	Clustering Jerárquico Objetivos Utilizar scikit-learn para el clustering jerárquico. Crear dendogramas para visualizar los datos. Tabla de contenido 1. Clustering jerárquico - Aglomerativo A. Generando datos aleatorios
	B. Clustering aglomerativo C. Dendograma asociado para el clustering aglomerativo 2. Clustering para el dataset de vehículos A. Limpieza de datos B. Clustering utilizando Scipy C. Clustering utilizando scikit-learn
In [1]:	Clustering jerárquico- Aglomerativo Recuerde que el aglomerativo es el enfoque bottom up. Utilizaremos Complete Linkage como criterio de enlace. Nota: Puede intentar utilizar Avergare Linkage donde está Complete Linkage para ver la diferencia. import numpy as np import pandas as pd from scipy import ndimage from scipy.cluster import hierarchy from scipy.spatial import distance_matrix from matplotlib import pyplot as plt from sklearn import manifold, datasets from sklearn.cluster import AgglomerativeClustering from sklearn.datasets.samples_generator import make_blobs %matplotlib inline
	C:\Users\marco\Anaconda3\lib\site-packages\sklearn\utils\deprecation.py:143: FutureWar ning: The sklearn.datasets.samples_generator module is deprecated in version 0.22 and will be removed in version 0.24. The corresponding classes / functions should instead be imported from sklearn.datasets. Anything that cannot be imported from sklearn.datasets is now part of the private API. warnings.warn(message, FutureWarning) Generando datos aleatorios Generaremos un conjunto de datos usando la clase make_blobs. Los parámetros de entrada son:
<pre>In [2]: In [3]: Out[3]:</pre>	 n_samples: El número de puntos igualmente divididos entre clusters Seleccione un número entre 10 y 1500 centers: El número de centros a generar, o las ubicaciones centrales fijas Seleccione coordenadas x e y para generar los centros. Tenga 1-10 centros (ej. centers=[[1,1], [2,5]]) cluster_std: La desviación estándar de los clusters. Seleccione un número entre 0.5-1.5 (ul> Guarde el resultado en X1 e y1. X1, y1 = make_blobs (n_samples=50, centers=[[4,4], [-2, -1], [1, 1], [10,4]], cluster_: Gráfico de dispersión de los datos generados plt.scatter(X1[:, 0], X1[:, 1], marker='o') matplotlib.collections.PathCollection at 0x20288808700>
	Clustering aglomerativo La clase clustering aglomerativo requiere de 2 parámetros: • n_clusters: El número de clusters a formar así como el número de centroides a generar • El valor será: 4 • linkage: Qué criterio de enlace se utilizará. El criterio de enlace determina qué distancia se usará entre conjuntos de observaciones. The algorithm will merge the pairs of cluster that minimize this criterion. • El valor será: 'complete' • Nota: Se recomienda probar también "Average" El resultado se guardará en una variable llamada agglom
<pre>In [4]: In [5]: Out[5]: In [6]:</pre>	<pre>agglom = AgglomerativeClustering(n_clusters = 4, linkage = 'average') Ajustamos el modelo con X1 e y1 generados como se vio arriba. agglom.fit(X1,y1) AgglomerativeClustering(linkage='average', n_clusters=4) Veamos el cluster. # Creamos una figura de tamaño 6x4 pulgadas. plt.figure(figsize=(6,4)) # Estas 2 líneas se utilizan para escalar los puntos de datos # de otro modo los puntos quedarán muy dispersos entre sí # Creamos un tango mínimo y máximo para X1 x_min, x_max = np.min(X1, axis=0), np.max(X1, axis=0) # Obtenemos la distancia promedio para X1 X1 = (X1 - x min) / (x max - x min)</pre>
	<pre># Este loop despliega los puntos de datos for i in range(X1.shape[0]): # Reemplazamos los puntos de datos con sus respectivos valores de cluster # (ejemplo 0) y su color es codificado con un colormap plt.text(X1[i, 0], X1[i, 1], str(y1[i]),</pre>
In [7]:	Dendograma ascoiado para el clustering jerárquico aglomerativo Recuerde que la matriz de distancia contiene la distancia de cada punto a cada otro punto del dataset. Utilice la función distance_matrix que requiere 2 entradas. Use la matriz de características, X1, como esas entradas y guarde la matriz de distancia en una variable llamada dist_matrix. Recuerde que la los valores de distancia son simétricos, con una diagonal de 0's. Esta es una forma de verificar que los valores son correctos. dist_matrix = distance_matrix(X1, X1)
	print (dist_matrix) [[0.
