

Universidad Nacional Autónoma de México



Facultad de ingeniería

Materia: Minería de Datos

Equipo: 10

Proyecto 3

Alumnos:

- Beristain Carbajal Armando
- Barrera Peña Víctor Miguel
- Torrecilla Jimenez Aaron Israel

Fecha: 31/Marzo/2024



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.

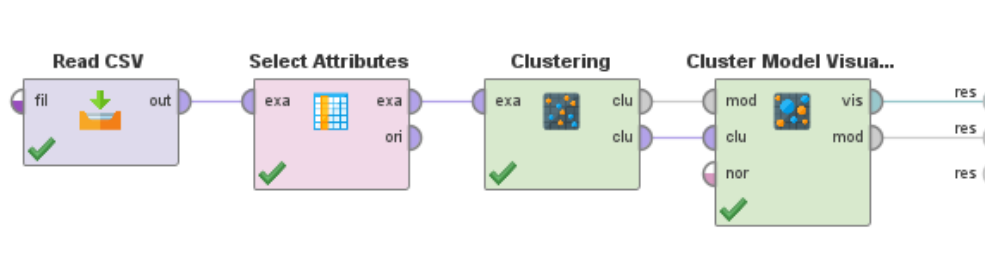
Como primera etapa se debe realizar una segmentación de tiendas en base a la Marca de bebidas vendidas.

1) Realizar los siguientes ejercicios utilizando:

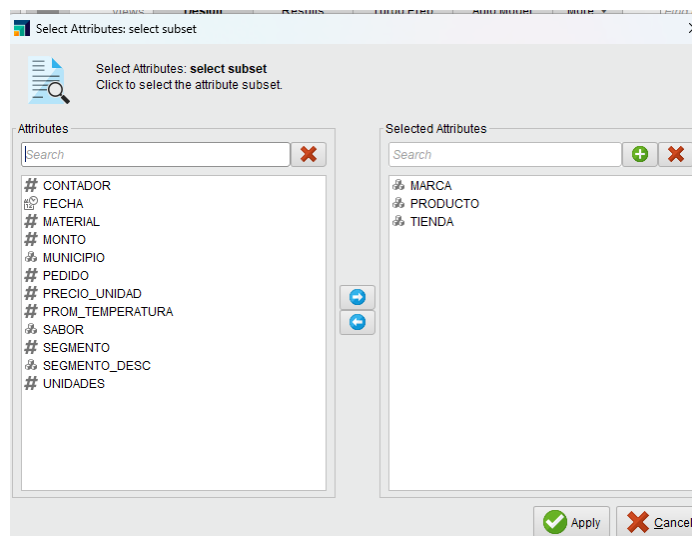
A) Con RapidMiner utilizar los algoritmos de K-means y DBScan

