

CEE 498 Applied Machine Learning - HW5

Rini Jasmine Gladstone (rjg7) and Maksymillian Podraza(podraza3)

April 17, 2020

1 Problem 1

1.1 Part a

We used the following image for our analysis.



We normalized the image so that the darkest pixel takes the value (0,0,0) and the lightest pixel takes the value (1,1,1).

```
1 image_paths = ["RobertMixed03.jpg", "smallstrelitzia.jpg", "  
    smallsunset.jpg"]  
2  
3 def load_images():  
4     image_array = []  
5     try:  
6         for image_path in image_paths:  
7             img = Image.open(image_path)  
8             image_array.append(img)  
9     except IOError:  
10        pass  
11    return image_array  
12
```

```

13 images = load_images()
14
15 working_image = images[0]
16
17 def image_matrix_create(image):
18     image_matrix = []
19     img_w, img_h = image.size
20     image_data = list(image.getdata())
21     for y in range(img_h):
22         image_matrix.append(image_data[y*img_w:(y+1)*img_w])
23     image_df = pd.DataFrame(image_matrix)
24     return image_df
25
26 def rgb_extract(image):
27     rgb = []
28     for i in range(3):
29         rgb.append(image.apply(lambda x: [y[i] for y in x]))
30     return rgb
31
32 def rgb_normalize(image):
33     normalized = []
34     for i in range(3):
35         x = image[i].values
36         min_max_scaler = preprocessing.MinMaxScaler()
37         x_scaled = min_max_scaler.fit_transform(x)
38         df = pd.DataFrame(x_scaled)
39         normalized.append(df)
40     return normalized
41
42 r = image_matrix_create(working_image)
43 rgb_matrix_list = rgb_extract(r)
44 rgb_normalized_list = rgb_normalize(rgb_matrix_list)

```

Then, we wrote the function for EM. We used K means to initialize the mean vector.

```

1 def EM_function(X,cov,k,niter):
2     weights = np.ones((k)) / k
3     kmeans = KMeans(n_clusters=k)
4     # fit kmeans object to data
5     kmeans.fit(X)
6     means = kmeans.cluster_centers_
7     eps=1e-8
8     likelihood = []
9     log_likelihoods = []
10    for step in range(niter):
11        likelihood = []
12        # Expectation step
13        for j in range(k):
14            likelihood.append(multivariate_normal.pdf(x=X, mean
15            =means[j], cov=cov[j]))
16            likelihood = np.array(likelihood)
17            assert likelihood.shape == (k, len(X))
18            b = []
19            # Maximization step
20            for j in range(k):
21                # use the current values for the parameters to
22                evaluate the posterior

```

```

21         # probabilities of the data have been generated by
each gaussian
22         b.append((likelihood[j] * weights[j]) / (np.sum([
likelihood[i] * weights[i] for i in range(k)], axis=0)+eps))
23         # update mean and variance
24         means[j] = np.sum(b[j].reshape(len(X),1) * X, axis
=0) / (np.sum(b[j]+eps))
25         # update the weights
26         weights[j] = np.mean(b[j])
27         assert cov.shape == (k, X.shape[1], X.shape[1])
28         assert means.shape == (k, X.shape[1])
29         log_likelihoods.append(np.log(np.sum([k*
multivariate_normal(means[i],cov[j]).pdf(X) for k,i,j in zip(
weights,range(len(means)),range(len(cov))]))))
30         if (step+1)%100==0 and (step+1)>=100:
31             print(step+1)
32         return means,weights,b,log_likelihoods,likelihood

