

Lab 4 : CLUSTERING Part 1

In this Lab you will have to write code for 2 clustering algorithms based on the mathematical theory :

1. K-means Clustering
2. Gaussian Mixture Model

You will then have to use these algorithms on a practical dataset and compare the results with the inbuilt algorithms present in scikit learn toolkit

Please use plots wherever possible to demonstrate the results

```
In [1]: import numpy as np
import matplotlib.pyplot as plt
from jupyterthemes import jtplot
jtplot.style(theme = "monokai", context = "notebook", ticks =
True, grid = False)
```

K-means Clustering

K-means clustering is a type of unsupervised learning, which is used when you have unlabeled data (i.e., data without defined categories or groups). The goal of this algorithm is to find groups in the data, with the number of groups represented by the variable K. The algorithm works iteratively to assign each data point to one of K groups based on the features that are provided.

Step 1 : Data Generation

Generate 2D gaussian data of 4 types each having 100 points, by taking appropriate mean and variance (example: mean : (0.5 0) (5 5) (5 1) (10 1.5), variance : Identity matrix)

```
In [2]: # write your code here
#generation of data
def dataset(num, mean, cov):
    data = np.random.multivariate_normal(mean, cov, num)
    return(data.T[0], data.T[1])

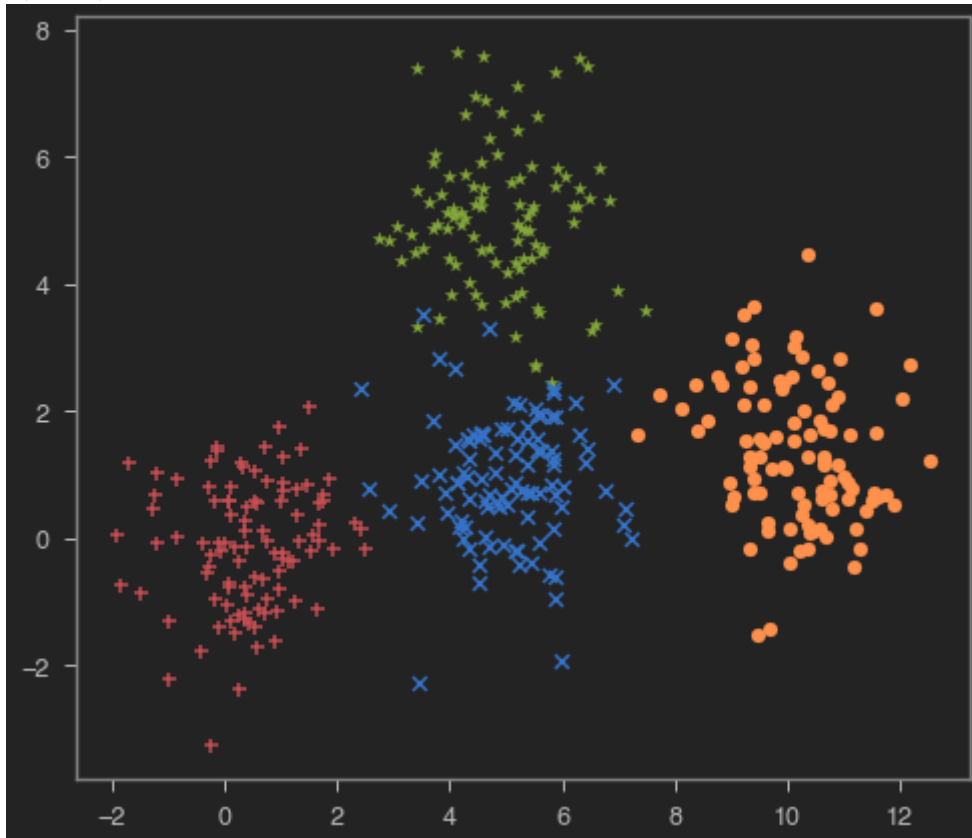
datax = np.array([])
datay = np.array([])
mu = np.array([[0.5, 0], [5, 5], [5, 1], [10, 1.5]])
cov = np.array([[1, 1], [1, 1], [1, 1], [1, 1]])
clr = ['r', 'g', 'b', 'y']
mrk = ['+', '*', 'x', 'o']
for i in range(0, 4):
```

```

mean = mu[i,:]
covariance = np.array([[cov[i,0],0],[0,cov[i,1]]])
x,y = dataset(100,mean,covariance)
datax = np.append(datax,x)
datay = np.append(datay,y)
plt.scatter(x,y,color = clr[i],marker = mrk[i])
print(datax.shape)

```

(400,)



Step 2 : Cluster Initialisation

Initialise K number of Clusters (Here, K=4)

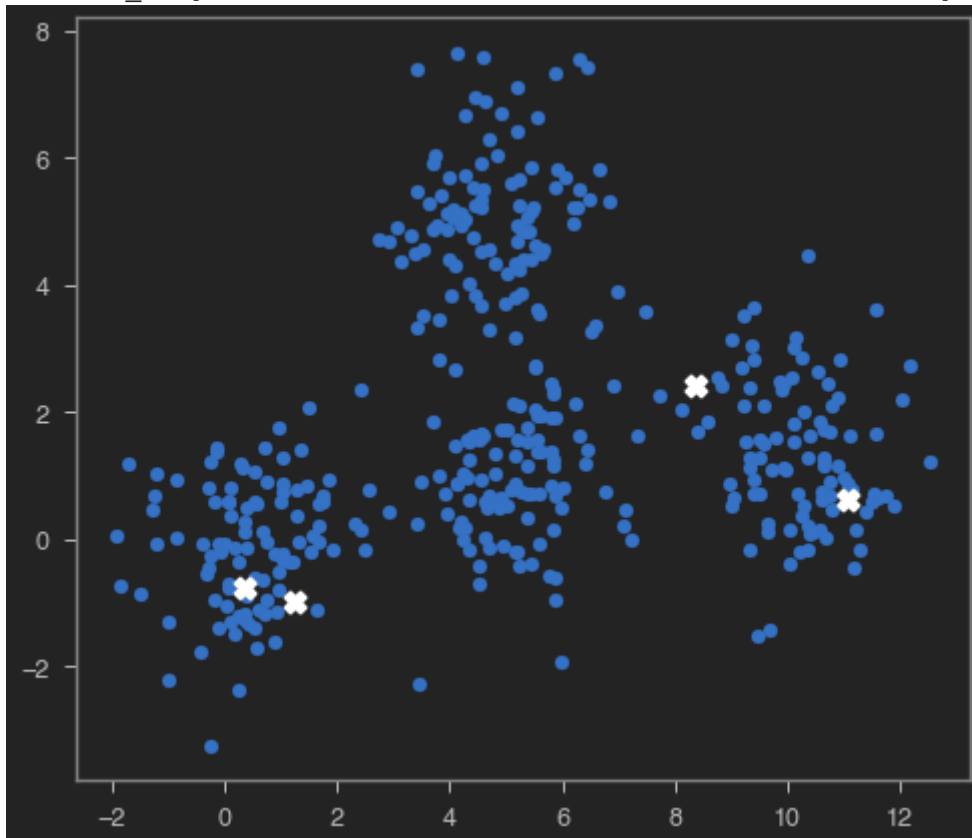
In [3]:

```

# write your code here
#defining clusters
ind = np.random.randint(0,399,4)
clstr_x = np.array([])
clstr_y = np.array([])
# clstr_x = np.random.choice(datax.flatten(), 4, replace=False)
# clstr_y = np.random.choice(datay.flatten(), 4, replace=False)
clstr_x = np.append(clstr_x,datax[ind])
clstr_y = np.append(clstr_y,datay[ind])
plt.figure()
plt.scatter(datax,datay)
plt.scatter(clstr_x,clstr_y,color = 'white',marker =
'x',linewidth=2,s=100)
print("Cluster_X:",clstr_x,"\nCluster_Y:",clstr_y)

