

Assignment Title: Unsupervised Learning with Dimensionality Reduction and Clustering

Authors: Koustab Ghosh¹ & Sujoy Kumar Biswas²

Affiliation:

1. Researcher, IDEAS-TIH, Indian Statistical Institute, Kolkata
2. Head of Research & Innovation, IDEAS-TIH, Indian Statistical Institute, Kolkata

Dated: Sep 07th, 2025

We shall work with the MNIST handwritten digits' image dataset. The details about the dataset is available [here](#).

We need the scikit-learn library to import the various machine learning models for our study.

```
from sklearn.datasets import load_digits
digits = load_digits()
digits.data.shape

(1797, 64)
```

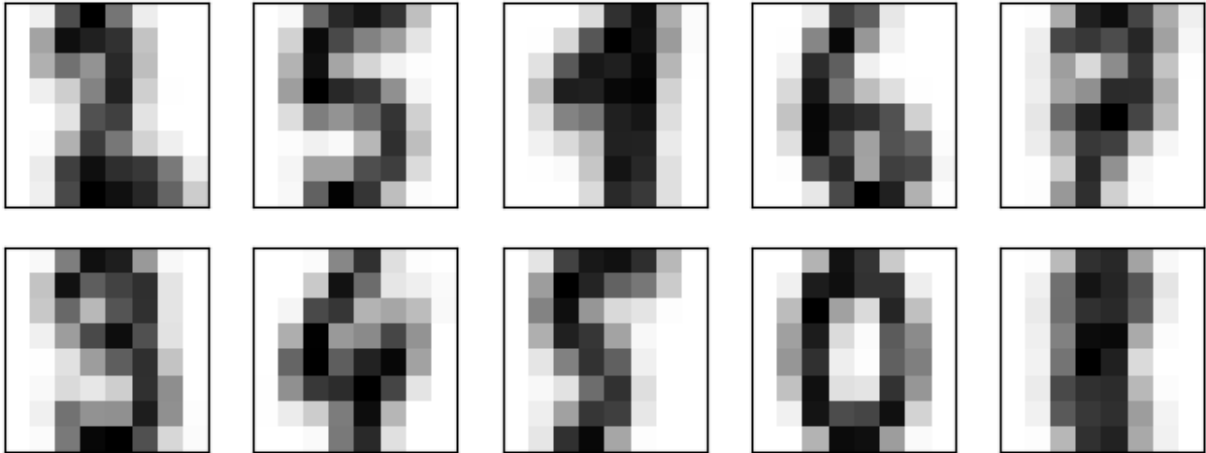
Question 1. Complete the following lines of code for K-Means clustering

```
from sklearn.cluster import KMeans
from matplotlib import pyplot as plt

# kmeans = ...
# clusters = ...
kmeans = KMeans(n_clusters=10, random_state=0)
clusters = kmeans.fit_predict(digits.data)
print(kmeans.cluster_centers_.shape)

(10, 64)

fig, ax = plt.subplots(2, 5, figsize=(8, 3))
centers = kmeans.cluster_centers_.reshape(10, 8, 8)
for axi, center in zip(ax.flat, centers):
    axi.set(xticks=[], yticks=[])
    axi.imshow(center, interpolation='nearest', cmap=plt.cm.binary)
```



We see that even without the labels, KMeans is able to find clusters whose centers are recognizable digits.

Next, we shall apply dimensionality reduction of MNIST handwritten datasets with PCA

```
from sklearn.decomposition import PCA
import matplotlib.pyplot as plt
import numpy as np
```

Question 2. Complete the following lines of code for

Step 1. Dimensionality reduction with PCA. The $8 \times 8 = 64$ dimensional data need to be reduced to 2-dimensional.

Step 2. K-means clustering should be done on the reduced dimensional data. Initial K value should be set to 10 like before.

Step 3. Visualization code is supplied below.

```
# data loading
data = digits.data

# PCA dimensionality reduction
model = PCA(n_components=2)
reduced_data = model.fit_transform(data)

# K-Means clustering with 10 clusters
kmeans = KMeans(n_clusters=10, random_state=0)
clusters = kmeans.fit_predict(reduced_data)
```

