# Project 4 - Color Image Compression Using Unsupervised Learning (Clustering)

Adrian Cross

November 25, 2019

across19@vols.utk.edu
Machine learning
COSC 522

## 1 Abstract

In this project an image has been input and the RGB value for each image has been gained. Using various clustering algorithms the original possible RGB variations has been reduced to various different numbers by using different clusters amounts. The results of these different algorithms and cluster amounts have been compared by directly looking at the images and by computing percentage differences between original and final color values. As expected the quality of images decreased with cluster number. It was also found that the quantifiable percentage error rates were similar but there was a vast difference in the actual images.

## 2 Introduction

An image is made up of rows and columns with a pixel making up each point in the image. Each pixel is made up of a variety of factors, including a color value for red, green and blue. For a full color image each color value has a total possible number of $2^{24}$. On some computers the total number of colors are reduced to $2^8$ or 256 colors. By reducing the number of colors the storage size of the image is reduced. This number can be further reduced which, to $2^7, 2^6$, etc., reduces the storage size of the image further. In this project the image shown in figure 1 has been reduced to various different sizes via various different algorithms and the result of this reduction has been compared.



Figure 1: Original input image

The reduction in colors is done via clustering algorithms. These algorithms start with the clusters assigned in specific places and then move the clusters dependent on the input data. The RGB coordinates of these clusters is

one of the new color values which each data point can take. For example if a cluster point was at $(R, G, B) = (154, 124, 123)$ then the pixels with RGB coordinates closest to that cluster, and not closer to any other cluster, will be assigned a new RGB value of $(154, 124, 123)$. The output image has been compared both via by directly looking at the output images, with reduced colors, and by looking at the percentage difference between the old color value and the new color value, outlined in section 3.3.

# 3  Experiments and Results

## 3.1  Input

In section 2 the input image was introduced (figure 1). By organizing this image into rows and columns, with the corresponding number of pixels in the rows and columns, the image can be converted to a data frame of size $120 * 120 = 14400$ with each data point being one pixel with a color value for red, green and blue (RGB). By plotting these RGB values on a graph we can see the distribution of colors in this image, as shown in image 2 by the blue data points.
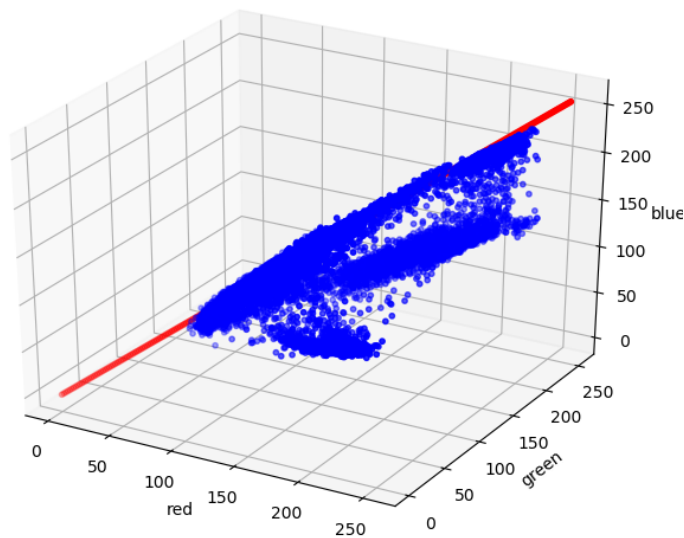


Figure 2: Pixel RGB values (blue) with 258 clusters distribution (red)

From this you can see the pixels have an RGB value between 0 and 256, concentrated around the center of the graph. This indicates the overall color of the image, for example if the image was mostly shades of red this image would show clustering which is biased towards the high red RGB value.

## 3.2  Initializing Clusters

After the RGB values have been inputted the initial cluster positions need to be defined. This can be done via a variety of means, such as random selection, but in this case the start positions have simply been defined as a straight line starting at $(R, G, B) = 0, 0, 0$ up to $(256, 256, 256)$ in even RGB steps [Steinbach et al., 2000]. The even steps are determined by the number of clusters used in each analysis. In this project a value of $256, 128, 64, 32, 16, 8, 4 and 2$ has been chosen for the cluster numbers. Therefore, for 256 clusters, the step inbetween each RGB value will simply be 1, as shown in figure 2 by the red colored line.

## 3.3  Data Comparison

In order to facilitate the comparison of the output image, two metrics have been used. The first is a simple visual test which is done by outputting the images with the new pixel color values. From this you can see how much the image changes with each compression amount. The second metric attempts to quantify the error between the source image and the output image. This is done by taking the input pixel RGB value and finding the percentage

difference to the new RGB value, as determined by the cluster center, as shown in equation 1.

$$\%diff = \frac{(RGB(actual) - RGB(Cluster))}{RGB(actual)} * 100 \tag{1}$$

This percentage difference is then plotted against number of clusters to see the difference in the trend for different number of clusters.

