

Kmeans on Geyser's Eruptions Segmentation: We'll first implement the kmeans algorithm on 2D data with 272 observations and 2 features. The data covers the waiting time between eruptions and the duration of the eruption in Yellowstone National Park, Wyoming, USA. We will try to find K subgroups within the data.

eruptions (float): Eruption time in minutes. waiting (int): Waiting time to next eruption.

```

1 # Modules
2 import matplotlib.pyplot as plt
3 from matplotlib.image import imread
4 import pandas as pd
5 import seaborn as sns
6 from sklearn.datasets.samples_generator import (make_blobs,
7                                                make_circles,
8                                                make_moons)
9 from sklearn.cluster import KMeans, SpectralClustering
10 from sklearn.preprocessing import StandardScaler
11 from sklearn.metrics import silhouette_samples, silhouette_score
12

```

```

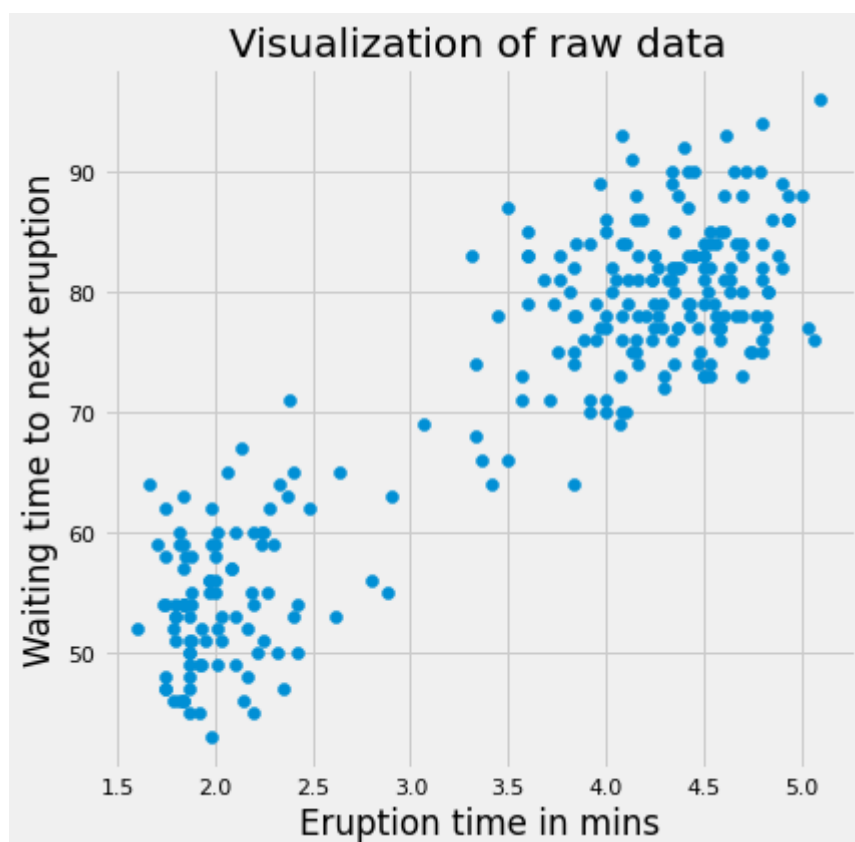
1 # Import the data
2 df = pd.read_csv('old_faithful.csv')

```

```

1 # Plot the data
2 plt.figure(figsize=(6, 6))
3 plt.scatter(df.iloc[:, 0], df.iloc[:, 1])
4 plt.xlabel('Eruption time in mins')
5 plt.ylabel('Waiting time to next eruption')
6 plt.title('Visualization of raw data');

