kmeans.predict` que nos muestran lo mismo: El cluster predicho para el punto $[0.5, 0]$ es: $[2]$, y para el punto $[0, -3]$ es: $[0]$

4. Conclusión

Nuestra experimentación demuestra que tanto KMeans como DBSCAN son herramientas útiles para analizar la distribución y agrupación de una población en función de sus características. El uso de métricas como el coeficiente de Silhouette nos permite evaluar y comparar la calidad de las clasificaciones obtenidas, lo que proporciona información valiosa para comprender mejor la estructura de los datos y las relaciones entre los individuos en el conjunto de datos de Villa Laminera

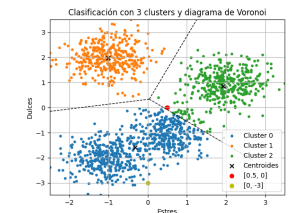


Figura 6: Predicción de la vecindad de los puntos

5. Anexo

5.1. Código implementado

```
1  """
2
3  Asignatura: Geometría computacional
4  Subgrupo: 1
5  Curso: 2023-2024
6  Alumno: Jiménez Poyatos, Pablo
7  Curso: 4 CC
8  Carrera: Grado en matemáticas.
9  Práctica 2. Diagrama de Voronói y clustering.
10
11  """
12
13  import os
14  import numpy as np
15  import matplotlib.pyplot as plt
16
17  from sklearn.cluster import KMeans, DBSCAN
18  from sklearn.metrics import silhouette_score
19  from scipy.spatial import Voronoi, voronoi_plot_2d
20
21
22
23
24  def preprocesamiento_datos(archive_name):
25      """
26      Función para cargar y preprocesar datos desde un archivo.
27
28      Parameters:
29          archive_name (str): Nombre del archivo que contiene los datos.
30
31      Return:
32          X (numpy.ndarray): Datos preprocesados.
33      """
34
35      ruta = os.getcwd()
36      archivo = os.path.join(ruta, archive_name)
37      X = np.loadtxt(archivo, skiprows=1)
38      return X
39
40  def calculo_silhouette(X, rango_inf, rango_sup):
41      """
42      Calcula el coeficiente de silhouette para un rango dado de número de
43      clusters.
44
45      Parameters:
46          X (numpy.ndarray): Datos.
47          rango_inf (int): Límite inferior del rango de número de clusters.
48          rango_sup (int): Límite superior del rango de número de clusters.
49
50      Return:
51          silhouette_scores (list): Lista de coeficientes de silhouette para cada
52          número de clusters en el rango dado.
53      """
54
55      silhouette_scores = []
56      for k in range(rango_inf, rango_sup + 1):
57          kmeans = KMeans(n_clusters=k, random_state=0).fit(X)
58          labels = kmeans.labels_
59          silhouette = silhouette_score(X, labels)
60          silhouette_scores.append(silhouette)
61      return silhouette_scores
62
```

```

63
64 def graficar(x,y,x_name, y_name, title_name, marker, style, metrics_l,
65             grid=True, legend=False):
66     """
67     Función para graficar datos.
68
69     Parámetros:
70         x (list): Valores del eje x.
71         y (list(list)): Valores del eje y (una lista de listas si se grafican
72                        múltiples líneas).
73         x_name (str): Nombre del eje x.
74         y_name (str): Nombre del eje y.
75         title_name (str): Título del gráfico.
76         marker (str): Tipo de marcador.
77         style (str): Estilo de línea.
78         metrics_l (list): Lista de nombres de las métricas.
79         grid (bool): Indica si mostrar la cuadrícula en el gráfico
80                     (por defecto True).
81         legend (bool): Indica si mostrar la leyenda en el gráfico
82                       (por defecto False).
83     """
84
85     plt.figure()
86     for i in range(len(y)):
87         plt.plot(x, y[i], marker=marker, linestyle=style, label=metrics_l[i])
88     plt.xlabel(x_name)
89     plt.ylabel(y_name)
90     plt.title(title_name)
91     plt.grid(grid)
92     if legend:
93         plt.legend()
94     plt.show()
95
96 def grafica_voronoi(kmeans, opt_k, X, labels, points, x_name, y_name, title):
97     """
98     Función para graficar un diagrama de Voronoi.
99
100    Parámetros:
101        kmeans: Modelo de KMeans entrenado.
102        opt_k (int): Número óptimo de clusters.
103        X (numpy.ndarray): Datos.
104        labels (numpy.ndarray): Etiquetas de los clusters.
105        points (list): Lista de puntos para destacar en el gráfico.
106        x_name (str): Nombre del eje x.
107        y_name (str): Nombre del eje y.
108        title (str): Título del gráfico.
109    """
110
111    plt.figure()
112    vor = Voronoi(kmeans.cluster_centers_)
113    voronoi_plot_2d(vor, show_vertices=False, show_points=False)
114    for i in range(opt_k):
115        plt.scatter(X[labels == i, 0], X[labels == i, 1], label=f'Cluster {i}',
116                    s=10)
117    plt.scatter(kmeans.cluster_centers_[ :, 0], kmeans.cluster_centers_[ :, 1],
118                marker='x', color='k', label='Centroides')
119
120    # Adding two red points
121    for i in range(len(points)):
122        plt.scatter([points[i][0][0]], [points[i][0][1]], color=points[i][1],
123                    label=str(points[i][0]))
124
125    plt.xlabel(x_name)
126    plt.ylabel(y_name)
127    plt.title(title)
128    plt.legend()

