

309551063 HW6 Report

a. code with detailed explanations

Part 1

the procedure of making gif image

```
94     def make_gif_image(cluster_idx_history, input_image, k, final_center, gif_name='test'):
95         """
96             Use the center to colorize the same cluster, and use imageio to make gif file.
97             :param cluster_idx_history: the cluster of each pixel in every iteration
98             :param input_image: the input image (shape: (10000, 3))
99             :param k: the number of clusters
100            :param final_center: the center after trained
101            :param gif_name: output gif file name
102            :return: None
103        """
104        # use imageio package to generate gif
105        images_frame = list()
106        for idx in cluster_idx_history:
107            image_frame = np.zeros(input_image.shape, dtype=np.uint8)
108            for c_idx in range(k):
109                # use the center to colorize the same cluster
110                image_frame[np.where(idx == c_idx)[0], :] = final_center[c_idx]
111                images_frame.append(image_frame.reshape(100, 100, 3))
112        imageio.mimsave(f'gif/kernel-k-means/{gif_name}.gif', images_frame)
113        # use *_tmp.gif file to show the final result in report
114        imageio.mimsave(f'gif/kernel-k-means/{gif_name}_tmp.gif', [images_frame[-1], images_frame[-1]])
```

kernel k-means clustering procedure

```
8     def kernel_function(spatial_a, spatial_b, color_a, color_b, theta=None):
9         """
10            Calculate the kernel noticed in the spec.
11            $k(x, x') = e^{\gamma_s ||S(x) - S(x')||^2} \times e^{\gamma_c ||C(x) - C(x')||^2}$
12            :param spatial_a: The spatial information of x
13            :param spatial_b: The spatial information of x'
14            :param color_a: The color information of x
15            :param color_b: The color information of x'
16            :param theta: The hyper-parameters used in kernel.
17                A list, format: [\gamma_s, \gamma_c]. (default: [1e-5, 1e-5])
18            :return: the mix kernel
19        """
20        if theta is None:
21            theta = np.array([1e-5, 1e-5]).reshape(-1, 1)
22            spatial_information = cdist(spatial_a, spatial_b, 'sqeuclidean')
23
24            color_information = cdist(color_a, color_b, 'sqeuclidean')
25            return np.exp(-theta[0] * spatial_information) * np.exp(-theta[1] * color_information)
```

$$k(x, x') = e^{\gamma_s ||S(x) - S(x')||^2} \times e^{\gamma_c ||C(x) - C(x')||^2}$$

```

27     def initialize_label(input_image, k, method='random'):
28         """
29             Initialize the label of each pixel.
30             :param input_image: the input image (shape: (10000, 3))
31             :param k: the number of clusters
32             :param method: There are two initialized method can choose, {'random', 'k-means++'}.
33             (Default: 'random')
34             :return: return the initialized label of each pixel
35         """
36         np.random.seed(666)
37         if method == 'random':
38             # random to initialize the label
39             return np.random.randint(k, size=len(input_image))
40         elif method == 'k-means++':
41             """

```

```

42         k-means++ algorithm
43         step1: sample a pixel randomly to be center
44         step2: calculate the minimum distance D(x) between pixel and the nearest center
45         step3: choose a new center based on the probability of  $\frac{D(x)^2}{\sum D(x)^2}$ 
46         repeat step2 and step 3 until get enough center
47         """
48         center = list()
49         center.append(np.random.randint(len(input_image)))
50         total_distance = kernel_k_means_square_distance_point_point()
51         while len(center) != k:
52             distance = np.min(total_distance[center], axis=0)
53             prob = distance / sum(distance)
54             center.append(np.random.choice(len(input_image), 1, p=list(prob))[0])
55         # Find the the nearest center of each pixel to get label
56         cluster_idx = np.argmin(total_distance[center], axis=0)

```

```

57         return cluster_idx
58     else:
59         # Undefined method
60         raise ValueError('Unknown initialize label method!')

```

```

63     def kernel_k_means_square_distance_point_point():
64         """
65             Calculate the square distance based on the kernel.
66              $\|\phi(x_i) - \phi(x_j)\|^2 = \phi(x_i)\phi(x_i) - 2\phi(x_i)\phi(x_j) + \phi(x_j)\phi(x_j) = k(x_i, x_i) - 2k(x_i, x_j) + k(x_j, x_j)$ 
67             :return: the square distance between point and point (shape: (10000, 10000))
68         """
69         tmp = total_kernel.diagonal().reshape(-1, 1)
70         return tmp - 2 * total_kernel + tmp

```

$$\|\phi(x_i) - \phi(x_j)\| = \phi(x_i)\phi(x_i) - 2\phi(x_i)\phi(x_j) + \phi(x_j)\phi(x_j) = k(x_i, x_i) - 2k(x_i, x_j) + k(x_j, x_j)$$

```

74     def kernel_k_means_distance_point_center(target_idx):
75         """
76             Use the formula in slide P22 to calculate the square distance between point and the center
77             $||\phi(x_j)-\mu^{\phi}k||=\\
78             ||\phi(x_j)-\frac{1}{|C_k|}\sum_{n=1}^N\alpha_{kn}\phi(x_n)||\\
79             =k(x_j,x_j)-\frac{1}{|C_k|}\sum_n\alpha_{kn}k(x_j,x_n)+\\
80             \frac{1}{|C_k|^2}\sum_p\sum_q\alpha_{kp}\alpha_{kq}k(x_p,x_q)$
81             :param target_idx: the index of the specified center
82             :return: the square distance between point and the center
83         """
84         first_term = total_kernel.diagonal()
85         second_term = 2 / len(target_idx)
86         alpha = np.zeros(total_kernel.shape[0])
87         alpha[target_idx] = 1
88         second_term = second_term * np.dot(total_kernel.T, alpha)
89         third_term = 1 / (len(target_idx) ** 2)

90         third_term = third_term * np.sum(total_kernel[np.ix_(target_idx, target_idx)])
91     return first_term - second_term + third_term