```



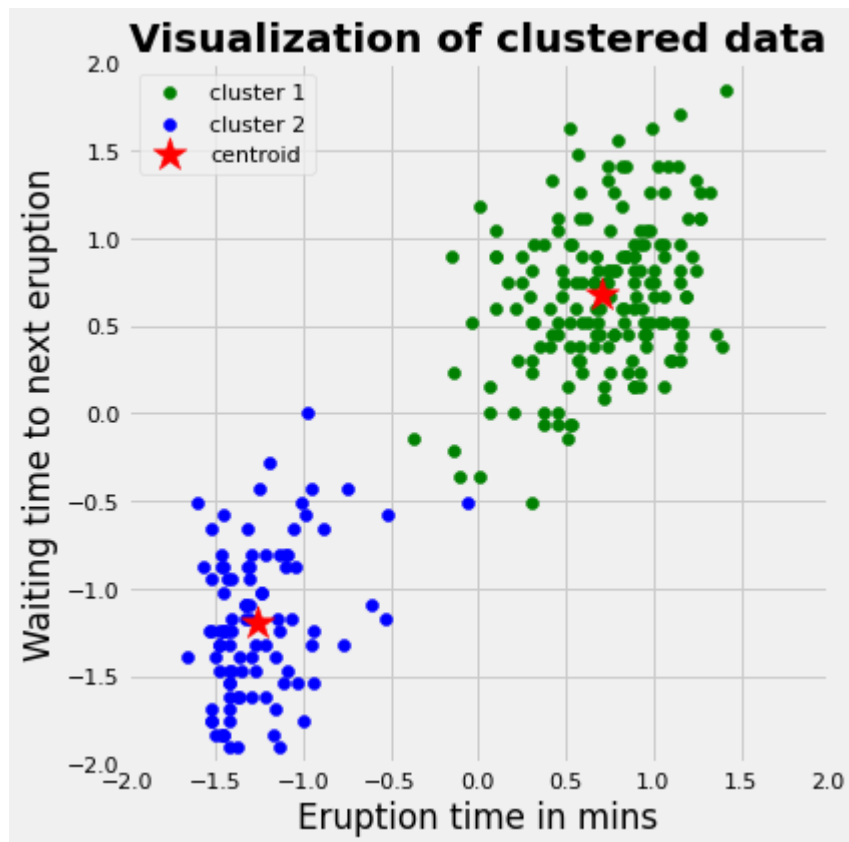
We'll use this data because it's easy to plot and visually spot the clusters since its a 2-dimension data. Let's standardize the data first and run the kmeans algorithm on the standardized data with K=2.

```
1
2 # Standardize the data
3 X_std = StandardScaler().fit_transform(df)
4

1
2 # Run local implementation of kmeans
3 km = Kmeans(n_clusters=2, max_iter=100)
4 km.fit(X_std)
5 centroids = km.centroids
6

1 # Plot the clustered data
2 fig, ax = plt.subplots(figsize=(6, 6))
3 plt.scatter(X_std[km.labels == 0, 0], X_std[km.labels == 0, 1],
4             c='green', label='cluster 1')
5 plt.scatter(X_std[km.labels == 1, 0], X_std[km.labels == 1, 1],
6             c='blue', label='cluster 2')
7 plt.scatter(centroids[:, 0], centroids[:, 1], marker='*', s=300,
8             c='r', label='centroid')
9 plt.legend()
10 plt.xlim([-2, 2])
11 plt.ylim([-2, 2])
12 plt.xlabel('Eruption time in mins')
13 plt.ylabel('Waiting time to next eruption')
14 plt.title('Visualization of clustered data', fontweight='bold')
15 ax.set_aspect('equal');
```



