

Clustering

I. Introduction

Clustering analysis is a unsupervised learning technique that is applied on an unlabelled database to identify sets of objects that belong to groups of similar properties or interests. It is an application used in many fields including image analysis, bioinformatics, social networks, social science, and so on...

The objective in this project is to apply two types of clustering analysis:

- a hard-clustering analysis, the K-means model.
- a type of soft-clustering analysis, the Gaussian mixture model.

The difference between hard and soft clustering is that, in one case, every single data point is assigned to only one cluster and boundaries are clearly visible, while in the other case, a weight is applied to each data point telling how confident we are that it belongs to a given cluster, and as such boundaries are more fuzzy.

The most popular database for education purposes, the Iris dataset [1], will be used for the analysis. It contains 150 datapoints each having 4 different attributes and belonging to one of three different classes. The goal will be to identify the class of each datapoint using a clustering model.

[1] <https://archive.ics.uci.edu/ml/index.php>

II. Methodology

1. K-means model

As mentionned above, K-means will seek to assign each datapoint to one of the K clusters. In the case of unlabelled datasets, the value of K is arbitrary, and several runs of the model with various values of K might be needed before reaching an appropriate partitionning of the data.

For each datapoint x_i , we associate a cluster assignment $c_i \in \{1, \dots, K\}$.

And each cluster in $\{1, \dots, K\}$ is represented by a mean vector μ_k , called a centroid, that defines the center of a cluster.

We are going to find the centroids such that all the datapoints assigned to a cluster are the closest to that same cluster centroid (we will use the squared Euclidian distance of x_i to μ_i). We want to minimize the following objective function:

$$\mu^*, c^* = \underset{\mu, c}{\operatorname{argmin}} \sum_{i=1}^n \sum_{k=1}^K \mathbb{1}\{c_i = k\} \|x_i - \mu_k\|^2$$

To solve this, we are going to split the model variables into two different sets. One set is going to contain all the vectors μ_k , and the other set will contain all the cluster assignments c_i .

We are first going to randomly initialize our centroids. Then we are going to iterate, updating each c_i and μ_k back and forth between the two following steps until the objective function stops changing:

- *Step 1:* Update each c_i by setting c_i to be the index of the cluster that the point x_i is closest to, in the Euclidian sense:

$$c_i = \underset{k}{\operatorname{argmin}} \|x_i - \mu_k\|^2$$

- *Step 2:* Update each μ_k by taking the average of all of the data assigned to a particular centroid. So for the k^{th} centroid we simply average all of the data that was assigned to the k centroid according to the most recent update of each c_i :

$$\mu_k = \frac{1}{N_k} \sum_{i=1}^n x_i \mathbb{1}\{c_i = k\} \quad \text{with} \quad N_k = \sum_{i=1}^n \mathbb{1}\{c_i = k\}$$

2. Gaussian mixture model

The Gaussian mixture model is a probabilistic soft-clustering model where the weight vector ϕ_i is a probability distribution of data point x_i being assigned to cluster c_i .

GMM works by first defining:

- A prior distribution π_k on the k^{th} cluster.
- A likelihood distribution on the i^{th} observation x_i given it comes from the k^{th} cluster, that we write as a Gaussian distribution of mean μ_k and covariance Σ_k : $N(x_i | \mu_k, \Sigma_k)$.

The goal of the GMM is to maximize the log of the marginal likelihood of each datapoint x_i over π , μ and Σ in order to find those 3 parameters:

$$\mathcal{L} = \sum_{i=1}^n \ln p(x_i | \pi, \mu, \Sigma) = \sum_{i=1}^n \ln \sum_{k=1}^K \pi_k N(x_i | \mu_k, \Sigma_k)$$

We first initialize the parameters such as:

- π is initialized as a uniform distribution in the interval $[0, 1]$.
- Each μ_k is randomly initialized on one of the x_i datapoint location.
- Each Σ_k is initialized as an identity matrix.

