

Kmeans算法的K值最优解

非凡wang咖 2019-04-02

画

赵雷 - 赵小雷



上篇文章为大家介绍了我们常用的聚类算法Kmenas算法，也为大家整理了一点小案例，今天为大家继续分享我们Kmenas算法，对Kmenas算法来说，如何确定簇数K值是一个至关重要的问题，为了解决这个问题，通常会选用探索法，即给定不同的k值下，对比某些评估指标的变动情况，进而选择一个比较合理的k值，在这我们上篇文章给大家推荐了三种方法（簇内离差平方和拐点法，轮廓系数法和间隔统计量法），接下来我们分别看看这三种方法是如何实现的：

Kmenas算法基础公式：

$$J = \sum_{i=1}^k \sum_{j \in c_k} (x^{(j)} - \mu_i)^2$$

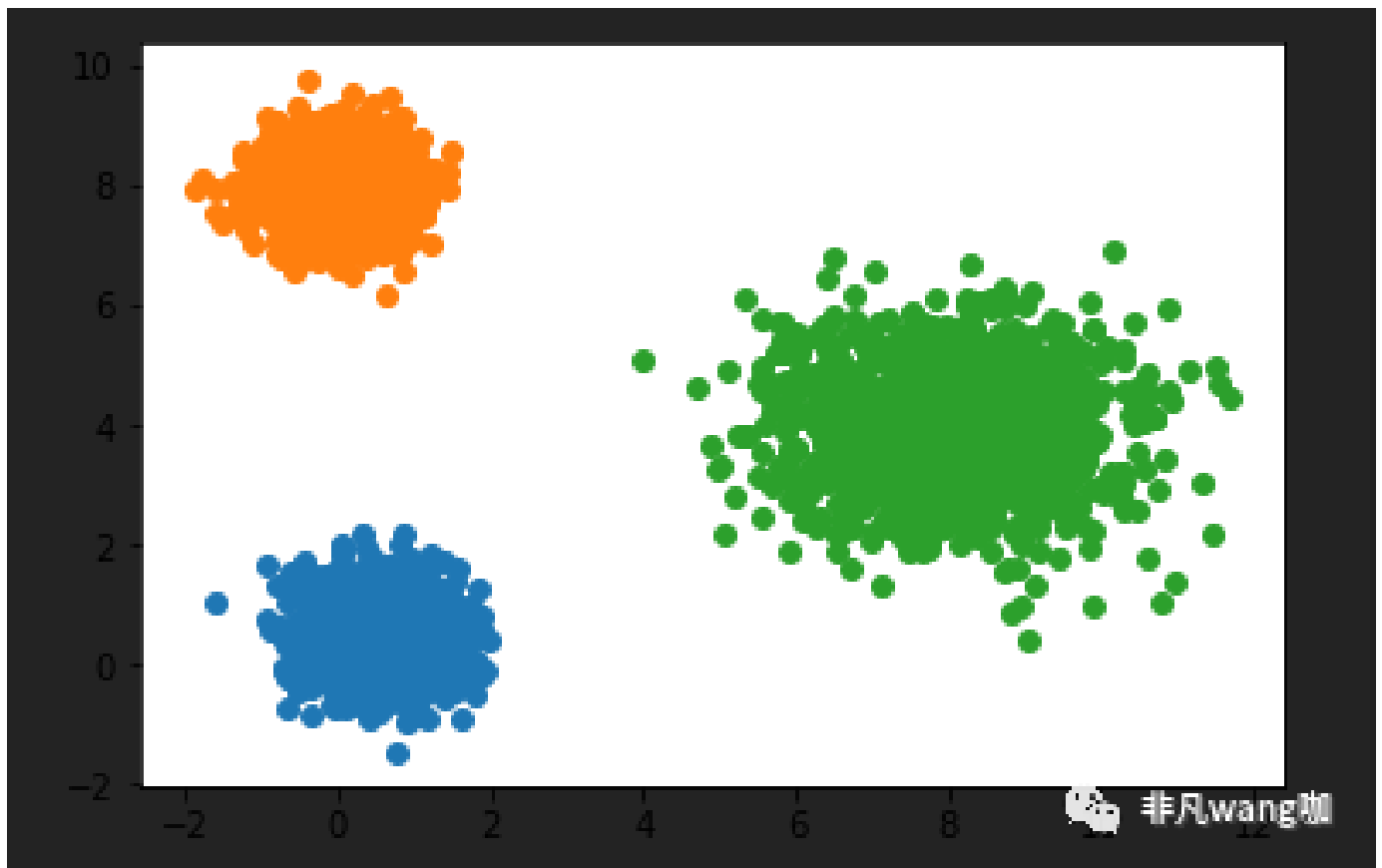
非凡wang咖

拐点法

簇内离差平方和拐点法的思想很简单，就是在不同的k值下计算簇内离差平方和，然后通过可视化的方法找到“拐点”所对应的k值，J为Kmeans算法的目标函数，随着簇数量的增加，簇中的样本量会越来越少，进而导致目标函数J的值也会越来越小，通过可视化方法，重点关注的是斜率的变化，当斜率由大突然变小时，并且之后的斜率变化缓慢，则认为突然变化的点就是寻找的目标点，因为继续随着簇数K的增加，聚类效果不再有大的变化。

接下来我们就验证这个方法，随机生成三组二元正态分布数据，首先基于该数据绘制散点图，如下代码：

```
1 import pandas as pd
2 import numpy as np
3 import matplotlib.pyplot as plt
4 from sklearn.cluster import KMeans
5
6 #随机生成三组二元正态分布随机数
7 np.random.seed(1234)
8 mean1 = [0.5,0.5]
9 cov1 = [[0.3,0],[0,0.3]]
10 x1,y1 = np.random.multivariate_normal(mean1,cov1,1000).T
11
12 mean2 = [0,8]
13 cov2 = [[0.3,0],[0,0.3]]
14 x2,y2 = np.random.multivariate_normal(mean2,cov2,1000).T
15
16 mean3 = [8,4]
17 cov3 = [[1.5,0],[0,1]]
18 x3,y3 = np.random.multivariate_normal(mean3,cov3,1000).T
19
20 #绘制三组数据的散点图
21 plt.scatter(x1,y1)
22 plt.scatter(x2,y2)
23 plt.scatter(x3,y3)
24 plt.show()
```



