

```

1 #####
2 # TITLE: ex3 machine learning BIU
3 # WRITER: Raphael Haehnel
4 # DATE: 8/12/2022
5 #####
6
7 import os
8 import matplotlib.pyplot as plt
9 import numpy as np
10 import functools
11 import scipy.io
12 import shutil
13 from tqdm import tqdm
14
15 LIMIT_ITER = 100
16
17 def generate_data(n, c, dim):
18     means = 12*(np.random.rand(c, dim) - 0.5)
19     means = means.repeat(n, axis=0)
20
21     data = means + np.random.randn(n*c, dim)
22     col = np.zeros((n*c, 1))
23     data = np.append(data, col, axis=1)
24     return data
25
26
27 def generate_means(k, dim):
28     means = 12*(np.random.rand(k, dim) - 0.5)
29     label = np.array([np.array(range(k))]).T
30     return np.append(means, label, axis=1)
31
32
33 def generate_means_pixels(k, dim):
34     means = np.random.rand(k, dim)
35     label = np.array([np.array(range(k))]).T
36     return np.append(means, label, axis=1)
37
38
39 def show_plot(data, mu, i, save, dim):
40     if dim == 2:
41         plt.scatter(x=data[:, 0], y=data[:, 1], c=data[:, 2], s=20)
42         plt.scatter(x=mu[:, 0], y=mu[:, 1], c=mu[:, 2],
43                     marker="*", edgecolors="black", s=100)
44         plt.xlim([-10, 10])
45         plt.ylim([-10, 10])
46     if dim == 3:
47         fig = plt.figure()
48         ax = fig.add_subplot(projection='3d')
49
50         ax.scatter(data[:, 0], data[:, 1], data[:, 2], c=data[:, 3], s=20)
51         ax.scatter(mu[:, 0], mu[:, 1], mu[:, 2], c=mu[:, 3],
52                     marker="*", edgecolors="black", s=100)
53
54     if save:
55         if not os.path.exists('./output'):
56             os.mkdir('./output')
57         plt.savefig(f'./output/{i}')
58         plt.close()
59     else:
60         plt.show()
61
62
63 def duplicate_mu(data, mu, dim):
64     n_data = data.shape[0]
65     n_mu = mu.shape[0]
66
67     # Duplicate the values of mu along the vertical axis
68     mu_v = np.apply_along_axis(functools.partial(
69         np.repeat, repeats=n_data, axis=0), axis=0, arr=mu[:, :dim])
70
71     # Split the array into multiple sub-arrays vertically
72     return np.vsplit(mu_v, n_mu)
73
74
75 def assignment(data, mu, dim):
76     mu_duplicated = duplicate_mu(data, mu, dim)

```

```

79 x = data[:, :dim]
80 result = np.zeros((1, data.shape[0]))
81
82 for i in tqdm(range(mu.shape[0])):
83     diff = x-mu_duplicated[i]
84     result = np.vstack([result, np.apply_along_axis(
85         lambda a: np.sum(a**2), 1, diff)])
86 result = result[1:]
87 arg_min = np.argmin(result, 0)
88
89 data[:, dim] = arg_min
90 cost = np.sum(np.min(result, 0))
91
92 return data, cost
93
94
95 def helper_sum(data, k, dim):
96     result = np.zeros((k, dim))
97     result[int(data[dim])] = data[:dim]
98
99     return result
100
101
102 def cout_data_labels(data, k, dim):
103
104     count = np.zeros((1, k))
105     for row in data:
106         count[0, int(row[dim])] += 1
107
108     return count
109
110
111 def centroid_update(data, k, dim):
112     cout_labels = cout_data_labels(data, k, dim)
113     new_mu = np.apply_along_axis(
114         functools.partial(helper_sum, k=k, dim=dim), 1, data)
115
116     new_mu = np.sum(new_mu, axis=0)
117     new_mu = np.divide(new_mu, cout_labels.repeat(dim, 0).T, out=np.zeros_like(
118         new_mu), where=cout_labels.repeat(dim, 0).T != 0)
119     label = np.array([np.array(range(k))]).T
120     new_mu = np.append(new_mu, label, axis=1)
121
122     return new_mu
123
124
125 def run_clustering(data, k, mu, dim, save):
126
127     total_cost = np.array([])
128
129     for i in range(LIMIT_ITER):
130         data, cost = assignment(data, mu, dim)
131         total_cost = np.append(total_cost, cost)
132         new_mu = centroid_update(data, k, dim)
133         if np.array_equal(new_mu, mu):
134             break
135         mu = new_mu
136         show_plot(data, mu, i, save=save, dim=dim)
137
138     return data, total_cost
139
140
141 def success_rate(y_train, y_output):
142     return np.count_nonzero(y_output == y_train)/len(y_output)*100
143
144
145 def normalize(array):
146     return array / array.max()
147
148
149 def retriive_data(max_n):
150     mat = scipy.io.loadmat('mnist_all.mat')
151
152     x0 = mat.get('train0')[:max_n, :]
153     x1 = mat.get('train1')[:max_n, :]
154     x2 = mat.get('train2')[:max_n, :]
155     x3 = mat.get('train3')[:max_n, :]
156     x4 = mat.get('train4')[:max_n, :]
157     x5 = mat.get('train5')[:max_n, :]

