

机器学习

Lecture 9

KMeans

一、 算法

Algorithm 1 KMeans**Input:** 数据集 $D = \{\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_m\}$; 聚类簇数 k

```

1: 初始化  $k$  个均值向量  $\{\mu_1, \mu_2, \dots, \mu_k\}$ 
2: repeat
3:   for  $i = 1, 2, \dots, m$  do
4:     for  $j = 1, 2, \dots, k$  do
5:       Compute  $d_{ij} = d(\mathbf{x}_i, \mu_j)$ 
6:     end for
7:     确定第  $i$  个样本的簇标记  $\lambda_i = \arg \min_j d(ij)$ 
8:     将  $\mathbf{x}_i$  划分到对应的簇:  $C_{\lambda_j} = C_{\lambda_j} \cup \{\mathbf{x}_i\}$ 
9:   end for
10:  for  $j = 1, 2, \dots, k$  do
11:    Update  $\mu_j = \frac{1}{|C_{\lambda_j}|} \sum_{\mathbf{x}_i \in C_{\lambda_j}} \mathbf{x}_i$ 
12:  end for
13: until 所有均值向量均未被更新

```

二、 理论

Theorem 1. K -means 算法等价于对如下目标函数进行优化

$$J = \frac{1}{m} \sum_{i=1}^m d(\mathbf{x}_i, \mu_j)$$

当采用欧式距离时, $\min J$ 等价于

$$\min \frac{1}{2} \sum_{i=1}^m \sum_{j=1}^k r_{ij} \|\mathbf{x}_i - \mu_j\|^2 \quad (1)$$

其中 $r_{ij} = 0, 1, \sum_{j=1}^k r_{ij} = 1$ 使用 BCD 算法对 (1) 式进行求解, 即先在固定 r 的情况下更新 μ , 再在固定 μ 的情况下更新 r

$$\begin{aligned} \nabla_{\mu_j} l(R, \mu) &= -2 \sum_{i=1}^m r_{ij} (\mathbf{x}_i - \mu_j) = 0 \\ \therefore \mu_j &= \frac{\sum_{i=1}^m r_{ij} \mathbf{x}_i}{\sum_{i=1}^m r_{ij}} = \frac{1}{m} \sum_{i=1}^m r_{ij} \mathbf{x}_i \end{aligned}$$

$$\begin{aligned}
 r_{ij} &= \arg \min_{r_{ij}} \sum_{i=1}^m \sum_{j=1}^k r_{ij} \|\mathbf{x}_i - \mu_j\|^2 \\
 &= \arg \min_{r_{ij}} \sum_{j=1}^k r_{ij} \|\mathbf{x}_i - \mu_j\|^2 \\
 &= \begin{cases} 1 & \arg \min_{1 \leq j \leq k} \|\mathbf{x}_i - \mu_j\|^2 \\ 0 & \text{otherwise} \end{cases}
 \end{aligned}$$

三、 K-Means 迭代过程

构造中心为 $[-1,1],[2,-2],[-2,-3]$ ，标准差为 0.6，样本数量为 100 的数据集。经过三轮迭代，算法收敛，结果如图1所示。

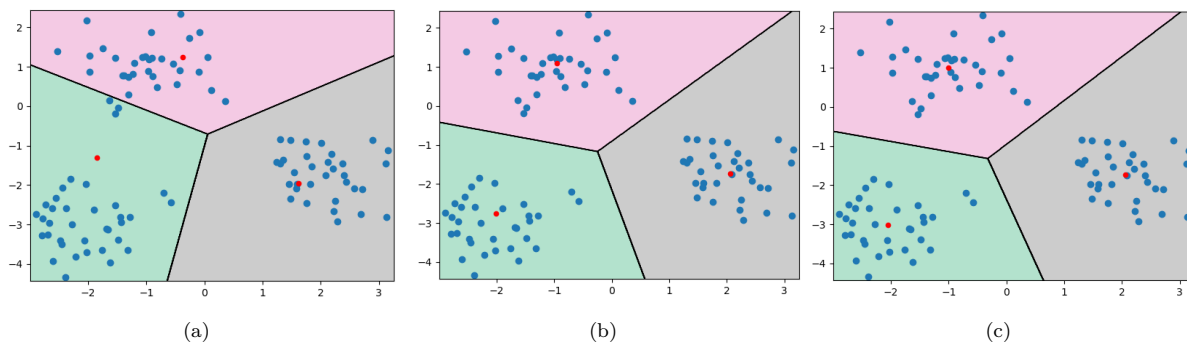


图 1: K-Means 迭代过程.

四、 通过手肘法挑选 k

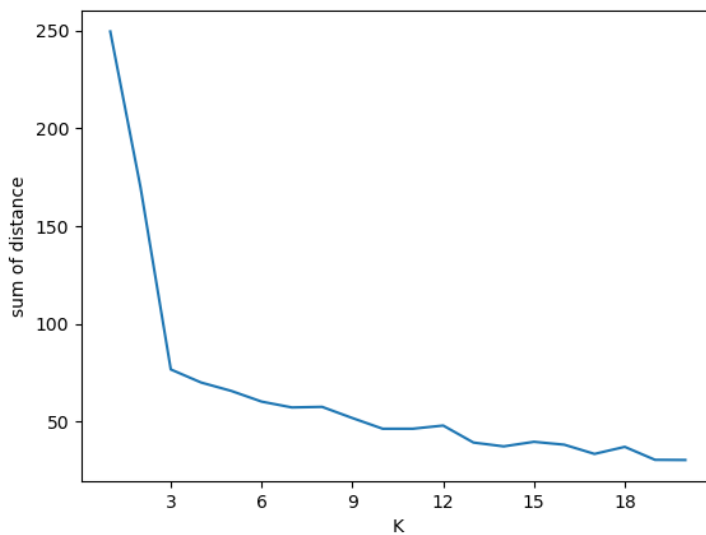


图 2: 手肘法.