Algoritmo de Kmeans

Diseño de Diagrama



Hacemos la selección de atributos para la segmentación



Para el algoritmo hacemos la selección de 3 clusters

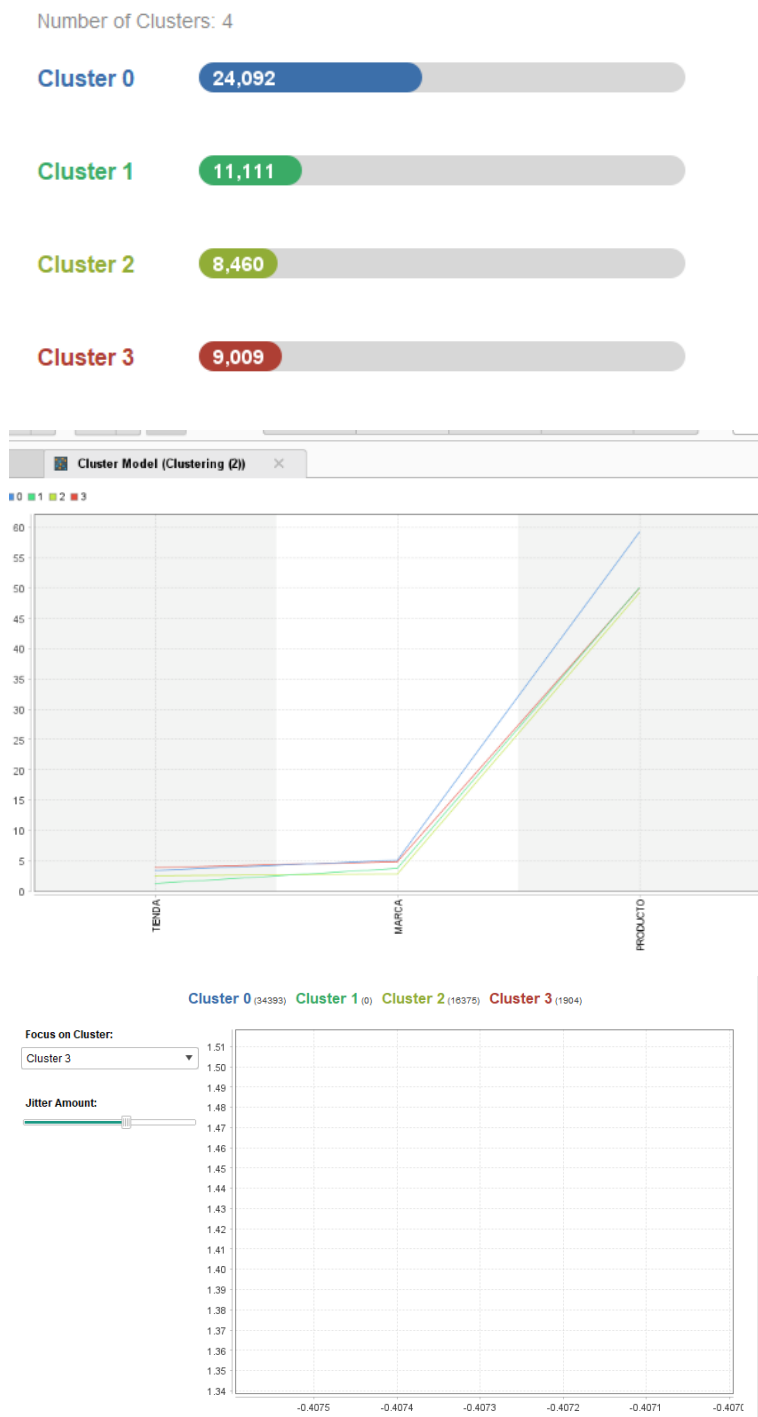
☒ add cluster attribute

☐ add as label

☐ remove unlabeled

k

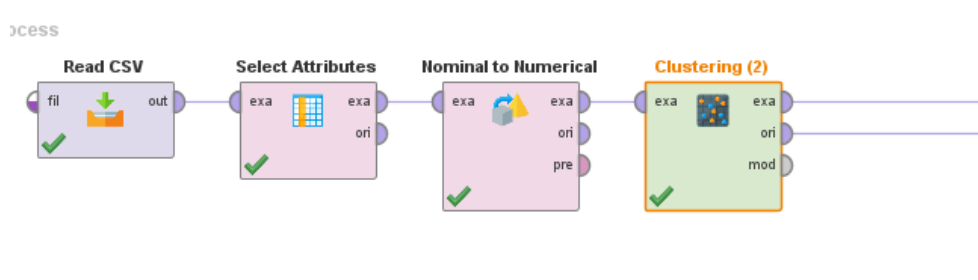
Resultados del modelo



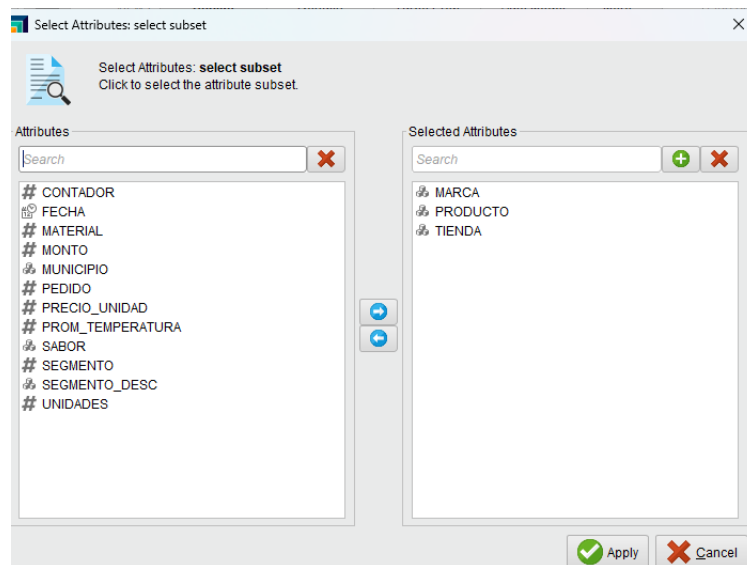
Aunque se colocó en escala automática y se moviera hasta los valores que deberían aparecer en la gráfica no se mostraba lo que se requería, y ello solo paso en la graficación, no se sabe porque, sin embargo la limpieza y acomodo de los datos es correcta.

