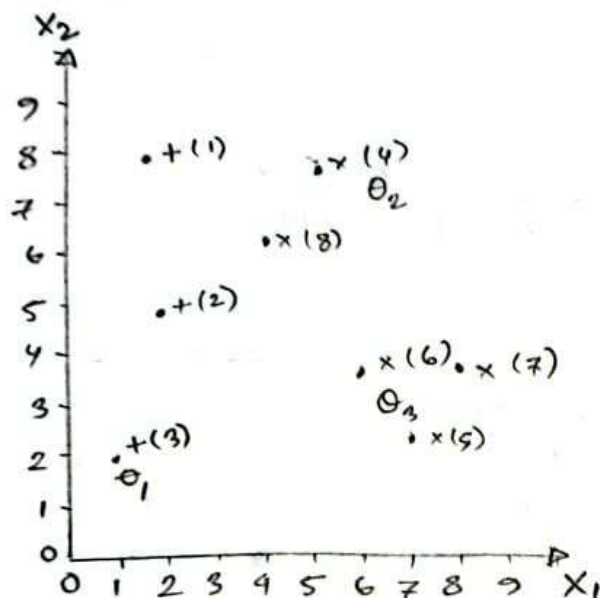## Forgy initialization:



initial cluster centers:

$$\theta_1 = x^{(3)}$$
$$\theta_2 = x^{(4)}$$
$$\theta_3 = x^{(6)}$$

Notation: $x^{(i)} = (x_1^{(i)}, x_2^{(i)})$
$$\theta_j = (\theta_{j1}, \theta_{j2})$$

## 1. Iteration

$$M_1 = M_2 = M_3 = \{ \}$$

The data points used for the initialization are easy to assign

$$M_1 = \{x^{(3)}\}, \quad M_2 = \{x^{(4)}\}, \quad M_3 = \{x^{(6)}\}$$

Let's compute the distance of each remaining data points to all cluster centre using the squared euclidian distance:

$$d(x^{(i)}, \theta_i) := \| x^{(i)} - \theta_i \|_2^2 = (x_1^{(i)} - \theta_{i1})^2 + (x_2^{(i)} - \theta_{i2})^2$$

| $d(x^{(i)}, \theta_k)$ | $\theta_1 = (1,2)$ | $\theta_2 = (5,8)$ | $\theta_3 = (6,4)$ | $x^{(i)}$'s assigned to clusters |
|---|---|---|---|---|
| $x^{(1)} = (2,8)$ | 37 | 9 | 32 | $\theta_2$ |
| $x^{(2)} = (2,5)$ | 10 | 18 | 17 | $\theta_1$ |
| $x^{(5)} = (7,3)$ | 37 | 29 | 2 | $\theta_3$ |
| $x^{(7)} = (8,4)$ | 53 | 25 | 4 | $\theta_3$ |
| $x^{(8)} = (4,7)$ | 34 | 2 | 13 | $\theta_2$ |