```

We set the covariance matrix (cov) to be identity matrix and ran the following code setting k = 10, 20 and 50.

```

1 rvalues=rgb_normalized_list[0].values.flatten()
2 gvalues=rgb_normalized_list[1].values.flatten()
3 bvalues=rgb_normalized_list[2].values.flatten()
4 list_of_tuples = list(zip(rvalues, gvalues, bvalues))
5 data=pd.DataFrame(list_of_tuples, columns = ['R', 'G', 'B'])
6 #Number of clusters
7 k=50
8 #Initializing the covariance matrix
9 cov = []
10 for i in range(k):
11     cov.append(np.eye(data.shape[1]))
12 cov = np.array(cov)
13 means,weights,b,log_likelihoods,likelihood=EM_function(data,cov,k
,200)
14 #Plotting the log-likelihood
15 plt.plot(log_likelihoods)
16 likelihood = np.array(likelihood)
17 #Finding the closest cluster for each pixel
18 predictions = np.argmax(likelihood, axis=0)
19 data['cluster']=predictions
20 data['mean_R']=np.zeros(data.shape[0])
21 data['mean_G']=np.zeros(data.shape[0])
22 data['mean_B']=np.zeros(data.shape[0])
23 for i in range(data.shape[0]):
24     data.at[i,'mean_R']=means[data.at[i,'cluster']][0]
25     data.at[i,'mean_G']=means[data.at[i,'cluster']][1]
26     data.at[i,'mean_B']=means[data.at[i,'cluster']][2]
27 for i in range(k):
28     data['weight_c'+str(i)]=b[i]

```

We ran the code for 200 iterations.

The following are the images obtained by replacing each pixel with the mean of its cluster center.

