

Clustering

Clustering is a type of unsupervised learning in which we try to find the features that are similar and group them. This group is referred to as a cluster. We mostly use L2 distance as a metric to cluster the features.

Kmeans and Kmeans++

Kmeans

Kmeans is one of the simplest and most popular clustering algorithm that uses distance between feature points to create clusters. Normal kmeans clustering randomly assigns K points in feature space as cluster centers. Then i-th cluster is assigned a point if, the distance between the point and cluster i is minimum as compared to all other cluster centers.

Once all points are clustered, mean of each cluster is calculated which becomes new cluster centers. This process is repeated till new cluster and cluster in previous iteration are almost similar to each other or certain number of iterations has been completed.

Kmeans++

Kmeans++ is similar to Kmeans with the only difference in the procedure to initialize starting cluster centers. In Kmeans++ instead of randomly initializing K cluster centers, we randomly initialize 1 cluster center. To initialize rest of the cluster centers following steps are performed

- First choose minimum distance of all the points with all the cluster centers.
- From the chosen minimum distances we choose the point with the maximum distance as the another cluster center.
- Repeat above two steps till we have acquired K cluster centers.

Kmeans++ ensures more intelligent initialization of cluster centers hence leads to better results as compared to kmeans. Below experiment justifies this statement



In the above result we can see that kmeans is incorrectly clustering the data points, whereas clusters created by kmeans++ looks correct.

Applications-

Kmeans++ has various applications like identifying crime localities, insurance fraud detection, customer segmentation etc. In this document we will look into following three application:

- Breast Cancer Detection
- Image compression
- Facial feature detection

Breast cancer detection

In below experiment we use breast cancer dataset present in sklearn to categorize feature into benign and malignant cancer type

