
```

1 import numpy as np
2 from matplotlib import pyplot as plt
3 plt.rcParams['figure.figsize'] = (12, 6)
4 plt.style.use('ggplot')
5 np.linalg.norm(np.array([8,10]) - np.array([3,10.5]))

```

```
5.024937810560445
```

```
1 np.sqrt((8-3)**2+(10-10.5)**2)
```

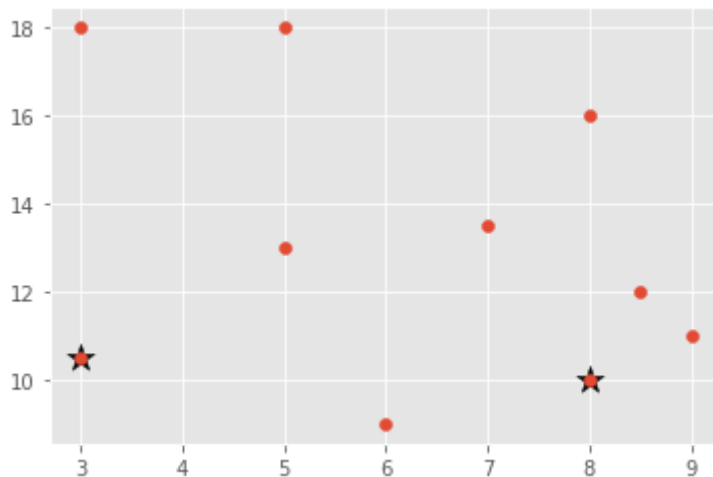
```
5.024937810560445
```

```

1 data = np.array([[8,10],[3,10.5],[7,13.5],[5,18],[5,13],[6,9],[9,11],[3,18],[8.5,
2 C = np.array([[8,10],[3,10.5]])
3
4 # Gráfica
5 fig = plt.figure()
6 plt.scatter(C[0][0], C[0][1], marker='*', s=200, c='#050505')
7 plt.scatter(C[1][0], C[1][1], marker='*', s=200, c='#050505')
8 plt.scatter(data[:, 0], data[:, 1])
9

```

 <matplotlib.collections.PathCollection at 0x7f2f39789dd0>



```

1 distances = []
2 clusters = np.zeros(len(data))
3
4 def dist(a, b, ax=1):
5     return np.linalg.norm(a - b, axis=ax)
6
7 for i in range(len(data)):
8     distance = dist(data[i], C)
9     distances.append(distance)
10    cluster = np.argmin(distances)
11    clusters[i] = cluster

```

```

12
13 print(clusters)
14 distances
15

[0. 1. 0. 1. 1. 0. 0. 1. 0. 0.]
[array([0.          , 5.02493781]),
 array([5.02493781, 0.          ]),
 array([3.64005494, 5.          ]),
 array([8.54400375, 7.76208735]),
 array([4.24264069, 3.20156212]),
 array([2.23606798, 3.35410197]),
 array([1.41421356, 6.02079729]),
 array([9.43398113, 7.5          ]),
 array([2.06155281, 5.70087713]),
 array([6.          , 7.43303437])]

1 for i in range(len(C)):
2     points = [data[j] for j in range(len(data)) if clusters[j] == i]
3     C[i] = np.mean(points, axis=0)
4
5 print(C)

[[ 7.75      11.91666667]
 [ 4.        14.875      ]]

1 # Gráfica
2 fig = plt.figure()
3 plt.scatter(C[0][0], C[0][1], marker='*', s=200, c='#050505')
4 plt.scatter(C[1][0], C[1][1], marker='*', s=200, c='#050505')
5 plt.scatter(data[:, 0], data[:, 1])
6

```

<matplotlib.collections.PathCollection at 0x7f2f39707a10>



```

1 '''
2 =====
3 scikit-learn
4

```

```

4 =====
5 '''
6 import numpy as np
7 import pandas as pd
8 from matplotlib import pyplot as plt
9 plt.rcParams['figure.figsize'] = (12, 6)
10 plt.style.use('ggplot')

```

```

1 # Conjunto de datos xclara
2 data = pd.read_csv('http://dicyg.fi-c.unam.mx:8080/lalo/pypcd/presentaciones/xcla
3 print(data.shape)
4 data.head()