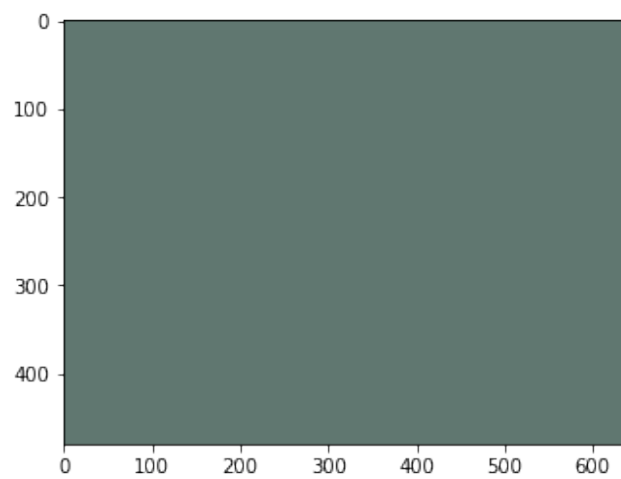


Figure 1: 10 Clusters

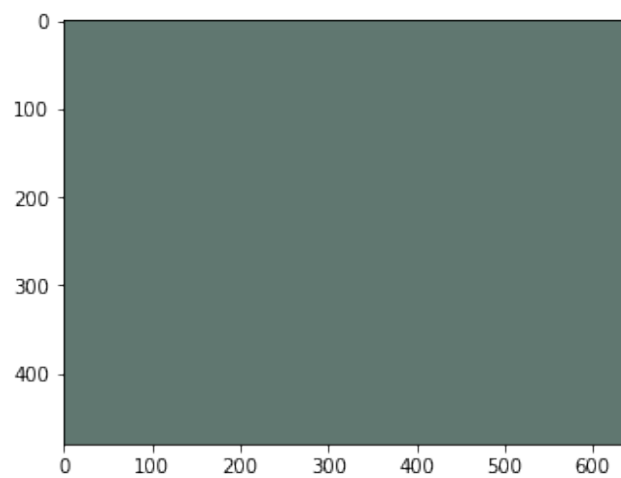


Figure 2: 20 Clusters

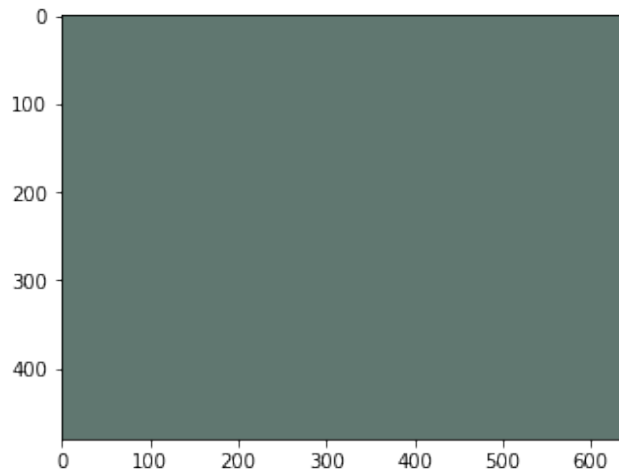


Figure 3: 50 Clusters

The code for generating the above images is as shown below.

```

1 def rescale_matrix(matrix, mode="int"):
2     new_matrix = matrix
3     for col in range(working_w):
4         new_matrix.loc[:, col] *= 255
5
6     if mode == "int":
7         new_matrix.astype(int)
8     return new_matrix
9
10
11 def print_matrix_image(original_image, image):
12     im2 = Image.new(original_image.mode, original_image.size)
13     img = image.values.flatten()
14     im2.putdata(img)
15     plt.imshow(im2)
16
17 def matricize_list(image, img_w=working_w, img_h=working_h):
18     image_matrix = []
19     for y in range(img_h):
20         image_matrix.append(image[y*img_w:(y+1)*img_w])
21     image_df = pd.DataFrame(image_matrix)
22     return image_df
23     return image_df
24
25
26 matrix_mean_R = matricize_list(list(data["mean_R"]))
27 matrix_mean_G = matricize_list(list(data["mean_G"]))
28 matrix_mean_B = matricize_list(list(data["mean_B"]))
29
30 matrix_mean_rescaled_R = rescale_matrix(matrix_mean_R)
31 matrix_mean_rescaled_G = rescale_matrix(matrix_mean_G)
32 matrix_mean_rescaled_B = rescale_matrix(matrix_mean_B)
33

```

```

34 matrix_mean_R = matrix_mean_R.astype(int)
35 matrix_mean_G = matrix_mean_G.astype(int)
36 matrix_mean_B = matrix_mean_B.astype(int)
37 matrix_mean_RGB = pd.DataFrame(np.rec.fromarrays((matrix_mean_R.
    values, matrix_mean_G.values, matrix_mean_B.values)).tolist(),
    columns=matrix_mean_R.columns,
    index=matrix_mean_R.index)
38
39
40
41 print_matrix_image(working_image, matrix_mean_RGB)

```

We observe from the figures that all the cluster centers have same mean, i.e. they drift together and hence the mean image is just a single color. Its the same for k=10,20 and 50

1.2 Part b

We constructed the figures showing the weights linking each pixel to each cluster center. Following is the code which executes this.

```

1 fig = plt.figure(figsize=(12, 12))
2 gs = gridspec.GridSpec(5, 2)
3 gs.update(wspace=0.25, hspace=0.25)
4
5 for i in range(10):
6     ax = plt.subplot(gs[i])
7     a = data["weight_c"+str(i)].values.reshape(working_h,working_w)
8     plt.axis('off')
9     plt.title("Cluster "+ str(i+1))
10    ax.set_aspect('equal')
11    plt.imshow(a)

```

The images generated are as shown below.

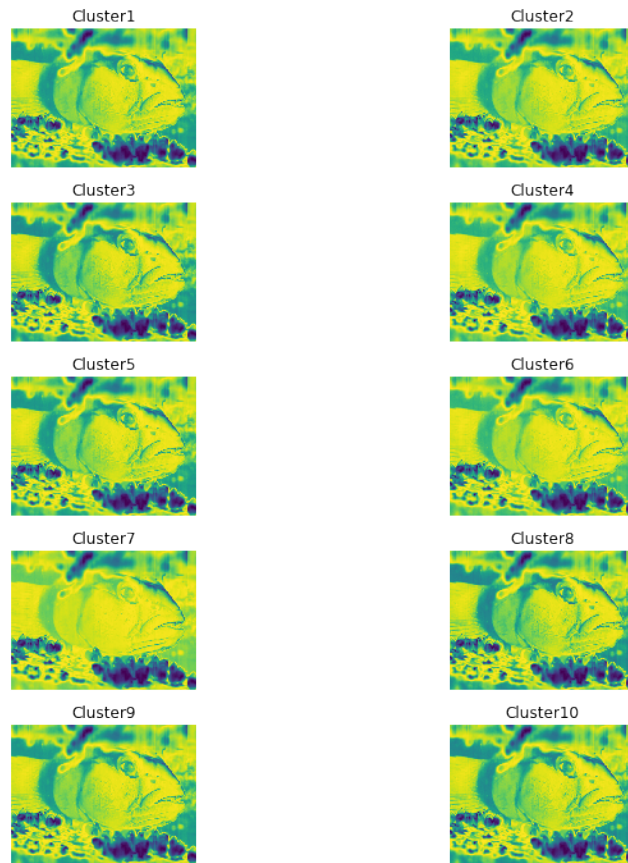


Figure 4: Cluster 1

We can observe that all the weight maps look similar, which is the expected behavior as at this scale, all the pixels are essentially one cluster and hence the weights are the same.

