	11.1 Gerando dados
ո [3]։	<pre>import numpy as np</pre>
	<pre>import matplotlib.pyplot as plt  # Gerando dados aleatórios para testar o algoritmo  1_x1 = np.random.rand(100) + 10  1_y1 = np.random.rand(100) + 10  1_x2 = np.random.rand(100) + 13  1_y2 = np.random.rand(100) + 13</pre>
	<pre>l_x3 = np.random.rand(100) + 13 l_y3 = np.random.rand(100) + 10  fig, ax = plt.subplots(1,1) ax.scatter(l_x1, l_y1, 5 ,color = "red") ax.scatter(l_x2, l_y2, 5 ,color = "blue")</pre>
	ax.scatter(1_x3, 1_y3, 5 ,color = "black") plt.show()  14.0 - 13.5 -
	13.0 - 12.5 - 12.0 - 11.5 - 11.0 -
	Vemos que de fato os dados são bem 'separáveis' em 3 conjuntos distintos. Abaixo vamos salvar o conjunto criado em um único array
n [4]:	<pre>X = [] for i in range(0, len(l_x1)):     X.append( np.array((l_x1[i],l_y1[i])) )</pre>
	<pre>X.append( np.array((1_x2[i],1_y2[i])) ) X.append( np.array((1_x3[i],1_y3[i])) )  X = np.array(X)  11.2 Calculando os grupos</pre>
n [5]:	Usamos o KMeans para calcular os grupos aos quais os dados pertencem. Para isso precisamos passar o número de grupos que queremos utilizar (no caso abaixo 3). O retorno é um array com o mesmo tamanho do conjunto de dados de entrada, em que cada elemento contém um número referente ao cluster que o dado pertence.  # Usando o algoritmo para ajustar os dados usando 3 clusters
ıt[5]:	<pre>from sklearn.cluster import KMeans as km kmeans = km(n_clusters = 3, random_state = 0).fit(X)  # O retorno é um array com o numero do cluster a que a amostra foi associada kmeans.labels_ array([1, 0, 2, 1,</pre>
	0, 2, 1, 0,
	2, 1, 0, 2])  Agora podemos separar os dados, cada um em seu cluster. Usando a indexação boolena no vetor X para coletar todos os valores em contrata de la coletar todos em contrata de la coletar todos os valores em contrata de la coletar todos os valores em contrata de la coletar todos en contrata de la coletar todos en coletar todos os valores em contrata de la coletar todos en contrata de la coletar todos en
	labels == 1, 2, Em seguida plotamos os dados em um gráfico (scatterplot) novamente, e conseguimos retornar ao mesmo gráfico que fizemos após a geração dos dados. O que mostra que todos os dados ficaram de fato em seus respectivos grupos.  Ainda, usamos o método cluster_centers, que contém um array bidimensional (neste caso) com os centróides dos grupos para plotá-los também.
n [6]:	<pre>cond = kmeans.labels_ == 0 C1 = X[cond]  cond2 = kmeans.labels_ == 1 C2 = X[cond2]  cond3 = kmeans.labels == 2</pre>
	C3 = X[cond3]  'Plotando os dados no plano, cada um com uma cor' fig, ax = plt.subplots(1,1) plt.scatter(C1[:,0], C1[:,1],5, color = "red") # Todas as linhas da coluna 0, # Todas as linhas da coluna plt.scatter(C2[:,0], C2[:,1],5, color = "blue") # Todas as linhas da coluna 0, # Todas as linhas da coluna plt.scatter(C3[:,0], C3[:,1],5, color = "black") # Todas as linhas da coluna 0, # Todas as linhas
ut[6]:	<pre>plt.scatter(kmeans.cluster_centers_[:, 0], kmeans.cluster_centers_[:, 1],40, c = 'black', label = 'Centroi plt.legend() kmeans.cluster_centers_ array([[13.52598415, 13.45670901],</pre>
	[13.51595468, 10.49950257]])  14.0 Centroids  13.0 12.5 1
	12.0 - 11.5 - 11.0 - 10.5 -
	11.3 Definindo o número de grupos ( <i>clusters</i> )
	Mas como determinar o número de clusters? No exemplo acima definimos 3, pois já sabíamos de antemão que esse era um bom número, dado que o conjunto foi gerado com 3 grupos. Em problemas reais, no entanto, não teremos nenhuma informações a respeit desse número.  Podemos usar uma abordagem visual chamada <b>gráfico de cotovelo (elbow plot)</b> . Para isso, usamos uma medida de desempenho do agrupamento chamada WCSS(Within-Cluster Sum of Square), ou soma dos quadrados intra-grupos. Esse medida calcula a soma da