```

```
(3000, 2)
```

	v1	v2
0	2.072345	-3.241693
1	17.936710	15.784810
2	1.083576	7.319176
3	11.120670	14.406780
4	23.711550	2.557729

```

1 # Gráfica
2 f1 = data['V1'].values
3 f2 = data['V2'].values
4 X = np.array(list(zip(f1, f2)))
5 plt.scatter(f1, f2, c='black', s=7)

```

```
<matplotlib.collections.PathCollection at 0x7f2f321e9610>
```



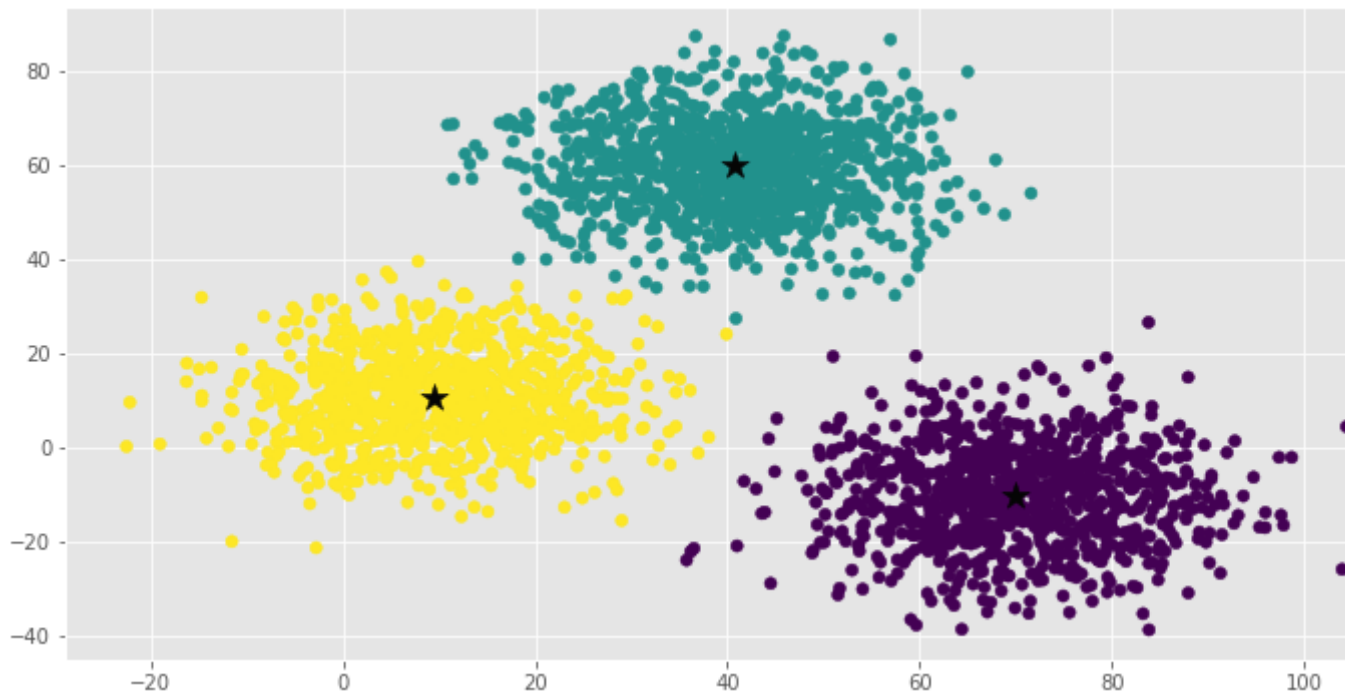
```

1
2 from sklearn.cluster import KMeans
3
4 # Número de grupos
5 kmeans = KMeans(n_clusters=3)
6 # Ajuste
7 kmeans = kmeans.fit(X)
8 # Etiquetas de cada clase
9 y = kmeans.predict(X)
10 # Centroides
11 C_skl = kmeans.cluster_centers_

1 fig, ax = plt.subplots()
2 ax.scatter(X[:, 0], X[:, 1], c=y)
3 ax.scatter(C_skl[:, 0], C_skl[:, 1], marker='*', s=200, c='#050505')