Algoritmo de DBSCAN:

Diseño del Diagrama



Se hace la selección de atributos para la segmentación



Open in [Turbo Prep](#) [Auto Model](#) [Interactive Analysis](#) Filter (52,672 / 52,672 examples): [all](#)

show the data in a table	score(clust...	score(clust...	score(clust...	score(clust...	score(clust...	score(clust...	score(clust...
1	0	1.414	1.414	2	1.414	1.414	1.414
2	1.414	0	1.414	2	1.414	1.414	1.414
3	1.414	1.414	0	2	1.414	1.414	1.414
4	2	2	2	0	2	2	2
5	1.414	1.414	1.414	2	0	1.414	1.414
6	1.414	1.414	0	2	1.414	1.414	1.414
7	1.414	1.414	1.414	2	1.414	0	1.414
8	1.414	1.414	1.414	2	1.414	1.414	0
9	2	2	2	2.449	2	1.414	2
10	2	1.414	2	2.449	2	2	2
11	2	2	2	2.449	2	2	1.414
12	2	2	2	2.449	2	2	2
13	2	2	2	2.449	2	2	2

ExampleSet (52,672 examples,596 special attributes,290 regular attributes)

PYTHON:

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

En la columna pedido existe un error en la línea 21876, tenemos la combinación de dos valores de diferente formato

21875	1000157892	100000740
21876	1000febrero	100005084
21877	1000febrero	100005084
21878	1000febrero	100005084
21879	1000febrero	100005084
21880	1000febrero	100005084
21881	1000febrero	100005084
21882	1000febrero	100005084
21883	1000febrero	100005084
21884	1000febrero	100005084
21885	1000febrero	100005084
21886	1000221720	100007101

Al mapearlos, se elimina esta inconsistencia y se numera como un valor más, sólo que cuando se realice un análisis es necesario tener una tabla de relación entre el n

