# Bangladesh University of Engineering and Technology

**EEE 6608**

**Machine Learning & Pattern Recognition**

*Final assignement*: Implementing K-means clustering alogrithm from scratch

Submitted by,
Name: Md. Al-Imran Abir
Student ID: 0421062549

Date of Submission: April 17, 2022

# 1 Code:

```python
# Libraries

import numpy as np
import pandas as pd

from sklearn.datasets import load_digits
from sklearn.metrics import calinski_harabasz_score, silhouette_score
from sklearn.cluster import KMeans

import matplotlib.pyplot as plt

#---------- Load data--------------

digits = load_digits()
features = digits.data
labels = digits.target
print(f"Feauters -> type:{type(features)}; shape:{features.shape}")
print(f"Labels   -> type:{type(labels)}; shape:{labels.shape}")

#----------------------------------------
#----------------PCA-------------------
#----------------------------------------

# sample average
mu = np.mean(features, axis=0)
print(f"Shape of mean:{mu.shape}")

# mean centered sample
X = features - mu
# print(np.mean(X, axis=0))

# the covariance matrix
A = np.cov(X, rowvar=False)
print(f"shape of covariance matrix:{A.shape}")
# eigenvalues and right eigenvectors
eig_value, eig_vect = np.linalg.eig(A)

print(f"eigenValues: {eig_value.shape}; type: {type(eig_value)}")
print(f"eigenValues: {eig_vect.shape}; type: {type(eig_vect)}")


# find the value of K for which 95% of energy is retained
k = 1
```

```python
44    energy = 0
45    while energy <= 0.95:
46        eig_value_k = eig_value[:k]
47        energy = eig_value_k.sum()/eig_value.sum()
48        k += 1
49    #     print(energy)
50    print(f"k:{k} -> energy:{energy}")
51
52    # choose the first k eigen vectors
53    U_k = eig_vect[:, :k]
54    print(f"shape of projection matrix: {U_k.shape}")
55
56    # Projection
57    X_proj = X.dot(U_k)
58    print(f"projected matrix shape: {X_proj.shape}")
59
60
61    #-------------------------------------------------------
62    #---------------- K-means clustering-------------------
63    #-------------------------------------------------------
64
65    def compute_cost(X, centroids, cluster):
66        m = X.shape[0]
67        sum = 0
68        for i in range(m):
69            sum += (X[i] - centroids[int(cluster[i])]) @ (X[i] - centroids[int(cluster[i])]).T
70        cost = sum/m
71        return cost
72
73    def my_kmeans(X, k, seed=41):
74        diff = 1
75        np.random.seed(seed)
76        cluster_labels = np.zeros(X.shape[0]) #cluster label for each sample
77        rand_idx = np.random.choice(len(X), k, replace=False)
78        #Randomly choosing Centroids
79        centroids = X[rand_idx, :] #Step 1
80
81        while diff:
82            for i, row in enumerate(X):
83                min_dist = float('inf')
84                for idx, centroid in enumerate(centroids):
85                    d = np.sqrt((row - centroid) @ (row-centroid).T)
86                    if min_dist > d:
87                        min_dist = d
88                        cluster_labels[i] = idx
89            new_centroids = pd.DataFrame(X).groupby(by=cluster_labels).mean().values
```

```python
            if np.count_nonzero(centroids-new_centroids) == 0:
                diff = 0
            else:
                centroids = new_centroids
    cost = compute_cost(X, centroids, cluster_labels)
    cal_score = calinski_harabasz_score(X, cluster_labels)
    return centroids, cluster_labels, cost, cal_score

#-------------------------------------------------
#------------------ With PCA ---------------------
#-------------------------------------------------

def plot_cost(cost, k):
    plt.plot(range(2, k), cost)
    plt.show()

# #### Variation wrt No of Clusters

K = 15
cost = np.zeros(K-2)
sil_score = np.zeros(K-2)
cal_score = np.zeros(K-2)
centroids_lst = []
labels = np.zeros((K-1, X.shape[0]))

for k in range(2, K):
    centroids, labels[k-2], cost[k-2], cal_score[k-2] = my_kmeans(X_proj, k)
    centroids_lst.append(centroids)
    sil_score[k-2] = silhouette_score(X_proj, labels[k-2])
    print(f"k={k} -> cost:{cost[k-2]:.3f};
        calinski:{cal_score[k-2]:.3f}; sil:{sil_score[k-2]:.3f}")

# print(f"silhouette_score:{sil_score}\nCost: {cost}\ncalinski_harabasz_score: {cal_score}")
plot_cost(cost, K)

### Plot cost vs clusters
plt.figure()
plt.plot(range(2, K), cost)
plt.savefig("pca_cost_vs_clusters.png")
plt.xlabel("No of clusters")
plt.ylabel("Cost")
plt.show()

### Plot Calinski-Harabasz Index vs clusters
plt.figure()
```