```

<matplotlib.collections.PathCollection at 0x7f2f19e231d0>



```

1 # ¿Cómo saber el valor inicial para K?

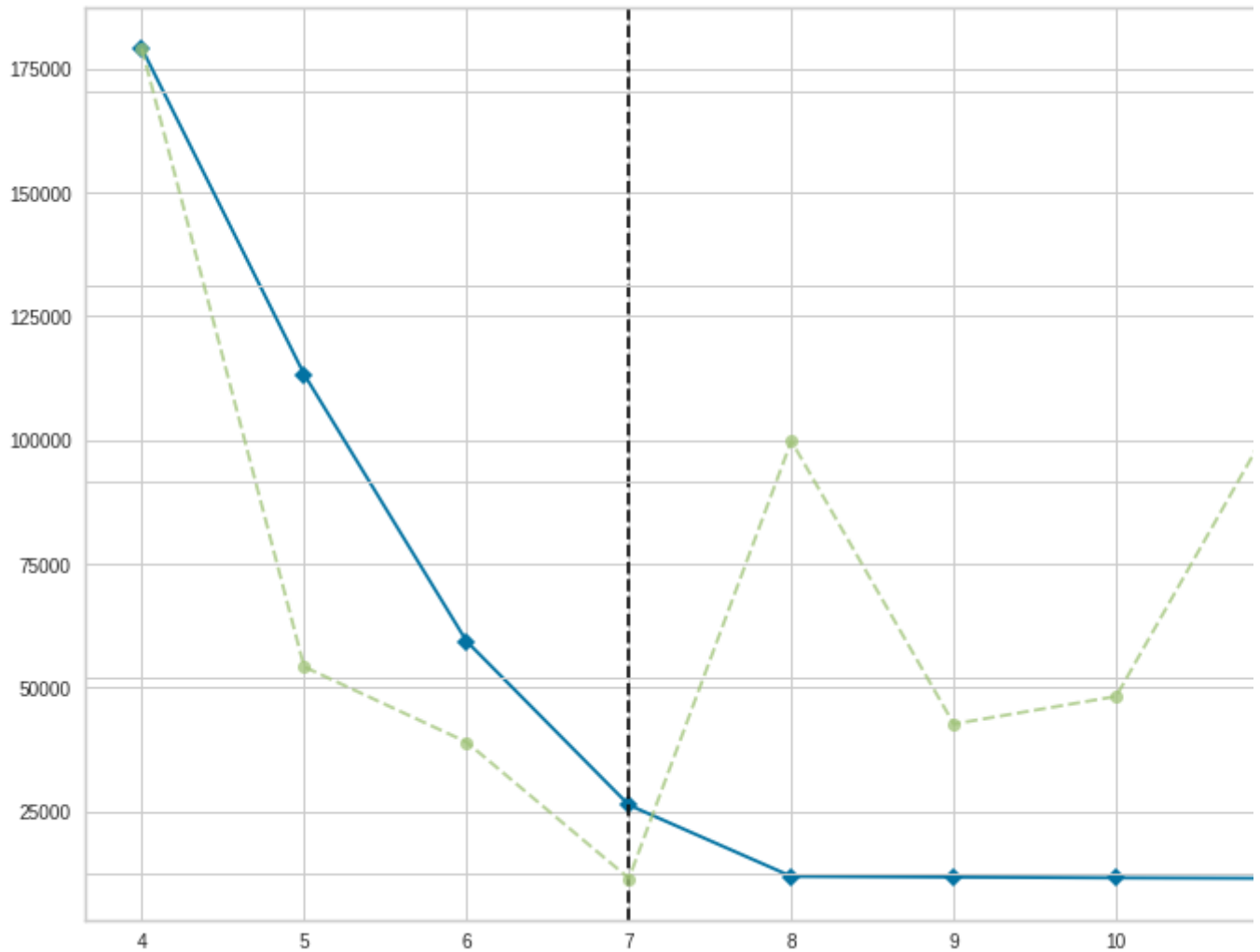
1 # Elección de k con la gráfica de codo KElbowVisualizer
2 # https://www.scikit-yb.org/en/latest/api/cluster/elbow.html
3 # By default, the scoring parameter metric is set to distortion,
4 # which computes the sum of squared distances from each point to its assigned cen
5 from sklearn.cluster import KMeans
6 from sklearn.datasets import make_blobs
7 from yellowbrick.cluster import KElbowVisualizer
8
9 X, y = make_blobs(n_samples=1000, n_features=12, centers=8, random_state=42)

```

```

10
11 visualizer = KElbowVisualizer(KMeans(), k=(4,12), timings=True)
12 visualizer.fit(X)
13 plt.xlabel('Número de grupos')
14 plt.ylabel('Distorsión')
15 plt.show()

