

Федеральное государственное бюджетное образовательное учреждение высшего
профессионального образования
Московский государственный технический университет имени Н.Э. Баумана

Домашняя работа №2
«Кластерный анализ методом k-средних»
по курсу:
«Моделирование»

Выполнил:
студент группы ИУ9-82
Иванов Георгий

Проверила:
Домрачева А.Б.

Москва, 2019

Содержание

1	Цель задачи	3
2	Теоретические сведения	4
2.1	Кластерный анализ (кластеризация)	4
2.2	Метод k-средних	4
3	Практическая реализация	5
4	Результаты	9
5	Вывод	10

1 Цель задачи

Имеются данные с информацией о ВИЧ-зараженных в различных странах за временной промежуток с 2000 года по 2014 год.

Необходимо исследовать связь между оценкой страны по ВИЧ (для каждой страны оценка вычисляется в зависимости от населения, уровня страны и количество зараженных) и количеством ВИЧ-зараженных с помощью кластерного анализа.

2 Теоретические сведения

2.1 Кластерный анализ (кластеризация)

Кластерный анализ — многомерная статистическая процедура, выполняющая сбор данных, содержащих информацию о выборке объектов, и затем упорядочивающая объекты в сравнительно однородные группы. Внутри каждой группы должны оказаться «похожие» объекты, а объекты разных группы должны быть как можно более отличны.

2.2 Метод k-средних

Метод k-средних наиболее популярный метод кластеризации. Целью этого метода является разделение m наблюдений (из пространства R^n) на k кластеров ($k \leq m$), при этом каждое наблюдение относится к тому кластеру, к центру (центроиду) которого оно ближе всего. В качестве меры близости используется Евклидово расстояние.

Основная идея заключается в том, что на каждой итерации перевычисляется центр масс для каждого кластера, полученного на предыдущем шаге, затем векторы разбиваются на кластеры вновь в соответствии с тем, какой из новых центров оказался ближе по выбранной метрике.

Алгоритм завершается, когда на какой-то итерации не происходит изменения внутрикластерного расстояния. Это происходит за конечное число итераций, так как количество возможных разбиений конечного множества конечно, а на каждом шаге суммарное квадратичное отклонение V уменьшается, поэтому заикливание невозможно.

3 Практическая реализация

```
1  from collections import defaultdict
2  from random import uniform
3  from math import sqrt
4  from matplotlib import pyplot as plt
5
6  import pandas as pd
7  import random
8
9
10 def point_avg(points):
11     """
12     Accepts a list of points, each with the same number of dimensions.
13     NB. points can have more dimensions than 2
14
15     Returns a new point which is the center of all the points.
16     """
17     dimensions = len(points[0])
18
19     new_center = []
20
21     for dimension in range(dimensions):
22         dim_sum = 0 # dimension sum
23         for p in points:
24             dim_sum += p[dimension]
25
26         # average of each dimension
27         new_center.append(dim_sum / float(len(points)))
28
29     return new_center
30
31
32 def update_centers(data_set, assignments):
33     """
34     Accepts a dataset and a list of assignments; the indexes
35     of both lists correspond to each other.
36
37     Compute the center for each of the assigned groups.
38
39     Return `k` centers where `k` is the number of unique assignments.
40     """
41     new_means = defaultdict(list)
42     centers = []
43     for assignment, point in zip(assignments, data_set):
44         new_means[assignment].append(point)
```

```

45
46     for points in new_means.values():
47         centers.append(point_avg(points))
48
49     return centers
50
51
52 def assign_points(data_points, centers):
53     """
54     Given a data set and a list of points between other points,
55     assign each point to an index that corresponds to the index
56     of the center point on it's proximity to that point.
57     Return a an array of indexes of centers that correspond to
58     an index in the data set; that is, if there are N points
59     in `data_set` the list we return will have N elements. Also
60     If there are Y points in `centers` there will be Y unique
61     possible values within the returned list.
62     """
63     assignments = []
64     for point in data_points:
65         shortest = float("+inf") # positive infinity
66         shortest_index = 0
67         for i in range(len(centers)):
68             val = distance(point, centers[i])
69             if val < shortest:
70                 shortest = val
71                 shortest_index = i
72         assignments.append(shortest_index)
73     return assignments
74
75
76 def distance(a, b):
77     """
78     """
79     dimensions = len(a)
80
81     _sum = 0
82     for dimension in range(dimensions):
83         difference_sq = (a[dimension] - b[dimension]) ** 2
84         _sum += difference_sq
85     return sqrt(_sum)
86
87
88 def generate_k(data_set, k):
89     """
90     Given `data_set`, which is an array of arrays,