```

```

129 plt.grid(True)
130
131 plt.xlim(min(X[:, 0]), max(X[:, 0])) # Ajusta los límites del eje x
132 plt.ylim(min(X[:, 1]), max(X[:, 1])) # Ajusta los límites del eje y
133
134 plt.show()
135
136
137 def grafica_Algoritmo(labels, X, n_clusters, title, x_title, y_title,
138                       core_samples_mask = None):
139     """
140     Función para graficar clusters encontrados por un algoritmo de clustering.
141
142     Parámetros:
143         labels (numpy.ndarray): Etiquetas de los clusters.
144         X (numpy.ndarray): Datos.
145         n_clusters (int): Número de clusters.
146         title (str): Título del gráfico.
147         x_title (str): Nombre del eje x.
148         y_title (str): Nombre del eje y.
149         core_samples_mask (numpy.ndarray): Máscara de muestras centrales (opc).
150     """
151
152     unique_labels = set(labels)
153     colors = [plt.cm.Spectral(each)
154               for each in np.linspace(0, 1, len(unique_labels))]
155
156     plt.figure()
157     for k, col in zip(unique_labels, colors):
158         if k == -1:
159             # Black used for noise.
160             col = [0, 0, 0, 1]
161
162         class_member_mask = (labels == k)
163
164         if core_samples_mask is not None:
165             xy = X[class_member_mask & core_samples_mask]
166             plt.plot(xy[:, 0], xy[:, 1], 'o', markerfacecolor=tuple(col),
167                     markeredgecolor='k', markersize=5)
168
169             xy = X[class_member_mask & ~core_samples_mask]
170             plt.plot(xy[:, 0], xy[:, 1], 'o', markerfacecolor=tuple(col),
171                     markeredgecolor='k', markersize=3)
172         else:
173             xy = X[class_member_mask]
174             plt.plot(xy[:, 0], xy[:, 1], 'o', markerfacecolor=tuple(col),
175                     markeredgecolor='k', markersize=5)
176     plt.xlabel(x_title)
177     plt.ylabel(y_title)
178     plt.title(title + str(n_clusters))
179     plt.show()
180
181
182 def calculos(metric, min_samples, X, epsilons):
183     """
184     Calcula el coeficiente de Silhouette y determina el número estimado de
185     clusters utilizando DBSCAN para diferentes
186     valores de epsilon (eps).
187
188     Parámetros:
189     metric : str
190         La métrica de distancia a utilizar. Debe ser compatible con las
191         métricas aceptadas por scikit-learn.
192     min_samples : int
193         El número mínimo de muestras requeridas para formar un cluster.
194     X : array-like, shape (n_samples, n_features)