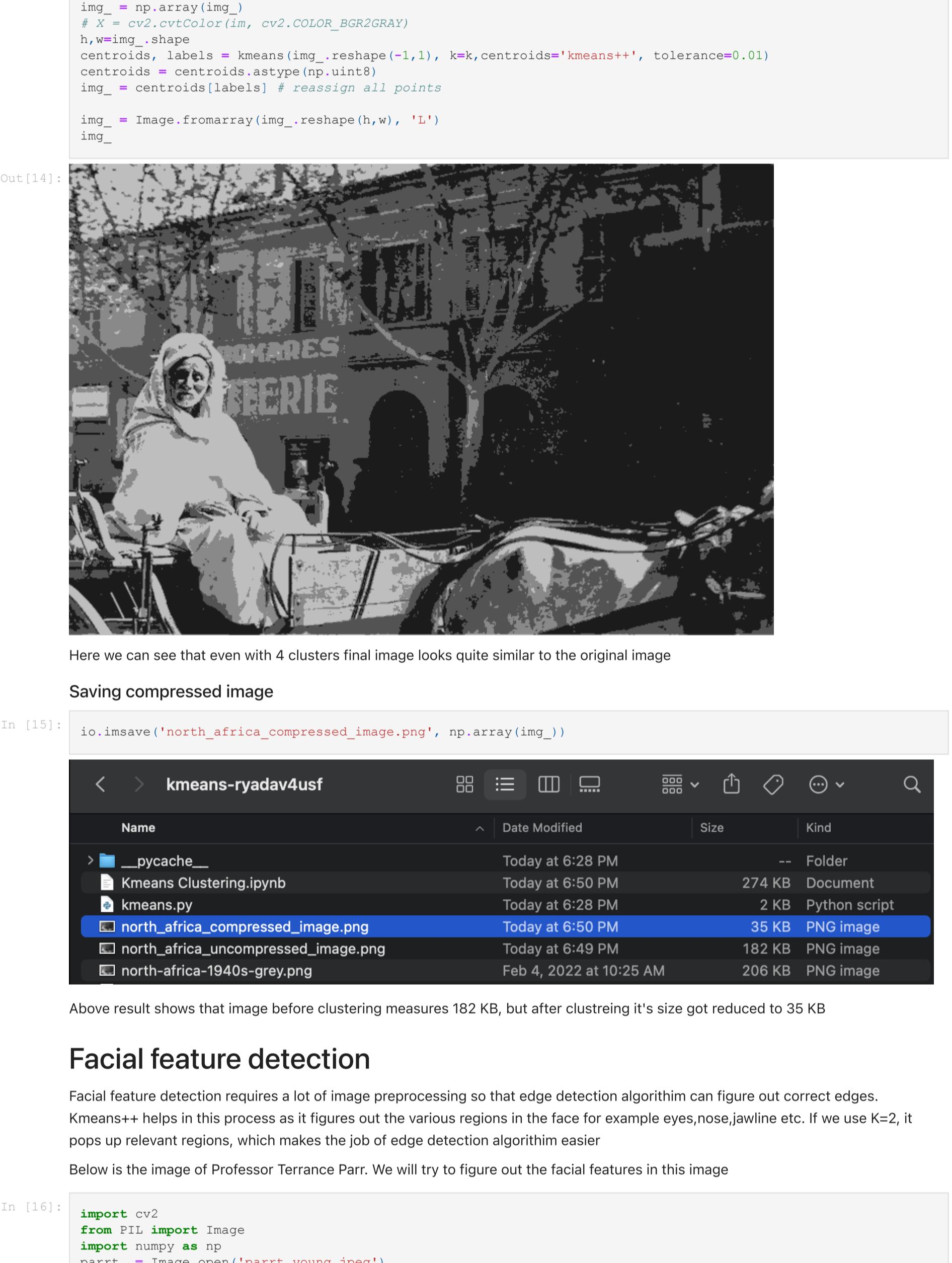
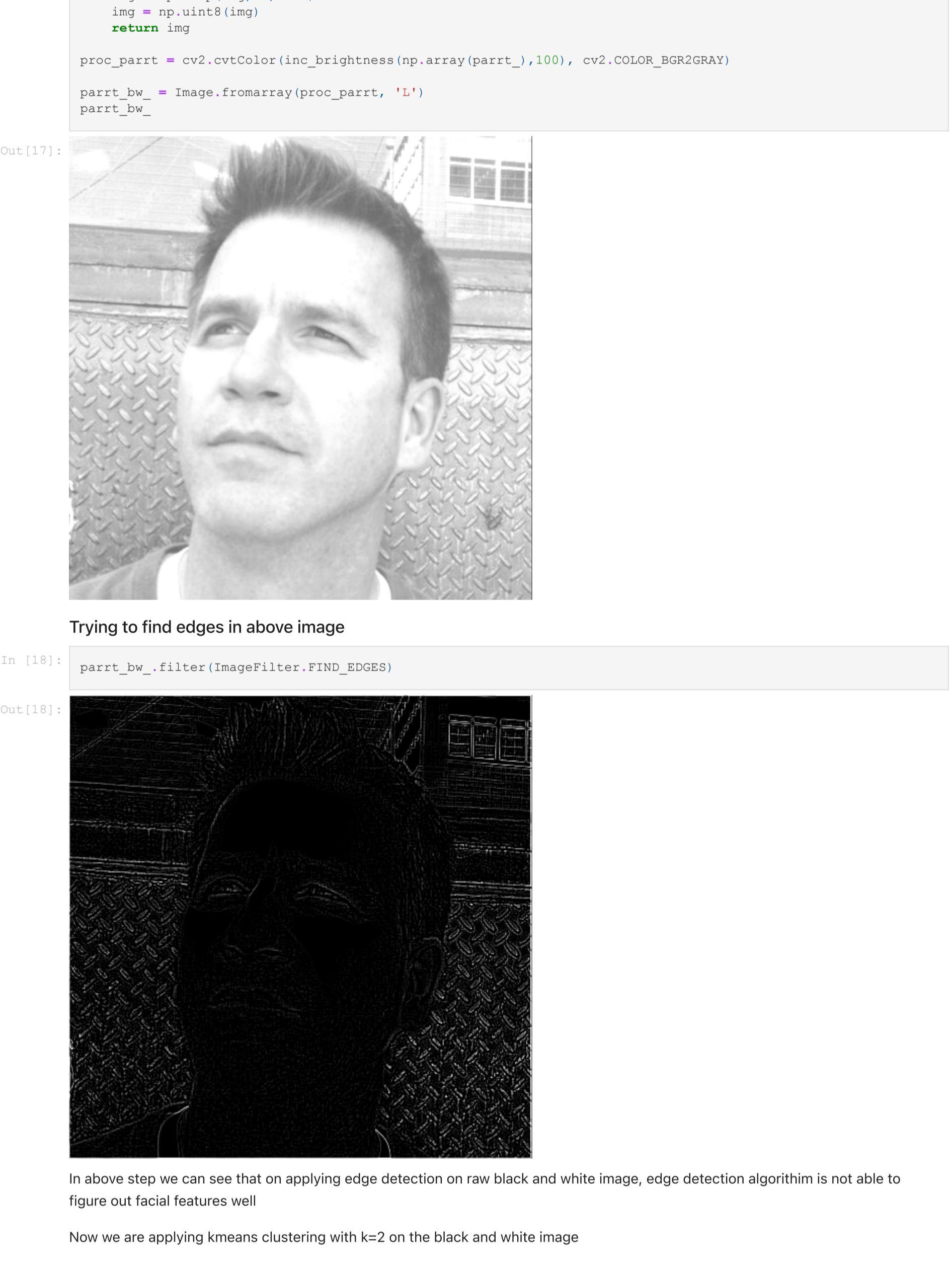