```



```

1 # KMeans in depth
2 # https://github.com/jakevdp/PythonDataScienceHandbook/blob/master/notebooks/05.1

```

1

1

1

1

1

```

1 # bibliotecas
2 import pandas as pd
3 import numpy as np
4 import matplotlib.pyplot as plt

1 # datos
2 #https://archive.ics.uci.edu/ml/machine-learning-databases/00292/Wholesale%20cust
3 url = 'https://bit.ly/2COHM14'
4 data = pd.read_csv(url)
5 data.head()
6

```

	Channel	Region	Fresh	Milk	Grocery	Frozen	Detergents_Paper	Delicassen
0	2	3	12669	9656	7561	214	2674	1338
1	2	3	7057	9810	9568	1762	3293	1776
2	2	3	6353	8808	7684	2405	3516	7844
3	1	3	13265	1196	4221	6404	507	1788
4	2	3	22615	5410	7198	3915	1777	5185

```

1 from sklearn.preprocessing import normalize
2 data_scaled = normalize(data)
3 data_scaled = pd.DataFrame(data_scaled, columns=data.columns)
4 data_scaled.head()

```

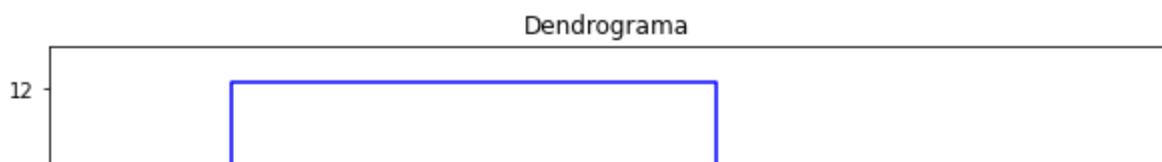
	Channel	Region	Fresh	Milk	Grocery	Frozen	Detergents_Paper	Delicassen
0	0.000112	0.000168	0.708333	0.539874	0.422741	0.011965	0.149505	0.074
1	0.000125	0.000188	0.442198	0.614704	0.599540	0.110409	0.206342	0.111
2	0.000125	0.000187	0.396552	0.549792	0.479632	0.150119	0.219467	0.489
3	0.000065	0.000194	0.856837	0.077254	0.272650	0.413659	0.032749	0.115
4	0.000079	0.000119	0.895416	0.214203	0.284997	0.155010	0.070358	0.205

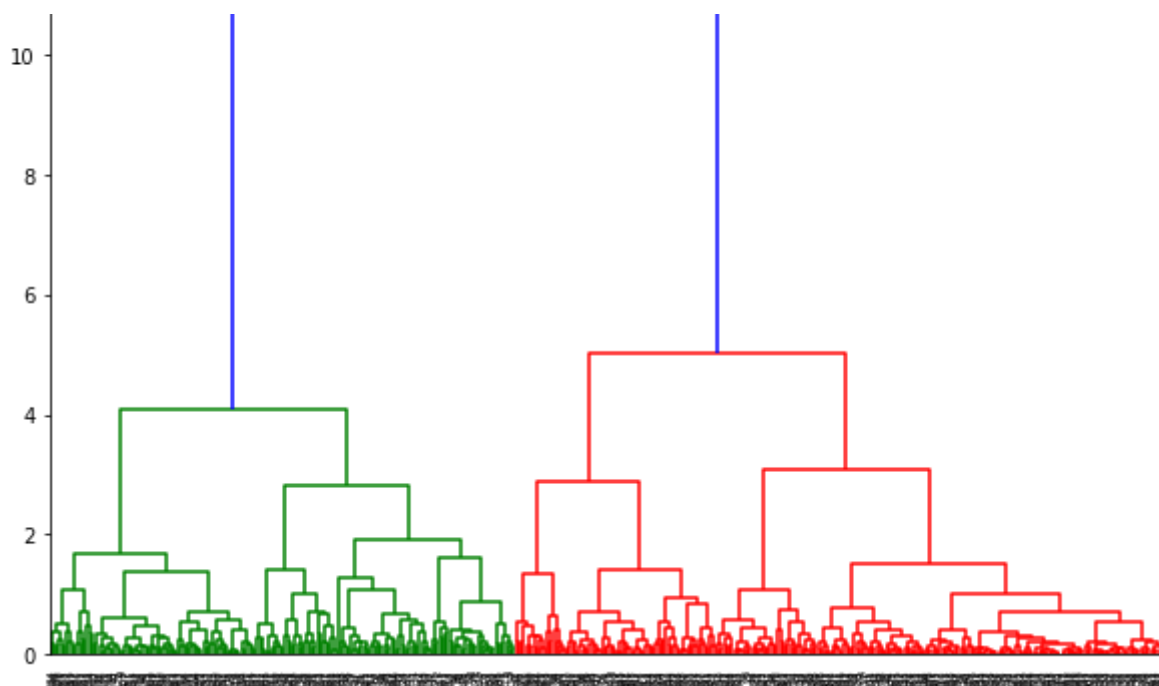
```

1 import scipy.cluster.hierarchy as shc
2 plt.figure(figsize=(10, 7))

3 plt.title("Dendrograma")
4 dend = shc.dendrogram(shc.linkage(data_scaled, method='ward'))