```

$$\|\phi(x_j) - \mu^\phi k\| = \|\phi(x_j) - \frac{1}{|C_k|} \sum_{n=1}^N \alpha_{kn} \phi(x_n)\| = k(x_j, x_j) - \frac{2}{|C_k|} \sum_n \alpha_{kn} k(x_j, x_n) + \frac{1}{|C_k|^2} \sum_p \sum_q \alpha_{kp} \alpha_{kq} k(x_p, x_q)$$

```

116     def kernel_k_means(input_image, gif_name, k=2, max_iter=100, tol=1e-4, method='random'):
117         """
118             kernel k-means algorithm
119             step1: use initialized method to get the cluster of each pixel
120             step2: calculate the distance between pixel and center
121             step3: update the cluster of each pixel based on the nearest center
122             repeat step2 and step3 until convergence or reach the max iteration
123             :param input_image: the input image (shape: (10000, 3))
124             :param gif_name: the output gif file name
125             :param k: the number of clusters (default: 2)
126             :param max_iter: the maximum iteration (default: 100)
127             :param tol: the tolerance to check convergence (default: 1e-4)
128             :param method: the center initialized method (default: 'random')
129             :return: the cluster of each pixel
130         """
131         old_cluster_distance = np.array(

```

```

132         [0 for _ in range(k)] * len(input_image)).reshape(-1, len(input_image))
133     cluster_idx = initialize_label(input_image, k, method)
134     # record the cluster of each pixel in every iteration to output gif
135     cluster_idx_history = [cluster_idx]
136     final_center = list()
137     final_iter_count = 0
138     for iter_count in range(max_iter):
139         cluster_distance = list()
140         for c_idx in range(k):
141             cluster_distance.append(
142                 kernel_k_means_distance_point_center(np.where(cluster_idx == c_idx)[0]))
143             cluster_distance = np.array(cluster_distance)
144             cluster_idx = np.argmin(cluster_distance, axis=0)
145             cluster_idx_history.append(cluster_idx)
146             if np.linalg.norm(old_cluster_distance - cluster_distance) < tol:
147                 final_iter_count = iter_count + 1

```

```

148         break
149     else:
150         old_cluster_distance = cluster_distance
151     # find the final center
152     for c_idx in range(k):
153         final_center.append(np.mean(input_image[np.where(cluster_idx == c_idx)[0], :], axis=0))
154     make_gif_image(cluster_idx_history, input_image, k, final_center,
155                     gif_name=f'{gif_name}_{final_iter_count}')
156     return cluster_idx

```

```

159 ► if __name__ == '__main__':
160     # use argparse to get all argument
161     parser = argparse.ArgumentParser()
162     parser.add_argument("-i", "--image", default='image1.png',
163                         help="input image: image1.png or image2.png (Default: image1.png)")
164     parser.add_argument("-m", "--method", default='k-means++',
165                         help="initialize method: k-means++ or random (Default: k-means++)")
166     parser.add_argument("--theta1", default=1e-5, type=float,
167                         help="hyper-parameters of spatial information (Default: 1e-5)")
168     parser.add_argument("--theta2", default=1e-5, type=float,
169                         help="hyper-parameters of color information (Default: 1e-5)")
170     parser.add_argument("-k", "--cluster", default=2, type=int,
171                         help="the number of clusters (Default: 2)")
172     args = parser.parse_args()
173     initialize_method = args.method
174     image_name = args.image

```

```

175     cluster = args.cluster
176     theta1 = args.theta1
177     theta2 = args.theta2
178     # read image
179     image = cv2.imread(image_name).reshape(-1, 3)
180     spatial = list()
181     for x in range(0, 100):
182         for y in range(0, 100):
183             spatial.append([x, y])
184     spatial = np.array(spatial)
185     output_name = f'{image_name}_{cluster}_{initialize_method}_{theta1}_{theta2}'
186     # use spatial and image information to get kernel
187     total_kernel = kernel_function(spatial, spatial, image, image, theta=[theta1, theta2])
188     kernel_k_means(image, gif_name=output_name, k=cluster, method=initialize_method)
189

```

Brief summary of kernel k-means clustering procedure:

- step1: use `argparse` package to get the parameters
- step2: use spatial information and color information to get kernel
- step3: use initialized method to get the cluster of each pixel
- step4: calculate the distance between pixel and center
- step5: update the cluster of each pixel based on the nearest center
(repeat step4 and step5 until convergence or reach the max iteration)
- step6: generate gif image file

spectral clustering (both normalized cut and ratio cut) clustering procedure

```

31     def initialize_label(input_image, k, method='random'):
32         """
33             Initialize the label of each pixel.
34             :param input_image: the input image (shape: (10000, 3))
35             :param k: the number of clusters
36             :param method: There are two initialized method can choose, {'random', 'k-means+'}.
37             (Default: 'random')
38             :return: return the initialized label of each pixel
39         """
40         np.random.seed(666)
41         if method == 'random':
42             return np.random.randint(k, size=len(input_image))
43         elif method == 'k-means+':
44             """
45                 k-means++ algorithm
46                 step1: sample a pixel randomly to be center

```

```

47     step2: calculate the minimum distance D(x) between pixel and the nearest center
48     step3: choose a new center based on the probability of  $\frac{D(x)^2}{\sum D(x)^2}$ 
49     repeat step2 and step 3 until get enough center
50     """
51
52     center = list()
53     center.append(np.random.randint(len(input_image)))
54     while len(center) != k:
55         distance = np.min(total_distance[center], axis=0)
56         prob = distance / sum(distance)
57         center.append(np.random.choice(len(input_image), 1, p=list(prob))[0])
58     # Find the the nearest center of each pixel to get label
59     cluster_idx = np.argmin(total_distance[center], axis=0)
60
61     return cluster_idx
62 else:
63     raise ValueError('Unknown method!')

```

```

87 def get_laplacian(method='ratio_cut'):
88     """
89     The weighted adjacency matrix  $W=(w_{ij})_{i,j=1,\dots,n}$ 
90     The degree:  $d_i=\sum_{j=1}^n w_{ij}$ 
91     The degree matrix D: the diagonal matrix with the degrees  $d_1,\dots,d_n$  on the diagonal.
92     The unnormalized Laplacian matrix L:  $D - W$ 
93     The normalized Laplacian matrix L:  $D^{-1/2}LD^{-1/2}$ 
94     (Source: A Tutorial on Spectral Clustering)
95     :param method: The cut method used in spectral clustering.
96     There are two methods can choose: {'normalized_cut', 'ratio_cut'} (Default: 'ratio_cut')
97     :return: The Laplacian matrix
98     """
99
100    # unknown cut method
101    if method not in ['normalized_cut', 'ratio_cut']:
102        raise ValueError('Unknown cut method!')
103    degree = np.sum(total_kernel, axis=0)

```

The weighted adjacency matrix $W = (w_{ij})_{i,j = 1, \dots, n}$

The degree: $d_i = \sum_{j=1}^n w_{ij}$

The degree matrix D: the diagonal matrix with the degrees d_1, \dots, d_n on the diagonal.