```

```
Cluster_X: [ 0.36122971  8.3522239   1.2388255  11.05168213]
Cluster_Y: [-0.77162826  2.42640321 -0.97241674  0.62156327]
```



Step 3 : Cluster assignment and re-estimation Stage

a) Using initial/estimated cluster centers (mean μ_i) perform cluster assignment.

b) Assigned cluster for each feature vector (X_j) can be written as:

$$\arg \min_i \|C_i - X_j\|_2, \quad 1 \leq i \leq K, \quad 1 \leq j \leq N$$

c) Re-estimation: After cluster assignment, the mean vector is recomputed as,

$$\mu_i = \frac{1}{N_i} \sum_{j \in i^{th} cluster} X_j$$

where N_i represents the number of datapoints in the i^{th} cluster.

d) Objective function (to be minimized):

$$Error(\mu) = \frac{1}{N} \sum_{i=1}^K \sum_{j \in i^{th} cluster} \|C_i - X_j\|_2$$

In [4]:

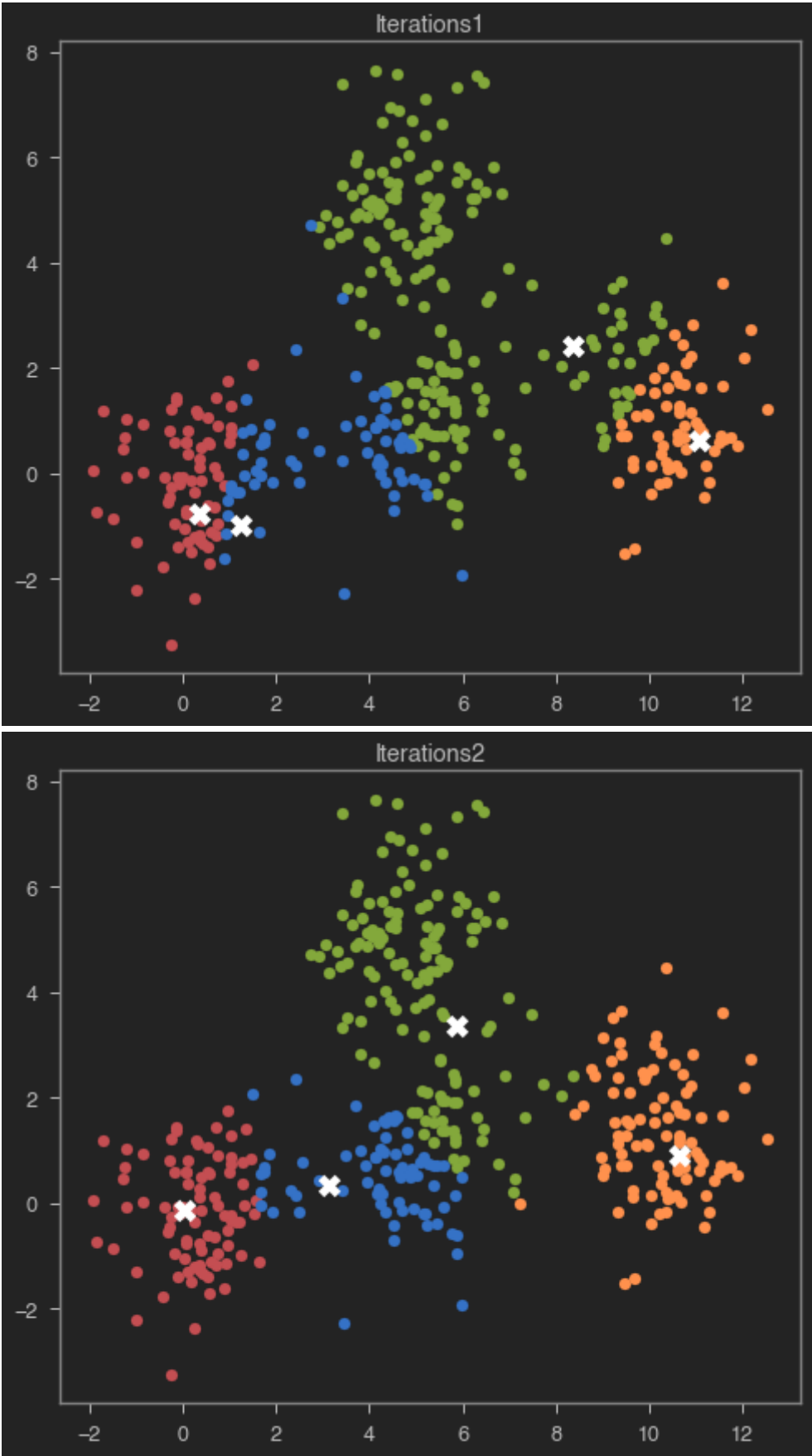
```
# write your code here
Error = np.zeros(15)
for num in range(15):
    clstr_ass = {
        0 : [],[],
        1 : [],[],
```

