

Homework 6 – Machine Learning

REPORT

Max LEMASQUERIER
0845034

I – CODE

```
1  import numpy as np
2  import matplotlib.pyplot as plt
3  import imageio
4  from scipy.spatial.distance import *
5  from tqdm import tqdm
```

In these first line, I'm just importing the useful libraries.

```
7  def read_input(image):
8      im = imageio.imread(image)
9      im=im.reshape((-1,3))
10     return im
```

Then here is the function that I used to load the image data into a variable. During my work, I've focused on the *image1.png* and not on the image 2, because the computation were already quite long to have good results with the image 1, I just had no time to do the same computations on the image 2 as well. But it should work all the same for both images.

```
12  def compute_rbf_kernel(image, gamma1, gamma2):
13      kernel = np.zeros((image.shape[0], image.shape[0]))
14      temp=[]
15      for i in range(100):
16          for j in range(100):
17              temp.append([i,j])
18      temp = np.asarray(temp)
19      temp2 = np.exp(-gamma1*cdist(temp,temp))
20      color = cdist(image,image)
21      kernel = np.multiply(temp2,np.exp(-gamma2*color))
22      return kernel
```

```

24 def initialization(data, k):
25     """
26     Function to initialize the means of the different clusters
27     and the classification first step.
28     """
29     n = data.shape[0]
30     initialize_method = "2"
31     means = np.random.rand(k,3)*255
32     if initialize_method == "1": # RANDOM INIT
33         classif_prec = np.random.randint(k, size=n)
34     elif initialize_method == "2": # Structured init : every 2 columns
35         classif_prec = []
36         for i in range(n):
37             if i % 2 == 1:
38                 classif_prec.append(0)
39             else:
40                 classif_prec.append(1)
41         classif_prec = np.asarray(classif_prec)
42     elif initialize_method == "3": # More accurate method, based on the random mean
43         classif_prec = np.zeros(n, dtype=np.int)
44         temp = np.zeros(n)
45         null_vector = np.zeros([1, 2])
46         for i in range(0, n):
47             temp[i] = np.linalg.norm(data[i,:] - null_vector[0,:])
48         mean_temp = np.mean(temp)
49         for i in range(0, n):
50             if temp[i] >= mean_temp:
51                 classif_prec[i] = 0
52             else:
53                 classif_prec[i] = 1
54     return means, np.asarray(classif_prec)

```

This is the function I used to compute the kernel matrix of the image, including as asked in the statement : a part concerning the spatial distance of the pixels, and another part concerning the color similarities, that I described as just a spatial distance between 2 different rgb vectors.

This part is about the initialization of the clustering. It turned out that it is quite important, and changes the results I could obtain by changing the initialization method. I thought it would end up being the same but in fact not really. So in this function I initialize the means of the different clusters, and a first classification of the pixels, which is quite important for how the algorithm will classify the pixels next. There are 3 different methods here, my favorite is the random one because it really starts from nothing to get to some quite impressive clustering. But I think the third method is much closer from an accurate description of the image.

```

56 def deuxieme_terme(data, kernel_data, classification, data_number, cluster_number, k):
57     result = 0
58     number_in_cluster = 0
59     for i in range(0, data.shape[0]):
60         if classification[i] == cluster_number:
61             number_in_cluster += 1
62     if number_in_cluster == 0:
63         number_in_cluster = 1
64     for i in range(0, data.shape[0]):
65         if classification[i] == cluster_number:
66             result += kernel_data[data_number][i]
67     return -2 * (result / number_in_cluster)

```

This function is computing the second term of the euclidean distance from a point to the center of the different clusters. This is important to determine afterwards, which cluster we will attribute to every pixel. And below is a similar function to compute the third term of this distance :

```

69 def troisieme_terme(kernel_data, classification, k):
70     """
71     Function to compute the third term of the euclidean distance
72     from a point to the center of the different clusters.
73     """
74     temp = np.zeros(k)
75     temp1 = np.zeros(k)
76     for i in range(0, classification.shape[0]):
77         temp[classification[i]] += 1
78     for i in range(0, k):
79         for p in range(0, kernel_data.shape[0]):
80             for q in range(p + 1, kernel_data.shape[1]):
81                 if classification[p] == i and classification[q] == i:
82                     temp1[i] += kernel_data[p,q]
83     for i in range(0, k):
84         if temp[i] == 0:
85             temp[i] = 1
86         temp1[i] /= (temp[i] ** 2)
87     return temp1

