

Assignment-6: K-Means Clustering

Rezwan-Ul-Alam (ID: 2011659042)
Md. Nur Alam Jowel (ID: 2012355042)
Raian Ruku (ID: 2013409642)

Email addresses:

rezwan.alam1@northsouth.edu
alam.jowel@northsouth.edu
raian.ruku@northsouth.edu

I. INTRODUCTION

K-Means clustering is a popular unsupervised machine learning algorithm used to partition a dataset into K clusters. Each cluster is represented by its centroid, and the objective is to minimize the variance within each cluster. In this assignment, we explore the application of K-Means clustering on a 2D dataset and its use in image compression. We will discuss the methodology, evaluation, results, and discussions derived from our experiments.

II. METHOD

The methodology for this assignment involves two main parts:

1. Clustering a 2D dataset.
2. Using K-Means for image compression.

III. CLUSTERING A 2D DATASET

At first, we loaded a 2D dataset from a .npz file and applied the K-Means algorithm to cluster the data points. Here are the steps below:

A. Data Loading

- a. First, we load the 2D data from a file. This data consists of points in a 2-dimensional space. Here is the code:

```
1 X = np.load("/content/drive/MyDrive/data/kmeans2d.npz")
2 print("First five elements of X are:\n", X[:5])
3 print('The shape of X is:', X.shape)
```

- b. The first five elements and the shape of the dataset are printed to understand its structure and size.

```
1 First five elements of X are:
2 [[1.84207953  4.6075716 ]
3  [5.65858312  4.79996405]
4  [6.35257892  3.2908545 ]
5  [2.90401653  4.61220411]
6  [3.23197916  4.93989405]]
7 The shape of X is: (300, 2)
```

B. Initialization

a. Randomly select initial centroids from the data points. This step is crucial as the choice of initial centroids can influence the final clusters. b. The function `kmeans_initial_centroids` randomly permute the dataset indices and selects the first K points

as the initial centroids. The code of the function is given below:

```

1 def kmeans_initial_centroids(X, K):
2     randidx = np.random.permutation(X.shape[0])
3     centroids = X[randidx[:K]]
4     return centroids

```

IV. ITERATION

1. Execute the K-Means algorithm for a fixed number of iterations.

Here is the code:

```

1 def find_closest_centroids(X, centroids):
2     m = centroids.shape[0]
3     n = X.shape[0]
4     idx = np.zeros(X.shape[0], dtype=int)
5     for i in range(n):
6         distance = []
7         for j in range(m):
8             norm = np.linalg.norm(X[i] - centroids[j])
9             distance.append(norm)
10        idx[i] = np.argmin(distance)
11    return idx
12
13 def compute_centroids(X, idx, K):
14     m, n = X.shape
15     centroids = np.zeros((K, n))
16     for k in range(K):
17         points = X[idx == k]
18         centroids[k] = np.mean(points, axis=0)
19    return centroids
20
21 def kmeans(X, initial_centroids, max_itr):
22     m, n = X.shape
23     K = initial_centroids.shape[0]
24     centroids = initial_centroids
25     previous_centroids = centroids
26     idx = np.zeros(m)
27     for i in range(max_itr):
28         idx = find_closest_centroids(X, centroids)
29         centroids = compute_centroids(X, idx, K)
30    return centroids, idx

```

2. find_closest_centroids: For each data point, this function calculates the distance to each centroid and assigns the point to the nearest centroid. Here below the code:

```

1 def find_closest_centroids(X, centroids):
2
3     m = centroids.shape[0]
4     n = X.shape[0]
5     idx = np.zeros(X.shape[0], dtype=int)
6     for i in range(n):
7         distance = []
8         for j in range(m):
9             norm = np.linalg.norm(X[i] - centroids[j])
10            distance.append(norm)
11
12        idx[i] = np.argmin(distance)
13    return idx

```

3. compute_centroids: This function updates the centroids by calculating the mean position of all points assigned to each centroid. Here below the code:

```

1 compute_centroids(X, idx, K):
2     m, n = X.shape
3     centroids = np.zeros((K, n))
4
5     for k in range(K):
6         points = X[idx == k]
7         centroids[k] = np.mean(points, axis = 0)
8     return centroids

```

4. kmeans: This function repeatedly assigns points to the nearest centroid and recalculates centroids for a set number of iterations. Here below the code:

```

1 def kmeans(X, initial_centroids, max_itr):
2
3     m, n = X.shape
4     K = initial_centroids.shape[0]
5     centroids = initial_centroids
6     previous_centroids = centroids
7     idx = np.zeros(m)
8     for i in range(max_itr):
9         idx = find_closest_centroids(X, centroids)
10        centroids = compute_centroids(X, idx, K)
11    return centroids, idx

```

V. VISUALIZATION

- Plot the clustered data with centroids for different values of K.
- We use matplotlib to create scatter plots of the clustered data points for different numbers of clusters (K=3 to K=6). Each plot shows the data points colored by their assigned cluster and the centroids marked with red 'x' symbols.

