

Proyecto 3.1 Cluster

Una empresa de distribución de bebidas de diferentes marcas realiza la distribución y venta a tiendas pequeñas con diferentes surtidos de acuerdo solicitud del dueño del negocio. El propósito que se requiere es sugerir algunas bebidas que no le solicita para conocer su comportamiento de venta, para ello se deben seleccionar marcas que sea más adecuadas al lugar y tamaño de la Tienda.

Como primera etapa se debe realizar una segmentación de tiendas en base a la Marca de bebidas vendidas.

1) Realizar los siguientes ejercicios utilizando:

1. Con RapidMiner utilizar los algoritmos de K-means y DBScan
2. Con Python incluir el método de Elbows para determinar el número de clusters y el algoritmo de K-means
3. Con Python realizar la segmentación con DBScan
4. Se puede incluir algún algoritmo adicional que permita ver el comportamiento de los datos
5. Realizar la limpieza y Transformación que sea requerida de los datos enviados

2) Identificar el comportamiento de cada segmento y dar una explicación

Nota: Utilizar el set de datos Venta_bebidas_Demo, se recomienda primero elegir una muestra para validar los algoritmos.

Solución

2) Con Python incluir el método de Elbows para determinar el número de clusters y el algoritmo de K-means

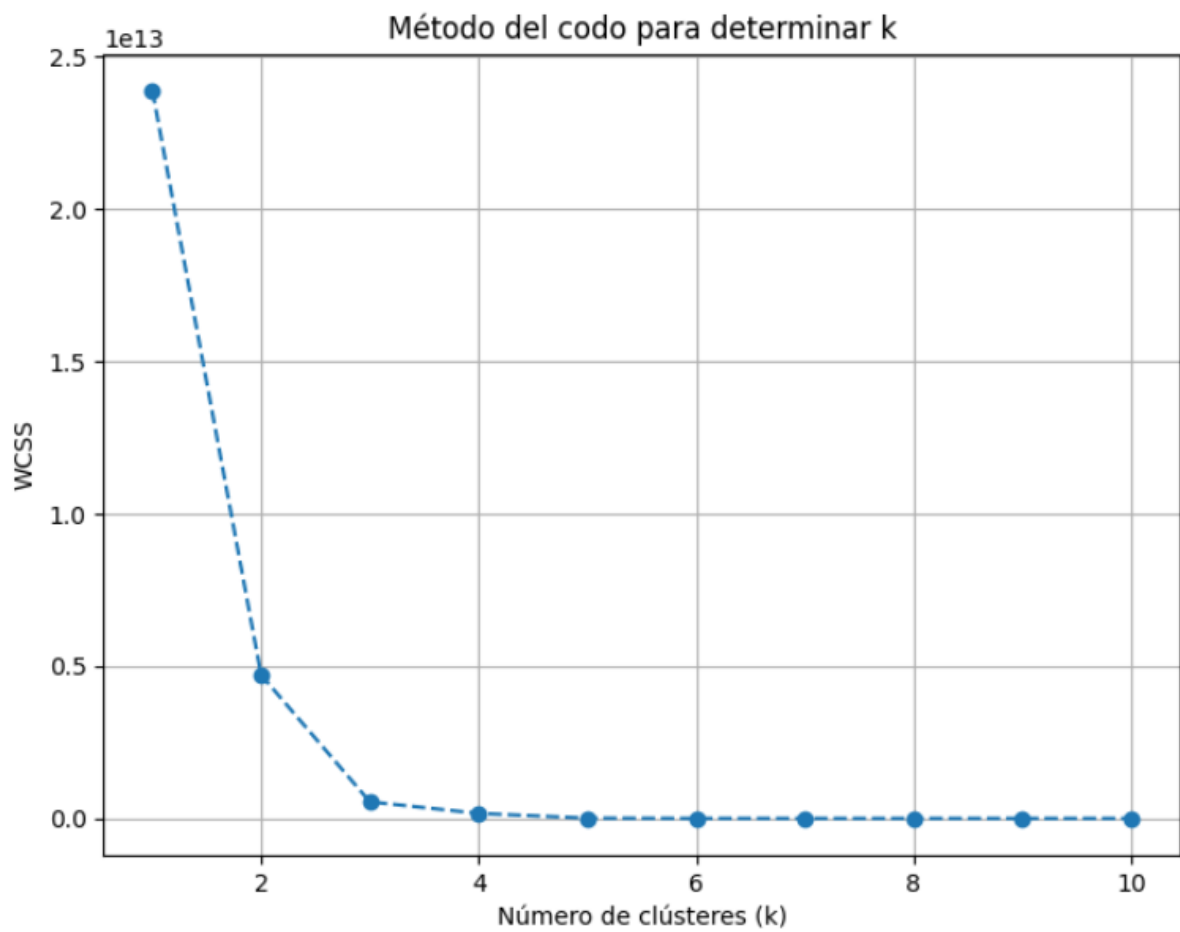
Elbows:

```
1 import pandas as pd
2 import numpy as np
3 import matplotlib.pyplot as plt
4 from sklearn.cluster import KMeans
5
6 # Cargar los datos desde el archivo CSV (ajusta el nombre del archivo según
  tu caso)
7 data = pd.read_csv("Ventas_Bebidas_Demo.csv")
8
```