如上图，虚拟出来的数据呈现出三个簇，接下来基于这个虚拟数据，使用拐点法绘制簇的个数与总的簇内离差平方和之间的折线图，确定最终的k值，代码如下：

```

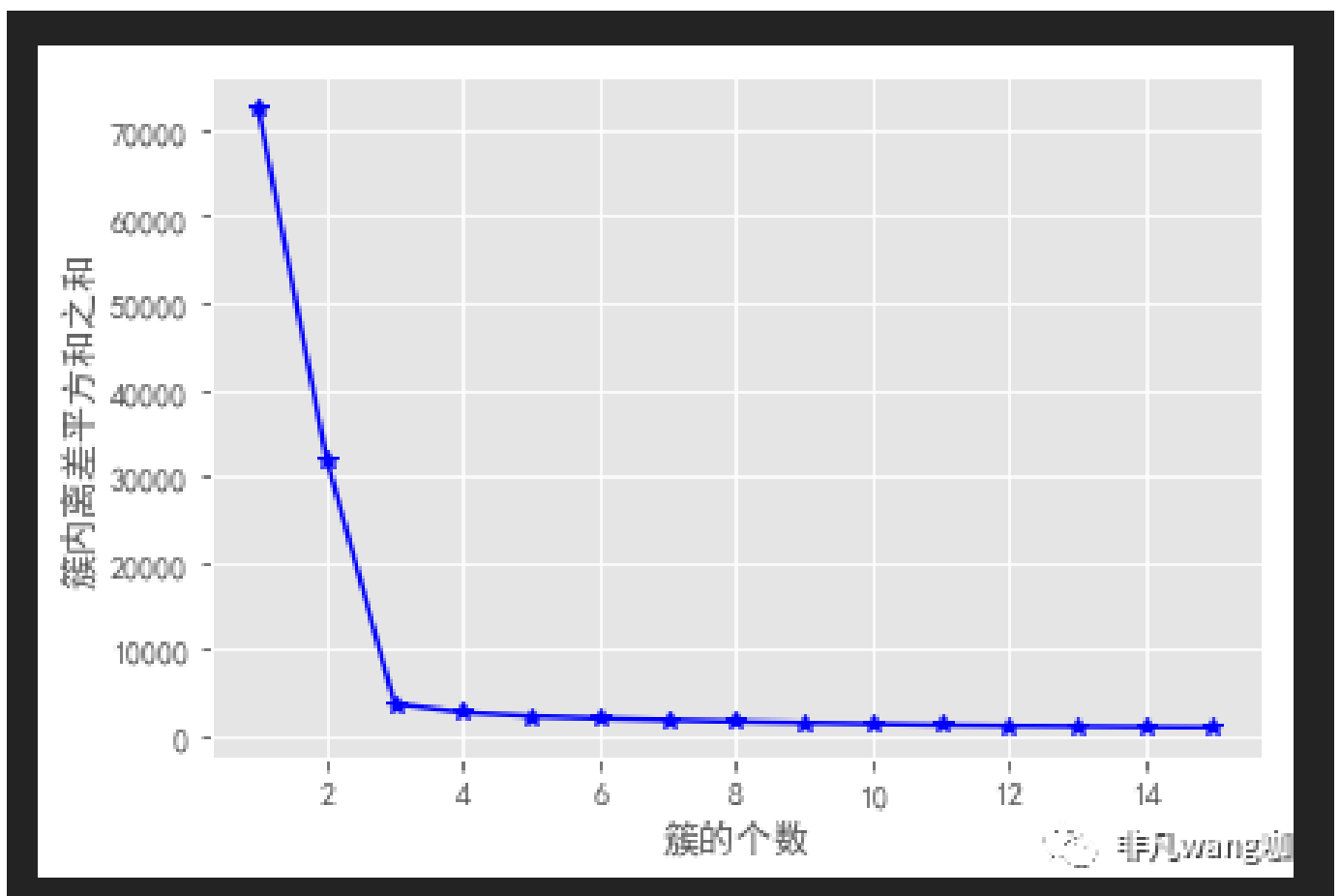
1  #构造自定义函数，用于绘制不同的k值和对应总的簇内离差平方和的折线图
2  def k_SSE(X,clusters):
3      #选择连续的K种不同的值
4      K= range(1,clusters+1)
5      #构建空列表用于存储总的簇内离差平方和
6      TSSE= []
7      for k in K:
8          #用于存储各个簇内离差平方和
9          SSE = []
10         kmeans = KMeans(n_clusters=k)
11         kmeans.fit(X)
12         #返回簇标签
13         labels = kmeans.labels_
14         #返回簇中心
15         centers = kmeans.cluster_centers_
16         #计算各簇样本的离差平方和，并保存到列表中
17         for label in set(labels):
18             SSE.append(np.sum((X.loc[labels==label,]-centers[label,:])**2))
19         #计算总的簇内离差平方和

```

```

20     TSSE.append(np.sum(SSE))
21     #中文和负号正常显示
22     plt.rcParams['font.sans-serif'] = 'SimHei'
23     plt.rcParams['axes.unicode_minus'] = False
24     #设置绘画风格
25     plt.style.use('ggplot')
26     #绘制K的个数与TSSE的关系
27     plt.plot(K, TSSE, 'b*-')
28     plt.xlabel('簇的个数')
29     plt.ylabel('簇内离差平方和之和')
30     #
31     plt.show()
32     #将三组数据集汇总到数据框中
33     X = pd.DataFrame(np.concatenate([np.array([x1,y1]),np.array([x2,y2]),np.array([x3,y3])]),
34                       columns=['x','y'])
35     k_SSE(X, 15)

