

Using K-means for color compression

```
In [1]: 1 import numpy as np
        2 import pandas as pd
        3
        4 %pylab inline
        5 import plotly.graph_objects as go
        6
        7 from sklearn.cluster import KMeans
```

%pylab is deprecated, use %matplotlib inline and import the required libraries.

Populating the interactive namespace from numpy and matplotlib

Read in data The "data" in this case is not a pandas dataframe. It's a photograph, which we'll convert into a numerical array.

```
In [2]: 1 # Read in a photo
        2 img = plt.imread('images_(79)[1].jpeg')
```

```
In [3]: 1 img
```

```
Out[3]: array([[220, 189, 125],
               [220, 186, 122],
               [220, 185, 117],
               ...,
               [206, 135, 133],
               [201, 158, 152],
               [191, 169, 158]],

            [[221, 187, 124],
             [219, 185, 121],
             [218, 183, 115],
             ...,
             [211, 131, 130],
             [205, 152, 146],
             [197, 164, 155]],

            [[217, 181, 121],
             [218, 180, 117],
             [217, 180, 112],
             ...,
             [219, 125, 125],
             [211, 142, 137],
             [206, 156, 149]],

            ...,

            [[ 1,  1,  1],
             [255, 255, 255],
             [255, 255, 255],
             ...,
             [255, 255, 255],
             [255, 255, 255],
             [255, 255, 255]],

            [[ 1,  1,  1],
             [255, 255, 255],
             [255, 255, 255],
             ...,
             [255, 255, 255],
             [255, 255, 255],
             [255, 255, 255]],

            [[ 1,  1,  1],
             [255, 255, 255],
             [255, 255, 255],
             ...,
             [255, 255, 255],
             [255, 255, 255],
             [255, 255, 255]]], dtype=uint8)
```

```
In [4]: 1 # Display the photo and its shape
        2 print(img.shape)
        3 plt.imshow(img)
        4 plt.axis('off');
```

(280, 463, 3)



```
In [5]: 1 # Reshape the image so that each row represents a single pixel
        2 # defined by three values: R, G, B
        3 img_flat = img.reshape(img.shape[0]*img.shape[1], 3)
        4 img_flat[:5, :]
```

```
Out[5]: array([[220, 189, 125],
               [220, 186, 122],
               [220, 185, 117],
               [225, 189, 115],
               [230, 193, 113]], dtype=uint8)
```

```
In [6]: 1 # Create a pandas df with r, g, and b as columns
        2 img_flat_df = pd.DataFrame(img_flat, columns = ['r', 'g', 'b'])
        3 img_flat_df.head()
```

Out[6]:

	r	g	b
0	220	189	125
1	220	186	122
2	220	185	117
3	225	189	115
4	230	193	113

```
In [7]: 1 # Create 3D plot where each pixel in the `img` is displayed in its actual
2 trace = go.Scatter3d(x = img_flat_df.r,
3                     y = img_flat_df.g,
4                     z = img_flat_df.b,
5                     mode='markers',
6                     marker=dict(size=1,
7                                 color=['rgb({}, {}, {})'.format(r,g,b) for
8                                     in zip(img_flat_df.r.values,
9                                             img_flat_df.g.values,
10                                            img_flat_df.b.values)],
11                                     opacity=0.5))
12
13 data = [trace]
14
15 layout = go.Layout(margin=dict(l=0,
16                                r=0,
17                                b=0,
18                                t=0),
19                    )
20
21 fig = go.Figure(data=data, layout=layout)
22 fig.update_layout(scene = dict(
23                     xaxis_title='R',
24                     yaxis_title='G',
25                     zaxis_title='B'),
26                 )
27 fig.show()
```



Cluster the data: k = 1

In [8]:

```
1 # Instantiate the model
2 kmeans = KMeans(n_clusters=1, random_state=42).fit(img_flat)
```

```
In [9]: 1 # Copy `img_flat` so we can modify it
2 img_flat1 = img_flat.copy()
3
4 # Replace each row in the original image with its closest cluster center
5 for i in np.unique(kmeans.labels_):
6     img_flat1[kmeans.labels_==i,:] = kmeans.cluster_centers_[i]
7
8 # Reshape the data back to (640, 480, 3)
9 img1 = img_flat1.reshape(img.shape)
10
11 plt.imshow(img1)
12 plt.axis('off');
```