```



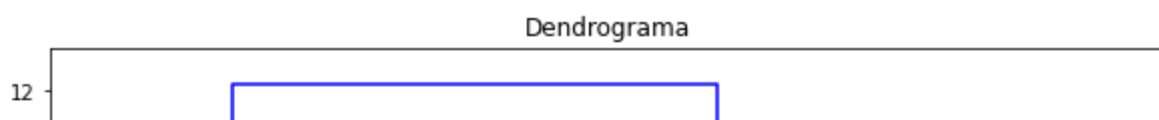


```

1 plt.figure(figsize=(10, 7))
2 plt.title("Dendrograma")
3 dend = shc.dendrogram(shc.linkage(data_scaled, method='ward'))
4 plt.axhline(y=6, color='r', linestyle='--')

```

<matplotlib.lines.Line2D at 0x7f4513572c90>



```

1 from sklearn.cluster import AgglomerativeClustering
2 cluster = AgglomerativeClustering(n_clusters=2, affinity='euclidean',

```

```

3                                     linkage='ward')
4 cluster.fit_predict(data_scaled)

array([1, 1, 1, 0, 0, 1, 0, 1, 1, 1, 1, 0, 0, 0, 0, 1, 1, 0, 0, 0, 0,
       0, 1, 0, 0, 0, 0, 1, 0, 0, 1, 0, 0, 1, 1, 0, 1, 1, 0, 0, 0, 1, 1,
       1, 1, 1, 1, 1, 1, 0, 1, 0, 1, 0, 1, 1, 1, 0, 1, 0, 1, 1, 1, 0, 1,
       1, 0, 1, 0, 0, 0, 0, 0, 1, 0, 0, 1, 0, 1, 0, 1, 1, 0, 0, 1, 1, 0,
       0, 0, 0, 1, 0, 1, 1, 1, 0, 0, 0, 1, 1, 1, 0, 0, 0, 1, 1, 1, 1,
       0, 1, 0, 0, 0, 0, 0, 1, 0, 0, 0, 0, 0, 1, 0, 0, 0, 0, 1, 0, 0, 0,
       0, 0, 0, 0, 1, 1, 0, 1, 0, 0, 0, 0, 0, 1, 0, 0, 0, 0, 0, 1, 0, 1,
       0, 1, 1, 0, 1, 1, 1, 0, 0, 1, 1, 1, 1, 1, 0, 0, 1, 1, 1, 1, 1, 1,
       0, 0, 0, 1, 0, 0, 1, 1, 1, 0, 0, 1, 1, 1, 0, 0, 0, 1, 0, 0, 0, 1,
       0, 0, 1, 1, 0, 1, 1, 1, 0, 1, 1, 1, 0, 1, 0, 1, 1, 1, 0, 1, 0,
       0, 1, 0, 0, 0, 0, 0, 0, 1, 0, 0, 1, 0, 1, 0, 1, 0, 0, 0, 0, 0,
       0, 0, 1, 1, 0, 0, 0, 0, 0, 1, 0, 0, 1, 0, 0, 0, 0, 0, 0, 0, 1,
       1, 1, 1, 0, 1, 0, 0, 1, 1, 0, 1, 1, 0, 0, 0, 0, 1, 0, 0, 0, 0,
       0, 0, 0, 0, 1, 0, 1, 1, 0, 0, 0, 0, 1, 1, 0, 1, 1, 1, 1, 1, 0,
       0, 1, 0, 0, 1, 0, 0, 1, 0, 1, 0, 1, 1, 0, 0, 0, 0, 0, 1, 0, 0,
       0, 1, 0, 1, 0, 0, 0, 0, 0, 0, 1, 1, 1, 1, 0, 1, 1, 0, 1, 1,
       1, 1, 0, 1, 0, 1, 1, 1, 0, 0, 1, 0, 0, 1, 0, 0, 0, 0, 0, 1, 0,
       0, 0, 1, 0, 1, 1, 0, 0, 0, 0, 1, 0, 1, 0, 0, 0, 0, 1, 1, 0, 0,
       1, 0, 0, 0, 0, 0, 0, 0, 0, 0, 1, 0, 0, 0, 1, 1, 0, 0, 1, 1, 1,
       1, 0, 1, 0, 0, 0, 0, 0, 1, 0, 1, 0, 1, 0, 0, 1, 0, 0, 0, 1, 0, 1])

1 plt.figure(figsize=(10, 7))
2 plt.scatter(data_scaled['Milk'], data_scaled['Grocery'], c=cluster.labels_)