1.3 Part c

We changed the covariance to $0.1I$ as well as $0.0025I$ and ran the code.

```

1 rvalues=rgb_normalized_list[0].values.flatten()
2 gvalues=rgb_normalized_list[1].values.flatten()
3 bvalues=rgb_normalized_list[2].values.flatten()
4 list_of_tuples = list(zip(rvalues, gvalues, bvalues))
5 data=pd.DataFrame(list_of_tuples, columns = ['R', 'G', 'B'])
6 #Number of clusters
7 k=50
8 #Initializing the covariance matrix
9 cov = []
10 for i in range(k):

```

```

11     cov.append(0.1*np.eye(data.shape[1]))
12 cov = np.array(cov)
13 means, weights, b, log_likelihoods, likelihood=EM_function(data, cov, k
    ,200)
14 #Plotting the log-likelihood
15 plt.plot(log_likelihoods)
16 likelihood = np.array(likelihood)
17 #Finding the closest cluster for each pixel
18 predictions = np.argmax(likelihood, axis=0)
19 data['cluster']=predictions
20 data['mean_R']=np.zeros(data.shape[0])
21 data['mean_G']=np.zeros(data.shape[0])
22 data['mean_B']=np.zeros(data.shape[0])
23 for i in range(data.shape[0]):
24     data.at[i, 'mean_R']=means[data.at[i, 'cluster']][0]
25     data.at[i, 'mean_G']=means[data.at[i, 'cluster']][1]
26     data.at[i, 'mean_B']=means[data.at[i, 'cluster']][2]
27 for i in range(k):
28     data['weight_c'+str(i)]=b[i]

```

We ran the code for 200 iterations.

The following are the images obtained by replacing each pixel with the mean of its cluster center.

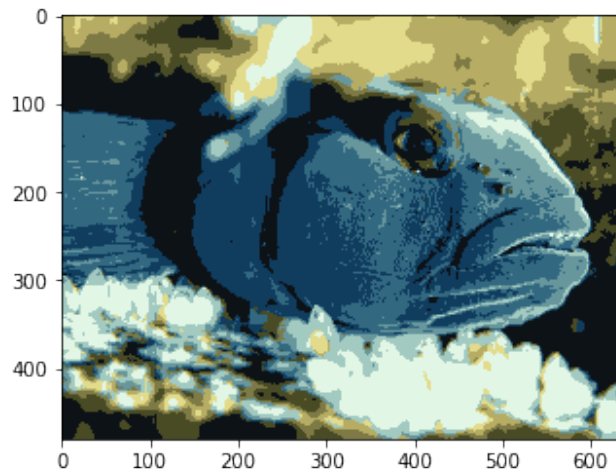


Figure 5: 10 Clusters (0.0025 I)

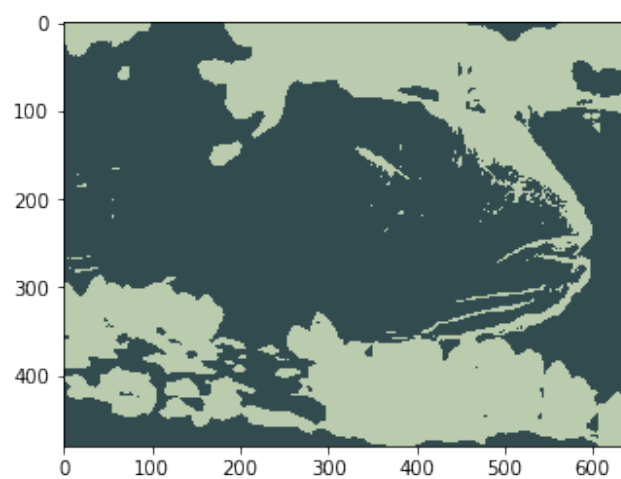


Figure 6: 10 Clusters (0.1 I)

