

EDS ASSIGNMENT 06 & Minor Project

Prof. Pushpmala Shinde
Mam

Project Members

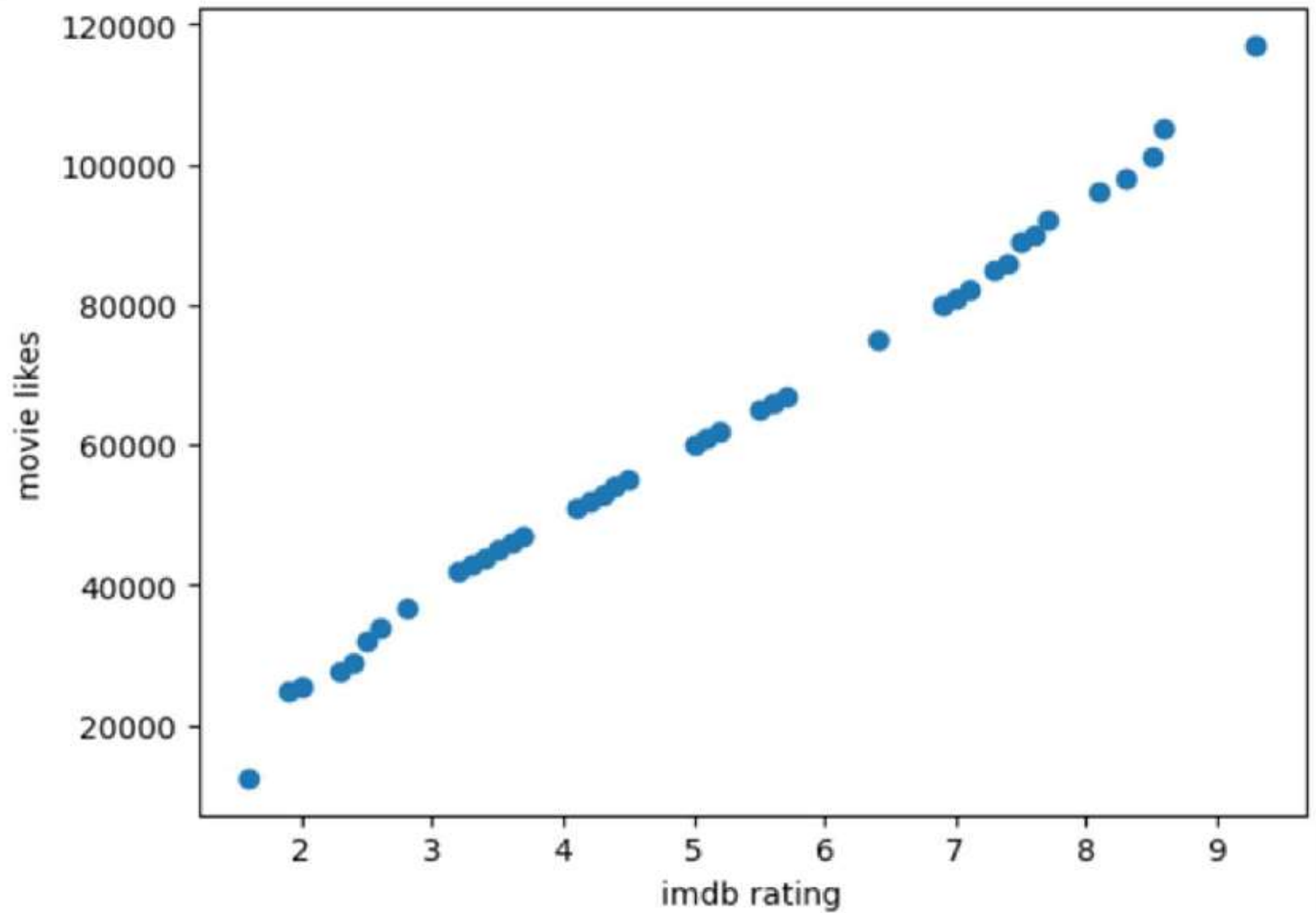
- 1. Rohan Budgujar 502**
- 2. Varun Balbudhe 503**
- 3. Abhay Isal 519**