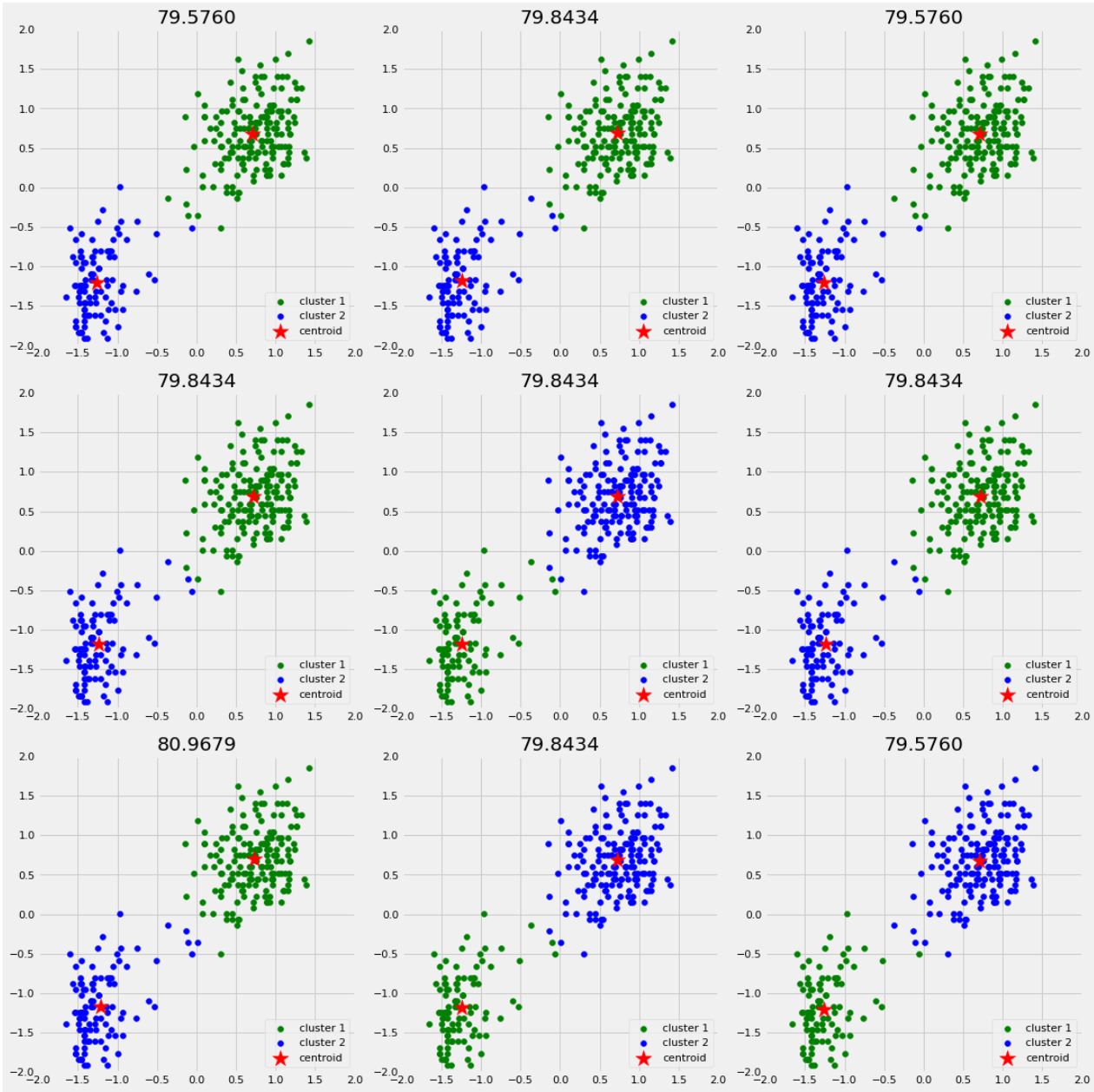


```

1
2 n_iter = 9
3 fig, ax = plt.subplots(3, 3, figsize=(16, 16))
4 ax = np.ravel(ax)
5 centers = []
6 for i in range(n_iter):
7     # Run local implementation of kmeans
8     km = Kmeans(n_clusters=2,
9                 max_iter=3,
10                 random_state=np.random.randint(0, 1000, size=1))
11     km.fit(X_std)
12     centroids = km.centroids
13     centers.append(centroids)
14     ax[i].scatter(X_std[km.labels == 0, 0], X_std[km.labels == 0, 1],
15                  c='green', label='cluster 1')
16     ax[i].scatter(X_std[km.labels == 1, 0], X_std[km.labels == 1, 1],
17                  c='blue', label='cluster 2')
18     ax[i].scatter(centroids[:, 0], centroids[:, 1],
19                  c='r', marker='*', s=300, label='centroid')
20     ax[i].set_xlim([-2, 2])
21     ax[i].set_ylim([-2, 2])
22     ax[i].legend(loc='lower right')
23     ax[i].set_title(f'{km.error:.4f}')
24     ax[i].set_aspect('equal')
25 plt.tight_layout();
26