$$\Rightarrow \mu_1 = \{x^{(3)}, x^{(2)}\}$$
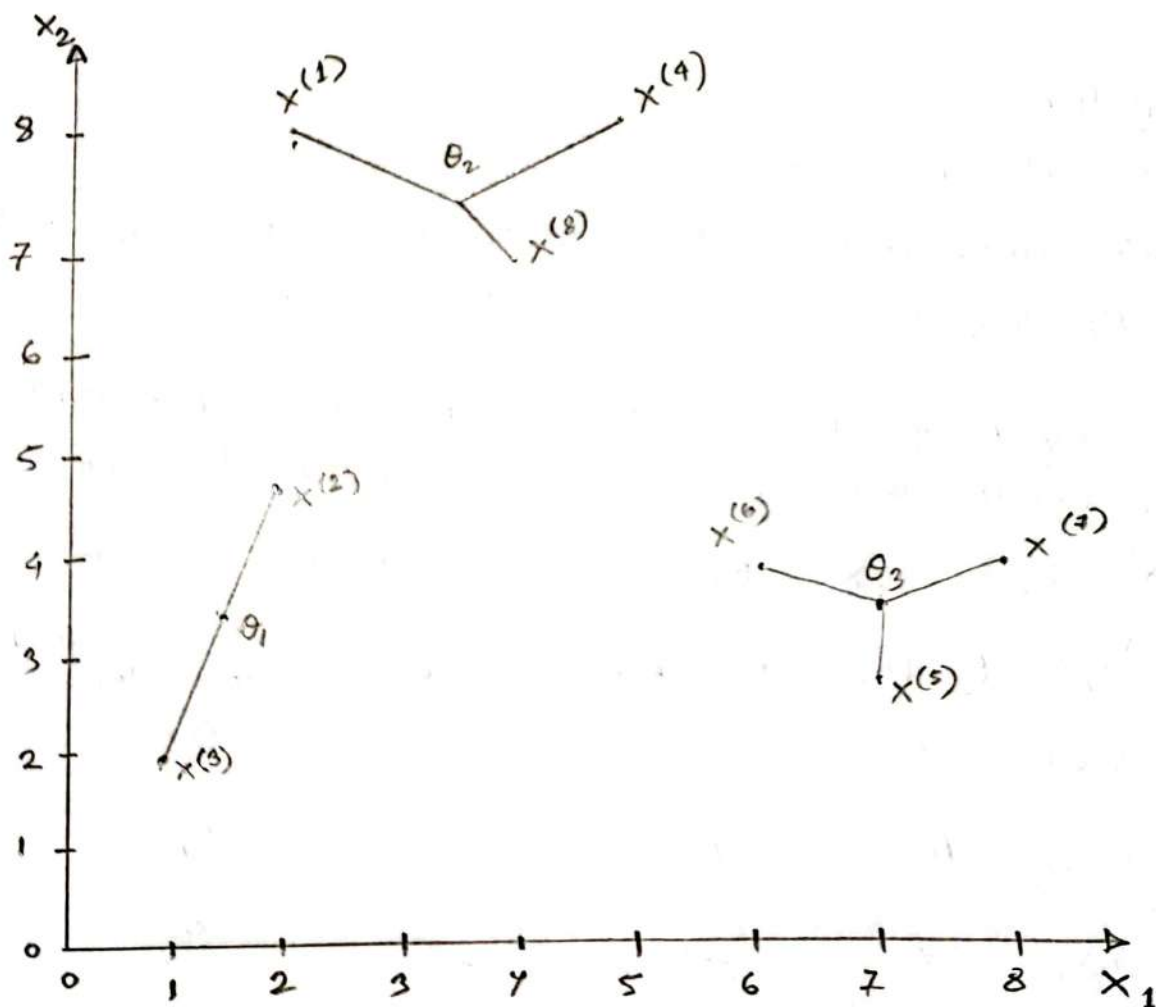
$$\mu_2 = \{x^{(4)}, x^{(1)}, x^{(8)}\}$$

$$\mu_3 = \{x^{(6)}, x^{(5)}, x^{(7)}\}$$

Compute new cluster centers: $\quad \theta_k = \dfrac{1}{|\mu_k|} \displaystyle\sum_{x^{(i)} \in \mu_k} x^{(i)}$

$$\theta_1 = \tfrac{1}{2}\left(x^{(3)} + x^{(2)}\right) = (1.5, 3.5)$$

$$\theta_2 = \tfrac{1}{3}\left(x^{(4)} + x^{(1)} + x^{(8)}\right) = \left(\tfrac{11}{3}, \tfrac{23}{3}\right)$$

$$\theta_3 = \tfrac{1}{3}\left(x^{(6)} + x^{(5)} + x^{(7)}\right) = \left(7, \tfrac{11}{3}\right)$$