```

```

91     find the minimum and maximum for each coordinate, a range.
92     Generate `k` random points between the ranges.
93     Return an array of the random points within the ranges.
94     """
95     centers = []
96     dimensions = len(data_set[0])
97     min_max = defaultdict(int)
98
99     for point in data_set:
100         for i in range(dimensions):
101             val = point[i]
102             min_key = 'min_%d' % i
103             max_key = 'max_%d' % i
104             if min_key not in min_max or val < min_max[min_key]:
105                 min_max[min_key] = val
106             if max_key not in min_max or val > min_max[max_key]:
107                 min_max[max_key] = val
108
109     for _k in range(k):
110         rand_point = []
111         for i in range(dimensions):
112             min_val = min_max['min_%d' % i]
113             max_val = min_max['max_%d' % i]
114
115             rand_point.append(uniform(min_val, max_val))
116
117         centers.append(rand_point)
118
119     return centers
120
121
122 def k_means(dataset, k):
123     k_points = generate_k(dataset, k)
124     assignments = assign_points(dataset, k_points)
125     old_assignments = None
126     while assignments != old_assignments:
127         new_centers = update_centers(dataset, assignments)
128         old_assignments = assignments
129         assignments = assign_points(dataset, new_centers)
130
131     result = dict()
132     for assign_num, point in zip(assignments, dataset):
133         if assign_num not in result:
134             result[assign_num] = [[point[i]] for i in range(len(point))]
135         else:
136             for i in range(len(point)):

```

```

137         result[assign_num][i].append(point[i])
138     return result, new_centers
139
140
141 def visualize_clusters(clusters, information):
142     """
143     Visualizes the first 2 dimensions of the data as a 2-D scatter plot.
144     """
145     plt.figure()
146     colors = ["#%06x" % random.randint(0, 0xFFFFFF) for _ in range(len(clusters))]
147     for num_cluster in clusters:
148         plt.plot(clusters[num_cluster][0], clusters[num_cluster][1], 'o',
149                 ↪ color=colors[num_cluster], markersize=2)
150
151     for num_cluster in clusters:
152         for i in range(0, len(clusters[num_cluster][0]), 111):
153
154             point = (clusters[num_cluster][0][i], clusters[num_cluster][1][i])
155
156             for j in range(len(information)):
157
158                 if point[0] == information[j][0] and point[1] == information[j][1]:
159                     dx = random.randint(-1, 1)
160                     dy = random.randint(-1, 1)
161
162                     plt.annotate(information[j][2], xy=(clusters[num_cluster][0][i],
163                     ↪ clusters[num_cluster][1][i]),
164                                 xytext=(clusters[num_cluster][0][i] + (dx * 40),
165                                 ↪ clusters[num_cluster][1][i] + (dy * 6)))
166
167                     break
168
169     plt.show()
170
171 fixed_df = pd.read_csv('aids.csv', sep=',', encoding='utf8')
172
173
174 cases = [int(case.replace(",", "")) for case in fixed_df['Cases'].to_list()]
175 rates = [rate for rate in fixed_df['Rate'].to_list()]
176 population = [p.replace(",", "").replace('\\', "'") for p in
177               ↪ fixed_df['Population'].to_list()]

```



```

179 clusters, centers = k_means(list(zip(cases, rates)), 2)
180 visualize_clusters(clusters, list(zip(cases, rates, population)))

```

4 Результаты

Примеры результаты работы программы приведены на рисунках для разного количества кластеров (рис. 1-3).