```





As the graph above shows that we only ended up with two different ways of clusterings based on  $c$  with the lowest sum of squared distance.

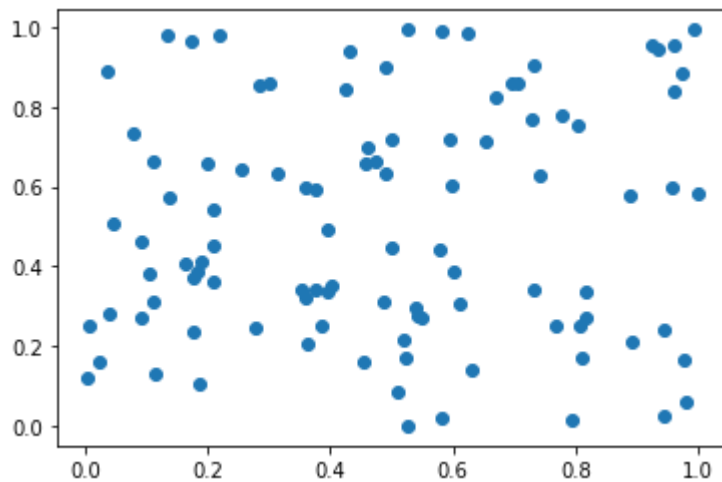
## Analysis

```
1 import matplotlib.pyplot as plt
2 %matplotlib inline
3 import numpy as np
4 from sklearn.cluster import KMeans
5 import time
```

```
1 np.random.seed(100)
2 X = np.random.rand(100, 2)
3 X_test = np.random.rand(10, 2)
```

```
1 plt.scatter(X[:,0],X[:,1], label='True Position')
```

 <matplotlib.collections.PathCollection at 0x7f5fc9e81898>



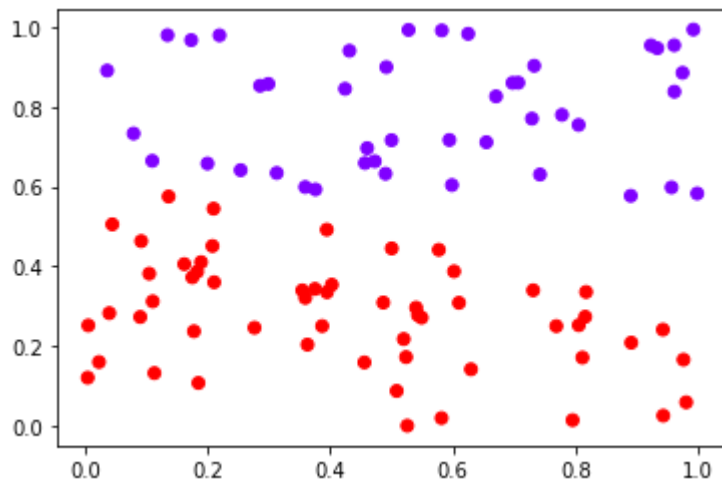
Number of Cluster = 2

```
1 start_time = time.time()
2 kmeans = KMeans(n_clusters=2, random_state=3425)
3 kmeans.fit(X)
4 end_time = time.time()
5 print("Time Taken : "+str(end_time-start_time))
```