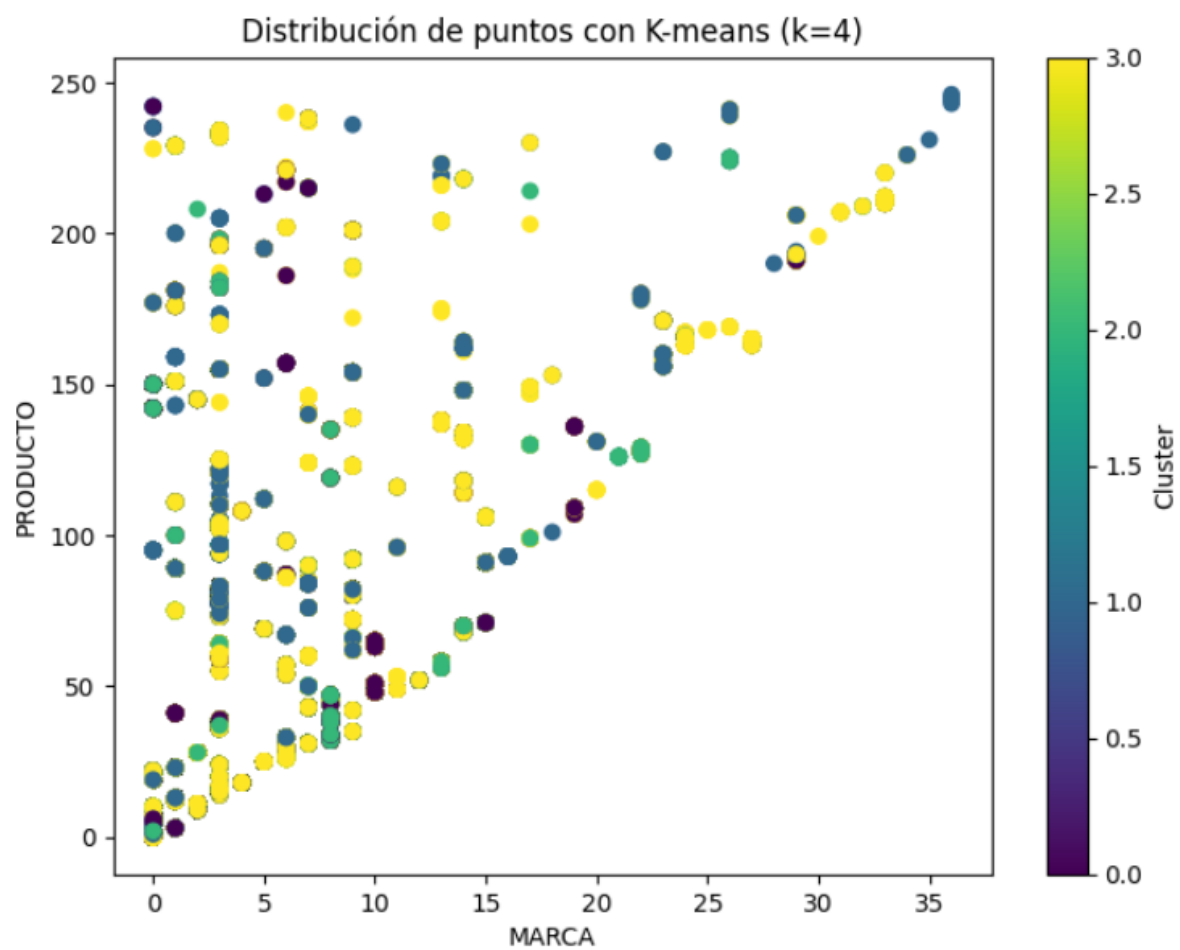
```

9 # Crear un diccionario para mapear MARCA y PRODUCTO a valores numéricos
10 marca_mapping = {marca: i for i, marca in enumerate(data["MARCA"].unique())}
11 producto_mapping = {producto: i for i, producto in
12     enumerate(data["PRODUCTO"].unique())}
13
14 # Reemplazar los valores categóricos con los valores numéricos
15 data["MARCA"] = data["MARCA"].map(marca_mapping)
16 data["PRODUCTO"] = data["PRODUCTO"].map(producto_mapping)
17
18 # Seleccionar las variables relevantes (MARCA, PRODUCTO, TIENDA)
19 selected_data = data[["MARCA", "PRODUCTO", "TIENDA"]]
20
21 # Determinar el número óptimo de clústeres utilizando la regla del codo
22 wcss = []
23 for k in range(1, 11):
24     kmeans = KMeans(n_clusters=k, init='k-means++', random_state=42)
25     kmeans.fit(selected_data)
26     wcss.append(kmeans.inertia_)
27
28 # Graficar la curva del codo
29 plt.figure(figsize=(8, 6))
30 plt.plot(range(1, 11), wcss, marker='o', linestyle='--')
31 plt.xlabel('Número de clústeres (k)')
32 plt.ylabel('WCSS')
33 plt.title('Método del codo para determinar k')
34 plt.grid()
35 plt.show()
36
37 # Seleccionar el valor óptimo de k basado en la gráfica
38 optimal_k = 3 # Supongamos que el codo está en k=3
39 kmeans_final = KMeans(n_clusters=optimal_k, init='k-means++',
40     random_state=42)
41 selected_data['cluster'] = kmeans_final.fit_predict(selected_data)