## 3.4  K-Means

The first algorithm tested is the k-means algorithm. This algorithm works by [University of Tennessee]:

1. Begin with an arbitrary assignment of samples to clusters or begin with an arbitrary set of cluster centers and assign samples to nearest clusters

2. Compute the sample mean of each cluster

3. Reassign each sample to the cluster with the nearest mean

4. If the classification of all samples has not changed, stop; else go to step 2.

   figure 3 shows how these cluster centers have moved from their original position in figure 2.
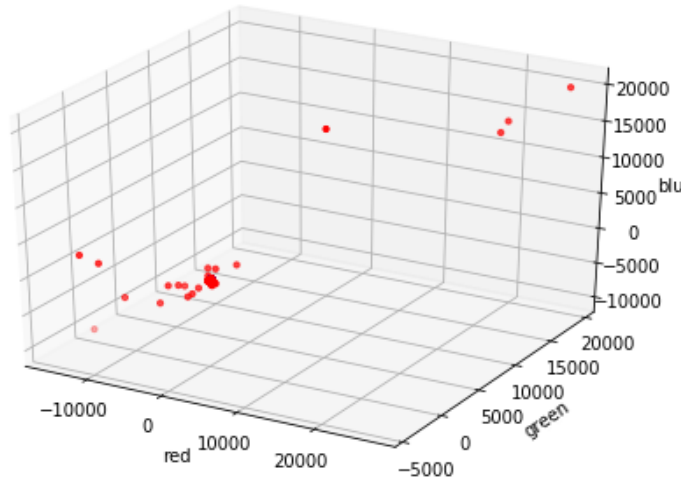


Figure 3: Final cluster position for k-means algorithm

By completing this algorithm over several iterations the final images are found, as shown in figure 9. By computing the percentage difference outlined in 3.3 the different percentage differences for each cluster number has been found and plotted in figure 4.
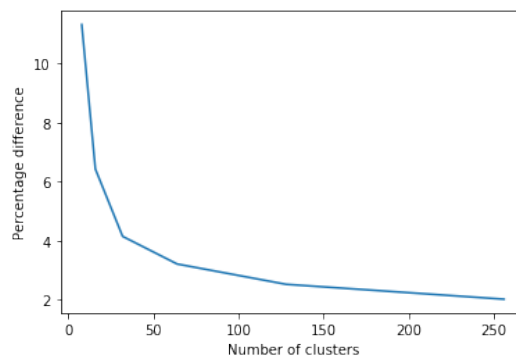


Figure 4: Quality reduction with number of clusters for k-means algorithm

3

(a) 256 clusters


(b) 8 clusters


(c) 4 clusters


(d) 2 clusters

Figure 5: Image quality reduction for k-means algorithm for different number of clusters

## 3.5 Winner Take All

The winner-take-all algorithm starts the same way as the k-means algorithm by first initializing the cluster centers. After this each pixel $x$ finds the nearest cluster center which is then defined as the winner. The winning cluster center $\omega_a$ is then modified using equation 2.

$$\omega_a^{new} = \omega_a^{old} + \epsilon(x - \omega_a^{old}) \tag{2}$$

Where $\epsilon$ is the learning parameter which is approximately equal to $\epsilon$ 0.01. By completing this algorithm over several iterations the final images are found, as shown in figure 9. By computing the percentage difference outlined in 6 the different percentage differences for each cluster number has been found and plotted in figure 6.
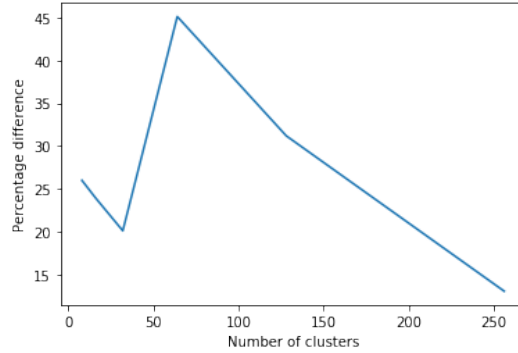