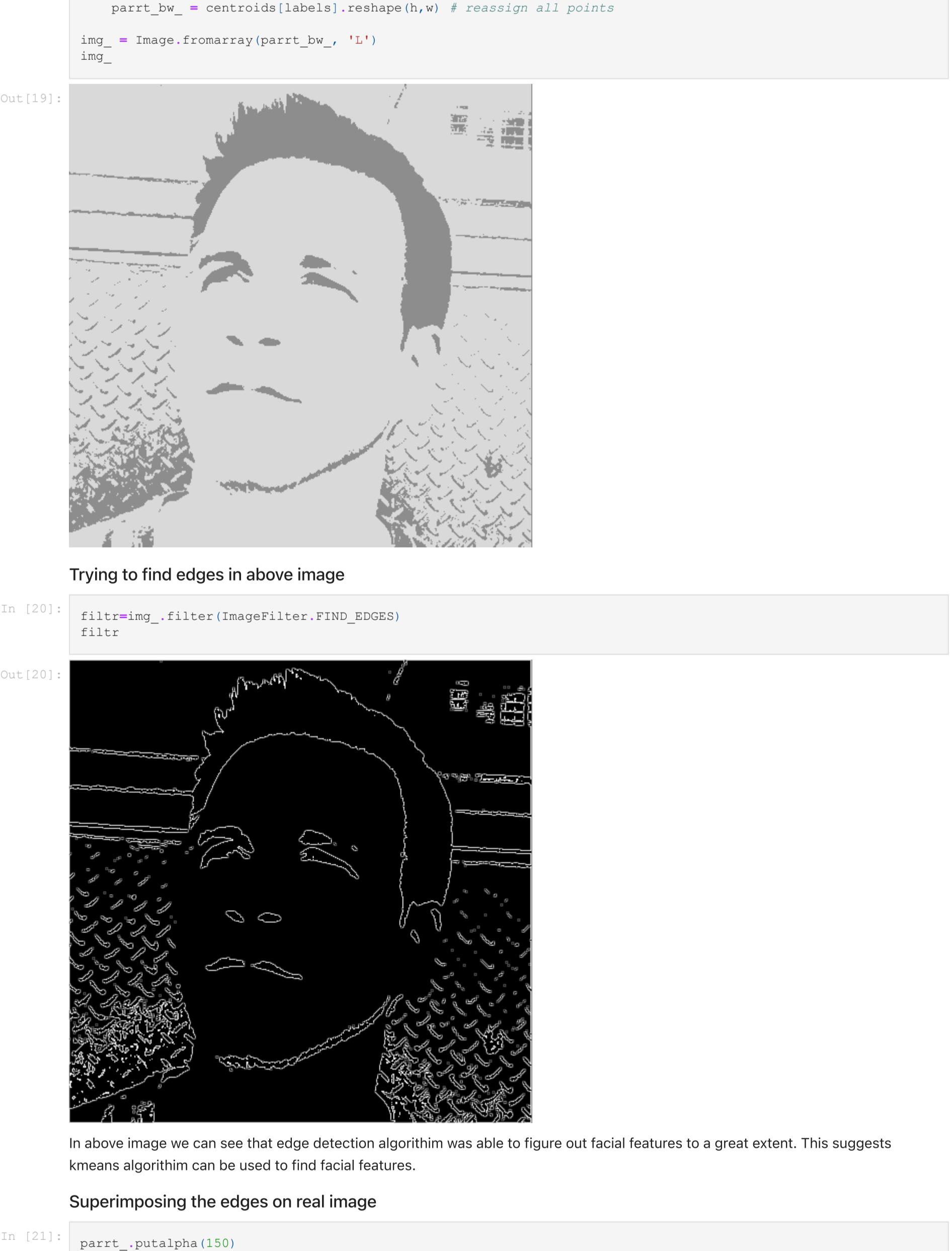
Image Compression

Kmean++ clustering is a good way to compress an image without much compromising on the quality of image. In the below example we are trying to cluster an image on just 4 clusters and compare the size of initial image and image after clustering