```

Next is the classifier :

```

89 def classifier(data, kernel_data, means, classification):
90     """
91     Attribute a cluster to every pixel of an image, based on the kernel results
92     """
93     temp_classification = np.zeros([data.shape[0]], dtype=np.int)
94     third_term = troisieme_terme(kernel_data, classification, means.shape[0])
95     for i in tqdm(range(0, data.shape[0])):
96         temp = np.zeros([means.shape[0]], dtype=np.float32) # temp size: k
97         for j in range(0, means.shape[0]):
98             temp[j] = deuxieme_terme(data, kernel_data, classification, i, j, means.shape[0]) + third_term[j]
99         temp_classification[i] = np.argmin(temp)
100     return temp_classification

```

As explained above, it uses the other functions in order to compute the full euclidean distance, and takes the minimum value of a corresponding cluster to attribute it to a data point, which is here a pixel.

```

102 def nb_erreur(classification, classif_prec):
103     """
104     Count the number of pixels that have change the cluster between an iteration
105     """
106     error = 0
107     for i in range(0, classification.shape[0]):
108         error += np.absolute(classification[i] - classif_prec[i])
109     return error

```

This function is counting the number of pixels that left a cluster to join a new one, on every iteration so that we can stop when every pixel is fixed, and not changing it's belonging to a cluster anymore.


```

123 def display_clusters(k, data, means, classification, iteration, filename):
124     title = "Kernel-K-Means Iteration-" + str(iteration)
125     plt.clf()
126     plt.suptitle(title)
127     plt.imshow(classification.reshape((100,100)), cmap='gray')
128     plt.show()

```

This function is just a displayer of figure. Once I classified every pixel to different clusters, I can plot a map of the assignment to see whether it looks like the original picture or not. I will show some results after.

```

131 def kkmeans(data, kernel_data, filename):
132     k = 2 # cluster number
133     means, classif_prec = initialization(data, k) # INITIALIZE everything
134     iteration, erreur_prec = 1, 0
135     display_clusters(k, data, means, classif_prec, iteration, filename) # display the initial assignment of clusters
136     classification = classifier(data, kernel_data, means, classif_prec)
137     error = nb_erreur(classification, classif_prec)
138     for i in range(50): # Limit to 50 iteration, after that if it did not converge already, it will stop.
139         display_clusters(k, data, means, classification, iteration, filename)
140         iteration += 1
141         classif_prec = classification
142         classification = classifier(data, kernel_data, means, classif_prec)
143         error = nb_erreur(classification, classif_prec)
144         print(error)
145         if error == erreur_prec:
146             break
147         erreur_prec = error
148     means = update(data, means, classification) # update the clusters mean to have the final centroids
149     display_clusters(k, data, means, classification, iteration, filename) # display the final assignment.

```

This is the main function, not the hardest one, because it just runs the other functions. So we just do the usual process, initialization, classification and iterate until it converges or reaches a max iteration number.

