

Programming Assignment 3

Part A:

1. Misclassified 71 emails which is 27 percent.
2. Train on different numbers of training examples:
 - a. 50 – misclassified 132 which is 50 percent.
 - b. 100 – misclassified 125 which is 48 percent.
 - c. 400 – misclassified 117 which is 45 percent.

3. Code

```
4. #!/usr/bin/env python2
5. # -*- coding: utf-8 -*-
6. """
7. Created on Wed Oct 11 20:21:17 2017
8.
9. @author: rditljtd
10. """
11.
12. import numpy as np
13. from scipy import misc
14. from scipy import sparse as sps
15. import matplotlib.pyplot as plt
16. from math import *
17.
18. #find probability of email being spam given that it contains a particular word
19. #this equals the probability of seeing this word given that it is a spam * proba
    bility of it being seen overall / probability of it being spam
20.
21. #find probability of email being non-
    spam given that it contains a particular word
22. #this equals the probability of seeing this word given that the email is non-
    spam *
23.     #probability of seeing this word overall
24.     #/probability of email being non-spam
25.
26. # Load the labels for the training set
27. train_labels = np.loadtxt('train-labels-50.txt',dtype=int)
28.
29. # Get the number of training examples from the number of labels
30. numTrainDocs = train_labels.shape[0]
31.
32. # This is how many words we have in our dictionary
33. numTokens = 2500
34.
35. # Load the training set feature information
36. M = np.loadtxt('train-features-50.txt',dtype=int)
37.
38. # Create matrix of training data
39. train_matrix = sps.csr_matrix((M[:,2], (M[:,0], M[:,1])),shape=(numTrainDocs,num
    Tokens))
40.
41. #train_labels[i] = ith document label: spam or not spam
42.
43. #train_matrix[i,:] = ith document
44.
45. #train_matrix[i:, j] = jth word in ith document
46.
```

```

47. #print train_matrix[69]
48.
49. print train_matrix.shape[1]
50.
51. spam_prob = (sum(x == 1 for x in train_labels))/(numTrainDocs*(1.0))
52.
53. nonspam_prob = (sum(x == 0 for x in train_labels))/(numTrainDocs*(1.0))
54.
55. print spam_prob, nonspam_prob
56.
57. #create variable for sum of all non spam emails
58. num_nonspam = (sum(x==0 for x in train_labels))
59. #create variable for sum of all spam emails
60. num_spam = (sum(x==1 for x in train_labels))
61.
62. #create array for probability all words in spam and non spam [spam_prob, nonspam
    _prob]
63. prob_for_each_word = []
64.
65.
66. #add up the instances of this word [spam, nonspam]
67. sum_for_this_word = [0, 0]
68.
69. #create array for how many times each word was foud [[spam, nonspam]] i = word
70. found_for_all_emails = []
71.
72. #loop through each email
73. for j in range (0, train_matrix.shape[0]):
74.
75.     #array of words in email [instances]
76.     words_in_email = train_matrix[j].toarray()[0]
77.     #print train_matrix[j],
78.
79.
80.
81.     #loop through each word in dictionary and see if it is in email
82.     for k in range (0, len(words_in_email)):
83.
84.         if j == 0:
85.             #create array for if this word was found [spam, nonspam]
86.             found_for_all_emails.append([0, 0])
87.
88.             #if email is non spam and there is more than zero instances of this word
89.
90.             if (train_labels[j] == 0 and words_in_email[k] > 0):
91.
92.                 #increment the sum for this word for non spam emails
93.                 found_for_all_emails[k][1] += 1
94.
95.                 #else if the email is spam and there is more than zero instances of this
96.                 word
97.                 elif (train_labels[j] == 1 and words_in_email[k] > 0):
98.
99.                     #print words_in_email[k]
100.                     #increment the sum for this word for spam emails
101.                     found_for_all_emails[k][0] += 1
102.
103.                     #print found_for_all_emails
104.                     #sum_for_this_word = sum(found_for_all_emails[:1])
105.                     #add the probability for this word to array

```

```

105.         #create an array to hold the probability for each word [spam, nonspam] i
    = word
106.         prob_for_each_word = []
107.
108.         #create array to hold the overall probability for each word [num] i = wo
rd
109.         overall_prob_for_each_word = []
110.
111.         for i in range (0, numTokens):
112.
113.             prob_for_each_word.append([0,0])
114.             overall_prob_for_each_word.append(0)
115.
116.             #divide by the total number of spam emails
117.             prob_for_each_word[i][0] = (found_for_all_emails[i][0]/(num_spam*(1.
0)))
118.
119.             #divide by the total number of nonspam emails
120.             prob_for_each_word[i][1] = (found_for_all_emails[i][1]/(num_nonspam*
(1.0)))
121.
122.             overall_prob_for_each_word[i] = (found_for_all_emails[i][0] + found_
for_all_emails[i][1])/((num_nonspam + num_spam)*(1.0))
123.
124.
125.         #print prob_for_each_word
126.         #print overall_prob_for_each_word
127.
128.
129.
130.         # Load the labels for the training set
131.         test_labels = np.loadtxt('test-labels.txt',dtype=int)
132.
133.         # Get the number of training examples from the number of labels
134.         numTestDocs = test_labels.shape[0]
135.
136.         # Load the training set feature information
137.         N = np.loadtxt('test-features.txt',dtype=int)
138.
139.         # Create matrix of training data
140.         test_matrix = sps.csr_matrix((N[:,2], (N[:,0], N[:,1])),shape=(numTestDo
cs,numTokens))
141.
142.         prob_for_all_test_emails = []
143.
144.         #iterate through each email in the test set
145.         for i in range (0, (test_matrix.shape[0])):
146.
147.             #array of words in email [instances]
148.             words_in_test_email = test_matrix[i].toarray()[0]
149.             #print train_matrix[j],
150.
151.             prob_spam_given_all_words = 0
152.             prob_nonspam_given_all_words = 0
153.
154.             #loop through each word in dictionary and see if it is in email
155.             for k in range (0, len(words_in_test_email)):
156.
157.                 prob_spam_given_word = 0
158.                 prob_nonspam_given_word = 0
159.