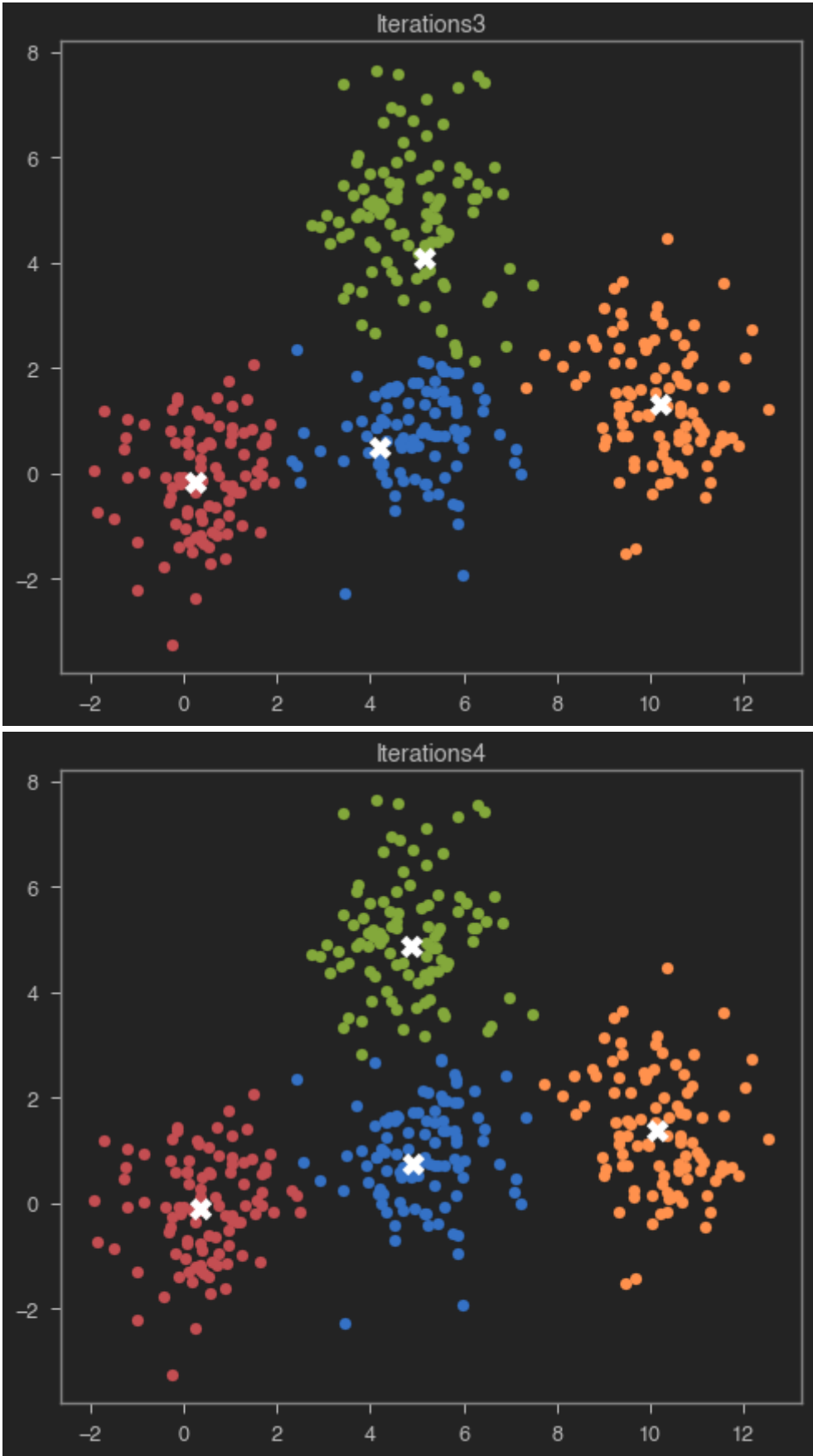
```

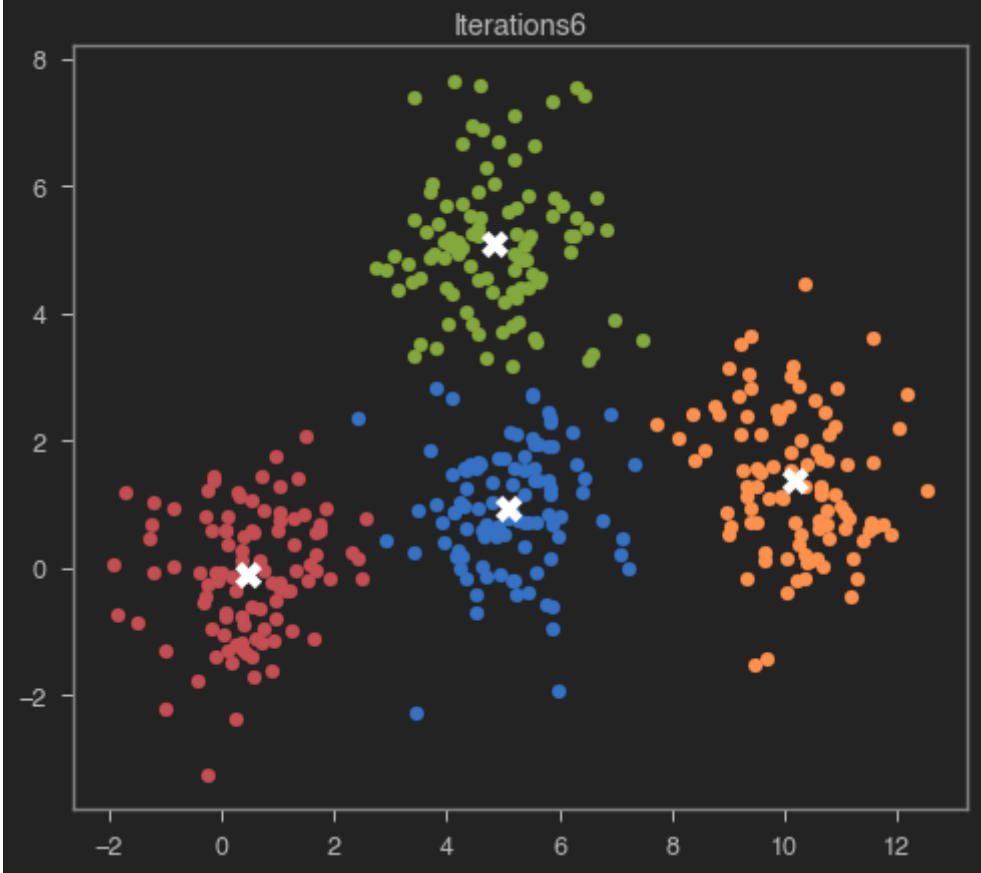
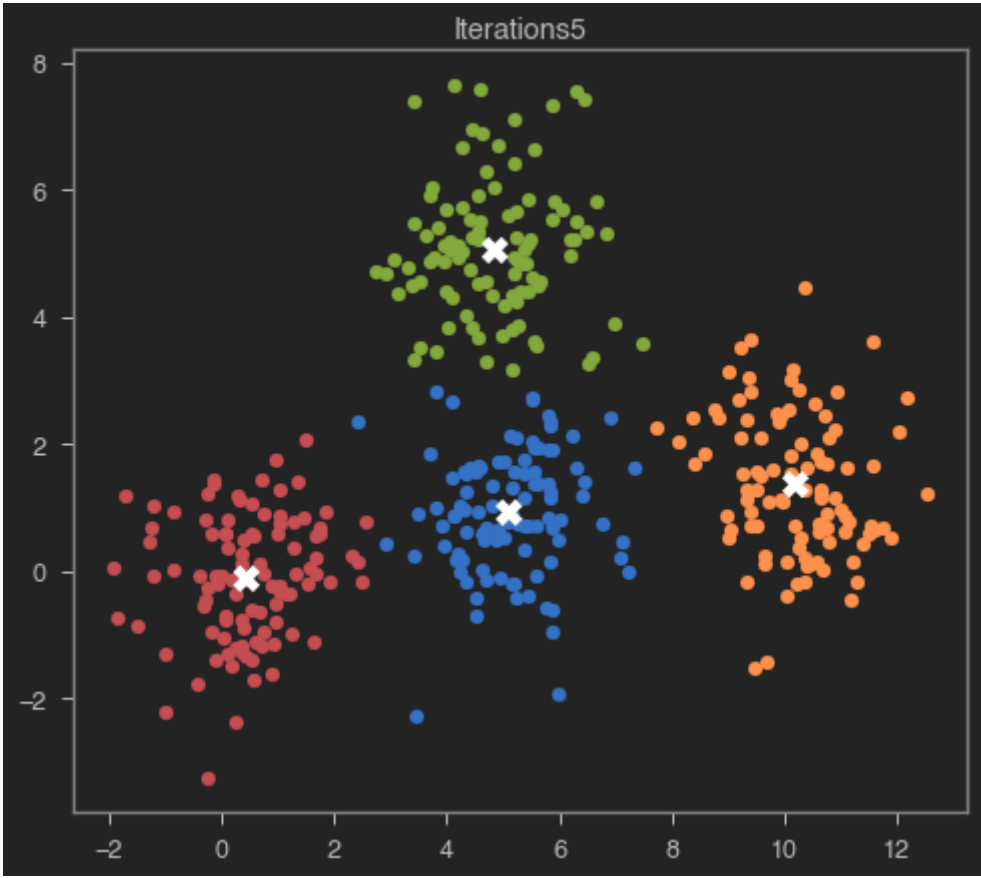
        2 : [[],[ ]],
        3 : [[],[ ]]
    }
    distx = np.array([ ])
    disty = np.array([ ])
    for i in range(0,4):
        distx = np.append(distx,(clstr_x[i] - datax)**2)
        disty = np.append(disty,(clstr_y[i] - datay)**2)
        dist = distx+disty
    Error[num] = 4*np.mean((dist)**0.5)
    dist = np.reshape(dist,(4,400))
    dist = dist.T
    min_d = np.min(dist,axis=1)
    inx = np.argmin(dist,axis=1)
    Error[num] = np.mean((min_d)**0.5)
    clr = ['r','g','b','y']
    mrk = ['+','*','x','o']
    for i in range(len(inx)):
        clstr_ass[inx[i]] =
np.append(clstr_ass[inx[i]],np.array([[datax[i]],
[datay[i]]]),axis=1)
    plt.figure()
    for i in range(4):
        plt.scatter(clstr_ass[i][0,:],clstr_ass[i]
[1,:],color=clr[i])
        plt.scatter(clstr_x[i],clstr_y[i],marker =
'x',color='white',s=100,linewidths=5)
    plt.title("Iterations"+str(num+1))
    for i in range(4):
        clstr_x[i] = np.mean(clstr_ass[i][0,:])
        clstr_y[i] = np.mean(clstr_ass[i][1,:])
plt.figure()
plt.plot(Error,color='white')
# print(clstr_ass[1][0,:].size)
# print("Cluster_X:",clstr_x,"\nCluster_Y:",clstr_y)

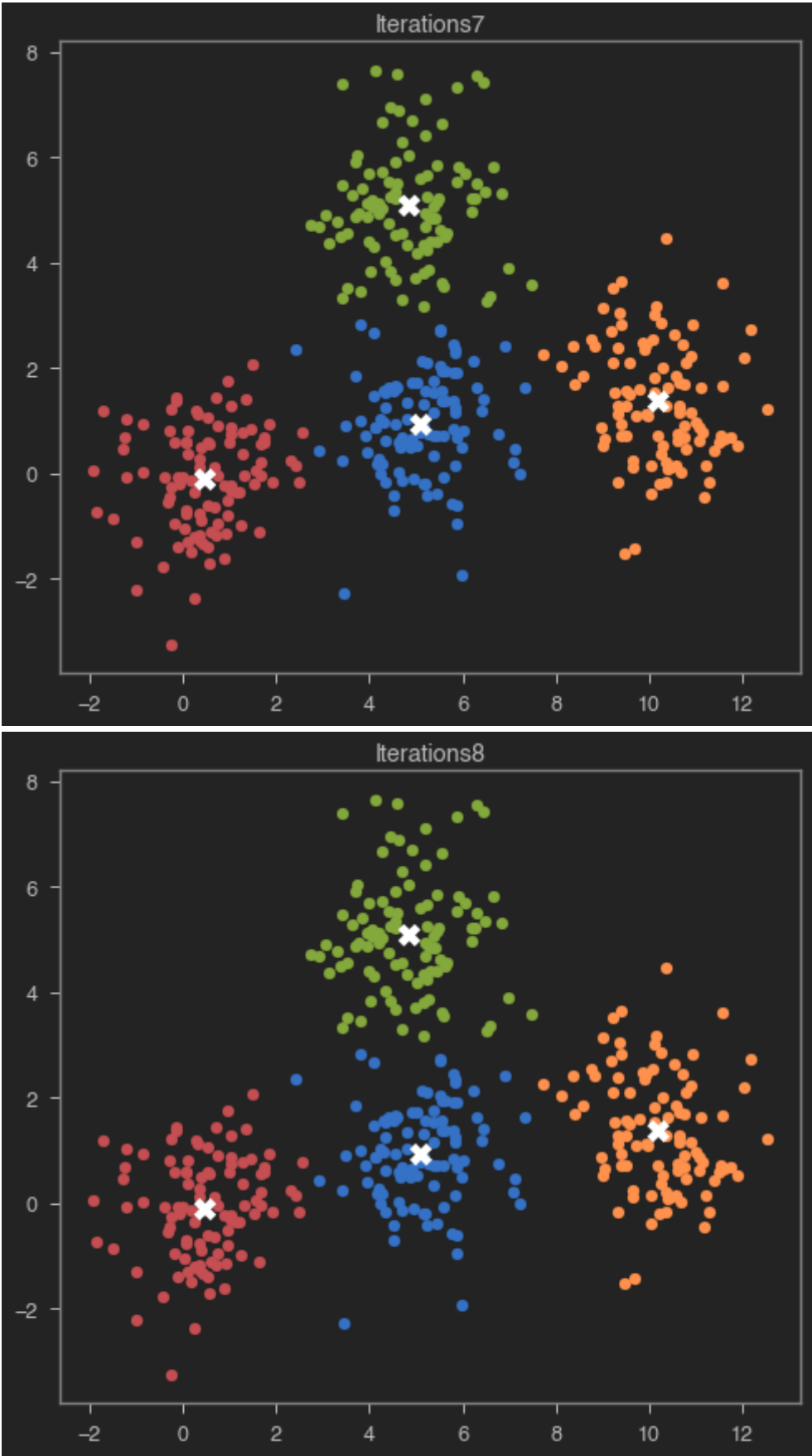
```