3

```
136    plt.plot(range(2, K), cal_score)
137    plt.xlabel("No of clusters")
138    plt.ylabel("Calinski-Harabasz Index")
139    plt.savefig("pca_cal_score_vs_clusters.png")
140    plt.show()
141
142    ### Plot Silhouette Coefficient vs clusters
143    plt.figure()
144    plt.plot(range(2, K), sil_score)
145    plt.xlabel("No of clusters")
146    plt.ylabel("Silhouette Coefficient")
147    plt.savefig("pca_cost_vs_sil_score.png")
148    plt.show()
149
150
151    # Variation wrt No of Samples--------------------
152
153
154    ki = 10
155    per = np.linspace(0.1, 1, num=10)
156    z = per.shape[0]
157    cost_sam = np.zeros(z)
158    sil_score_sam = np.zeros(z)
159    cal_score_sam = np.zeros(z)
160    centroids_lst_sam = []
161    labels_lst_sam = []
162    for i, j in enumerate(per):
163        idx = np.random.choice(range(X_proj.shape[0]), int(j*X_proj.shape[0]), replace=False)
164        X_per = X_proj[idx,:]
165    #      print(i, j, X_per.shape[0])
166        centroids_sam, labels_sam, cost_sam[i], cal_score_sam[i] = my_kmeans(X_per, ki)
167        centroids_lst_sam.append(centroids_sam)
168        labels_lst_sam.append(labels_sam)
169        sil_score_sam[i] = silhouette_score(X_per, labels_lst_sam[i])
170        print(f"j={j:.2f} -> cost:{cost_sam[i]:.3f};
171            calinski:{cal_score_sam[i]:.3f}; sil:{sil_score_sam[i]:.3f}")
172
173    plt.plot(per, sil_score_sam)
174
175    ### Plot cost vs sample size
176    plt.figure()
177    plt.plot(per, cost_sam)
178    plt.savefig("pca_cost_vs_per_sample.png")
179    plt.xlabel("% of total sample")
180    plt.ylabel("Cost")
```

```python
181    plt.show()
182
183    ### Plot Calinski-Harabasz Index vs sample size
184    plt.figure()
185    plt.plot(per, cal_score_sam)
186    plt.savefig("pca_cal_vs_per_sample.png")
187    plt.xlabel("% of total sample")
188    plt.ylabel("Calinski-Harabasz Index")
189    plt.show()
190
191    ### Plot Silhouette Coefficient vs smaple size
192    plt.figure()
193    plt.plot(per, sil_score_sam)
194    plt.savefig("pca_sim_vs_per_sample.png")
195    plt.xlabel("% of total sample")
196    plt.ylabel("Silhouette Coefficient")
197    plt.show()
198
199    # Different Initialization----------------------
200
201    seeds = [21, 75, 84, 12, 51]
202    ls = len(seeds)
203    cost_in = np.zeros(ls)
204    cal_in = np.zeros(ls)
205    sil_in = np.zeros(ls)
206    for ij, seed in enumerate(seeds):
207        labels_in, cost_in[ij], cal_in[ij] = my_kmeans(X_proj, 10, seed=seed)[1:4]
208        sil_in[ij] = silhouette_score(X_proj, labels_in)
209        print(f" {seed} -> cost:{cost_in[ij]:.3f};
210            cal:{cal_in[ij]:.3f}; sil:{sil_in[ij]:.3f}")
211
212    #-----------------------------------------------------
213    #------------------- Without PCA --------------------
214    #-----------------------------------------------------
215
216    # Variation wrt No of Clusters ---------------------
217
218    K_w = 15
219    cost_w = np.zeros(K_w-2)
220    sil_score_w = np.zeros(K_w-2)
221    cal_score_w = np.zeros(K_w-2)
222    centroids_lst_w = []
223    labels_w = np.zeros((K_w-1, X.shape[0]))
224
225    for k in range(2, K_w):
```