```

```

160.             #if email is non spam and there is more than zero instances of t
            his word
161.             if (words_in_test_email[k] > 0):
162.                 #prob_nonspam_given_word = float(((prob_for_each_word[k][1])
                *(1.0))*overall_prob_for_each_word[k])/((nonspam_prob)*1.0)
163.                 #prob_spam_given_word = float(((prob_for_each_word[k][0])*(1
                .0))*overall_prob_for_each_word[k])/((spam_prob)*1.0)
164.                 if (prob_for_each_word[k][1] > 0):
165.                     prob_nonspam_given_word = float(math.exp(math.log(prob_f
                or_each_word[k][1]))) + float(math.exp(math.log(overall_prob_for_each_word[k])))
                - float(math.exp(math.log(nonspam_prob)))
166.                 if (prob_for_each_word[k][0] > 0):
167.                     prob_spam_given_word = float(math.exp(math.log(prob_for_
                each_word[k][0]))) + float(math.exp(math.log(overall_prob_for_each_word[k]))) -
                float(math.exp(math.log(spam_prob)))
168.
169.                 if (prob_spam_given_word > 0):
170.                     #print prob_spam_given_all_words,
171.                     #prob_spam_given_all_words += float((prob_spam_given_word*(1
                .0)))
172.                     #print prob_spam_given_all_words
173.                     prob_spam_given_all_words += prob_spam_given_word
174.
175.                 if (prob_nonspam_given_word > 0):
176.                     #print prob_nonspam_given_all_words,
177.                     #prob_nonspam_given_all_words += float((prob_nonspam_given_w
                ord*(1.0)))
178.                     prob_nonspam_given_all_words += prob_nonspam_given_word
179.
180.                     #prob_spam_given_all_words += float(math.exp(math.log(spam_prob)))
181.                     #prob_nonspam_given_all_words += float(math.exp(math.log(nonspam_pro
                b)))
182.                     prob_for_all_test_emails.append([prob_spam_given_all_words, prob_non
                spam_given_all_words])
183.
184.             print prob_for_all_test_emails
185.
186.             classify_each_test_email = []
187.             spam_count = 0
188.             nonspam_count = 0
189.             incorrect = []
190.             print len(prob_for_all_test_emails)
191.             for k in range(0, len(prob_for_all_test_emails)):
192.                 if prob_for_all_test_emails[k][0] > prob_for_all_test_emails[k][1]:
193.                     classify_each_test_email.append(1)
194.                     spam_count += 1
195.                     if(test_labels[k] != classify_each_test_email[k]):
196.                         incorrect.append([k, classify_each_test_email[k], test_label
                s[k]])
197.                 elif prob_for_all_test_emails[k][1] > prob_for_all_test_emails[k][0]
                :
198.                     classify_each_test_email.append(0)
199.                     nonspam_count += 1
200.                     if(test_labels[k] != classify_each_test_email[k]):
201.                         incorrect.append([k, classify_each_test_email[k], test_label
                s[k]])
202.                 else:
203.                     classify_each_test_email.append('unknown')
204.                     incorrect.append([k, classify_each_test_email[k], test_labels[k]
                ])

```

```

205.
206.     print classify_each_test_email
207.     print "spam count: " + str(spam_count) + " spam percentage: " + str(float(
208.         spam_count*100)/(len(classify_each_test_email)))
209.     print "nonspam count: " + str(nonspam_count) + " nonspam percentage: " +
210.         str(float((nonspam_count*100)/(len(classify_each_test_email))))
211.     print "number misclassified: " + str(len(incorrect)) + " percentage mis
212.         classified: " + str(float(((len(incorrect))*100)/(len(classify_each_test_email)))
213.         )
214.
215.     #def nBayes(train_matrix, train_labels):
216.     #     print 'called nBayes'
217.     #     #reset num_unique_words_in_spam ratio and wordsInSpam
218.     #
219.     #     #create variable for number of unique words in spam emails
220.     #     num_unique_words_in_spam = 0
221.     #
222.     #     #create variable for number of unique words in nonspam emails
223.     #     num_unique_words_in_nonspam = 0
224.     #
225.     #     #total number of words in spam emails
226.     #     totalWordsInSpam = 0
227.     #     totalWordsInNonSpam = 0
228.     #     wordsInSpam = {}
229.     #     wordsInNonSpam = {}
230.     #
231.     #     #loop through each email
232.     #     for i in range(0, train_matrix.shape[0]):
233.     #         #get array with words and # of instance of these words
234.     #         #two dim array with first being word id and second being # of i
235.     #         #instances
236.     #         wordsInEmail = train_matrix[i].toarray()[0]
237.     #         #print wordsInEmail
238.     #         #print "-----"
239.     #         #for each possible word in the email
240.     #         for j in range(0, len(wordsInEmail)):
241.     #
242.     #             if (wordsInEmail[j] == 0):
243.     #                 continue
244.     #             #if the email is a spam email
245.     #             if train_labels[i] == 1:
246.     #
247.     #                 #increment summation value for a unique word in this sp
248.     #                 #am email
249.     #                 num_unique_words_in_spam = num_unique_words_in_spam+1
250.     #                 #add the # of instances of this word to the number of w
251.     #                 #ords total in the email
252.     #                 totalWordsInSpam = totalWordsInSpam + wordsInEmail[j]
253.     #
254.     #                 #if this word has already been seen in a spam email
255.     #                 if (j in wordsInSpam.keys()):
256.     #                     #add the # of instances of this word to the # of in
257.     #                     #stances of this word in all spam emails
258.     #                     wordsInSpam[j] += wordsInEmail[j]
259.     #
260.     #                 #word has not been seen in a spam email yet

```

```

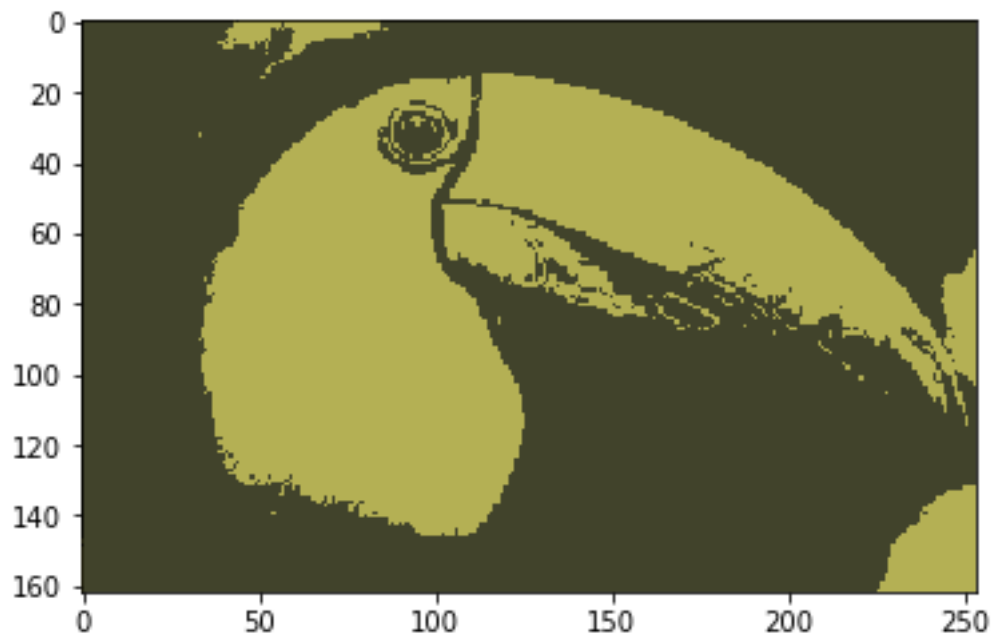
258.         #         else:
259.         #         #set the # of instances of this word in a spam email
        l to the # of instances in this email
260.         #         wordsInSpam[j] = wordsInEmail[j]
261.         #
262.         #         #if email is non-spam
263.         #         if train_labels[i] == 0:
264.         #         #increment summation value for a unique word in this non
        spam email
265.         #         num_unique_words_in_nonspam = num_unique_words_in_nonsp
        am+1
266.         #         #add the # of instances of this word to the number of w
        ords total in the email
267.         #         totalWordsInNonSpam = totalWordsInNonSpam + wordsInEmail[j]
268.         #
269.         #         #if this word has already been seen in a non spam email
270.         #         if (j in wordsInNonSpam.keys()):
271.         #         #add the # of instances of this word to the # of in
        stances of this word in all non spam emails
272.         #         wordsInNonSpam[j] += wordsInEmail[j]
273.         #
274.         #         #word has not been seen in a non spam email
275.         #         else:
276.         #         #set the # of instances of this word in non spam em
        ails to the # of instances in this email
277.         #         wordsInNonSpam[j] = wordsInEmail[j]
278.         #
279.         #         #I think what actually needs to be done here is that for each unique
        word in the dictionary, calculate the percentage of spam emails that contain t
        hat word
280.         #         #then calculate the # of non spam emails that contain that word. Mu
        ltiply the percentages for every word in
281.         #
282.         #
283.         #         #dictionary of words and # of instances in spam emails
284.         #         print wordsInSpam
285.         #         #dictionary of words and # of instances in non spam emails
286.         #         print wordsInNonSpam
287.         #         #of unique words in spam email
288.         #         print num_unique_words_in_spam
289.         #         #of unique words in non spam email
290.         #         print num_unique_words_in_nonspam
291.         #         #total # of words in spam emails
292.         #         print totalWordsInSpam
293.         #         #total # of words in non spam emails
294.         #         print totalWordsInNonSpam
295.         #
296.         #         #return # of unique words in spam + 1 / total # of words in spam em
        ails + # of words in dictionary,
297.         #         #return # of unique words in non spam +1 / total # of words in non
        spam emails + # of words in dictionary
298.         #         return (num_unique_words_in_spam+1)/((totalWordsInSpam+train_matrix
        .shape[1])*1.0), (num_unique_words_in_nonspam+1)/((totalWordsInNonSpam+train_mat
        rix.shape[1])*1.0)
299.         #
300.         #print nBayes(train_matrix, train_labels)