Loss: $\quad L = \displaystyle\sum_{k=1}^{3} \sum_{x^{(i)} \in \mu_k} d(x^{(i)}, \theta_k) = \sum_{k=1}^{3} \sum_{x^{(i)} \in \mu_k} \left((x_1^{(i)} - \theta_{k_1})^2 + (x_2^{(i)} - \theta_{k_2})^2\right)$

$$d(x^{(2)}, \theta_1) = 2.5 \qquad d(x^{(1)}, \theta_2) \approx 2.89 \qquad d(x^{(5)}, \theta_3) \approx 0.44$$

$$d(x^{(3)}, \theta_1) = 2.5 \qquad d(x^{(4)}, \theta_2) \approx 1.89 \qquad d(x^{(6)}, \theta_3) \approx 1.11$$

$$d(x^{(8)}, \theta_2) \approx 0.56 \qquad d(x^{(7)}, \theta_3) \approx 1.11$$

$L = 13.0$

2. Iteration

$\mu_1 = \mu_2 = \mu_3 = \{\ \}$

| $d(x^{(i)}, \theta_k)$ | $\theta_1 = (1.5, 3.5)$ | $\theta_2 = (\frac{11}{3}, \frac{23}{3})$ | $\theta_3 = (7, \frac{11}{3})$ | $x^{(i)}$ is assigned to cluster |
|---|---|---|---|---|
| $x^{(1)} = (2,8)$ | 20.50 | 2.89 | 43.78 | $\theta_2$ |
| $x^{(2)} = (2,5)$ | 2.50 | 9.89 | 26.78 | $\theta_1$ |
| $x^{(3)} = (1,2)$ | 2.50 | 39.22 | 38.78 | $\theta_1$ |
| $x^{(4)} = (5,8)$ | 32.50 | 1.89 | 22.78 | $\theta_2$ |
| $x^{(5)} = (7,3)$ | 30.50 | 32.89 | 0.44 | $\theta_3$ |
| $x^{(6)} = (6,4)$ | 20.50 | 18.89 | 1.11 | $\theta_3$ |
| $x^{(7)} = (8,4)$ | 42.50 | 32.22 | 1.11 | $\theta_3$ |
| $x^{(8)} = (4,7)$ | 18.50 | 0.56 | 20.11 | $\theta_2$ |

$\Rightarrow$ each datapoint is already assigned to its closed cluster center

$\Rightarrow \mu_1 = \{x^{(3)}, x^{(2)}\}$

$\mu_2 = \{x^{(4)}, x^{(1)}, x^{(8)}\}$

$\mu_3 = \{x^{(6)}, x^{(5)}, x^{(7)}\}$

remains unchanged

$\Rightarrow L = 13.0$ unchanged

# Statistical data analysis

## Problem Sheet 9

### Exercise 2

```
1 # necessary libraries
2 from PIL import Image
3 import numpy as np
4 import os
5 import matplotlib.pyplot as plt
```

Loading the Image and converting into ndarray

```
1 file = os.path.abspath("/content/Rubix_cube_ps9.jpg")
2 path = os.path.dirname(file)
```

```
1 arr = np.array(Image.open(file))
2 print("The image and pixel has shape: " + str(arr.shape))
3
```

```
    The image and pixel has shape: (410, 400, 3)
```

*So, we can see the resolution of the image is 410x400 where each pixel has RGB values. Which means each pixel is a combination of three values in range of 0-255 .*