Then the iteration is done in two steps until convergence of the objective function \mathcal{L} :

- Step 1: For each data point, the posterior probability $\phi_i(k)$ of coming from each of the k different clusters is calculated:

$$\phi_i(k) = \frac{\pi_k N(x_i | \mu_k, \Sigma_k)}{\sum_{j=1}^K \pi_j N(x_i | \mu_j, \Sigma_j)}, \quad \text{for } k = 1, \dots, K$$

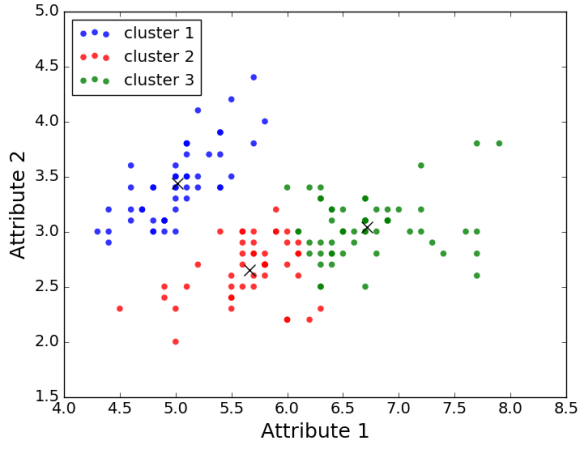
- Step 2: By defining $n_k = \sum_{i=1}^n \phi_i(k)$, we update the model parameters as follow:

$$\pi_k = \frac{n_k}{n}, \quad \mu_k = \frac{1}{n_k} \sum_{i=1}^n \phi_i(k) x_i, \quad \Sigma_k = \frac{1}{n_k} \sum_{i=1}^n (x_i - \mu_k)(x_i - \mu_k)^T$$

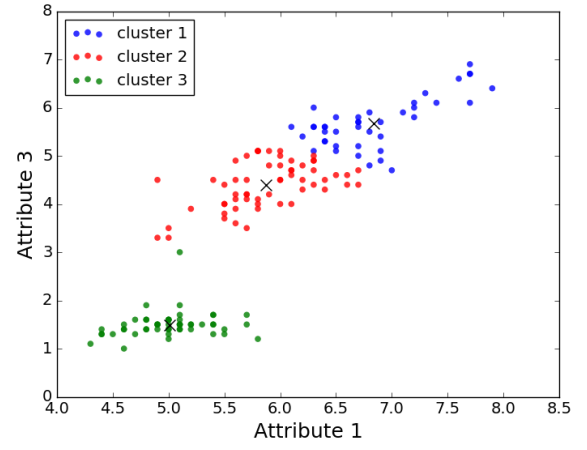
III. Results

1. K-means model

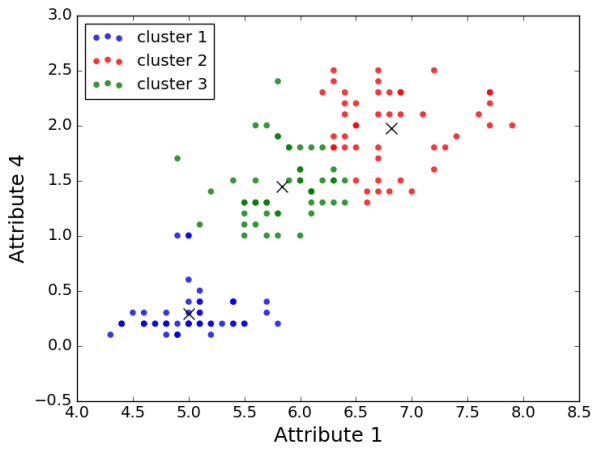
In order to visually see the partitionning and centroids location, only two of the four original attributes were used (see Figure 1). For each couple of attributes, comparisons were made between predictions and real class types.