 Time Taken : 0.030069828033447266

```
1 plt.scatter(X[:,0],X[:,1], c=kmeans.labels_, cmap='rainbow')
```

 <matplotlib.collections.PathCollection at 0x7f5fc6ec0ef0>

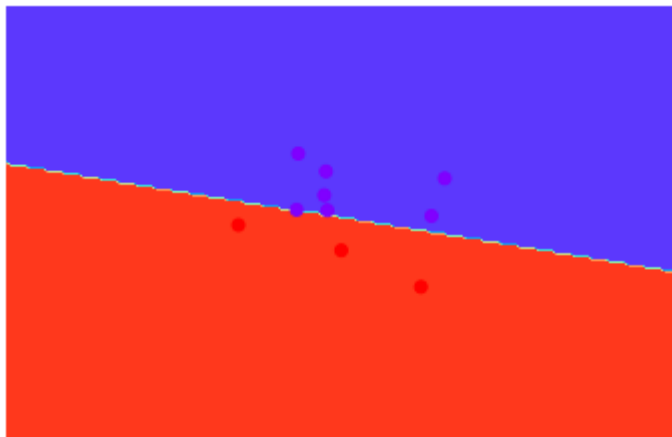


```

1 x_min, x_max = X[:, 0].min() - 1, X[:, 0].max() + 1
2 y_min, y_max = X[:, 1].min() - 1, X[:, 1].max() + 1
3 xx, yy = np.meshgrid(np.arange(x_min, x_max, 0.02),
4                       np.arange(y_min, y_max, 0.02))
5
6 # here "model" is your model's prediction (classification) function
7 Z = kmeans.predict(np.c_[xx.ravel(), yy.ravel()])
8
9 # Put the result into a color plot
10 Z = Z.reshape(xx.shape)
11 plt.contourf(xx, yy, Z, cmap='rainbow')
12 plt.axis('off')
13
14 # Plot also the training points
15 plt.scatter(X_test[:,0],X_test[:,1], c=kmeans.predict(X_test), cmap='rainbow')

```

 <matplotlib.collections.PathCollection at 0x7f5fc6eacf60>



```
1 centroid = kmeans.predict(X_test)
```

```

2 sse = 0
3 for i in range(len(X_test)):
4     temp = np.linalg.norm(X_test[i]-kmeans.cluster_centers_[centroid[i]])
5     temp = temp*temp
6     sse += temp
7 print(sse)

```

1.0645816624472537

### Similarly Doing for n\_clusters=3

```

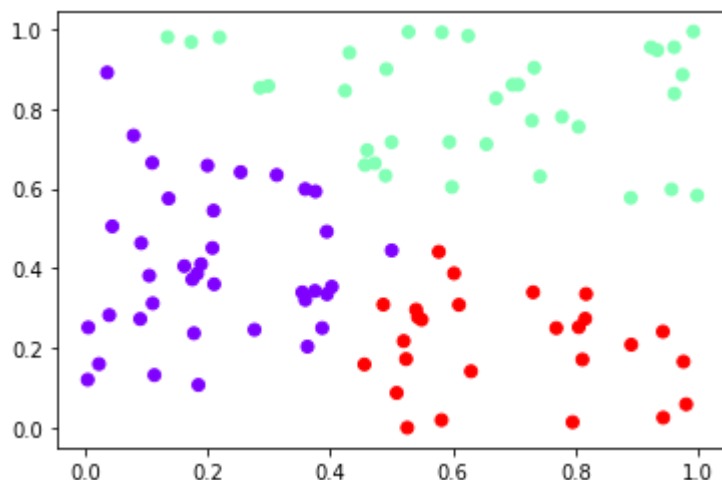
1 start_time = time.time()
2 kmeans = KMeans(n_clusters=3, random_state=3425)
3 kmeans.fit(X)
4 end_time = time.time()
5 print("Time Taken : "+str(end_time-start_time))