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

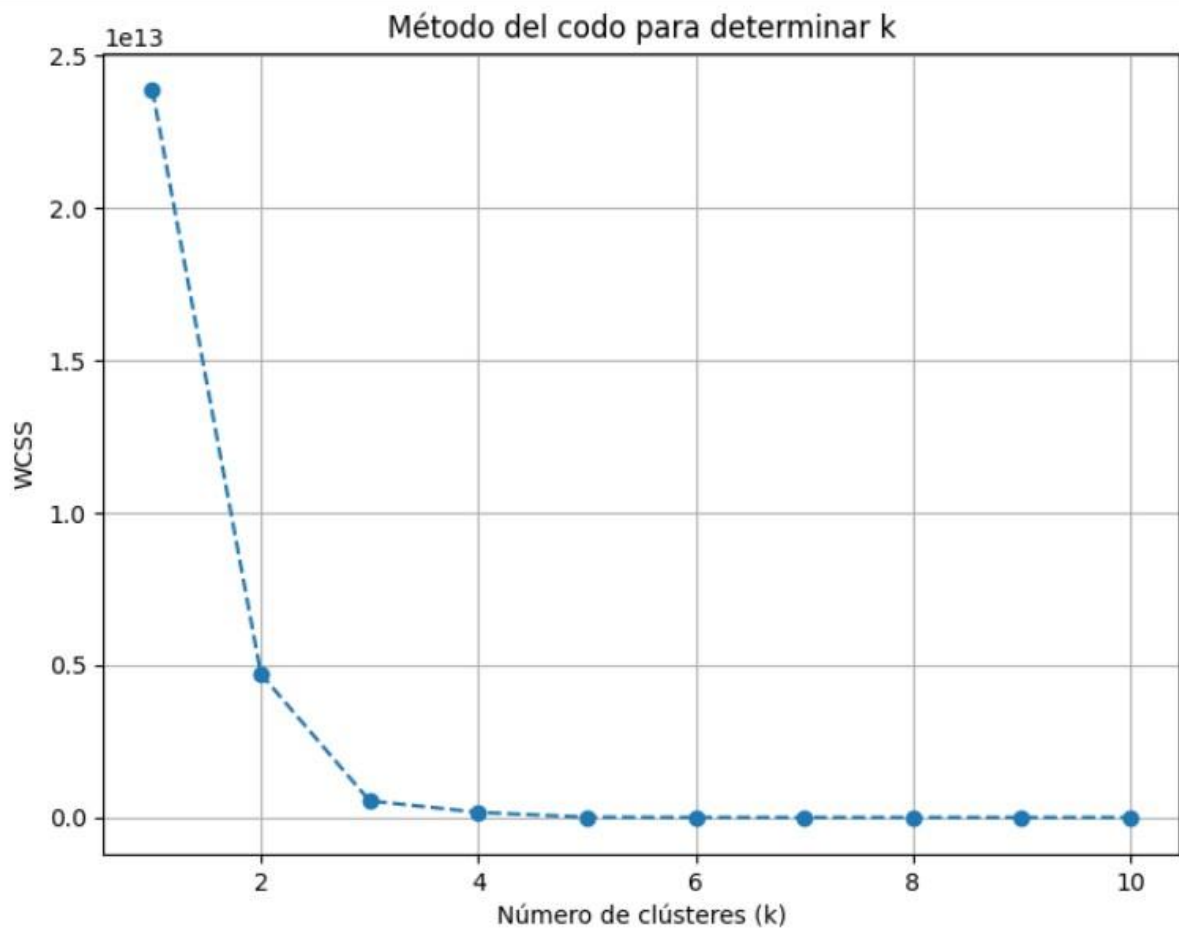
Nota: Utilizar el set de datos Venta_bebidas_Demo, se recomienda primero elegir una muestra