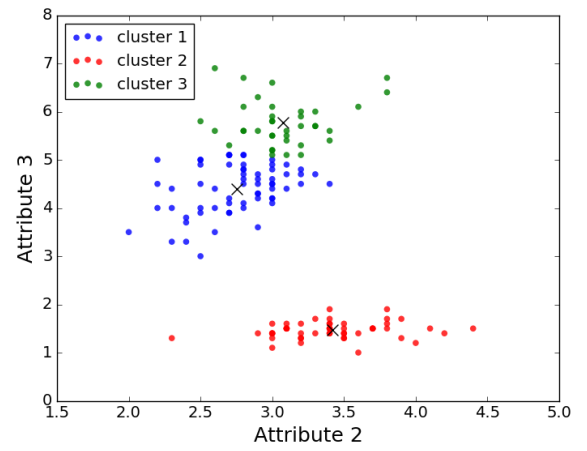
(a) Attributes 1 and 2



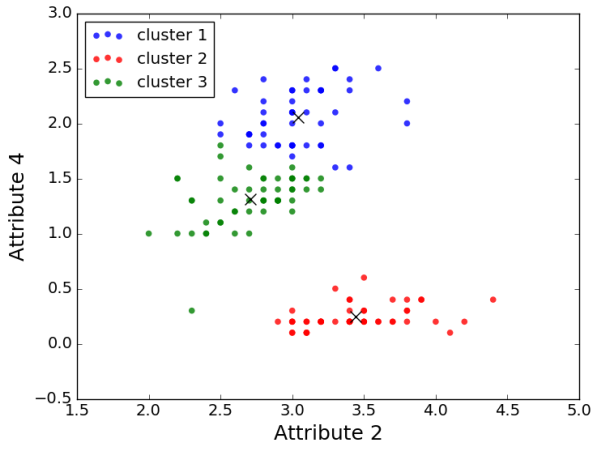
(b) Attributes 1 and 3



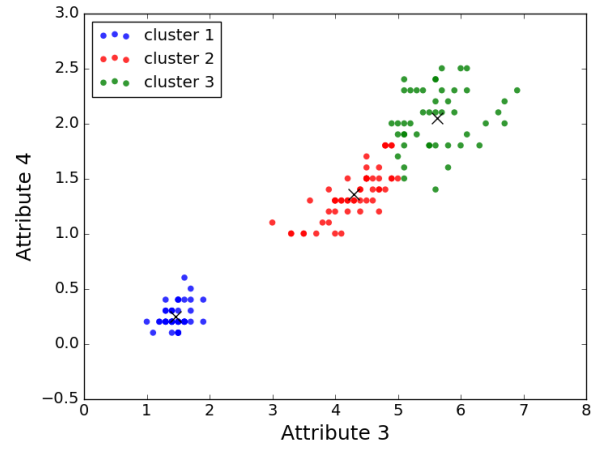
(c) Attributes 1 and 4



(d) Attributes 2 and 3



(e) Attributes 2 and 4



(f) Attributes 3 and 4

Figure 1: Various clusters plots with centroids for 2 of the 4 attributes

Attributes	Cluster 1	Cluster 2	Cluster 3
1 & 2	100%	76%	68%
1 & 3	100%	94%	74%
1 & 4	100%	82%	74%
2 & 3	100%	100%	74%
2 & 4	98%	92%	88%
3 & 4	100%	96%	88%
1 to 4	100%	96%	72%