```

```

158 x6 = mat.get('train6')[:max_n, :]
159 x7 = mat.get('train7')[:max_n, :]
160 x8 = mat.get('train8')[:max_n, :]
161 x9 = mat.get('train9')[:max_n, :]
162 X_train = np.concatenate((x0, x1, x2, x3, x4, x5, x6, x7, x8, x9), axis=0)
163
164 y0 = 0*np.ones(len(x0), dtype=int)
165 y1 = 1*np.ones(len(x1), dtype=int)
166 y2 = 2*np.ones(len(x2), dtype=int)
167 y3 = 3*np.ones(len(x3), dtype=int)
168 y4 = 4*np.ones(len(x4), dtype=int)
169 y5 = 5*np.ones(len(x5), dtype=int)
170 y6 = 6*np.ones(len(x6), dtype=int)
171 y7 = 7*np.ones(len(x7), dtype=int)
172 y8 = 8*np.ones(len(x8), dtype=int)
173 y9 = 9*np.ones(len(x9), dtype=int)
174 y_train = np.concatenate((y0, y1, y2, y3, y4, y5, y6, y7, y8, y9), axis=0)
175
176 x0_test = mat.get('test0')[:max_n, :]
177 x1_test = mat.get('test1')[:max_n, :]
178 x2_test = mat.get('test2')[:max_n, :]
179 x3_test = mat.get('test3')[:max_n, :]
180 x4_test = mat.get('test4')[:max_n, :]
181 x5_test = mat.get('test5')[:max_n, :]
182 x6_test = mat.get('test6')[:max_n, :]
183 x7_test = mat.get('test7')[:max_n, :]
184 x8_test = mat.get('test8')[:max_n, :]
185 x9_test = mat.get('test9')[:max_n, :]
186 X_test = np.concatenate((x0_test, x1_test, x2_test, x3_test,
187                          x4_test, x5_test, x6_test, x7_test, x8_test, x9_test),
axis=0)
188
189 y0_test = 0*np.ones(len(x0_test), dtype=int)
190 y1_test = 1*np.ones(len(x1_test), dtype=int)
191 y2_test = 2*np.ones(len(x2_test), dtype=int)
192 y3_test = 3*np.ones(len(x3_test), dtype=int)
193 y4_test = 4*np.ones(len(x4_test), dtype=int)
194 y5_test = 5*np.ones(len(x5_test), dtype=int)
195 y6_test = 6*np.ones(len(x6_test), dtype=int)
196 y7_test = 7*np.ones(len(x7_test), dtype=int)
197 y8_test = 8*np.ones(len(x8_test), dtype=int)
198 y9_test = 9*np.ones(len(x9_test), dtype=int)
199 y_test = np.concatenate((y0_test, y1_test, y2_test, y3_test,
200                          y4_test, y5_test, y6_test, y7_test, y8_test, y9_test),
axis=0)
201
202 return normalize(X_train), normalize(y_train), normalize(X_test),
normalize(y_test)
203
204
205 if __name__ == "__main__":
206
207     k = 10
208     c = k
209
210     if os.path.exists('./output'):
211         shutil.rmtree('./output')
212
213     X_train, y_train, X_test, y_test = retrieve_data(max_n=100)
214
215     #####
216     # To use the MNIST data, use this code:
217     n = X_train.shape[0]
218     dim = X_train.shape[1]
219     data = np.append(X_train, np.array([y_train]).T, axis=1)
220     mu = generate_means_pixels(k, dim)
221     # To generate your own data, use this code:
222     # n = 100
223     # dim = 3
224     # data = generate_data(n, c, dim)
225     # mu = generate_means(k, dim)
226     #####
227
228     output, total_cost = run_clustering(data, k, mu, dim, save=True)
229     result_success = success_rate(y_train, output[:, -1])
230     print(f'Success rate = {result_success}%')
231     plt.plot(total_cost, marker="o")
232     plt.show()
233

```