The unnormalized Laplacian matrix L : $D - W$

The normalized Laplacian matrix L : $D^{-1/2}LD^{-1/2}$

(Source: [A Tutorial on Spectral Clustering](#)).

```

112     def plot_eigen_space(u, cluster_idx, final_center, k, file_name='test'):
113         """
114             Plot the eigen space of graph Laplacian.
115             If k is 2, output the 2D scatter plot.
116             If k is 3, output the 3D scatter plot.
117             If k is bigger than 3, output the 1D scatter plot of each dimension.
118             :param u: the first k eigenvectors(the eigenvectors corresponding to the k smallest eigenvalues)
119             :param cluster_idx: the cluster of each pixel
120             :param final_center: the center after trained
121             :param k: the number of clusters
122             :param file_name: output jpg file name
123             :return: None
124         """
125         if k == 2:
126             for c_idx in range(k):
127                 idx = np.where(cluster_idx == c_idx)[0]
128                     plt.scatter(u[idx, 0], u[idx, 1], c=tuple(final_center[c_idx] / 255))
129             elif k == 3:
130                 fig = plt.figure()
131                 ax = fig.gca(projection='3d')
132                 for c_idx in range(k):
133                     idx = np.where(cluster_idx == c_idx)[0]
134                     ax.scatter(u[idx, 0], u[idx, 1], u[idx, 2], c=tuple(final_center[c_idx] / 255))
135             else:
136                 fig, axs = plt.subplots(k, figsize=(15, 15))
137                 count = 0
138                 for i in range(k):
139                     axs[i].set_title(f'eigenspace - the {i} dimension')
140                     for c_idx in range(k):
141                         idx = np.where(cluster_idx == c_idx)[0]
142                         axs[i].scatter(range(count, count+len(idx)), u[idx, i],
143                                         c=tuple(final_center[c_idx] / 255))
144                         count = count + len(idx)
145             plt.tight_layout()
146             plt.savefig(f'gif/spectral-clustering/{file_name}_eigen_space.jpg')

```

```

149     def k_means(input_image, u, gif_name, k=2, max_iter=100, tol=1e-4, method='random'):
150         """
151             k-means algorithm
152             step1: use initialized method to get the cluster of each pixel
153             step2: calculate the distance between pixel and center
154             step3: update the cluster of each pixel based on the nearest center
155             repeat step2 and step3 until convergence or reach the max iteration
156             :param input_image: the input image (shape: (10000, 3))
157             :param u: the first k eigenvectors of Laplacian
158             :param gif_name: the output gif file name
159             :param k: the number of clusters (default: 2)
160             :param max_iter: the maximum iteration (default: 100)
161             :param tol: the tolerance to check convergence (default: 1e-4)
162             :param method: the center initialized method (default: 'random')
163             :return: the cluster of each pixel
164         """
165
166         old_cluster_distance = np.array([0 for _ in range(k)] * len(u)).reshape(-1, len(u))      ▲1 ✘3
167         cluster_idx = initialize_label(u, k, method)
168         center = list()
169         for _ in range(k):
170             center.append(u[np.random.randint(len(u))])
171         cluster_idx_history = [cluster_idx]
172         final_center = list()
173         final_iter_count = 0
174         for iter_count in range(max_iter):
175             cluster_distance = list()
176             for c_idx in range(k):
177                 cluster_distance.append(cdist(u, center[c_idx].reshape(1, -1), 'euclidean').flatten())
178             cluster_distance = np.array(cluster_distance)
179             cluster_idx = np.argmin(cluster_distance, axis=0)
180             cluster_idx_history.append(cluster_idx)
181             if np.linalg.norm(old_cluster_distance - cluster_distance) < tol:
182                 final_iter_count = iter_count + 1
183                 break
184             else:
185                 old_cluster_distance = cluster_distance
186                 for c_idx in range(k):
187                     center[c_idx] = np.mean(u[np.where(cluster_idx == c_idx)[0], :], axis=0)
188                 for c_idx in range(k):
189                     final_center.append(np.mean(input_image[np.where(cluster_idx == c_idx)[0], :], axis=0))
190                 make_gif_image(cluster_idx_history, input_image, k, final_center,
191                               gif_name=f'{gif_name}_{final_iter_count}')
192                 plot_eigen_space(u, cluster_idx, final_center, k, file_name=f'{gif_name}_{final_iter_count}')
193         return cluster_idx

```

```

195 |     def spectral_clustering(input_image, k, gif_name,
196 |                               spectral_clustering_method='ratio_cut', initialize_method='random'):
197 |         """
198 |             :param input_image: the input image (shape: (10000, 3))
199 |             :param k: the number of clusters
200 |             :param gif_name: the output gif file name
201 |             :param spectral_clustering_method: The cut method used in spectral clustering.
202 |                 There are two methods can choose: {'normalized_cut', 'ratio_cut'} (Default: 'ratio_cut')
203 |             :param initialize_method: There are two initialized method can choose, {'random', 'k-means++'}.
204 |                 (Default: 'random')
205 |             :return: None
206 |         """
207 |         laplacian = get_laplacian(spectral_clustering_method)
208 |         # compute eigenvectors
209 |         if os.path.exists(f'{args.image}_{args.cut_method}_{theta1}_{theta2}_a.npy'):
210 |             a = np.load(f'{args.image}_{args.cut_method}_{theta1}_{theta2}_a.npy')

211 |             b = np.load(f'{args.image}_{args.cut_method}_{theta1}_{theta2}_b.npy')
212 |             print('npy file loaded')
213 |         else:
214 |             a, b = np.linalg.eigh(laplacian)
215 |             np.save(f'{args.image}_{args.cut_method}_{theta1}_{theta2}_a', a)
216 |             np.save(f'{args.image}_{args.cut_method}_{theta1}_{theta2}_b', b)
217 |             # find the first k eigenvectors(the eigenvectors corresponding to the k smallest eigenvalues)
218 |             # and drop the first column
219 |             u = b[:, a.argsort()[1:k+1]]
220 |             if spectral_clustering_method == 'normalized_cut':
221 |                 # if spectral clustering method is 'normalized_cut', we need to normalize the rows to norm 1
222 |                 u = u / cdist(u, np.array([0] * k).reshape(1, -1), 'euclidean')
223 |             return k_means(input_image=input_image, u=u, gif_name=gif_name, k=k, method=initialize_method)

224 ► if __name__ == '__main__':
225 |     # use argparse to get all argument
226 |     parser = argparse.ArgumentParser()
227 |     parser.add_argument("-i", "--image", default='image1.png',
228 |                         help="input image: image1.png or image2.png (Default: image1.png)")
229 |     parser.add_argument("--initialize_method", default='k-means++',
230 |                         help="initialize method: k-means++ or random (Default: k-means++)")
231 |     parser.add_argument("--cut_method", default='normalized_cut',
232 |                         help="clustering method: normalized_cut or ratio_cut (Default: normalized_cut)")
233 |     parser.add_argument("--theta1", default=0.001, type=float,
234 |                         help="hyper-parameters of spatial information (Default: 0.001)")
235 |     parser.add_argument("--theta2", default=0.001, type=float,
236 |                         help="hyper-parameters of color information (Default: 0.001)")
237 |     parser.add_argument("-k", "--cluster", default=2, type=int,
238 |                         help="the number of clusters (Default: 2)")
239 |     args = parser.parse_args()