So what happened? Well, let's run through the K-means steps:

We randomly placed our centroid in the colorspace. We assigned each point to its nearest centroid. Since there was only one centroid, all points were assigned to the same centroid, and thus to the same cluster. We updated the centroid's location to the mean location of all of its points. Again, since there is only a single centroid, it updated to the mean location of every point in the image. Repeat until the model converges. In this case, it only took one iteration for the model to converge. We then updated each pixel's RGB values to be the same as the centroid's. The result is the image of our tulips when every pixel is replaced with the average color. The average color of this photo was brown—all the colors muddled together.

We can verify this for ourselves by manually calculating the average for each column in the flattened array. This will give us the average R value, G value, and B value.

```
In [10]: 1 # Calculate mean of each column in the flattened array
2 column_means = img_flat.mean(axis=0)
3
4 print('column means: ', column_means)
```

```
column means: [179.55814563 111.90674946  79.88253625]
```

Now, we can compare this to what the K-means model calculated as the final location of its one centroid.

```
In [11]: 1 print('cluster centers: ', kmeans.cluster_centers_)  
cluster centers:  [[179.55814563 111.90674946  79.88253625]]
```

They're the same! Now, let's return to the 3-D rendering of our data, only this time we'll add the centroid.

```
In [12]: 1 # Create 3-D plot where each pixel in the `img` is displayed in its actual
2         trace = go.Scatter3d(x = img_flat_df.r,
3                               y = img_flat_df.g,
4                               z = img_flat_df.b,
5                               mode='markers',
6                               marker=dict(size=1,
7                                           color=[ 'rgb({}, {}, {})' .format(r,g,b) for
8                                               r,g,b in zip(img_flat_df.r.values,
9                                                         img_flat_df.g.values,
10                                                         img_flat_df.b.values)
11                                           opacity=0.5))
12
13         data = [trace]
14
15         layout = go.Layout(margin=dict(l=0,
16                                         r=0,
17                                         b=0,
18                                         t=0))
19
20         fig = go.Figure(data=data, layout=layout)
21
22         # Add centroid to chart
23         centroid = kmeans.cluster_centers_[0].tolist()
24
25         fig.add_trace(
26             go.Scatter3d(x = [centroid[0]],
27                           y = [centroid[1]],
28                           z = [centroid[2]],
29                           mode='markers',
30                           marker=dict(size=7,
31                                       color=[ 'rgb(125.79706706,78.8178776,42.580901
32                                       opacity=1))
33
34         )
35         fig.update_layout(scene = dict(
36                             xaxis_title='R',
37                             yaxis_title='G',
38                             zaxis_title='B'),
39                             )
40         fig.show()
```




We can see the centroid as a large circle in the middle of the colorspace. (If you can't, just click on the image and spin/zoom it.) Notice that this is the "center of gravity" of all the points in the graph.

Now let's try something else. Let's refit a K-means model to the data, this time using $k = 3$. Take a moment to consider what you might expect to result from this. Go through the steps of what the model is doing like we did above. What colors are you likely to see?

Cluster the data: $k = 3$

```
In [13]: 1 # Instantiate k-means model for 3 clusters
          2 kmeans3 = KMeans(n_clusters=3, random_state=42).fit(img_flat)
          3
          4 # Check the unique values of what's returned by the .labels_ attribute
          5 np.unique(kmeans3.labels_)
```

```
Out[13]: array([0, 1, 2])
```

The `.cluster_centers_` attribute returns an array where each element represents the coordinates of a centroid (i.e., their RGB values). We'll use these coordinates as we did previously to generate the colors that are represented by our centroids.