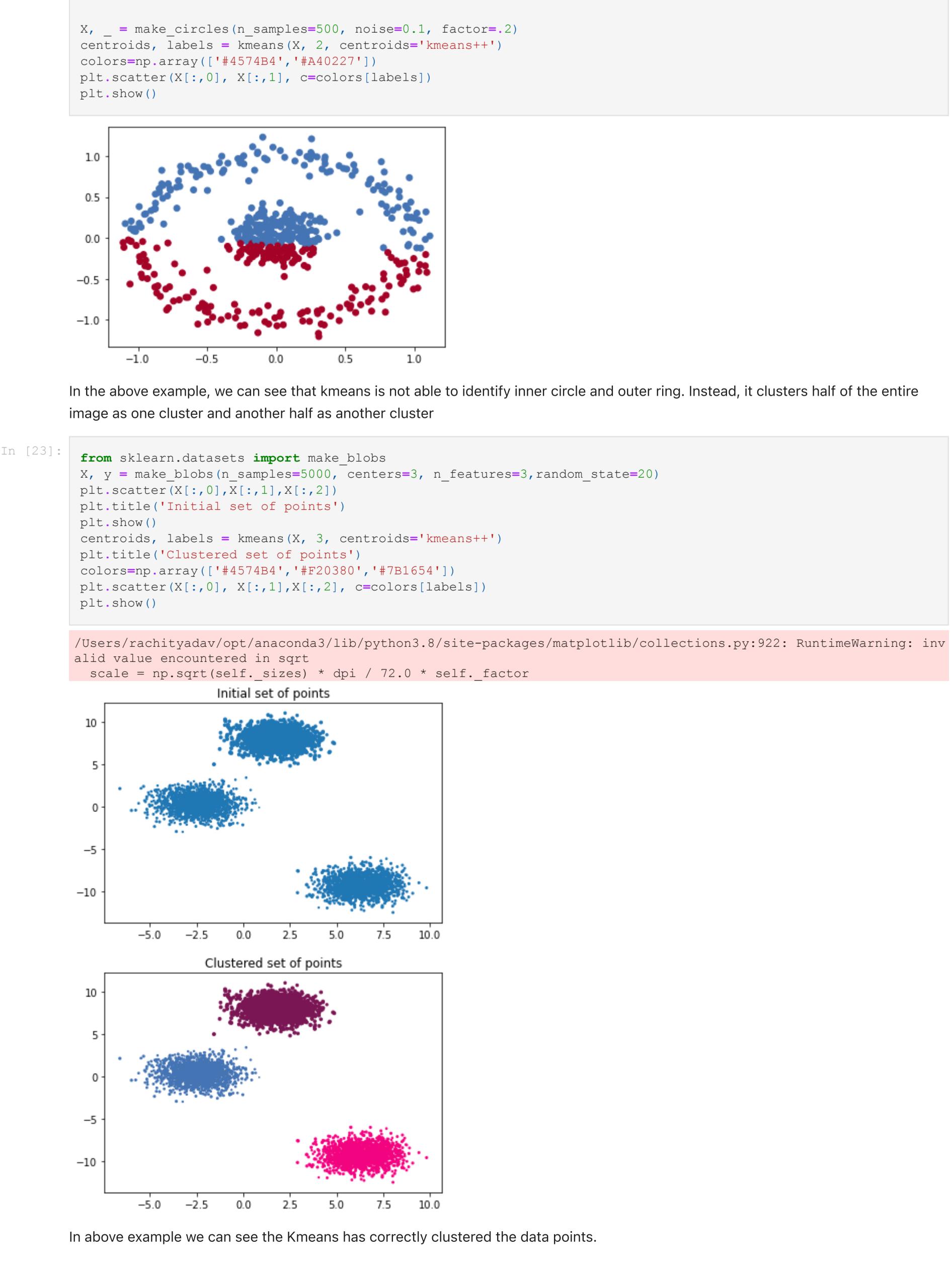
Since we are using only 4 clusters final image might not be as sharp as initial image, but on increasing number of clusters, we can achieve sharper image at the cost of image size



Saving compressed image



Converting the image to gray scale and increasing the brightness



Limitation of Kmeans:

Even though Kmeans has lot of applications in practical life, it also has few limitations

a. Not good with huge data - K-means requires distances to be calculated from cluster center to each and every point. This process becomes computationally expensive when we are dealing with huge amount of data. Hence Kmeans is preferred only with small amount of data

b. Works with numeric data - Since Kmeans uses distance between points in feature space, we generally don't get good result with categorical features. It's example is demonstrated below. In this example we try to compare clustering result in an image with non spherical clusters vs an image with spherical clusters

This suggests performance of Kmeans is best when the shape of cluster is kind of spherical

Conclusion

Kmeans is one of the most popular clustering algorithm and is usually the first choice while implementing clustering. It initially chooses random centroids then iteratively shifts them and eventually arrives at a set of centroids that are almost the same location as the centroid in previous iteration. All features that are closest to a specific centroid forms a cluster

Kmeans works well only when we have kind of spherical clusters in the data. If the cluster present is not spherical then kmeans does not give correct result. Also, kmeans requires number of clusters to be pre defined. Therefore, we can say that kmeans clustering doesn't work in all scenarios, but in the scenarios where it does work, kmeans gives really good results.