Figure 6: Quality reduction with number of clusters for winner take all algorithm
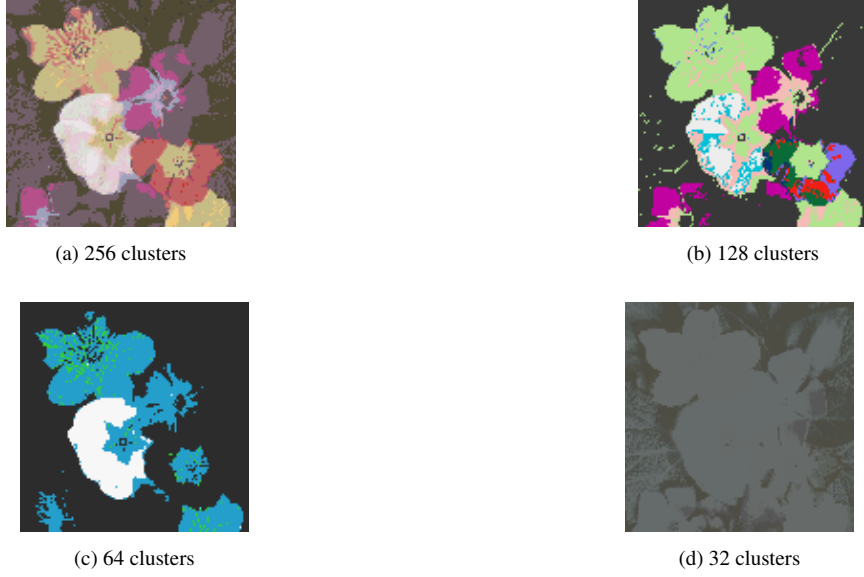
4

(a) 256 clusters


(b) 128 clusters


(c) 64 clusters


(d) 32 clusters

Figure 7: Image quality reduction for winner take all algorithm for different number of clusters

## 3.6 Kohonen

The Kohonen feature map (NN) is a modification of the winner take all algorithm. It works by assuming a problem dependent topological distance exists between a pair of cluster centers. This means that when the winning cluster center is updated so are its neighbours, dependent on this topological distance. Therefore the algorithm works the same way as winner take all, except equation 2 is modified so that it becomes equation 3.

$$\omega_r^{k+1} = \omega_r^k + \epsilon(k) * \Phi(k)(x - \omega_r^k) \tag{3}$$

Where $\omega_r$ are the cluster centers and epsilon is constantly adjusted so that as the number of iterations increases the learning rate decreases as it is becoming more stable, as shown in equation 4

$$\epsilon(k) = \epsilon_{max} * \left(\frac{\epsilon_{min}}{\epsilon_{max}}\right)^{k/k_{max}} \tag{4}$$

and the topological $\Phi$ distance is defined by equation 5

$$\Phi(k) = exp(-\frac{||g_{\omega_r} - frac\epsilon_{min}||^2}{2\sigma^2}) \tag{5}$$

Where $g_{\omega_r}$ is the coordinate of the cluster centers and $g_{\omega_{winner}}$ is the coordinate of the winner.
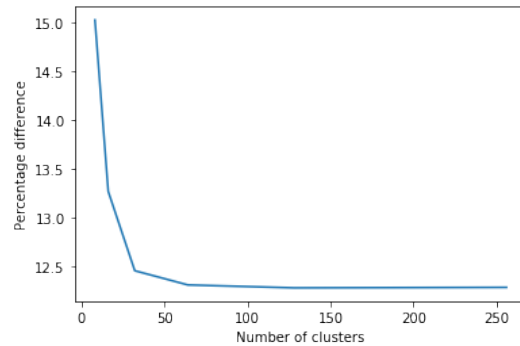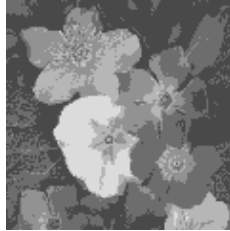


Figure 8: Quality reduction with number of clusters for Kohonen algorithm

5

(a) 64 clusters



(b) 32 clusters



(c) 8 clusters



(d) 2 clusters

Figure 9: Image quality reduction for Kohonen algorithm for different number of clusters

## 4  Conclusion

From the difference in images you can see a stark difference in the way the different algorithms alter the images. In the k-means algorithm the image quality is decreased and the final image is much different than for the winner take all algorithm, however there is still an overall decrease in quality.

While the images are graphically different, the actual reduction in quality, shown by the quantification of these errors using percentage differences, is very similar across each method. Therefore it can be concluded that while each algorithm does reduce both the storage size and percentage error of the image in similar ways, the actual visual quality is changed in different ways. Further study could be done to investigate this further, possibly by using image recognition tools, such as neural networks, in order to quantify the visual differences of the images between different number of clusters, as well as different algorithms.

