# Unsupervised Learning - Coding Practice Questions

CM52054: Foundational Machine Learning
*Practice set with fully worked answers*

1) **Implement k-means from scratch.**

   Implement k-means clustering with:

   - centroid initialisation by sampling *K distinct data points*,
   - Euclidean distance,
   - stopping when assignments stop changing *or* centroid movement is below a tolerance.

   Return (`centroids, labels, wcss`), where

   $$\text{WCSS (Within-Cluster Sum of Squares)} = \sum_{i=1}^{N} \|x_i - c_{\ell_i}\|_2^2.$$

   **Answer:**

```python
import numpy as np

def kmeans_from_scratch(X, K, max_iter=100, tol=1e-6, seed=0):
    """
    Simple, easy-to-follow k-means (NumPy only).

    Returns
    -------
    centroids : (K, M) array
    labels : (N,) array of ints in {0,...,K-1}
    wcss : float, sum of squared distances to assigned centroid
    """
    # ---------- Step 0: basic setup ----------
    X = np.asarray(X, dtype=float)
    N, M = X.shape
    rng = np.random.default_rng(seed)

    # ---------- Step 1: initialise centroids ----------
    # Pick K random data points as initial centroids
    centroids = X[rng.choice(N, size=K, replace=False)].copy()

    # ---------- Step 2: repeat assignment + update ----------
    for _ in range(max_iter):
        # (A) Assignment: compute distance to each centroid, pick
            nearest
        # dists_sq[i, k] = ||X[i] - centroids[k]||^2
        dists_sq = np.zeros((N, K))
        for k in range(K):
            diff = X - centroids[k] # (N, M)
            dists_sq[:, k] = np.sum(diff**2, axis=1)

        labels = np.argmin(dists_sq, axis=1) # (N,)
```

```
32
33          # (B) Update: recompute each centroid as mean of its assigned
                points
34          new_centroids = centroids.copy()
35          for k in range(K):
36              points_k = X[labels == k]
37              if len(points_k) > 0:
38                  new_centroids[k] = points_k.mean(axis=0)
39              else:
40                  # If a cluster becomes empty, re-pick a random data
                      point
41                  new_centroids[k] = X[rng.integers(0, N)]
42
43          # (C) Convergence: stop if centroids barely move
44          shift = np.linalg.norm(new_centroids - centroids)
45          centroids = new_centroids
46          if shift < tol:
47              break
48
49      # ---------- Step 3: compute WCSS ----------
50      # Sum of squared distances of each point to its assigned
            centroid
51      wcss = 0.0
52      for i in range(N):
53          diff = X[i] - centroids[labels[i]]
54          wcss += float(np.sum(diff**2))
55
56      return centroids, labels, wcss
57
58
59  if __name__ == "__main__":
60      X = np.array([[0.0, 0.0],
61                    [0.2, 0.1],
62                    [3.0, 3.0],
63                    [3.1, 2.9],
64                    [10.0, 10.0]])
65      centroids, labels, wcss = kmeans_from_scratch(X, K=2, seed=42)
66      print("Centroids:\n", centroids)
67      print("Labels:", labels)
68      print("WCSS:", wcss)
```

2) **Best-of-$n$ restarts for k-means (mitigate non-determinism).**

Implement a wrapper kmeans_best_of_n that:

- runs k-means n_init times with different seeds,
- returns the solution with the *lowest* WCSS.

**Answer:**

```
1  import numpy as np
2
3  def kmeans_best_of_n(X, K, n_init=10, max_iter=100, tol=1e-6,
       base_seed=0):
4      """
5      Run k-means multiple times and choose the result with the lowest
          WCSS.
6
```

```
7      Parameters
8      ----------
9      X : array-like, shape (N, M)
10     K : int
11     n_init : int
12         Number of independent initialisations.
13     max_iter, tol : forwarded to kmeans_from_scratch
14     base_seed : int
15         Base seed; each run uses base_seed + i.
16
17     Returns
18     -------
19     best_centroids : np.ndarray, shape (K, M)
20     best_labels : np.ndarray, shape (N,)
21     best_wcss : float
22     """
23     X = np.asarray(X, dtype=float)
24
25     best_centroids = None
26     best_labels = None
27     best_wcss = np.inf
28
29     # Run k-means with different seeds, inducing different initial
           centroid choices.
30     for i in range(n_init):
31         seed = base_seed + i
32
33         # Call the k-means implementation from Q1.
34         centroids, labels, wcss = kmeans_from_scratch(
35             X, K, max_iter=max_iter, tol=tol, seed=seed
36         )
37
38         # Keep the solution with the smallest objective value.
39         if wcss < best_wcss:
40             best_wcss = wcss
41             best_centroids = centroids
42             best_labels = labels
43
44     return best_centroids, best_labels, float(best_wcss)
45
46
47  if __name__ == "__main__":
48     rng = np.random.default_rng(0)
49
50     # Create two clusters for demonstration
51     X1 = rng.normal(loc=(0, 0), scale=0.5, size=(50, 2))
52     X2 = rng.normal(loc=(5, 5), scale=0.5, size=(50, 2))
53     X = np.vstack([X1, X2])
54
55     centroids, labels, wcss = kmeans_best_of_n(X, K=2, n_init=20,
           base_seed=100)
56     print("Best WCSS:", wcss)
```

3) **Compute "good clustering" metrics (intra vs inter).**

Implement a function clustering_quality_metrics(X, labels, centroids) that computes:

(a) average intra-cluster distance: mean $\|x_i - c_{\ell_i}\|_2$,

(b) minimum inter-centroid distance: $\min_{k \neq k'} \|c_k - c_{k'}\|_2$.