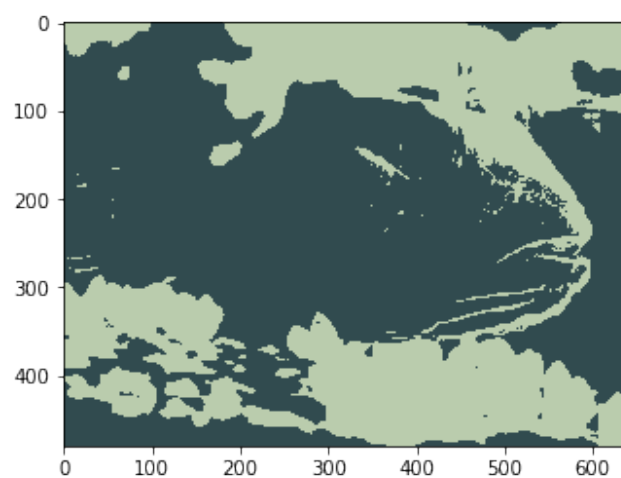


Figure 7: 20 Clusters (0.1 I)

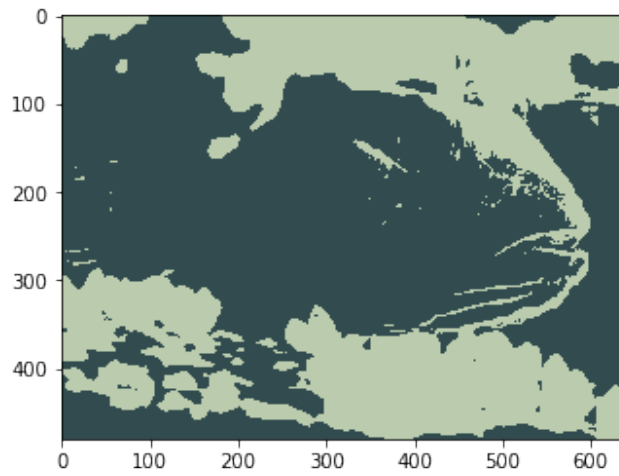


Figure 8: 50 Clusters (0.1 I)

Following are the new set of weight maps.