## A  Code

### A.1  Algorithm File

```
1  #reads in image and outputs dataframe for clusters and pixels
2  import sys
3  import math
4  import timeit
5  import pandas as pd
6  import numpy as np
7  import matplotlib.pyplot as plt
8  from mpl_toolkits.mplot3d import Axes3D
9  from skimage import io
10
11 #custom modules
12 import find_closest_cluster
13 import plot_rgb
14 import k_means
15 import winner_take_all
16 import fig_io
17
18 def run_k_means(n_clusters,pixel_data):
19     count_max=10
20     clusters= fig_io.initialize_clusters(n_clusters)
21     clusters=k_means.initialize(clusters.copy())
22     pixel_data_comp=pixel_data
23     counter=0
24     limit_reached=False
```

```
25      while(limit_reached==False):
26          counter+=1
27          pixel_data_comp=pixel_data
28
29          #reset
30          clusters=k_means.reset_clusters(clusters.copy())
31          clusters.reset_index(drop=True,inplace=True)
32          pixel_data=find_closest_cluster.reset_closest_cluster(pixel_data.copy())
33          pixel_data.reset_index(drop=True,inplace=True)
34
35          #algorithm
36          pixel_data=find_closest_cluster.find_all_closest_cluster(pixel_data.copy(),clusters.
    copy())
37          clusters=k_means.sample_mean(pixel_data,clusters.copy())
38          clusters=k_means.reassign_cluster_pos(clusters.copy())
39
40          #stop conditions
41          if pixel_data_comp.equals(pixel_data)==True:
42              limit_reached=True
43              print('algorithm has converged stopping....')
44          print('counter='+str(counter))
45          if counter==count_max:
46              limit_reached=True
47              print('algorithm has reached '+str(counter)+' iterations stopping....')
48
49
50      print('completed for '+str(n_clusters)+' clusters')
51      pixel_data=find_closest_cluster.find_all_closest_cluster(pixel_data.copy(),clusters.copy()
    )
52      pixel_data=find_closest_cluster.update_pixels_to_closest_cluster(pixel_data.copy(),
    clusters.copy())
53
54      plot_rgb.plot_rgb_cluster_pix(clusters.copy(),pixel_data,'k_means_'+str(n_clusters))
55      plot_rgb.save_file(clusters.copy(),pixel_data.copy(),'k_means_'+str(n_clusters))
56
57
58
59
60  def run_winner_take_all(n_clusters,pixel_data,learning_parameter):
61      count_max=5
62      clusters=fig_io.initialize_clusters(n_clusters)
63      clusters=winner_take_all.initialize(clusters.copy())
64      pixel_data_comp=pixel_data
65      counter=0
66      limit_reached=False
67
68      while(limit_reached==False):
69          counter+=1
70          pixel_data_comp=pixel_data
71
72          #reset
73          clusters=winner_take_all.reset_clusters(clusters.copy())
74          clusters.reset_index(drop=True,inplace=True)
75          pixel_data=find_closest_cluster.reset_closest_cluster(pixel_data.copy())
76          pixel_data.reset_index(drop=True,inplace=True)
77
78          #algorithm
79          pixel_data=find_closest_cluster.find_all_closest_cluster(pixel_data.copy(),clusters.
    copy())
80          clusters=winner_take_all.pixel_shift(pixel_data.copy(),clusters,learning_parameter)
81
82          #stop conditions
83          if pixel_data_comp.equals(pixel_data)==True:
84              limit_reached=True
85              print('algorithm has converged stopping....')
86          print('counter='+str(counter))
87          if counter==count_max:
88              limit_reached=True
89              print('algorithm has reached '+str(counter)+' iterations stopping....')
90
91      print('completed for '+str(n_clusters)+' clusters')
92      pixel_data=find_closest_cluster.find_all_closest_cluster(pixel_data.copy(),clusters.copy()
    )
```

```
 93     pixel_data=find_closest_cluster.update_pixels_to_closest_cluster(pixel_data.copy(),
        clusters.copy())
 94
 95     plot_rgb.plot_rgb_cluster_pix(clusters.copy(),pixel_data,'wta_'+str(n_clusters))
 96     plot_rgb.save_file(clusters.copy(),pixel_data.copy(),'wta_'+str(n_clusters))
 97
 98
 99
100
101 def run_kohonen(n_clusters,pixel_data,eps_max,eps_min):
102     count_max=5
103     clusters=fig_io.initialize_clusters(n_clusters)
104     clusters=winner_take_all.initialize(clusters.copy())
105     pixel_data_comp=pixel_data
106     counter=0
107     limit_reached=False
108
109     while(limit_reached==False):
110         counter+=1
111         pixel_data_comp=pixel_data
112
113         #reset
114         clusters=winner_take_all.reset_clusters(clusters.copy())
115         clusters.reset_index(drop=True,inplace=True)
116         pixel_data=find_closest_cluster.reset_closest_cluster(pixel_data.copy())
117         pixel_data.reset_index(drop=True,inplace=True)
118
119         #algorithm
120         pixel_data=find_closest_cluster.find_all_closest_cluster(pixel_data.copy(),clusters.
        copy())
121         winner_take_all.kohonen(pixel_data=pixel_data.copy(),clusters=clusters.copy(),
        epsilon_max=eps_max,epsilon_min=eps_min,k=counter,k_max=count_max)
122
123         #stop conditions
124         print('counter='+str(counter))
125         if counter==count_max:
126             limit_reached=True
127             print('algorithm has reached '+str(counter)+' iterations stopping....')
128
129     print('completed for '+str(n_clusters)+' clusters')
130     pixel_data=find_closest_cluster.find_all_closest_cluster(pixel_data.copy(),clusters.copy()
        )
131     pixel_data=find_closest_cluster.update_pixels_to_closest_cluster(pixel_data.copy(),
        clusters.copy())
132
133     plot_rgb.plot_rgb_cluster_pix(clusters.copy(),pixel_data,'kohonen_'+str(n_clusters))
134     plot_rgb.save_file(clusters.copy(),pixel_data.copy(),'kohonen_'+str(n_clusters))
135
136
137
138 pixel_data=fig_io.read_in_fig('flowersm.ppm')
139 for n_clusters in [256,128,64,32,16,8,4,2]:
140     # run_k_means(n_clusters,pixel_data)
141     #run_winner_take_all(n_clusters,pixel_data,0.05)
142     run_kohonen(n_clusters,pixel_data,2,0.5)
```