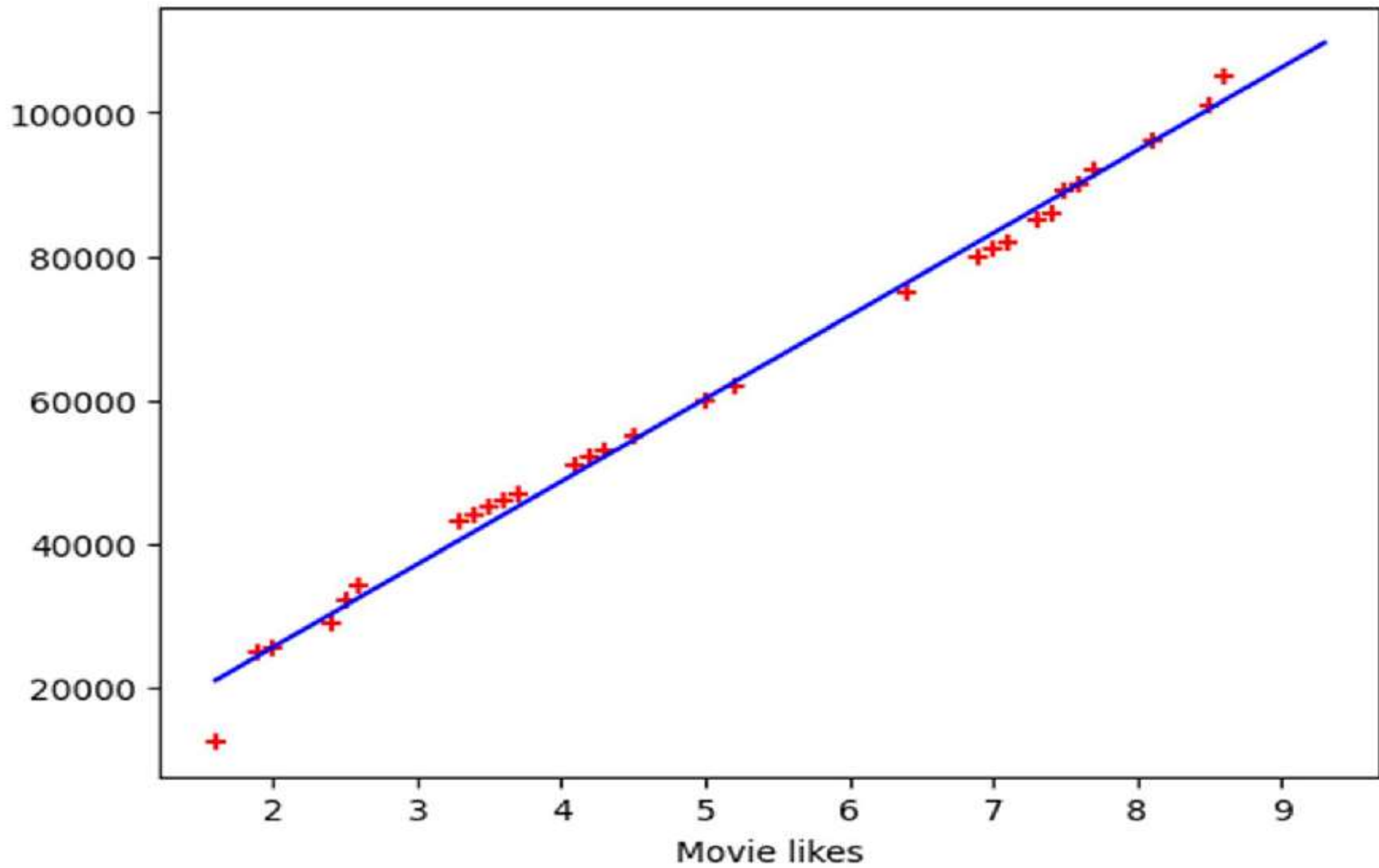
ESSENTIALS OF DATA SCIENCE

Using linear regression

```
1  import numpy as np
2  import pandas as pd
3  import matplotlib.pyplot as plt
4  from sklearn import linear_model
5  from sklearn.model_selection import train_test_split
6  df = pd.read_csv("/content/movie_data.csv")
7
8  #data cleaning df.dropna(inplace=True)
9  df.reset_index(drop=True, inplace=True)
10 df1 = df.head(40)
11 # print(df1)
12 plt.scatter(df1['imdb_score'], df1['movie_likes'])
13 plt.xlabel('imdb rating')
14 plt.ylabel('movie likes')
```



```
16 X = np.array(df1[['imdb_score']]).reshape(-1,1)
17 Y = np.array(df1[['movie_likes']]).reshape(-1,1)
18 X_train,X_test,Y_train,Y_test = train_test_split(X,Y,test_size = 0.25)
19
20 # create linear regression object
21 reg = linear_model.LinearRegression()
22 reg.fit(X_train, Y_train) #training the model
23 # predicting movie likes using the testing dataset on the trained model
24 reg.predict(X_test)
25 # plotting linear regression line
26 plt.scatter(X_train, Y_train, color='red', marker='+')
27 plt.xlabel('IMDB')
28 plt.xlabel('Movie likes')
29 plt.plot(df1['imdb_score'],
30 reg.predict(df1[['imdb_score']]),
31 color='blue')
```