II – RESULTS

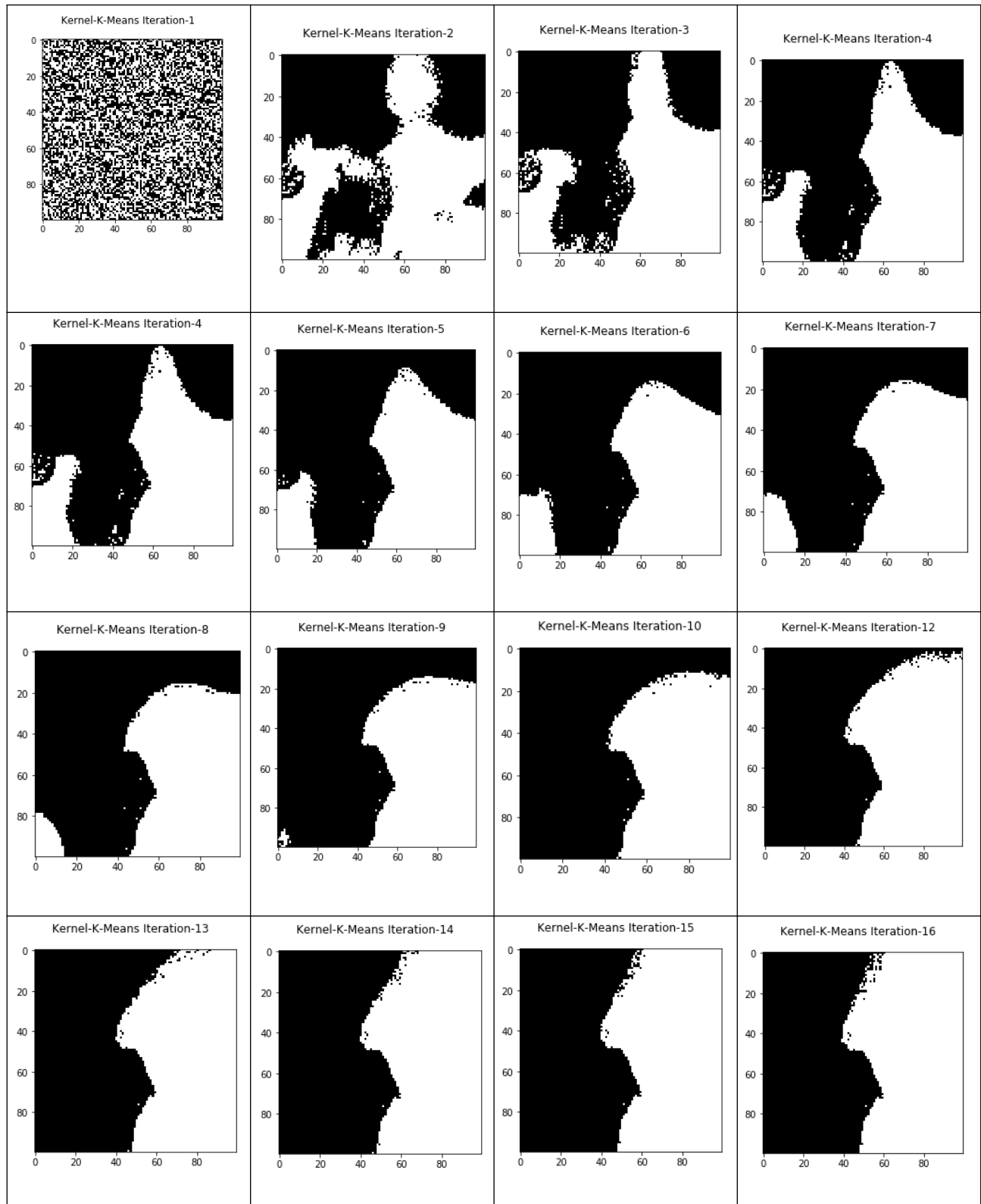
From what I understood, the spectral is the same as k mean, but using the output cut and ratio cut as input for the spectral clustering.

So I will only introduce the results I had with the kernel k-means algorithm. Every iteration was about 5 minute to compute because the third term of the euclidean distance was quite long to compute, and the classification step was therefore quite long. Plus I had to loop over 10 000 two times which also takes a while.

Here is our base image :