## A.2   Image Output

```
 1 #output_images to images folder
 2 import sys
 3 import math
 4 import timeit
 5 import pandas as pd
 6 import numpy as np
 7 import matplotlib.pyplot as plt
 8 import matplotlib.image as mpimg
 9 import scipy.misc as smp
10 from PIL import Image
11 from mpl_toolkits.mplot3d import Axes3D
12 from skimage import io
13
```

```
14  #custom modules
15  import find_closest_cluster
16  import plot_rgb
17  import k_means
18  import fig_io
19
20
21  for i in [256,128,64,32,8,4,2]:
22      pixel_data = pd.read_csv("data/k_means_"+str(i)+"_pixel_data.csv") #k_mean
23      fig_io.read_out_fig(pixel_data.copy(),'flowersm.ppm','k_means_'+str(i)+'clusters')
24
25      pixel_data = pd.read_csv("data/wta_"+str(i)+"_pixel_data.csv") #wta
26      fig_io.read_out_fig(pixel_data.copy(),'flowersm.ppm','wta_'+str(i)+'clusters')
27
28      pixel_data = pd.read_csv("data/kohonen_"+str(i)+"_pixel_data.csv") #kohonen
29      fig_io.read_out_fig(pixel_data.copy(),'flowersm.ppm','kohonen_'+str(i)+'clusters')
```

## A.3   Evaluate Algorithms

```
1   #evaluates data
2   import sys
3   import math
4   import timeit
5   import pandas as pd
6   import numpy as np
7   import matplotlib.pyplot as plt
8   import matplotlib.image as mpimg
9   import scipy.misc as smp
10  from PIL import Image
11  from mpl_toolkits.mplot3d import Axes3D
12  from skimage import io
13
14  #custom modules
15  import find_closest_cluster
16  import plot_rgb
17  import k_means
18  import fig_io
19
20  def r_perc_diff(pixel_data): #returns the mean percentage difference
21      perc_diff=(pixel_data['red']-pixel_data['r_cluster'])/pixel_data['red']
22      perc_diff=abs(perc_diff*100)
23      perc_diff=perc_diff.mean()
24      return perc_diff
25
26  def g_perc_diff(pixel_data): #returns the mean percentage difference
27      perc_diff=(pixel_data['green']-pixel_data['g_cluster'])/pixel_data['green']
28      perc_diff=abs(perc_diff*100)
29      perc_diff=perc_diff.mean()
30      return perc_diff
31
32  def b_perc_diff(pixel_data): #returns the mean percentage difference
33      perc_diff=(pixel_data['blue']-pixel_data['b_cluster'])/pixel_data['blue']
34      perc_diff=abs(perc_diff*100)
35      perc_diff=perc_diff.mean()
36      return perc_diff
37
38  def output_graphs(algorithm):
39      cluster_no=[256,128,64,32,16,8]
40      perc_diff_r=[]
41      perc_diff_g=[]
42      perc_diff_b=[]
43
44      for i in cluster_no:
45          pixel_data = pd.read_csv("data/"+algorithm+"_"+str(i)+"_pixel_data.csv")
46          perc_diff_r.append(r_perc_diff(pixel_data))
47          perc_diff_g.append(g_perc_diff(pixel_data))
48          perc_diff_b.append(b_perc_diff(pixel_data))
49      perc_diff=[(perc_diff_r+perc_diff_g+perc_diff_b)/3 for perc_diff_r,perc_diff_g,perc_diff_b
          in zip(perc_diff_r,perc_diff_g,perc_diff_b)]
50      plot_rgb.plot_perc_diff(cluster_no,perc_diff_r,'red',algorithm)
51      plot_rgb.plot_perc_diff(cluster_no,perc_diff_g,'green',algorithm)
```

```
52        plot_rgb.plot_perc_diff(cluster_no,perc_diff_b,'blue',algorithm)
53        plot_rgb.plot_perc_diff(cluster_no,perc_diff,'all',algorithm)
54
55
56   output_graphs('k_means')
57   output_graphs('wta')
58   output_graphs('kohonen')
```