5

```
226      centroids_w, labels_w[k-2], cost_w[k-2], cal_score_w[k-2] = my_kmeans(X, k)
227      centroids_lst_w.append(centroids_w)
228      sil_score_w[k-2] = silhouette_score(X, labels_w[k-2])
229      print(f"k={k} -> cost:{cost_w[k-2]}; calinski:{cal_score_w[k-2]};
230          sil:{sil_score_w[k-2]}")
231
232  # print(f"silhouette_score:{sil_score}\nCost: {cost}\ncalinski_harabasz_score: {cal_score}")
233  plot_cost(cost_w, K)
234
235  plt.figure()
236  plt.plot(range(2, K_w), cost_w)
237  plt.savefig("no_pca_cost_vs_clusters.png")
238  plt.xlabel("No of clusters")
239  plt.ylabel("Cost")
240  plt.show()
241
242  plt.figure()
243  plt.plot(range(2, K_w), cal_score_w)
244  plt.savefig("no_pca_cal_vs_clusters.png")
245  plt.xlabel("% of total sample")
246  plt.ylabel("Calinski-Harabasz Index")
247  plt.show()
248
249  plt.figure()
250  plt.plot(range(2, K_w), sil_score_w)
251  plt.savefig("no_pca_sil_vs_clusters.png")
252  plt.xlabel("% of total sample")
253  plt.ylabel("Silhouette Coefficient")
254  plt.show()
255
256  # Variation wrt No of Samples------------------------
257
258  ki_w = 10
259  per_w = np.linspace(0.1, 1, num=10)
260  z_w = per_w.shape[0]
261  cost_sam_w = np.zeros(z_w)
262  sil_score_sam_w = np.zeros(z_w)
263  cal_score_sam_w = np.zeros(z_w)
264  centroids_lst_sam_w = []
265  labels_lst_sam_w = []
266  for i, j in enumerate(per_w):
267      idx_w = np.random.choice(range(X.shape[0]), int(j*X.shape[0]), replace=False)
268      X_per_w = X[idx_w,:]
269  #     print(i, j, X_per.shape[0])
270      centroids_sam_w, labels_sam_w, cost_sam_w[i], cal_score_sam_w[i] = my_kmeans(X_per_w, ki_w)
```

```
271        centroids_lst_sam_w.append(centroids_sam_w)
272        labels_lst_sam_w.append(labels_sam_w)
273        sil_score_sam_w[i] = silhouette_score(X_per_w, labels_lst_sam_w[i])
274        print(f"j={j:.2f} -> cost:{cost_sam_w[i]:.3f};
275            calinski:{cal_score_sam_w[i]:.3f}; sil:{sil_score_sam_w[i]:.3f}")
276
277    plt.figure()
278    plt.plot(per_w, cost_sam_w)
279    plt.savefig("no_pca_cost_vs_per_sample.png")
280    plt.xlabel("% of total sample")
281    plt.ylabel("Cost")
282    plt.show()
283
284    plt.figure()
285    plt.plot(per_w, cal_score_sam_w)
286    plt.savefig("no_pca_cal_vs_per_sample.png")
287    plt.xlabel("% of total sample")
288    plt.ylabel("Calinski-Harabasz Index")
289    plt.show()
290
291    plt.figure()
292    plt.plot(per_w, sil_score_sam_w)
293    plt.savefig("no_pca_sil_vs_per_sample.png")
294    plt.xlabel("% of total sample")
295    plt.ylabel("Silhouette Coefficient")
296    plt.show()
297
298    # Different Initialization--------------------------
299
300    seeds_w = [48, 14, 97, 62, 53]
301    ls_w = len(seeds_w)
302    cost_in_w = np.zeros(ls_w)
303    cal_in_w = np.zeros(ls_w)
304    sil_in_w = np.zeros(ls_w)
305    for ij, seed in enumerate(seeds_w):
306        labels_in_w, cost_in_w[ij], cal_in_w[ij] = my_kmeans(X, 10, seed=seed)[1:4]
307        sil_in_w[ij] = silhouette_score(X, labels_in_w)
308        print(f" {seed} -> cost:{cost_in_w[ij]:.3f};
309            cal:{cal_in_w[ij]:.3f}; sil:{sil_in_w[ij]:.3f}")
```
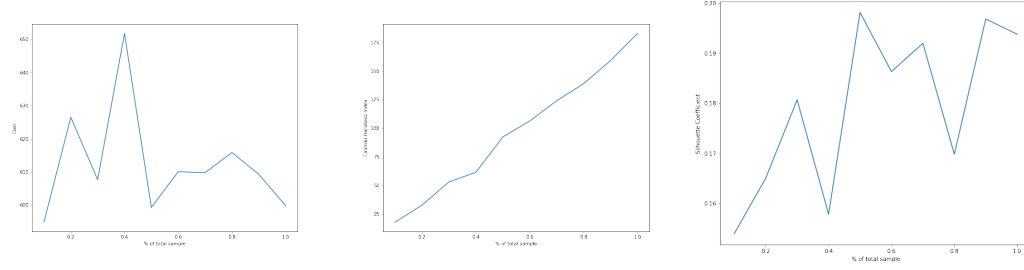
7