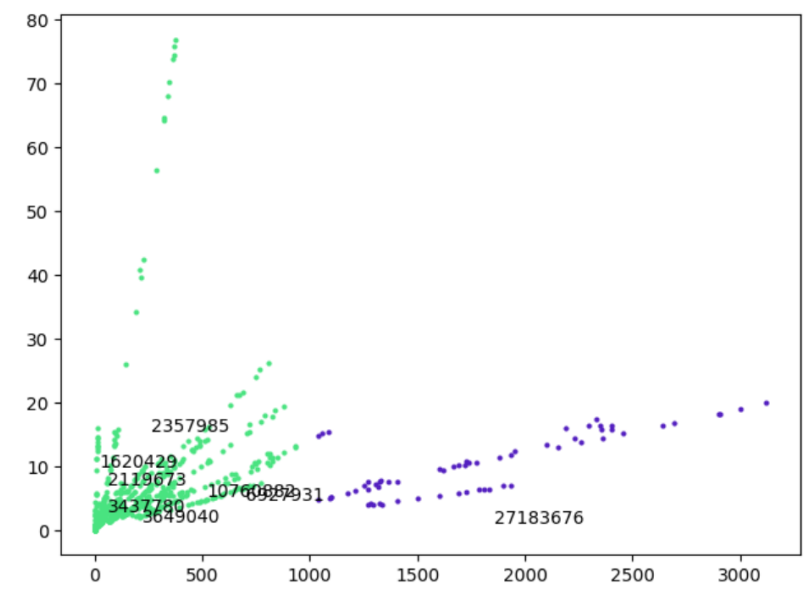


Рис. 1: 2 кластера

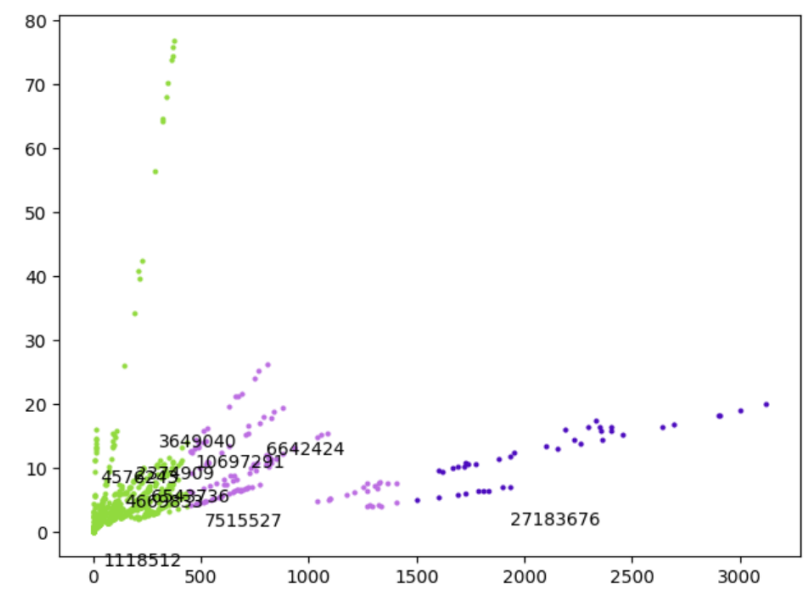


Рис. 2: 3 кластера

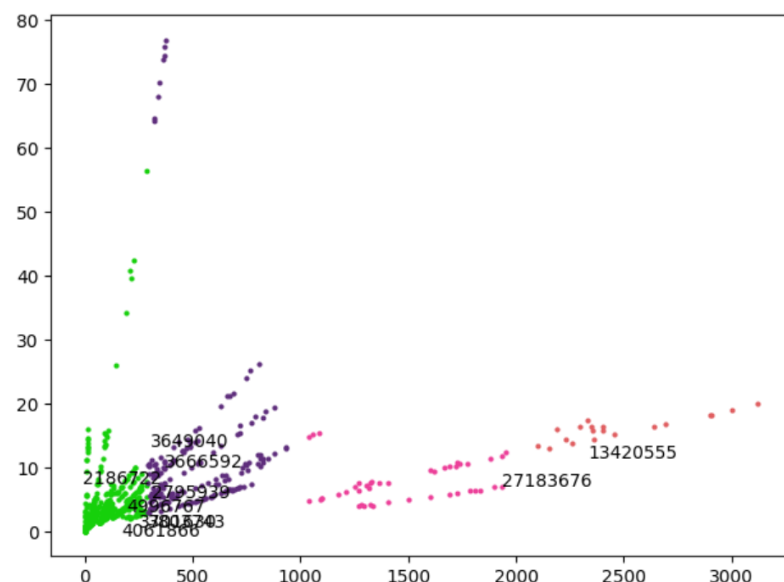


Рис. 3: 4 кластера

5 Вывод

На рисунке 1 более четко можно сформулировать гипотезу о то, что в мире есть тенденция увеличения числа ВИЧ-зараженных, при этом в развивающихся странах число более резко возрастает, нежели чем в развитых странах.