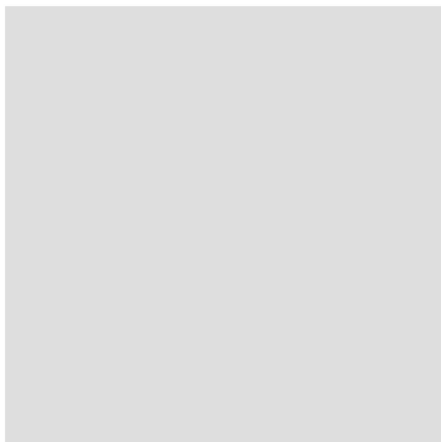
```
In [14]: 1 # Assign centroid coordinates to `centers` variable
          2 centers = kmeans3.cluster_centers_
          3 centers
```

```
Out[14]: array([[210.43027568, 103.15660618,  65.82105407],
                 [238.60707255, 224.70164054, 177.11527525],
                 [ 79.27359372,  42.74099099,  32.04084267]])
```

Now we'll create a helper function to easily display RGB values as color swatches, and use it to check the colors of the model's centroids.

```
In [15]: 1 # Helper function that creates color swatches
          2 def show_swatch(RGB_value):
          3     '''
          4     Takes in an RGB value and outputs a color swatch
          5     '''
          6     R, G, B = RGB_value
          7     rgb = [[np.array([R,G,B]).astype('uint8')]]
          8     plt.figure()
          9     plt.imshow(rgb)
         10     plt.axis('off');
```

```
In [16]: 1 # Display the color swatches  
2 for pixel in centers:  
3     show_swatch(pixel)
```



Hopefully, you hypothesized that we'd see similar colors as a result of a 3-cluster model. If you examine the original image of the tulips, it's apparent that there are generally three dominant colors: reds, greens, and golds/yellows, which is very close to what the model returned.

Just as before, let's now replace each pixel in the original image with the RGB value of the centroid to which it was assigned.

```

In [17]: 1 # Helper function to display our photograph when clustered into k clusters
2 def cluster_image(k, img=img):
3     '''
4     Fits a K-means model to a photograph.
5     Replaces photo's pixels with RGB values of model's centroids.
6     Displays the updated image.
7
8     Args:
9         k:      (int)          - Your selected K-value
10        img:    (numpy array) - Your original image converted to a numpy array
11
12    Returns:
13        The output of plt.imshow(new_img), where new_img is a new numpy array
14        where each row of the original array has been replaced with the \
15        coordinates of its nearest centroid.
16    '''
17
18    img_flat = img.reshape(img.shape[0]*img.shape[1], 3)
19    kmeans = KMeans(n_clusters = k, random_state = 42).fit(img_flat)
20    new_img = img_flat.copy()
21
22    for i in np.unique(kmeans.labels_):
23        new_img[kmeans.labels_ == i, :] = kmeans.cluster_centers_[i]
24
25    new_img = new_img.reshape(img.shape)
26
27    return plt.imshow(new_img), plt.axis('off');
28

```

Generate image when k=3

```
cluster_image(3);
```

We now have a photo with just three colors. Each pixel's RGB values correspond to the values of its nearest centroid.

We can return once more to our 3-D colorspace. This time, we'll re-color each dot in the colorspace to correspond with the color of its centroid. This will allow us to see how the K-means algorithm clustered our data spatially.

Again, don't concern yourself so much with the code. Feel free to skip down to the graph.

```
In [18]: 1 # Just to get an understanding of what the data structures look like
2
3 print(kmeans3.labels_.shape)
4 print(kmeans3.labels_)
5 print(np.unique(kmeans3.labels_))
6 print(kmeans3.cluster_centers_)
```

```
(129640,)
[1 1 1 ... 1 1 1]
[0 1 2]
[[210.43027568 103.15660618  65.82105407]
 [238.60707255 224.70164054 177.11527525]
 [ 79.27359372  42.74099099  32.04084267]]
```

```
In [19]: 1 # Create a new column in the df that indicates the cluster number of each
2 # (as assigned by Kmeans for k=3)
3 img_flat_df['cluster'] = kmeans3.labels_
4 img_flat_df.head()
```

Out[19]:

	r	g	b	cluster
0	220	189	125	1
1	220	186	122	1
2	220	185	117	1
3	225	189	115	1
4	230	193	113	1

```
In [20]: 1 # Create helper dictionary to map RGB color values to each observation in
2 series_conversion = {0: 'rgb' +str(tuple(kmeans3.cluster_centers_[0])),
3                       1: 'rgb' +str(tuple(kmeans3.cluster_centers_[1])),
4                       2: 'rgb' +str(tuple(kmeans3.cluster_centers_[2])),
5                       }
6 series_conversion
```

```
Out[20]: {0: 'rgb(210.43027567550294, 103.1566061821888, 65.82105407130598)',
1: 'rgb(238.60707254830987, 224.7016405395384, 177.11527524605327)',
2: 'rgb(79.27359371564982, 42.74099099098083, 32.04084267194157)'}
```

```
In [21]: 1 # Replace the cluster numbers in the 'cluster' col with formatted RGB values
          2 # (made ready for plotting)
          3 img_flat_df['cluster'] = img_flat_df['cluster'].map(series_conversion)
          4 img_flat_df.head()
```

