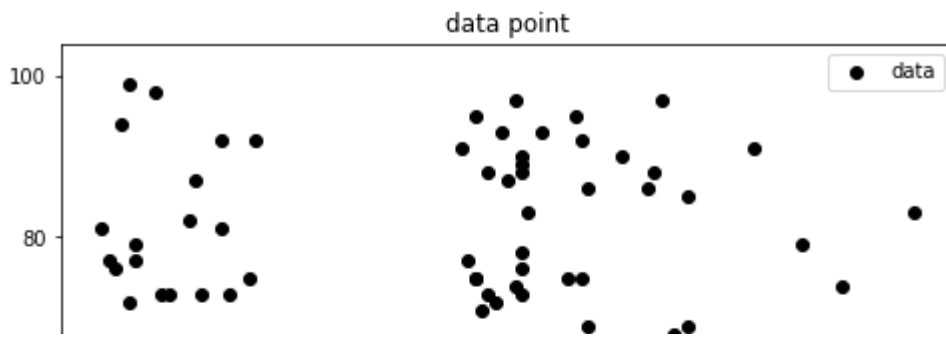# K-means clustering

## 1. Data

- the data are given by the file data-kmeans.csv
- the data consist of a set of points $\{(x_i, y_i)\}_{i=1}^{n}$ where $z_i = (x_i, y_i)$ denotes a 2-dimensional point in the cartesian coordinate and $n$ is given as $200$

load the data from the files

```
1   import numpy as np
2   import matplotlib.pyplot as plt
3   import pandas as pd
4   import random as rd
5
6   path = '/content/drive/My Drive/ML_Assignment/data/data-kmeans.csv'
7   dataset = pd.read_csv(path)
8   data = dataset.values
9   x_data = data[:,0] # x
10  y_data = data[:,1] # y
```

Plot the data points

```
1   fig_1 = plt.figure(figsize = (8,8))
2   plt.scatter(x_data, y_data, c='k', label='data')
3   plt.title('data point')
4   plt.legend()
5   plt.show()
6   fig_1.savefig('data point.png')
```

data point

## 2. Loss

- the loss function $\mathcal{L}(C_1, C_2, \cdots, C_k, \mu_1, \mu_2, \cdots, \mu_k)$ with a given number of clusters $k$ for a set of data $\{z_i\}_{i=1}^n$ is defined by:

$$\mathcal{L}(C_1, C_2, \cdots, C_k, \mu_1, \mu_2, \cdots, \mu_k) = \frac{1}{n} \sum_{i=1}^n \|z_i - \mu_{l(z_i)}\|_2^2 = \frac{1}{n} \sum_{j=1}^k \sum_{z_i \in C_j} \|z_i - \mu_j\|_2^2$$

- $l(z) = k$ is a label function that defines a label $k$ of point $z$
- $C_k$ denotes a set of points $\{z_i | l(z_i) = k\}$ of label $k$
- $\mu_k$ denotes a centroid of points in $C_k$

define a function to compute a initial centroid

```
1  def init_centroid(k):
2     centroids = np.array([]).reshape(2,0)
3     for i in range(k):
4        rand = rd.randint(0, 200)
5        centroids = np.c_[centroids, data[rand]]
6
7     return centroids.T
```

define a function to compute a distance between two points $a$ and $b$

```
1  def compute_distance(data, c):
2
3      dist = np.array([]).reshape(200,0)
4
5      # distance between data and cluster
6      for i in range(5):
7         i_dist = np.sqrt(np.sum((data - c[i,:])**2, axis=1))
8         dist = np.c_[dist, i_dist]
9
10     return dist
```

```
1  def compute_centroid_distrance(c):
2      dist = []
3      # distance between data and cluster
4      for i in range(5):
5         i_dist = np.sqrt(np.sum((c[i,:])**2))
6         dist.append(i_dist)
```

```
7
8        return dist
```