In [8]:	Guarde el resultado en la variable Z Z = hierarchy.linkage(dist_matrix, 'complete') <ipython-input-8-3814b774a052>:1: ClusterWarning: scipy.cluster: The symmetric non-neg ative hollow observation matrix looks suspiciously like an uncondensed distance matrix Z = hierarchy.linkage(dist_matrix, 'complete') El clustering jerárquico es típicamente visualizado como un dendograma tal cual se muestra en la celda siguiente. Cada combinación es representada por una línea horizontal. La coordenada y de la línea horizontal es ls similaridad de los 2 clusters combinados, donde las ciudades son vistas como clusters singleton. Al movernos desde la capa inferior hacia el nodo superior, el dendograma nos permite reconstruir la historia de combinaciones que resultaron en el cluster mostrado. Luego, guardaremos el dendograma en una variable llamada dendro. Al hacerlo, el dendograma será desplegado. Utilizaremos dendogram de la clase hierarchy y le pasaremos el parámetro:</ipython-input-8-3814b774a052>
In [9]:	<pre>• Z dendro = hierarchy.dendrogram(Z)</pre>
In [10]:	Práctico Utilicemos average como linkage para ver cómo cambia el dendograma. Z = hierarchy.linkage(dist_matrix, 'average') dendro = hierarchy.dendrogram(Z) <ipython-input-10-90bcb3750f27>:1: ClusterWarning: scipy.cluster: The symmetric non-ne gative hollow observation matrix looks suspiciously like an uncondensed distance matrix Z = hierarchy.linkage(dist_matrix, 'average') 4.0 3.5 3.0 3.0</ipython-input-10-90bcb3750f27>
	Clustering en el dataset de vehículos Imagine que un fabricante de automóviles ha desarrollado prototipos para un nuevo vehículo. Antes de introducir el modelo nuevo, el fabricante quiere determinar qué vehículos existentes en el mercado son los más parecidos al prototipo; esto es, cómo los vehículos pueden ser agrupados, qué grupo es el más similar
In [11]:	a este modelo y por tanto contra qué modelos competirá. Nuestro objetivo es usar métodos de clustering para encontrar los clusters más distintivos de vehículos. Resumirá los vehículos existentes y ayudará a los fabricantes a tomar una decisión sobre el suministro de nuevos modelos. Descargando los datos #!wget -0 cars_clus.csv https://cf-courses-data.s3.us.cloud-object-storage.appdomain.cloud/import urllib.request url = 'https://cf-courses-data.s3.us.cloud-object-storage.appdomain.cloud/IBMDevelope;
Out[11]: In [12]:	
Out[12]:	print ("Shape of dataset: ", pdf.shape) pdf.head(5) Shape of dataset: (159, 16) manufact model sales resale type price engine_s horsepow wheelbas width length curb_wgt 0 Acura Integra 16.919 16.360 0.000 21.500 1.800 140.000 101.200 67.300 172.400 2.639 1 Acura TL 39.384 19.875 0.000 28.400 3.200 225.000 108.100 70.300 192.900 3.517 2 Acura CL 14.114 18.225 0.000 null 3.200 225.000 106.900 70.600 192.000 3.470 3 Acura RL 8.588 29.725 0.000 42.000 3.500 210.000 114.600 71.400 196.600 3.850 4 Audi A4 20.397 22.255 0.000 23.990 1.800 150.000 102.600 68.200 178.000 2.998 Limpieza de dataset eliminando las filas que contienen valores null:
<pre>In [13]: Out[13]:</pre>	<pre>print ("Shape of dataset before cleaning: ", pdf.size) pdf[['sales', 'resale', 'type', 'price', 'engine_s',</pre>
In [14]:	Audi A4 20.397 22.255 0.0 23.99 1.8 150.0 102.6 68.2 178.0 2.998 Audi A6 18.780 23.555 0.0 33.95 2.8 200.0 108.7 76.1 192.0 3.561 Selección de características featureset = pdf[['engine_s', 'horsepow', 'wheelbas', 'width', 'length', 'curb_wgt', Normalización MinMaxScaler transforma características escalando cada característica en un rango dado. Es por defecto
<pre>In [15]: Out[15]:</pre>	(0,1). Esto es, el estimador escala y traslada cada característica individualmente de modo que queda entre 0 y 1. from sklearn.preprocessing import MinMaxScaler x = featureset.values #returns a numpy array min_max_scaler = MinMaxScaler() feature_mtx = min_max_scaler.fit_transform(x) feature_mtx [0:5] array([[0.11428571, 0.21518987, 0.18655098, 0.28143713, 0.30625832,
In [16]:	Clustering usando Scipy En esta parte utilizaremos el paquete Scipy. Primero, calculemos la matriz de distancia. import scipy leng = feature_mtx.shape[0] D = scipy.zeros([leng,leng]) for i in range(leng):
In [17]:	<pre>- single - complete - average - weighted - centroid Utilizaremos complete import pylab import scipy.cluster.hierarchy Z = hierarchy.linkage(D, 'complete') <ipython-input-17-8655000d21de>:3: ClusterWarning: scipy.cluster: The symmetric non-ne</ipython-input-17-8655000d21de></pre>
In [18]:	gative hollow observation matrix looks suspiciously like an uncondensed distance matrix Z = hierarchy.linkage(D, 'complete') Esencialmente, el clustering jerárquico NO requiere un número pre-establecido de clusteres. Sin embargo, en algunas aplicaciones, queremos una partición de clusteres disjuntos como en el "flat clustering". Entonces, puede utilizar una línea de corte: from scipy.cluster.hierarchy import fcluster max_d = 3 clusters = fcluster(Z, max_d, criterion='distance') clusters array([1, 5, 5, 6, 5, 4, 6, 5, 5, 5, 5, 5, 5, 4, 4, 5, 1, 6,
<pre>In [19]: Out[19]:</pre>	5, 5, 5, 4, 2, 11, 6, 6, 5, 6, 5, 1, 6, 6, 10, 9, 8, 9, 3, 5, 1, 7, 6, 5, 3, 5, 5, 4, 4, 3, 2, 6, 6, 5, 4, 2, 1, 6, 5, 2, 7, 5, 5, 5, 4, 4, 3, 2, 6, 6, 5, 7, 4, 7, 6, 6, 5, 3, 5, 5, 6, 5, 4, 4, 1, 6, 5, 5, 5, 6, 4, 5, 4, 1, 6, 5, 6, 6, 5, 5, 5, 7, 7, 7, 7, 2, 2, 1, 2, 6, 5, 1, 1, 1, 7, 8, 1, 1, 6, 1, 1], dtype=int32) También puede determinar el número de clusters directamente: from scipy.cluster.hierarchy import fcluster k = 5 clusters = fcluster(Z, k, criterion='maxclust') clusters array([1, 3, 3, 3, 3, 2, 3, 3, 3, 3, 3, 3, 2, 2, 3, 1, 3, 3, 3, 3, 2, 1,