Data Visualization

We shall visualize the reduced dimension and cluster data (overlapped) with the help of the following code snippet.

```

# Step size of the mesh. Decrease to increase the quality of the VQ.
h = 0.02 # point in the mesh [x_min, x_max]x[y_min, y_max].

# Plot the decision boundary. For that, we will assign a color to each
x_min, x_max = reduced_data[:, 0].min() - 1, reduced_data[:, 0].max()
+ 1
y_min, y_max = reduced_data[:, 1].min() - 1, reduced_data[:, 1].max()
+ 1
xx, yy = np.meshgrid(np.arange(x_min, x_max, h), np.arange(y_min,
y_max, h))

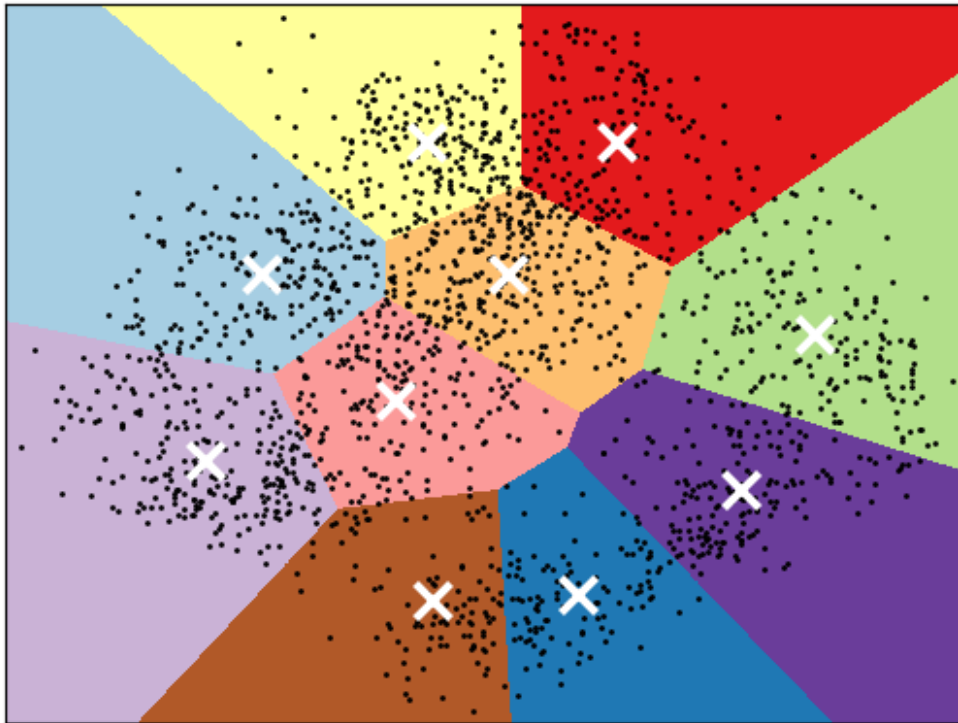
# Obtain labels for each point in mesh. Use last trained model.
clusters = kmeans.predict(np.c_[xx.ravel(), yy.ravel()])

# Put the result into a color plot
clusters = clusters.reshape(xx.shape)
plt.figure(1)
plt.clf()
plt.imshow(
    clusters,
    interpolation="nearest",
    extent=(xx.min(), xx.max(), yy.min(), yy.max()),
    cmap=plt.cm.Paired,
    aspect="auto",
    origin="lower",
)

plt.plot(reduced_data[:, 0], reduced_data[:, 1], "k.", markersize=2)
# Plot the centroids as a white X
centroids = kmeans.cluster_centers_
plt.scatter(
    centroids[:, 0],
    centroids[:, 1],
    marker="x",
    s=169,
    linewidths=3,
    color="w",
    zorder=10,
)
plt.title(
    "K-means clustering on the digits dataset (PCA-reduced data)\n"
    "Centroids are marked with white cross"
)
plt.xlim(x_min, x_max)
plt.ylim(y_min, y_max)
plt.xticks(())
plt.yticks(())
plt.show()

```

K-means clustering on the digits dataset (PCA-reduced data)
Centroids are marked with white cross



Question 3.

Find a high dimensional dataset of your choice. Show how you load the dataset. Do the basic exploratory data analysis to become familiar with the dataset.

Question 4.

Next, the objective would be to reduce the dimension of your dataset and do the clustering on it. Complete the following code for clustering in an object-oriented manner. Do the exact process as above: PCA dimension reduction followed by clustering.

A template code is provided below for your guidance.

```
class YourDataClustering:
    def __init__(self, n_clusters=3):
        self.n_clusters = n_clusters
        self.data = _____
        self.labels = _____
        self.kmeans = _____
        self.scaled_data = _____

    def load_data(self):
        """Load the Iris dataset"""
```

```

        iris = _____._____()
        self.data = iris._____
        return _____

    def preprocess_data(self):
        """Standardize the dataset"""
        scaler = _____._____()
        self.scaled_data = scaler._____(_____)
        return _____

    def apply_kmeans(self):
        """Apply KMeans clustering"""
        self.kmeans = _____._____ (n_clusters=self.n_clusters,
random_state=42)
        self.labels = self.kmeans._____(_____)
        return _____

    def evaluate_clusters(self):
        """Compute silhouette score"""
        score = _____._____ (_____, _____)
        print(f"Silhouette Score: {score:.3f}")
        return _____

    def visualize_clusters_matplotlib(self):
        """Visualize clustering result using Matplotlib"""
        plt.scatter(_____[ :, 0], _____[ :, 1], c=_____,
cmap='viridis')
        plt.title("KMeans Clustering on Iris Dataset (Matplotlib)")
        plt.xlabel("_____")
        plt.ylabel("_____")
        plt.show()

    def visualize_clusters_opencv(self):
        """Visualize clustering result using OpenCV"""
        canvas = np.ones((_____, _____, 3), dtype=np.uint8) * 255
        colors = [(255, 0, 0), (0, 255, 0), (0, 0, 255)]
        scaled = (_____[ :, :2] * 100 + 250).astype(int)

        for i, point in enumerate(_____):
            cv2.circle(canvas, tuple(point), 5, colors[_____ % 3],
-1)

        cv2.imshow("KMeans Clustering (OpenCV)", _____)
        cv2.waitKey(0)
        cv2.destroyAllWindows()

```

The following code executes all the parts of the complete system.

```
# Step 1: Create clustering object
```

```
clustering = _____(n_clusters=3)
```

```
# Step 2: Load dataset
```

```
data = clustering._____()
```

```
# Step 3: Preprocess dataset
```

```
scaled_data = clustering._____()
```

```
# Step 4: Apply KMeans clustering
```

```
labels = clustering._____()
```

```
# Step 5: Evaluate clusters
```

```
score = clustering._____()
```

```
# Step 6: Visualize with Matplotlib
```

```
clustering._____()
```

```
# Step 7: Visualize with OpenCV
```

```
clustering._____()
```