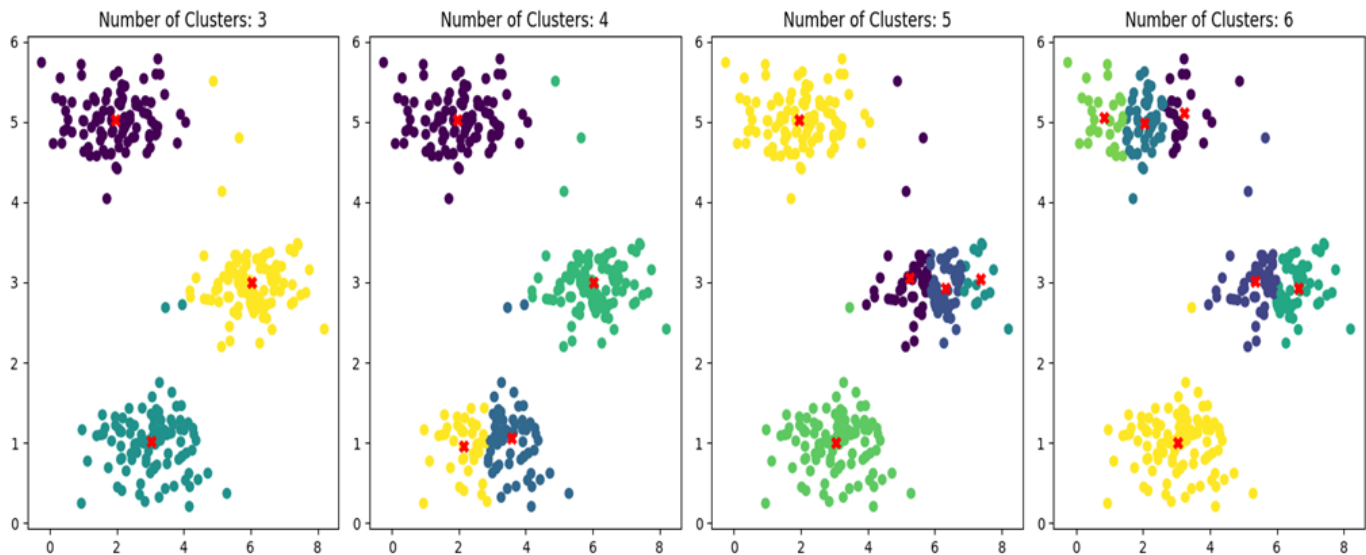
Here is the code:

```

1 maximum_clusters = 6
2 minimum_clusters = 3
3
4 fig, axs = plt.subplots(1, maximum_clusters - 2, figsize=(15, 5))
5 max_iteration = 10
6 #different numbers of clusters
7 for idx, clusters in enumerate(range(minimum_clusters, maximum_clusters + 1)):
8     initial_centroid = X[np.random.choice(range(len(X)), size=clusters)]
9     centroids, labels = kmeans(X, initial_centroid, max_iteration)
10    axs[idx].scatter(X[:, 0], X[:, 1], c=labels)
11    axs[idx].scatter(centroids[:, 0], centroids[:, 1], marker='x', c='r', linewidths=3)
12    axs[idx].set_title(f"Number_of_Clusters:_{clusters}")
13
14 plt.tight_layout()
15 plt.show()

```

Here is the visualization output:



VI. IMAGE COMPRESSION

Image compression using K-Means involves treating each pixel as a data point in a 3D space (RGB color space). Here below the steps:

A. Data Preparation

- Load and preprocess the image by normalizing pixel values to the range $[0, 1]$.
- The image is loaded and displayed to understand its structure. We normalize the pixel values and reshape the image into a matrix where each row represents a pixel's RGB values.