	distancia de todos os pontos ao centróide de seus grupos. A idéia é que a medida que o número de grupos aumenta, essa distância d ser reduzida. Porém a partir de um determinado número a redução passa a ser insignificante, de forma que esse é o ponto em que aumentar o número de grupos não influencia mais no desempenho do agrupamento.  O atributo .inertia_ fornece o WCCS. Com isso podemos calcular os grupos para diversos valores de $k$ , salvando os valores de WC
n [7]:	<pre>em uma lista, em seguida plotamos os valores e obtemos o gráfico de cotovelo.  l_inertia = [] for k in range(1,10):     kmeans = km(n_clusters = k, random_state = 0).fit(X)     l_inertia.append(kmeans.inertia_) l_inertia</pre>
	<pre># Agora plotamos os valores de inertia e vemos como o gráfico de cotovelo fica: fig, ax = plt.subplots(1,1) fig.set_size_inches(10,6) ax.plot(l_inertia, marker = "o", color = "orange") plt.show()</pre>
	1200 -
	600 -
	200 - 0 1 2 3 4 5 6 7 8
	Pela análise do gráfico percebemos que o cotovelo está localizado no segundo grupo (como começa em zero, no terceiro), o que cono com o esperado.  11.4 Prevendo novos valores
[151	Com o modelo estimado, podemos prever em quais grupos novos conjuntos de dados seriam atribuidos, usando a função predict exemplo abaixo demonstra:  # Estimando à qual cluster o ponto [10.3, 13.1] seria atribuido novo_dado = [[10.3, 13.1]]
	<pre>grupo_novo_dado = kmeans.predict(novo_dado) print("O grupo do novo dado é : ",grupo_novo_dado[0])  O grupo do novo dado é : 1  11.5 Cuidados</pre>
	<ul> <li>11.5.1 Tratando dados não numéricos (categóricos/ordinais)</li> <li>O método k-means só funciona com dados numéricos, de forma que se existem atributos categóricos se faz necessário realizar uma transformações nos dados. Considere o DataFrame abaixo, com a coluna C2 e C3 do tipo categórico.</li> </ul>
[15]:	C1 = [.4,.5,.8,.9] #Float C2 = ["Tipo1", "Tipo2", "Tipo4"]
	C3 = ["SP", "SP", "SC", "MG"] C4 = [10,20,20,30] dic = {"C1": C1,
t[15]:	C3 = ["SP", "SP", "SC", "MG"]  C4 = [10,20,20,30]  dic = {"C1": C1,
t[15]:	C3 = ["SP", "SP", "SC", "MG"]  C4 = [10,20,20,30]  dic = {"C1": C1,
t[15]: [25]:	C3 = ["SP", "SP", "SC", "MG"] C4 = [10,20,20,30] dic = {"C1": C1,
	C3 = ["SP", "SP", "SC", "MG"] C4 = [10,20,20,30] dic = {"C1": C1,
[25]:	C3 = ["sgr", "sgr", "scr", "wG"] C4 = [10, 20, 20, 30] dic = ("C1": C1,
[25]: t[25]:	C3 = [19.5]***, "S9C**, "S9C**, "S9C**] C4 = [10.2, 20, 30] dic =   "C1***; C1,
[25]: t[25]:	C3 =   Tipstranser
[25]: t[25]:	CS =   T99P*, T9P*, 20,3 ct   Cd =   10,2 y,2 y,3 ct   Cd   Cd   Cd   Cd   Cd   Cd   Cd   C
[25]: t[25]: t[26]:	C3
[25]: t[25]: t[26]:	cos = (15.9%, 25.0%) **sex** *
[25]: t[25]: t[26]:	C1   C2   C3   C3   C3   C3   C4
[25]: t[25]: t[26]:	content of the conten
[25]: t[25]: t[26]:	Call = 170,190,190,190   100,190   1
[25]: t[25]: t[26]:	C = 12,712,712,712,712,713,713,713,713,713,713,713,713,713,713
[25]: t[25]: t[26]:	13
[25]: t[25]: t[26]: [34]:	12
[25]: t[25]: t[26]: [34]:	Side of Section 1997 (1997) 19
[25]: t[25]: t[26]: [34]:	Company   Comp
[25]: t[25]: [34]: [34]:	Colling 1971 - 1
[25]: t[25]: [34]: [34]:	College Control Contro
[25]: t[25]: t[26]: t[34]: [34]:	Compare   Comp
[25]: t[25]: t[26]: t[34]: [34]:	Compare   Comp
[25]: t[25]: t[26]: t[34]: [34]:	See 1 May 19 12 13 1
[25]: t[25]:  [26]: t[26]: t[47]:	The second control of the control of
[25]: t[25]:  [26]: t[26]: t[47]:	Control   Cont
[25]: t[25]:  [26]: t[26]: t[47]:	Company   Comp
[25]: t[25]:  [26]: t[26]: t[47]:	Company   Comp
[25]: t[25]:  [26]: t[26]: t[47]:	Commission of the commission o
[25]: t[25]:  [26]: t[26]: t[47]:	Comment of the commen
[25]: t[25]:  [26]: t[26]: t[47]:	Comments of the comments of th
[25]: t[25]:  [26]: t[26]: t[47]:	The control of the co
[25]: t[25]:  [26]: t[26]: t[47]:	See that the second of the sec
[25]: t[25]:  [26]: t[26]: t[47]:	See the control of th