para validar los algoritmos.

Solución

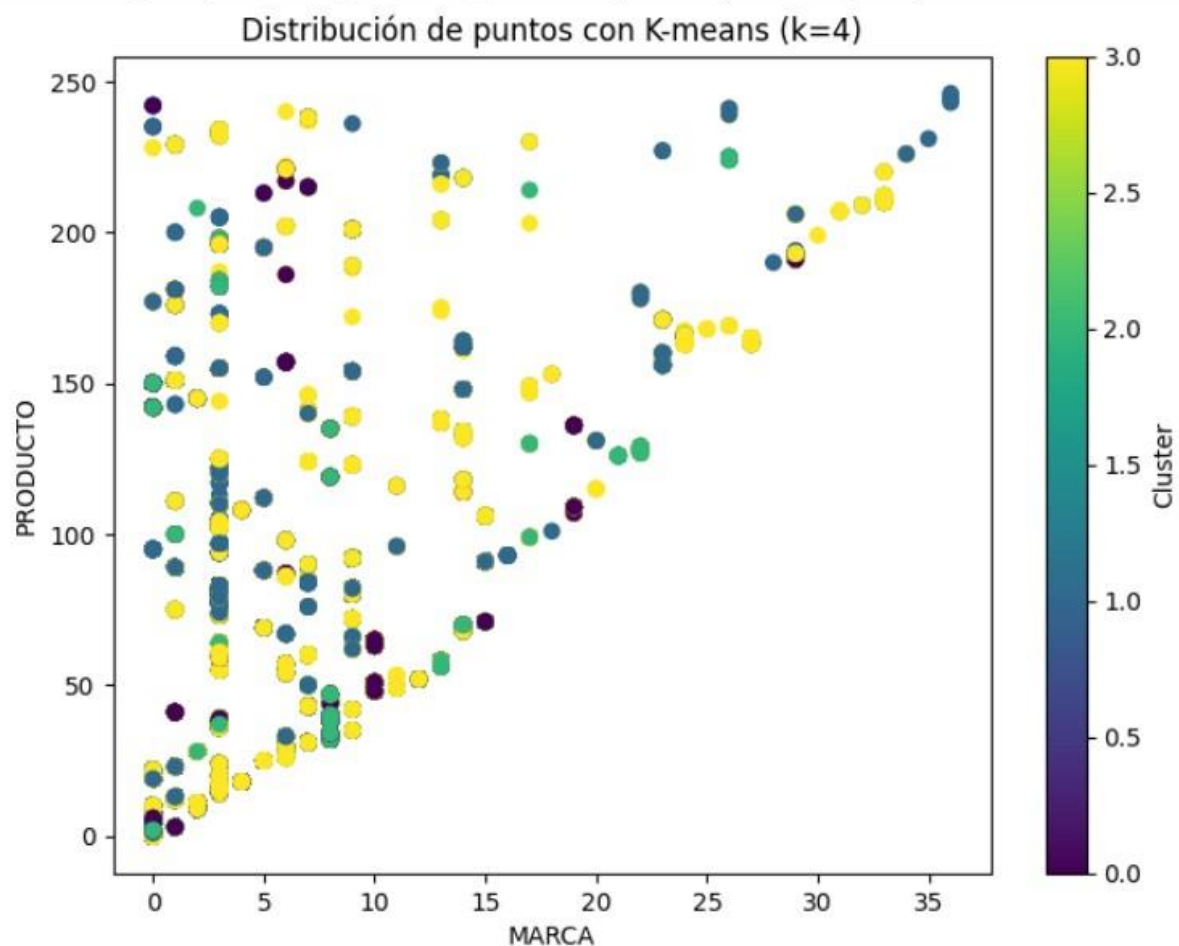
2) Con Python incluir el método de Elbows para determinar el número de clusters y el algoritmo de K-means