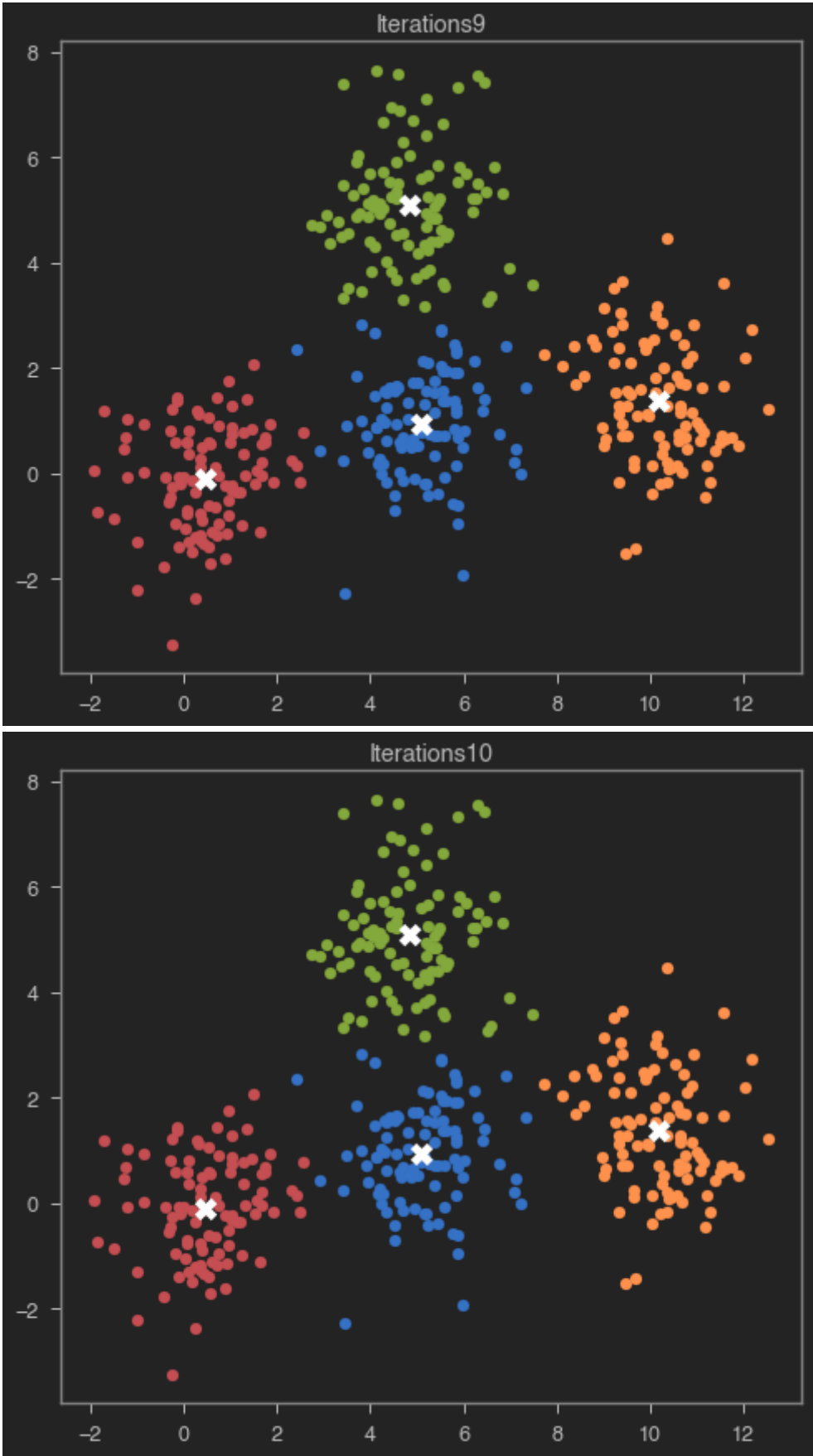
Out[4]: []