```

```

242     image_name = args.image
243     cluster = args.cluster
244     theta1 = args.theta1
245     theta2 = args.theta2
246     image = cv2.imread(image_name).reshape(-1, 3)
247     spatial = list()
248     for x in range(0, 100):
249         for y in range(0, 100):
250             spatial.append([x, y])
251     spatial = np.array(spatial)
252     total_kernel = kernel_function(spatial, spatial, image, image, theta=[theta1, theta2])
253     output_name = f'{args.image}_{args.cut_method}_{args.initialize_method}_{cluster}_{theta1}_{theta2}'
254     spectral_clustering(image, cluster, output_name,
255                         spectral_clustering_method=args.cut_method,
256                         initialize_method=args.initialize_method)

```

Brief summary of spectral clustering procedure:

- step1: use `argparse` package to get the arguments
- step2: use spatial information and color information to get kernel
- step3: calculate Laplacian
The unnormalized Laplacian matrix $L : D - W$
The normalized Laplacian matrix $L : D^{-1/2}LD^{-1/2}$
- step4: compute the first k eigenvectors of L
- step5: if spectral clustering method is `normalized_cut`, we need to normalize the rows to norm 1
- step6: Use the eigenvectors as data points to do k-means algorithm
- step7: generate gif image file and plot the eigenspace of graph Laplacian

Part 2

If we want to try more clusters, such as 4 clusters, we can run the program with argument `-k`.

- Example:

```

$ python kernel-k-means.py -k 4
$ python spectral_clustering.py -k 4

```

Part 3

For the initialization, I wrote a function `initialize_label` in order to get the initialized label of each pixel(data point). I implemented two methods, `random` and `k-means++`.

- `random`
 - random initialized label of each pixel
- `k-means++`
 - step1: sample a pixel randomly to be center
 - step2: calculate the minimum distance $D(x)$ between pixel and the nearest center
 - step3: choose a new center based on the probability of $\frac{D(x)^2}{\sum D(x)^2}$
repeat step2 and step 3 until get enough center

The main difference of `kernel k-means` and `spectral clustering` in `k-means++` initialization is the distance calculation.

- `kernel k-means` : use kernel to calculate distance
- `spectral clustering` : use eigenvectors to calculate distance

Part 4

In order to plot the eigenspace of graph Laplacian, I wrote a function `plot_eigen_space`.

- If k is 2, output the 2D scatter plot.
- If k is 3, output the 3D scatter plot.
- If k is bigger than 3, output the 1D scatter plot of each dimension.

b. experiments settings and results & discussion

Part 1 & Part 2 & Part 3 results

In kernel k-means, the hyper-parameters setting: theta1=1e-5, theta2=1e-5

(file: gif/kernel-k-means/report/{image_name}_{k}_{initialization_method}_{theta1}_{theta2}_{iteration_times}.gif)

- image1 with kernel k-means



original



method: random (k=2, iteration: 9)



method: k-means++ (k=2, iteration: 9)



method: random (k=3, iteration: 11)



method: k-means++ (k=3, iteration: 10)



method: random (k=4, iteration: 23)



method: k-means++ (k=4, iteration: 24)

- image2 with kernel k-means



original



method: random (k=2, iteration: 13)



method: k-means++ (k=2, iteration: 12)



method: random (k=3, iteration: 20)



method: k-means++ (k=3, iteration: 17)



method: random (k=4, iteration: 60)



method: k-means++ (k=4, iteration: 20)

In spectral clustering, the hyper-parameters setting: theta1=0.001, theta2=0.001

(file: gif/spectral-clustering/report/{image_name}_{cut_method}_{initialization_method}_{k}_{theta1}_{theta2}_{iteration_times}.gif)

- image1 with spectral clustering(normalized cut)



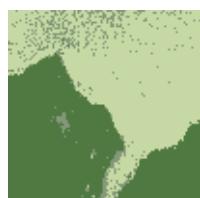
original



method: random (k=2, iteration: 6)



method: k-means++ (k=2, iteration: 5)



method: random (k=3, iteration: 16)



method: k-means++ (k=3, iteration: 8)



method: random (k=4, iteration: 15)



method: k-means++ (k=4, iteration: 11)

- image1 with spectral clustering(ratio cut)



original



method: random (k=2, iteration: 6)



method: k-means++ (k=2, iteration: 6)



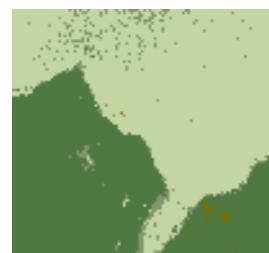
method: random (k=3, iteration: 6)



method: k-means++ (k=3, iteration: 6)

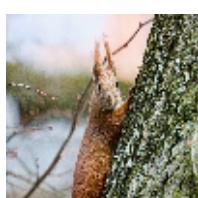


method: random (k=4, iteration: 28)



method: k-means++ (k=4, iteration: 16)

- image2 with spectral clustering(normalized cut)