Out[21]:

	r	g	b	cluster
0	220	189	125	rgb(238.60707254830987, 224.7016405395384, 177...
1	220	186	122	rgb(238.60707254830987, 224.7016405395384, 177...
2	220	185	117	rgb(238.60707254830987, 224.7016405395384, 177...
3	225	189	115	rgb(238.60707254830987, 224.7016405395384, 177...
4	230	193	113	rgb(238.60707254830987, 224.7016405395384, 177...

```
In [22]: 1 # Replot the data, now showing which cluster (i.e., color) it was assigned
2
3 trace = go.Scatter3d(x = img_flat_df.r,
4                     y = img_flat_df.g,
5                     z = img_flat_df.b,
6                     mode='markers',
7                     marker=dict(size=1,
8                                 color=img_flat_df.cluster,
9                                 opacity=1))
10
11 data = trace
12
13 layout = go.Layout(margin=dict(l=0,
14                                r=0,
15                                b=0,
16                                t=0))
17
18 fig = go.Figure(data=data, layout=layout)
19 fig.show()
```



You may be thinking to yourself that you would have clustered the data differently based on the distribution of points that you saw in the first 3-D plot. For example, why is there a sharp line that separates red and green, when there doesn't appear to be any empty space there in the data?

You're not incorrect. Even though there's no such thing as "wrong" clustering, some ways can definitely be better than others.

You'll notice in the original 3-D rendering that there are long bands—not round balls—of clustered data. K-means works best when the clusters are more circular, because it tries to minimize distance from point to centroid. It may be worth trying a different clustering algorithm if you want to cluster a long, narrow, continuous band of data. (More on these later!)

Nonetheless, K-means successfully compresses the colors of this photograph. This process can be applied for any value of k . Here's the output of each photo for $k = 2-10$.

Cluster the data: $k = 2-10$


```
In [23]: 1 # Helper function to
2 def cluster_image_grid(k, ax, img=img):
3     '''
4     Fits a K-means model to a photograph.
5     Replaces photo's pixels with RGB values of model's centroids.
6     Displays the updated image on an axis of a figure.
7
8     Args:
9         k:      (int)          - Your selected K-value
10        ax:      (int)          - Index of the axis of the figure to plot to
11        img:      (numpy array) - Your original image converted to a numpy array
12
13    Returns:
14        A new image where each row of the ori array has been replaced with the
15        coordinates of its nearest centroid.
16    '''
17    img_flat = img.reshape(img.shape[0]*img.shape[1], 3)
18    kmeans = KMeans(n_clusters=k, random_state=42).fit(img_flat)
19    new_img = img_flat.copy()
20
21    for i in np.unique(kmeans.labels_):
22        new_img[kmeans.labels_==i, :] = kmeans.cluster_centers_[i]
23
24    new_img = new_img.reshape(img.shape)
25    ax.imshow(new_img)
26    ax.axis('off')
27
28    fig, axs = plt.subplots(3, 3)
29    fig = matplotlib.pyplot.gcf()
30    fig.set_size_inches(9, 12)
31    axs = axs.flatten()
32    k_values = np.arange(2, 11)
33    for i, k in enumerate(k_values):
34        cluster_image_grid(k, axs[i], img=img)
35        axs[i].title.set_text('k=' + str(k))
```

k=2



k=3



k=4



k=5



k=6



k=7



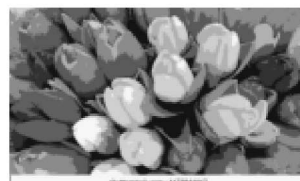
k=8



k=9



k=10



Notice that it becomes increasingly difficult to see the difference between the images each time a color is added. This is a visual example of something that happens with all clustering models, even if the data is not an image that you can see. As you group the data into more and more clusters, additional clusters beyond a certain point contribute less and less to your understanding of your data.

In []:

1