Table 1: Rate of accurate predictions obtained with K-means

2. Gaussian mixture model

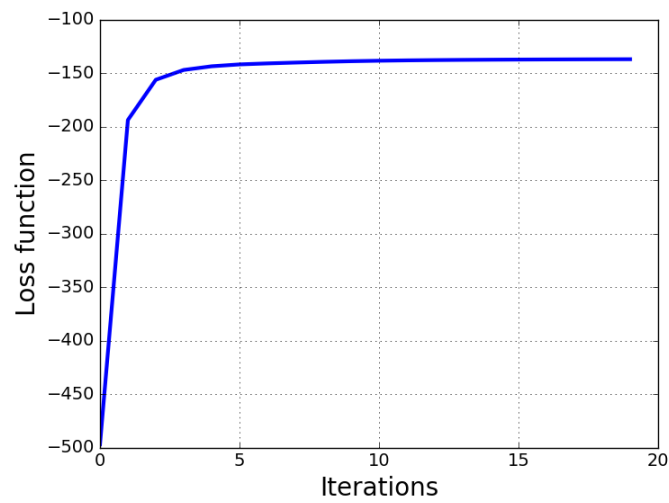


Figure 2: Objective function \mathcal{L} of the GMM over iteration time

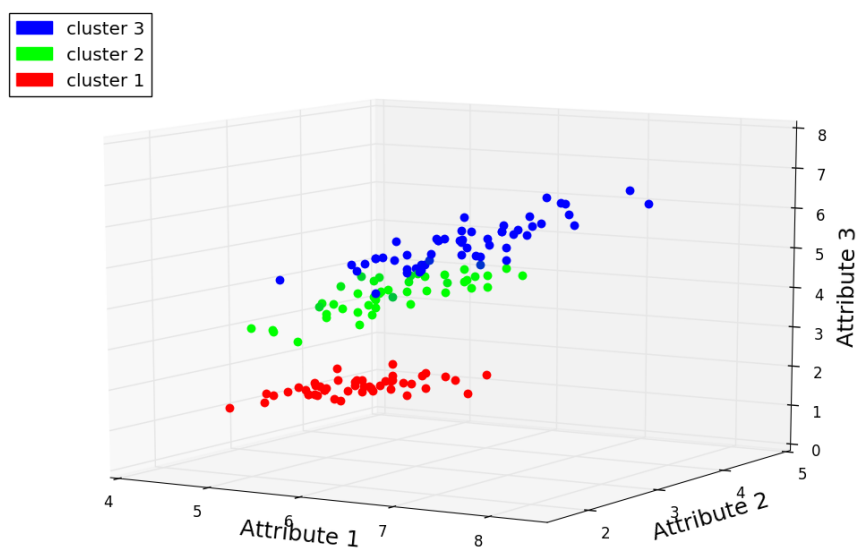


Figure 3: 3D plot of the clusters using all attributes with the GMM

IV. Python Code

```
1 from __future__ import division
2 import numpy as np
3 import sys
4 import random as rd
5 from matplotlib import pyplot as plt
6 import matplotlib.patches as mpatches
7 from mpl_toolkits.mplot3d import Axes3D
8 import matplotlib.ticker as ticker
9 from scipy.stats import multivariate_normal as MvN
10
11 """
12 This script should be launched from the console using the following command line:
13 python script_file X_data Label_data nb_cluster Clustering_type nb_iter
14 ex:    > python Clustering.py X.csv Y.csv 3 Kmeans 50
15 or:    > python Clustering.py X.csv Y.csv 3 GMM 50
16 Where:
17 .script_file is the name of this python file
18 .X_data is a .csv file of the data nodes coordinates
19 there should be one column per dimension, so for instance 2 columns if we are in 2D
20 .Label_data is a .csv file that should contain the known labels corresponding
21 to the nodes in the X_data file
22 .nb_cluster is the number of clusters we want to partition the data into (5 is the default)
23 .Clustering_type is the type of clustering we want to apply:
24 Kmeans if we want to apply hard-clustering
25 GMM if we want to apply soft-clustering
26 .nb_iter is the number of iterations we are going to use
27 """
28
29
30 class Kmeans:
31     def __init__(self, nb_clust, nb_iter):
32         self.X = np.genfromtxt(sys.argv[1], delimiter=",")
33         self.Y = np.genfromtxt(sys.argv[2], delimiter=",")
34         self.nb_clust = nb_clust
35         self.nb_iter = nb_iter
36
37         # Random initialisation of the mean  $\mu = (\mu_1, \dots, \mu_K)$ 
38         self.mu = np.zeros((nb_clust, len(self.X[0])))
39         # print(self.mu.shape)
40         for k in range(self.nb_clust):
41             self.mu[k] = rd.choice(self.X)
42
43         self.ci = np.zeros((len(self.X)))
44         self.nk = np.zeros((nb_clust))
45
46     def KMeans(self):
47         LX = len(self.X) # nb of datapoints
48
```

```

49     ### do the following steps for self.nb_iter iterationsb ###
50
51     for iter in range(self.nb_iter):
52         centerslist = []
53         print('iteration %s' % iter)
54
55         # calculate for each cluster the measured Euclidian distance between
56         # and select for c_i the cluster with the minimal distance
57         count = 0
58         for xi in self.X:
59             val1 = 1e10
60             clust = -1
61             for k in range(self.nb_clust):
62                 mu = self.mu[k]
63                 val2 = np.dot(xi - mu, xi - mu)
64                 # if count == 133:
65                 # print('xi: %s, mu: %s, xi-mu: %s' % (xi, mu, xi - mu))
66