```

Time Taken : 0.033385276794433594

```
1 plt.scatter(X[:,0],X[:,1], c=kmeans.labels_, cmap='rainbow')
```

<matplotlib.collections.PathCollection at 0x7f5fc6e0ef60>

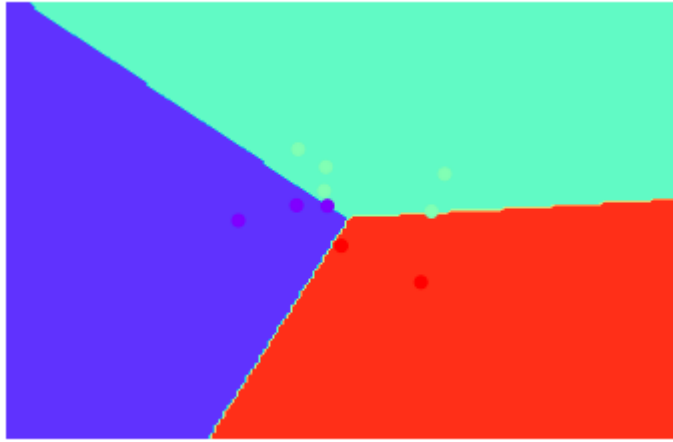


```

1 x_min, x_max = X[:, 0].min() - 1, X[:, 0].max() + 1
2 y_min, y_max = X[:, 1].min() - 1, X[:, 1].max() + 1
3 xx, yy = np.meshgrid(np.arange(x_min, x_max, 0.02),
4                       np.arange(y_min, y_max, 0.02))
5
6 # here "model" is your model's prediction (classification) function
7 Z = kmeans.predict(np.c_[xx.ravel(), yy.ravel()])
8
9 # Put the result into a color plot
10 Z = Z.reshape(xx.shape)
11 plt.contourf(xx, yy, Z, cmap='rainbow')
12 plt.axis('off')
13
14 # Plot also the training points
15 plt.scatter(X_test[:,0],X_test[:,1], c=kmeans.predict(X_test), cmap='rainbow')

```

↳ <matplotlib.collections.PathCollection at 0x7f5fc6df7b38>



```
1 centroid = kmeans.predict(X_test)
2 sse = 0
3 for i in range(len(X_test)):
4     temp = np.linalg.norm(X_test[i]-kmeans.cluster_centers_[centroid[i]])
5     temp = temp*temp
6     sse += temp
7 print(sse)
```

↳ 0.7406067549399893

Similarly Doing for n\_clusters=4

```
1 start_time = time.time()
2 kmeans = KMeans(n_clusters=4, random_state=3425)
3 kmeans.fit(X)
4 end_time = time.time()
5 print("Time Taken : "+str(end_time-start_time))
```

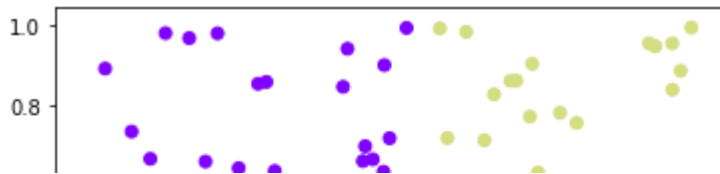
↳ Time Taken : 0.03960251808166504

```
1 plt.scatter(X[:,0],X[:,1], c=kmeans.labels_, cmap='rainbow')
```