```

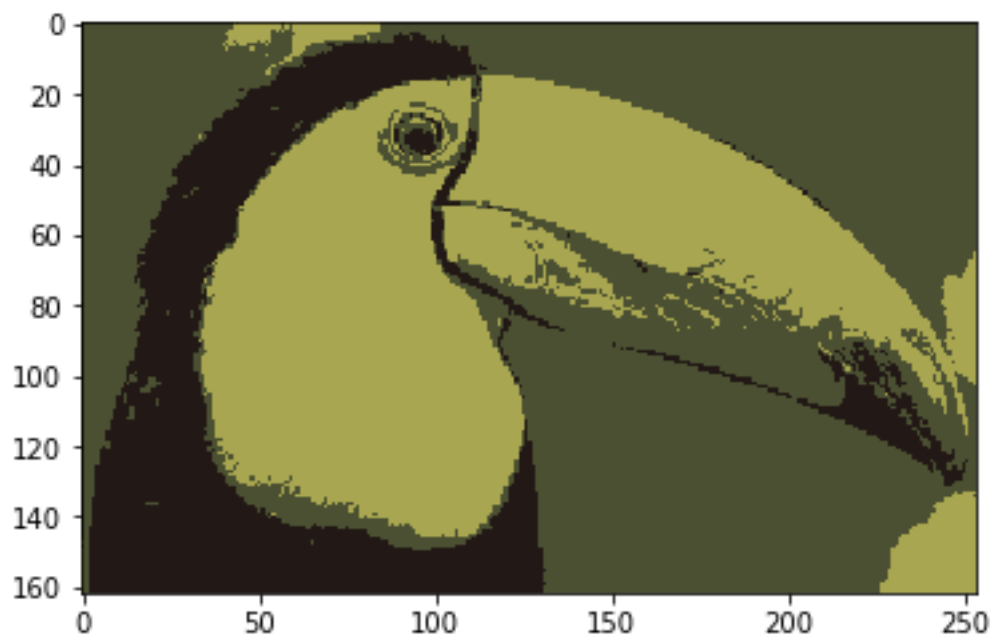
Part B:

1. Images for $k = 2..15$

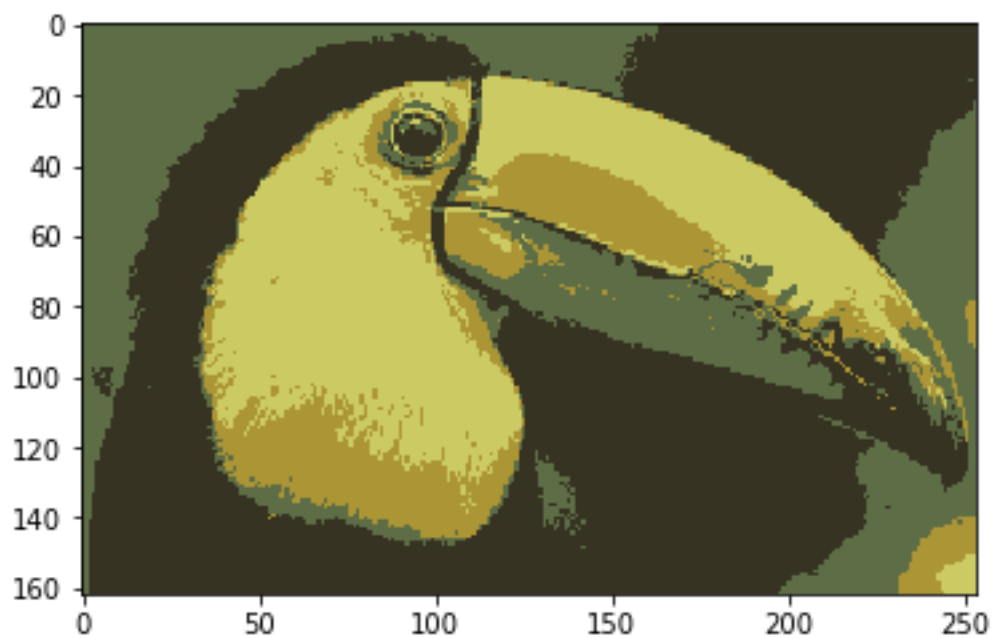
a. 2



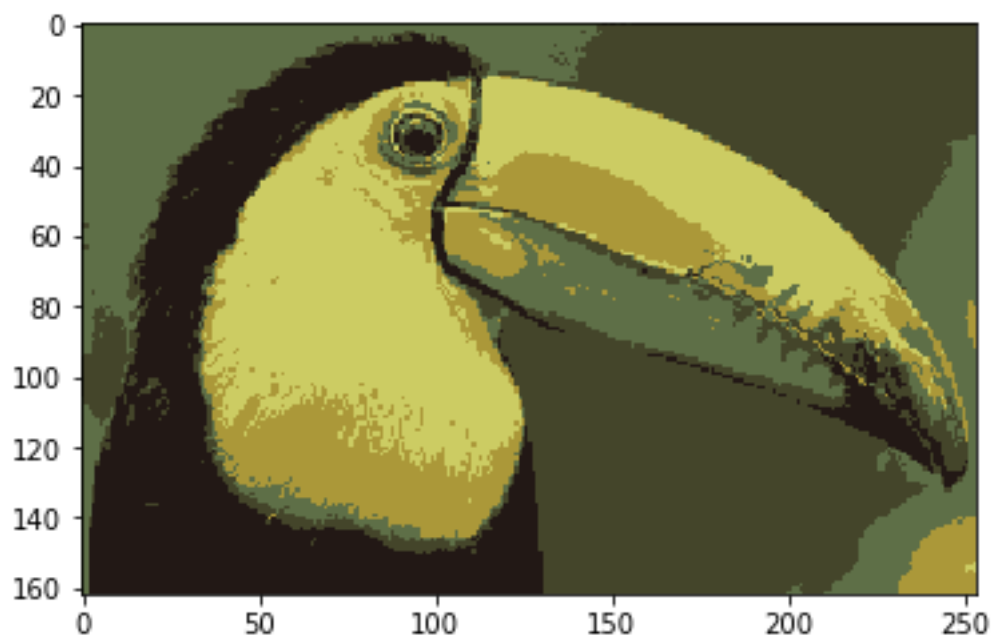
b. 3



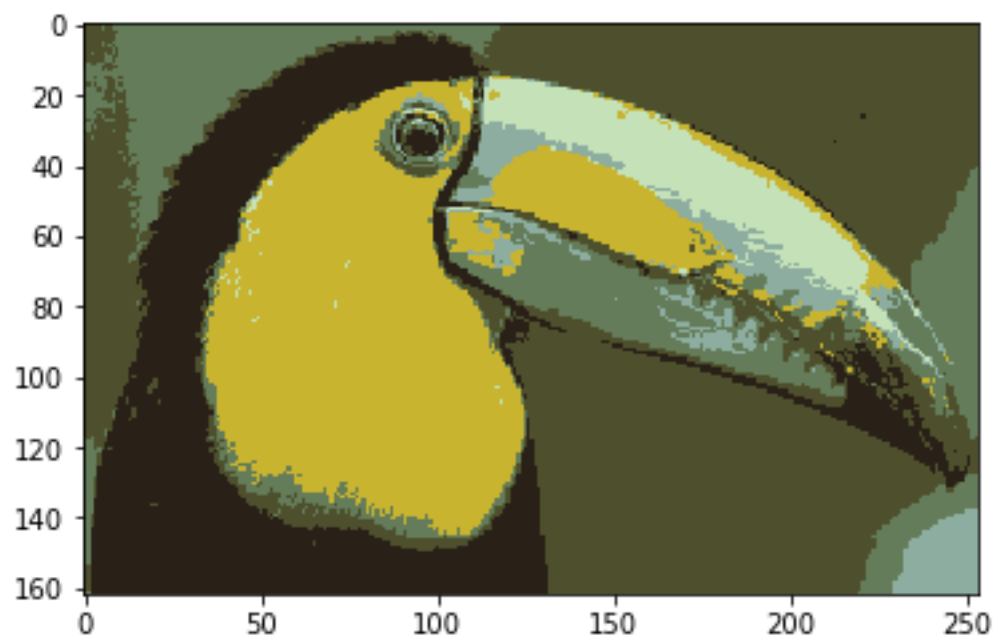
c. 4



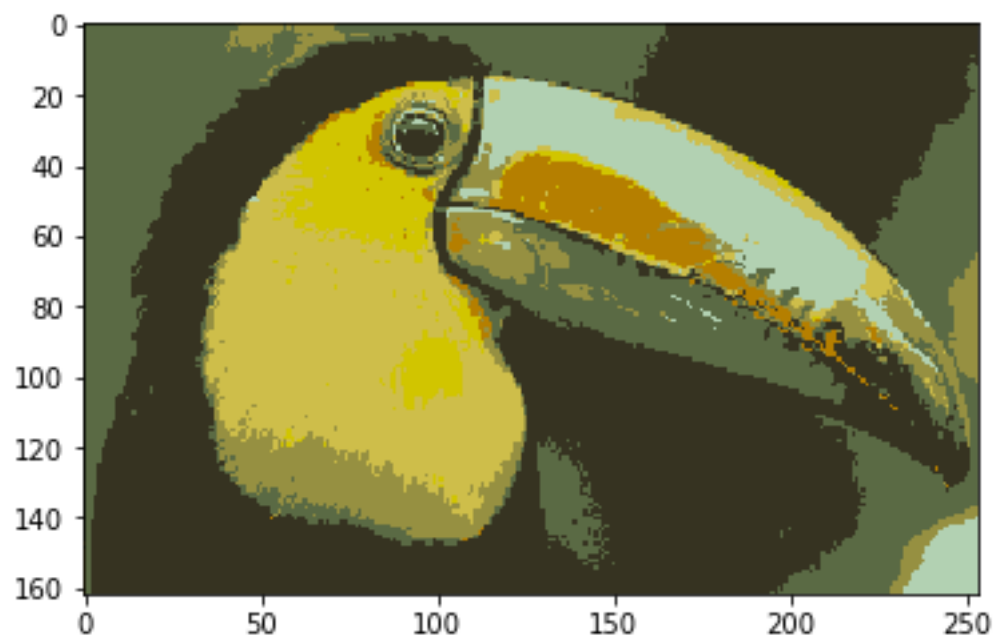
d. 5



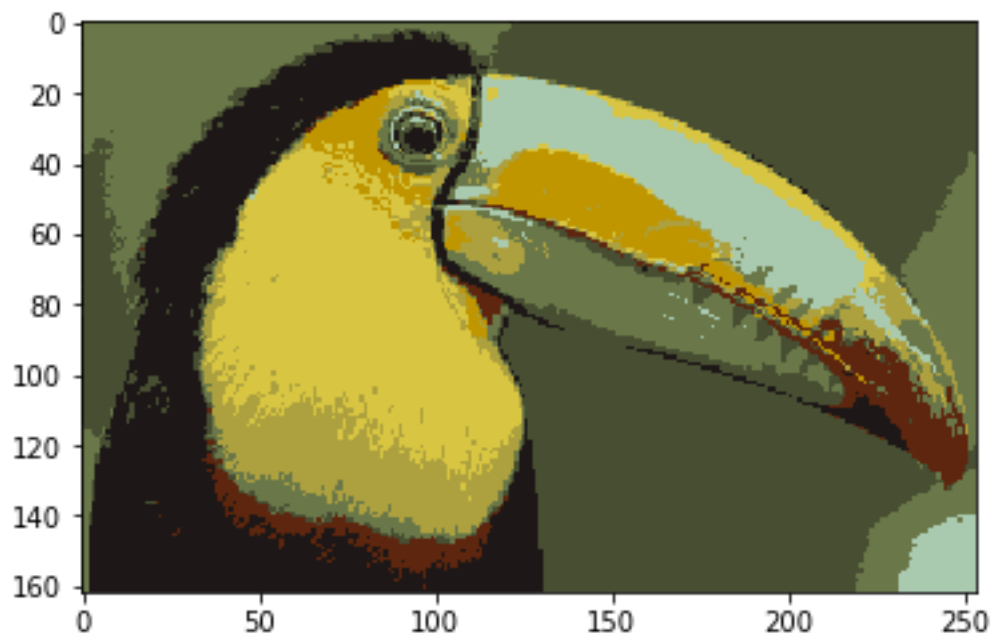