67                 if val2 < val1:
68                     val1 = val2
69                     clust = k
70                     # if count == 133:
71                     # print('count: %s, k:%s, np.dot:%s' % (count, k, val2))
72             self.ci[count] = clust
73             count += 1
74
75         # calculate numbers of indicators c_i for each cluster k
76         for k in range(self.nb_clust):
77             nbk = 0
78             for i in range(LX):
79                 if self.ci[i] == k:
80                     nbk += 1
81             self.nk[k] = nbk
82             if self.nk[k] == 0:
83                 self.nk[k] = 1
84         # print('nk: %s'%self.nk)
85
86         # Update the mean mu = (mu_1,...,m_K)
87         for k in range(self.nb_clust):
88             val = np.zeros((len(self.X[0])))
89             count = 0
90             for i in range(LX):
91                 if self.ci[i] == k:
92                     val += self.X[i]
93             self.mu[k] = val / self.nk[k]
94             centerslist.append(self.mu[k])
95
96         # save centroids into file for each iteration
97         filename = "centroids-" + str(iter + 1) + ".csv" # "i" would be each
98         np.savetxt(filename, centerslist, delimiter=",")
99         # print(Kmeans.check_prediction(self, self.ci))
100
101     return self.mu, self.ci

```

```

102
103     def check_prediction(self, ci):
104         LX = len(self.X)  # nb of datapoints
105         cTT = 50
106         ct1 = 0
107         ct2 = 0
108         ct3 = 0
109         print(ci)
110         for i in range(LX):
111             if 0 <= i < cTT:
112                 if i == 0:
113                     val = ci[i]
114                     if ci[i] == val:
115                         ct1 += 1
116             elif cTT <= i < 2 * cTT:
117                 if i == cTT:
118                     val = ci[i]
119                     if ci[i] == val:
120                         ct2 += 1
121             elif 2 * cTT <= i < 3 * cTT:
122                 if i == 2 * cTT:
123                     val = ci[i]
124                     if ci[i] == val:
125                         ct3 += 1
126
127         return ct1, ct2, ct3
128
129     def plot_data(self, ci, mu, nb):
130         x1 = []
131         y1 = []
132         x2 = []
133         y2 = []
134         x3 = []
135         y3 = []
136         x4 = []
137         y4 = []
138         x5 = []
139         y5 = []
140         mux = []
141         muy = []
142
143         # get mean
144         for k in range(len(mu)):
145             mux.append(mu[k][0])
146             muy.append(mu[k][1])
147
148         ## plot datapoints per cluster
149         colors = ('b', 'r', 'g', 'c', 'y')
150         groups = ("cluster 1", "cluster 2", "cluster 3", "cluster 4", "cluster 5")
151         for i in range(len(self.X)):
152             if ci[i] == 0:
153                 x1.append(self.X[i][0])
154                 y1.append(self.X[i][1])