## A.4  k-means Algorithm File

```
1    import numpy as np
2
3    def initialize(clusters): #initialize the clusters with empty totals
4        clusters['red_total']=np.zeros(shape=len(clusters),dtype=int)
5        clusters['green_total']=np.zeros(shape=len(clusters),dtype=int)
6        clusters['blue_total']=np.zeros(shape=len(clusters),dtype=int)
7        clusters['n_total']=np.zeros(shape=len(clusters),dtype=int)
8        clusters['red_mean']=np.nan
9        clusters['green_mean']=np.nan
10       clusters['blue_mean']=np.nan
11       return clusters
12
13   def reset_clusters(clusters): #reset clusters to 0 total
14       clusters['red_total']=0
15       clusters['green_total']=0
16       clusters['blue_total']=0
17       clusters['n_total']=0
18       clusters['red_mean']=0
19       clusters['green_mean']=0
20       clusters['blue_mean']=0
21       return clusters
22
23   def sample_mean(pixel_data,clusters): #Finds the rgb sample mean
24       for i in range(len(pixel_data)):
25           cid=int(pixel_data.closest_id.iloc[i])
26           clusters.red_total.iloc[cid]+=pixel_data.red.iloc[i]
27           clusters.green_total.iloc[cid]+=pixel_data.green.iloc[i]
28           clusters.blue_total.iloc[cid]+=pixel_data.blue.iloc[i]
29           clusters.n_total.iloc[cid]+=1
30
31       clusters['red_mean']=clusters['red_total']/clusters['n_total']
32       clusters['green_mean']=clusters['green_total']/clusters['n_total']
33       clusters['blue_mean']=clusters['blue_total']/clusters['n_total']
34       return clusters
35
36   def eliminate_empty_clusters(clusters): #removes the clusters with no local pixel_data points
37       clusters = clusters[clusters.n_total != 0]
38       return clusters
39
40   def reassign_cluster_pos(clusters): #reassigns the position of clusters to their means
41       clusters['red']=clusters['red_mean']
42       clusters['green']=clusters['green_mean']
43       clusters['blue']=clusters['blue_mean']
44       return clusters
45
46   def check_clusters_size_limit(clusters,limit): #checks the number of clusters left, returns
         true if required cluster size is reached
47       if len(clusters)<=limit:
48           return True
49       else:
50           return False
```