```



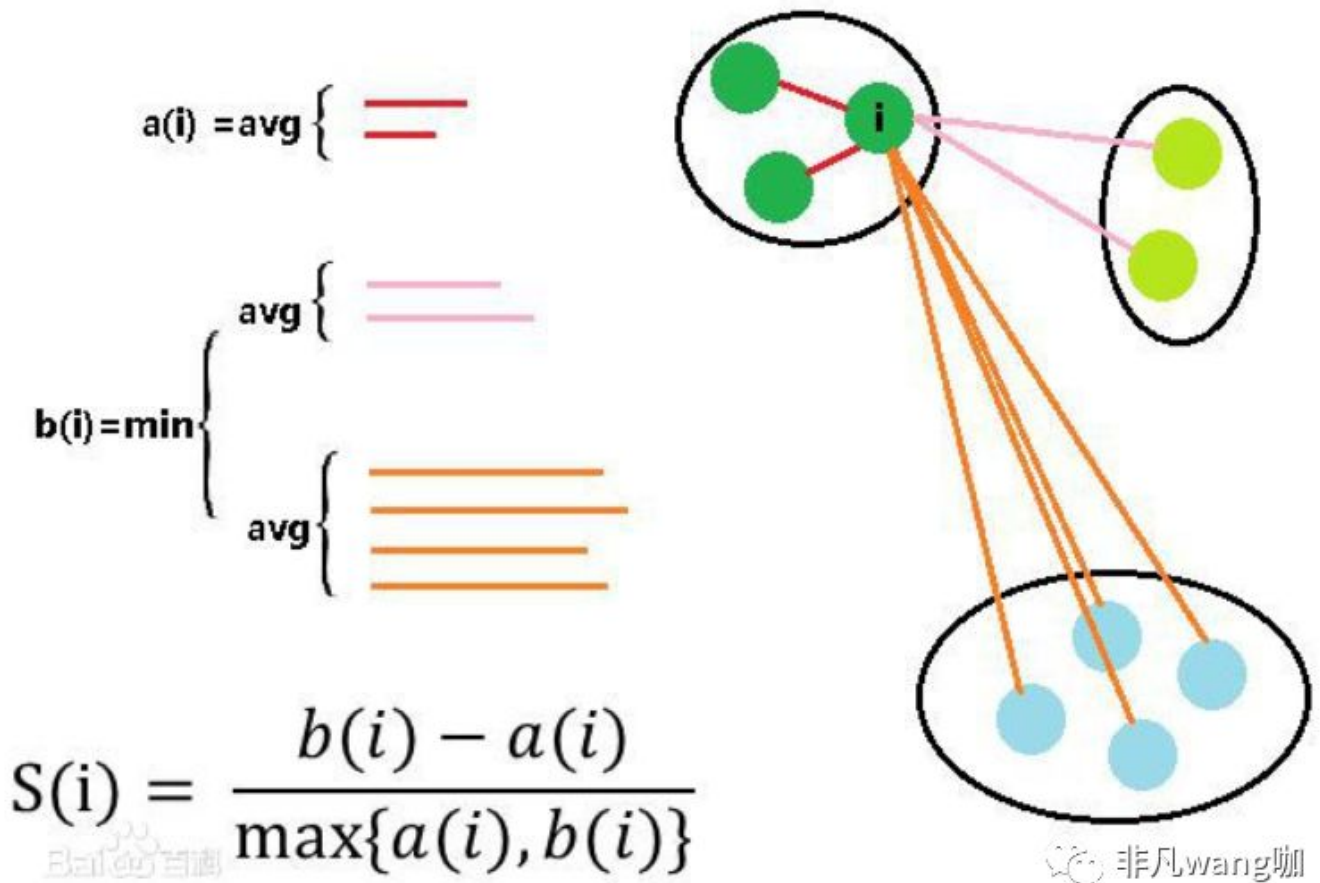
如上图，当簇的个数为3时，形成了一个明显的“拐点”，因为K值从1到3时，折线的斜率都比较大，但是k值为4时斜率突然就降低了很多，并且之后的簇对应的斜率都变动很小，所以，合理的k值应该为3，与虚拟数据集的三个簇相吻合。

轮廓系数法

$$s(i) = \frac{b(i) - a(i)}{\max\{a(i), b(i)\}} \quad s(i) = \begin{cases} 1 - \frac{a(i)}{b(i)}, & a(i) < b(i) \\ 0, & a(i) = b(i) \\ \frac{b(i)}{a(i)} - 1, & a(i) > b(i) \end{cases}$$

啥都不说先上公式，该方法综合考虑了簇的密集性和分散性两个信息，如果数据集被分割为理想的K个簇，那么对应的簇内样本会很密集，而簇间样本会很分散。上述公式中a(i)体现了簇内的密集性，代表样本i与同簇内其他样本点距离的平均值；b(i)反映了簇间的分散性，他的计算过程是样本i与其他非同簇样本点距离的平均值，然后从平均值中挑选出最小值。

通过公式可知当S(i)接近于-1时，说明样本i分配的不合理，需要将其分配到其他簇中；当S(i)近似为0时，说明样本i落在了模糊地带，即簇的边界处，当S(i)近似为1时，说明样本i的分配是合理的。