```

```

155         elif ci[i] == 1:
156             x2.append(self.X[i][0])
157             y2.append(self.X[i][1])
158         elif ci[i] == 2:
159             x3.append(self.X[i][0])
160             y3.append(self.X[i][1])
161         elif ci[i] == 3:
162             x4.append(self.X[i][0])
163             y4.append(self.X[i][1])
164         elif ci[i] == 4:
165             x5.append(self.X[i][0])
166             y5.append(self.X[i][1])
167
168     g1 = (x1, y1)
169     g2 = (x2, y2)
170     g3 = (x3, y3)
171     g4 = (x4, y4)
172     g5 = (x5, y5)
173     data = (g1, g2, g3, g4, g5)
174
175     # Create plot
176     fig = plt.figure()
177     #ax = fig.add_subplot(1, 1, 1, axisbg="1.0") #python 2.7
178     ax = fig.add_subplot(1, 1, 1)
179     for data, color, group in zip(data, colors[0:nb], groups[0:nb]):
180         x, y = data
181         ax.scatter(x, y, alpha=0.8, c=color, edgecolors='none', s=30, label=group)
182     plt.scatter(mux, muy, c='k', marker='x', s=100)
183     plt.legend(loc=2)
184     plt.show()
185
186
187 class GMM:
188     def __init__(self, nb_clust, nb_iter):
189         self.X = np.genfromtxt(sys.argv[1], delimiter=",")
190         self.nb_clust = nb_clust
191         self.nb_iter = nb_iter
192         self.nk = np.zeros((nb_clust))
193
194         # Random initialisation of the mean mu = (mu_1,...,m_K)
195         self.mu = np.zeros((nb_clust, len(self.X[0])))
196         for k in range(self.nb_clust):
197             self.mu[k] = rd.choice(self.X)
198
199         # Initialisation of the sigma_k matrix as identity matrix
200         self.sigma = np.zeros((nb_clust, len(self.X[0]), len(self.X[0])))
201         for k in range(self.nb_clust):
202             self.sigma[k] = np.matrix(np.identity(len(self.X[0])), copy=False)
203
204         # Initialisation of the conditional posterior probability distribution
205         self.phi = np.zeros((len(self.X), nb_clust))
206
207         # Initialisation of pi as a uniform distribution

```



```

208         self.pi = np.random.uniform(low=0.0, high=1.0, size=nb_clust)
209
210     def EMGMM(self):
211         LX = len(self.X) # nb of datapoints
212         L_save = [] # list saving loss function at each iteration
213
214         ### do the following steps for self.nb_iter iterationsb ###
215
216         for iter in range(self.nb_iter):
217             print("iteration: %s" % (iter + 1))
218             # Initialize objective function at the beginning of each iteration
219             L = 0
220
221             ## E-step: generate phi posterior probability
222             # print("E-step")
223             # calculate for each datapoint x_i the sum of pi_k*N_k(mu_k,Sigma_k)
224             sum_pi_j_N = np.zeros((len(self.X)))
225             for i in range(LX):
226                 val = 0
227                 for k in range(self.nb_clust):
228                     val += self.pi[k] * MvN.pdf(self.X[i], mean=self.mu[k], cov=self.sigma[k])
229                     # print(MvN.pdf(self.X[i], mean=self.mu[k], cov=self.sigma[k]))
230                 sum_pi_j_N[i] = val
231
232             # calculate phi & objective function
233             L = 0 # Initialize objective function at the beginning of each iteration
234             for i in range(LX):
235                 val2 = 0
236                 for k in range(self.nb_clust):
237                     val = self.pi[k] * MvN.pdf(self.X[i], mean=self.mu[k], cov=self.sigma[k])
238                     val2 += val
239                     self.phi[i][k] = val / sum_pi_j_N[i]
240                 L += np.log(val2)
241             L_save.append(L)
242
243             ## M-step
244             # print("M-step")
245             # update empirical distribution pi_k using expected nb of pts n_k computed
246             for k in range(self.nb_clust):
247                 val = 0
248                 for i in range(LX):
249                     val += self.phi[i][k]
250                 self.nk[k] = val
251                 if self.nk[k] == 0:
252                     self.nk[k] = 1
253
254             for k in range(self.nb_clust):
255                 self.pi[k] = self.nk[k] / LX
256
257             # update mean mu_k
258             for k in range(self.nb_clust):
259                 val = np.zeros((len(self.X[0])))
260                 for i in range(LX):