In [20]:	5, 3, 3, 3, 3, 3, 1, 3, 3, 4, 4, 4, 4, 2, 3, 1, 3, 3, 3, 2, 3, 2, 4, 3, 4, 1, 3, 3, 3, 2, 1, 1, 3, 3, 1, 3, 3, 3, 3, 2, 2, 2, 1, 3, 3, 3, 3, 2, 3, 3, 3, 3, 2, 3, 3, 3, 3, 3, 3, 3, 3, 3, 3, 3, 3, 3,
	Ford F-Series 1 Clodge Ram Picksp 1 1 Clodge Ram Vers 1 Teyer Legedition 1 Clodge Ram Wagon 1 Clodge Ram Wagon 1 Clodge Marw Wagon 1 Clodge Wiper O Cleverold Metro 0 Metsubshi 900007 Ford Mustang 0 Trysta & diumer 1 Perche Carrae Carlog 0 Trysta & Grand Ford 1 Perche Carrae Carlog 0 Clodge Accord 0 Perche Carrae Carlog 0 Clodge Status 0 Perche Carrae Carlog 0 Perche Carrae Carlog 0 Perche Carrae Carlog 0 Perche Carlog 0 Perche Carrae Carlog 0 Perche Carlog 0 Perche Carrae Carlog 0 P
	[Referedes-Benz E-Class 0] [Internal State of Class of Cl
	[Mercury Grand Marquis 0] [Ford Crown Victoria 0] [Lincoln Town car 0] [Ford Windstar 1] [Dodge Dakota 1] [Toyota Corolla 0] [Chevrolet Prizm 0] [Missubishi Mirage 0] [Saturn SC 0] [Hyundai Accent 0] [Hyundai Accent 0] [Jeep Wrangler 1] [Volkswagen Cabrio 0] [Toyota RAV4 1] [Volkswagen GTi 0] [Volkswagen GTi 0] [Isyota Celica 0] [Nissan Sentra 0] [Saturn SW 0] [Ford Escort 0] [Toyota Tacoma 1] [Chevrolet Cavalier 0] [Volkswagen jetta 0] [Volkswagen jetta 0]
In [21]:	Clustering utilizando scikit-learn
<pre>In [22]: Out[22]:</pre>	 Ward minimiza la suma de las diferencias cuadráticas entre todos los clusters. Es un enfoque de varianza mínima y en este sentido es similar a la función objetivo de k-means pero atacado con un enfoque de clustering aglomerativo. Maximun o complete linkage minimiza la distancia máxima entre observaciones de pares de clusters. Average minimiza el promedio de distancias entre observaciones de pares de clusters. agglom = AgglomerativeClustering(n_clusters = 6, linkage = 'complete') agglom.fit(feature_mtx) agglom.labels_ array([1, 2, 2, 1, 2, 3, 1, 2, 2, 2, 2, 2, 2, 3, 3, 2, 1, 1, 2, 2, 2, 2, 5, 1, 4, 1, 1, 2, 1, 2, 1, 2, 1, 1, 1, 5, 0, 0, 0, 3, 2, 1, 2, 1, 2, 3, 2, 3,
In [23]: Out[23]:	<pre>4, 1, 1, 2, 1, 2, 1, 1, 1, 5, 0, 0, 0, 3, 2, 1, 2, 1, 2, 3, 2, 3, 0, 3, 0, 1, 1, 1, 2, 3, 1, 1, 1, 2, 1, 1, 2, 2, 2, 3, 3, 3, 3, 1, 1, 1, 2, 1, 2, 2, 1, 1, 2, 3, 2, 3, 1, 2, 3, 5, 1, 1, 2, 3, 2, 1, 3, 2, 3, 1, 1, 2, 1, 1, 2, 2, 2, 1, 1, 1, 1, 1, 1, 1, 1, 2, 1, 1, 2, 0, 1, 1, 1, 1], dtype=int64) Agregamos un nuevo campo al dataframe para mostrar el cluster de cada fila: pdf['cluster_'] = agglom.labels_ pdf.head()</pre>
In [24]:	Acura RL 8.588 29.725 0.0 42.00 3.5 210.0 114.6 71.4 196.6 3.850 3 Audi A4 20.397 22.255 0.0 23.99 1.8 150.0 102.6 68.2 178.0 2.998 4 Audi A6 18.780 23.555 0.0 33.95 2.8 200.0 108.7 76.1 192.0 3.561 import matplotlib.cm as cm n_clusters = max(agglom.labels_)+1 colors = cm.rainbow(np.linspace(0, 1, n_clusters)) cluster_labels = list(range(0, n_clusters)) # Create a figure of size 6 inches by 4 inches.