```

<matplotlib.collections.PathCollection at 0x7f451ba13550>



```

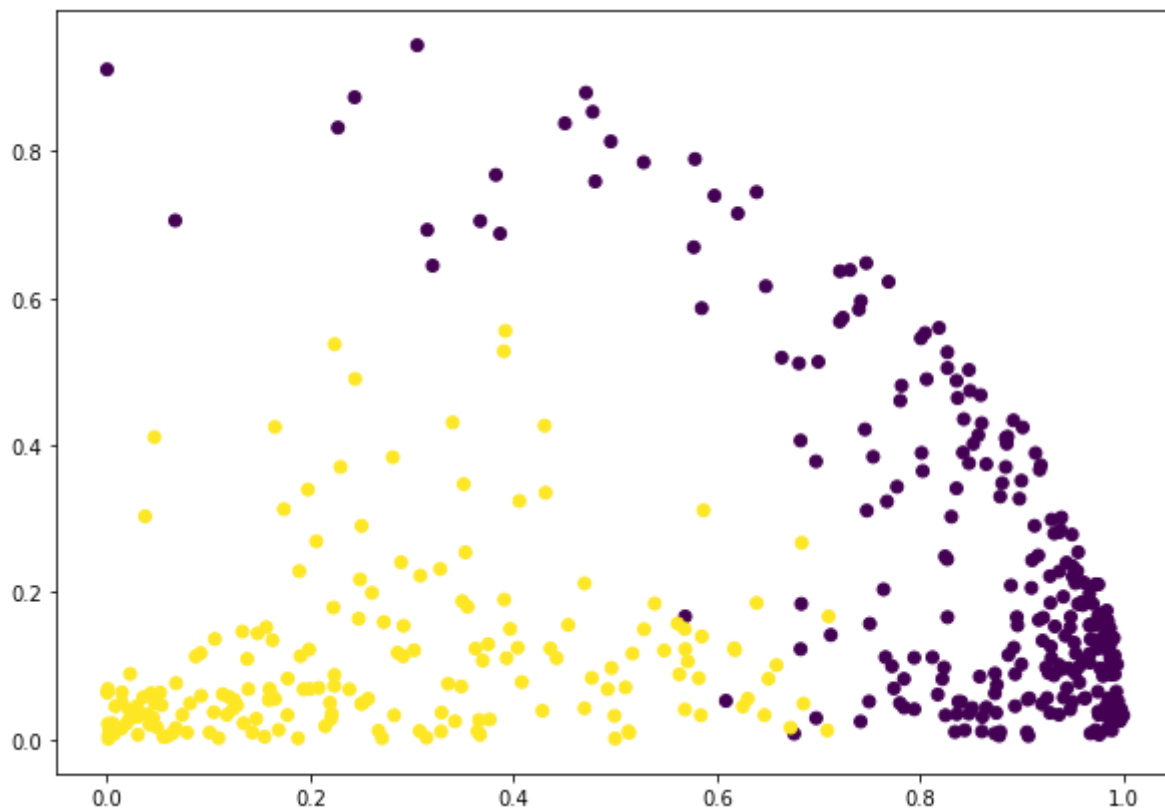
1 plt.figure(figsize=(10, 7))
2 plt.scatter(data_scaled['Milk'], data_scaled['Grocery'], c=cluster.labels_)

```



```
2 plt.scatter(data_scaled['Fresh'], data_scaled['Frozen'], c=cluster.labels_)
```

```
<matplotlib.collections.PathCollection at 0x7f450bfe0410>
```