define a function to compute a centroid from a given set of points $Z$

```
1    def compute_centroid(cluster):
2        center = np.array([]).reshape(2,0)
3        # centroid of a set of points in Z
4        for i in range(5):
5            idx = (cluster[:,2]==i)
6            i_center = np.mean(data[idx],axis=0)
7            center = np.c_[center, i_center]
8        return center.T
```

define a function to compute the loss with a set of clusters $C$ and a set of centroids $M$

```
1    def compute_loss(cluster, centroids):
2        loss_list = []
3        loss = 0
4
5        for i in range(5):
6            idx = (cluster[:,2]==i)
7            i_loss = np.sqrt(np.sum((data[idx] - centroids[i,:])**2))
8            loss += i_loss
9
10       loss = loss / len(cluster)
11       return loss
```

# 3. Optimization

- the label $l(z)$ of each point $z$ is determined by: $l(z) = \arg\min_k \|z - \mu_k\|_2^2$
- the centroid $\mu_i$ of cluster $k$ is determined by: $\mu_k = \frac{\sum_{z_i \in C_k} z_i}{|C_k|}$

define a function to determine the label of point $z$ with a set of centroids $M$

```
1    def compute_label(dist):
2
3        argmin_label = np.argmin(dist, axis=1) #label of point z with a set of centroids M#
4        label = np.c_[data, argmin_label]
5
6        return label
```

# 4. Clustering

- initialise labels $l(z_i)$ for point $z_i$ for all $i$ randomly

- optimise the loss function with respect to the centroids and the clusters in an alternative way
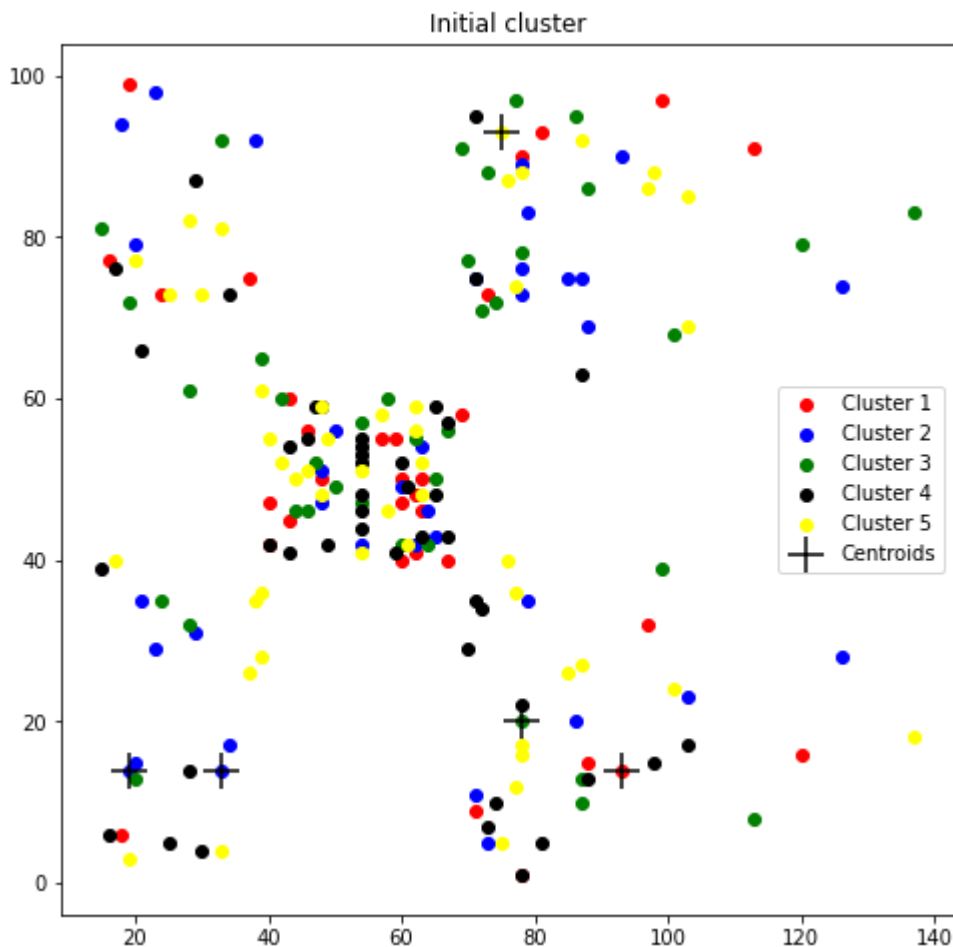- set the number of clusters $k = 5$