```

```

261         # print(self.phi[i][k], self.X[i])
262         val += self.phi[i][k] * self.X[i]
263         self.mu[k] = val / self.nk[k]
264
265     # update covariance sigma_k
266     for k in range(self.nb_clust):
267         val = np.zeros((len(self.X[0]), len(self.X[0])))
268         for i in range(LX):
269             diff_iT = self.X[i] - self.mu[k]
270             diff_i = diff_iT[np.newaxis, :].T
271             aa = np.zeros((len(self.X[0]), len(self.X[0])))
272             bb = np.zeros((len(self.X[0]), len(self.X[0])))
273             bb[0] = diff_iT
274             for j in range(len(self.X[0])):
275                 aa[j][0] = diff_i[j]
276
277             val += self.phi[i][k] * np.dot(aa, bb)
278         self.sigma[k] = val / self.nk[k]
279
280     # save pi, phi, mu & sigma into files for each iteration
281     filename = "pi-" + str(iter + 1) + ".csv"
282     np.savetxt(filename, self.pi, delimiter=",")
283     filename = "phi-" + str(iter + 1) + ".csv"
284     np.savetxt(filename, self.phi, delimiter=",")
285     filename = "mu-" + str(iter + 1) + ".csv"
286     np.savetxt(filename, self.mu, delimiter=",") # this must be done at
287
288     for k in range(self.nb_clust): # k is the number of clusters
289         filename = "Sigma-" + str(k + 1) + "-" + str(
290             iter + 1) + ".csv" # this must be done 5 times (or the number
291         np.savetxt(filename, self.sigma[k], delimiter=",")
292
293     # save objective function
294     np.savetxt("Loss_function.csv", L_save, delimiter=",")
295
296     # Plot objective function
297     plt.plot(L_save, linewidth=3)
298     axes = plt.gca()
299     # axes.set_ylim([-0.5e6, 0])
300     plt.ticklabel_format(axis='y', style='sci') # , scilimits=(-2, 2))
301     plt.grid()
302     plt.xlabel('Iterations', fontsize=20)
303     plt.ylabel('Loss function', fontsize=20)
304     axes.tick_params(labelsize=14)
305     plt.show()
306
307     return self.pi, self.mu, self.sigma, self.phi
308
309 def plot_GMM(self, phi, nb):
310
311     # create two vectors with data point coordinates (we are considering to k
312     x = []
313     y = []

