

МОСКОВСКИЙ ГОСУДАРСТВЕННЫЙ ТЕХНИЧЕСКИЙ
УНИВЕРСИТЕТ ИМЕНИ Н. Э. БАУМАНА
Факультет информатики и систем управления
Кафедра теоретической информатики и компьютерных технологий

Лабораторная работа №14
по курсу «Информационный поиск»
«Агломеративная кластеризация»

Лабораторная работа №15
по курсу «Информационный поиск»
«K-means»

Выполнил:
студент группы ИУ9-21М
Беляев А. В.

Проверила:
Лукашевич Н. В.

Москва 2019

1 Цель работы

Для заданного набора точек построить агломеративную кластеризацию по методу single-link и complete-link.

Для этого же набора точек и заданных двух изначальных точек построить кластеризацию методом K-means, выбрав в качестве начальных кластера 2 указанные ранее точки.

Сравнить методы.

2 Агломеративная кластеризация

2.1 Ход работы

```
1 import matplotlib.pyplot as plot
2 import numpy as np
3 import scipy.cluster.hierarchy as sch
4
5
6 POINTS = [
7     (0.6, 1.9),
8     (1.8, 1.6),
9     (2.7, 2.0),
10    (3.0, 2.1),
11    (3.0, 2.6),
12    (3.1, 4.5),
13    (3.8, 0.6),
14    (4.2, 2.7)
15 ]
16
17
18 class Cluster:
19     def __init__(self):
20         self.points = []
21
22     def add_point(self, point: tuple):
23         self.points.append(point)
24
25     def coordinates(self):
26         return list(point[0] for point in self.points), list(point[1]
27             ↪ for point in self.points)
28
29     def __repr__(self):
30         return self.__str__()
31
32     def __str__(self):
33         return f'Cluster {self.points}'
```

```

34
35 def single_link(distance_function):
36     linkage = sch.linkage(distance_function, method='single')
37     return sch.fcluster(linkage, criterion='maxclust', t=2)
38
39
40 def complete_link(distance_function):
41     linkage = sch.linkage(distance_function, method='complete')
42     return sch.fcluster(linkage, criterion='maxclust', t=2)
43
44
45 def main():
46     distance_function = sch.distance.pdist(np.array(POINTS),
47     ↪ metric='cosine')
48     clusterized = single_link(distance_function)
49     # clusterized = complete_link(distance_function)
50
51     # (-1) for each index since numeration after `fclutser` starts from
52     ↪ 1
53     clusterized = [cluster_indx - 1 for cluster_indx in clusterized]
54
55     clusters_num = max(clusterized) + 1
56     clusters = [Cluster() for _ in range(clusters_num)]
57
58     i = 0
59     for cluster_index in clusterized:
60         clusters[cluster_index].add_point(POINTS[i])
61         i += 1
62
63     for c in clusters:
64         xs, ys = c.coordinates()
65         plot.scatter(xs, ys)
66     plot.show()
67
68 if __name__ == '__main__':
69     main()

```

Результаты кластеризации по методу single-link представлены на Рисунке 1. Complete-link – на Рисунке 2. Используемая мера близости - косинусная мера. Функция `fcluster` пакета `scipy` позволяет разбить данное множество на t кластеров. В данном случае $t = 2$.