Here is the code below:

```

1 # Divide by 255 so that all values are in the range 0 - 1
2 original_img = original_img / 255
3
4 # Reshape the image into an m x 3 matrix where m = number of pixels
5 # (in this case m = 128 x 128 = 16384)
6 # Each row will contain the Red, Green and Blue pixel values
7 # This gives us our dataset matrix X_img that we will use K-Means on.
8
9 X_img = np.reshape(original_img, (original_img.shape[0] * original_img.shape[1], 3))
10 X_img.shape

```

B. Initialization

Randomly select initial centroids from the pixel values.

Code:

```

1 K = 8
2 max_iteration = 10
3
4 initial_centroids = kmeans_initial_centroids(X_img, K)
5
6 centroids, label = kmeans(X_img, initial_centroids, max_iteration)

```

C. Reconstruction

Replace each pixel's color with its corresponding centroid's color to compress the image.

Code:

```

1 X_recovered = centroids[label, :]
2 X_recovered = np.reshape(X_recovered, original_img.shape)

```

D. Visualization

a. Compare the original and compressed images.

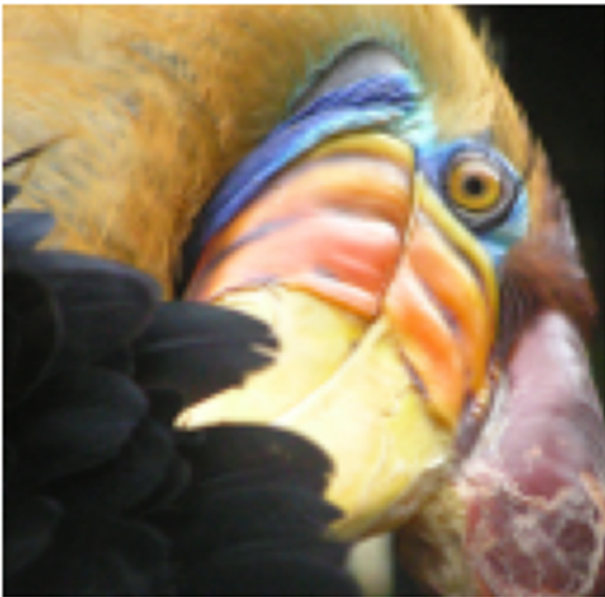
b. The original and compressed images are displayed side by side for visual comparison. The compressed image uses a reduced color palette, effectively demonstrating the compression.

Here is the code below:

```
1 # Display original image
2 fig, ax = plt.subplots(1,2, figsize=(8,8))
3 plt.axis('off')
4
5 ax[0].imshow(original_img*255)
6 ax[0].set_title('Original')
7 ax[0].set_axis_off()
8
9 # Display compressed image
10 ax[1].imshow(X_recovered*255)
11 ax[1].set_title('Compressed_with_%d_colours'%K)
12 ax[1].set_axis_off()
```

Here is the visualization output:

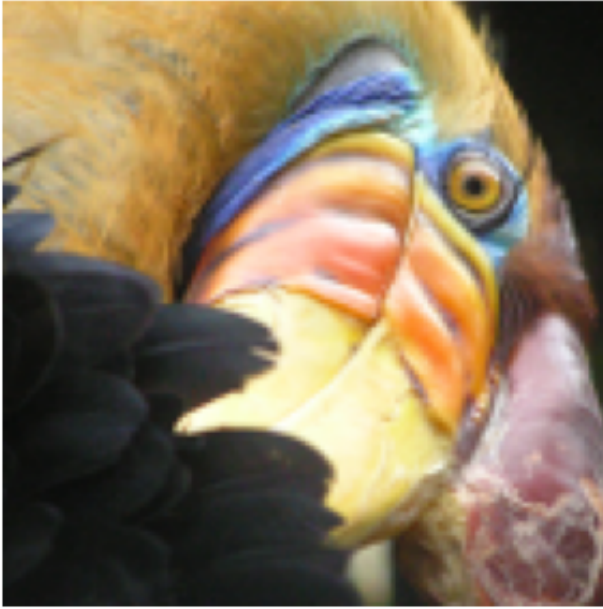
Original



Compressed with 8 colours



Original



Compressed with 32 colours



VII. EVALUATION

A. *Clustering*

To evaluate the clustering results, we plotted the clustered data points and centroids for different values of K (from 3 to 6). This visual inspection helps us understand how well the data is partitioned.

B. *Image Compression*

For image compression, we evaluated the results by visually comparing the original and compressed images. The quality of compression is judged by how closely the compressed image resembles the original while reducing the number of colors.

VIII. RESULTS

A. *2D Dataset Clustering*

Here below are the results of clustering the 2D dataset for different values of K :

1. $K=3$:

- a. Three distinct clusters with clear separation.
- b. Centroids are well positioned at the center of each cluster.