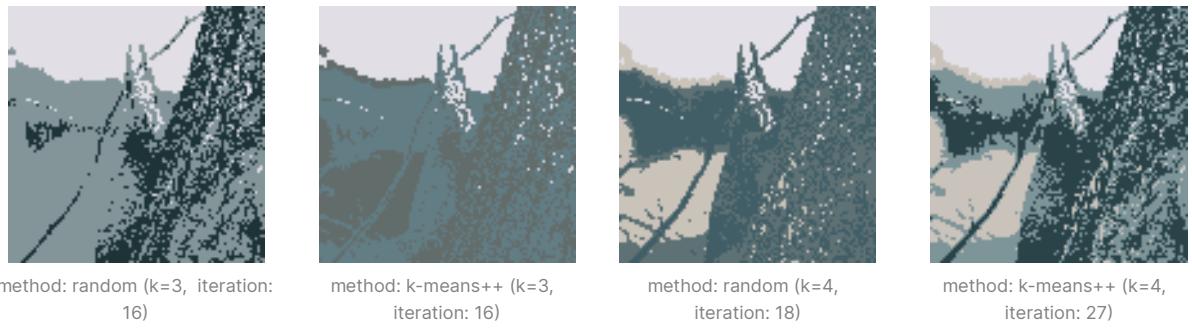
original



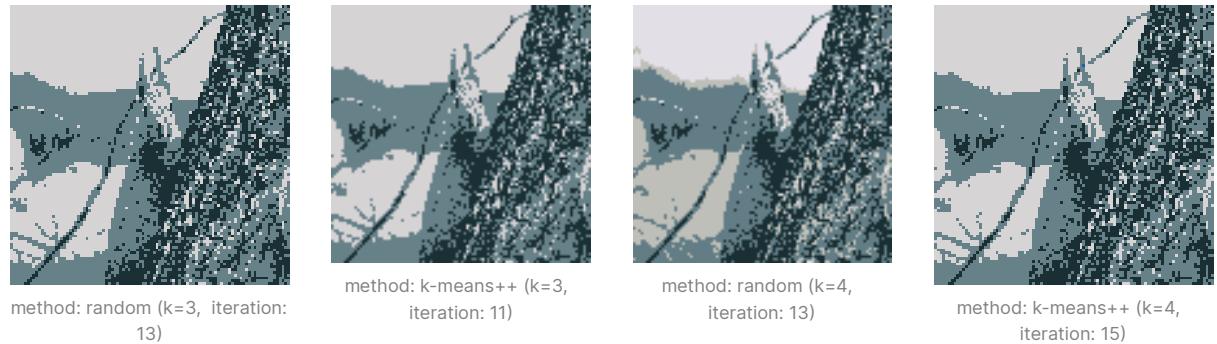
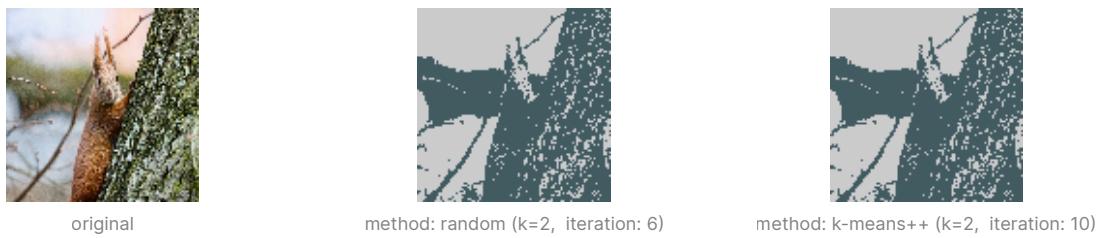
method: random (k=2, iteration: 17)



method: k-means++ (k=2, iteration: 6)



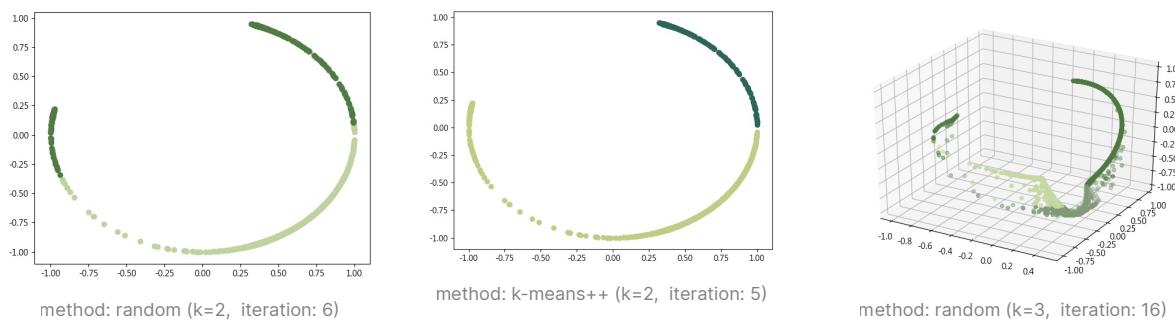
- image2 with spectral clustering(ratio cut)

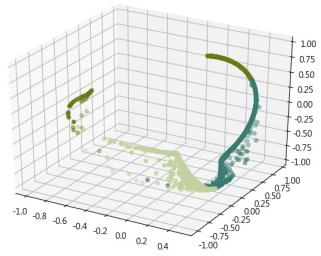


Part4 results

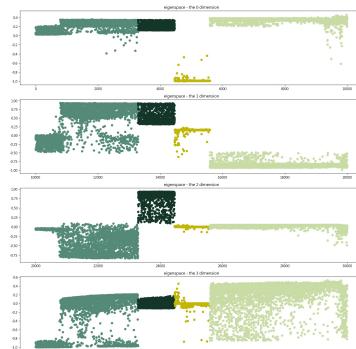
(file: gif/spectral-clustering/report/{image_name}_{cut_method}_{initialization_method}_{k}_{theta1}_{theta2}_{iteration_times}_eigen_space.gif)

- image1 with spectral clustering(normalized cut)

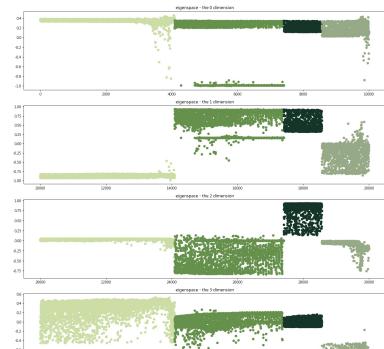




method: k-means++ (k=3, iteration: 8)

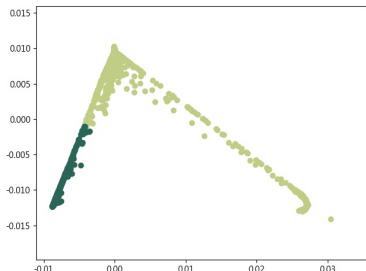


method: random (k=4, iteration: 15)

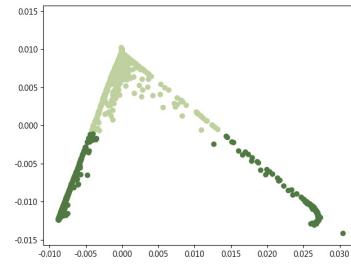


method: k-means++ (k=4, iteration: 11)