1

1

1

1

1

1

1

```

1 # Bibliotecas
2 import numpy as np
3 import matplotlib.pyplot as plt
4 import pandas as pd

1 # Dataset
2 dataset = pd.read_csv('http://dicyg.fi-c.unam.mx:8080/lalo/pypcd/presentaciones/M
3 X = dataset.iloc[:, [3, 4]].values
4 dataset.head()

```

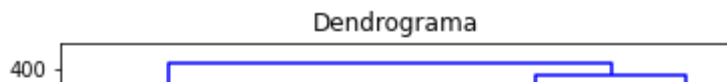


	CustomerID	Genre	Age	Annual Income (k\$)	Spending Score (1-100)
0	1	Male	19	15	39
1	2	Male	21	15	81
2	3	Female	20	16	6
3	4	Female	23	16	77
4	5	Female	31	17	40

```

1 # Dendrograma para determinar el número óptimo de grupos
2 import scipy.cluster.hierarchy as sch
3 dendrograma = sch.dendrogram(sch.linkage(X, method = 'ward'))
4 plt.title('Dendrograma')
5 plt.xlabel('Clientes')
6 plt.ylabel('Distancia euclideana')
7 plt.show()

```



```

1 from sklearn.cluster import AgglomerativeClustering
2 hc = AgglomerativeClustering(n_clusters = 5, affinity = 'euclidean',
3                               linkage = 'ward')
4 v hc = hc.fit_predict(X)

```

die 200

```
<matplotlib.collections.PathCollection at 0x7fbb7d03dd50>
```





```

1 import numpy as np
2 import matplotlib.pyplot as plt
3 import seaborn as sns
4 sns.set()
5 from sklearn.mixture import GaussianMixture
6 from sklearn.datasets import make_blobs

```

+ Code

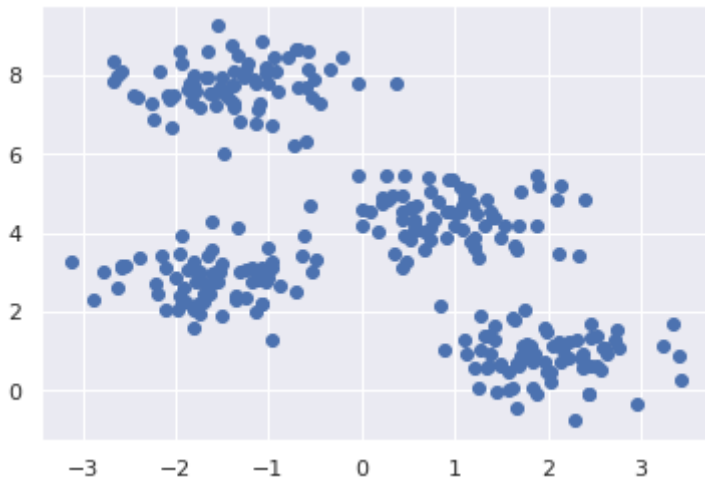
+ Text

```

1 X, y = make_blobs(n_samples=300, centers=4, cluster_std=0.6, random_state=0)
2 plt.scatter(X[:,0],X[:,1])

```

<matplotlib.collections.PathCollection at 0x7faf2360d890>

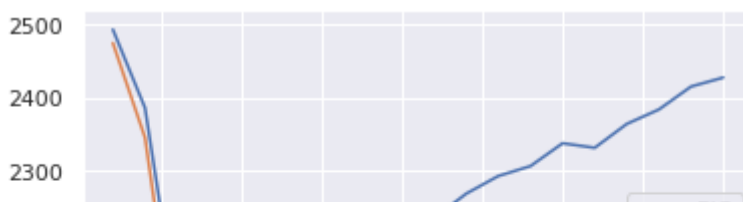


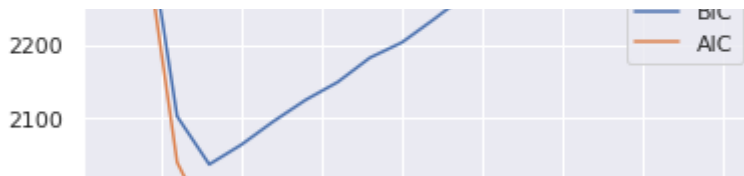
```

1 # The optimal number of clusters (K) is the value that minimizes the
2 # Akaike information criterion (AIC) or the Bayesian information criterion (BIC)
3 # https://en.wikipedia.org/wiki/Akaike\_information\_criterion
4 # https://en.wikipedia.org/wiki/Bayesian\_information\_criterion
5 n_clusters = np.arange(1, 21)
6 models = [GaussianMixture(n, covariance_type='full',
7                             random_state=0).fit(X) for n in n_clusters]
8 plt.plot(n_clusters, [m.bic(X) for m in models], label='BIC')
9 plt.plot(n_clusters, [m.aic(X) for m in models], label='AIC')
10 plt.legend(loc='best')
11 plt.xlabel('n_clusters')