## A.5  WTA and Kohonen Algorithm File

```
1    import numpy as np
2
3    def initialize(clusters): #initialize the clusters with empty totals
4        clusters['red_temp']=np.zeros(shape=len(clusters),dtype=int)
5        clusters['green_temp']=np.zeros(shape=len(clusters),dtype=int)
6        clusters['blue_temp']=np.zeros(shape=len(clusters),dtype=int)
```

```python
 7      clusters['phi']=np.zeros(shape=len(clusters))
 8      return clusters
 9
10 def reset_clusters(clusters): #reset clusters to 0 total
11      clusters['red_temp']=0
12      clusters['green_temp']=0
13      clusters['blue_temp']=0
14      return clusters
15
16 def pixel_shift(pixel_data,clusters,epsilon): #shifts the cluster position
17      for i in range(len(pixel_data)):
18          red_shift=epsilon*(pixel_data['red']-clusters['red'])
19          green_shift=epsilon*(pixel_data['green']-clusters['green'])
20          blue_shift=epsilon*(pixel_data['blue']-clusters['blue'])
21          clusters['red']=clusters['red']+red_shift
22          clusters['green']=clusters['green']+green_shift
23          clusters['blue']=clusters['blue']+blue_shift
24      return clusters
25
26
27 def kohonen(pixel_data,clusters,epsilon_max,epsilon_min,k,k_max):
28      epsilon=epsilon_max*(epsilon_min/epsilon_max)**(k/k_max)
29
30      for i in range(len(pixel_data)):
31          phi=0
32          std_dev=0
33          cid=pixel_data.closest_id.iloc[i]
34
35          for j in range (0,len(clusters)):
36              std_dev_r=((clusters['red']-clusters.red.iloc[j])**2).sum()/(len(clusters)-1)
37              std_dev_g=((clusters['green']-clusters.red.iloc[j])**2).sum()/(len(clusters)-1)
38              std_dev_b=((clusters['blue']-clusters.red.iloc[j])**2).sum()/(len(clusters)-1)
39
40              a=(clusters.red.iloc[cid]-clusters.red.iloc[j])**2+(clusters.green.iloc[cid]-
    clusters.green.iloc[j])**2+(clusters.blue.iloc[cid]-clusters.blue.iloc[j])**2
41              std_dev=(std_dev_r+std_dev_g+std_dev_b)/3
42              phi=np.exp(a/(2*std_dev)) #calculated phi value for shift
43              clusters.phi.iloc[j]=phi
44
45          red_shift=clusters['phi']*epsilon*(pixel_data['red']-clusters['red'])
46          green_shift=clusters['phi']*epsilon*(pixel_data['green']-clusters['green'])
47          blue_shift=clusters['phi']*epsilon*(pixel_data['blue']-clusters['blue'])
48          clusters['red']=clusters['red']+red_shift
49          clusters['green']=clusters['green']+green_shift
50          clusters['blue']=clusters['blue']+blue_shift
51      return clusters
52
53
54
55
56
57
58 def kohonen1234(pixel_data,clusters,epsilon_max,epsilon_min,k,k_max): #shifts the cluster
       position
59      epsilon=epsilon_max*(epsilon_min/epsilon_max)**(k/k_max)
60      red_shift=(pixel_data['red']-clusters['red'])
61      green_shift=(pixel_data['green']-clusters['green'])
62      blue_shift=(pixel_data['blue']-clusters['blue'])
63
64      for j in range (0,len(clusters)):
65          std_dev_r=((clusters['red']-clusters.red.iloc[j])**2).sum()/(len(clusters)-1)
66          std_dev_g=((clusters['green']-clusters.red.iloc[j])**2).sum()/(len(clusters)-1)
67          std_dev_b=((clusters['blue']-clusters.red.iloc[j])**2).sum()/(len(clusters)-1)
68          std_dev=(std_dev_r+std_dev_g+std_dev_b)/3
69          for i in range(0,len(clusters)):
70              print(clusters)
71              a=(clusters.red.iloc[i]-clusters.red.iloc[j])**2+(clusters.green.iloc[i]-clusters.
    green.iloc[j])**2+(clusters.blue.iloc[i]-clusters.blue.iloc[j])**2
72              print('red',clusters.red.iloc[i]-clusters.red.iloc[j])
73              print('green',clusters.green.iloc[i]-clusters.green.iloc[j])
74              print('blue',clusters.blue.iloc[i]-clusters.blue.iloc[j])
75
76              print(a)
```

```
77          print(std_dev)
78
79          phi=np.exp(a/(2*std_dev))
80          clusters.red_temp.iloc[i]=red_shift
81          clusters.red_temp.iloc[i]=phi
82          clusters.red_temp.iloc[i]=epsilon
83          clusters.red_temp.iloc[i]=clusters.red.iloc[i]+red_shiftclusters.red.iloc[i]*phi*
    epsilon
84          clusters.green_temp.iloc[i]=clusters.green.iloc[i]+green_shift*phi*epsilon
85          clusters.blue_temp.iloc[i]=clusters.blue.iloc[i]+blue_shift*phi*epsilon
86
87      clusters['red']=clusters['red_temp']
88      clusters['green']=clusters['green_temp']
89      clusters['blue']=clusters['blue_temp']
90      return clusters
```

## A.6  Input/Output Image File

```
1  #reads in picture and outputs rgb and clusters with closest cluster for each data point
2  #outputs clusters.csv for cluster locations
3  #outputs pixel_data.csv for each data point with rgb and closest cluster and closest cluster
      dist
4
5  import sys
6  import math
7  import timeit
8  import pandas as pd
9  import numpy as np
10 import matplotlib.pyplot as plt
11 from sklearn.svm import SVC
12 from sklearn import svm
13 from skimage import io
14 from PIL import Image
15
16 def read_in_fig(figure):
17     image = io.imread(figure)
18     io.imshow(image)
19     io.show()
20     rows = image.shape[0]
21     cols = image.shape[1]
22
23     #reads in rgb list from image
24     r=[]#length of each should be 120*120
25     g=[]
26     b=[]
27     c=[]
28     cdist=[]
29     for line in image:
30         for pixel in line:
31             temp_r,temp_g,temp_b=pixel
32             r.append(temp_r)
33             g.append(temp_g)
34             b.append(temp_b)
35             c.append(np.NaN)
36             cdist.append(np.NaN)
37     data = pd.DataFrame({'red': r,'green': g,'blue': b,'closest_id':c,'closest_dist':cdist})
38     return data
39
40 def read_out_fig(pixel_data,image_input_name,name):#reads out image to same size as read in
      image in image_name
41     image = io.imread(image_input_name)
42     rows = image.shape[0]
43     cols = image.shape[1]
44     data = np.zeros( (rows,cols,3), dtype=np.uint8 )
45     count=0
46
47     for row in range(0,rows):
48         for column in range(0,cols):
49             data[row,column] = [pixel_data.r_cluster.iloc[count],pixel_data.g_cluster.iloc[
    count],pixel_data.b_cluster.iloc[count]]  #inputs pixel color for specific pixel
50             count+=1
```

```
51      img = Image.fromarray( data )
52      img.save('/home/across/UTK_PhD/Machine_learning_fall_2019/project_4/images/'+name+'.png')
53
54
55  def initialize_clusters(n): #initialize clusters
56      cluster = pd.DataFrame(index=range(0,n),columns=['red','green','blue'])
57      cluster['red']=np.linspace(0,256,n)
58      cluster['green']=np.linspace(0,256,n)
59      cluster['blue']=np.linspace(0,256,n)
60      return cluster
```