Visualise the initial condition of the point labels

```
1   k = 5    # set the number of clusters
2   n = len(data)
3   max_iter = 50
4   centroids = init_centroid(k)
5   labels = np.random.randint(low=0, high=k, size=n)
6   result = np.c_[data, labels]
```

```
1    fig_2 = plt.figure(figsize = (8,8))
2    color=['red','blue','green', 'black', 'yellow']
3    label=['Cluster 1','Cluster 2','Cluster 3', 'Cluster 4', 'Cluster 5']
4    for i in range(k):
5        idx = (result[:,2]==i)
6        plt.scatter(x_data[idx],y_data[idx], c=color[i],label=label[i])
7    plt.scatter(centroids[:,0],centroids[:,1],s=300, c='k', marker='+', label='Centroids')
8
9    plt.title('Initial cluster')
10   plt.legend()
11   plt.show()
12   fig_2.savefig('Initial cluster.png')
```


Initial cluster

```
1    def k_means_clustering(max_iter, data, centroids):
```

```
2        loss_iters = [] # record the loss values
3        centroid_iters = []
4
5        for i in range(max_iter):
6            dist = compute_distance(data, centroids)
7            cluster = compute_label(dist)
8            centroids = compute_centroid(cluster)
9            loss = compute_loss(cluster, centroids)
10           c_dist = compute_centroid_distrance(centroids)
11           loss_iters.append(loss)        # save the current loss value
12           centroid_iters.append(c_dist)
13
14       return cluster, centroids, loss_iters, centroid_iters
```

```
1    final_result, final_c, loss_iter, c_iter = k_means_clustering(max_iter, data, centroids)
```
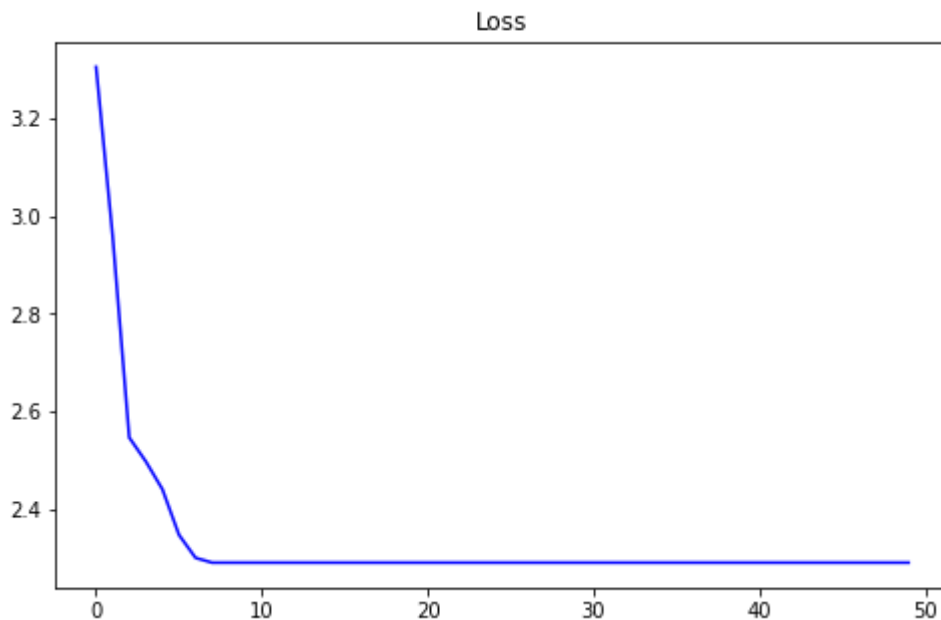
## Plot the loss curve

```
1    # Plot the loss curve
2    fig_3 = plt.figure(figsize = (8,5))
3    plt.plot(np.array(range(max_iter)),loss_iter, c = 'b')
4    plt.title('Loss')
5    plt.show()
6    fig_3.savefig('Loss.png')
```