```



El optimo podría ser 3 ó 4, ya que es donde se observa el codo.

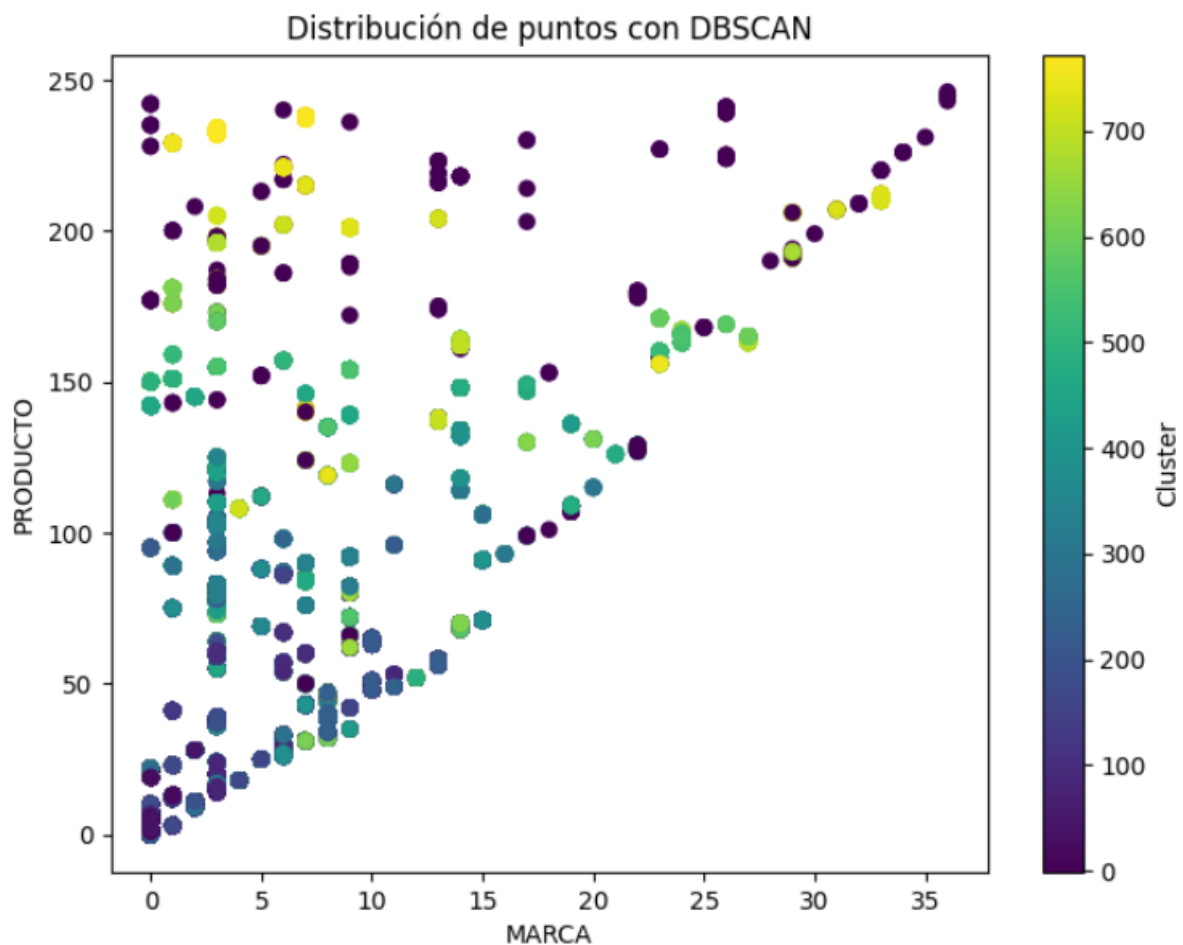


Dbscan:

```

1 import pandas as pd
2 import numpy as np
3 import matplotlib.pyplot as plt
4 from sklearn.cluster import DBSCAN
5
6 # Cargar los datos desde el archivo CSV (ajusta el nombre del archivo según
  tu caso)
7 data = pd.read_csv("Ventas_Bebidas_Demo.csv")
8
9 # Crear un diccionario para mapear MARCA y PRODUCTO a valores numéricos
10 marca_mapping = {marca: i for i, marca in enumerate(data["MARCA"].unique())}
11 producto_mapping = {producto: i for i, producto in
  enumerate(data["PRODUCTO"].unique())}
12
13 # Reemplazar los valores categóricos con los valores numéricos
14 data["MARCA"] = data["MARCA"].map(marca_mapping)
15 data["PRODUCTO"] = data["PRODUCTO"].map(producto_mapping)
16
17 # Seleccionar las variables relevantes (MARCA, PRODUCTO, TIENDA)
18 selected_data = data[["MARCA", "PRODUCTO", "TIENDA"]]
19
20 # Aplicar DBSCAN con parámetros eps=0.5 y min_samples=5
21 dbscan = DBSCAN(eps=0.5, min_samples=5)
22 selected_data['cluster'] = dbscan.fit_predict(selected_data)
23
24 # Visualizar la distribución de puntos
25 plt.figure(figsize=(8, 6))
26 plt.scatter(selected_data["MARCA"], selected_data["PRODUCTO"],
  c=selected_data["cluster"], cmap='viridis')
27 plt.xlabel('MARCA')
28 plt.ylabel('PRODUCTO')
29 plt.title('Distribución de puntos con DBSCAN')
30 plt.colorbar(label='Cluster')
31 plt.show()
32

```