```

```

195     La matriz de datos de entrada.
196     epsilons : array-like
197     Lista de valores de epsilon (eps) para probar.
198
199     Retorna:
200     lista : list
201     Lista de coeficientes de Silhouette para cada valor de epsilon probado.
202     labels : array-like, shape (n_samples,)
203     Etiquetas de cluster asignadas por DBSCAN.
204     core_samples_mask : array-like, shape (n_samples,)
205     Máscara de muestras centrales identificadas por DBSCAN.
206     n_clusters : int
207     Número estimado de clusters identificados por DBSCAN.
208     """
209
210     lista = []
211     for epsilon in epsilons:
212         db = DBSCAN(eps=epsilon, min_samples=min_samples, metric=metric).fit(X)
213         labels = db.labels_
214         silhouette = silhouette_score(X, labels)
215         lista.append(silhouette)
216
217     max_silhouette = max(lista)
218     optimal_epsilon = epsilons[lista.index(max_silhouette)]
219     print(f"Mayor coef. de Silhouette para DBSCAN con métrica '{metric}':",
220           max_silhouette)
221     print(f"Eps óptimo para DBSCAN con métrica '{metric}':",
222           optimal_epsilon)
223
224     db = DBSCAN(eps=optimal_epsilon, min_samples=10, metric=metric).fit(X)
225     core_samples_mask = np.zeros_like(db.labels_, dtype=bool)
226     core_samples_mask[db.core_sample_indices_] = True
227     labels = db.labels_
228     n_clusters = len(set(labels)) - (1 if -1 in labels else 0)
229     print(f"Número estimado de clusters para DBSCAN con métrica '{metric}':",
230           n_clusters)
231     return (lista, labels, core_samples_mask, n_clusters)
232
233
234
235
236
237 if __name__ == '__main__':
238
239     #Preprocesamiento de datos
240     X = preprocesamiento_datos("Personas_de_villa_laminera.txt")
241     plt.plot(X[:, 0], X[:, 1], 'ro', markersize=2)
242     plt.xlabel('Estrés')
243     plt.ylabel('Dulces')
244     plt.title('Población Villa Laminera')
245     plt.grid(True)
246     plt.show()
247
248
249     ### APARTADO 1 ###
250
251     # Calcular coeficiente de Silhouette para diferentes valores de k
252     valores_s = calculo_silhouette(X, 2, 15)
253     graficar(range(2, 16), [valores_s], 'Número de Vecindades (k)',
254             'Coeficiente de Silhouette (s)',
255             'Coeficiente de Silhouette en función de k', 'o', '-',
256             [None], True)
257
258     # Clasificar los datos con el número óptimo de vecindades
259     max_sil_kmeans = max(valores_s)
260     optimal_k = valores_s.index(max_sil_kmeans) + 2

```

```

261     print("Mayor coeficiente de Silhouette para KMeans:", max_sil_kmeans)
262     print("Número óptimo de vecindades (k) para KMeans:", optimal_k)
263
264
265     kmeans = KMeans(n_clusters=optimal_k, random_state=0).fit(X)
266     labels = kmeans.labels_
267
268     # Graficar la clasificación resultante con el diagrama de Voronoi
269     points=[]
270     grafica_voronoi(kmeans, optimal_k, X, labels, points, 'Estres', 'Dulces',
271                     f'Clasificación con {optimal_k} clusters y Voronoi')
272
273     grafica_Algoritmo(labels, X, optimal_k,
274                       'Numero de clusters usando Kmeans: ', 'Estres', 'Dulces')
275
276
277
278     ### APARTADO 2 ###:
279
280     eps = np.arange(0.1, 0.4, 0.01) # Rango de umbrales de distancia
281     min_samples = 10 # Número mínimo de elementos en una vecindad
282     metrics_list = ['euclidean', 'manhattan']
283
284     # Calculamos Silhouette para cada combinación de parámetros
285     valores_s = {}
286     for metric in metrics_list:
287         l,label, core_mask, n_cluster= calculos(metric, min_samples, X, eps)
288         valores_s[metric] = l
289         grafica_Algoritmo(label, X, n_cluster,
290                           f'Estimación de clusters para DBSCAN con metrica {metric}: ',
291                           'Estres', 'Dulces' , core_mask)
292
293     # Gráfica comparativa
294     graficar(eps, list(valores_s.values()), 'Umbral de Distancia (Eps)',
295             'Coeficiente de Silhouette (s)',
296             'Coeficiente de Silhouette en función Eps para DBSCAN',
297             None, None, metrics_list, True, True)
298
299
300
301     ### APARTADO 3 ###:
302
303     # Coordenadas de los puntos
304     a = [0.5, 0]
305     b = [0, -3]
306     punto_a = np.array([a])
307     punto_b = np.array([b])
308
309     # Predicción de los puntos utilizando el modelo KMeans
310     cluster_a = kmeans.predict(punto_a)
311     cluster_b = kmeans.predict(punto_b)
312
313     points = [(a,'r'),(b,'y')]
314     grafica_voronoi(kmeans, optimal_k, X, labels, points, 'Estres', 'Dulces',
315                     f'Clasificación con {optimal_k} clusters y Voronoi')
316
317     print(f"El cluster predicho para el punto {a} es:", cluster_a)
318     print(f"El cluster predicho para el punto {b} es:", cluster_b)

```

Programa 1: Entrega2.DiagramaVoronoiClustering.py