# 2 Results

## 2.1 With PCA

### 2.1.1 Variation of clustering performance with number of samples

Cost, Calinski-Harabasz Index, and Silhouette Coefficients while varying the sample size are shown in fig. 1a, fig. 1b, and fig. 1c respectively.
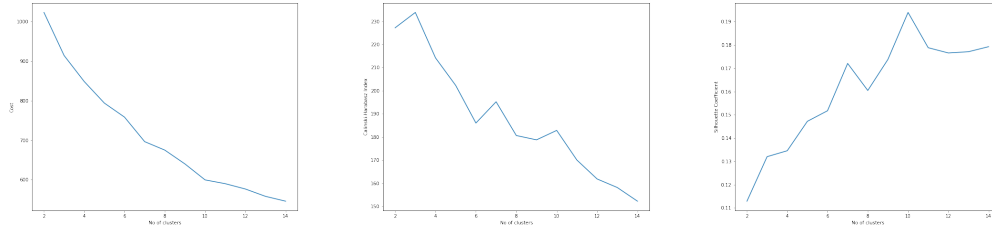


(a) Cost vs no of sample plot    (b) Calinski-Harabasz Index vs no of sample    (c) Silhouette Coefficient vs no of samples plot

Figure 1: Variation of clustering performance with number of samples (with PCA)

### 2.1.2 Variation of clustering performance with number of clusters K

Cost, Calinski-Harabasz Index, and Silhouette Coefficients while varying the no of clusters are shown in fig. 2a, fig. 2b, and fig. 2c respectively.



(a) Cost vs no of clusters plot    (b) Calinski-Harabasz Index vs no of clusters    (c) Silhouette Coefficient vs no of clusters

Figure 2: Variation of clustering performance with number of clusters K (with PCA)

### 2.1.3 Variation of clustering performance with different initializations of the cluster centers

Different initialization was performed by setting the seed of NumPy pseudo-random number generator (`numpy.random.seed(seed)`) to different numbers. The corresponding performance metrics
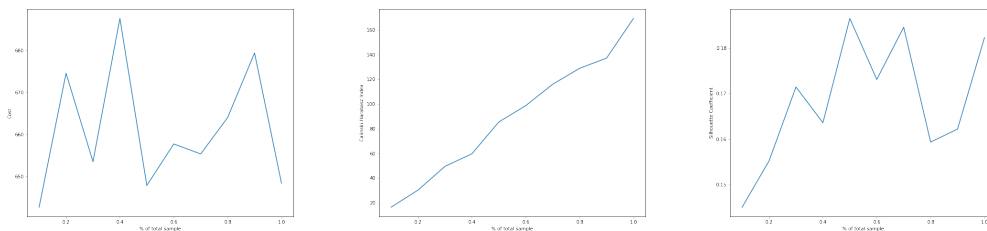
are shown in table 1.

Table 1: Variation of clustering performance with different initializations of the cluster centers (with PCA)

| Seed | Cost | Calinski-Harabasz Index | Silhouette Coefficient |
|------|------|-------------------------|------------------------|
| 21 | 603.331 | 180.672 | 0.195 |
| 75 | 619.554 | 170.742 | 0.182 |
| 84 | 643.371 | 157.071 | 0.160 |
| 12 | 624.812 | 167.634 | 0.179 |
| 51 | 602.507 | 181.190 | 0.198 |

## 2.2 Without PCA

### 2.2.1 Variation of clustering performance with number of samples

Cost, Calinski-Harabasz Index, and Silhouette Coefficients while varying the sample size are shown in fig. 3a, fig. 3b, and fig. 3c respectively.