	<pre>plt.figure(figsize=(16,14)) for color, label in zip(colors, cluster_labels): subset = pdf[pdf.cluster_ == label] for i in subset.index:</pre>
	as value-mapping will have precedence in case its length matches with *x* & *y*. Plea se use the *color* keyword-argument or provide a 2-D array with a single row if you in tend to specify the same RGB or RGBA value for all points. *c* argument looks like a single numeric RGB or RGBA sequence, which should be avoided as value-mapping will have precedence in case its length matches with *x* & *y*. Plea se use the *color* keyword-argument or provide a 2-D array with a single row if you in tend to specify the same RGB or RGBA value for all points. *c* argument looks like a single numeric RGB or RGBA sequence, which should be avoided as value-mapping will have precedence in case its length matches with *x* & *y*. Plea se use the *color* keyword-argument or provide a 2-D array with a single row if you in tend to specify the same RGB or RGBA value for all points. *c* argument looks like a single numeric RGB or RGBA sequence, which should be avoided as value-mapping will have precedence in case its length matches with *x* & *y*. Plea se use the *color* keyword-argument or provide a 2-D array with a single row if you in tend to specify the same RGB or RGBA value for all points. *c* argument looks like a single numeric RGB or RGBA sequence, which should be avoided as value-mapping will have precedence in case its length matches with *x* & *y*. Plea se use the *color* keyword-argument or provide a 2-D array with a single row if you in tend to specify the same RGB or RGBA value for all points.
Out[24]:	*c* argument looks like a single numeric RGB or RGBA sequence, which should be avoided as value-mapping will have precedence in case its length matches with *x* & *y*. Plea se use the *color* keyword-argument or provide a 2-D array with a single row if you in tend to specify the same RGB or RGBA value for all points. Text(0, 0.5, 'mpg') Clusters duster0 duster1 duster2 duster3 duster4 duster5
	25 - 25 - 25 - 26 - 27 - 28 - 28 - 28 - 29 - 20 -
In [25]:	Estamos viendo la distribución de cada cluster en el gráfico de dispersión, pero no es claro donde está el centroide de cada cluster. Más aún, hay 2 tipos de vehículos en nuestro dataset, "truck (camión)" (valor de 1 en la columna type) y "car". Por lo tanto, los usamos para distinguir las clases y resumir el clúster. Primero contamos el número de casos en cada grupo: pdf.groupby(['cluster_', 'type']) ['cluster_'].count()
Out[25]:	cluster_ type 0
<pre>In [26]: Out[26]:</pre>	<pre>agg_cars = pdf.groupby(['cluster_','type'])['horsepow','engine_s','mpg','price'].mean agg_cars <ipython-input-26-fb63ecec89ff>:1: FutureWarning: Indexing with multiple keys (implici tly converted to a tuple of keys) will be deprecated, use a list instead. agg_cars = pdf.groupby(['cluster_','type'])['horsepow','engine_s','mpg','price'].mea n() horsepow engine_s mpg price cluster_ type 0 1.0 211.666667 4.483333 16.166667 29.024667</ipython-input-26-fb63ecec89ff></pre>
	1 0.0 146.531915 2.246809 27.021277 20.306128 1.0 145.000000 2.580000 22.200000 17.009200 2 0.0 203.111111 3.303704 24.214815 27.750593 1.0 182.090909 3.345455 20.181818 26.265364 3 0.0 256.500000 4.410000 21.500000 42.870400 1.0 160.571429 3.071429 21.428571 21.527714 4 0.0 55.000000 1.000000 45.000000 9.235000 5 0.0 365.666667 6.233333 19.3333333 66.010000