```

Text(0.5, 0, 'n_clusters')

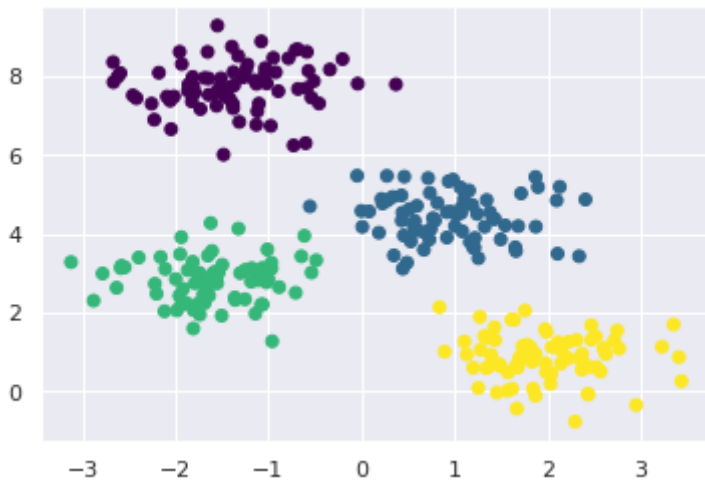




```
1 gmm = GaussianMixture(n_components=4)
2 gmm.fit(X)
```

```
GaussianMixture(n_components=4)
```

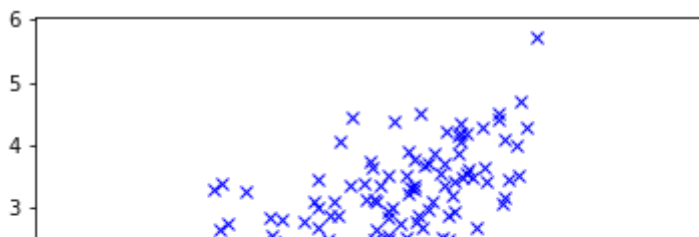
```
1 labels = gmm.predict(X)
2 plt.scatter(X[:, 0], X[:, 1], c=labels, cmap='viridis');
```

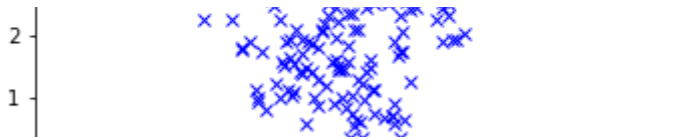


```
1 import numpy as np
2 import matplotlib.pyplot as plt
3 from sklearn.mixture import GaussianMixture
```

```
1 # https://bit.ly/3Bic3iu
2 X_train = np.load('data.npy')
```

```
1 plt.plot(X_train[:,0], X_train[:,1], 'bx')
2 plt.axis('equal')
3 plt.show()
```





```

1 gmm = GaussianMixture(n_components=2)
2 gmm.fit(X_train)
3
4 print("Medias: \n", gmm.means_)
5 print('\n')
6 print("Covarianzas: \n",gmm.covariances_)
7

```

Medias:

```

[[3.04641134 3.10654272]
 [1.60718016 1.35251723]]

```

Covarianzas:

```

[[[ 0.83656079  0.37865596]
 [ 0.37865596  0.72727426]]

 [[ 0.74995307 -0.5010097 ]
 [-0.5010097   0.74377694]]]

```

```

1 X, Y = np.meshgrid(np.linspace(-1, 6), np.linspace(-1,6))
2 XX = np.array([X.ravel(), Y.ravel()]).T
3 Z = gmm.score_samples(XX)
4 Z = Z.reshape((50,50))
5
6 plt.contour(X, Y, Z)
7 plt.scatter(X_train[:, 0], X_train[:, 1])
8 plt.show()
9

```



1



1



1

u



1

1

1

1

1

1

1

1

1

1

1

1

1

1

1

1

1

1