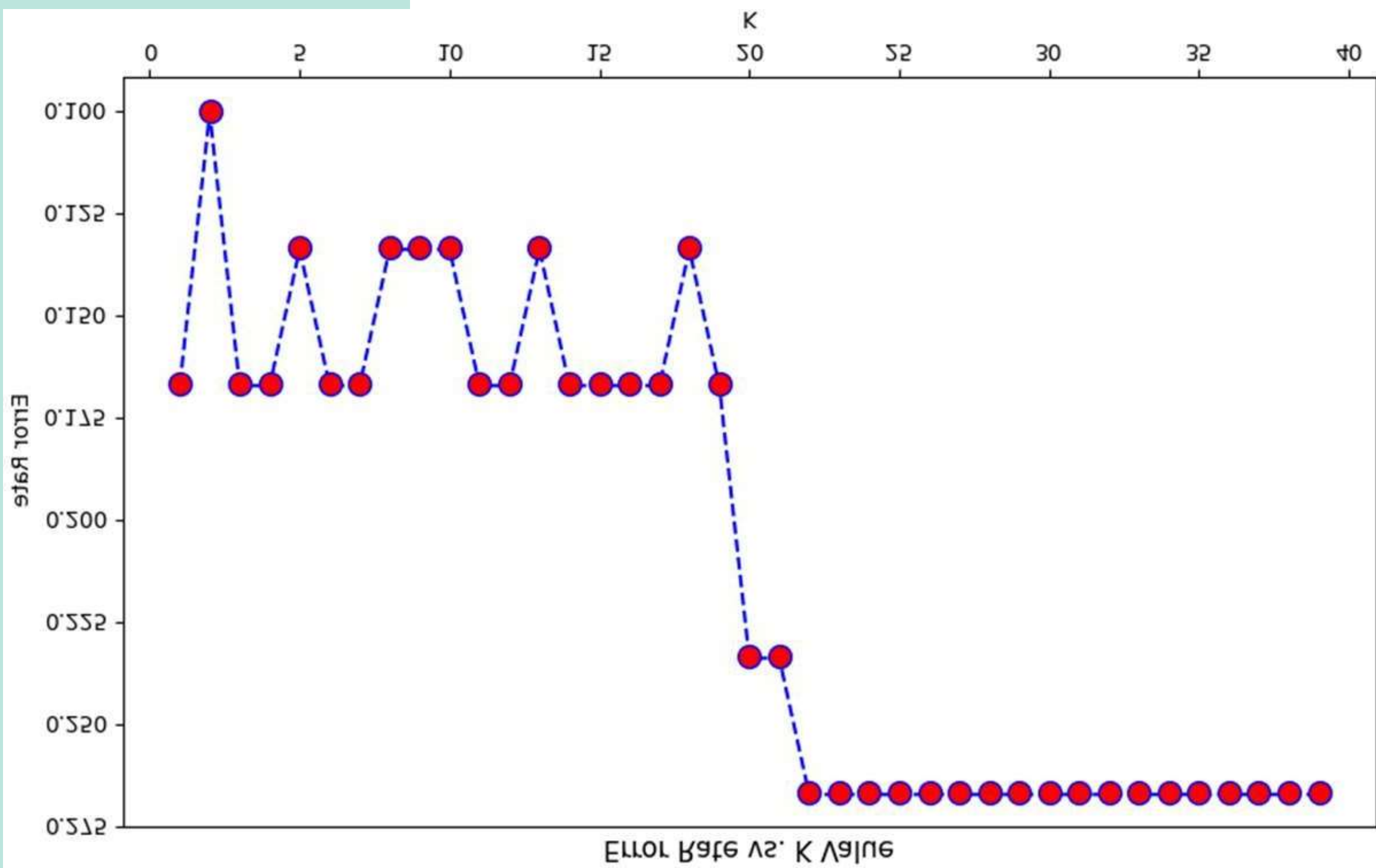


Using KNN

```
1  import pandas as pd
2  import seaborn as sns
3  import matplotlib.pyplot as plt
4  import numpy as np
5  df = pd.read_csv("prostate.csv")
6  df.head()
7  from sklearn.preprocessing import StandardScaler
8  scaler = StandardScaler()
9  scaler.fit(df.drop('Target', axis=1))
10 scaled_features = scaler.transform(df.drop('Target', axis=1))
11
```

```
12
13 df_feat = pd.DataFrame(scaled_features, columns=df.columns[:-1])
14 df_feat.head()
15 from sklearn.metrics import classification_report
16 from sklearn.neighbors import KNeighborsClassifier
17 from sklearn.model_selection import train_test_split
18 X_train, X_test, \
19 y_train, y_test = train_test_split(scaled_features,
20 df['Target'], test_size=0.30)
21 # Remember that we are trying to come up # with a model to predict whether
22 # someone will Target or not. # We'll start with k = 1.
23 knn = KNeighborsClassifier(n_neighbors=1)
24 knn.fit(X_train, y_train)
25 pred = knn.predict(X_test)
26 # Predictions and Evaluations
27 # Let's evaluate our KNN model !
28 print(confusion_matrix(y_test, pred))
29 print(classification_report(y_test, pred))
30 error_rate = []
31 # Will take some time for i in range(1, 40):
32 knn = KNeighborsClassifier(n_neighbors=i)
33 knn.fit(X_train, y_train)
34 pred_i = knn.predict(X_test)
35 error_rate.append(np.mean(pred_i != y_test))
```

```
29 print(classification_report(y_test, pred))
30 error_rate = []
31 # Will take some time
32 for i in range(1, 40):
33     knn = KNeighborsClassifier(n_neighbors=i)
34     knn.fit(X_train, y_train)
35     pred_i = knn.predict(X_test)
36     error_rate.append(np.mean(pred_i != y_test))
37
38 plt.figure(figsize=(10, 6))
39 plt.plot(range(1, 40), error_rate, color='blue', linestyle='dashed', marker='o',
40         markerfacecolor='red', markersize=10)
41
42 plt.title('Error Rate vs. K Value')
43 plt.xlabel('K')
44 plt.ylabel('Error Rate')
45 plt.show()
```