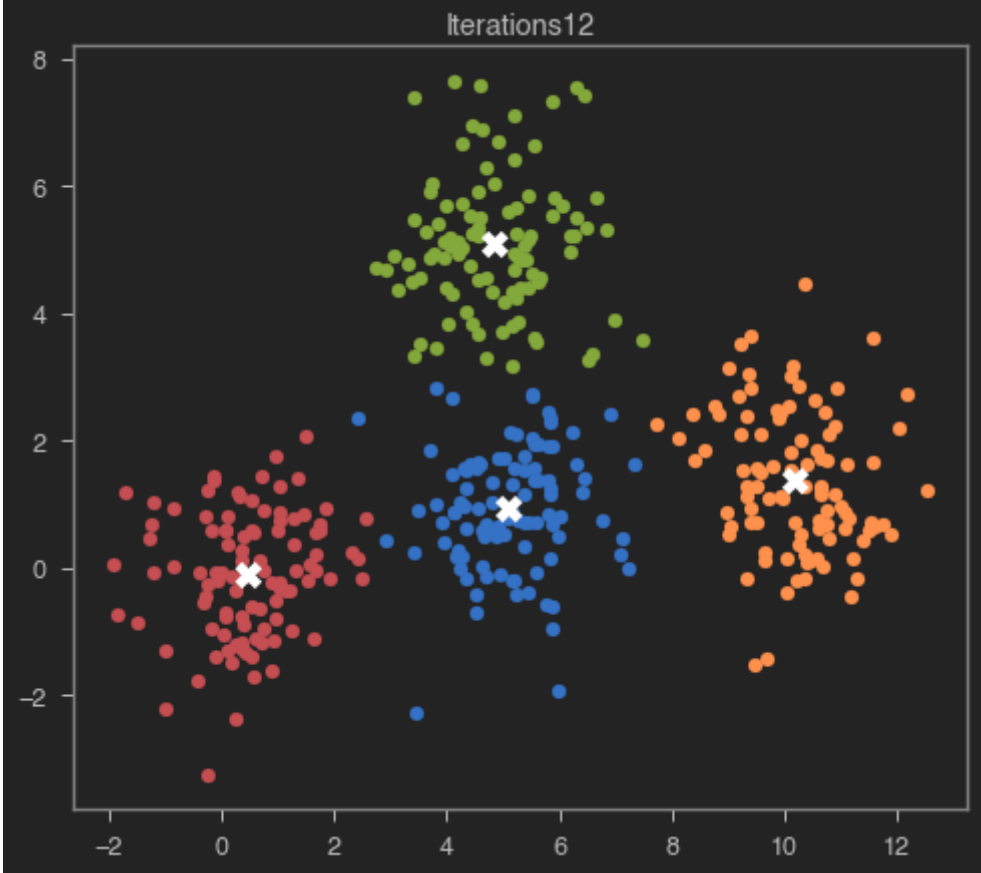
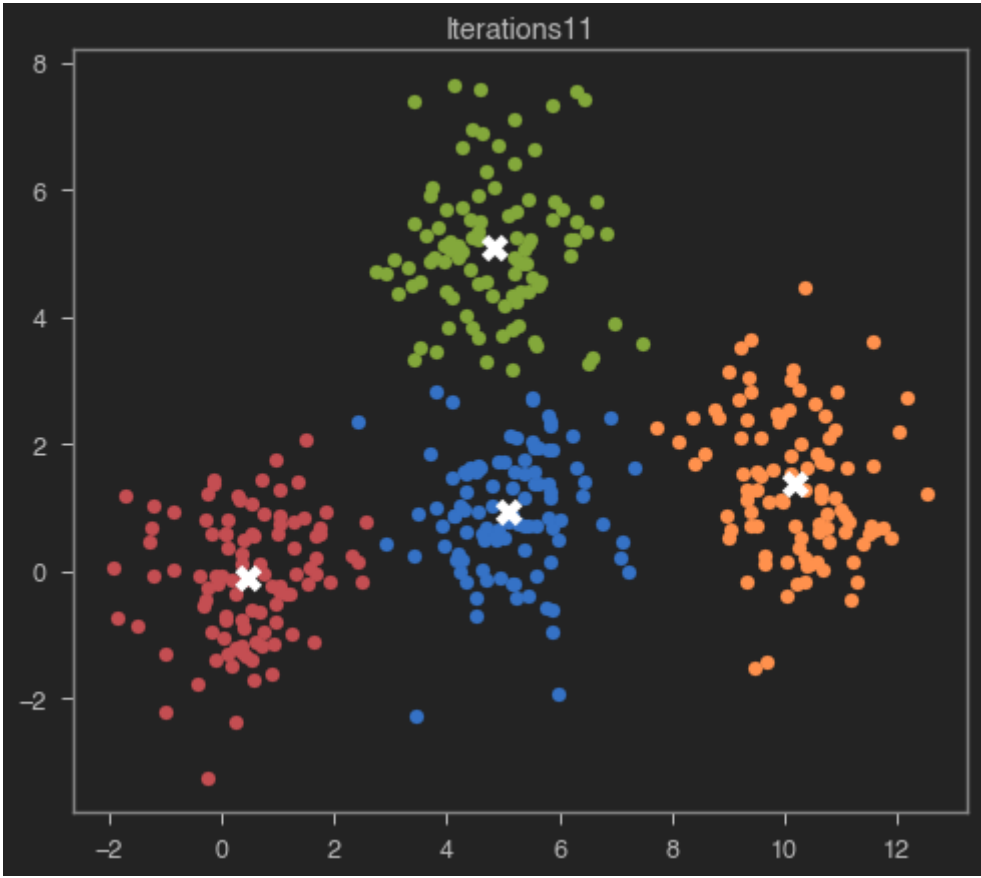