## Implementing the k-means algorithm.

```
1 def getCenter(X, R):
2     # initializing shapes
3     le, wi, n = X.shape
4     k = R.shape[-1]
5     # initializing Means
6     Miu = np.zeros((k, n))
7
8     # converting R shape into 2d
9     R = R.reshape((le*wi, k))
10    # converting X as 2d
11    X = X.reshape((le*wi, n))
12
13    # iterate over each cluster
14    for i in range(k):
15        # get the pixels that belong to cluster i and calculate their mean value
16        Miu[i] = X[np.where(R[:, i] == 1)].mean(axis=0)
17
18    return Miu
```

```
1 def getCluster(X, Miu):
2     # initializing shapes
3     le, wi, n = X.shape
4     k = Miu.shape[0]
5     # initializing R in 2D shape
6     R = np.zeros((le*wi, k))
7
8     # convert the image ndarray into 2d shape
9     X = arr.reshape((le*wi,n))
10
11    # calculating the interval
```

```
12      interv = np.linalg.norm(np.repeat(X, k, axis=0) - np.tile(Miu, (le * wi, 1)), axis=1, keepdims=True)
13
14      # converting interval array shape
15      interv = interv.reshape(le*wi, k)
16
17      # finding ic as the index of the closest mean
18      ic = np.argmin(interv, axis=1)
19
20      # plugging 1 to the closest mean
21      R[np.arange(le*wi),ic] = 1
22      # converting R into 3d shape
23      R = R.reshape((le,wi,k))
24
25      return R
26
```

```
 1 def kmeans(X, k):
 2      # initializing shapes
 3      le, wi, n = X.shape
 4      y = np.zeros([le,wi,])
 5      # initializing R as 3d shape
 6      R = np.zeros((le,wi,k))
 7      i = 0
 8      verbose = 0
 9      maxiter = 200
10      tolerance = 1e-2
11
12      # Initializing means randomly
13      np.random.seed(9)
14      Miu = np.random.uniform(low=np.min(X), high=np.max(X), size=(k, n))
15
16      # iterate until convergence or maxiter
17      while True:
18          # stopping criterion maxiter
19          if i>maxiter:
20              break
21          i += 1
22
23          if verbose == 1:
24              print(Miu)
25
26          try:
27              # calculating the cluster pixel for
28              newR = getCluster(X, Miu)
29              # if newR ~ R, then break
30              assert np.abs(newR - R).sum() < tolerance*newR.size
31              break
32          except AssertionError:
33              # if R is different than newR, then plugin the newR
34              R = newR
35              # calculating the new cluster means
36              Miu = getCenter(X, R)
37
38      # calculate y as the cluster pixel
39      y = np.argmax(R.reshape((le*wi, k)), axis=1).reshape((le, wi))
40
41      return (y, Miu)
42
```