(a) Cost vs no of samples plot   (b) Calinski-Harabasz Index vs no of samples   (c) Silhouette Coefficient vs no of smaples
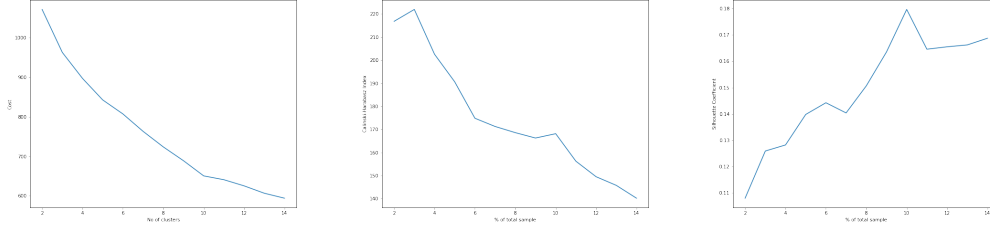
Figure 3: Variation of clustering performance with number of samples (without PCA)

### 2.2.2 Variation of clustering performance with number of clusters K

Cost, Calinski-Harabasz Index, and Silhouette Coefficients while varying the no of clusters are shown in fig. 4a, fig. 4b, and fig. 4c respectively.

### 2.2.3 Variation of clustering performance with different initializations of the cluster centers

Different initialization was performed by setting the seed of NumPy pseudo-random number generator (`numpy.random.seed(seed)`) to different numbers. The corresponding performance metrics are shown in table 2.

(a) Cost vs no of clusters plot

(b) Calinski-Harabasz Index vs no of clusters

(c) Silhouette Coefficient vs no of clusters

Figure 4: Variation of clustering performance with number of clusters K (without PCA)

Table 2: Variation of clustering performance with different initializations of the cluster centers (without PCA)

| Seed | Cost | Calinski-Harabasz Index | Silhouette Coefficient |
|------|------|--------------------------|------------------------|
| 48 | 650.183 | 168.357 | 0.188 |
| 14 | 650.176 | 168.361 | 0.188 |
| 97 | 682.802 | 150.829 | 0.163 |
| 62 | 650.852 | 167.980 | 0.187 |
| 53 | 659.397 | 163.230 | 0.173 |

10

# 3   Discussion

In this project, the K-means clustering algorithm was implemented from scratch. Also, principal component analysis (PCA) was used for dimensionality reduction. We investigated the variation in the performances of the K-means algorithm with respect to different aspects.

## 3.1   Principal Component Analysis

We implemented PCA using the `numpy.linalg.eig()` function which computes the eigenvalues and right eigenvectors of a square matrix. As per the instruction, we also retained 95% energy in the reduced "k" dimensional space. The energy can be calculated by summing the corresponding eigenvalues.

The original dataset had 64 features. To capture 95% of the original energy, we needed only 30 principal features. These 30 features retained 95.48% of the total energy.

## 3.2   K-Means Clustering

Then we performed K-means clustering on both the given data samples and the reduced featured samples. We varied the no of clusters, sample size and observed three metrics, namely the loss function, Calinski-Harabasz Index and Silhouette Coefficient. We also observed the effect of initialization.

### 3.2.1   With PCA

At first, we performed K-means clustering on the reduced dimensional space. As we varied the size of the datasize (i.e. no of samples), we found that Calinski-Harabasz Index increases monotonically with no of samples (fig. 1b) and Silhouette Coefficient has an increasing trend (fig. 1c). But the cost function has no regular shape with respect to the no of sample (fig. 1a).

With the increase in no of clusters, we found that the cost function decreases monotonically as expected (fig. 2a). But there is no clear "knee" in the plot from which we may reach a conclusion about the optimal no of clusters. However, the Silhouette Coefficient (fig. 2c) is highest for 10 clusters and also the Calinski-Harabasz Index (fig. 2b) has a higher value at 10 than the neighbors. So, we may conclude that there are 10 clusters in the dataset. This is also supported by our prior knowledge about the dataset (the dataset has 10 different classes).

We also observe the effect of different initialization on the performance metrics (table 1). Depending on the initialization the performance metrics could vary in a large scale. However for our case, the performance metrics didn't vary that much.

### 3.2.2   Without PCA

Using the full feature space, we found also more or less the same results as with PCA (reduced dimensional space, fig. 3, fig. 4). The performance metrics didn't vary that much. This proves the fact that using PCA, we can achieve similar performances as using the main dataset with the added benefit of reduced computational resources.