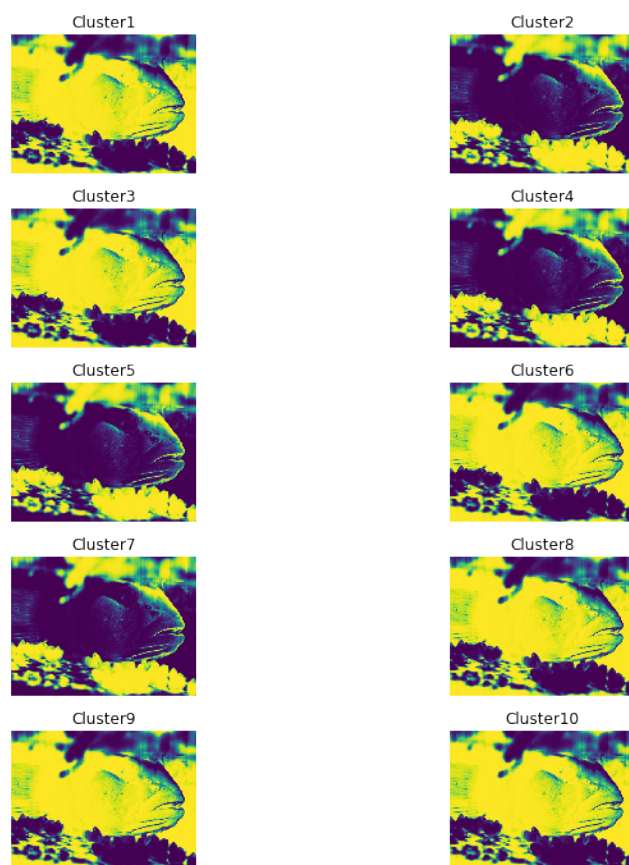


Figure 9: Covariance = $0.1 \mathbf{I}$

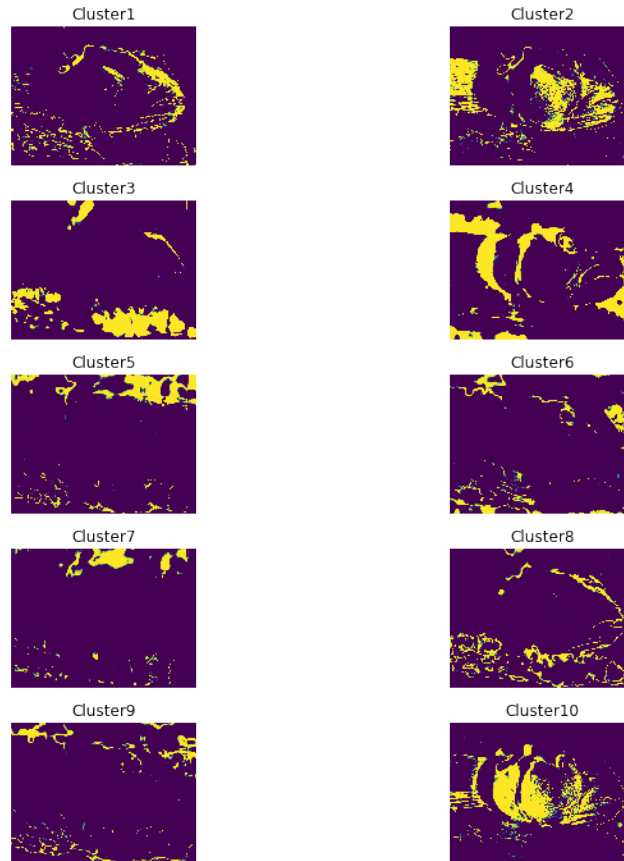


Figure 10: Covariance = 0.0025 I

We can observe that for covariance = 0.1I, the weights are not so different. But as we change the covariance to 0.0025I, we observe variation in weights image across the clusters, which means that the cluster centers do not drift together for covariance = 0.0025I. This also resulted in a better recovery of image by replacing the pixels with the mean value (Figure 5).

1.4 Part d

We estimated the covariance of pixel values by assuming that pixels are normally distributed and clustered them using EM. The code for this is as given below.

```

1 rvalues=rgb_normalized_list[0].values.flatten()
2 gvalues=rgb_normalized_list[1].values.flatten()
3 bvalues=rgb_normalized_list[2].values.flatten()
4 list_of_tuples = list(zip(rvalues, gvalues, bvalues))
5 data=pd.DataFrame(list_of_tuples, columns = ['R', 'G', 'B'])
6 #Number of clusters

```

```

7 k=50
8 #Initializing the covariance matrix
9 cov = []
10 for i in range(k):
11     cov.append(np.array(data.cov()))
12 cov = np.array(cov)
13 means, weights, b, log_likelihoods, likelihood=EM_function(data, cov, k
    ,200)
14 #Plotting the log-likelihood
15 plt.plot(log_likelihoods)
16 likelihood = np.array(likelihood)
17 #Finding the closest cluster for each pixel
18 predictions = np.argmax(likelihood, axis=0)
19 data['cluster']=predictions
20 data['mean_R']=np.zeros(data.shape[0])
21 data['mean_G']=np.zeros(data.shape[0])
22 data['mean_B']=np.zeros(data.shape[0])
23 for i in range(data.shape[0]):
24     data.at[i, 'mean_R']=means[data.at[i, 'cluster']][0]
25     data.at[i, 'mean_G']=means[data.at[i, 'cluster']][1]
26     data.at[i, 'mean_B']=means[data.at[i, 'cluster']][2]
27 for i in range(k):
28     data['weight_c'+str(i)]=b[i]

```

We ran the code for 200 iterations.

The following are the images obtained by replacing each pixel with the mean of its cluster center.