Using k-means cluster all the pixels of an image into k clusters and assign each pixel the color represented by its nearest cluster center.
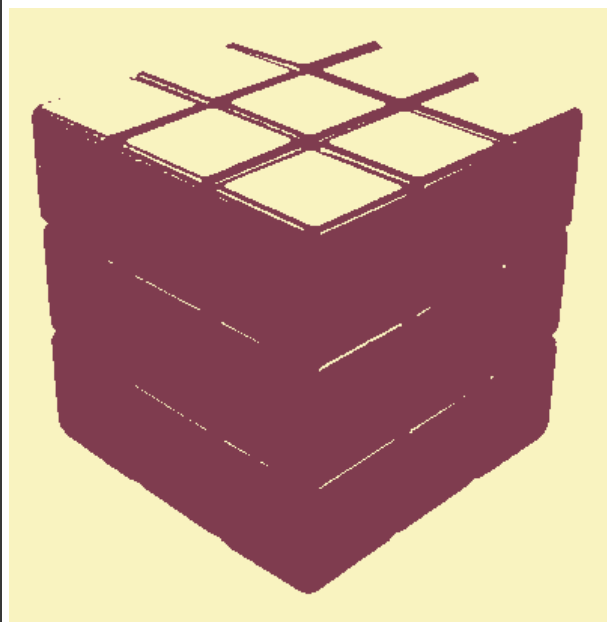
```
 1 def getImage(X, k):
```
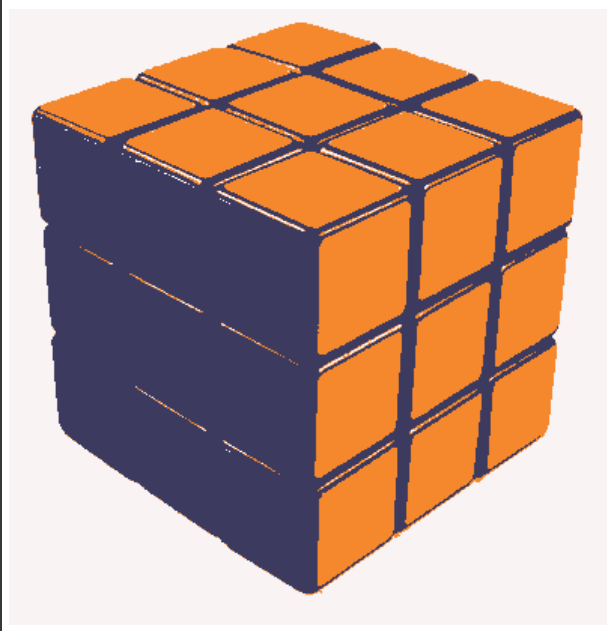
```
2      # initializing segmented image
3      seImage = np.zeros(X.shape)
4      # run k-means
5      y, Miu = kmeans(X, k)
6
7      # assigning each pixel the color of the nearest cluster center
8      for i in range(k):
9          seImage[np.where(y == i)] = Miu[i]
10
11     # generating the segmented Image
12     data = Image.fromarray(np.uint8(seImage))
13
14     return data
```
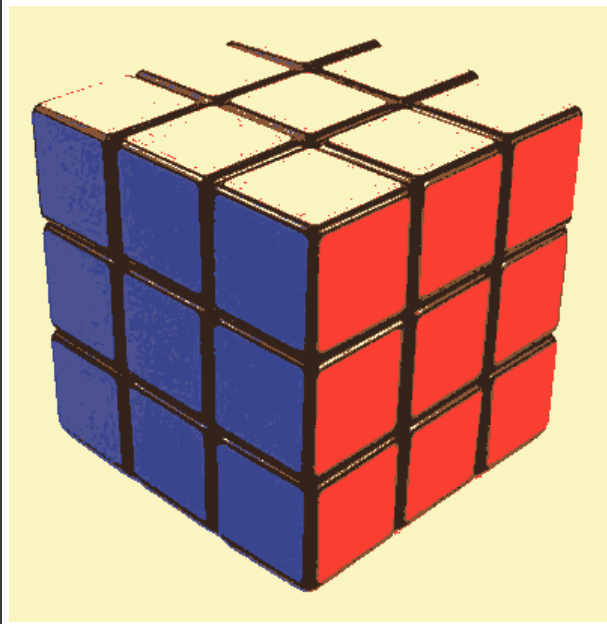
```
1 # 2 clusters
2 image = getImage(arr, 2)
3 image
```
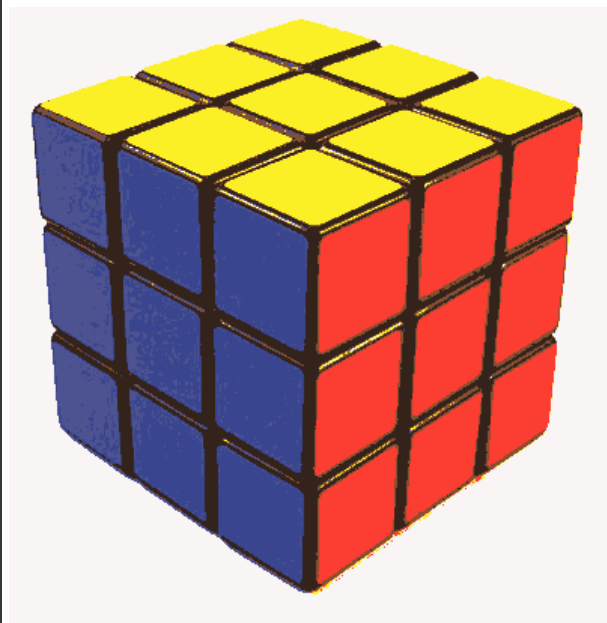


```
1 # 3 clusters
2 image = getImage(arr, 3)
3 image
```

```
1 # 6 clusters
2 image = getImage(arr, 6)
3 image
```



```
1 # 7 clusters
2 image = getImage(arr, 7)
3 image
```



**These are different clusters where each pixel color is taken from its nearest cluster center.**