

```

1 # k-Means demo
2 import numpy as np
3 import matplotlib.pyplot as plt

1 # Set the number of data vectors (n) and the dimension of the data space (m)
2 n = 12 # Example 1
3 # n = 1000 # random data
4 m = 2
5
6 # Set the number of clusters (k)
7 k = 5 # changed from 4 to 5 for q1 part b
8
9 # Initialize the data - either as in Example 1 or using random data
10 XData = np.load("blobs.npy")
11 n = len(XData) # Note: Added for Q1 part b

1 # Create data structures to store the (randomly selected) representative vectors for cluster (c)
2 c = np.vstack([np.random.uniform(-1,1,k), np.random.uniform(-1,1,k)]).T #np.zeros(shape=(k, m))

```

## ✓ Question 1a part ii

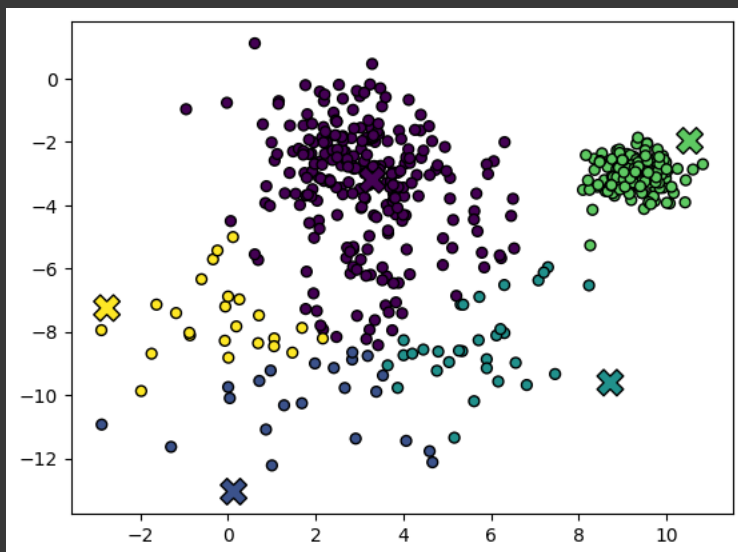
```

1 # Question 1a part ii:
2
3 # Pick random point for first rep vec
4 c[0] = XData[np.random.choice(n)]
5
6
7 for i in range(1, k):
8     # Find distance between each data point and its nearest rep vector
9     nearest_dist = np.array([min(np.linalg.norm((c[j] - xD), ord=2) for j in range(i)) for xD in XData])
10
11     # Pick next rep vec as the point with the highest dist
12     next_rep_index = np.argmax(nearest_dist)
13     c[i] = XData[next_rep_index]

1 # Create a data structure to store closest representative vector for each data point
2 closestCluster = np.zeros(n)
3 # Assign each data vector to the new, closest cluster
4 for d in range(n):
5     # Store the coordinates of the current data vector
6     xD = XData[d, :]
7
8     # Set the minimum distance tracker to be a very large number
9     sqDistMin = 1e16
10
11     # Find the closest representative vector (cluster) to the current data vector
12     for i in range(k):
13         sqDist = np.linalg.norm(c[i, :] - xD, ord=2)
14
15         # If the distance is less than the current min, assign the
16         # current data vector to this cluster
17         if sqDist < sqDistMin:
18             closestCluster[d] = i
19             sqDistMin = sqDist

1 # Update the assignments of the data vectors to their new clusters
2 IndexSet = closestCluster.astype(int)
3
4 # Plot the data
5 plt.scatter(XData[:, 0], XData[:, 1], s=32, c=IndexSet, cmap='viridis', edgecolors='k', marker='o')
6 plt.scatter(c[:, 0], c[:, 1], s=200, c=np.linspace(1, k, k), cmap='viridis', edgecolors='k', marker='X')
7 plt.show()

```



```

1 # Create data structures to store the representative vectors from the previous iteration (cPrev)
2 cPrev = np.copy(c)
3
4 # The Alternating Minimization Scheme
5 doneFlag = False
6
7 # Keep alternating updates to representative vectors and cluster assignments until representative vectors no longer change their locations
8 while not doneFlag:
9     # Update the representative vectors in each cluster via the centroid formula
10    for i in range(k):
11        # Find the indices for all data vectors currently in cluster i
12        ClusterIndices = np.where(IndexSet == i)[0]
13
14        # Find the number of data vectors currently in cluster i
15        NumVecsInCluster = len(ClusterIndices)
16
17        # Create a data structure to store representative vector for the current cluster
18        c[i, :] = np.zeros(m)
19
20        # Update cluster vector using the centroid formula
21        for j in range(NumVecsInCluster):
22            c[i, :] += XData[ClusterIndices[j], :] / NumVecsInCluster
23
24    # Plot the updated representative vectors for each cluster
25    plt.scatter(XData[:, 0], XData[:, 1], s=64, c=IndexSet, cmap='viridis', edgecolors='k', marker='o')
26    plt.scatter(c[:, 0], c[:, 1], s=200, c=np.linspace(1, k, k), cmap='viridis', edgecolors='k', marker='X')
27    plt.show()
28
29    # Now reassign all data vectors to the closest representative vector (cluster)
30    # Create a data structure to store closest representative vector for each data point
31    closestCluster = np.zeros(n)
32
33    # Reassign each data vector to the new, closest cluster
34    for d in range(n):
35        # Store the coordinates of the current data vector
36        xD = XData[d, :]
37
38        # Set the minimum distance tracker to be a very large number
39        sqDistMin = 1e16
40
41        # Find the closest representative vector (cluster) to the current data vector
42        for i in range(k):
43            sqDist = np.linalg.norm(c[i, :] - xD, ord=2)
44
45            # If the distance is less than the current min, assign the
46            # current data vector to this cluster
47            if sqDist < sqDistMin:
48                closestCluster[d] = i
49                sqDistMin = sqDist
50
51    # Update the assignments of the data vectors to their new clusters
52    IndexSet = closestCluster.astype(int)
53
54    # Plot the data and the updated representative vectors

```

```
54 # Plot the data and the updated representative vectors
55 plt.scatter(XData[:, 0], XData[:, 1], s=64, c=IndexSet, cmap='viridis', edgecolors='k', marker='o')
56 plt.scatter(c[:, 0], c[:, 1], s=200, c=np.linspace(1, k, k), cmap='viridis', edgecolors='k', marker='X')
57 plt.show()
58
59 # Terminate the alternating scheme if the representative vectors are unaltered
60 # relative to the previous iteration
61 if np.array_equal(c, cPrev):
62     doneFlag = True
63 else:
64     cPrev = np.copy(c)
```