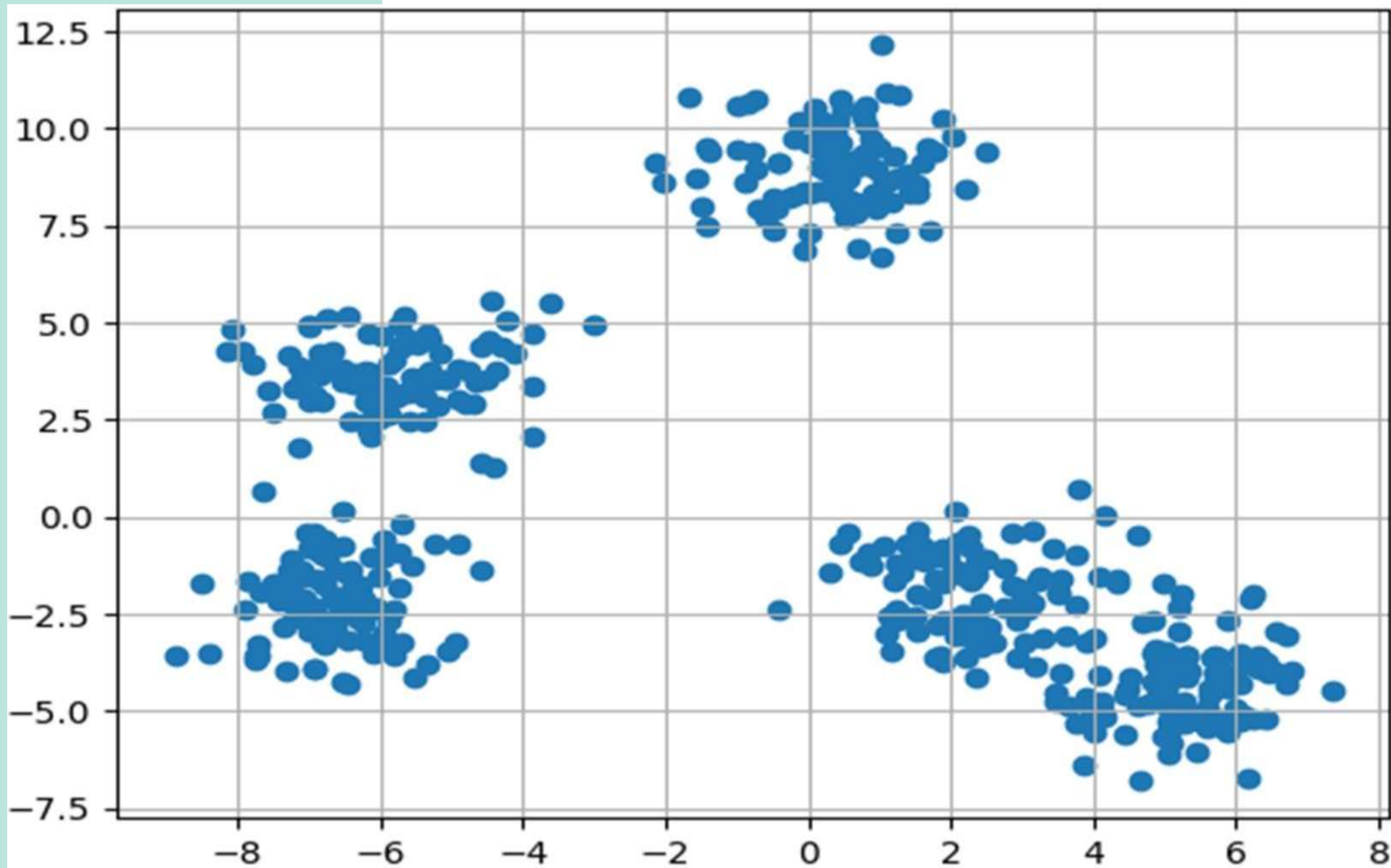



```
47 # FIRST A QUICK COMPARISON TO OUR ORIGINAL K = 1
48 knn = KNeighborsClassifier(n_neighbors = 1)
49
50 knn.fit(X_train, y_train)
51 pred = knn.predict(X_test)
52 print('WITH K = 1')
53 print('Confusion Matrix')
54 print(confusion_matrix(y_test, pred))
55 print('Classification Report')
56 print(classification_report(y_test, pred))
57
58 # NOW WITH K = 10
59 knn = KNeighborsClassifier(n_neighbors = 10)
60 knn.fit(X_train, y_train)
61 pred = knn.predict(X_test)
62 print('WITH K = 10')
63 print('Confusion Matrix')
64 print(confusion_matrix(y_test, pred))
65 print('Classification Report')
66 print(classification_report(y_test, pred))
```

Using K Means

```
1 import matplotlib.pyplot as plt
2 from sklearn.datasets import make_blobs
3 X,y = make_blobs(n_samples = 500,n_features = 2,centers = 5,random_state = 23)
4 fig = plt.figure(0)
5 plt.grid(True)
6 plt.scatter(X[:,0],X[:,1])
7 plt.show()
```