↳



<matplotlib.collections.PathCollection at 0x7f5fc6d5be48>

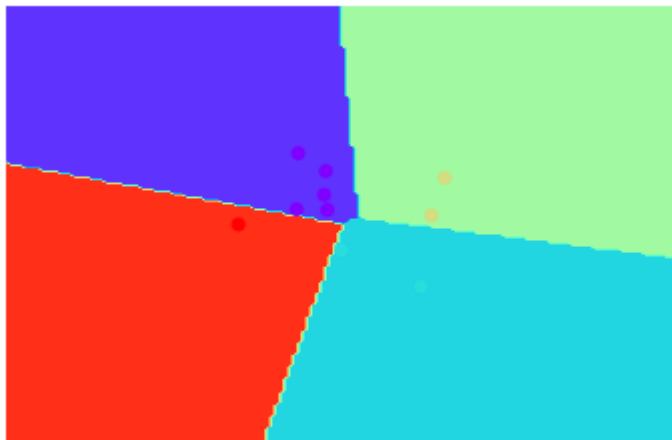


```

1 x_min, x_max = X[:, 0].min() - 1, X[:, 0].max() + 1
2 y_min, y_max = X[:, 1].min() - 1, X[:, 1].max() + 1
3 xx, yy = np.meshgrid(np.arange(x_min, x_max, 0.02),
4                       np.arange(y_min, y_max, 0.02))
5
6 # here "model" is your model's prediction (classification) function
7 Z = kmeans.predict(np.c_[xx.ravel(), yy.ravel()])
8
9 # Put the result into a color plot
10 Z = Z.reshape(xx.shape)
11 plt.contourf(xx, yy, Z, cmap='rainbow')
12 plt.axis('off')
13
14 # Plot also the training points
15 plt.scatter(X_test[:,0],X_test[:,1], c=kmeans.predict(X_test), cmap='rainbow')

```

↳ <matplotlib.collections.PathCollection at 0x7f5fc6cbf9e8>



```

1 centroid = kmeans.predict(X_test)
2 sse = 0
3 for i in range(len(X_test)):
4     temp = np.linalg.norm(X_test[i]-kmeans.cluster_centers_[centroid[i]])
5     temp = temp*temp
6     sse += temp
7 print(sse)

```

↳ 0.39793603053158416

Similarly Doing for n\_clusters=5

```

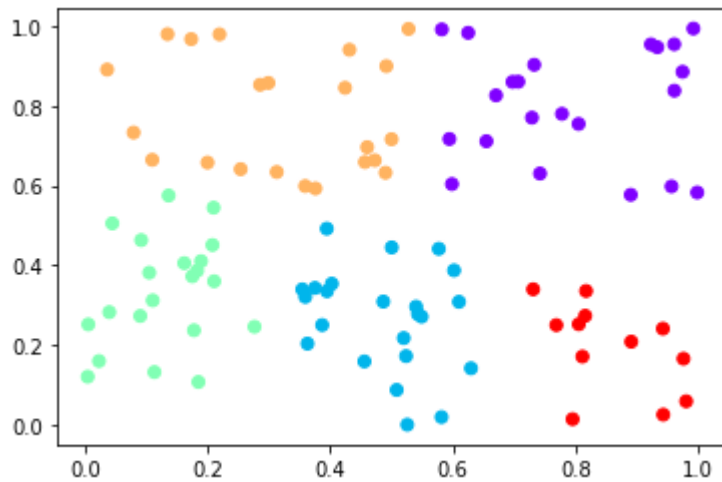
1 start_time = time.time()
2 kmeans = KMeans(n_clusters=5, random_state=3425)
3 kmeans.fit(X)
4 end_time = time.time()
5 print("Time Taken : "+str(end_time-start_time))

```

Time Taken : 0.04885458946228027

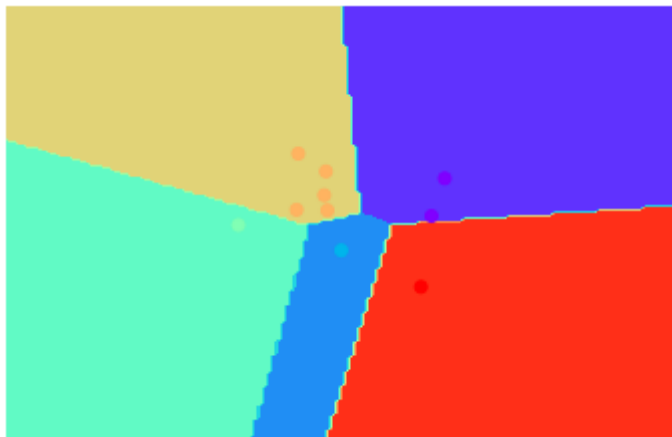
```
1 plt.scatter(X[:,0],X[:,1], c=kmeans.labels_, cmap='rainbow')
```