e. 6



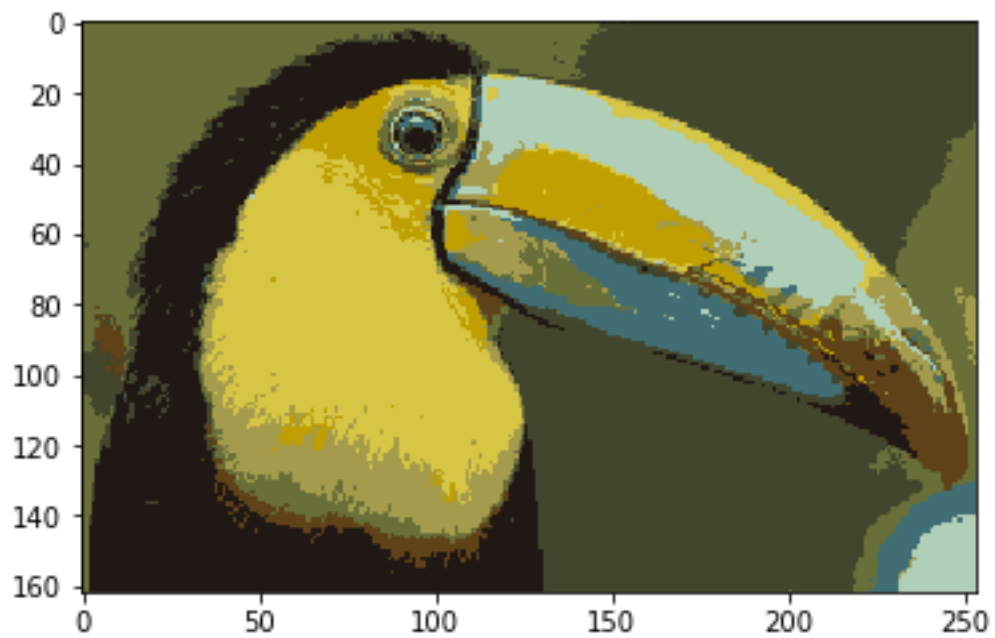
f. 7



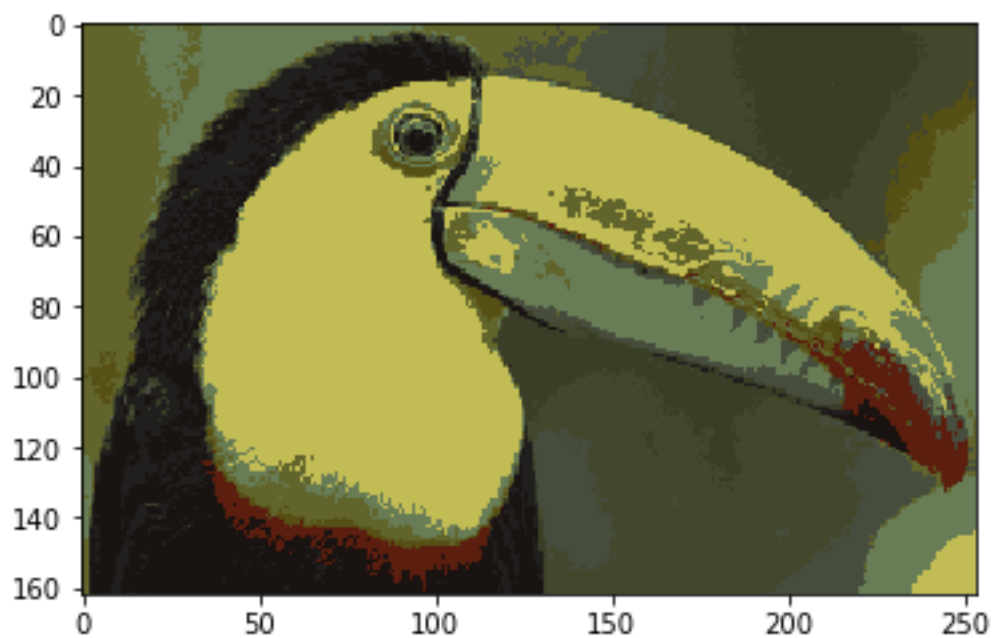
g. 8



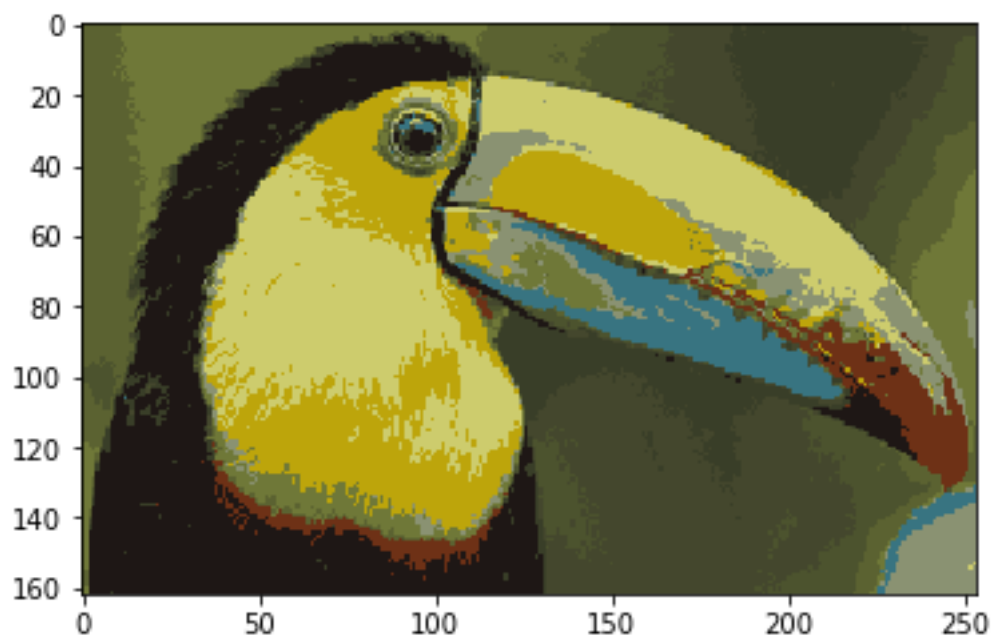
h. 9



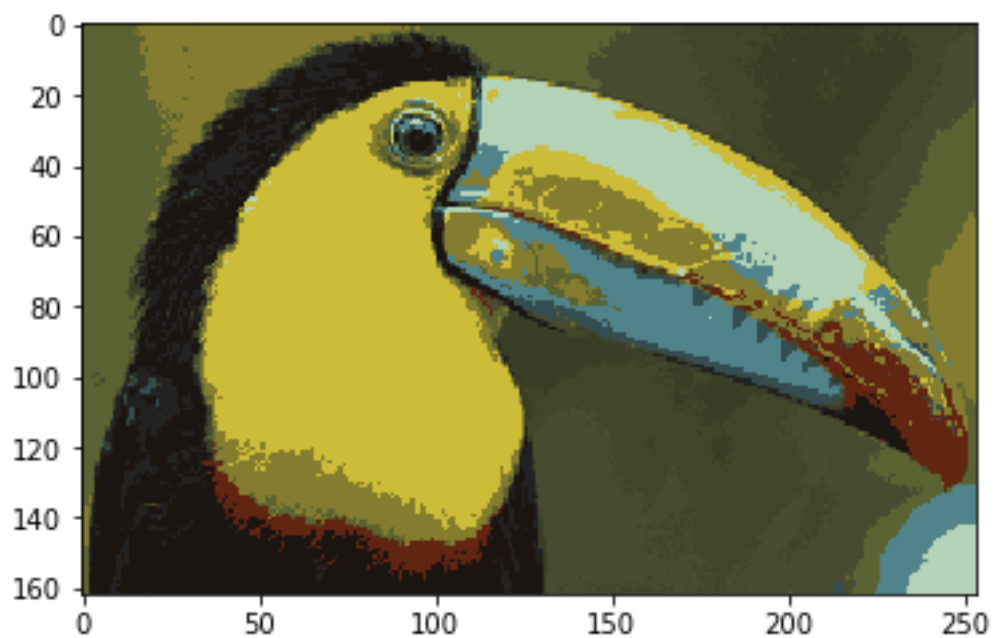
i. 10



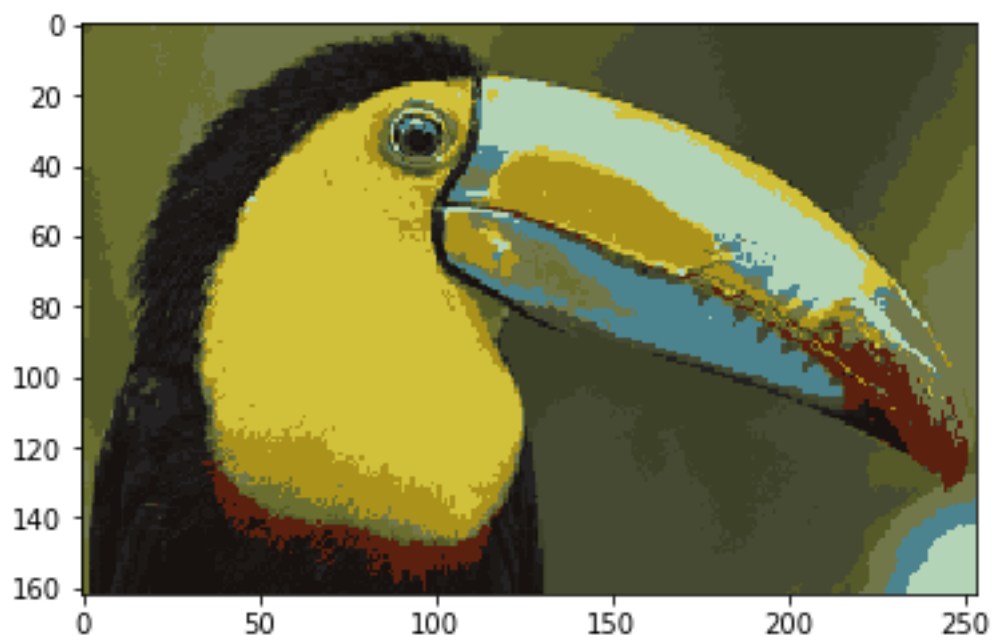
j. 11