**Answer:**

```python
import numpy as np

def clustering_quality_metrics(X, labels, centroids):
    """
    Simple, easy-to-follow clustering metrics:

    1) avg_intra: average distance from each point to its assigned
        centroid
    2) min_inter: minimum distance between any two centroids
    """
    X = np.asarray(X, dtype=float)
    labels = np.asarray(labels, dtype=int)
    centroids = np.asarray(centroids, dtype=float)

    N = X.shape[0]
    K = centroids.shape[0]

    # ------------------------------
    # 1) Average intra-cluster distance
    # ------------------------------
    # For each point i, compute distance to its assigned centroid.
    total = 0.0
    for i in range(N):
        c = centroids[labels[i]] # centroid of point i
        diff = X[i] - c
        dist = np.sqrt(np.sum(diff**2)) # Euclidean distance
        total += float(dist)
    avg_intra = total / N

    # ------------------------------
    # 2) Minimum inter-centroid distance
    # ------------------------------
    # Compute distance between every pair of centroids and take the
        smallest.
    min_inter = float("inf")
    for i in range(K):
        for j in range(i + 1, K):
            diff = centroids[i] - centroids[j]
            dist = np.sqrt(np.sum(diff**2))
            if dist < min_inter:
                min_inter = float(dist)

    return avg_intra, min_inter


if __name__ == "__main__":
    rng = np.random.default_rng(0)
    X = rng.normal(size=(20, 2))

    # Example: two fixed centroids and nearest-centroid labels
    centroids = np.array([[0.0, 0.0], [1.0, 1.0]])
    labels = np.zeros(X.shape[0], dtype=int)
    for i in range(X.shape[0]):
```

```
52        d0 = np.sum((X[i] - centroids[0])**2)
53        d1 = np.sum((X[i] - centroids[1])**2)
54        labels[i] = 0 if d0 <= d1 else 1
55
56     avg_intra, min_inter = clustering_quality_metrics(X, labels,
          centroids)
57     print("Average intra distance:", avg_intra)
58     print("Minimum inter-centroid distance:", min_inter)
```

4) **Optional: Agglomerative hierarchical clustering (single linkage) for small datasets.**
   Implement agglomerative clustering with single linkage:

   - start with each point as its own cluster,
   - repeatedly merge the pair of clusters with minimum single-link distance,
   - return a merge history list.

   Assume $N$ is small; clarity is more important than speed.

   **Answer:**

```
1  import numpy as np
2
3  def euclidean(a, b):
4      """Euclidean distance between two vectors a and b."""
5      return float(np.sqrt(np.sum((a - b) ** 2)))
6
7  def single_link_distance(X, cluster_a, cluster_b):
8      """
9      Single-link distance between two clusters:
10     the minimum distance between any point in cluster_a and any
          point in cluster_b.
11     """
12     best = float("inf")
13     for i in cluster_a:
14         for j in cluster_b:
15             d = euclidean(X[i], X[j])
16             if d < best:
17                 best = d
18     return best
19
20 def agglomerative_single_linkage(X):
21     """
22     Easy-to-follow agglomerative clustering (single linkage).
23
24     Returns
25     -------
26     merges : list of (a_id, b_id, dist, new_id)
27         a_id and b_id are the cluster IDs merged at distance dist,
28         producing a new cluster with ID new_id.
29     """
30     X = np.asarray(X, dtype=float)
31     N = X.shape[0]
32
33     # Start with N singleton clusters: {0}, {1}, ..., {N-1}
34     clusters = {i: {i} for i in range(N)} # cluster_id -> set of
          point indices
35     next_id = N
36     merges = []
```

```
37
38     # Keep merging until only one cluster remains
39     while len(clusters) > 1:
40         ids = list(clusters.keys())
41
42         # Find the closest pair of clusters
43         best_dist = float("inf")
44         best_pair = None
45
46         for i in range(len(ids)):
47             for j in range(i + 1, len(ids)):
48                 a_id = ids[i]
49                 b_id = ids[j]
50                 d = single_link_distance(X, clusters[a_id], clusters[
                        b_id])
51                 if d < best_dist:
52                     best_dist = d
53                     best_pair = (a_id, b_id)
54
55         # Merge the closest pair
56         a_id, b_id = best_pair
57         new_cluster = clusters[a_id] | clusters[b_id]
58
59         new_id = next_id
60         next_id += 1
61
62         # Record this merge event
63         merges.append((a_id, b_id, best_dist, new_id))
64
65         # Update active clusters
66         del clusters[a_id]
67         del clusters[b_id]
68         clusters[new_id] = new_cluster
69
70     return merges
71
72
73 if __name__ == "__main__":
74     X = np.array([
75         [0.0, 0.0],
76         [0.1, 0.0],
77         [2.0, 2.0],
78         [2.1, 2.0],
79     ])
80
81     merges = agglomerative_single_linkage(X)
82     for a_id, b_id, dist, new_id in merges:
83         print(f"Merged {a_id} and {b_id} at dist={dist:.3f} -> new
                cluster {new_id}")
```

6