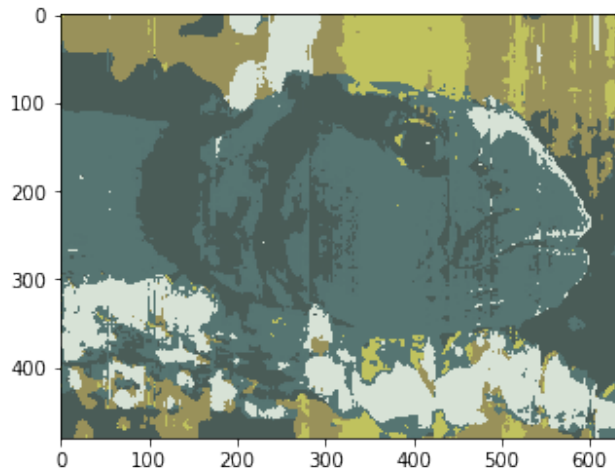


Figure 11: 10 Clusters

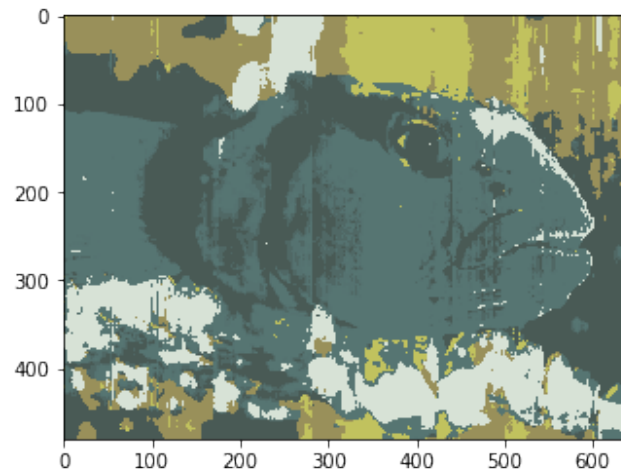


Figure 12: 20 Clusters

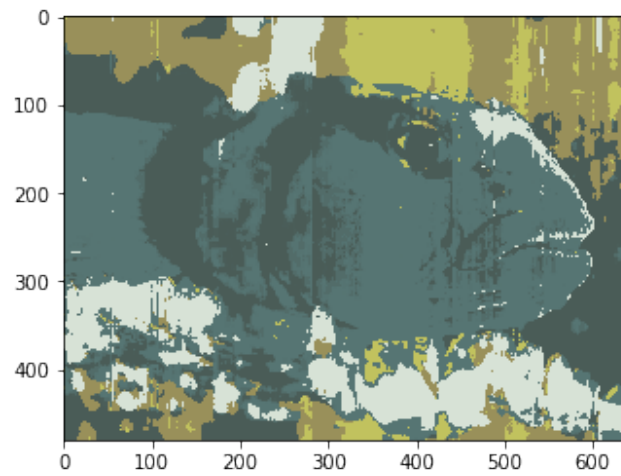


Figure 13: 50 Clusters

Following are the new set of weight maps.

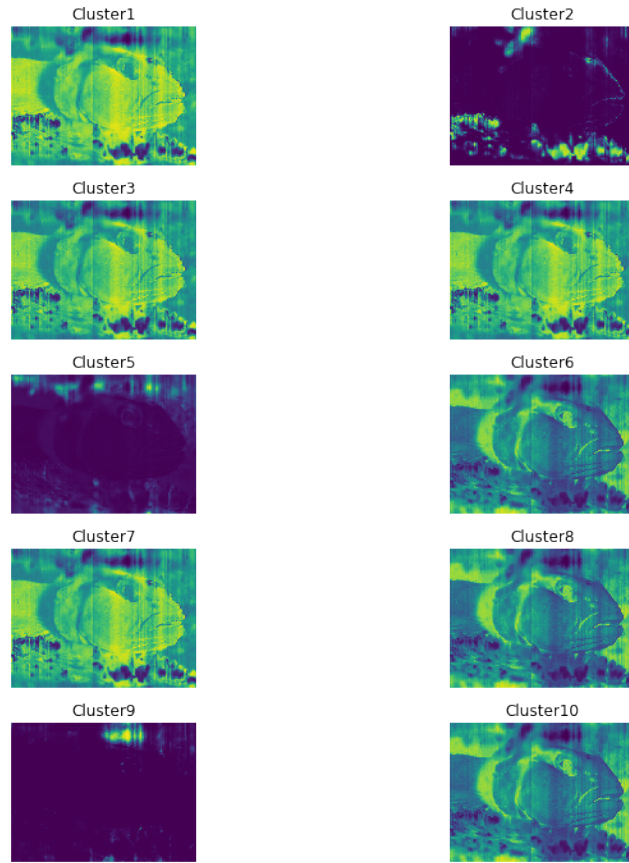


Figure 14:

We observe that the centers do not drift, and the maps are different from each other. The mean image also has many colors and could recover the original image to an extent. This is very different from part a and b, where the centers drift together and the mean image is a single color.