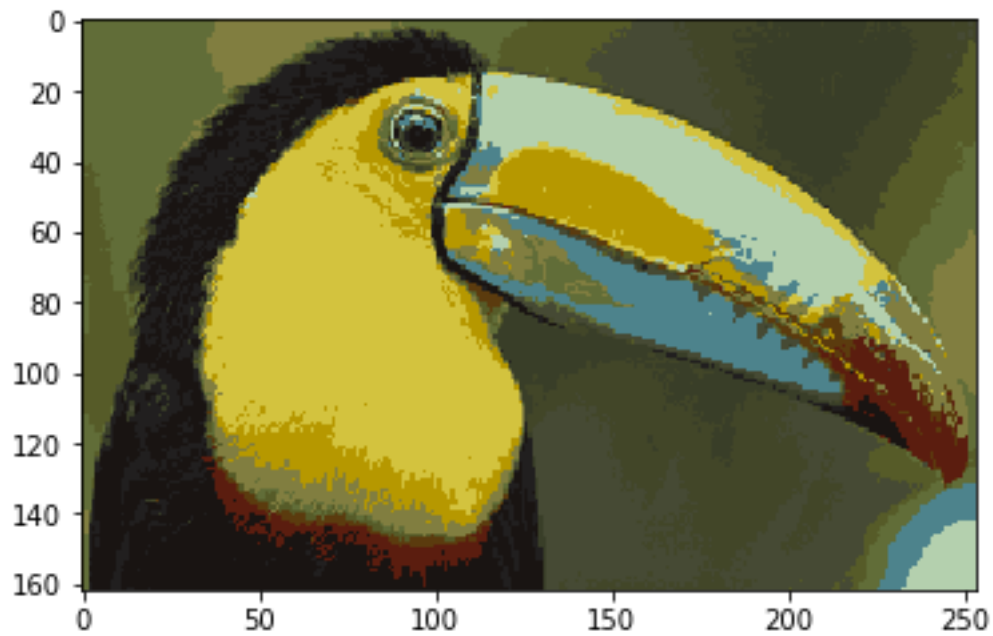
k. 12



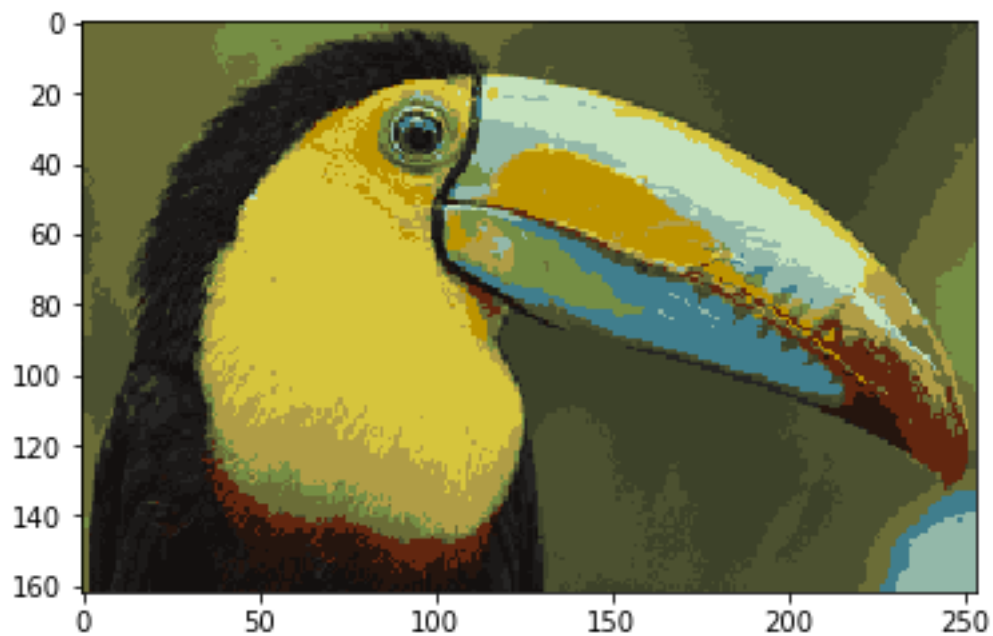
l. 13



m. 14



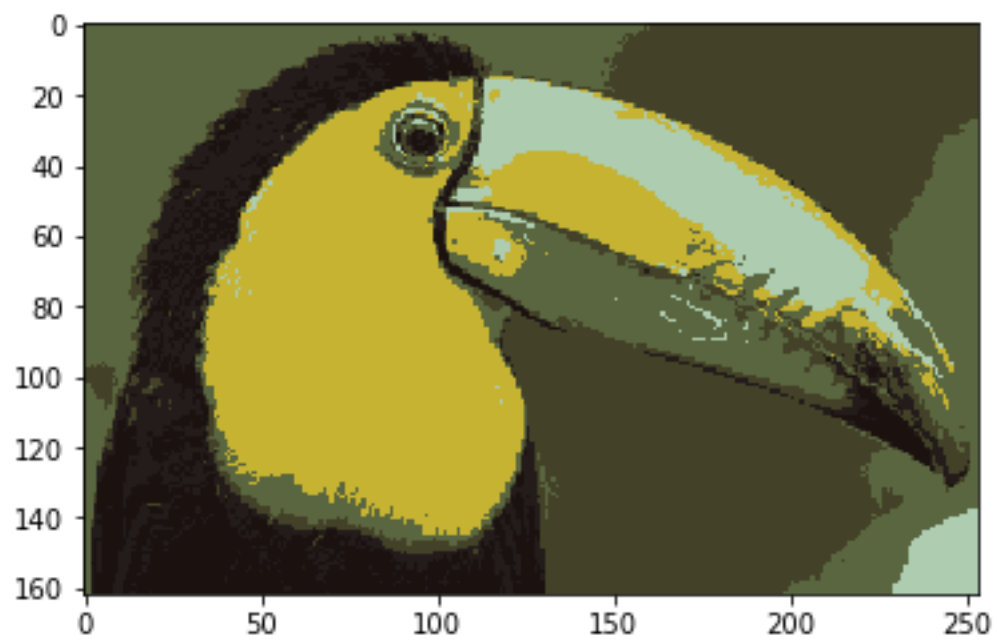
n. 15



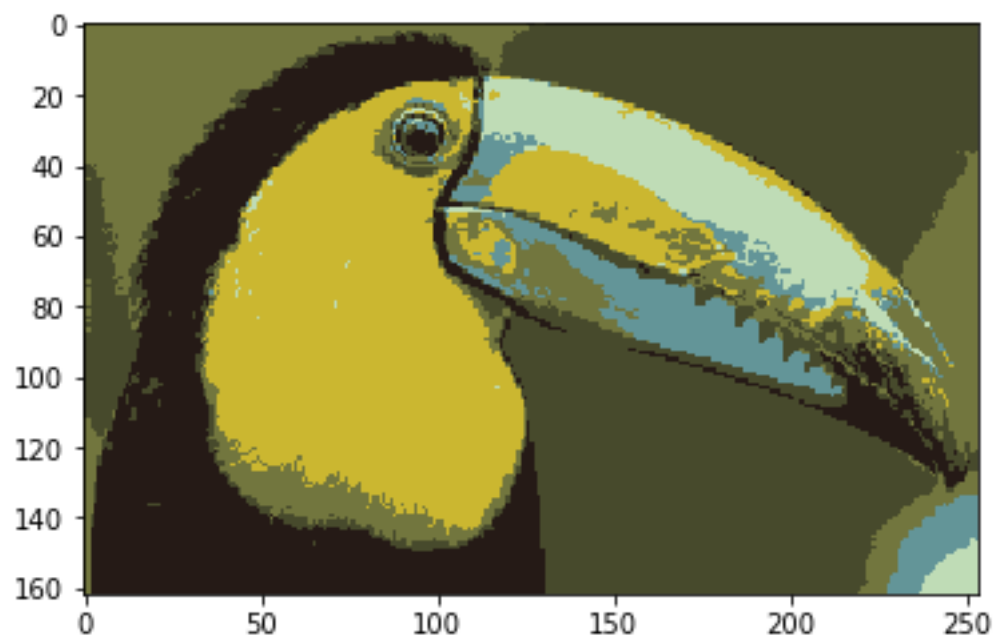
Overall the images improve as k increases. Depending on the centroid