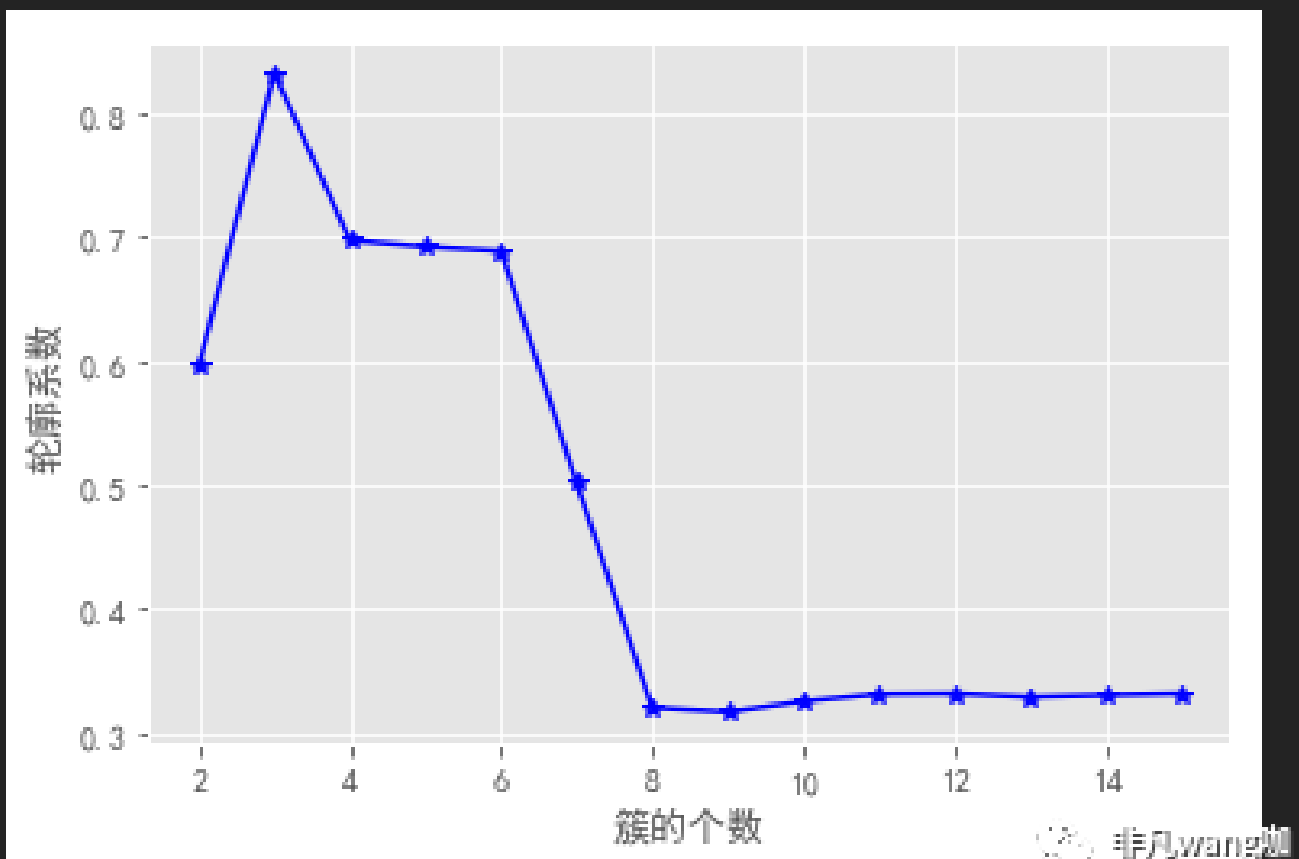
上图是百度百科对该方法的解释：

<https://baike.baidu.com/item/轮廓系数/17361607?fr=aladdin>

接下来我们就看看如何用轮廓系数解决我们的k取值问题，由于轮廓系数计算较复杂，所以我们直接使用sklearn中的metrics中的silhouette_score方法，需要注意的是该方法需要接受的聚类簇数必须大于等于2.代码如下：

```
1 from sklearn import metrics
2 #构造自定义函数
3 def k_silhouette(X,clusters):
4     K = range(2,clusters+1)
5     #构建空列表，用于存储不同簇数下的轮廓系数
6     S = []
7     for k in K:
8         kmeans = KMeans(n_clusters=k)
9         kmeans.fit(X)
10        labels = kmeans.labels_
11        #调用子模块metrics中的silhouette_score函数，计算轮廓系数
12        S.append(metrics.silhouette_score(X,labels,metric='euclidean'))
13    #设置绘图风格
14    plt.rcParams['font.sans-serif'] = 'SimHei'
```

```
15 plt.rcParams['axes.unicode_minus'] = False
16 # 设置绘画风格
17 plt.style.use('ggplot')
18 # 绘制K的个数与轮廓系数的关系
19 plt.plot(K,S, 'b*-')
20 plt.xlabel('簇的个数')
21 plt.ylabel('轮廓系数')
22 plt.show()
23 k_silhouette(X,15)
```



如上图，利用之前构造的虚拟数据，绘制了不同K值下对应的轮廓系数图，当k取值为3时轮廓系数最大，且比较接近于1，说明应该把虚拟数据聚为3类比较合理。

间隔统计法


2000年Hastie等人提出了间隔统计量法（Gap Statistic方法），该方法可以适用与任何聚类算法，公式如下：

最简单的方法是使用类内样本点之间的欧式距离来表示，记为 D_k ， D_k 越小聚类的紧支性越好。Ref