Algoritmo de Agrupamiento Jerárquico (Hierarchical Agglomerative Clustering - HAC)

```

1  import pandas as pd
2  import numpy as np
3  import matplotlib.pyplot as plt
4  from sklearn.cluster import AgglomerativeClustering
5
6  # Cargar los datos desde el archivo CSV (ajusta el nombre del archivo según
   tu caso)
7  data = pd.read_csv("Ventas_Bebidas_Demo.csv")
8
9  # Crear un diccionario para mapear MARCA y PRODUCTO a valores numéricos
10 marca_mapping = {marca: i for i, marca in enumerate(data["MARCA"].unique())}
11 producto_mapping = {producto: i for i, producto in
   enumerate(data["PRODUCTO"].unique())}
12
13 # Reemplazar los valores categóricos con los valores numéricos
14 data["MARCA"] = data["MARCA"].map(marca_mapping)
15 data["PRODUCTO"] = data["PRODUCTO"].map(producto_mapping)
16
17 # Seleccionar las variables relevantes (MARCA, PRODUCTO, TIENDA)
18 selected_data = data[["MARCA", "PRODUCTO", "TIENDA"]]
19
20 # Aplicar HAC con parámetros n_clusters=4 y linkage='ward'
21 hac = AgglomerativeClustering(n_clusters=4, linkage='ward')
22 selected_data['cluster'] = hac.fit_predict(selected_data)
23

```

```

24 # Visualizar la distribución de puntos
25 plt.figure(figsize=(8, 6))
26 plt.scatter(selected_data["MARCA"], selected_data["PRODUCTO"],
27             c=selected_data["cluster"], cmap='viridis')
28 plt.xlabel('MARCA')
29 plt.ylabel('PRODUCTO')
30 plt.title('Distribución de puntos con HAC')
31 plt.colorbar(label='Cluster')
32 plt.show()

```