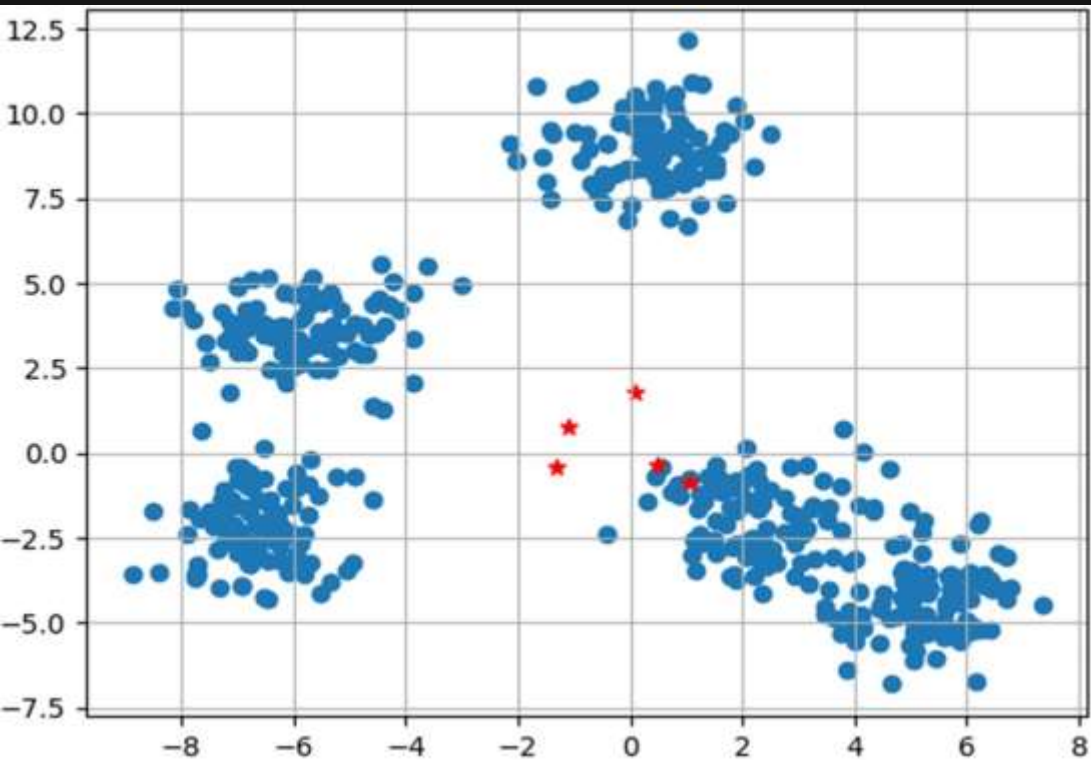
8





```
9 import numpy as np
10 clusters = {}
11 np.random.seed(23)
12
13 for idx in range(k):
14     center = 2*(2*np.random.random((X.shape[1],))-1)
15     points = []
16     cluster = {
17         'center' : center,
18         'points' : []
19     }
20     clusters[idx] = cluster
21
22
23 clusters
```

```
{0: {'center': array([0.06919154, 1.78785042]), 'points': []},
 1: {'center': array([ 1.06183904, -0.87041662]), 'points': []},
 2: {'center': array([-1.11581855,  0.74488834]), 'points': []},
 3: {'center': array([-1.33144319, -0.43023013]), 'points': []},
 4: {'center': array([ 0.47220939, -0.35227962]), 'points': []}}
```

```
23 clusters
24 plt.scatter(X[:,0],X[:,1])
25 plt.grid(True)
26 for i in clusters:
27     center = clusters[i]['center']
28     plt.scatter(center[0],center[1],marker = '*',c = 'red')
29 plt.show()
30
```



#Implementing E step

```
def assign_clusters(X, clusters):  
    for idx in range(X.shape[0]):  
        dist = []  
        curr_x = X[idx]  
         for i in range(k):  
            dis = distance(curr_x, clusters[i]['center'])  
            dist.append(dis)  
            curr_cluster = np.argmin(dist)  
            clusters[curr_cluster]['points'].append(curr_x)  
    return clusters
```

#Implementing the M-Step

```
def update_clusters(X, clusters):  
    for i in range(k):  
        points = np.array(clusters[i]['points'])  
        if points.shape[0] > 0:  
            new_center = points.mean(axis = 0)  
            clusters[i]['center'] = new_center
```

```
clusters[i]['points'] = []  
return clusters
```

```
def pred_cluster(X, clusters):  
    pred = []  
    for i in range(X.shape[0]):  
        dist = []  
        for j in range(k): dist.append(distance(X[i], clusters[j]['center']))  
        pred.append(np.argmin(dist))  
    return pred
```

```
62 clusters = assign_clusters(X,clusters)
63 clusters = update_clusters(X,clusters)
64 pred = pred_cluster(X,clusters)
65
66 plt.scatter(X[:,0],X[:,1],c = pred)
67 ✓ for i in clusters:
68     center = clusters[i]['center']
69     plt.scatter(center[0],center[1],marker = '^',c = 'red')
70     plt.show()
71
```