```

```

314         for i in range(len(self.X)):
315             x.append(self.X[i][0])
316             y.append(self.X[i][3])
317
318         # adjust phi for colors
319         col = phi[:]
320         for i in range(len(self.X)):
321             for k in range(nb):
322                 if col[i][k] < 0.01:
323                     col[i][k] = 0
324                 elif col[i][k] > 1:
325                     col[i][k] = 1
326                 # else:
327                 #     col[i][k] *= 1
328                 # col[i][k] = int(round(col[i][k]))
329
330         np.savetxt("colors.csv", col, delimiter=",")
331
332         groups_colors = {'cluster 1': 'r', 'cluster 2': [0, 1, 0], 'cluster 3':
333         # Create plot
334         fig = plt.figure()
335         # ax = fig.add_subplot(1, 1, 1, axisbg="1.0") # python 2.7
336         ax = fig.add_subplot(1, 1, 1)
337         # handles, labels = ax.get_legend_handles_labels()
338         # ax.legend(colors, groups)
339         patchList = []
340         for key in groups_colors:
341             data_key = mpatches.Patch(color=groups_colors[key], label=key)
342             patchList.append(data_key)
343
344         ax.legend(handles=patchList, loc=2)
345
346         for i in range(len(self.X)):
347             ax.scatter(x[i], y[i], alpha=0.8, c=[col[i][0], col[i][1], col[i][2],
348             s=40) # , label=group)
349
350         plt.xlabel('Attribute 1', fontsize=18)
351         plt.ylabel('Attribute 4', fontsize=18)
352         ax.tick_params(labelsize=14)
353
354         plt.show()
355
356
357     def plot_GMM_3D(self, phi, nb):
358         # create two vectors with data point coordinates (we are considering to b
359         x = []
360         y = []
361         z = []
362
363         for i in range(len(self.X)):
364             x.append(self.X[i][0])
365             y.append(self.X[i][1])
366             z.append(self.X[i][2])

```

```

367
368     # adjust phi for colors
369     col = phi[:]
370     for i in range(len(self.X)):
371         for k in range(nb):
372             if col[i][k] < 0.01:
373                 col[i][k] = 0
374             elif col[i][k] > 1:
375                 col[i][k] = 1
376
377
378     np.savetxt("colors.csv", col, delimiter=",")
379
380     groups_colors = {'cluster 1': 'r', 'cluster 2': [0, 1, 0], 'cluster 3':
381     # Create plot
382     fig = plt.figure()
383     ax = fig.add_subplot(111, projection='3d')
384
385     patchList = []
386     for key in groups_colors:
387         data_key = mpatches.Patch(color=groups_colors[key], label=key)
388         patchList.append(data_key)
389
390     ax.legend(handles=patchList, loc=2)
391
392     for i in range(len(self.X)):
393         ax.scatter(x[i], y[i], z[i], zdir='z', s=50, c=[col[i][0], col[i][1],
394
395     ax.set_xlabel('Attribute 1', fontsize=18)
396     ax.set_ylabel('Attribute 2', fontsize=18)
397     ax.set_zlabel('Attribute 3', fontsize=18)
398     tick_spacing = 1.0
399     ax.xaxis.set_major_locator(ticker.MultipleLocator(tick_spacing))
400     ax.yaxis.set_major_locator(ticker.MultipleLocator(tick_spacing))
401     ax.zaxis.set_major_locator(ticker.MultipleLocator(tick_spacing))
402     ax.tick_params(labelsize=12)
403
404     plt.show()
405
406
407 if __name__ == "__main__":
408     try:
409         assert len(sys.argv) > 5, "Missing input arguments"
410         iter = int(sys.argv[5]) # nb of iterations
411         nb = int(sys.argv[3])   # nb of clusters
412         Cl_type = sys.argv[4]   # type of clustering
413         if Cl_type == "Kmeans":
414             mu, ci = Kmeans(nb, iter).KMeans()
415             ct1, ct2, ct3 = Kmeans(nb, iter).check_prediction(ci)
416             Kmeans(nb, iter).plot_data(ci, mu, nb)
417             print("ct1 = %s; ct2 = %s; ct3 = %s" % (ct1, ct2, ct3))
418         elif Cl_type == "GMM":
419             pi, mu, sigma, phi = GMM(nb, iter).EMGMM()

```

```
420         GMM(nb, iter).plot_GMM(phi, nb)
421         GMM(nb, iter).plot_GMM_3D(phi, nb)
422     else:
423         print("Clustering type not known")
424 except Exception as e:
425     print(e)
```