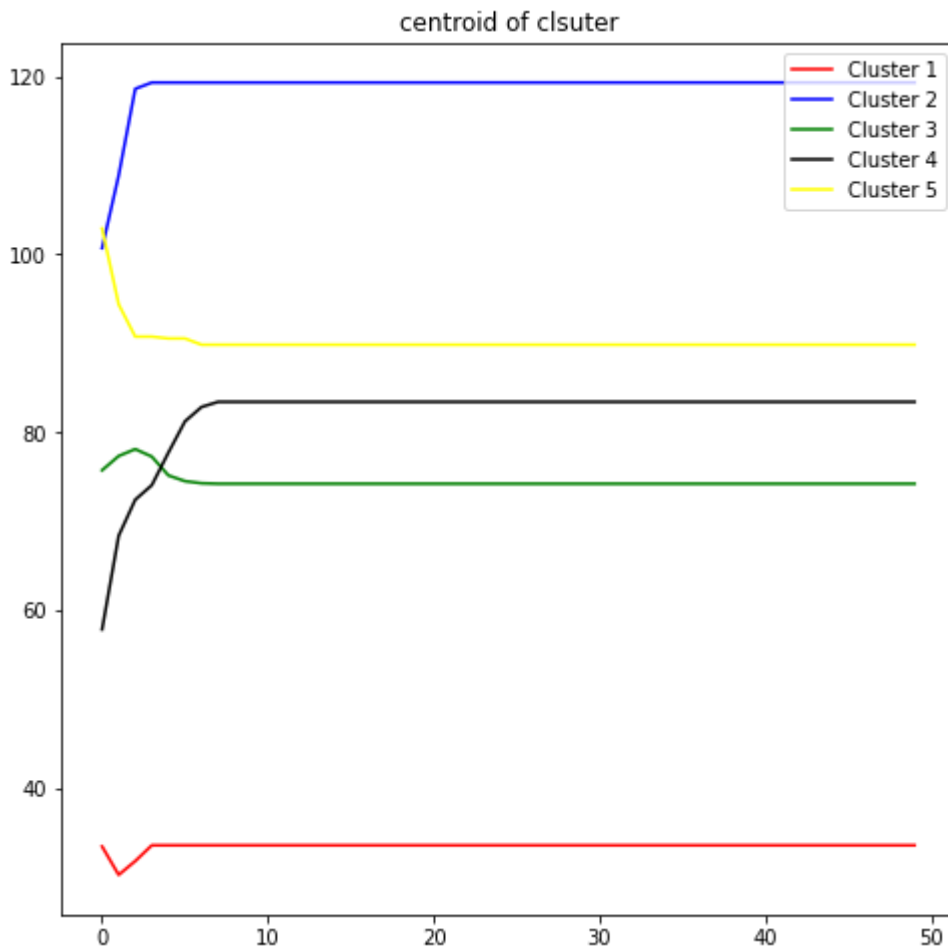## Plot the centroid of each clsuter

```
1    # Plot the centroid of each clsuter
2    fig_4 = plt.figure(figsize = (8,8))
3    color=['red','blue','green', 'black', 'yellow']
4    label=['Cluster 1','Cluster 2','Cluster 3', 'Cluster 4', 'Cluster 5']
5    np_c_iter = np.array(c_iter)
6
7    for i in range(k):
8        idx = (final_result[:,2]==i)