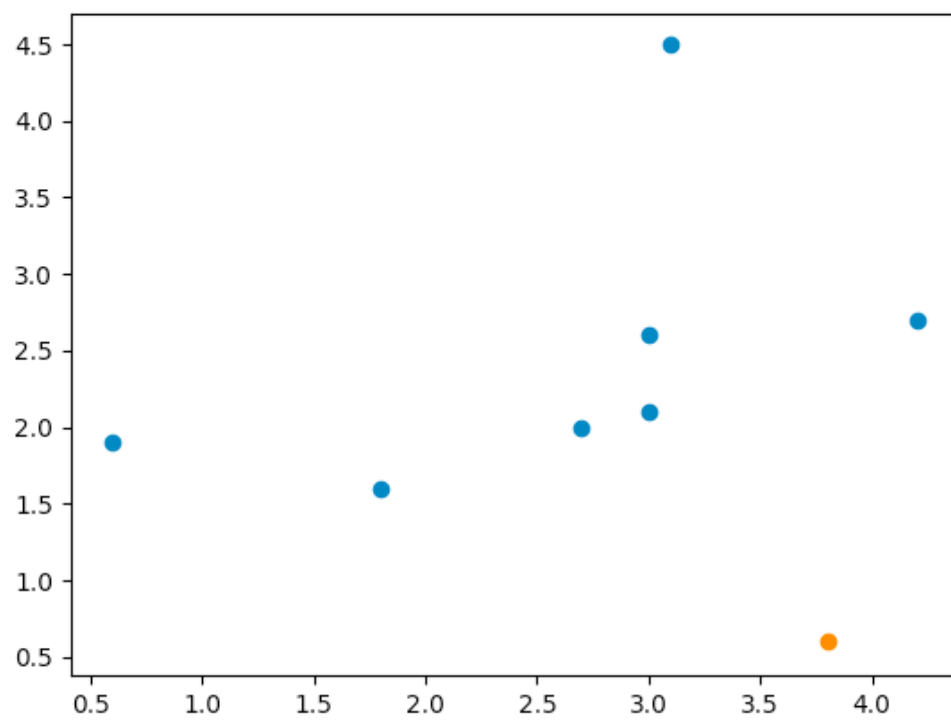


Рис. 1: Single-link

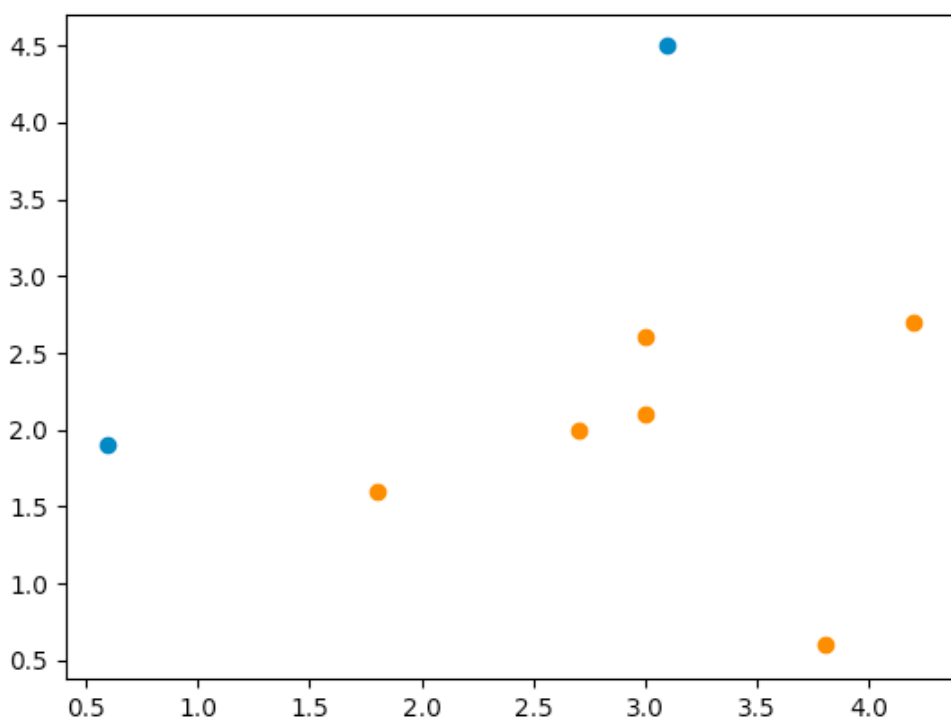


Рис. 2: Complete-link

3 Кластеризация K-means

3.1 Ход работы

```
1 import math
2 import matplotlib.pyplot as plot
3 from sklearn.metrics.pairwise import cosine_similarity
4
5
6 POINTS = [
7     (0.6, 1.9),
8     (1.8, 1.6),
9     (2.7, 2.0),
10    (3.0, 2.1),
11    (3.0, 2.6),
12    (3.1, 4.5),
13    (3.8, 0.6),
14    (4.2, 2.7)
15 ]
16
17 INIT_CLUSTER_CENTERS = [
18     (1.8, 1.6),
19     (3.0, 2.6)
20 ]
21
22
23 class Point:
24     def __init__(self, x: float, y: float):
25         self.x = x
26         self.y = y
27         self.cluster_index = None
28
29     def to_vector(self):
30         return [self.x, self.y]
31
32     def __repr__(self):
33         return self.__str__()
34
35     def __str__(self):
36         return f'[{self.x} {self.y}]'
37
38
39 class Cluster:
40     def __init__(self, center: Point):
41         self.points = []
```

```

42         self.center = center
43
44     def add_point(self, p: Point):
45         self.points.append(p)
46
47     def remove_point(self, p: Point):
48         if p in self.points:
49             self.points.remove(p)
50
51     def update_center(self):
52         if 0 != len(self.points):
53             x_avg = sum(p.x for p in self.points) / len(self.points)
54             y_avg = sum(p.y for p in self.points) / len(self.points)
55             self.center = Point(x_avg, y_avg)
56
57     def coordinates(self):
58         return list(point.x for point in self.points), list(point.y for
59             ↪ point in self.points)
60
61     def __repr__(self):
62         return self.__str__()
63
64     def __str__(self):
65         return f'Cluster {len(self.points)}: {self.center}'
66
67 def reshuffle(points: list, clusters: list) -> bool:
68     def cosine_dist(a: Point, b: Point):
69         return cosine_similarity([a.to_vector()], [b.to_vector()])
70
71     # not used in favor of cosine_sim
72     def euclidean_dist(a: Point, b: Point):
73         x = math.fabs(a.x - b.x)
74         y = math.fabs(a.y - b.y)
75         return math.sqrt(x**2 + y**2)
76
77     # algorithm converges when there are no more changes in clusters
78     smth_has_changed = False
79     for p in points:
80
81         min_dist = float('inf')
82         new_cluster_index = 0
83
84         # which cluster the point should be added to
85         for j in range(len(clusters)):