clusters, the image improves variably, but as a whole, they improve noticeably in the images above. We can see this as we scroll down through them.

2. $k = 6$ examples

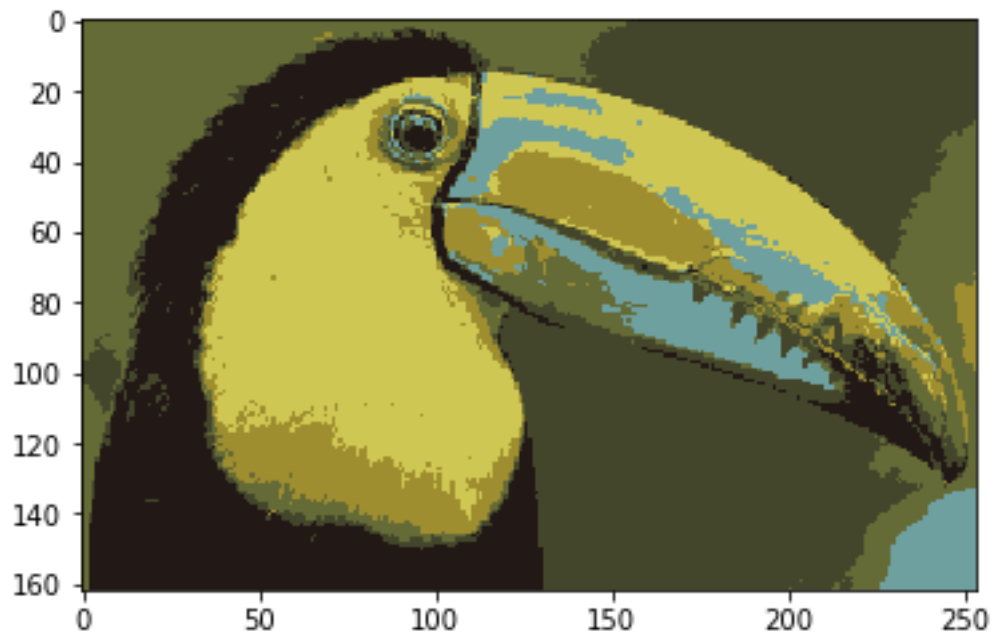
a.



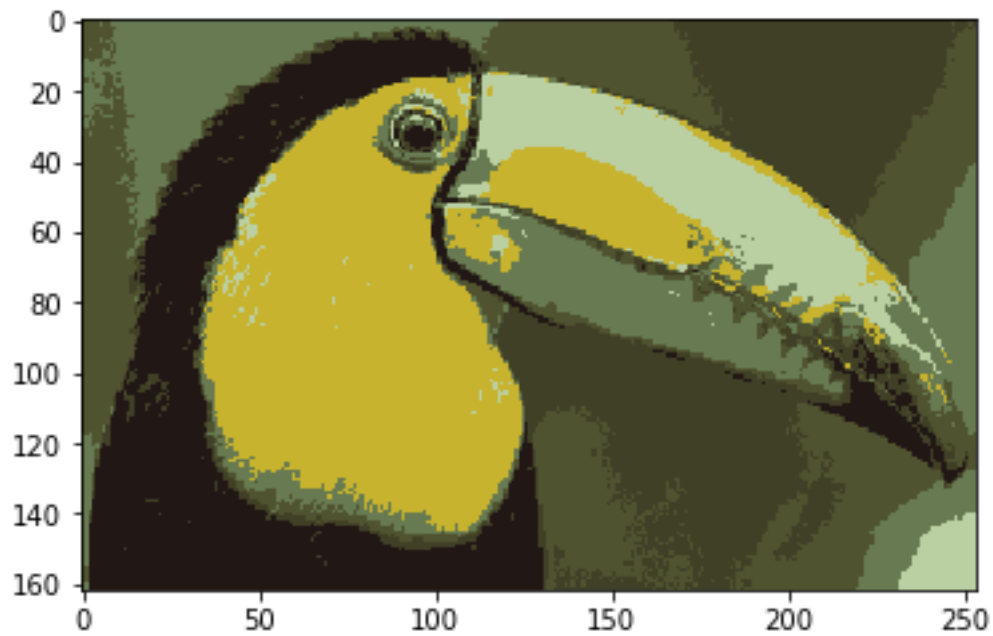
b.



c.