<matplotlib.collections.PathCollection at 0x7f5fc6bb7d30>



```
1 x_min, x_max = X[:, 0].min() - 1, X[:, 0].max() + 1
2 y_min, y_max = X[:, 1].min() - 1, X[:, 1].max() + 1
3 xx, yy = np.meshgrid(np.arange(x_min, x_max, 0.02),
4                       np.arange(y_min, y_max, 0.02))
5
6 # here "model" is your model's prediction (classification) function
7 Z = kmeans.predict(np.c_[xx.ravel(), yy.ravel()])
8
9 # Put the result into a color plot
10 Z = Z.reshape(xx.shape)
11 plt.contourf(xx, yy, Z, cmap='rainbow')
12 plt.axis('off')
13
14 # Plot also the training points
15 plt.scatter(X_test[:,0],X_test[:,1], c=kmeans.predict(X_test), cmap='rainbow')
```

<matplotlib.collections.PathCollection at 0x7f5fc6af8b00>



```
1 centroid = kmeans.predict(X_test)
```

```

2 sse = 0
3 for i in range(len(X_test)):
4     temp = np.linalg.norm(X_test[i]-kmeans.cluster_centers_[centroid[i]])
5     temp = temp*temp
6     sse += temp
7 print(sse)

```

0.30321005588436034

### Kmeans on Image Compression:

In this part, we'll implement kmeans to compress an image. The image that we'll be working on is located at `./img/IMG_6115.jpeg`. At each location we would have 3 8-bit integers that specify the red, green, and blue intensity values. Our goal is to compress and represent (compress) the photo using those 30 colors only. To pick which colors to use, we'll use K-Means. We'll treat every pixel as a data point. That means reshape the image from height x width x channels to (height x width) x channels.  $128 \times 128 \times 3 = 49152$  data points in 3-dimensional space which are the intensity of RGB. Doing so will allow us to find centroids for each cluster and would significantly reduce the size of the image by a factor of 6. The original image is 3,763,584 bits; however, the new compressed image would be  $30 \times 24 + 128 \times 128 \times 4 = 627,984$  bits that we'll be using centroids as a lookup for pixels' colors and that would reduce the size of each pixel to 24 bits.

```

1
2 # Read the image
3 img = imread('IMG_6115.jpeg')
4 img_size = img.shape
5
6 # Reshape it to be 2-dimension
7 X = img.reshape(img_size[0] * img_size[1], img_size[2])
8
9 # Run the Kmeans algorithm
10 km = KMeans(n_clusters=30)
11 km.fit(X)
12
13 # Use the centroids to compress the image
14 X_compressed = km.cluster_centers_[km.labels_]
15 X_compressed = np.clip(X_compressed.astype('uint8'), 0, 255)
16
17 # Reshape X_recovered to have the same dimension as the original image 128 * 128 * 3
18 X_compressed = X_compressed.reshape(img_size[0], img_size[1], img_size[2])
19
20 # Plot the original and the compressed image next to each other
21 fig, ax = plt.subplots(1, 2, figsize = (12, 8))
22 ax[0].imshow(img)
23 ax[0].set_title('Original Image',)
24 ax[1].imshow(X_compressed)
25 ax[1].set_title('Compressed Image with 30 colors')
26 for ax in fig.axes:
27     ax.axis('off')
28 plt.tight_layout()
29

```