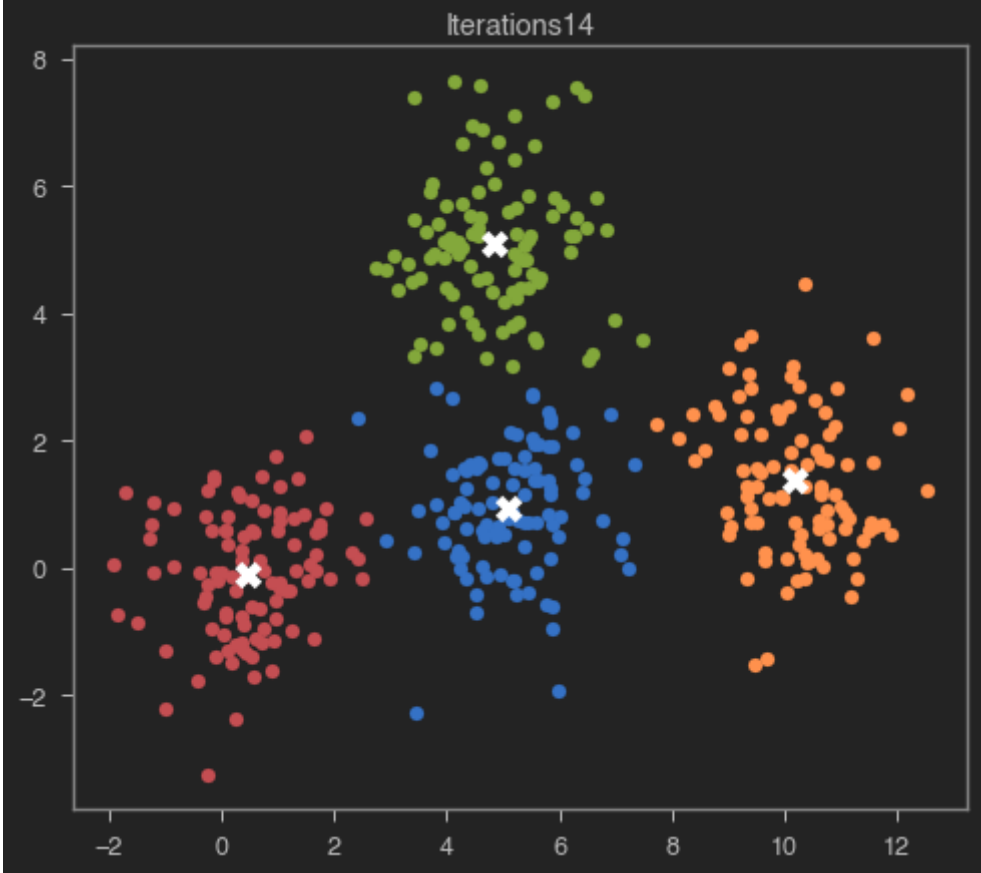
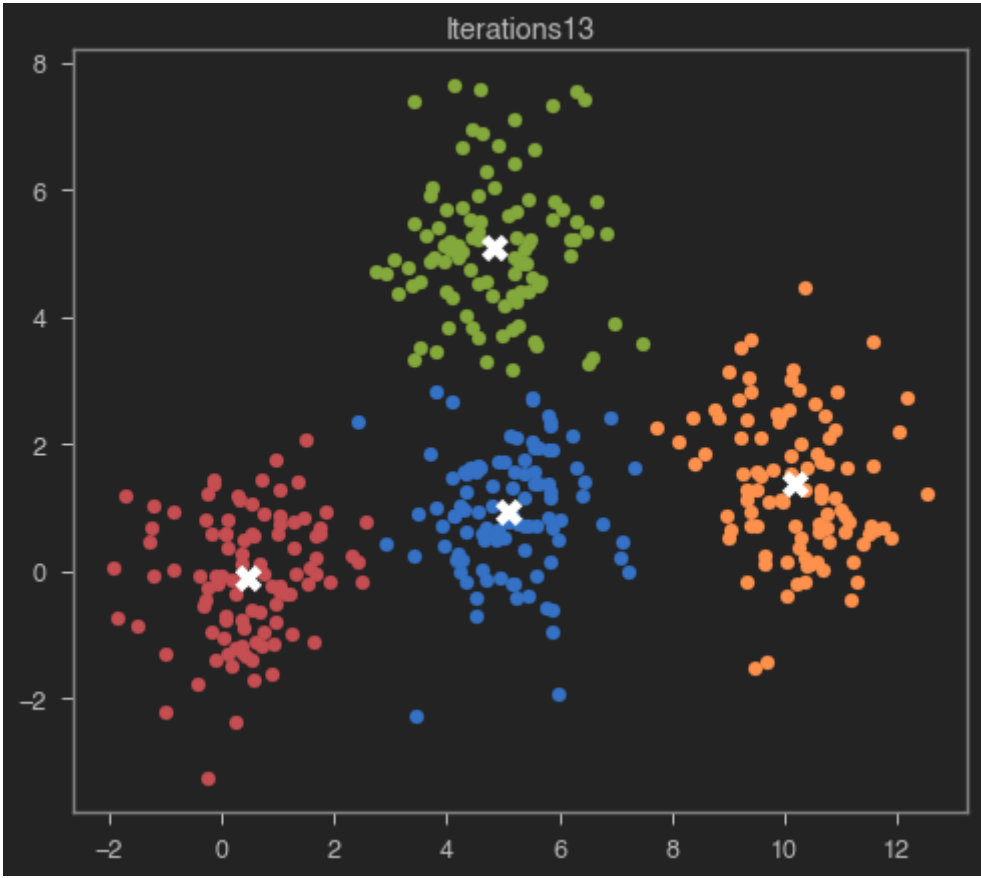


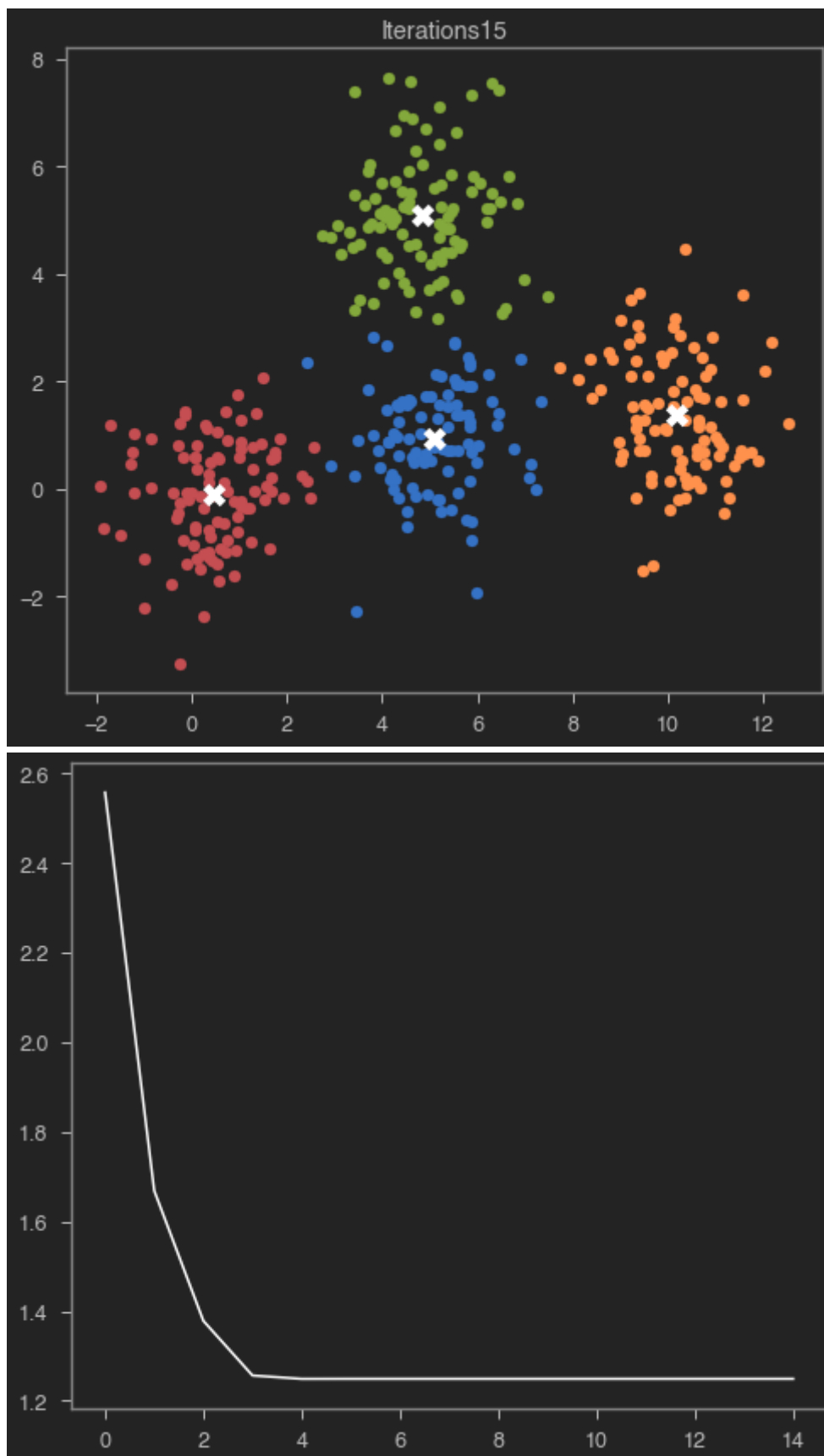












Step 4 : Performance metric

Compute Homogeneity score and Silhouette coefficient using the information given below

Homogeneity score : A clustering result satisfies homogeneity if all of its clusters contain only data points which are members of a single class. This metric is independent of the absolute values of the labels: a permutation of the class or cluster label values won't change the score value in any way.