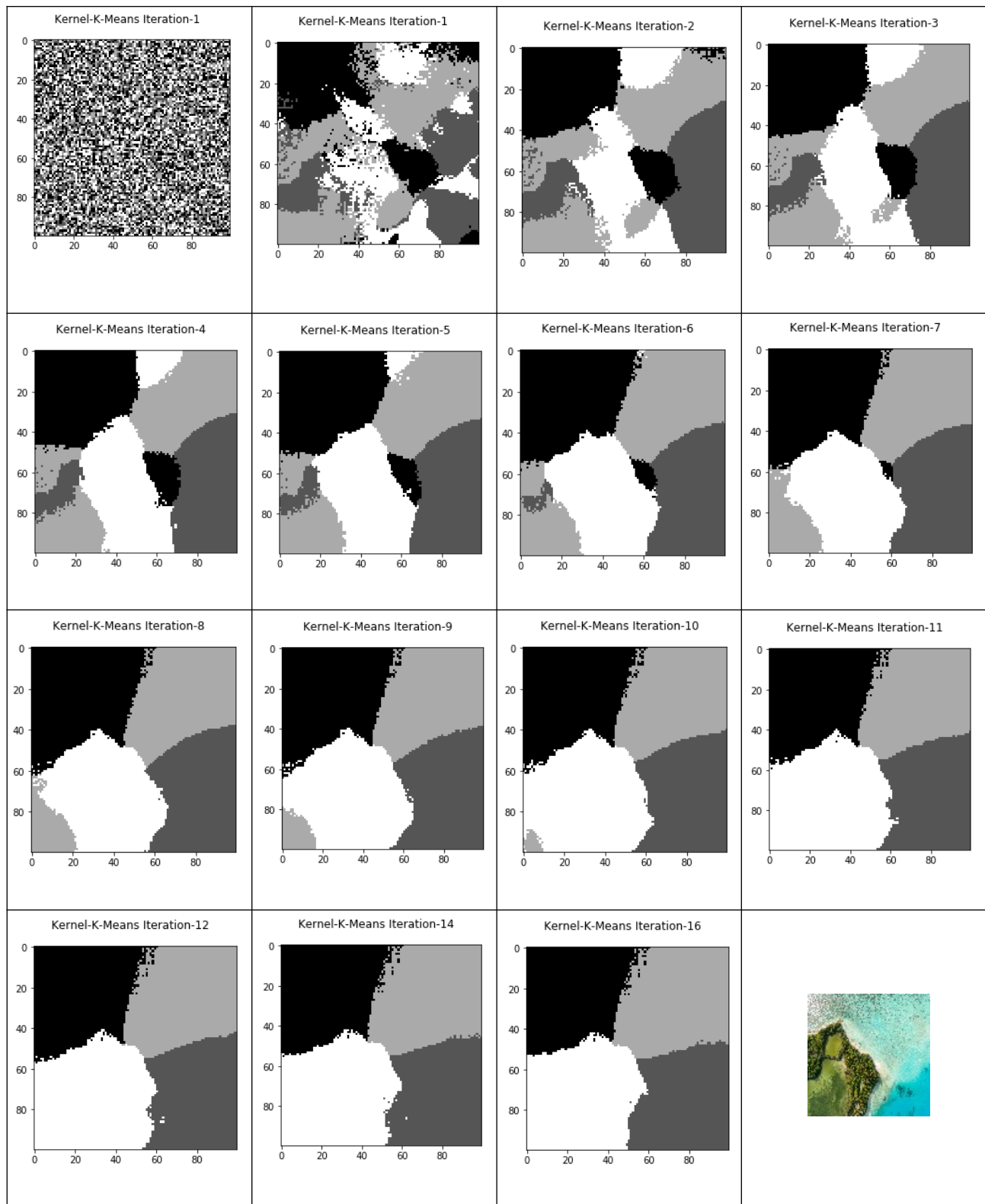


First, with a random initialization and 2 clusters :



