

```

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

1 # n = 1000 # random data
2
3 m = 2
4
5 # Set the number of clusters (k)
6 k = 5
7
8 XData = np.load("blobs.npy")
9 n = len(XData) # is 500 for blobs.npy
10
11 num_realizations = 10

```

Question 1 part c:

```

1 # Random initialization
2 def rand_init(X, k):
3     c = np.vstack([np.random.uniform(-1,1,k), np.random.uniform(-1,1,k)]).T #np.zeros(shape=(k, m))
4     return c
5
6 # k++ initialization
7 def k_plus(X, k):
8     c = np.vstack([np.random.uniform(-1,1,k), np.random.uniform(-1,1,k)]).T
9     c[0] = XData[np.random.choice(n)]
10
11     for i in range(1, k):
12         # Find distance between each data point and its nearest rep vector
13         nearest_dist = np.array([min(np.linalg.norm((c[j] - xD), ord=2) for j in range(i)) for xD in XData])
14
15         # Pick next rep vec as the point with the highest dist
16         next_rep_index = np.argmax(nearest_dist)
17         c[i] = XData[next_rep_index]
18     return c
19
20 # Function calculates coherence
21 def coherence(X, centroids, labels):
22     coherence = 0
23     for i, centroid in enumerate(centroids):
24         indices = np.where(labels == i)[0]
25         coherence += np.sum(np.linalg.norm((X[indices] - centroid), axis=1, ord=2))
26     return coherence

```

```

1 random_coherence = [] # coherence of random init per trial
2 k_plus_coherence = [] # coherence of k++ init per trial
3
4 for _ in range(num_realizations): # Loops 10x in our case
5
6     # Random init:
7     random_centroids = rand_init(XData, k)
8     closestCluster = np.zeros(len(XData))
9     for d in range(len(XData)):
10         xD = XData[d, :]
11         sqDistMin = 1e16
12         for i in range(k):
13             sqDist = np.linalg.norm(random_centroids[i, :] - xD, ord=2)
14             if sqDist < sqDistMin:
15                 closestCluster[d] = i
16                 sqDistMin = sqDist
17     IndexSet = closestCluster.astype(int) # our 'labels' param to be passed
18     random_coherence.append(coherence(XData, random_centroids, IndexSet))
19
20     # k++ init
21     k_plus_centroids = k_plus(XData, k)
22     closestCluster = np.zeros(len(XData))
23     for d in range(len(XData)):
24         xD = XData[d, :]
25         sqDistMin = 1e16
26         for i in range(k):
27             sqDist = np.linalg.norm(k_plus_centroids[i, :] - xD, ord=2)
28             if sqDist < sqDistMin:

```

```

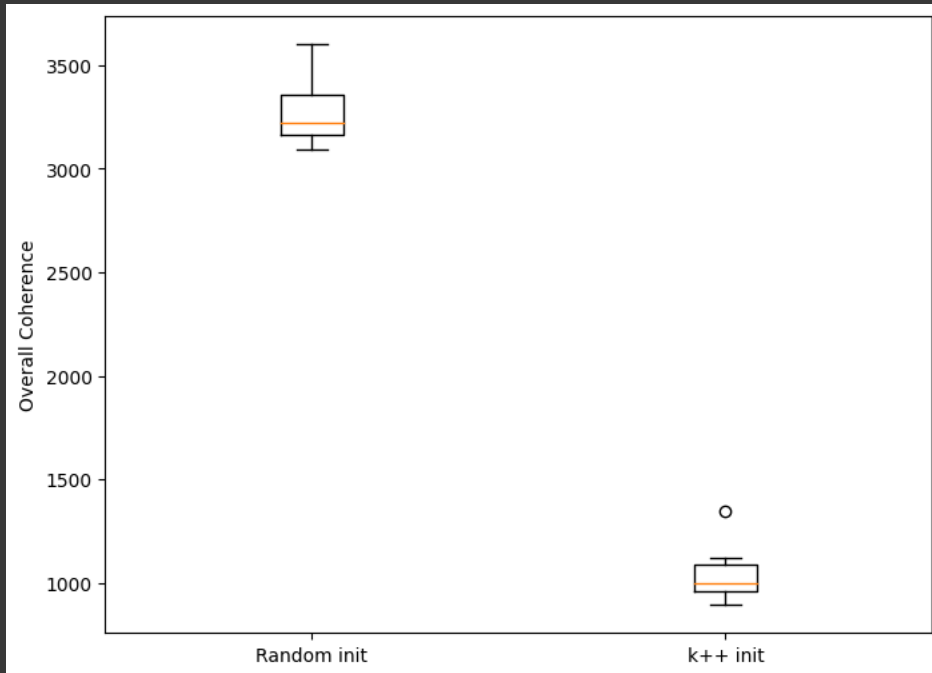
29         closestCluster[d] = i
30         sqDistMin = sqDist
31     IndexSet = closestCluster.astype(int) # labels param
32     k_plus_coherence.append(coherence(XData, k_plus_centroids, IndexSet))
33

```

```

1 # Plots
2 plt.figure(figsize=(8, 6))
3 plt.boxplot([random_coherence, k_plus_coherence], labels=['Random init', 'k++ init'])
4 plt.ylabel('Overall Coherence')
5 plt.show()

```



```

1 random_avg_coherence = np.mean(random_coherence)
2 k_plus_avg_coherence = np.mean(k_plus_coherence)
3 print("random init average coherence: ", random_avg_coherence)
4 print("k++ init average coherence: ", k_plus_avg_coherence)

```

```

random init average coherence: 3279.542290246182
k++ init average coherence: 1039.5696601560392

```

Of course, k++ initialization had a much lower coherence level, which makes sense, because the nodes are initialized to lie on an exact data point, while being spread out from one another. However, the downside is the runtime of the initialization itself. At least in my implementation, a triple nested for loop was used, which is very inefficient.