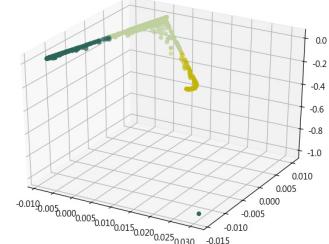
- image1 with spectral clustering(ratio cut)



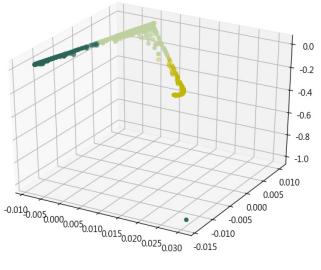
method: random (k=2, iteration: 6)



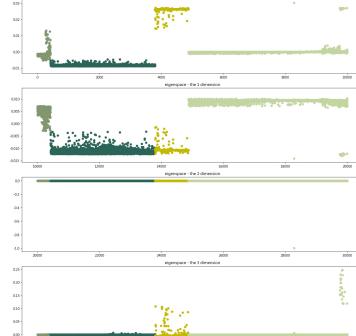
method: k-means++ (k=2, iteration: 6)



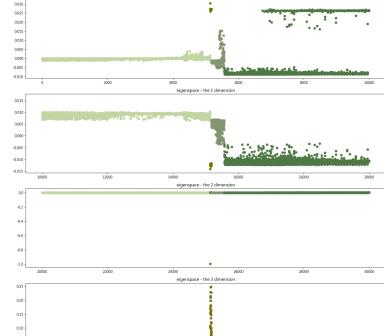
method: random (k=3, iteration: 6)



method: k-means++ (k=3, iteration: 6)

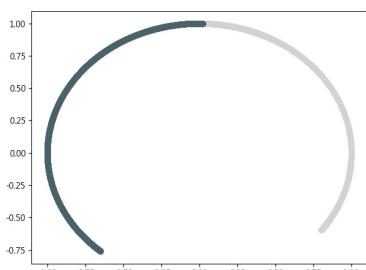


method: random (k=4, iteration: 28)

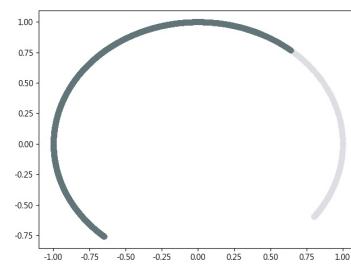


method: k-means++ (k=4, iteration: 16)

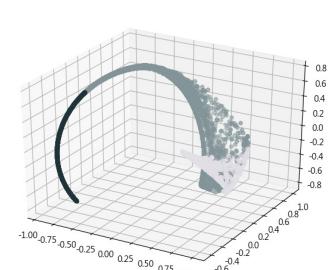
- image2 with spectral clustering(normalized cut)



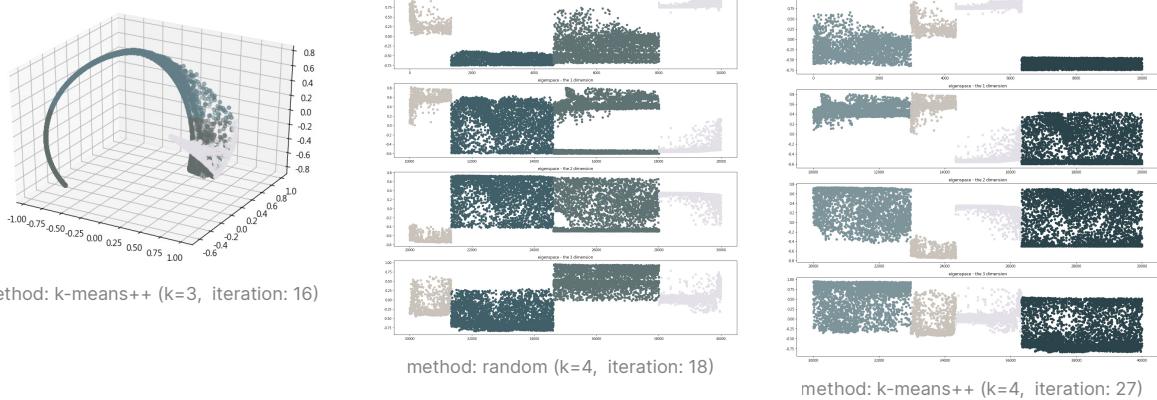
method: random (k=2, iteration: 17)



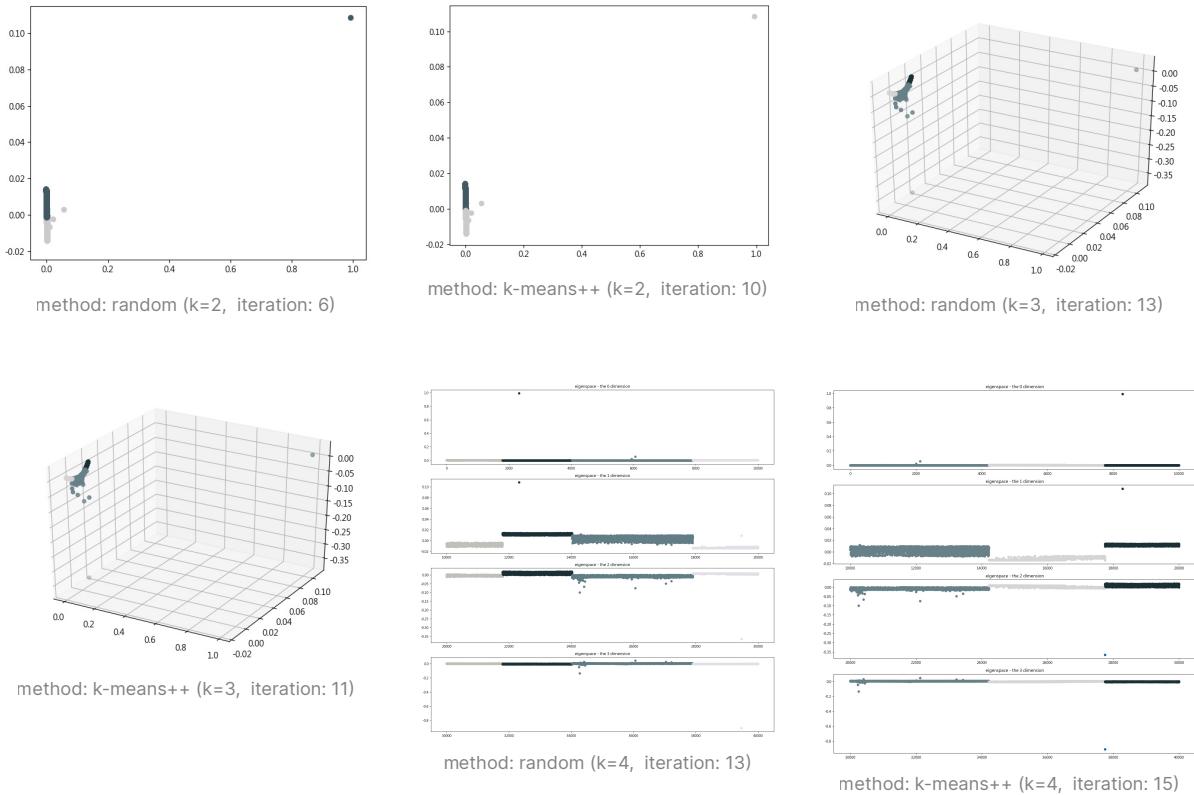
method: k-means++ (k=2, iteration: 6)