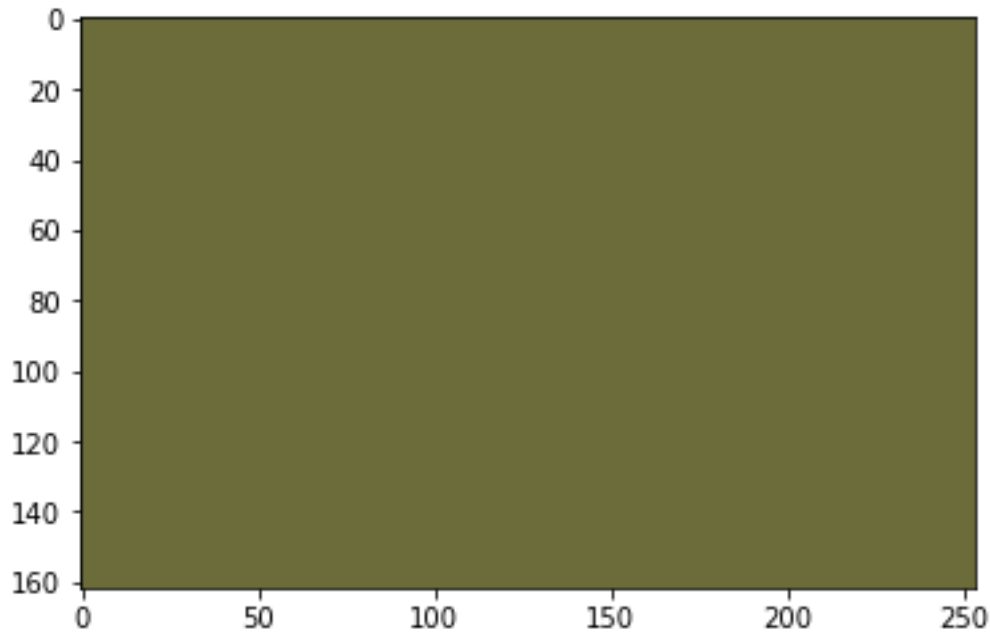


d.



We can see by looking at these images that the centroid clusters are different between each of them. They are close in most of them, but if we examine closer we can see pixels with slightly different RGB values and centroids.

3. $k = 1$
 - a.



It takes one iteration to get the average color for the image, which is returned as the only RGB value. Since there is only one centroid, it becomes the average for all the points immediately, and all the points fit to it since there is no other centroid.

4. Code

```
5. #!/usr/bin/env python2
6. # -*- coding: utf-8 -*-
7. """
8. Created on Wed Oct 18 22:01:20 2017
9.
10. @author: rditljtd
11. """
12. import numpy as np
13. from scipy import misc
14. from scipy import sparse as sps
15. import matplotlib.pyplot as plt
16. from collections import namedtuple
17. from math import sqrt
18. import random
19. try:
20.     import Image
21. except ImportError:
22.     from PIL import Image
23.
24. #make array of RGB values for each pixel in image
25. A = misc.imread('b_small.tiff', mode='RGB')
26.
27. plt.imshow(A)
28.
29. num_of_centroids = 16
30.
31. #instantiate array of random 16 pixels (used for cluster centroids)
32. random_16 = []
33.
```



```

34. #loop 16 times to actually choose the random 16 pixels for cluster centroids
35. for i in range (0, num_of_centroids):
36.
37.     #generate random x-value
38.     random_x = random.randint(0, 253)
39.     #generate random y-value
40.     random_y = random.randint(0, 161)
41.
42.     #add random pixel to array of 16 random centroids
43.     random_16.append(A[random_y][random_x])
44.
45. #set number of iterations
46. iterations = 50
47.
48.
49. #print out the centroids
50. #print random_16
51.
52. #print out the last pixel in the image
53. print A[161][253]
54.
55. #create an array for all pixels that is equal to the array of RGB values for the
    image
56. all_pixels = A
57.
58. #create an array for associating pixels to centroids
59. pixels_to_centroids = []
60.
61. #loop the number of iterations
62. for a in range(0, iterations):
63.     print a+1,
64.
65.     #loop through all y-values
66.     for i in range (0, len(A)):
67.
68.         #loop through all x-values
69.         for j in range (0, len(A[0])):
70.
71.             #print out this pixel's RGB values
72.             #print A[i][j]
73.
74.             #set this_pixel variable equal to this pixel's RGB values
75.             this_pixel = A[i][j]
76.
77.             #set distance = to maximum distance possible
78.             distance = 3*253
79.
80.             #create a variable for closest centroid
81.             pixel_to_centroid = [0, 0, 0, [0, 0, 0]]
82.
83.             #loop through the cluster centroids
84.             for k in range(0, len(random_16)):
85.
86.                 #set this_distance equal to 0
87.                 this_distance = 0
88.
89.                 #loop through the red, green, blue values for this centroid
90.                 for l in range (0, len(random_16[0])):
91.
92.                     #calculate the distance between the r, g, or b value for thi
                        s centroid and this pixel

```

```

93.         this_distance += abs(int(random_16[k][1]) -
    int(this_pixel[1]))
94.
95.         #if the distance between the RGB values for this centroid is less
    s than all other centroids
96.         if (this_distance < distance):
97.
98.             #set the distance to the minimum distance calculated
99.             distance = this_distance
100.
101.             #set the closest centroid
102.             pixel_to_centroid = [int(k), int(i), int(j), A[i][j]
    ]
103.
104.             #set this pixel's centroid to closest
105.             pixels_to_centroids.append(pixel_to_centroid)
106.
107.             #all_pixels[i][j] = this_pixel
108.
109.
110.         #create variable for number of values associated with each centroid
111.
112.         num_assoc_16 = []
113.
114.         #create variable for sum of values associated with each centroid
115.         sum_16 = []
116.
117.         #create variable for new centroids
118.         average_16 = []
119.
120.         #loop through centroids
121.         for i in range(0, len(random_16)):
122.             num_assoc_16.append(0)
123.             sum_16.append([0, 0, 0])
124.
125.         #loop through each pixel
126.         for j in range(0, len(pixels_to_centroids)):
127.
128.             #if this pixel is associated with this cluster
129.             if (pixels_to_centroids[j][0] == i):
130.
131.                 #number of points associated with this centroid
132.                 num_assoc_16[i] += 1
133.                 #sum up the R values
134.                 sum_16[i][0] += int(pixels_to_centroids[j][3][0])
135.                 #sum up the G values
136.                 sum_16[i][1] += int(pixels_to_centroids[j][3][1])
137.                 #sum up the B values
138.                 sum_16[i][2] += int(pixels_to_centroids[j][3][2])
139.
140.         for i in range(0, len(sum_16)):
141.             average_16.append([sum_16[i][0] / num_assoc_16[i],
142.                 sum_16[i][1] / num_assoc_16[i],
143.                 sum_16[i][2] / num_assoc_16[i]])
144.
145.         #set centroids as the averages
146.         if (np.array_equal(np.array(random_16), np.array(average_16))):
147.             break
148.         random_16 = average_16
149.
150.         B = A

```

```
150.
151.     #loop through all pixels
152.     for a in range (0, len(pixels_to_centroids)):
153.
154.         i=pixels_to_centroids[a][1]
155.         j=pixels_to_centroids[a][2]
156.
157.         k=pixels_to_centroids[a][0]
158.         B[i][j] = random_16[k]
159.
160.     plt.imshow(B)
161.     plt.savefig('kmeans-' + str(iterations) + '.png')
162.
163.     #loop through y-values
164.     #for i in range (0, len(all_pixels)):
165.
166.         #loop through x-values
167.         #for j in range (0, len(all_pixels[0])):
168.
169.             #loop through centroids
170.             #for k in range (0, len(random_16)):
171.
172.                 #if this pixel is closest to this centroid
173.                 #if all_pixels[i][j] = random_16[k]
174.
```