手肘法的核心思想是：随着聚类数 k 的增大，样本划分会更加精细，每个簇的聚合程度会逐渐提高，那么误差平方和 SSE 自然会逐渐变小。并且，当 k 小于真实聚类数时，由于 k 的增大会大幅增加每个簇的聚合程度，故 SSE 的下降幅度会很大，而当 k 到达真实聚类数时，再增加 k 所得到的聚合程度回报会迅速变小，所以 SSE 的下降幅度会骤减，然后随着 k 值的继续增大而趋于平缓，也就是说 SSE 和 k 的关系图是一个手肘的形状，而这个肘部对应的 k 值就是数据的真实聚类数。

从图2中可以看出， $k = 3$ 为一个拐点，因此 k 取 3 时效果最优，这与构造的数据相吻合。

五、 代码实现

```
1 import numpy as np
2 from numpy.linalg import norm
3 import matplotlib.pyplot as plt
4 from tqdm import tqdm
5 from matplotlib.ticker import MaxNLocator
6
7 class KMeans:
8     def __init__(self, x, k, resolution=1000, plot=True):
9         self.x = np.array(x)
10        self.k = k
11        self.m = len(x)
12        self.meanVectors = self.x[np.random.choice(self.m,k,replace=False)]
13        self.clusters = [[]for i in range(self.k)]
14        self.resolution = resolution
15        self.plot = plot
16
17    def train(self):
18        '''
19        KMeans 迭代训练
20        :return: 均值向量
21        '''
22        cnt = 0
23        while True:
24            # ----- #
25            # 初始化所有簇
26            # ----- #
27            self.clusters = [[]for i in range(self.k)]
28            # ----- #
29            # 计算最近簇
30            # ----- #
31            for i in range(self.m):
32                d = float('inf')
33                index = 0
34                for j in range(self.k):
35                    d_ij = norm(self.x[i]-self.meanVectors[j])
36                    if d_ij < d:
37                        d = d_ij
38                        index = j
39                self.clusters[index].append(self.x[i])
```

```
40     # ----- #
41     # 更新均值向量
42     # ----- #
43     done = True
44     for i in range(self.k):
45         if len(self.clusters[i])!=0:
46             newMeanVector = np.average(np.array(self.clusters[i]),axis=0)
47         else:
48             newMeanVector = self.meanVectors[i]
49         if not np.array_equal(newMeanVector,self.meanVectors[i]):
50             self.meanVectors[i] = newMeanVector
51         done = False
52     # print('# ----- #')
53     # print(self.meanVectors)
54     # ----- #
55     # 绘制决策边界
56     # ----- #
57     if self.plot:
58         self.plot_decision_boundaries(iteration=cnt,resolustion=self.resolustion)
59     cnt += 1
60     # ----- #
61     # 判断结束条件 退出循环
62     # ----- #
63     if done:
64         return self.meanVectors
65
66     def predict(self, x):
67         '''
68         预测
69         :param x: 某个给定向量
70         :return: 距离x最近的均指向量的编号
71         '''
72         d = float('inf')
73         index = 0
74         for i in range(self.k):
75             d_ij = norm(x-self.meanVectors[i])
76             if d_ij < d:
77                 d = d_ij
78                 index = i
79         return index
80
81     def select_k(self, start=1, end=20, plot=False):
82         '''
83         手肘法确定k的取值
84         :return:
85         '''
86         all_distance = []
87         for k in range(start,end+1):
88             self.__init__(self.x, k, plot=plot)
89             self.train()
```

```
90     all_distance.append(self.sumDistance())
91     plt.plot(range(1,len(all_distance)+1),all_distance)
92     plt.ylabel('sum of distance')
93     plt.xlabel('K')
94     plt.gca().xaxis.set_major_locator(MaxNLocator(integer=True))
95     plt.show()
96
97
98     def sumDistance(self):
99         '''
100         计算点到聚点的距离之和
101         :return:
102         '''
103         sumdistance = 0
104         for i in range(self.m):
105             d = float('inf')
106             for j in range(self.k):
107                 d_ij = norm(self.x[i]-self.meanVectors[j])
108                 if d_ij<d:
109                     d = d_ij
110             sumdistance += d
111         return sumdistance
112
113     def plot_decision_boundaries(self, resolution=1000, iteration=0):
114         '''
115         绘制决策边界
116         :param resolution: 网格密度
117         :param iteration: 迭代次数
118         :return:
119         '''
120         mins = self.x.min(axis=0) - 0.1
121         maxs = self.x.max(axis=0) + 0.1
122         xx, yy = np.meshgrid(np.linspace(mins[0], maxs[0], resolution),
123                               np.linspace(mins[1], maxs[1], resolution))
124         grid = np.c_[xx.ravel(), yy.ravel()]
125         predict = []
126         for i in tqdm(grid):
127             predict.append(self.predict(i))
128         predict = np.array(predict)
129         predict = predict.reshape(xx.shape)
130         # print(predict)
131
132         plt.contourf(predict, extent=(mins[0], maxs[0], mins[1], maxs[1]),
133                      cmap='Pastel2')
134
135         plt.contour(predict, extent=(mins[0], maxs[0], mins[1], maxs[1]),
136                    linewidths=1, colors='k')
137
138         plt.scatter(self.x[:,0],self.x[:,1])
139         plt.scatter(self.meanVectors[:,0],self.meanVectors[:,1],s=20,c='r')
```

```
140     # plt.show()
141     plt.savefig('kmeans_{}.png'.format(iteration))
142     plt.close()
```