```

```

86         curr_center = clusters[j].center
87         distance = cosine_dist(p, curr_center)
88         if distance <= min_dist:
89             new_cluster_index = j
90             min_dist = distance
91
92     if p.cluster_index != new_cluster_index:
93         # remove point from old cluster
94         if p.cluster_index is not None:
95             clusters[p.cluster_index].remove_point(p)
96
97         # append point to new cluster
98         clusters[new_cluster_index].add_point(p)
99         p.cluster_index = new_cluster_index
100         smth_has_changed = True
101
102     return smth_has_changed
103
104
105 def k_means_clusterize(points: list, initial_centers: list):
106     clusters = [Cluster(initial_center) for initial_center in
107         ↪ initial_centers]
108
109     not_converged = True
110     while not_converged:
111         # assign points to nearest clusters
112         not_converged = reshuffle(points, clusters)
113
114         # update cluster centers
115         [c.update_center() for c in clusters]
116         print(clusters)
117
118     return clusters
119
120 def main():
121     points = [Point(p[0], p[1]) for p in POINTS]
122     centers = [Point(p[0], p[1]) for p in INIT_CLUSTER_CENTERS]
123
124     clusters = k_means_clusterize(points, centers)
125
126     # draw clusters
127     for c in clusters:
128         xs, ys = c.coordinates()
129         plot.scatter(xs, ys)

```

```

130
131     # draw cluster centers
132     cxs, cys = list(c.center.x for c in clusters), list(c.center.y for c
↵      in clusters)
133     plot.scatter(cxs, cys, color='black')
134     plot.show()
135
136
137 if __name__ == '__main__':
138     main()

```

Метод сошелся за 2 итерации. Результаты кластеризации представлены на Рисунке 3. Цетры кластеров обозначены черным цветом.

Используемая мера близости - косинусная мера.

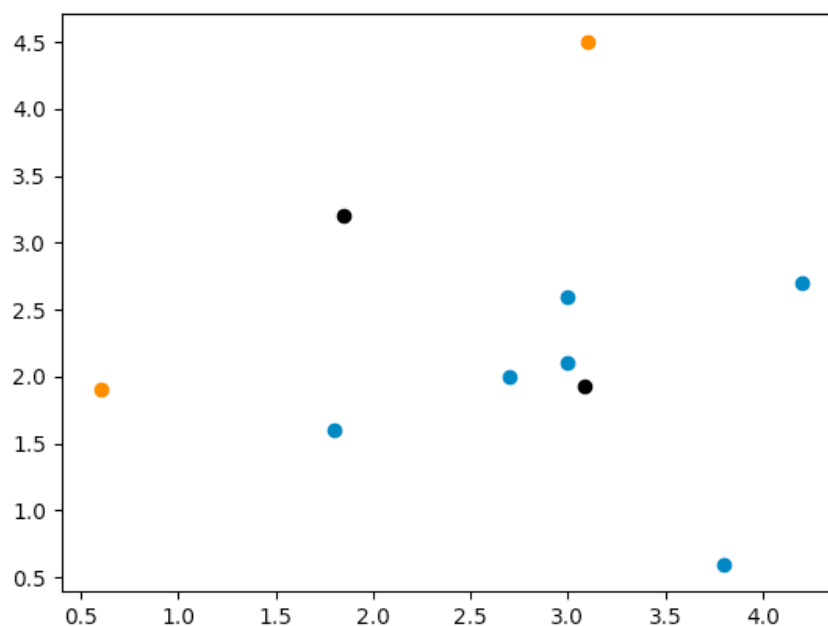


Рис. 3: K-means

4 Выводы

Были реализованы 2 способа кластеризации – агломеративная (с помощью single-link и complete-link) и K-means.

Метод кластеризации K-means предельно прост и понятен и может быть реализован вручную за адекватное время.

В то же время агломеративная кластеризация требует больших усилий для ручной реализации. Поэтому был использован пакет `scipy`. Главным недостатком подобных пакетов (`scipy`, `scikit` и прочие) является то, что реализация метода скрыта

глубоко в коде (в данном случае пакет представляет лишь обертку над реализацией на языке *C*) и не может быть изучена, а управление методами реализовано посредством большого количества параметров. Таким образом, для простейшего использования метода необходимо изучить все параметры, а затем, ввиду их большого количества, подбирать их практически случайным образом. При этом, параметры позволяют лишь «подогнать» получаемый результат под ожидаемый в конкретном случае результат и, с большой долей вероятности, на другом наборе данных работать не будут.

После определенного количества попыток, подобные параметры для методов кластеризации из пакета *scipy* были подобраны и кластеризация прошла успешно. Тенденция к «широким» и «узким» кластерам отражена на соответствующих рисунках агломеративной кластеризации. Она некоторым образом коррелирует с кластеризацией K-means.