```python
1  import os
2  import shutil
3  import matplotlib.pyplot as plt
4  import numpy as np
5  import functools
6  import scipy.stats as sc
7  import sklearn.metrics as skl
8
9  MIN_X = -10
10 MAX_X = 10
11 MIN_Y = MIN_X
12 MAX_Y = MAX_X
13 LIMIT_ITER = 100
14
15
16 def generate_specific_means():
17     n = 100
18     c = 2
19
20     means = np.array([[-3, 1.5], [-3, -1.5]])
21     sigma = np.array([[6, 0], [0, 0.8]])
22     x = np.random.rand(100, 2)
23
24     x = np.matmul(x, sigma)
25     x1 = x + means[0]
26     x2 = x + means[1]
27
28     data = np.vstack([x1, x2])
29     data_real = data
30
31     label_default = np.zeros((n*c, 1))
32     data = np.append(data, label_default, axis=1)
33
34     label_real = np.repeat(np.array(range(c)), repeats=n).reshape(-1, 1)
35     data_real = np.append(data_real, label_real, axis=1)
36
37     return data, data_real
38
39
40 def generate_data(n: int, c: int):
41     """
42     Generate two matrices of random samples. The dimensions of the matrices are (nc,
   d+1)
43     :param n: The number of samples
44     :param m: The number of generated cluster
45     :return: data: data with empty label
46             data_real: data with correct label
47     """
48     means = 12*(np.random.rand(c, 2) - 0.5)
49     means = means.repeat(n, axis=0)
50
51     data = means + np.random.randn(n*c, 2)
52     data_real = data
53
54     label_default = np.zeros((n*c, 1))
55     data = np.append(data, label_default, axis=1)
56
57     label_real = np.repeat(np.array(range(c)), repeats=n).reshape(-1, 1)
58     data_real = np.append(data_real, label_real, axis=1)
59
60     return data, data_real
61
62
63 def generate_means(k):
64     """
65     Generate a matrix of labeled random means. The dimension of the matrix is (k,
   d+1)
66     :param k: The number of means
67     :return: A matrix of labeled random means
68     """
69     means = 12*(np.random.rand(k, 2) - 0.5)
70     label = np.array([np.array(range(k))]).T
71     return np.append(means, label, axis=1)
72
73
74 def generate_sigmas(k):
75     """
76     Generate k identical matrices of sigma
```

```python
77          :param k: The number of gaussians
78          :return: A list of matrices
79          """
80          return np.repeat(np.eye(2)[None, ...], 2, axis=0)
81
82
83   def generate_alphas(k):
84          return np.ones(k)
85
86
87   def show_plot(data, mu=None, i=0, gaussian_list=None, alpha_list=None, save=False,
     show=True):
88
89          plt.scatter(x=data[:, 0], y=data[:, 1], c=data[:, 2], s=5)
90          if mu is not None:
91              plt.scatter(x=mu[:, 0], y=mu[:, 1], c=mu[:, 2],
92                          marker="*", edgecolors="black", s=100)
93          if gaussian_list is not None:
94              N = 200
95              X = np.linspace(MIN_X, MAX_X, N)
96              Y = np.linspace(MIN_Y, MAX_Y, N)
97              Z = draw_gaussian(X, Y, gaussian_list, alpha_list)
98              plt.contour(X, Y, Z)
99
100         plt.xlim([MIN_X, MAX_X])
101         plt.ylim([MIN_Y, MAX_Y])
102
103         if save:
104             if not os.path.exists('./output'):
105                 os.mkdir('./output')
106             plt.savefig(f'./output/{i}')
107             plt.close()
108         else:
109             if show == True:
110                 plt.show()
111
112         if gaussian_list is not None:
113             return Z
114
115
116  def duplicate_mu(data, mu):
117         n_data = data.shape[0]
118         n_mu = mu.shape[0]
119
120         # Duplicate the values of mu along the vertical axis
121         mu_v = np.apply_along_axis(functools.partial(
122             np.repeat, repeats=n_data, axis=0), axis=0, arr=mu[:, :2])
123
124         # Split the array into multiple sub-arrays vertically
125         return np.vsplit(mu_v, n_mu)
126
127
128  def assignment(data, mu):
129
130         mu_duplicated = duplicate_mu(data, mu)
131         x = data[:, :2]
132         result = np.zeros((1, data.shape[0]))
133
134         for i in range(mu.shape[0]):
135             diff = x-mu_duplicated[i]
136             result = np.vstack([result, np.apply_along_axis(
137                 lambda a: a[0]**2+a[1]**2, 1, diff)])
138         result = result[1:]
139         arg_min = np.argmin(result, 0)
140
141         data[:, 2] = arg_min
142         cost = np.sum(np.min(result, 0))
143
144         return data, cost
145
146
147  def helper_sum(data, k):
148         result = np.zeros((k, 2))
149         result[int(data[2])] = [data[0], data[1]]
150         return result
151
152
153  def get_linked_data(data, k):
154
```

```python
155        sum_data_linked = np.zeros((1, k))
156        for row in data:
157            sum_data_linked[0, int(row[2])] += 1
158
159        return sum_data_linked
160
161
162 def centroid_update(data, k):
163        linked_data = get_linked_data(data, k)
164        new_mu = np.apply_along_axis(functools.partial(helper_sum, k=k), 1, data)
165