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
  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 # Determinar el número óptimo de clústeres utilizando la regla del codo
21 wcss = []
22 for k in range(1, 11):
23     kmeans = KMeans(n_clusters=k, init='k-means++', random_state=42)
24     kmeans.fit(selected_data)
25     wcss.append(kmeans.inertia_)
26
27 # Graficar la curva del codo
28 plt.figure(figsize=(8, 6))
29 plt.plot(range(1, 11), wcss, marker='o', linestyle='--')
30 plt.xlabel('Número de clústeres (k)')
31 plt.ylabel('WCSS')
32 plt.title('Método del codo para determinar k')
33 plt.grid()
34 plt.show()
35
36 # seleccionar el valor óptimo de k basado en la gráfica
37 optimal_k = 3 # Supongamos que el codo está en k=3
38 kmeans_final = KMeans(n_clusters=optimal_k, init='k-means++',
  random_state=42)
39 selected_data['cluster'] = kmeans_final.fit_predict(selected_data)
```



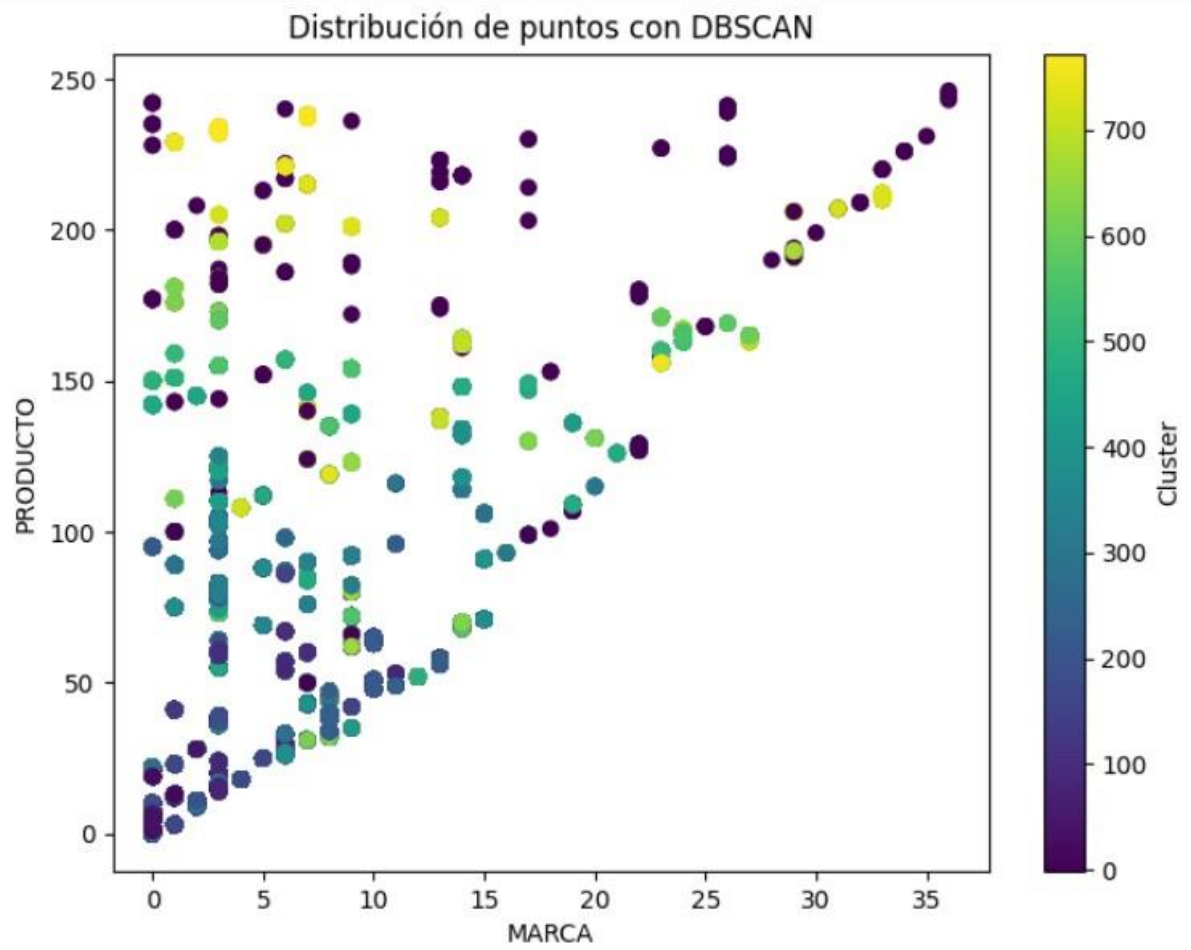
El óptimo podría ser 3 ó 4, ya que es donde se observa el codo.



DBSCAN:

La implementación se logrará mediante el siguiente código:

```
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
```

El resultado es similar a DBSCAN, por ello se considera correcto, que signa las mismas tendencias en ambos algoritmos.

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
```

Muestra un poco más de dispersión este último algoritmo, sin embargo, se mantiene la misma tendencia, esto fue retomado con la implementación de la biblioteca, por ello lo único que se tuvo que arreglar, fue arreglar los datos para que funcionen los algoritmos, en este caso funciona igual que K-means. Las transformaciones realizadas están vinculadas a conversión numérica de Marca y Producto haciendo un mapeo, para todos los algoritmos era necesario de otra manera se ocasionaba un error.

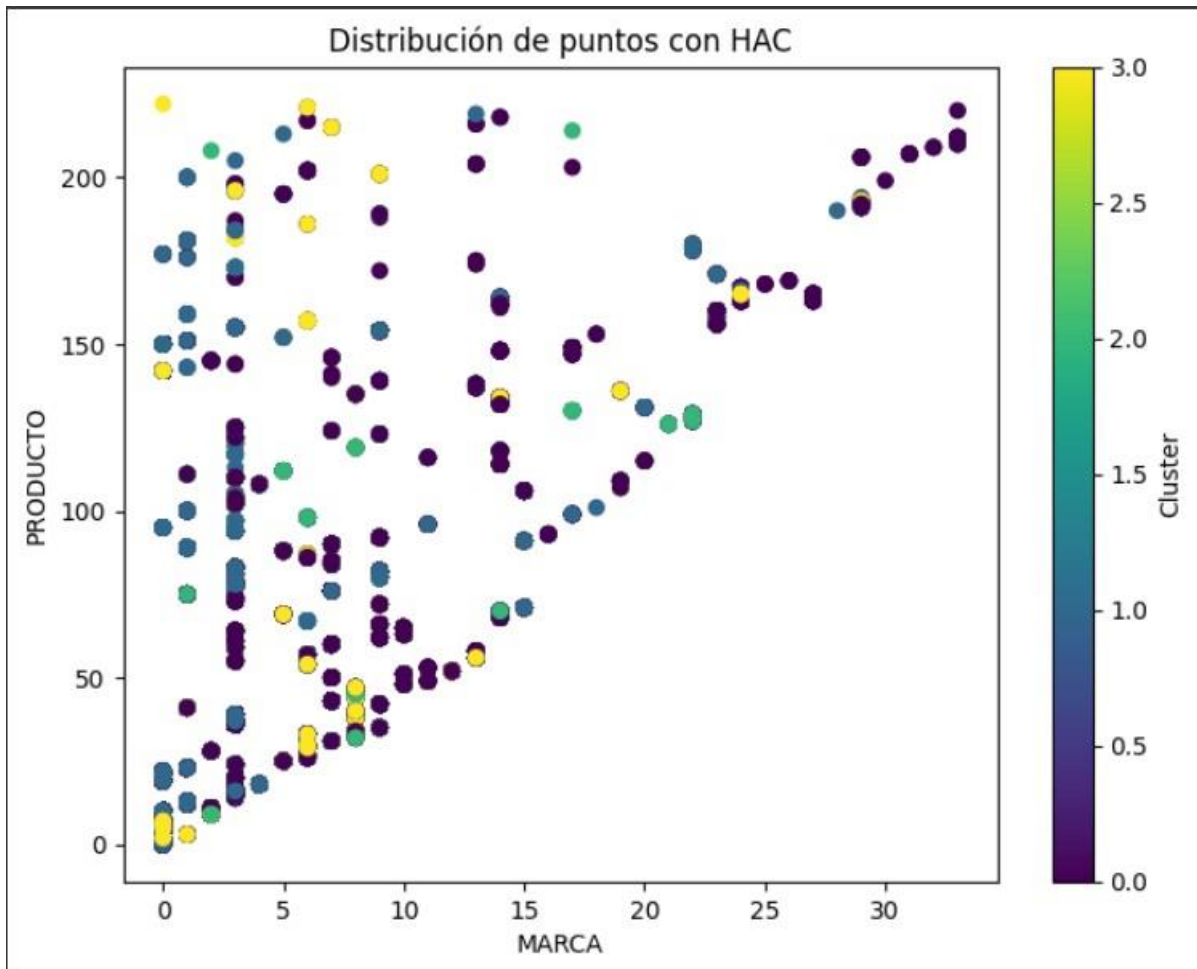
Limpiezas realizadas, ninguna, transformaciones realizadas:

- Mapeo Marca a número entero.
- Producto mapeado a número entero.

```

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()

```



El gráfico sigue la tendencia de los anteriores gráficos usando algoritmos ya conocidos, por lo cual este también se considera que está correcto, la diferencia notable de este algoritmo es que no se le tiene que especificar el valor de K, ya que no lo necesita.

Conclusiones

- **Barrera Peña Víctor Miguel:** Se cumplió las actividades encomendadas, no se tuvo algún error al realizarla, las comprobaciones se lograron al determinar que mediante las implementaciones de los algoritmos el resultado era idéntico, sólo se tuvo que investigar cuando se obtenía alguna clase de error, la búsqueda de
- **Torrecilla Jimenez Aaron Israel:** Las actividades realizadas en el proyecto de Minería de Datos cumplen con los objetivos establecidos. Se implementaron con éxito los algoritmos de K-means y DBScan, tanto en RapidMiner como en Python, con resultados consistentes entre ambos enfoques. La determinación del número óptimo de clusters se realizó mediante el método del codo, sugiriendo 3 o 4 clusters para el conjunto de datos. Se compararon los resultados de DBSCAN y K-means, encontrando similitudes en las tendencias identificadas. Se exploró también el algoritmo de Agrupamiento Jerárquico (HAC), observando resultados similares a K-means, aunque con una ligera dispersión mayor. Se llevaron a cabo transformaciones de datos, como el mapeo de marcas y productos a números enteros, para garantizar la compatibilidad con los algoritmos utilizados. El proyecto proporciona una sólida base para tomar decisiones en la distribución y venta de bebidas por parte de la empresa.
- **Beristain Carbajal Armando:** Mediante el caso practico de la venta de bebidas de un conjunto de tiendas logramos aplicar con claridad los diferentes algoritmos abordados en las horas teóricas , mediante el uso de K-Means y DBScan segmentamos el tipo de tiendas con ciertas necesidades y demandas deacuerdo a sus condiciones y características donde comercian , tras explorar los algoritmos en Python y Rapidminer , llegamos a las mismas conclusiones y resultados