Silhouette coefficient :

$a(x)$: Average distance of x to all other vectors in same cluster

$b(x)$: Average distance of x to the vectors in other clusters. Find minimum among the clusters

$$s(x) = \frac{b(x) - a(x)}{\max(a(x), b(x))}$$

Silhouette coefficient (SC) :

$$SC = \frac{1}{N} \sum_{i=1}^N s(x)$$

In [5]:

```
# write your code here
a_x = np.array([])
b_x = np.array([])
means = np.array([])

for i in range(4):
    for j in range(len(clstr_ass[i][0,:])):
        b = np.array([])
        a_x = np.append(a_x, np.mean(np.sqrt((clstr_ass[i][0,:]-
clstr_ass[i][0,j])**2 + (clstr_ass[i][1,:]-clstr_ass[i]
[1,j])**2)))

        b_x = np.append(b_x, np.mean(np.sqrt((clstr_ass[(i+1)%4]
[0,:]-clstr_ass[i][0,j])**2 + (clstr_ass[(i+1)%4][1,:]-
clstr_ass[i][1,j])**2)))

        b_x = np.append(b_x, np.mean(np.sqrt((clstr_ass[(i+2)%4]
[0,:]-clstr_ass[i][0,j])**2 + (clstr_ass[(i+2)%4][1,:]-
clstr_ass[i][1,j])**2)))

        b_x = np.append(b_x, np.mean(np.sqrt((clstr_ass[(i+3)%4]
[0,:]-clstr_ass[i][0,j])**2 + (clstr_ass[(i+3)%4][1,:]-
clstr_ass[i][1,j])**2)))

        b = np.append(b, np.min(b_x))
s = (b-a_x)/(np.maximum(b,a_x))
s = np.mean(s)
print(s)
```

0.2854734349833749

Gaussian Mixture Models Clustering

Gaussian mixture model is an unsupervised machine learning method. It summarizes a multivariate probability density function with a mixture of Gaussian probability distributions as its name suggests. It can be used for data clustering and data mining. In this lab, GMM is used for clustering.

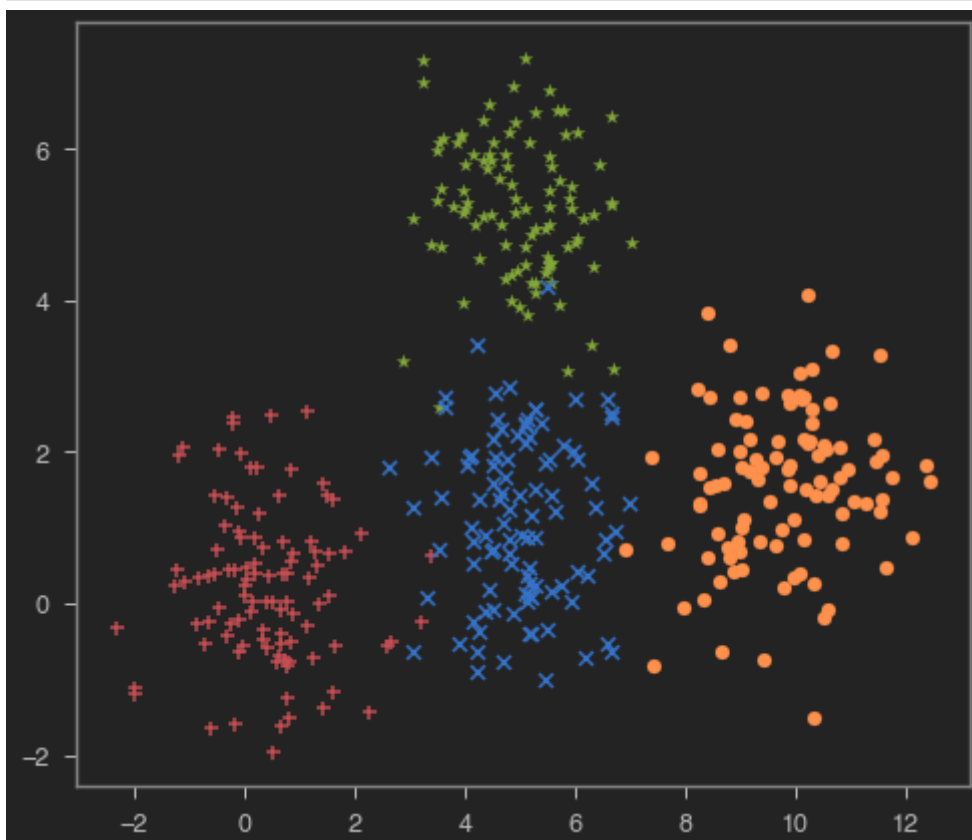
Step 1: Data generation

a) Follow the same steps as in K-means Clustering to generate the data

In [6]:

```
# write your code here
def dataset(num,mean,cov):
    data = np.random.multivariate_normal(mean, cov, num)
    return(data.T[0],data.T[1])

datax = np.array([])
datay = np.array([])
mu = np.array([[0.5,0],[5,5],[5,1],[10,1.5]])
cov = np.array([[1,1],[1,1],[1,1],[1,1]])
clr = ['r','g','b','y']
mrk = ['+','*','x','o']
for i in range(0,4):
    mean = mu[i,:]
    covariance = np.array([[cov[i,0],0],[0,cov[i,1]]])
    x,y = dataset(100,mean,covariance)
    datax = np.append(datax,x)
    datay = np.append(datay,y)
    plt.scatter(x,y,color = clr[i],marker = mrk[i])
```



Step 2. Initialization

- a) Mean vector (randomly any from the given data points) (μ_k)
- b) Covariance (initialize with (identity matrix)*max(data)) (Σ_k)
- c) Weights (uniformly) (w_k), with constraint: $\sum_{k=1}^K w_k = 1$

In [7]:

```
def initialization(data,K):

    # write your code here

    return theta
```

File "/var/folders/ls/mdl2q7410zj32bxvfvf078ff40000gn/T/ipykernel_20495/226370346.py", line 5

```
    return theta
    ^
```

IndentationError: expected an indented block

Step 3: Expectation stage

$$\gamma_{ik} = \frac{w_k P(x_i | \Phi_k)}{\sum_{k=1}^K w_k P(x_i | \Phi_k)}$$

where,

$$\Phi_k = \{\mu_k, \Sigma_k\}$$

$$\theta_k = \{\Phi_k, w_k\}$$

$$w_k = \frac{N_k}{N}$$

$$N_k = \sum_{i=1}^N \gamma_{ik}$$

$$P(x_i | \Phi_k) = \frac{1}{(2\pi)^{d/2} |\Sigma_k|^{1/2}} e^{-(x_i - \mu_k)^T \Sigma_k^{-1} (x_i - \mu_k)}$$

In []:

```
# E-Step GMM
from scipy.stats import multivariate_normal

def E_Step_GMM(data,K,theta):

    # write your code here

    return responsibility
```

Step 4: Maximization stage

$$a) w_k = \frac{N_k}{N}, \text{ where } N_k = \sum_{i=1}^N \gamma_{ik}$$

$$b) \mu_k = \frac{\sum_{i=1}^N \gamma_{ik} x_i}{N_k}$$

$$c) \Sigma_k = \frac{\sum_{i=1}^N \gamma_{ik} (x_i - \mu_k)(x_i - \mu_k)^T}{N_k}$$

Objective function(maximized through iteration):

$$L(\theta) = \sum_{i=1}^N \log \sum_{k=1}^K w_k P(x_i | \Phi_k)$$

```
In [ ]: # M-STEP GMM

def M_Step_GMM(data, responsibility):

    # write your code here

    return theta, log_likelihood
```

Step 5: Final run (EM algorithm)

a) Initialization

b) Iterate E-M untill $L(\theta_n) - L(\theta_{n-1}) \leq th$

c) Plot and see the cluster allocation at each iteration

```
In [ ]: log_l=[]
Itr=50
eps=10**(-14) # for threshold
clr=['r','g','b','y','k','m','c']
mrk=['+','*','x','o','.','<','p']

K = 4 # no. of clusters

theta=initialization(data,K)

for n in range(Itr):

    responsibility=E_Step_GMM(data,K,theta)

    cluster_label=np.argmax(responsibility,axis=1) #Label Points

    theta,log_likh=M_Step_GMM(data,responsibility)
```



```

log_l.append(log_likhd)

plt.figure()
for l in range(K):
    id=np.where(cluster_label==l)
    plt.plot(data[id,0],data[id,1],'.',color=clr[l],marker=mrk[l])
Cents=theta[0].T
plt.plot(Cents[:,0],Cents[:,1],'X',color='k')
plt.title('Iteration= %d' % (n))

if n>2:
    if abs(log_l[n]-log_l[n-1])<eps:
        break

plt.figure()
plt.plot(log_l)

```

Step 6 : Performance metric

Compute Homogeneity score and Silhouette coefficient using the information given below

Homogeneity score : A clustering result satisfies homogeneity if all of its clusters contain only data points which are members of a single class. This metric is independent of the absolute values of the labels: a permutation of the class or cluster label values won't change the score value in any way.