2. $K=4$:

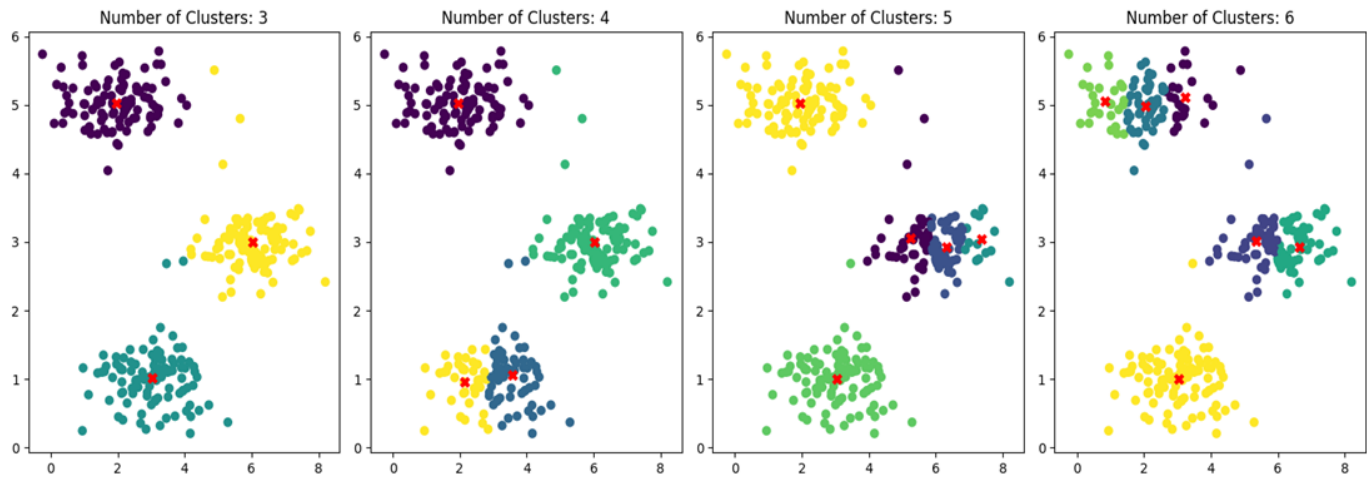
- a. Four clusters with good separation.
- b. Centroids effectively represent each cluster.

3. $K=5$:

- a. Five clusters with increasing granularity.
- b. Centroids continue to accurately represent the clusters.

4. $K=6$:

- a. Six clusters with more detailed partitioning.
- b. Centroids maintain good cluster representation.



B. Image Compression

The image compression results for $K=8$ is shown below:

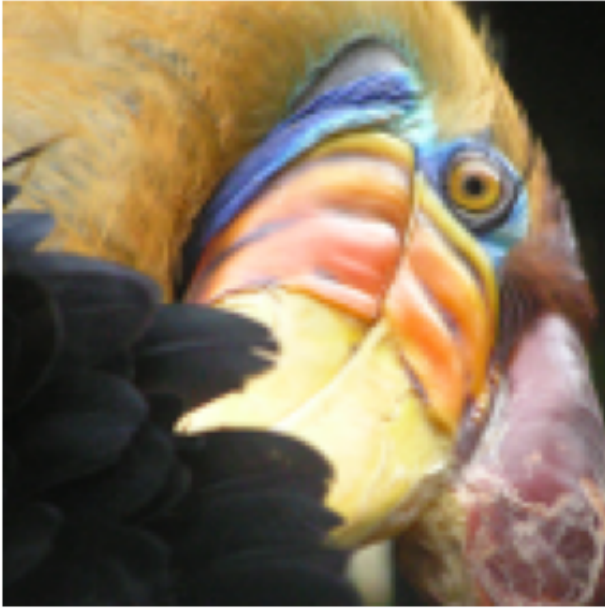
1) Original Image:

- a. Full color range with 128x128 pixels.
- b. Shape: (128, 128, 3).

2) Compressed Image::

- a. Reduced color palette with $K=16$.
- b. Visual quality is high, closely resembling the original.

Original



Compressed with 16 colours



IX. DISCUSSION:

1. The K-Means algorithm effectively clusters 2D data. For the 2D dataset, increasing the number of clusters provides finer granularity and better representation of data points.
2. In image compression, K-Means significantly reduces the number of colors while maintaining visual quality, demonstrating its utility in practical applications.