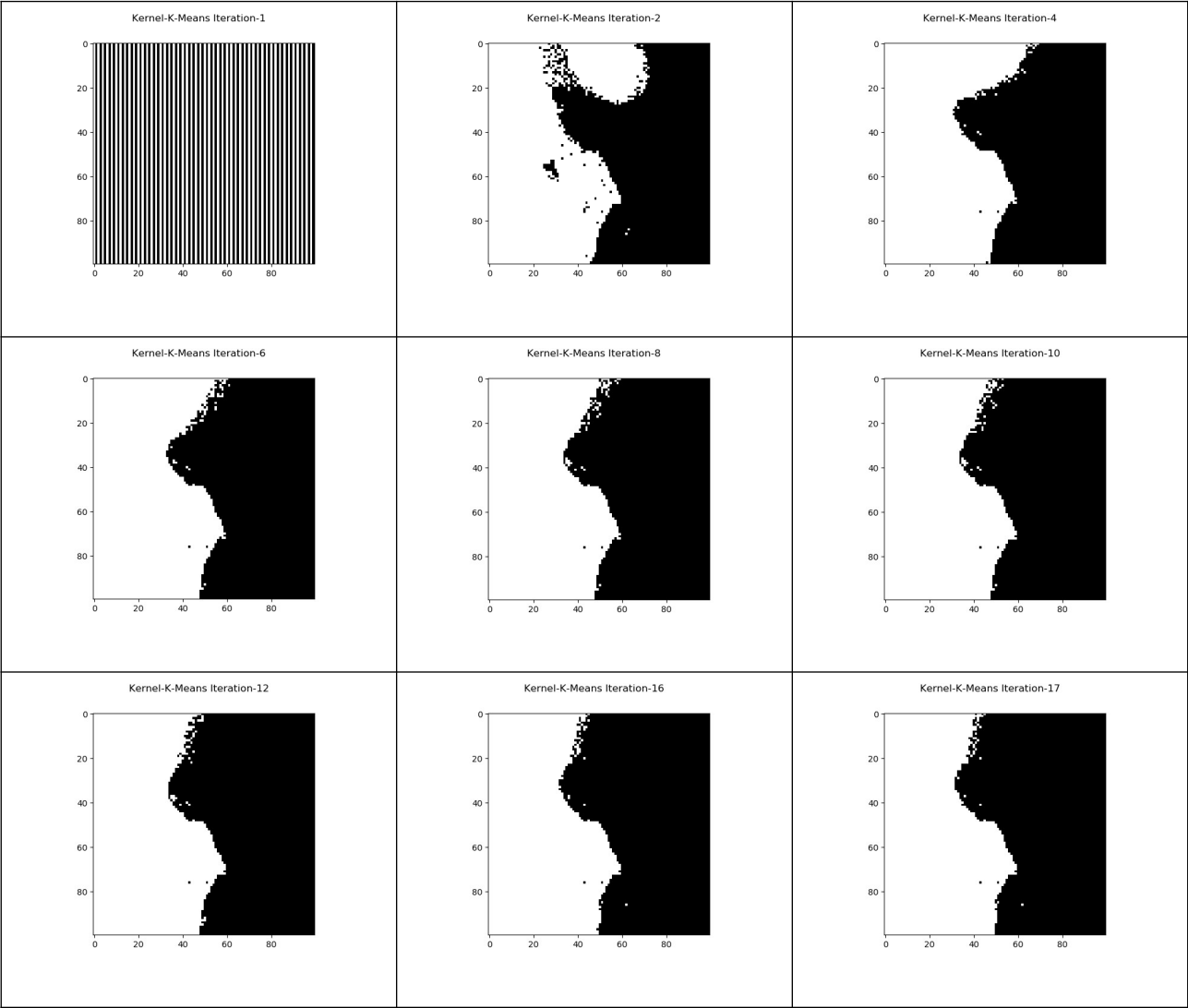
You can see the gif as well which should be in the compressed file. We see a nice evolution after iterating.

With the same initialization with 4 clusters :