method: random (k=3, iteration: 16)

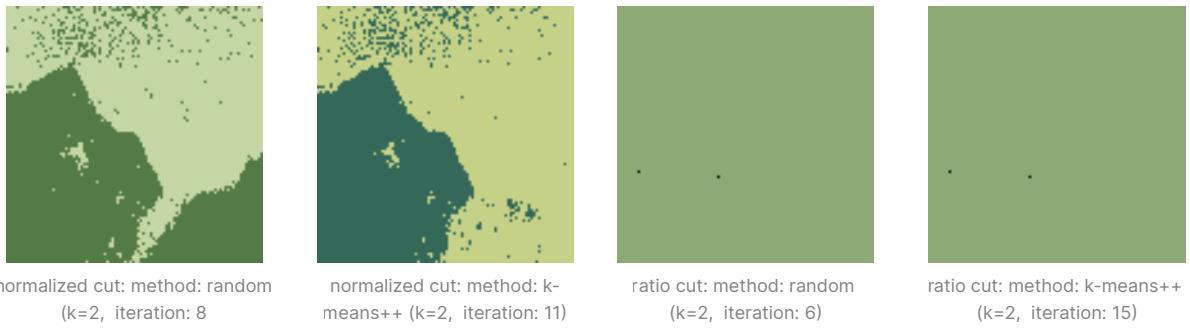


- image2 with spectral clustering(ratio cut)



Part 1 discussion

When $k=2$, we can see that the results in image2 of all three method(kernel k-means, spectral clustering(normalized cut), and spectral clustering(ratio cut)) are quite similar. But the results in image1, there are some different in three method. In kernel k-means, we can see that the green dots in the sea are successfully classified into the same cluster with island. But in spectral clustering (both normalized cut and ratio cut), we can see that the results are quite different. The island and the sea are classified into the same cluster. I think the reason is that the hyper-parameters are different in kernel k-means and spectral clustering. When the hyper-parameters are setting to the same(theta1=1e-5, theta2=1e-5), the results are showed below.

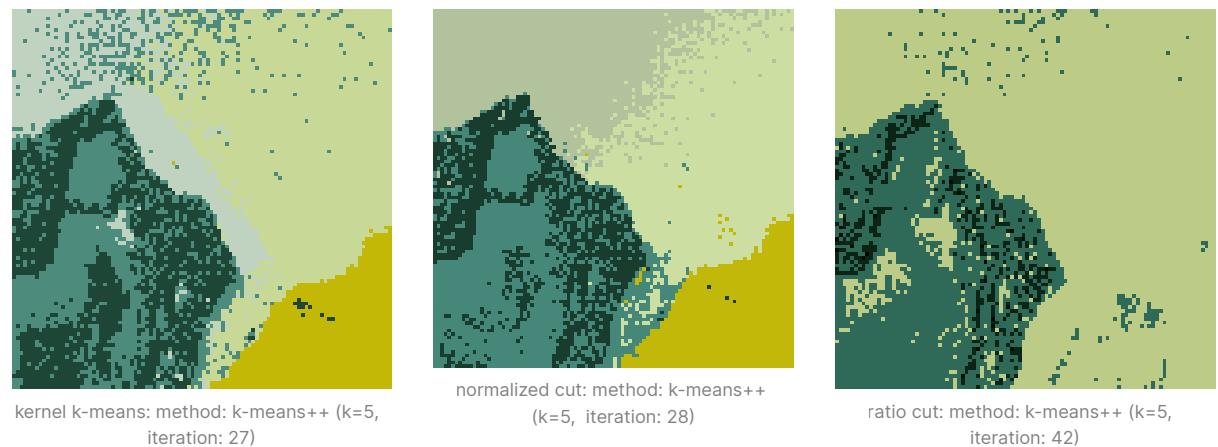


We can see that spectral clustering with normalized cut and with k-means++ initialization method can get the same clustering results with kernel k-means. But the result in spectral clustering with ratio cut are very strange, there are only two points clustered into the same cluster. So I think the hyper-parameters chosen are very important.

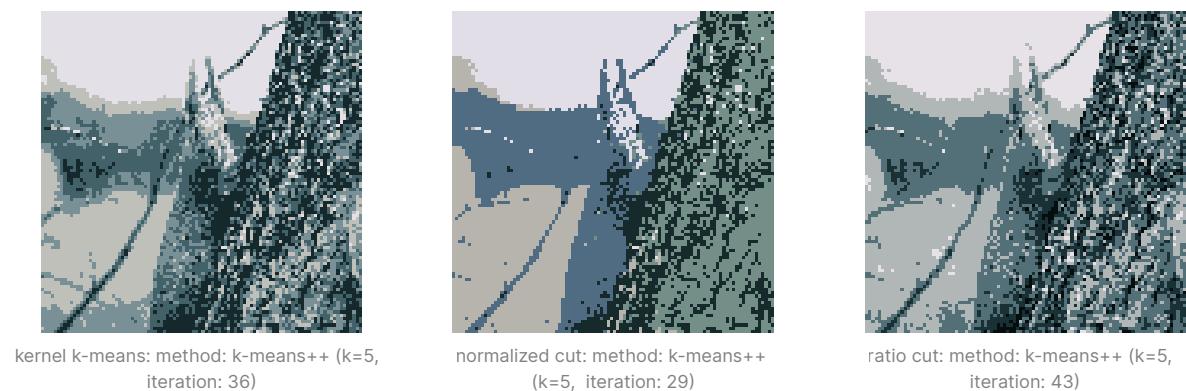
Part 2 discussion

When try more clusters, we need to more iterations to converge. We can see that when we use higher clusters, we can get more accurate clustering results. But it may overfit when k is too large. For example, when k=5, the results are below.