## A.7 Output Data File

```
1   import matplotlib.pyplot as plt
2   from mpl_toolkits.mplot3d import Axes3D
3
4   def plot_rgb_cluster_pix(cluster,pixel,name): #plots rgb for cluster and pixels
5       cluster_fig = plt.figure()
6       ax=Axes3D(cluster_fig)
7       ax.scatter(cluster['red'],cluster['green'],cluster['blue'],color='red',s=10)
8       ax.scatter(pixel['red'],pixel['green'],pixel['blue'],color='blue',s=10)
9       ax.set_xlabel('red')
10      ax.set_ylabel('green')
11      ax.set_zlabel('blue')
12      plt.savefig('/home/across/UTK_PhD/Machine_learning_fall_2019/project_4/plots/'+name+'.png'
        )
13      print('rgb plot saved with cluster and pixels')
14
15  def save_file(cluster,pixel_data,name):
16      cluster.to_csv('/home/across/UTK_PhD/Machine_learning_fall_2019/project_4/data/'+name+'
        _clusters.csv',index=False)
17      pixel_data.to_csv('/home/across/UTK_PhD/Machine_learning_fall_2019/project_4/data/'+name+'
        _pixel_data.csv',index=False)
18
19  def plot_rgb_cluster(cluster,pixel,name): #plots rgb for cluster and pixels
20      cluster_fig = plt.figure()
21      ax=Axes3D(cluster_fig)
22      ax.scatter(cluster['red'],cluster['green'],cluster['blue'],color='red',s=10)
23      ax.set_xlabel('red')
24      ax.set_ylabel('green')
25      ax.set_zlabel('blue')
26      plt.savefig('/home/across/UTK_PhD/Machine_learning_fall_2019/project_4/plots/'+name+'
        cluster.png')
27      print('rgb plot saved with clusters')
28
29  def plot_perc_diff(cluster_no,perc_diff,color_name,algorithm):
30      plt.figure()
31      plt.plot(cluster_no,perc_diff)
32      plt.xlabel('Number of clusters')
33      plt.ylabel('Percentage difference')
34      plt.show()
35      plt.savefig('/home/across/UTK_PhD/Machine_learning_fall_2019/project_4/plots/'+algorithm+'
        _'+color_name+'.png')
```

## A.8 Find Closest Cluster File

```
1   import numpy as np
2   import pandas as pd
3
4   def find_coord_closest_cluster(coordinate, cluster): #find closest cluster for single
        coordinate
5       coord_comp=(coordinate['red']-cluster['red'])**2+(coordinate['green']-cluster['green'])
        **2+(coordinate['blue']-cluster['blue'])**2
6       minid=int(coord_comp.idxmin())
7       mindist=np.sqrt(coord_comp.min())
8
9       return minid,mindist
10
11  def find_all_closest_cluster(data_pixel,cluster): #find closest cluster for all coordinates
12      for i in range(len(data_pixel)):
```

```
13          data_pixel.closest_id.iloc[i], data_pixel.closest_dist.iloc[i]=
        find_coord_closest_cluster(data_pixel.iloc[i],cluster.copy())
14
15      return data_pixel
16
17  def update_pixels_to_closest_cluster(pixel_data,clusters): #updates pixel to the rgb value of
        closest cluster
18      pixel_data['closest_id']=pixel_data['closest_id'].astype(int)
19      r=clusters.red.iloc[pixel_data['closest_id']].copy()
20      g=clusters.green.iloc[pixel_data['closest_id']].copy()
21      b=clusters.blue.iloc[pixel_data['closest_id']].copy()
22      r=r.tolist()
23      g=g.tolist()
24      b=b.tolist()
25      cluster_addition = pd.DataFrame(list(zip(r, g,b)), columns =['r_cluster', 'g_cluster','
        b_cluster'])
26      pixel_data=pd.concat([pixel_data,cluster_addition],axis=1)
27      return pixel_data
28
29  def reset_closest_cluster(data_pixel):
30      data_pixel['closest_id']=0
31      data_pixel['closest_dist']=0
32
33      return data_pixel
```

# References

Michael Steinbach, George Karypis, Vipin Kumar, et al. A comparison of document clustering techniques. In *KDD workshop on text mining*, volume 400, pages 525–526. Boston, 2000.

University of Tennessee. Machine learning cosc 522 lecture 15. `http://web.eecs.utk.edu/~hqi/cosc522/`. Accessed: 11/23/2019.