9        plt.plot(np.array(range(max_iter)), np_c_iter[:, i], c=color[i], label=label[i])
```

```
 9     plt.plot(np.array(range(max_iter)), np_c_iter[:,i], c=color[i],label=label[i])
10
11     plt.title('centroid of clsuter')
12     plt.legend(loc = 'upper right')
13     plt.show()
14     fig_4.savefig('centroid of clsuter.png')
```
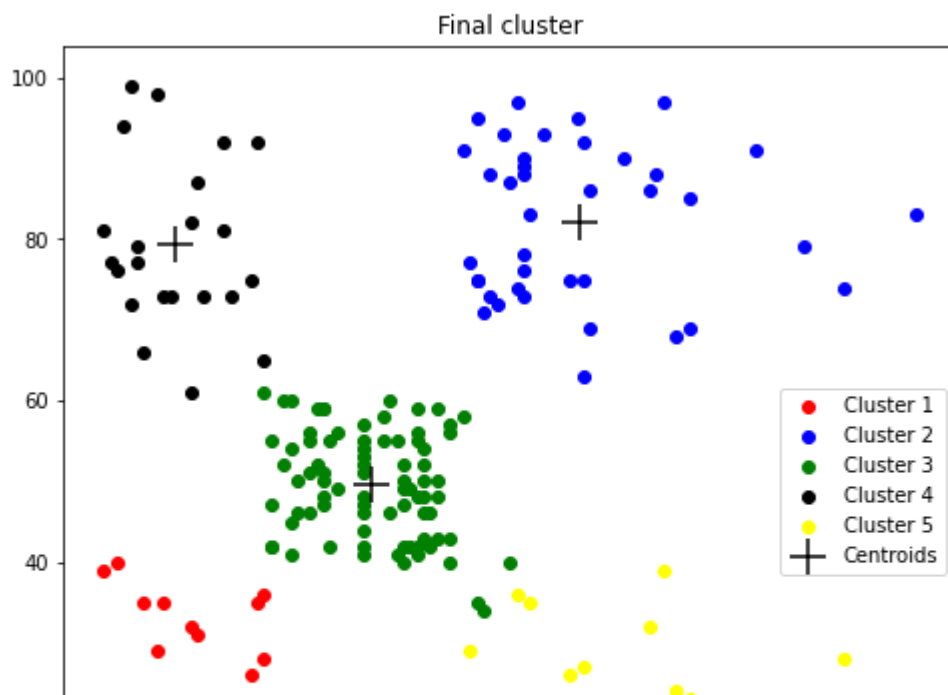