- image1



- image2



There are some strange in the results with k=5. For example, in kernel k-means with k=5, there is a weird cluster in the left corner. In conclusion, we can see that the clustering results with k=5 is not better than the clustering result with k=4.

Part 3 discussion

We know that the initialization of k-means clustering is highly dependent to the final result, so the initialization of center is important. In part3, I implemented two methods, `random` and `k-means++`. Because we want to get the same result in each time, we use `np.random.seed` to fix the random seed. Overall, `k-means++` converges faster than `random` in particular when k is higher. There is an example in image2 with kernel k-means, when k=4, `random` needs 60 iterations but `k-means++` only needs 20 iterations. I think the reason is that the initialization center in `k-means++` is much closer to final centers.

But sometimes `k-means++` does not perform well as we expected. An example in image2 with spectral clustering(normalized cut), when k=2, although the iteration times in `k-means++` is much smaller than `random`, the clustering result is `random` better than `k-means++`. I think the reason is that `k-means++` is trapped in the local minimum.

Part 4 discussion

Although we have same eigenvectors in the same setting(k=2, theta1=0.001, theta2=0.001), the clustering results are influenced by the k-means initialization method. So it would cause different results.

Also, we can see that there is often have extreme values in ratio cut, but normalized cut doe not have any extreme values. I think the reason is that we have already normalize the rows to norm 1. In conclusion, the normalize cut of results are more stable because it does not affected by extreme values.

c. observations and discussion

- In spectral clustering, we need to solve the eigenvalue problem. First of all, I use `np.linalg.eig` to get the eigenvalues and eigenvectors, but it is very slow. After some survey, `np.linalg.eigh` uses a faster algorithm when the matrix is symmetric. We know that our kernel function is symmetric, so I used `np.linalg.eigh` and save the eigenvalues and eigenvectors to speed up.
 - calculation time used `np.linalg.eig`: 527.00s
 - calculation time used `np.linalg.eigh`: 82.72s
- Hyper-parameters tuning: at first, I wrote a function `silhouette_score` to tune the hyper-parameters in kernels(theta1 and theta2). `silhouette_score` is one of the evaluation in unsupervised learning, the range from -1 to +1. The higher silhouette score means the better clustering model. The calculation and the code are shown below.

$$\begin{aligned}
 a(i) &= \frac{1}{|C_i| - 1} \sum_{j \in C_i, i \neq j} d(i, j) \\
 b(i) &= \min_{k \neq i} \frac{1}{|C_k|} \sum_{j \in C_k} d(i, j) \\
 s(i) &= \frac{b(i) - a(i)}{\max(a(i), b(i))} \text{ if } |C_i| > 1 \\
 \text{silhouette-score} &= \text{mean}(s)
 \end{aligned}$$

```

94     def silhouette_score(k, cluster_idx):
95         """
96             one of the evaluation in unsupervised learning, the range from -1 to +1.
97             The higher silhouette score means the better clustering model.
98             $a(i)=\frac{1}{|C_i|-1}\sum_{j \in C_i, i \neq j} d(i,j)$
99             $b(i)=\min_{k \neq i} \frac{1}{|C_k|}\sum_{j \in C_k} d(i,j)$
100            $s(i)=\frac{b(i)-a(i)}{\max(a(i), b(i))} \text{ if } |C_i| > 1$"
101            $silhouette-score = mean(s)$
102            :param k: the number of clusters
103            :param cluster_idx: the cluster of each pixel
104            :return: silhouette score
105        """
106        a = list()
107        tmp = list()
108        distance = np.triu(kernel_k_means_square_distance_point() ** 0.5)

```

```

109        for c_idx in range(k):
110            a.append(np.sum(distance[np.where(cluster_idx == c_idx)[0]]) /
111                      (len(np.where(cluster_idx == c_idx)[0]) - 1))
112            tmp.append(np.sum(distance[np.where(cluster_idx == c_idx)[0]]) /
113                      (len(np.where(cluster_idx == c_idx)[0])))
114        b = list()
115        for c_idx in range(k):
116            tmp2 = tmp.copy()
117            tmp2[c_idx] = np.inf
118            b.append(np.min(tmp2))
119        s = list()
120        for c_idx in range(k):
121            if len(np.where(cluster_idx == c_idx)[0]) > 1:
122                s.append((b[c_idx] - a[c_idx]) / max(a[c_idx], b[c_idx]))
123            else:
124                s.append(0)
125        return np.mean(s)

```

After tuning, we can see the best results below.

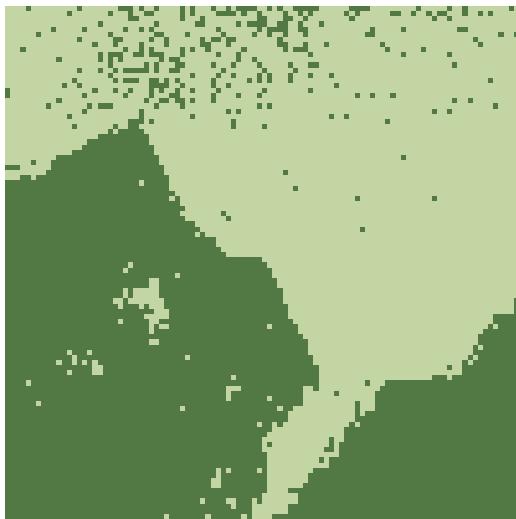


image1 with kernel k-means: method: k-means++ (k=2, iteration: 8, theta1: $2^{-19.5}$, theta2: $2^{-15.5}$)



image2 with kernel k-means: method: k-means++ (k=2, iteration: 19, theta1: 2^{-14} , theta2: 2^{-19})

- We can see that image1 results are clustered perfectly, but image2 results are very bad. I think the problem is in the `silhouette_score` calculation. I use kernel distance to replace $d(i, j)$ in the equation, but kernel distance are highly dependant to theta1 and theta2. When we change the theta1 and theta2, the kernel distance will be also changed. There is no same distance can be used, which means the criterion of each hyper-parameters are not the same. So finally, I abandon this method because I think it is not work in hyper-parameters tuning.
- In clustering, there is an important question: How to determine the optimal number of clusters? In our case, we have already showed all results of k=2~5. As part2 discussion, we know k=4 is the best choice in our case. But there are some popular techniques to find the best k, for example: elbow method, and Silhouette Curve.