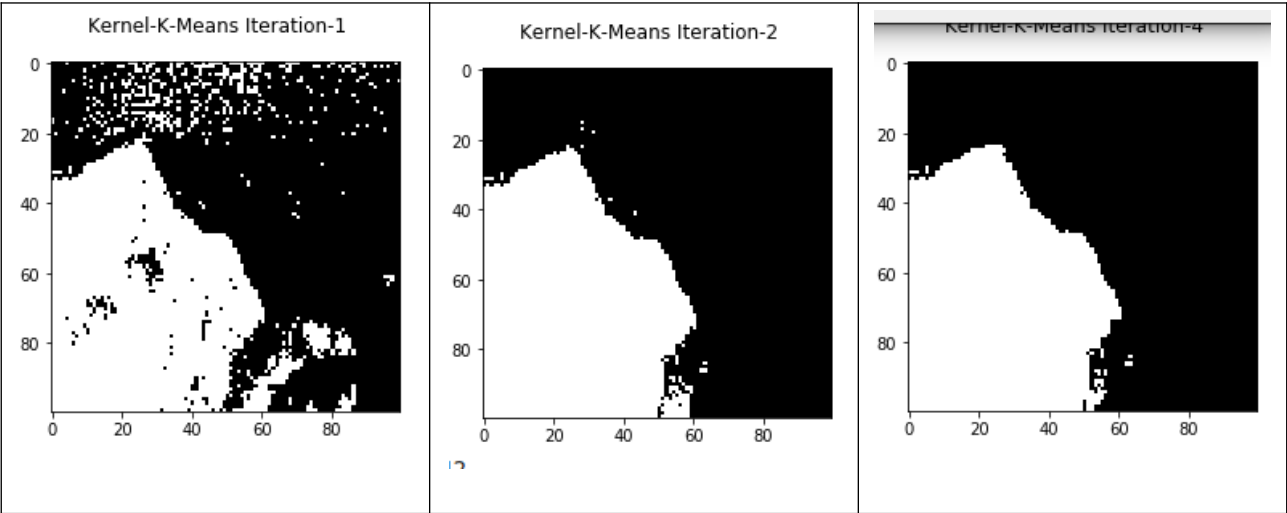


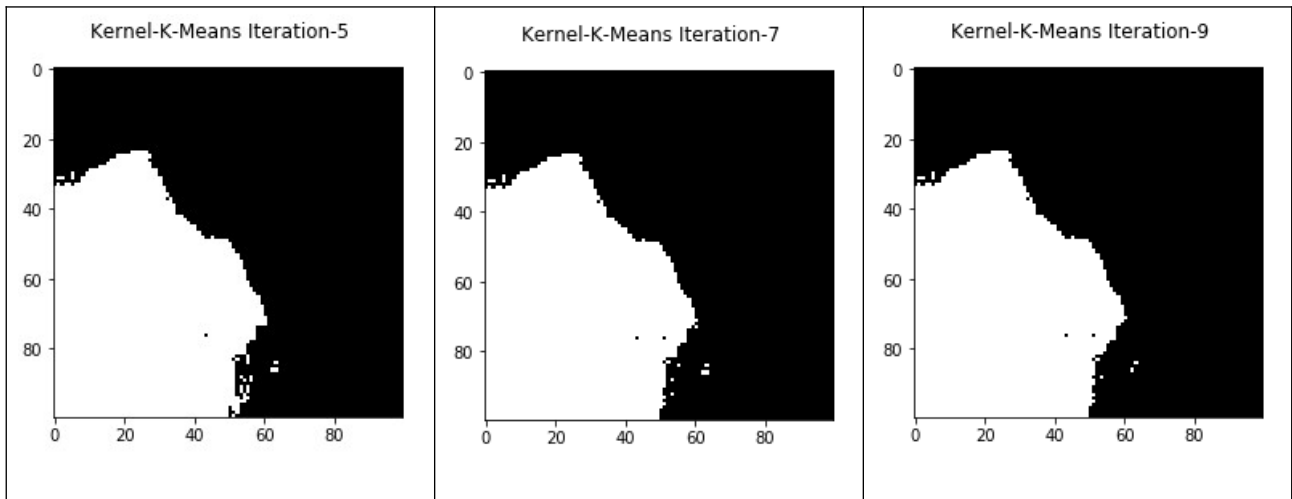
We recognize in a better way in my opinion the different regions of the initial picture. Very happy about this result.

Then with different init method :



And :





So for clustering with 2 regions only, this method is way better than the previous ones. We clearly identify the land and the sea here.

The spectral clustering aims to take the eigenvalues of the kernel matrix in order to reduce the dimensionality of the problem. I didn't have time to fully run the spectral clustering because I still have to debug a few stuff. But I missed some time, so I can't discuss about the eigenspace of graph Laplacian.

The idea was to compute the Laplacian, then to compute its eigenvalues and eigenvectors before using a very similar to k-means algorithm.