Silhouette coefficient :

$a(x)$: Average distance of x to all other vectors in same cluster

$b(x)$: Average distance of x to the vectors in other clusters. Find minimum among the clusters

$$s(x) = \frac{b(x)-a(x)}{\max(a(x),b(x))}$$

Silhouette coefficient (SC) :

$$SC = \frac{1}{N} \sum_{i=1}^N s(x)$$

In []: # write your code here

GMM v/s K-means

(a) Generate Data to show shortcomings of Kmeans and advantage of GMM over it

(b) Perform GMM on the same data and justify how it is better than K-means in that particular case

(c) Verify the same using performance metrics

In []:

```
# write your code here
```

Practical Use Case : K-means Clustering

For this exercise we will be using the **IRIS FLOWER DATASET** and explore how K-means clustering is performing

IRIS Dataset consists of 50 samples from each of the three species of Iris flower (Iris Setosa, Iris Virginica and Iris Versicolor)

Four features were measured from each sample : Length of Sepals, Width of sepals, Length of Petals, Width of Sepals all in centimeters. Based on the combinations of these 4 features each flower was categorized into one of the 3 species

Steps :

(a) Convert the given iris.csv file into a Pandas Dataframe, then extract both feature vector and target vector

(b) Perform analysis of Dataset, Plot the following features : (Sepal Length vs Sepal Width), (Petal Length vs Petal Width)

(c) Next group the data points into 3 clusters using the above K-means Clustering algorithm and compare the performance against the true labels obtained by the target vector, Also explain the results using a Confusion matrix

(d) Next use scikit learn tool to perform K-means Clustering and compare the performance against the true labels obtained by the target vector, Also explain the results using a Confusion matrix

(e) Vary the Number of Clusters (K) and run K-means algorithm from 1-10 and find the optimal number of clusters

In [8]:

```
## write your code here

import numpy as np
import pandas as pd
import matplotlib as mpl
import matplotlib.pyplot as plt
from jupyterthemes import jtplot
jtplot.style(theme = "monokai", context = "notebook", ticks =
True, grid = False)
import warnings
warnings.filterwarnings("ignore")
# show plots inline
```

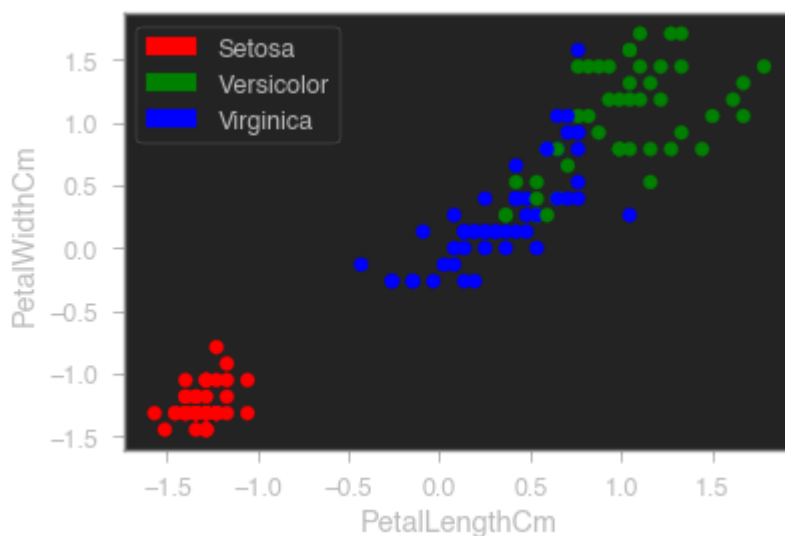
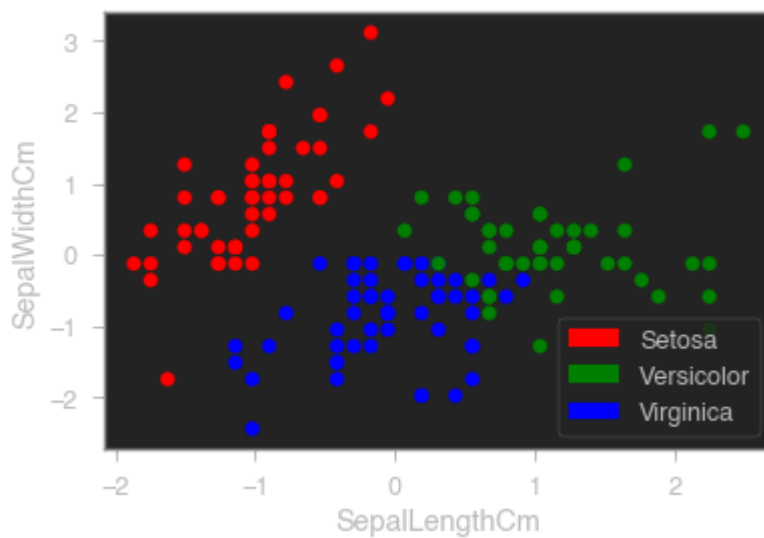
```
%matplotlib inline
data = pd.read_csv('Iris.csv')
data = data.drop('Id', axis=1) # get rid of the Id column - don't need it
data.sample(5)
X = data.iloc[:,0:4]
y = data.iloc[:, -1]

from sklearn import preprocessing
scaler = preprocessing.StandardScaler()
scaler.fit(X)
X_scaled_array = scaler.transform(X)
X_scaled = pd.DataFrame(X_scaled_array, columns = X.columns)
X_scaled.sample(5)

from sklearn.cluster import KMeans
nclusters = 3 # this is the k in kmeans
seed = 0
km = KMeans(n_clusters=nclusters, random_state=seed)
km.fit(X_scaled)
# predict the cluster for each data point
y_cluster_kmeans = km.predict(X_scaled)
y_cluster_kmeans

import matplotlib.patches as mpatches
red_patch = mpatches.Patch(color='red', label='Setosa')
green_patch = mpatches.Patch(color='green', label='Versicolor')
blue_patch = mpatches.Patch(color='blue', label='Virginica')
colors = np.array(['blue', 'red', 'green'])
plt.figure()
plt.scatter(X_scaled.iloc[:, 0],X_scaled.iloc[:,
1],c=colors[y_cluster_kmeans])
plt.xlabel("SepalLengthCm")
plt.ylabel("SepalWidthCm")
plt.legend(handles=[red_patch, green_patch, blue_patch])
plt.figure()
plt.scatter(X_scaled.iloc[:, 2],X_scaled.iloc[:,
3],c=colors[y_cluster_kmeans])
plt.xlabel("PetalLengthCm")
plt.ylabel("PetalWidthCm")
plt.legend(handles=[red_patch, green_patch, blue_patch])
plt.show()
```

```
from sklearn import metrics
score = metrics.silhouette_score(X_scaled, y_cluster_kmeans)
print("Coefficient:", score)
```



Coefficient: 0.4589717867018717

Practical Use Case : GMM

Steps :

- Convert the given iris.csv file into a Pandas Dataframe, then extract both feature vector and target vector
- Next group the data points into 3 clusters using the above GMM Clustering algorithm and compare the performance against the true labels obtained by the target vector, Also explain the results using a Confusion matrix
- Next use scikit learn tool to perform GMM Clustering and compare the performance against the true labels obtained by the target vector, Also explain the results using a Confusion matrix

In []: `# write your code here`