	 5 0.0 365.666667 6.233333 19.333333 66.010000 Tenemos 3 clusters principales con la mayoría de los vehículos en ellos: Cars: Cluster 1: mpg casi alto y baja potencia Cluster 2: buen mpg y horsepower, pero alto precio Cluster 3: mpg bajo, alto horsepower y alto precio Trucks: Cluster 1: mpg casi alto dentro de los camiones, bajo horsepower y precio
In [28]:	 Cluster 1: mpg casi alto dentro de los camiones, bajo horsepower y precio Cluster 2: mpg casi bajo, horsepower medio y alto precio Cluster 3: buen mpg y horsepower, bajo precio plt.figure (figsize=(16,10)) for color, label in zip(colors, cluster_labels): subset = agg_cars.loc[(label,),] for i in subset.index: plt.text(subset.loc[i][0]+5, subset.loc[i][2], 'type='+str(int(i)) + ', price-plt.scatter(subset.horsepow, subset.mpg, s=subset.price*20, c=color, label='clusteplt.legend() plt.title('Clusters') plt.xlabel('horsepow') plt.ylabel('mpg')
	c argument looks like a single numeric RGB or RGBA sequence, which should be avoided as value-mapping will have precedence in case its length matches with *x* & *y*. Plea se use the *color* keyword-argument or provide a 2-D array with a single row if you in tend to specify the same RGB or RGBA value for all points. *c* argument looks like a single numeric RGB or RGBA sequence, which should be avoided as value-mapping will have precedence in case its length matches with *x* & *y*. Plea se use the *color* keyword-argument or provide a 2-D array with a single row if you in tend to specify the same RGB or RGBA value for all points. *c* argument looks like a single numeric RGB or RGBA sequence, which should be avoided as value-mapping will have precedence in case its length matches with *x* & *y*. Plea se use the *color* keyword-argument or provide a 2-D array with a single row if you in tend to specify the same RGB or RGBA value for all points. *c* argument looks like a single numeric RGB or RGBA sequence, which should be avoided
Out[28]:	*c* argument looks like a single numeric RGB or RGBA sequence, which should be avoided as value—mapping will have precedence in case its length matches with *x* & *y*. Plea se use the *color* keyword-argument or provide a 2-D array with a single row if you in tend to specify the same RGB or RGBA value for all points. *c* argument looks like a single numeric RGB or RGBA sequence, which should be avoided as value—mapping will have precedence in case its length matches with *x* & *y*. Plea se use the *color* keyword-argument or provide a 2-D array with a single row if you in tend to specify the same RGB or RGBA value for all points. *c* argument looks like a single numeric RGB or RGBA sequence, which should be avoided as value—mapping will have precedence in case its length matches with *x* & *y*. Plea se use the *color* keyword-argument or provide a 2-D array with a single row if you in tend to specify the same RGB or RGBA value for all points. Text(0, 0.5, 'mpg') *Custers *duster0 *duster1 *duster2 *duster2 *duster3 *duster4 *duster3 *duster4 *duster5 *duster4 *duster5 *duster6 *duster6 *duster7 *duster7 *duster7 *duster9 *duster6 *duster6 *duster6 *duster7 *duster7 *duster7 *duster7 *duster7 *duster9 *d
In []:	25 - type=0, price=20k 25 - type=1, price=17k type=1, price=21k type=1, price=21k type=1, price=26k type=1, price=26k type=1, price=26k type=1, price=29k type=1, price=29k type=1, price=29k