$$D_k = \sum_{x_i \in C_k} \sum_{x_j \in C_k} \|x_i - x_j\|^2 = 2n_k \sum_{x_i \in C_k} \|x_i - \mu_k\|^2$$

标准化后：

$$W_k = \sum_{k=1}^K \frac{1}{2n_k} D_k$$

 非凡wang咖

详情参考地址：

https://blog.csdn.net/baidu_17640849/article/details/70769555

接下来我们构造自定义函数，绘制不同K值对应的间隙统计量折线图：

```

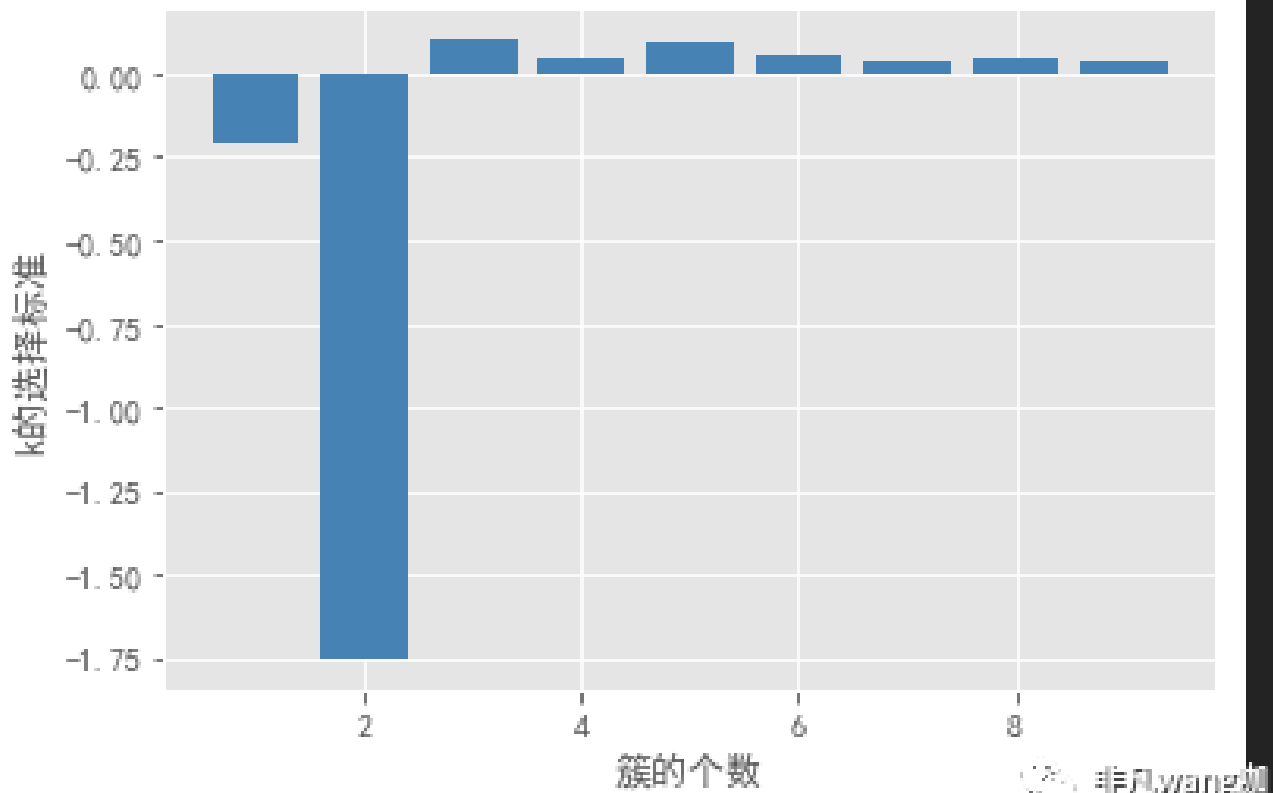
1  #自定义函数，计算簇内任意俩样本之间的欧式距离Dk
2  def short_pair_wise_D(each_cluster):
3      mu = each_cluster.mean(axis=0)
4      Dk = sum(sum((each_cluster-mu)**2*each_cluster.shape[0]))
5      return Dk
6  #自定义函数，计算簇内的Wk值
7  def compute_Wk(data, classfication_result):
8      Wk = 0
9      label_set = set(classfication_result)
10     for label in label_set:
11         each_cluster = data[classfication_result == label, :]
12         Wk = Wk + short_pair_wise_D(each_cluster)/(2.0*each_cluster.shape[0])
13     return Wk
14
15 # 计算GAP统计量
16 def gap_statistic(X, B=10, K=range(1,11), N_init = 10):
17     # 将输入数据集转换为数组
18     X = np.array(X)
19     # 生成B组参照数据
20     shape = X.shape
21     tops = X.max(axis=0)
22     bots = X.min(axis=0)
23     dists = np.matrix(np.diag(tops-bots))
24     rands = np.random.random_sample(size=(B,shape[0],shape[1]))