166        new_mu = np.sum(new_mu, axis=0)
167        new_mu = np.divide(new_mu, linked_data.repeat(2, 0).T, out=np.zeros_like(
168            new_mu), where=linked_data.repeat(2, 0).T != 0)
169        label = np.array([np.array(range(k))]).T
170        new_mu = np.append(new_mu, label, axis=1)
171
172        return new_mu
173
174
175 def cluster(data, mu, save):
176
177        total_cost = np.array([])
178
179        for i in range(16):
180
181            data, cost = assignment(data, mu)
182            total_cost = np.append(total_cost, cost)
183
184            new_mu = centroid_update(data, k)
185            if np.array_equal(new_mu, mu):
186                break
187            mu = new_mu
188            i += 1
189            show_plot(data, mu, i, save=save, show=False)
190
191        return data, mu, total_cost
192
193
194 def draw_gaussian(X, Y, gaussian_list:
    list[sc._multivariate.multivariate_normal_frozen], alpha_list):
195
196        X, Y = np.meshgrid(X, Y)
197        pos = np.dstack((X, Y))
198        Z_tot = np.zeros(X.shape)
199
200        for i in range(len(gaussian_list)):
201            Z_tot += alpha_list[i]*gaussian_list[i].pdf(pos)
202        return Z_tot
203
204
205 def generate_list_gaussians(means, sigma):
206
207        gaussian_list = []
208        sigma = np.array([[1, 0], [0, 1]])
209
210        for mu in means:
211            gaussian_list.append(sc.multivariate_normal(mu, sigma))
212
213        return gaussian_list
214
215
216 def expectation_step(gaussian: list[sc._multivariate.multivariate_normal_frozen],
    data, k, alpha_list):
217        x = data[:, :2]
218
219        # Sum for the denominator
220        divisor = 0
221        for l in range(k):
222            divisor += gaussian[l].pdf(x)*alpha_list[l]
223
224        w_list = []
225        for j in range(k):
226            w = gaussian[j].pdf(x)*alpha_list[j] / divisor
227            w = w.reshape((w.shape[0], 1))
228            w_list.append(w)
229
230        return w_list
231
```

```python
232
233 def minimalization_step(w_list, data):
234     x = data[:, :2]
235     gaussian_list = []
236     alpha_list = []
237
238     for w in w_list:
239         n = w.shape[0]
240
241         alpha = 1/n * np.sum(w)
242         mu = sum(w*x)/np.sum(w)
243         sigma = np.matmul((x-mu).T, w*(x-mu))/np.sum(w)
244
245         gaussian_list.append(sc.multivariate_normal(
246             mu, sigma, allow_singular=True))
247         # TODO numpy.linalg.LinAlgError: When `allow_singular is False`, the input
    matrix must be symmetric positive definite
248
249         alpha_list.append(alpha)
250
251     return gaussian_list, alpha_list
252
253
254 def gmm(gaussian_list, data, k, alpha_list):
255
256     # Show the data with the initialized gaussians
257     Z = show_plot(data, mu, -1, save=False,
258                   gaussian_list=gaussian_list, alpha_list=alpha_list)
259
260     # Run the Gaussian Mixture Method
261     for i in range(LIMIT_ITER):
262         w_list = expectation_step(gaussian_list, data, k, alpha_list)
263         gaussian_list, alpha_list = minimalization_step(w_list, data)
264
265         new_Z = show_plot(data, mu, i, save=True,
266                           gaussian_list=gaussian_list, alpha_list=alpha_list)
267         print(f'Iteration {i}', end="\r")
268         if np.allclose(Z, new_Z, rtol=1e-4, atol=1e-4):
269             break
270
271         Z = new_Z
272
273     print("")
274     w_as_array = np.array(w_list).squeeze()
275     return w_as_array
276
277
278 def get_label(w_array):
279     labels = np.argmax(w_array, axis=0)
280     return labels.reshape(labels.shape[0])
281
282
283 if __name__ == "__main__":
284
285     k = 2  # The number of clusters we'd like to find
286     n = 100  # The size of the data
287     c = k  # The number of generated cluster (we should define c=k)
288
289     if os.path.exists('./output'):
290         shutil.rmtree('./output')
291
292     data, data_real = generate_data(n, c)
293
294     # If you want to generate specific data, uncomment this
295     # data, data_real = generate_specific_means()
296
297     # Show only the data
298     show_plot(data_real, i=-2, save=False)
299
300     # Generate the data for the gmm algorithm
301     mu = generate_means(k)
302     sigma = generate_sigmas(k)
303     alpha_list = generate_alphas(k)
304     gaussian_list = generate_list_gaussians(mu[:, :2], sigma)
305
306     # Run the gmm algorithm
307     w_array = gmm(gaussian_list, data, k, alpha_list)
308     labels = get_label(w_array)
309
```

```python
    data_parsed = data
    data_parsed[:, 2] = labels
    show_plot(data_parsed, i=-2, save=False)

    score_gmm = skl.adjusted_rand_score(data_real[:, 2], labels)

    data, mu, total_cost = cluster(data, mu, save=False)

    score_kmeans = skl.adjusted_rand_score(data_real[:, 2], data[:, 2])

    print("gmm score: ", score_gmm)
    print("kmeans score: ", score_kmeans)

    print("eof")
```