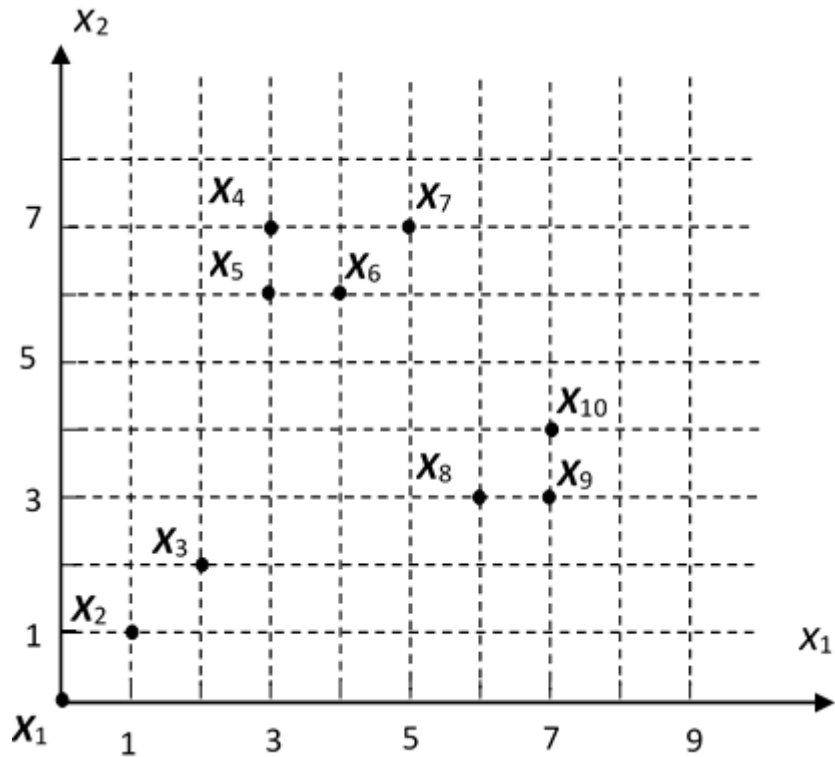


2016302580320 任思远 第二次作业

2.1



1. 选择 X_1 作为 Z_1 。
2. 选择距离最远的 X_7 作为 Z_2 。
3. 记 $M_i = \min\{\|X_i - Z_1\|, \|X_i - Z_2\|\}$, 经计算:

$$M_2 = \sqrt{2}$$

$$M_3 = 2\sqrt{2}$$

$$M_4 = 2$$

$$M_5 = \sqrt{5}$$

$$M_6 = \sqrt{2}$$

$$M_8 = \sqrt{17}$$

$$M_9 = \sqrt{20}$$

$$M_{10} = \sqrt{13}$$

$M_9 > \theta \|Z_1 - Z_2\|$, 选择 M_9 作为 Z_3 , 再无满足的点, 结束聚类中心的选择。

4. 划分:

$\{X_1, X_2, X_3\}, \{X_7, X_4, X_5, X_6\}, \{X_9, X_8, X_{10}\}$

2.2

k-means算法流程:

input: 样本集 $D = \{x_1, x_2, \dots, x_n\}$, 聚类簇数 k .

output: 聚类划分的集合。

算法流程:

1. 从 D 中随机选择 k 个作为初始均值向量 $\{\mu_1, \mu_2, \dots, \mu_k\}$.
2. loop until 均值向量不再变化 (或变化不大, 或其他近似条件)

令 $C_i = \emptyset$, $i = 1, \dots, k$ 表示每个聚类的集合,

for $j = 1, \dots, n$:

计算每个向量和均值向量的距离

$$d_{ji} = \text{dist}(x_j, \mu_i), \quad i = 1, \dots, k$$

选择最近的簇, 并加入 $\lambda_j = \text{argmin}_i \{d_{ji}\}$, $C_j \cup = x_j$

重新计算均值向量:

for $j = 1, \dots, k$:

$$\mu_j = \frac{1}{|C_j|} \sum_{x \in C_j} x$$

3. 退出 2 中的 loop 后, 得到的 $\{C_i\} \quad i = 1, \dots, k$ 即为聚类划分。

[C++ 实现 k-means](#) C++代码太长了, 这里放不下。

[Python 实现 k-means](#) Python代码:

```
1 import numpy as np
2 from functools import reduce
3
4
5 def k_means(cluster_num: int, dist, *data_set):
6     size = len(data_set)
7     assert size >= cluster_num > 0
8     s = set()
9     while len(s) != cluster_num:
10         s.add(np.random.random_integers(0, size - 1))
11
```

```

12     mean_vecs = [] # mean vectors in each cluster,
    initialized with random vector selected in data set.
13     for num in s:
14         mean_vecs.append(data_set[num])
15
16     C = [] # clusters
17     for i in range(cluster_num):
18         C.append([])
19
20     isChanged = True
21
22     '''
23     Divide each vector into one of the sets,
24     and calculate the new mean_vecs.
25
26     Loop until mean_vecs have no change.
27     '''
28     while isChanged:
29         isChanged = False
30         for c in C:
31             c.clear()
32
33         # for each vector in data set, divide it into
    one of the clusters.
34         for vec in data_set:
35             d = dist(vec, mean_vecs[0])
36             label = 0
37             for l in range(cluster_num):
38                 cur_dist = dist(vec, mean_vecs[l])
39                 if cur_dist < d:
40                     d = cur_dist

```

```
41         label = 1
42         C[label].append(vec)
43
44         # update mean_vecs, and set up `isChanged`
flag if needed
45         for label in range(cluster_num):
46             assert len(C[label]) > 0
47             mean_vec = reduce(lambda x, y: x + y,
C[label]) / len(C[label])
48             if dist(mean_vec, mean_vecs[label]) >
10e-2:
49                 mean_vecs[label] = mean_vec
50                 isChanged = True
51
52     return C
```