centroid of clsuter

## Plot the final clustering result

```
 1     fig_5 = plt.figure(figsize = (8,8))
 2     color=['red','blue','green', 'black', 'yellow']
 3     label=['Cluster 1','Cluster 2','Cluster 3', 'Cluster 4', 'Cluster 5']
 4     for i in range(k):
 5         idx = (final_result[:,2]==i)
 6         plt.scatter(x_data[idx],y_data[idx], c=color[i],label=label[i])
 7     plt.scatter(final_c[:,0],final_c[:,1],s=300, c='k', marker='+', label='Centroids')
 8
 9     plt.title('Final cluster')
10     plt.legend()
11     plt.show()
12     fig_5.savefig('Final cluster.png')
```
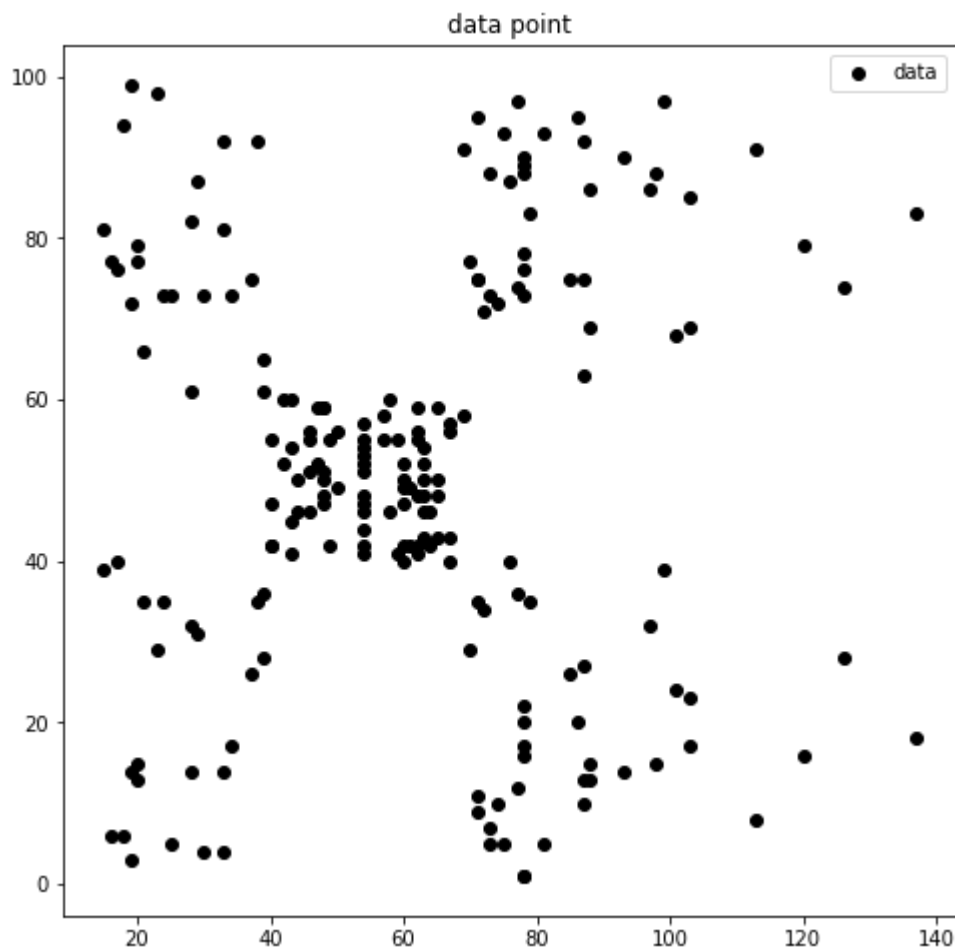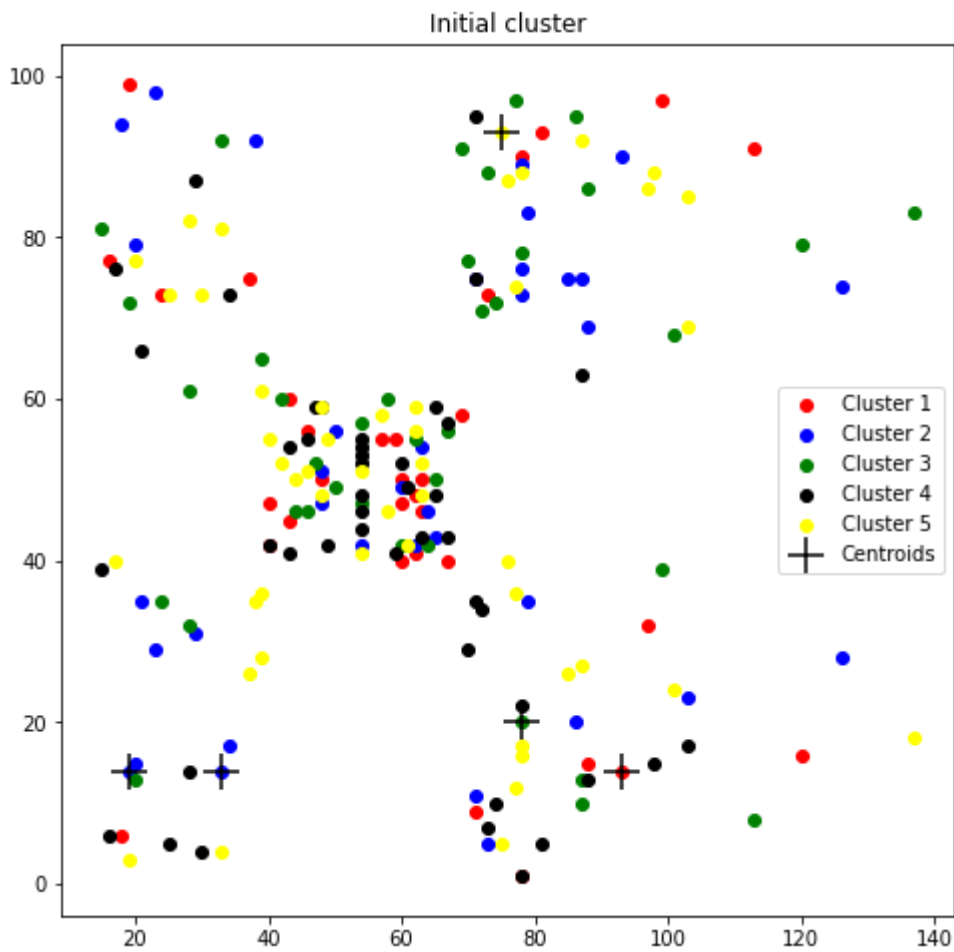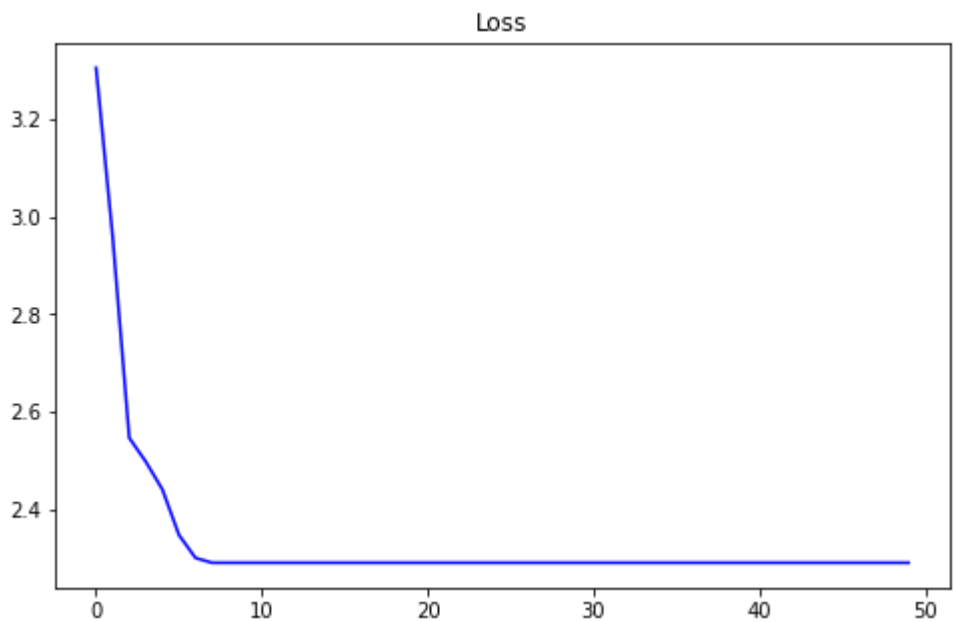
Final cluster

## 1. Plot the data points [1pt]

```
1   fig_1
```



data point

1    fig_2

### Initial cluster



Legend:
- Cluster 1 (red)
- Cluster 2 (blue)
- Cluster 3 (green)
- Cluster 4 (black)
- Cluster 5 (yellow)
- Centroids (+)

## 3. Plot the loss curve [5pt]

1    fig_3

### Loss

## 4. Plot the centroid of each clsuter [5pt]

```
1    fig_4
```

centroid of clsuter



## 5. Plot the final clustering result [5pt]

```
1    fig_5
```

Final cluster