```



```
25     for i in range(B):
26         rands[i,:,:] = rands[i,:,:]*dists+bots
27
28     # 自定义0元素的数组，用于存储gaps、wks和wkbs
29     gaps = np.zeros(len(K))
30     wks = np.zeros(len(K))
31     wkbs = np.zeros((len(K),B))
32     # 循环不同的k值，
33     for idxk, k in enumerate(K):
34         k_means = KMeans(n_clusters=k)
35         k_means.fit(X)
36         classification_result = k_means.labels_
37         # 将所有簇内的wk存储起来
38         wks[idxk] = compute_wk(X,classification_result)
39
40         # 通过循环，计算每一个参照数据集下的各簇wk值
41         for i in range(B):
42             Xb = rands[i,:,:]
43             k_means.fit(Xb)
44             classification_result_b = k_means.labels_
45             wkbs[idxk,i] = compute_wk(Xb,classification_result_b)
46
47     # 计算gaps、sd_ks、sk和gapDiff
48     gaps = (np.log(wkbs)).mean(axis = 1) - np.log(wks)
49     sd_ks = np.std(np.log(wkbs), axis=1)
50     sk = sd_ks*np.sqrt(1+1.0/B)
51     # 用于判别最佳k的标准，当gapDiff首次为正时，对应的k即为目标值
52     gapDiff = gaps[:-1] - gaps[1:] + sk[1:]
53
54     #设置绘图风格
55     plt.rcParams['font.sans-serif'] = 'SimHei'
56     plt.rcParams['axes.unicode_minus'] = False
57     #设置绘画风格
58     plt.style.use('ggplot')
59     # 绘制gapDiff的条形图
60     plt.bar(np.arange(len(gapDiff))+1, gapDiff, color = 'steelblue')
61     plt.xlabel('簇的个数')
62     plt.ylabel('k的选择标准')
63     plt.show()
64
```

```
65 # 自定义函数的调用
66 gap_statistic(X)
```



如上图，x轴代表了不同的簇数k，y轴代表k值选择的判断指标gapDiff，gapDiff首次出现正值时对应的k为3，所以对于虚拟的数据集来说，将其划分为三个簇是比较合理的。

以上就是Kmeans算法中簇数k的确定方法，嗯，是有点小小的难度😁

喜欢此内容的人还喜欢

山东栖霞金矿重大爆炸事故调查处理结果公布 45人被追责问责

中央纪委国家监委网站

新手上路，请别超速！

中华全国学联